



Data Independent Order Policy Enforcement: Limitations and Solutions

Sarisht Wadhwa
sarisht.wadhwa@duke.edu
Duke University
Durham, NC, United States

Luca Zanolini
luca.zanolini@ethereum.org
Ethereum Foundation
London, United Kingdom

Aditya Asgaonkar
aditya.asgaonkar@ethereum.org
Ethereum Foundation
Sunnyvale, CA, United States

Francesco D'Amato
francesco.damato@ethereum.org
Ethereum Foundation
Berlin, Germany

Chengrui Fang
Chengrui_Fang@zju.edu.cn
Zhejiang University
Hangzhou, China

Fan Zhang
f.zhang@yale.edu
Yale University
New Haven, CT, United States

Kartik Nayak
kartik@cs.duke.edu
Duke University
Durham, NC, United States

Abstract

Order manipulation attacks such as frontrunning and sandwiching have become an increasing concern in blockchain applications such as DeFi. To protect from such attacks, several recent works have designed order policy enforcement (OPE) protocols to order transactions fairly in a data-independent fashion. However, while the manipulation attacks are motivated by monetary profits, the defenses assume honesty among a significantly large set of participants. In existing protocols, if all participants are *rational*, they may be incentivized to collude and circumvent the order policy without incurring any penalty.

This work makes two key contributions. First, we explore whether the need for the honesty assumption is fundamental. Indeed, we show that it is *impossible* to design OPE protocols under some requirements when all parties are rational. Second, we explore the tradeoffs needed to circumvent the impossibility result. In the process, we propose a novel concept of rationally binding transactions that allows us to construct AnimaguSwap¹, the first content-oblivious Automated Market Makers (AMM) interface that is secure under rationality. We report on a prototype implementation of AnimaguSwap and performance evaluation results demonstrating its practicality.

CCS Concepts

• Security and privacy → Distributed systems security.

¹A key design in AnimaguSwap is that user orders may *transform* to a different direction—like the fictional creatures Animagi in Harry Potter—in order to achieve the desired game theoretic properties.



This work is licensed under a Creative Commons Attribution International 4.0 License.

Keywords

Blockchain, MEV, Cryptoeconomics

ACM Reference Format:

Sarisht Wadhwa, Luca Zanolini, Aditya Asgaonkar, Francesco D'Amato, Chengrui Fang, Fan Zhang, and Kartik Nayak. 2024. Data Independent Order Policy Enforcement: Limitations and Solutions. In *Proceedings of the 2024 ACM SIGSAC Conference on Computer and Communications Security (CCS '24)*, October 14–18, 2024, Salt Lake City, UT, USA. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3658644.3670367>

1 Introduction

Blockchains can provide a trustworthy platform for transacting and smart contract execution. Blockchain-powered finance applications, also known as DeFi, have grown to a market of more than \$46 billion² in value. However, despite the strong integrity and availability properties offered by blockchains, they do not protect the *ordering* of user transactions. As a result, order manipulation attacks — e.g., frontrunning attacks, sandwich attacks — are rampant, where an attacker listens for user transactions sent in public and strategically places her exploiting transactions around the victim to gain a profit. The profits earned through inserting and reordering transactions are referred to as Maximal Extractable Values (MEV) [13]. An estimated \$1.2B of MEV has been extracted as of the time of writing.³

To protect users from order manipulation attacks, an extensively explored direction [2, 6, 9, 10, 22–24, 28] is to design protocols that enforce certain “fair” transaction ordering policy. A popular approach is *data-independent* ordering, which guarantees that given a set of user transactions as input, the final ordering of them on the blockchain should be independent of the transaction content. For example, some fair ordering protocols [9, 22, 23] order user transactions based on the *time* they are received by a *committee* of parties. Content-oblivious ordering (e.g., [2, 6, 24, 28]) guarantees that user transactions are hidden from the committee who orders them, e.g., through encryption until after an ordering has been

²<https://defillama.com>

³<https://explore.flashbots.net/>

decided. In this case, transaction ordering may be based on any metadata, such as the ciphertext, the sender address, etc.

Both approaches can prevent an attacker from placing exploiting transactions before user transactions *after* having observed user transactions. However, all known data-independent ordering protocols share the same limitation: they only work under the strong assumption that enough parties running the protocol are honest. E.g., [22] assumes more than three-fourths of the participants are honest (for $\gamma = 1$, a parameter in their work).

Indeed, in a permissionless blockchain system where players are pseudonymous and can join and leave freely, the assumption that players are always honest is hard to justify. A much more palatable assumption is to assume rationality instead of honesty, i.e., instead of assuming parties are intrinsically honest, a rational party may take any action to maximize utility. In fact, the existence of MEV is tied to the rationality of the participants. Thus, the goal is to design a protocol so that following the protocol is incentive-compatible, which is significantly more challenging because all of the parties running the protocol may deviate from the protocols arbitrarily if doing so leads to a higher utility.

In this paper, we systematically investigate the design of data-independent ordering protocols in the presence of rational parties, asking two fundamental questions:

- (1) All known data-independent ordering protocols require some honesty assumption. Is that a limitation of existing solutions or something fundamental? We answer this question negatively by showing an impossibility result that not only are existing protocols insecure in the presence of rational parties, but a wide range of protocols compliant with the same specification also cannot be secure.
- (2) Given the impossibility, what tradeoffs must one make in order to realize a data-independent ordering protocol under the rationality assumption? We propose a novel concept called *rationality binding commitments* and present the first decentralized exchange construction, called AnimaguSwap, with a built-in data-independent ordering protocol under rationality.

1.1 Overview of results

1.1.1 Existing protocols are not secure. Intuitively, it is not hard to see how rational parties might lead to an insecure execution: in existing data-independent ordering protocols [2, 6, 9, 10, 22–24, 28], there is no way to retroactively verify whether the ordering output was indeed data-independent. Thus, if violating data independence increases parties' utility, all parties running the protocol to collude is a dominant strategy.

For fair ordering protocols, if enough parties collude, they can order transactions arbitrarily by lying about when transactions are received — an action that cannot be held accountable unless assuming a global trustworthy timestamping service (which is a strong assumption for applications we care about).

The situation is a bit trickier for content-oblivious ordering, as collusion *might be accountable*. For example, in schemes where user transactions are encrypted using threshold encryption (e.g., [10]), enough parties can reconstruct the decryption key if they collude.

However, this way of colluding may be accountable since the decryption key itself could serve as irrefutable proof of the fact that collusion has taken place.

This leads to a natural question for a protocol designer: can we leverage proof of collusion to design data-independent ordering protocols under rationality?

Answering this question negatively and identifying the conditions under which this is true is the crux of our first contribution. We observe that colluding parties do not necessarily need to decrypt the transaction or leave any proof of collusion whatsoever, by running the collusion algorithm in a way that the only outcome of collusion is a set of transactions that resemble benign user transactions while giving colluding parties a higher utility (e.g., with their frontrunning transactions inserted before the victim). We emphasize that the cost to collude between parties is very low since today's blockchain landscape is not so decentralized, and a few of the pools interacting with each other are all required to attack and collude. Further, once parties collude, they can profit for a longer duration, which decreases the amortized cost.

An order policy enforcement (OPE) framework and impossibilities. To prove this claim, we first present a generic framework in Section 4.1 that captures all known data-independent ordering protocols. Then, we show that in any concrete protocol Π following this framework, if violating the ordering policy increases parties' utility, there always exists a collusion protocol π with which parties can collude and violate the ordering policy with deniability: even after executing π , no participants of it can generate a cryptographic proof to incriminate any participants (including herself).

Section 4.2 presents the full proof.

1.1.2 New directions informed by the impossibility. Our impossibility proof not only shows the fundamental limitations of existing approaches in achieving security under rationality, but it also carves out avenues to improve. The impossibility critically relies on some assumptions about Π . First, users may go offline after sending one message (typically a transaction or a cryptographic commitment thereof). This is a desirable usability feature because users do not need to stay online. Consequently, once the user submits her transaction, the parties have the capability to retrieve it. Second, if the user sends a cryptographic commitment of her transaction, it is *binding* in that the commitment can only be opened to one transaction plaintext, which is a natural requirement so that transaction execution is unambiguous.

Designing protocols that violate these assumptions can circumvent the impossibility, but dispensing with them naively will lead to undesirable outcomes. For instance, if we require users to stay online, there exists a (somewhat trivial) solution where a user first sends a commitment to the parties running the ordering protocol, and then opens the commitment after the ordering is determined. This construction, while secure against collusion of parties, not only introduces a usability challenge for users but also potential problems when users refuse to open the commitment.

Introducing rationally binding commitments. Our next result is a novel way to relax the second assumption by introducing *rationality binding commitments*. A key observation from the impossibility proof is that if a user only sends one message and that

message binds to her transaction, then if enough parties collude, they can uniquely recover her transaction (and thus can frontrun it, for example), no matter what cryptographic protections are employed. (Since the user only sends one message, that message should enable recovery of some transactions; due to the binding property, the committee can recover the exact transaction the user committed to). Can we dispense with the binding property as a way to circumvent the impossibility? This seems paradoxical. After all, a user's transaction needs to be encoded somehow in the commitment, otherwise, the commitment may be opened against her will. Our answer is to replace binding with *rationally binding*, as follows.

We first require parties running the protocol to put down collateral (i.e., to *stake*) that can be confiscated (or slashed) for detected misbehavior. We call these parties stakers hereafter. Suppose one of the stakers is designated as the “flipper” (the meaning of the name will become clear momentarily). In order to create a rationally binding commitment to a transaction tx , the user samples a random bit $b \in \{0, 1\}$ and depending on this bit, creates a transaction that is either the one that the user intended (tx) or a related but different transaction (\bar{tx}), e.g., the other transaction must satisfy certain requirements that we will specify later for specific applications. The user sends b to the flipper in a deniable message [20, 29] and gets back an acknowledgment of the bit signed by the flipper. (If the flipper does not respond, the user can designate another flipper.) The user then shares the created transaction (which can be different from the one it intended) with the rest of the stakers. To open the commitment, the stakers reveal the shared transaction, and the flipper reveals b , and the transaction tx will be executed. Crucially, if the flipper reveals the wrong bit \bar{b} , the user can use the acknowledgment it received as evidence to slash the flipper.

From the user's point of view, assuming the penalty is appropriately set, a rational flipper will always reveal b , so tx will always be executed, similarly to the binding property. On the other hand, from the stakers' point of view, even if all parties collude, they cannot identify which transaction will be executed since the flipper might lie about b , and there is no way for the flipper to prove the correctness of b due to the use of deniable messaging. In fact, the protocol can be made such that lying about b is a dominant strategy for the flipper by carefully crafting \bar{tx} , which ensures that no stable collusion can be formed amongst the stakers.

In Section 5.1, we present AnimaguSwap, an Automated Market Makers (AMM) decentralized exchange that uses rationally binding commitments to defeat sandwich attacks, assuming buying and selling a token is equally likely. In our protocol, if user transaction tx sells a certain asset, then \bar{tx} is the reverse order, i.e., buying the same asset. If the stakers collude, they must still guess which will be executed (with no more than a $1/2$ probability of being correct). Thus, in expectation, it is not worthwhile to attempt sandwiching.

In Section 5.2, we provide a game theoretic analysis of AnimaguSwap and show that following the protocol specification is dominant for all the involved parties. To evaluate practicality, in Section 5.3, we implement a base AnimaguSwap as a smart contract and show that the key overhead in terms of the amount of gas is 1.3x compared to a typical insecure trade today. Moreover, since this cost does not depend on the number of stakers or the value of the transaction, this is already quite practical for high-value transactions. We then extend the result in Section 5.4 to consider scenarios

where the buying and selling of a token may not be equally likely. We further enhance the security of the game to cover repeated games for a more practical solution by ensuring the unlinkability between different games.

Contributions. In summary, this paper makes the following contribution:

- We present a framework that captures existing protocols for data-independent order policy enforcement (OPE), such as fair ordering and content-oblivious ordering protocols.
- We present an impossibility proof showing that a wide range of OPE protocols cannot be secure when all parties are rational
- We propose the notion of rationally binding commitments as a practical way to circumvent the impossibility. We present the first AMM interface construction AnimaguSwap, that can achieve data-independent ordering of user transactions in the presence of rational parties. We analyze the efficacy of AnimaguSwap by a game-theoretic proof in the presence of rational parties.
- We implement AnimaguSwap using a smart contract and show that the overhead of security is about 1.3x in gas cost compared to vanilla UniSwap; this will be practical for high-value transactions.

2 Related Work

Data-independent ordering protocols. As reviewed in Section 1, several works purpose to order transactions independent of their content as a way to reduce MEV [13]. Below is a non-exhaustive list of protocols that are covered by the framework (Section 4.1) and the impossibility theorem (Theorem 1).

The first category of protocols is fair ordering. Kelkar *et al.* [23] investigate a notion of *fair transaction ordering* for (permissioned) consensus protocols, which prevents adversarial manipulation of the ordering of transactions. The authors then formulate a new class of consensus protocols, called Aequitas, that achieve fair transaction ordering while also providing the usual consistency and liveness. Their findings have been later extended in permissionless settings [21]. Subsequently, Kelkar *et al.* [22] devised Themis, a (permissioned) consensus protocol that, along the same lines as [23], achieves fair transaction ordering while preventing a liveness issue in Aequitas. Cachin *et al.* [9] introduce a *differential order fairness* property and present a quick order-fair atomic broadcast protocol that guarantees payload message delivery in a differentially fair order. The protocol of Cachin *et al.* results in a more efficient protocol than the previous solutions, but it relies on a weaker form of validity property.

The second category of solutions is content-oblivious ordering. A popular idea (used by, e.g., [2, 6, 10, 28]) is to encrypt user transactions using a threshold public key encryption scheme so that the ordering of transactions is done based on the ciphertext. Fino [24] efficiently integrates threshold encryption and secret sharing to a DAG-based BFT protocol. Shutter, Osmosis, and Sikka [2, 10, 28] are examples of operational systems in this category.

The protocols in these works make the assumption of honest majority participation, e.g., a majority (or two-thirds) of the participants do not deviate from the specified protocol, even if such deviations are undetectable. Our work investigates ways to relax such assumptions.

MEV mitigation leveraging rationality. Platforms have emerged to auction off the opportunities to extract MEV so that MEV extraction is democratized [15, 16]. MEV auctions rely on the rationality of bidders (or builders) to maximize MEV extraction. Our solution (in Section 5) aims to achieve a different goal of reducing MEV.

Heimbach *et al.* [19] analyzed the sandwich game between an AMM trader and predatory bots and identified the optimal slippage tolerance a trader could set to disincentive bots from attacking while limiting the probability of execution failures. Their algorithm crucially relies on estimating the execution failure probability using historical data and thus cannot guarantee accuracy. Our solution is fundamentally different and does not have this limitation.

PROF [3] is a protocol that leverages the profit-maximizing nature of proposers to promote the inclusion of fairly ordered transactions (PROF defines fairness broadly as any order that follows a given policy). PROF is agnostic to specific transaction ordering protocols and, thus, is complementary to our solution. Note that PROF does not address ordering under rationality, though it suggested a TEE-based content-oblivious ordering protocol.

Lower bounds on MEV mitigation. Ferreira *et al.* [32] presents an impossibility result showing that for a class of liquidity pool exchanges (e.g., Uniswap), for any data-dependent ordering policy (called sequencing rule in [32]), there are always valid sequences in which the miners get risk-free profits. Their result leaves it open whether data-independent ordering policies can be enforced, which is our focus. (We show it is impossible in Theorem 1.)

Data dependent ordering policies. All of the discussion in this work is only pertinent to ordering policies that are *data independent*, i.e., policies that only rely on the metadata related to the transactions and not the transaction content themselves. [32] proposes a *data-dependent* sequencing rule that alternates between BUY and SELL orders to guarantee that user transactions are executed at a price as good as being executed at the beginning of the block (unless the miner does not gain anything from manipulating the ordering).

Moreover, [32] relies on the assumption that each block is created by a different miner, a questionable assumption in today's Ethereum ecosystem with Proposer-Builder Separation (PBS) [8], which our solution (Section 5.1) avoids.

3 Model and Problem Statement

Throughout this paper, we consider *data-independent transaction ordering protocols* run by a set of N parties called *stakers* $\mathcal{S} = \{s_1, \dots, s_N\}$. Such protocols process transactions submitted by *users* and output an ordered list of received transactions while ensuring that the ordering is independent of the transaction content. Examples include fair ordering based on receive order [9, 22, 23] and content-oblivious ordering (e.g., [2, 6, 24, 28]). The purpose of data-independent ordering is to prevent order manipulation attacks such as frontrunning attacks, sandwich attacks, etc. We refer readers to [33] for a survey of such attacks.

We assume all users, including stakers, to be *rational*, i.e., they act to maximize their utility function. To keep things simple, we assume that this utility function is the amount of monetary profit (in the number of *tokens*) that the party can make. If a staker s_i fails

to serve the role assigned to it or tries to deliberately deviate from the protocol, i.e., s_i is *Byzantine*, and a proof of this misbehavior is given, it loses a part of its stake (s_i gets *slashed*), and it might be removed from the system. A protocol specifies rules that provide rewards to stakers who complete certain tasks. We sometimes refer to users (and stakers) as *players* or *parties*.

Adversary model. Stakers are adversarial, and they may deviate from the protocol arbitrarily if doing so increases their utility (after counting the penalty, if any). Their goal is to tamper with the ordering process so that transactions are ordered to their advantage. For example, in receive order-based fair ordering protocols, stakers may collude and order a later transaction before an earlier one to facilitate a frontrunning attack; in content-oblivious ordering protocols, stakers may collude to decrypt user transactions and profit from the information thereof.

Problem statement. We ask two questions: First, existing data-independent ordering protocols are insecure under the above rationality assumption. Is this limitation fundamental or can it be mitigated? We answer this question negatively with an impossibility result. Second, given the impossibility, what relaxation of the problem can we make to obtain a practical data-independent ordering protocol under the above rationality assumption?

Notation. We denote the evaluation of a protocol using $(pub_o; (y_1, \dots, y_k) \leftarrow \text{prot}(pub_i; (x_1, \dots, x_k)))$. Here, there is a public input pub_i and private inputs (x_1, \dots, x_k) , resulting in a public output pub_o and private outputs (y_1, \dots, y_k) . Public inputs/outputs might be omitted if not applicable.

4 Impossibility of OPE under rationality

To study the common features of data-independent order policy enforcement (OPE) protocols [2, 6, 9, 10, 22–24, 28], we first present an abstract framework to capture the essence of aforementioned protocols with four sub-protocols (submit, process, order, reveal) and two predicates ShouldRelease and ShouldReveal. To aid understanding, we show how existing schemes can be mapped to our framework.

4.1 Framework for Order Policy Enforcement

Parties, transactions, and ordering policies. Our framework is run by *users*, who submit transactions, and a set of *stakers*, who execute the ordering protocol to order submitted transactions. Stakers' protocol can either be a component of a larger consensus protocol or a standalone protocol in parallel with the consensus (e.g., on layer 2).

Definition 1 (Data and Metadata). *A transaction tx_i can be considered to consist of two parts – metadata md_i and data $data_i$. Metadata is defined as the part of a transaction not given to the application (i.e., a smart contract) for execution. Data is defined as the part of a transaction that is required for application execution.*

Our framework defines a generic protocol to enforce a data-independent policy \mathcal{P} .

Definition 2 (Data-independent Policy). *A policy is defined as data-independent if it takes as input a set of metadata (one for each*

Framework for Order Policy Enforcement	
Initialization:	
1:	Each staker s_i runs init (possibly interactively with other stakers) to get $\text{param}_i := (\text{spri}_i, \text{spp}_i)$
2:	Each staker s_i publishes spp_i
3:	Each staker s_i sets $\text{state}_i := \emptyset$
Transaction submission:	
4:	Whenever initiated by a user u , stakers in S and u run (possibly interactively)
	$(\text{txid}; (\perp, \text{out}_1, \dots, \text{out}_N)) \leftarrow \text{submit}(\text{tx}, \text{inp}_1, \dots, \text{inp}_N)$
	where tx is user's input (her transaction) and inp_i is staker s_i 's input derived from param_i and state_i .
5:	Each staker s_i processes the metadata information and the transaction information and updates its state.
	$(\text{md}_i, \text{data}_i) \leftarrow \text{process}(\text{txid}, \text{out}_i, \text{state}_i)$
	$\text{state}_i \leftarrow \text{state}_i.\text{add}((\text{txid}, \text{md}_i, \text{data}_i))$
Transaction inclusion:	
6:	Whenever $\text{ShouldRelease}(s_i)$, stakers in S evaluate
	$(\text{tSeq} = (\bar{\text{tx}}_1, \dots, \bar{\text{tx}}_\ell); (\text{state}_1, \dots, \text{state}_N)) \leftarrow$
	$\text{order}(\text{state}_1, \dots, \text{state}_N)$
	where the order of $(\bar{\text{tx}}_1, \dots, \bar{\text{tx}}_\ell)$ is dependent only on $\text{md}_1, \dots, \text{md}_\ell$.
7:	Staker s_i adds tSeq to the blockchain.
Transaction revealing:	
8:	For each $k \in [\ell]$, when $\text{ShouldReveal}(\bar{\text{tx}}_k)$, stakers evaluate
	$(\text{tx}_k; (\text{state}_1, \dots, \text{state}_N)) \leftarrow \text{reveal}(\bar{\text{tx}}_k;$
	$(\text{state}_1, \text{spri}_1), \dots, (\text{state}_N, \text{spri}_N))$

Figure 1: A general framework that captures proposed ordering policy enforcement protocols [6, 9, 10, 22–24] using four protocols (submit, manipulate, order, reveal) and two predicates ShouldRelease , ShouldReveal .

transaction) and outputs one or more permutations of transactions associated with them, i.e., $\mathcal{P}(\text{md}_1, \dots, \text{md}_\ell) \subseteq \sigma(\ell)$, where $\sigma(\ell)$ is the set of all permutations of $(\bar{\text{tx}}_1, \dots, \bar{\text{tx}}_\ell)$.

Generally, each staker may have some different metadata for a given transaction, thus $\text{md}_i = (\text{md}_i^1, \dots, \text{md}_i^N)$ represents the metadata for transaction tx_i across all N stakers.

The framework. As shown in Fig. 1, the framework for order policy enforcement consists of four sub-protocols. These subprotocols are reactive in that they are activated when specific conditions are met and may execute in parallel to each other. We now describe the four subprotocols following the life cycle of a given transaction, although note that these subprotocols are reactive and may execute in parallel for different transactions.

- Stakers engage in an initialization protocol to generate a parameter $\text{param} = (\text{spri}, \text{spp})$ that consists of secret parameters spri and public ones spp . Initialization will also set a local variable, state_i — the set of pending transactions with metadata, to \emptyset .
- First, to send a transaction tx to a blockchain, the user runs the submit protocol with stakers. Specifically,

$$(\text{txid}; (\perp, \text{out}_1, \dots, \text{out}_N)) \leftarrow \text{submit}(\text{tx}, (\text{inp}_1, \dots, \text{inp}_N))$$

, where inp_i and out_i are the input (output) from (to) staker s_i , and txid is an id identifying the transaction. We do not restrict how submit may be realized, e.g., it can be realized as a non-interactive protocol where the user simply encrypts the transaction under stakers' public keys (in which case $\text{inp}_i = \text{pk}_i$); submit may also be implemented with an interactive Multi-Party Computation (MPC) protocol where the user engages in MPC protocol with stakers (in this case inp_i might be secret). At the end of submit, each staker s_i receives some information about tx in out_i , which will be used in later protocols. Note that not all stakers may be required to participate in submit; however, a minimum of t_s is required ($1 \leq t_s \leq N$). For the stakers that do not participate, the input and output are \perp .

Users are ephemeral, i.e., they may go offline after running submit, a usability feature enjoyed by most real-world systems[6, 9, 10, 22–24]. Consequently, $(\text{txid}; (\perp, \text{out}_1, \dots, \text{out}_N))$ together must contain enough information to recover tx , an observation that will play a critical role in our subsequent analysis. We discuss alternative protocols if this assumption does not hold in Section 5.

We also assume w.l.o.g. that non-staker users submit their transactions before a staker adds its own, considering all the information revealed to it by the non-staking users.

- Having finished the submit protocol for a given tx , a staker runs a local process function to capture any local state to be used in later sub-protocols, e.g., the time at which tx was received. Specifically, $(\text{md}_i, \text{data}_i) \leftarrow \text{process}(\text{txid}, \text{out}_i, \text{state}_i)$.
- The goal of an OPE protocol is to produce blocks with transactions ordered in a desirable way. In our framework, whenever predicate $\text{ShouldRelease}(s_i)$ is true, stakers will run the order protocol, with s_i being the leader if applicable, to order transactions and to output a sequence of transactions. Specifically, let $\text{tSeq} = (\bar{\text{tx}}_1, \dots, \bar{\text{tx}}_\ell)$

$$(\text{tSeq}; (\text{state}_1, \dots, \text{state}_N)) \leftarrow \text{order}(\text{state}_1, \dots, \text{state}_N)$$

where each staker inputs its local set of pending transactions (with any metadata captured in process). The output is a sequence of transactions to be added to the blockchain and an updated local variable (e.g., with transactions added to the block removed).

Note that, like in submit, not all stakers may be required to participate in order; however, a minimum of t_o is required ($1 \leq t_o \leq N$). For the stakers that do not participate, the input and output are \perp . These stakers would appropriately need to change state according to the on-chain published ordering of transactions.

This sub-protocol captures any multiparty computation mandated by an ordering protocol, e.g., fair ordering schemes generate the contents of the next block based on timestamps (or relative receiving orders) across all stakers.

- In some protocols, order only includes some cryptographic representation of transactions in the blockchain, and another step reveal is required to reveal the transaction plaintext so it can be executed. Whenever $\text{ShouldReveal}(B)$ is true, stakers will run reveal to reveal transactions in B .

Again, not all stakers may be required to participate in reveal; however, a minimum of t_r is required ($1 \leq t_r \leq N$).

We use $tSeq \models (tx_1, \dots, tx_\ell)$ to represent that if $tSeq$ is posted on-chain after order then the reveal execution would correspond to (tx_1, \dots, tx_ℓ) .

Requirements. To rule out trivial or impractical constructions, our framework makes the following assumptions.

First, we require $submit(tx, \cdot)$ to be binding to the given transaction tx in that if $(_, out_1, \dots, out_N) = submit(tx, \cdot)$, then $(tx; \cdot) \leftarrow reveal(\bar{tx}; \cdot)$. All practical blockchain systems do achieve this.

Second, we require that a submitted transaction is eventually included in the blockchain, and revealed, if applicable. This is the standard liveness property.

Third, we note that as expressed in the framework, the function $reveal()$ takes as input the output of the function $order()$ and the static private parameter in $spri$. Thus, we assume the protocols and the predicates in the interim do not affect the inputs to the function $reveal$, and thus, the function $reveal$ can be run any time after order (even before staker s_i adds the output block to the blockchain). This implies our framework does not apply to protocols that use cryptographic primitives that change state of a transaction between order and reveal such as by using time-locked encryption [27] or witness encryption [17]. These primitives are not widely used due to their practical limitations (e.g., it is hard to calibrate the timeout in time-lock encryption, and decrypting a timelock encrypted ciphertext requires constant computation; there is yet no practical witness encryption schemes[17]).

Examples. In [31, Appendix A], we show that our framework can capture OPE protocols based on DKG [6, 10], secret-sharing [24], as well as fair ordering protocols [9, 22, 23].

4.2 Delineating Impossibility Conditions for Data Independent Ordering

Existing data independent order policy enforcement (OPE) protocols order transactions under the assumption that a fraction (less than one-third or one-half) of stakers are Byzantine and the remaining stakers are honest. However, in practice, the motivation to introduce additional transactions, delete existing transactions, or to order transactions differently is to obtain higher monetary gains for the stakers. Thus, a model where all stakers are rational and maximizing their utility (in terms of monetary gains) captures the adversarial setting better. In this section, we analyze OPE protocols under such an adversary. In particular, we show that under some circumstances, there exists an attacking strategy where we can ensure that rational stakers *do not* follow the OPE protocol. The key challenge is in identifying the conditions under which this statement holds, and showing the resulting attacking strategy. Recalling the notations defined in Section 3, our result can be stated as follows:

Theorem 1. *Let Π be a protocol that follows the ordering policy enforcement framework (Fig. 1) to enforce a data-independent policy \mathcal{P} , and let \mathcal{S} be the set of rational stakers executing Π . Suppose there exists a sequence of transactions $tSeq = \{tx_1, \dots, tx_\ell\} \in \mathcal{P}(md_1, \dots, md_\ell)$ with max utility for some input stream $((md_1, data_1), \dots, (md_\ell, data_\ell))$. Moreover, let us assume that there exists a function $extract()$ known to all stakers in \mathcal{S} s.t. $tSeq' \models extract(tx_1, \dots, tx_\ell) \in \mathcal{P}(md'_1, \dots, md'_\ell)$ where tx_i corresponds to the reveal of \bar{tx}_i , for another set of valid*

md'_1, \dots, md'_ℓ ; such that the utility from publishing $tSeq'$ is more than the utility from publishing $tSeq$. Then, Π cannot enforce \mathcal{P} .

In other words, assuming MEV extraction is possible (i.e., $extract$ exists), data-independent ordering policies cannot be enforced by protocols following the ordering policy enforcement framework defined in Fig. 1. The necessary $extract$ function, in practice, can be an algorithm that uses a combination of techniques publicly known to stakers today and outputs the sequence that produces the highest utility.

To prove the above impossibility result we present an attacking protocol (Algorithm 1), and show that the stakers can present a different reality $tSeq'$ where no proof of malice can be obtained.

Suppose an adversarial set of stakers \mathcal{A} ($|\mathcal{A}| \geq \max(t_s, t_o, t_r)$), such that \mathcal{A} is able to run $submit$, $order$, $reveal$) want to attack, they will run Algorithm 1 using a protocol in a Trusted Execution Environment (such as Intel SGX) when $ShouldRelease(s_i)$ is true (and skip the honest protocol). Such an algorithm in TEE is described in [31, Appendix B]. Note that we use an algorithm that provides deniability to the stakers. Stakers in \mathcal{A} ($s_i \in \mathcal{A}$) will provide inputs to the TEE running Algorithm 1, which will release any output bit-by-bit to ensure all parties receive the output [7, Sec 5.4].

Note that all computations except the final outputs are hidden during the execution and not available to any party in the clear. Given ℓ received outputs (each one submitted by a user u_i for a transaction tx_i), and given a list $spri^a$ of inputs $spri_i$ of stakers $s_i \in \mathcal{A}$, an ordered list of transactions $tSeq = (\bar{tx}_1, \dots, \bar{tx}_\ell)$ is generated (Line 5). Then, the $reveal$ function is computed by the stakers in \mathcal{A} by passing as inputs the previously generated list of ordered transactions, the list $state^a$ of states $state_i$ of stakers $s_i \in \mathcal{A}$, and $spri^a$. Once the transactions tx_i are available, transaction signatures are verified in order to confirm that each member provided the correct input to the protocol. Next, the $extract$ function is run (Line 10) in order to introduce new transactions att_txn (Line 12), which are then submitted (Line 13) and added in the local state (Line 15). The resulting block containing MEV-extracting transactions is then published (Line 17).

At a high level, the above construction of an attacking protocol works because i) $tSeq'$ is more profitable for the stakers than $tSeq$, and thus they are incentivized to join the coalition and ii) no party can prove that the coalition of stakers was formed to violate the ordering policy, and thus cannot be penalized.

We prove this formally in [31, Appendix C]. We show an example attack that follows the attack protocol Algorithm 1 in [31, Appendix D].

5 OPE using Rational Binding Commitments

In the impossibility result in the previous section, we assumed that given a sequence of transactions $tSeq$, the parties have access to an $extract()$ function that provides a higher utility. For existing systems such as Ethereum, such MEV extraction strategies are known for sequences of transactions such as sandwich attacks [34], frontrunning [14, 25], arbitrage [14] etc. To make them work in the attack in Algorithm 1 (where $tSeq$ is available but not in the clear), we can create an $extract()$ circuit that attempts all known attack strategies and applies them to $tSeq$, and picks the best among them to produce a new sequence $tSeq'$.

Algorithm 1 Protocol for a set \mathcal{A} of stakers extracting an ordering with a higher utility (protocol for $s_i \in \mathcal{A}$)

```

1:  $state_j^a \leftarrow$  if  $s_j \in \mathcal{A}$  then  $state_j$  else  $\perp$ 
2:  $inp_j^a \leftarrow$  if  $s_j \in \mathcal{A}$  then  $inp_j$  else  $\perp$ 
3:  $spri_j^a \leftarrow$  if  $s_j \in \mathcal{A}$  then  $spri_j$  else  $\perp$ 
4: procedure  $ATTACK^K(state^a, spri^a, inp^a)$ 
5:    $(tSeq = (\bar{tx}_1, \dots, \bar{tx}_\ell), state^a) \leftarrow order(state^a)$ 
6:   for  $j \in \{1, \dots, \ell\}$  do
7:      $(tx_j; state^a) \leftarrow reveal(\bar{tx}_j, state^a, spri_j^a)$ 
8:    $B = (tx_1, \dots, tx_\ell)$ 
9:    $VerifySigs(B)$ 
10:   $att\_B \leftarrow extract(B)$ 
11:   $state' \leftarrow \perp$ 
12:  for  $att\_txn \in att\_B$  do
13:     $(txid; (\perp, out_1, \dots, out_N)) \leftarrow submit(att\_txn, inp^a)$ 
14:     $md_i, data_i \leftarrow process(txid, out_i, state'_i)$ 
15:     $state'_i \leftarrow state'_i.add((txid, md_i, data_i))$ 
16:   $(tSeq' = (tx'_1, \dots, tx'_{\ell'}); state') \leftarrow order(state')$ 
17:  return  $tSeq', state'_i$ 

```

$\triangleright state^a$ is a list of states $state_j$ for every $state_j \in \mathcal{A}$
 $\triangleright inp^a$ is a list of inputs inp_j for every $state_j \in \mathcal{A}$
 $\triangleright spri^a$ is a list of private inputs $spri_j$ for every s_j in \mathcal{A}
 \triangleright Executed when $ShouldRelease(s_i)$ is true
 \triangleright Validators in \mathcal{A} order ℓ transactions
 \triangleright Reveal the block earlier than protocol intended
 \triangleright Get MEV-extracting transactions
 \triangleright Replay extracted in the desired order
 \triangleright Add to state the MEV-extracting transactions
 \triangleright Publish the block containing the MEV-extracting transaction

Importantly, for such an attack to work, indeed, the $extract()$ function needs to have access to *all* the information about the transaction (e.g., having access to signed transactions that cannot change). What happens if some information could be withheld from the stakers? To understand this question, let us consider the following example. An ideal strategy to sandwich an AMM transaction $tx := \text{"Buy } x \text{ tokens of } X \text{ for } y \text{ tokens of } Y \text{ with a slippage of } s"$ is to produce a sequence $(tx_{BUY}^{attack}, tx, tx_{SELL}^{attack})$ so that the first attacking transaction tx_{BUY}^{attack} reduces the supply of token X for tx making it pay a higher price, and tx_{SELL}^{attack} extracts the sandwiching profit. However, if the attackers are unaware whether tx was a buy or sell transaction, or if it may be reversed with some probability (i.e., tx became selling token X for Y), then using the same attack can backfire and can result in losses for the attackers.

This idea leads to two natural questions. First, can we deviate from the framework to design a scheme that withholds some information from attacking stakers? Second, can we disincentivize attacks when the information is withheld?

In Section 5.1, we devise a strategy with *rationaly binding commitments* by creating an information asymmetry (e.g., only one party knows whether it is a BUY or a SELL transaction) between a specific staker \mathcal{F} (a flipper) and other stakers. In particular, the transaction can be modified after reveal has been invoked and \mathcal{F} is responsible to *complete* the transaction. The asymmetry of information allows a rational \mathcal{F} to improve its own utility at the expense of other stakers if the stakers choose to sandwich it. Consequently, this disincentivizes the other stakers to attack in the first place. We call this *rationaly binding* since the correctness of the transaction relies on \mathcal{F} being rational, which is a reasonable assumption. In this world, the client needs only to monitor the chain and hold \mathcal{F} accountable in case it observes \mathcal{F} does not complete the transaction correctly.

We can also rely on users or TEEs held by stakers to withhold some information; this information is only revealed during the reveal phase. We discuss how to disincentivize attacks when this is

possible in [31, Appendix E]. However, such a solution either requires the user to be online (which breaks the general ephemerality requirement) or needs additional assumptions, such as TEEs in the protocol. Since users do not have a stake, they typically tend to be ephemeral and this may cause liveness issues by not revealing their transactions.

5.1 AnimaguSwap

In this subsection, we describe a protocol design where some information is withheld from the attackers by a designated rational staker called flipper \mathcal{F} . This approach, as is, only works towards mitigating, and sometimes eliminating, sandwich attacks in constant product automated market makers (AMMs) like in Uniswap V2 [1]; though it can be easily extended to any constant function AMM. The key intuition is that if a set of stakers choose to sandwich a transaction, the protocol design allows the flipper to use its knowledge to gain a profit at the expense of those stakers. Thus, the binding property of the transaction relies on the flipper being rational. We first provide some background on an AMM and how sandwich attacks can be performed on transactions. Then, we present our protocol design and analyze it.

5.1.1 Background. An Automated Market Maker (AMM) such as Uniswap [1], Balancer [4], and Curve[12], uses automated algorithms to facilitate decentralized exchange of assets. AMMs set prices based on a mathematical formula based on the available liquidity of a given asset. In particular, in a Constant Product Market Maker, the product of the asset amounts in the liquidity pool is kept constant. Thus, if we have an AMM with two assets X and Y with quantities r_X and r_Y respectively, then $r_X * r_Y = k$ holds for some fixed value of k at all times.

When a user wants to trade one asset (X) for another (Y), they must deposit an amount of the first asset Δr_X and receive an appropriate amount of the second asset Δr_Y in return. Each transaction to the AMM is charged an additional fee, which we represent by f (e.g., $f = 0.3\%$ is a common value in practice). The

constraint becomes $(r_X + (1-f)\Delta r_X) * (r_Y - \Delta r_Y) = k$. Such a trade is SwapTokensForExactTokens in the Uniswap implementation and we represent it by Buy. Post the trade, the liquidity available would be $(r_X + \Delta r_X)$ and $(r_Y - \Delta r_Y)$ respectively (independent of the f).

A trade can be made with the fixed Δr_X amount in which to receive a fixed amount of first asset Δr_X , the user deposits an appropriate amount of second asset Δr_Y . Such a trade can be achieved by SwapExactTokensForTokens in Uniswap implementation and is referred to as Sell for the paper. The constraint for Sell is $(r_X - (1-f)\Delta r_X) * (r_Y + \Delta r_Y) = k$. Post the trade, the liquidity available would be $(r_X - \Delta r_X)$ and $(r_Y + \Delta r_Y)$ respectively.

Thus, given the current state of an AMM with r_X and r_Y tokens, a user can estimate Δr_Y received in exchange for depositing Δr_X or estimate Δr_X to deposit in exchange for receiving Δr_Y . However, if the state of the system changes due to some other transactions getting executed and affecting the liquidity pool, receiving Δr_Y for depositing Δr_X is not guaranteed. Thus, the system allows the user to specify a parameter expressed as a fraction called slippage s so that the number of tokens received by the user is not exact, e.g., $\geq (1-s)\Delta r_Y$. In other words, the user's transaction is specified as "Deposit $(1-f)\Delta r_X$ of X in exchange for $\geq (1-s)\Delta r_Y$ of Y ".

5.1.2 Sandwich Attack on Constant Product AMM. While slippage upper bounds users' loss, an attacker can still profit from user loss up to what is permitted by slippage by mounting *sandwich attacks*. This can be done by executing a transaction, depositing X , and receiving Y before the user's transaction (frontrunning). Once the user's transaction is executed, observe that the liquidity of Y has reduced further while it is the other way around for X . Thus, the attack can then run a reversed transaction, where the attacker sells the Y earned from the frontrunning transaction, in exchange for X . Such a transaction is called backrunning, and in an AMM, the attacker obtains a higher amount of X compared to what it had deposited in the frontrunning transaction. We refer interested readers to [31, Appendix F] for a mathematical analysis of the optimal frontrunning and backrunning parameters.

5.1.3 AnimaguSwap specification. We now present a protocol that can either reduce attacker gains or under some parameterizations, result in attacker losses, when sandwiching is attempted. As we have seen, in the frontrunning part of a sandwich attack, the attacker reduces the liquidity of the token that the user is interested in (token Y in our example). However, if the direction of the trade can be withheld from the attacker, then the attacker essentially has to guess one of the two directions. In situations where the attacker guesses incorrectly, it instead increases the liquidity of Y due to which the user can enjoy a much better trade and obtain $\Delta r'_Y > \Delta r_Y$ tokens of Δr_Y .

Our protocol is shown in [31, Appendix I, Fig. 11]. It generally follows the structure of the framework in Figure 1 except for a couple of aspects that we will describe later. Recall that \mathcal{F} refers to the flipper, a designated staker who would withhold the information from other stakers.

Transaction generation. Suppose the user intends to perform a trade from X to Y . This intent can be fulfilled in two ways: a "buy" transaction tx_{BUY} that buys Y or, equivalently, a "sell" transaction tx_{SELL} that sells X . With properly adjusted parameters, these two

transactions have the same execution outcome. Specifically, we write $tx_{SELL} = \text{SELL}(X, Y, \Delta r_X, \Delta r_Y, s, md)$, where Δr_X represents the number of tokens of X to be sold, in order to get maximum possible Y units, which is expected to be Δr_Y . The transaction would only go through if the number of tokens received $> (1-s)\Delta r_Y$. md represents any other metadata to be used by the transaction. Similarly, $tx_{BUY} = \text{BUY}(Y, X, \Delta r_Y, \Delta r_X, s, md)$. In our notation, the first parameter is $(\Delta r_X$ in case of SELL, and Δr_Y in case of BUY) is "exact" whereas the second parameter is determined by the first parameter and s .

The user first generates a random bit b to determine which transaction to use to fulfill its intended trade (note that the user is indifferent). Without loss of generality, we require that the user chooses the "buy" transaction tx_{BUY} if and only if $b = 0$. In the transaction metadata for tx_b , a hash of $v||w$ is included, where v and w are randomly generated numbers. This would be later used to allow slashing.

The key trick in AnimaguSwap is that the same coin decides if the user will "flip" the chosen transaction again. By flipping, we mean changing the polarity of the trade from selling asset X to buying asset X and vice versa, thereby creating a flipped transaction that is the opposite of the user's intent. Specifically, we require that the user flips the chosen transaction if and only if $b = 1$. We denote the transaction after the optional flipping as tx_b .

Following the same example where the user intends to trade from X to Y . If $b = 0$, the user will choose tx_{BUY} and does not flip, i.e., $tx_b = tx_{BUY}$. If $b = 1$, the user will choose tx_{SELL} and flips, i.e., $tx_b = \text{BUY}(X, Y, \Delta r_X, \Delta r_Y, s, md)$. Note that in this case $tx_b \neq tx_{BUY}$. Also, the committee always receives a "buy" transaction from the user, but the true intent is hidden in the flip bit. As we will detail in the next step, the user will submit tx_b to the committee and the flip bit b to a different staker called the flipper \mathcal{F} .

The second key trick is to disincentivize the flipper from revealing the flip bit (b), by having the user create another transaction $tx_{\mathcal{F}}$ which pays the flipper \mathcal{F} some amount of tokens if stakers attempt to sandwich tx_b but the direction of the sandwiched transaction is opposite. In particular, following the same example, if $b = 1$ and the committee creates a sandwich assuming $b = 0$, the user will earn $\Delta r'_Y > \Delta r_Y$. It can then pay the flipper $\Delta r'_Y - \Delta r_Y$ without decreasing its utility from a no-attack scenario. Similarly, if $b = 0$ and the committee assumes $b = 1$, then the user would swap $\Delta r'_Y < \Delta r_Y$ and pay the flipper $\Delta r_Y - \Delta r'_Y$. To represent it mathematically, the user pays the flipper $b(\Delta r'_Y - \Delta r_Y) + (1-b)(\Delta r_Y - \Delta r'_Y)$. Observe that obtaining Δr_Y is what the user expected; paying the remaining amount incentivizes \mathcal{F} . In scenarios where the polarity is guessed correctly, the flipper does not gain or lose money.

Transaction submission. During the transaction submission process, the user sends the bit b to the flipper. Importantly, the user does not sign this message, ensuring that the flipper cannot prove the polarity of the transaction to the other stakers. The bit b would later be revealed by the flipper to the blockchain by sending a signed message. What if the flipper cheats and presents an incorrect value? To ensure this does not happen, the flipper sends a signed message only to the user stating that it would reveal bit b corresponding to this transaction; if the flipper does otherwise, or does not reveal any value, then it can be slashed by the user based on this message.

However, one might argue that the flipper can forward a similar message to the stakers, and if this bit is incorrect, the stakers would be able to slash the flipper.

In order to safeguard against that, the user sends a random value v as an unsigned message to the flipper. It is crucial to ensure the deniability of the message sent by the user while maintaining the integrity of the message i.e., the message sent to the flipper could have been generated by the flipper itself. To ensure this, the user sends $m = (pk_u || b || v || txid)$ to the flipper encrypted under $pk_{flipper}$ using a hybrid public key encryption scheme (e.g., [5]). The message sent above could only be generated by a party who knows the correct random number v and the transaction ID $txid$. This ensures that no party except for the user and flipper (the two parties that know the content of the message) could have generated the message.

When returning the signed message to the user committing to b , it also includes v in the commitment. The user then generates another random number w , and uses $hash(v || w)$ in the transaction metadata. This ensures that only the user or a party with w can slash the flipper using the signed message the flipper sent, and thus, the flipper is free to sign any message it wants without risk of getting slashed.

Once both these steps succeed, the user secret-shares the (potentially flipped) transaction with the remaining stakers.

Transaction inclusion and reveal. The transaction inclusion process is straightforward. An accumulator value corresponding to the transaction is added to the chain whenever ShouldRelease predicate is true. Finally, the transaction content is revealed from the secret-shares when ShouldReveal is true.⁴ In this step, \mathcal{F} reveals the bit b too so that the correct transaction is revealed.

Pessimistic slashing. In case the flipper reveals bit \tilde{b} instead of b , then the user uses the signed message $(b, txid, v)_{\sigma(pk_{\mathcal{F}})}$ in addition to v and w to slash the flipper. The slashing rule gives us the following guarantees:

- **Correctness:** The user can only slash the flipper in case an incorrect bit is revealed, as slashing requires the user to show a signed flip bit different from what the flipper revealed. Correctness follows from the unforgeability of digital signatures.
- **Soundness:** If the flipper releases an incorrect bit, then the user can slash the flipper. Since the user has the signed message which contains the correct bit b , it acts as a commitment by the flipper and since both v and w are known to the user, the signature can be used to slash the flipper by showing the authenticity of v (revealing v and w , and verifying it against $hash(v || w)$).
- **Non-transferability:** The flipper cannot convince any party (other than the user) that $(b, txid, v)_{\sigma(pk_{\mathcal{F}})}$ can be used to slash the flipper. Note that this message can slash the flipper only if v is committed to by the metadata $md = hash(v || w)$. Since w is private to the user, md is a perfectly hiding commitment to v , so no party can verify that flipper's claimed v is committed to by md , following the definition of hiding.

⁴These predicates are abstract since their choice does not affect the design. In practice, one can replace these with predicates used by Shutter DKG [10], Ferveo [6], or Fino [24].

Observations. Here are a few observations related to this protocol. First, all known blockchains typically rely on accepting transactions that are signed only by the end users. This is the first protocol, to our knowledge, that includes a transaction where a portion of it (the bit b) is signed by a party (the flipper) other than the user. Second, a consequence of our approach is that, in the presence of a Byzantine flipper, the polarity of the executed transaction can be reversed. In practice, however, parties are sensitive to their utility, and thus, due to the existence of the slashing mechanism, a rational flipper would always reveal the correct bit. Thus, our protocol is only *rationaly binding* – this is the key aspect where we deviate from the requirement in the framework in Fig. 1. Third, since we expect the user to slash the flipper in case it deviates, the user cannot be ephemeral in the pessimistic case. The user needs to penalize the flipper within a reasonable timeframe (e.g., a few days). Finally, while the flipper can be any designated staker, a reasonable choice would be to have the staker that is expected to reveal the content of the transaction as the flipper. This ensures that the staker can reveal without waiting for inputs from other stakers.

5.2 AnimaguSwap Analysis

In this subsection, we will analyze AnimaguSwap detailed in Section 5.1.3. Our goal is to show that following protocol specifications is the dominant strategy for all the parties involved. For ease of analysis, we assume that the committee is colluding (e.g., through a collusion protocol such as Algorithm 1), and hence, treat the committee as a single party. We start the game after the user sends the flip bit to the flipper \mathcal{F} , and the transaction is secret shared with the committee C . For the analysis, C reconstructs the secret transaction sent to it. Also, the analysis would follow a single-shot game (i.e., the flipper and the committee are not repeated). In more practical scenarios, the game would be a multi-shot game. The reduction from a multi-shot game to a single-shot game is shown in [31, Appendix G].

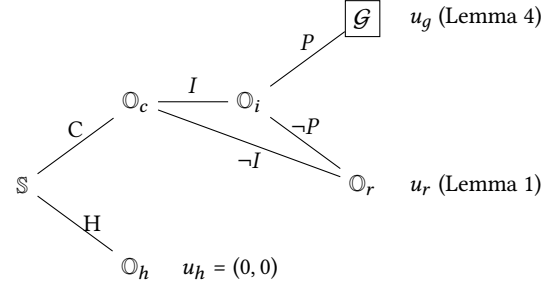
5.2.1 Game Setup. We analyze the game for a single AnimaguSwap transaction. The game underlying AnimaguSwap consists of two players - the flipper (\mathcal{F}), who receives one bit of information from the user on whether or not to flip the direction of the trade, and the committee (C), which receives the transaction.

Definition 3 (AnimaguSwap Game). We define the game as a tuple (N, A, O, μ, u) , where $N = \{\mathcal{F}, C\}$, $A = \{A_{\mathcal{F}}, A_C\}$ is the set of actions available to \mathcal{F} and C respectively, $O = \{\odot_h, \odot_r, \odot_s, \odot_{rs}\}$ ⁵ is the set of outcomes from the game, μ is the function that maps the set of actions to the outcome, and $u = \{u_h, u_r, u_s, u_{rs}\}$ is utility corresponding to outcomes.

The game states achievable through the game are $\mathbb{S} = \{\odot_c, \odot_i, \odot_h, \odot_r, \odot_s, \odot_{rs}\}$. The action space for the game is defined as $A_C = \{H$ (Honest), C (Collude), I (Invite \mathcal{F}), AI (Accept Information), AB (Anticipate Betray)) and $A_{\mathcal{F}} = \{P$ (Participate), Co (Cooperate), B (Betray)). The game tree Fig. 2 helps to understand action space better and also designs the function μ .

We assume that both players are rational, and have perfect information on strategy used by the other player, a concept used to find a Nash equilibrium of the game. We will first show such an

⁵ h stands for honest, r random, s sandwich, and rs reverse sandwich.



(a) Sequential Game

$\mathcal{F} \backslash C$	Accept information	Anticipate Betray
Co-operate	u_s (Lemma 2)	u_{rs} (Lemma 3)
Betray	u_{rs} (Lemma 3)	u_s (Lemma 2)

(b) \mathcal{G} - Simultaneous Game

Figure 2: Game Tree for actions taken in AnimaguSwap. Consists of a sequential game in which each of \mathcal{F} and C decide to participate or not and if both decide to participate, there exists a simultaneous game to decrypt the transaction correctly or incorrectly.

equilibrium of strategy in a simultaneous game (Fig. 2(b)), and then plug in the utility for the equilibrium strategy in the sequential game (Fig. 2(a)), to use iterative elimination of dominated strategies to find the best strategy across both games to earn the highest utility.

In this model, a player's utility is defined by the net number of tokens gained or lost in the game. Further, the action set for the game is limited, and any actions such as being involved in binding side contracts are out of the scope of this analysis. We will discuss the case with binding side contracts in [31, Appendix J]. There are two ways to reach outcomes O_{rs} and O_s , but the players' utility is only a function of the state, not the actions leading to the state (by definition of utility function). Specifically, in case the game reaches O_s , i.e., a successful sandwich attack due to either $\{Co, AI\}$ actions or $\{B, AB\}$, \mathcal{F} and C get a utility independent of actions taken to reach the state; In case the game reaches O_{rs} , i.e., unsuccessful sandwich attack due to either $\{Co, AB\}$ actions or $\{B, AI\}$, \mathcal{F} does not lose any utility due to the actions taken. As an example of utility sharing, after the action I a conditional bribe can be set to \mathcal{F} , which would only go through if the sandwich is successful.

From Section 5.1.1, without any attacker transaction, if the user's transaction was SELL then it would have followed

$$(r_X + (1 - f)\Delta r_X)(r_Y - \Delta r_Y^S) = r_X r_Y \quad (1)$$

Without any attacker transaction, if the user's transaction was BUY then it would have followed

$$(r_X - (1 - f)\Delta r_X)(r_Y + \Delta r_Y^B) = r_X r_Y \quad (2)$$

5.2.2 Analysis. Let us represent direction as a random variable chosen uniformly by the user from $\{BUY, SELL\}$. Without loss in generality, we can always represent the frontrunning transaction as a SELL transaction. In the case of BUY, the fee would be charged

from the other token, but the essence of the proof would remain the same. First, the frontrunning transaction would follow the constant product invariant.

$$(r_X + (1 - f)\Delta a_X)(r_Y - \Delta a_Y) = r_X r_Y \quad (3)$$

Now, after the frontrunning transaction, the victim transaction would follow. The transaction here could be a SELL transaction or BUY transaction, following the same set of parameters as C 's frontrunning transaction. If the user's transaction is SELL, then it would follow the constant product invariant.

$$(r_X + \Delta a_X + (1 - f)\Delta r_X)(r_Y - \Delta a_Y - \Delta r_Y^+) = (r_X + \Delta a_X)(r_Y - \Delta a_Y) \quad (4)$$

Also, the user's transaction would only be executed if

$$\Delta r_Y^+ > (1 - s)\Delta r_Y^S \quad (5)$$

The backrunning transaction would follow the constant product invariant with updated liquidity pools.

$$(r_X + \Delta a_X + \Delta r_X - \Delta a_X^+)(r_Y - \Delta a_Y - \Delta r_Y^+ + (1 - f)\Delta a_Y) = (r_X + \Delta a_X + \Delta r_X)(r_Y - \Delta a_Y - \Delta r_Y^+) \quad (6)$$

The profit would be given by

$$p^+ = \Delta a_X^+ - \Delta a_X \quad (7)$$

Since the sandwich is successful, $\Delta r_Y^+ < \Delta r_Y^S$. If the transaction is BUY, then it would follow Eqs. (8) to (11),

$$(r_X + \Delta a_X - (1 - f)\Delta r_X)(r_Y - \Delta a_Y + \Delta r_Y^-) = (r_X + \Delta a_X)(r_Y - \Delta a_Y) \quad (8)$$

$$\Delta r_Y^- > (1 - s)\Delta r_Y^B \quad (9)$$

$$(r_X + \Delta a_X - \Delta r_X - \Delta a_X^-)(r_Y - \Delta a_Y + \Delta r_Y^- + (1 - f)\Delta a_Y) = (r_X + \Delta a_X - \Delta r_X)(r_Y - \Delta a_Y + \Delta r_Y^-) \quad (10)$$

$$p^- = \Delta a_X^- - \Delta a_X \quad (11)$$

Since the sandwich is unsuccessful, $\Delta r_Y^- < \Delta r_Y^B$

LEMMA 1. *If the transaction's direction is uniformly distributed between $\{BUY, SELL\}$, and if C takes the action to cooperate (C), and either takes action to not invite \mathcal{F} ($-I$), or after taking action to invite \mathcal{F} (I), \mathcal{F} does not participate ($-P$), to reach output state O_r , then independent of the direction of trade C chooses, the utility of \mathcal{F} and C ($u_r = u_r(\mathcal{F}), u_r(C)$) is given by $(\frac{1}{2}(\Delta r_Y^B - \Delta r_Y^-), \frac{1}{2}(p^+ + p^-))$.*

PROOF. If C does not have information about the direction of the trade, then it can assume a direction among $\{BUY, SELL\}$. Without loss of generality, we represent the committee's frontrunning transaction in the form of a SELL($r_X, r_Y, \Delta a_X, \Delta a_Y, s, md$).

Since the direction of transaction is chosen at random from $\{BUY, SELL\}$, Eq. (7) and Eq. (11) govern the profit with probability 0.5 each, and the expected utility would be given by

$$u_r(C) = \frac{p^+ + p^-}{2} \quad (12)$$

Next, the utility for \mathcal{F} would be given from the AnimaguSwap protocol only in the case when C guesses the transaction direction incorrectly (and 0 in the other case).

$$u_r(\mathcal{F}) = \frac{1}{2}(\Delta r_Y^B - \Delta r_Y^-) \quad (13)$$

□

LEMMA 2. For game \mathcal{G} , if \mathcal{F} takes the action to co-operate (Co), and C accepts the information (AI), or \mathcal{F} takes the action to betray (B), and C anticipates betrayal (AB), then the utility is given by $u_s = (u_s(C) = \varepsilon, u_s(\mathcal{F}) = p^+ - \varepsilon)$, where $0 < \varepsilon < p^+$.

PROOF. The proof for the lemma follows Eq. (7). The profit gained from choosing the correct direction to sandwich would be shared between \mathcal{F} and C , regardless of the actions taken to reach the state. Thus, if C receives a utility of ε , then \mathcal{F} receives a utility of $p^+ - \varepsilon$, $0 < \varepsilon < p^+$. \mathcal{F} receives no utility directly from the AnimaguSwap protocol. □

LEMMA 3. For game \mathcal{G} , if \mathcal{F} takes the action to co-operate (Co), but C anticipates betrayal (AB), or \mathcal{F} takes the action to betray (B), but C accepts the information (AI), then the utility is given by $u_{rs} = (u_{rs}(C) = p^-, u_{rs}(\mathcal{F}) = \Delta r_Y^B - \Delta r_Y^-)$.

PROOF. Both sets of actions lead to a state where the committee inserts a frontrunning transaction with the incorrect direction. As stated in the setup, this would mean that \mathcal{F} has no utility from the game itself, and C loses utility governed by p^- (Eq. (11)). However, in accordance with the AnimaguSwap protocol, \mathcal{F} receives incentives from the protocol. This would be given by $\Delta r_Y^B - \Delta r_Y^-$. Thus $u_{rs} = (p^-, \Delta r_Y^B - \Delta r_Y^-)$. □

LEMMA 4. For game \mathcal{G} , if $\Delta r_Y^B - \Delta r_Y^- > p^+ - \varepsilon$ the Nash Equilibrium is governed by a mixed strategy for both \mathcal{F} and C , with \mathcal{F} betraying the committee with a probability of 0.5, and C anticipating betrayal with a probability of 0.5. The overall utility from game \mathcal{G} is given by $\left(\frac{\varepsilon + p^-}{2}, \frac{p^+ - \varepsilon + \Delta r_Y^B - \Delta r_Y^-}{2}\right)$

PROOF. To prove that the above strategy is a Nash Equilibrium, we reveal the strategy of each player to the other player and see if the strategy changes. From \mathcal{F} 's perspective, if it knows that C anticipates betrayal with a probability of 0.5, then the expected utility from betraying is $\frac{(u_s(\mathcal{F}) + u_{rs}(\mathcal{F}))}{2}$, whereas the utility from cooperating is $\frac{(u_s(\mathcal{F}) + u_{rs}(\mathcal{F}))}{2}$. From lemmas 2 and 3, both of these are equal, and thus \mathcal{F} does not have any additional utility from deviating from the strategy.

From C 's perspective, if it knows that \mathcal{F} betrays with a probability of 0.5, then the expected utility from anticipating betraying is $\frac{(u_s(C) + u_{rs}(C))}{2}$, whereas the utility from accepting the information is $\frac{(u_s(C) + u_{rs}(C))}{2}$. From lemmas 2 and 3, both of these are equal, and thus C does not have any additional utility from deviating from the strategy.

Thus, the given strategy is a Nash Equilibrium and by substituting the utilities from lemmas 2 and 3, the utility is given by $\left(\frac{\varepsilon + p^-}{2}, \frac{p^+ - \varepsilon + \Delta r_Y^B - \Delta r_Y^-}{2}\right)$. □

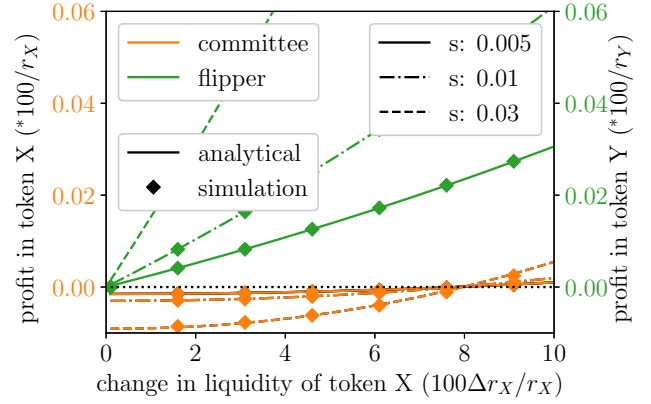


Figure 3: Gains of the committee and the flipper in a sandwich attack when using AnimaguSwap.

LEMMA 5. In the game state \odot_c , if the transaction's direction is uniformly distributed between $\{BUY, SELL\}$, and if $\Delta r_Y^B - \Delta r_Y^- > p^+ - \varepsilon$, it is strictly dominant for C to not invite \mathcal{F} ($-I$) over inviting \mathcal{F} (I).

PROOF. From Lemma 1, the utility of C from random choice is $u_r(C) = \frac{p^+ + p^-}{2}$. From Lemma 4, the utility of C from game \mathcal{G} is $\frac{\varepsilon + p^-}{2}$. Since $\varepsilon < p^+$, $\frac{p^+ + p^-}{2} > \frac{\varepsilon + p^-}{2}$. Thus, choosing the direction at random strictly dominates colluding with \mathcal{F} . □

Theorem 2. If the transaction's direction is uniformly distributed between $\{BUY, SELL\}$, $\Delta r_Y^B - \Delta r_Y^- > p^+$ and $\Delta a_X^+ + \Delta a_X^- - 2\Delta a_X < 0$, it is dominant for all parties in AnimaguSwap to follow the specification in Section 5.1.3.

PROOF. If $\Delta r_Y^B - \Delta r_Y^- > p^+$, then $\Delta r_Y^B - \Delta r_Y^- > p^+ - \varepsilon$ since $\varepsilon > 0$. Thus, from Lemma 5, it is strictly dominant for C to choose the direction of trade arbitrarily over colluding with \mathcal{F} . If $\Delta a_X^+ + \Delta a_X^- - 2\Delta a_X < 0$, then the expected profit from attacking AnimaguSwap by arbitrarily choosing a direction of trade is < 0 by plugging in p^+ and p^- in Lemma 1. Since taking honest actions leads to 0 utility, the committee would choose to take honest actions in AnimaguSwap. □

The chart (Fig. 3) represents the variation of the maximum utility with the user's transaction input (Δr_X) relative to the liquidity of X available (r_X), where slippage is set to be 0.005, 0.01, and 0.03 respectively. If both \mathcal{F} and C act honestly, they both receive no utility from the game. From the chart, we observe that $u_r(C) < 0$ until $\sim 7.9\%$ of the liquidity is traded. From this, we can conclude that C and consequently \mathcal{F} would act honestly unless the traded amount $> 7.9\%$ of the liquidity. 7.9% value is roughly when $\Delta a_X^+ + \Delta a_X^- - 2\Delta a_X = 0$.

To validate the analysis when C chooses the direction of trade arbitrarily, we simulate the attack in an AMM modeled after Uniswap v2, and calculate the expected gains for the attacker and the flipper when the attacker decides to sandwich the transaction across multiple values of s set by the user. Figure 3 shows that the analytical and simulation results are consistent.

Numerical example. As a concrete example, Table 2 in [31, Appendix H] shows the loss/gain of an attack against a simulated user trade sent through AnimaguSwap to Uniswap V2.

5.3 Base AnimaguSwap Evaluation

To evaluate the practicality of AnimaguSwap, we implemented the protocol (without the sub-routines for non-uniform distribution for {BUY, SELL} and repeated games) as a smart contract AS in 350 lines of Solidity. At a high level, AS is a middleware between users and an AMM, where users send commitments of transactions to AS following the AnimaguSwap protocol, stakers and flipper reveal user transactions, and finally, AS forwards the revealed transactions to the destination. In our prototype implementation, the destination is a fork of Uniswap V2 AMM. Specifically, the smart contract handles the following tasks:

- **Staking and slashing.** Stakers and the flipper deposit an appropriate amount of collateral to AS to join the system. When misbehavior is detected, users can submit evidence to AS to slash their stakes (c.f. Pessimistic Slashing [31, Appendix I, Figure 12]).
- **Commit.** To commit to transaction tx , a user runs *Generate Transaction* and *Transaction Submission* specified in Section 5.1.3. After interacting with stakers and flippers, the user calls `AS.commit` with the hash of tx_b .
- **Reveal and execution.** The reveal process follows *Transaction Revealing* in Section 5.1.3. Specifically, the stakers reconstruct the secret shared transaction off-chain and one of the stakers calls `AS.revealStaker` with the reconstructed transaction tx_b . `revealStaker` verifies the correctness of tx_b against the commitment. Then, Flipper calls `AS.revealFlipper` with flip bit b . With both b and tx_b , AS recovers the user's original intent and executes it. In our implementation, this triggers a call to Uniswap V2.

Off-chain parties (stakers and the flipper) are implemented in 400 lines of TypeScript. The code can be found at <https://anonymous.4open.science/r/AnimaguSwap-D31F/>.

Evaluation. The stakers' main task is to reconstruct user transactions from secret shares to open commitments and execution transactions. The flipper's task does not involve any costly computation. Compared to smart contract execution, off-chain computation is much more efficient (see, e.g., [6], for evaluation). The main performance metric, therefore, is the gas consumption of AS.

Table 1 shows the gas cost breakdown of executing a Uniswap trade from 1 wBTC to DAI, which costs 351k gas (\$4.6 at the time of writing). To compare, a typical Uniswap V2 trade costs about 150k gas (\$2).⁶ The strong protection of AnimaguSwap thus incurs a 1.3x overhead. Note that since transactions are reconstructed off-chain, the gas cost does not increase with the number of stakers. Also, note that the cost does not increase with the value of the transaction. Therefore, while gas usage can be potentially reduced further, our preliminary implementation is already quite practical for high-value transactions.

5.4 Non-uniform Distribution in {BUY, SELL}

In Theorem 2, we assumed for the result that there is an even spread between two trade directions, {BUY, SELL}. However, this

Function	Gas cost	Function	Gas cost
commit	66275	revealFlipper	46069
revealStaker	239342	complain	34862

Table 1: Cost for a AnimaguSwap call. It takes 1564706 gas to deploy (one-time cost). Complain is not on a critical path.

assumption might not hold in real-world scenarios. For instance, during the LUNA crash of May 2022, a trade involving ETH and LUNA is more likely to be selling LUNA, than buying it. In such scenarios, AnimaguSwap, as presented above, does not work, as the attacker can *guess* the flipper's bit based on public information such as market sentiment. Another problematic scenario is when the user only owns one of the assets in her transaction, so the attacker can deduce that the transaction must be selling that asset.

Formally, for a given pair of assets $Pair = (X, Y)$, if the probability of trading X to Y in a random transaction is different from that of the reverse direction, we say this asset pair is *biased*, and we denote the probability of the more probable direction with $P_{DIRPair} \geq 0.5$. (The probability of the less probable direction is $1 - P_{DIRPair}$.) We call $P_{DIRPair}$ the *bias* of the pair for short. We omit the subscription when the asset pair is clear from context. Now we present an enhancement to AnimaguSwap that can protect biased asset pairs.

Our idea is to obfuscate the user transaction (using techniques to be presented shortly) so that an attacker mounting sandwich attacks based on guessed transaction information will equally likely fail or succeed. 'Success' here means the attacker manages to profit from the sandwich attack, while 'failure' means the attacker loses. As long as the amounts gained or lost in a 'success' and 'failure' remain the same as in the gains and loss in Section 5.2, all the lemma statements and the theorem statement would follow. We present the subroutines integrated into AnimaguSwap in [31, Appendix I, Figure 12].

In order to get an equal probability of 'Success' and 'Failure', we take the following three steps: 1) remove any user dependency by creating a pool of users; 2) hide the asset being traded among multiple transactions in a way that no party can distinguish which asset is being traded, and 3) set the parameters so that the attacker can only lose utility from sandwiching incorrectly, but never gain any profit, if it chooses the wrong asset. With these three properties, we make it such that the attacker gains a profit when it gets everything correct, however, in multiple cases where it gets the trade direction incorrectly, the attacker loses. There also exist cases where the attacker neither wins nor loses (pays some transaction fee), but we assume the utility for such cases is 0.

To start with, we obfuscate the user identity by mixing it among a pool of decoy users so that the attacker cannot gain information about the user's transaction from the user's on-chain presence (e.g., if this user doesn't possess the asset X , then the attacker can infer that the user cannot only be buying X).

To select decoy users, the user randomly chooses users who collectively own all the assets being traded (the real asset pair and all the auxiliary asset pairs, which we will introduce next). The selection is made such that given all the assets (including auxiliary assets, which will be introduced soon) involved in the

⁶<https://etherscan.io/gastracker>

trade, selecting any user out of the pool gives no extra information about the trade to the attacker over any other user in the pool. For example, let's say the user submits a trade between USDC and ETH. It has USDC but not ETH. The user would create a pool of users (from the set of all blockchain users), such that one user has ETH but no other assets. This pool of users (in this case, 2 users are sufficient, but adding more users to have redundancy does not hurt) would be used to generate a ring signature [26], such that no party can distinguish which user amongst the pool of users created the transaction. The flipper escrows the true identity of the transaction creator and will be released in the same way as the flip bit.

However, in doing so, the flipper can release incorrect information about the user sending the transaction. i.e. it can create a transaction itself and generate a transaction, but since the identity is obscured, it releases the wrong identity. To prevent this, we need a way for a user to prove that the transaction is not its own. We do this by controlling the v variable used in metadata creation. It was introduced to hide the commitment of the commitment received from Flipper to the information it has. Instead of randomly generating it, the user uses $v = h(sk, txid)$, where sk is the secret key, and $txid$ is the block number. To any party without sk , it remains a completely random number, but any user can generate a zero-knowledge proof that it is not the user generating the transaction.

We also introduce the concept of auxiliary transactions. The idea behind the concept is to hide how much is being traded and which pair of assets is involved. Consider, for instance, the bias of the user's "real" transaction is $P_{DIR_{real}} > 0.5$. To achieve an equal probability of 'Success' and 'Failure', the user chooses N_{ASSET} different asset pairs with a skew of the direction of trade similar to $P_{DIR_i} > 0.5$ for each asset pair i , the value of which we will define shortly. Next, it creates N_{ASSET} auxiliary transactions with the newly chosen pairs, the same value as the original transaction, and chooses the direction of the trade with the same probability P_{DIR_i} . We refer to $P_{ASSET_{real}}$ as the probability that amongst a set of assets, real is the actual asset pair, and P_{ASSET_i} for introduced auxiliary asset pair i , such that $P_{ASSET_{real}} + \sum_i P_{ASSET_i} = 1$. For example, let's say for the ETH-USDC pair, there exists a probability of 0.6 for traders to buy ETH from USDC. The user would find another pair (e.g., WBTC-Tether) such that the attacker cannot differentiate whether this transaction was an ETH-USDC pair or the WBTC-Tether pair. We will show how this other pair (WBTC-Tether) is found after the claim statement. The slippage for this transaction is such that at the current price, the trade would fail, however, at any price better than the current price, the transaction would succeed. This way, the attacker can only 'Fail' if it chooses to attack the auxiliary assets, and not 'Succeed' even when the direction guessed is correct.

Introducing new auxiliary transactions would require the token being traded to be available to the user. However, the user may not own auxiliary assets to make transactions in the first place. One workaround is to borrow auxiliary assets from a decentralized lending platform in an indistinguishable way (indistinguishable which asset is being borrowed) from the attacker. A naïve solution to this is to loan all possible assets involved in the trade, with the smallest union of assets not involved in the trade over the user pool. e.g. If assets A and B are the real traded assets, and C, D, E, and F are used as auxiliary assets, and in the user pool, all users

have either asset G, H, or I, then a loan for A, B, C, D, E, F would be required keeping an asset G, H and I as collateral (i.e., a total of 18 loans). Most of these loans would fail since the user would not have the required collateral to obtain the loan.

However, with AnimaguSwap, we have an advantage that at the time of execution, we have complete knowledge, and hiding the access is not important. Thus, if we can program AnimaguSwap contract to involve generating transactions for loans, we can optimize the process to not involve unnecessary loan transactions.

In addition, the collateralized asset needs to be outside the set of assets involved in the trade (the main transaction and the auxiliary transaction). For example, in the ETH-USDC trade, if the user chooses the auxiliary pair as WBTC-Tether (Let's say WBTC is being bought more than Tether), then it would need to ensure that the auxiliary transaction is valid. If it does not have the asset being traded away (Tether), then the AnimaguSwap would issue a loan for Tether, keeping SUSHI (A third asset not involved in the trade) as collateral. This is important for the user ambiguity constraint since if the user does not have an asset for collateral, it gives the attacker information that the user may be more likely to use the asset it has as the traded asset. If the user is using a loan from an external asset, all the assets involved are equally likely, however, in case an asset is used that is in the set of traded assets, then that asset is more likely to be a traded asset (real asset, or auxiliary traded asset).

With the above-described changes to the AnimaguSwap specification ([31, Appendix I]), we can claim the following:

Claim 1. *Given a probability of the trade being $P_{DIR_{real}} > 0.5$ in one direction, if there exists a set of asset pairs with probability $P_{DIR_i} > 0.5$ of it being traded in a given direction, such that $\sum_i (P_{ASSET_i} * P_{DIR_i}) = 1 - 2P_{ASSET_{real}} * P_{DIR_{real}}$, where P_{ASSET_i} represents the probability of asset pair i to be the real asset pair amongst the set of asset pairs chosen for a pool of users, given that the user owns an asset not involved in the trade as collateral, then the subroutine described ensures that the probability of an attacker successfully gaining utility from sandwich is equal to probability of an attacker losing utility.*

PROOF. The protocol creates the following scenarios when sandwiched attacked - 1) the attacker chooses both the asset and direction of trade correctly (probability = $P_{ASSET_{real}} * P_{DIR_{real}}$); 2) chooses the asset correctly, but the direction is inverted (probability = $P_{ASSET_{real}} * (1 - P_{DIR_{real}})$); 3) chooses the asset incorrectly but the direction is the same as the direction chosen for the dummy transaction (probability = $\sum_i (P_{ASSET_i} * P_{DIR_i})$); and 4) chooses the asset as well as the direction of the dummy transaction incorrectly (probability = $\sum_i (P_{ASSET_i} * (1 - P_{DIR_i}))$). Now, the attacker loses capital in the second and fourth scenarios, i.e. with probability $P_{ASSET_{real}} * (1 - P_{DIR_{real}}) + \sum_i (P_{ASSET_i} * (1 - P_{DIR_i}))$ and gains capital in only the first scenario with probability $P_{ASSET_{real}} * P_{DIR_{real}}$. The auxiliary assets are chosen such that $\sum_i (P_{ASSET_i} * P_{DIR_i}) = 1 - 2P_{ASSET_{real}} * P_{DIR_{real}}$. Earlier $P_{DIR_i} > 0.5$ was chosen, and thus the restriction needs to be set when choosing the asset pairs (because the attacker would always sandwich the more probable direction after choosing the asset). \square

To complete the running example, if user wants to trade USDC to buy ETH, where there is a 0.6 probability ($P_{DIR_{real}}$) to buy ETH, and amongst all asset pairs traded in AMMs, USDC-ETH pair occurs 20% of the times, then it would choose another asset pair such that $0.2 * 0.6 = 0.2 * 0.4 + P_{ASSET_i} * P_{DIR_i}$. Thus, an asset that could be used by the user is the WBTC-Tether if it occurs 8.8% of time, and with probability 0.55 the direction is biased towards WBTC. (It could have been 10%, 0.6 probability, and so on as long as it satisfies the formula in the claim). After choosing the asset, the user creates the auxiliary transaction as described. According to our claim, if an attacker sandwiches the new transaction, then it would have the same probability for ‘success’ case (guessing USDC-ETH, and that the transaction is in the direction of ETH) and ‘failure’ case (guessing the wrong direction of trade for the guessed pair).

Using the claim with Theorem 2, we get the following theorem:

Theorem 3. *Given there exists a set of asset pairs with probability $P_{DIR_i} > 0.5$ of it being traded in a given direction, such that $\sum_i (P_{ASSET_i} * P_{DIR_i}) = 1 - 2P_{ASSET_{real}} * P_{DIR_{real}}$, where P_{ASSET_i} represents the probability of asset i to be the real asset amongst the set of assets chosen for a pool of users, given that the user owns an asset not involved in the trade as collateral, $\Delta r_Y^B - \Delta r_Y^- > p^+$ and $\Delta a_X^+ + \Delta a_X^- - 2\Delta a_X < 0$, it is dominant for all parties in AnimaguSwap to follow the specification in Section 5.1.3.*

6 Discussion and Future Work

On using primitives such as witness encryption, time lock encryption, or traceable secret sharing to circumvent Theorem 1. Our setup of Framework 1 assumes that the output of the order function is directly used as an input to the reveal function. This implies that a transaction can be revealed at any time after it is ordered so far as sufficiently many stakers participate. On the other hand, the use of cryptographic primitives such as Witness Encryption [17] and Time Lock Encryption [27] tie the reveal of transactions to satisfying some condition (e.g., the passage of time); thus, these primitives can be used to circumvent the impossibility result. The use of TEEs in [31, Appendix E] can be considered as an implementation of witness encryption assuming trusted hardware.

The notion of traceable secret sharing introduced by Goyal et al. [18] allows users to produce secret shares such that once the data is reconstructed, parties releasing their secret shares can be identified. However, our attack strategy in Algorithm 1 circumvents this by producing only the generated transactions as output.

On sending deniable messages. Recent studies [30] demonstrate that deniability may be compromised when keys are encumbered in a Trusted Execution Environment (TEE) such as Intel SGX or if a committee manages the flipper’s keys through a distributed key system. Consequently, users must verify that they are interacting with a single, unrestricted user as the flipper. This verification can be achieved by employing a Complete Knowledge Proof [20], which substantiates that a single user possesses unrestricted access to the information provided, thereby reinstating deniability. To use CK in practice, all flippers would require a CK certificate either obtained through a TEE (since the input to a TEE is public to the party that inputs it) or an ASIC-based proof (which can be generated in a

reasonable time only if the key is known to a party generating the proof) verified on-chain.

On lack of knowledge of real-world entities. Our impossibility results crucially rely on the inability of the protocol participants to distinguish whether two public keys belong to the same real-world entity or not. This is reasonable, especially in a permissionless setting. However, in practice, if we can perform an analysis of the flow of transactions across different keys and their uses, and derive intelligence based on these transactions (e.g., [11]), we can identify the existence of such attacks with the analysis acting as a ‘proof’.

On collusion between the user and the flipper. Even if the user and flipper collude, in AnimaguSwap, it is not possible for a user’s transaction to violate the soundness condition. During a collusion there are two cases that may arise: the user does not receive a commitment itself, or the flipper violates the commitment it shared. The protocol specification requires the user to collect a commitment before posting a transaction. If the commitment is not given to the user, and the user does not post a commitment to what it receives from the flipper, the flipper can refuse to share profits with the user and instead collude with the committee for a potentially larger share of the profit. Thus, if the user and flipper collude, and the user has the flipper’s commitment, the user can not only get better execution and share profits with the flipper but also slash the flipper.

On user acting as flipper. The user cannot be asked to withhold the information in AnimaguSwap. This is because the user cannot be trusted to release the information. As with user withholding information, if the user does not release a flip bit, it cannot be slashed and thus threatens the protocol’s liveness. Flipper can, however, be a user (which has a designated stake from being a Flipper) since the transactions being sent by the Flipper would only further discourage sandwich attacks since even when the committee is randomly predicting, the flipper can ensure that the prediction they choose is incorrect.

Acknowledgments

This work was supported in part by NSF award 2237814 and an Ethereum Research Grant.

References

- [1] Hayden Adams, Noah Zinsmeister, and Dan Robinson. 2020. *Uniswap v2 Core*. Technical Report. Tech. rep., Uniswap.
- [2] atom_crypto. 2022. The MEV Game of the Crypto Economy: Osmosis’ Threshold Encryption vs. SGX of Flashbot? <https://mirror.xyz/infnet.eth/SFjR1H1-RMnKoloPjqkxpauVPrLYGqLHQp1dY9FHvx4>.
- [3] Kushal Babel, Yan Ji, Ari Juels, and Mahimna Kelkar. 2023. PROF: Fair Transaction-Ordering in a Profit-Seeking World. <https://initc3org.medium.com/prof-fair-transaction-ordering-in-a-profit-seeking-world-b6dadd71f086>.
- [4] Balancer. [n. d.]. Balancer Docs. <https://docs.balancer.fi/reference/math/stable-math.html>.
- [5] Richard Barnes, Karthikeyan Bhargavan, Benjamin Lipp, and Christopher A. Wood. 2022. RFC 9180: Hybrid public key encryption. <https://datatracker.ietf.org/doc/rfc9180/>
- [6] Joseph Bebel and Dev Ojha. 2022. Ferveo: Threshold decryption for mempool privacy in BFT networks. *Cryptology ePrint Archive* (2022).
- [7] Iddo Bentov, Yan Ji, Fan Zhang, Lorenz Breidenbach, Philip Daian, and Ari Juels. 2019. Tesseract: Real-time cryptocurrency exchange using trusted hardware. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 1521–1538.

- [8] Vitalik Buterin. [n. d.]. State of research: Increasing censorship resistance of transactions under proposer/builder separation (PBS). https://notes.ethereum.org/@vbuterin/pbs_censorship_resistance.
- [9] Christian Cachin, Jovana Micić, Nathalie Steinhauer, and Luca Zanolini. 2022. Quick order fairness. In *Financial Cryptography and Data Security: 26th International Conference, FC 2022, Grenada, May 2–6, 2022, Revised Selected Papers*. Springer, 316–333.
- [10] Cducrest. 2022. Shutterized Beacon Chain. <https://ethresear.ch/t/shutterized-beacon-chain/12249>.
- [11] Chainalysis. [n. d.]. Chainalysis. <https://www.chainalysis.com/>.
- [12] Curve. [n. d.]. Understanding Curve v1 Curve Finance. <https://resources.curve.fi/base-features/understanding-curve>.
- [13] Philip Daian, Steven Goldfeder, Tyler Kell, Yunqi Li, Xueyuan Zhao, Iddo Bentov, Lorenz Breidenbach, and Ari Juels. 2020. Flash boys 2.0: Frontrunning in decentralized exchanges, miner extractable value, and consensus instability. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 910–927.
- [14] Philip Daian, Steven Goldfeder, Tyler Kell, Yunqi Li, Xueyuan Zhao, Iddo Bentov, Lorenz Breidenbach, and Ari Juels. 2020. Flash boys 2.0: Frontrunning in decentralized exchanges, miner extractable value, and consensus instability. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 910–927.
- [15] FlashBots. 2020. Flashbots Resource Document. <https://docs.flashbots.net/>.
- [16] flashbots. 2022. Mev-Boost GitHub. <https://github.com/flashbots/mev-boost>.
- [17] Sanjam Garg, Craig Gentry, Amit Sahai, and Brent Waters. 2013. Witness encryption and its applications. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*. 467–476.
- [18] Vipul Goyal, Yifan Song, and Akshayaram Srinivasan. 2021. Traceable secret sharing and applications. In *Advances in Cryptology—CRYPTO 2021: 41st Annual International Cryptology Conference, CRYPTO 2021, Virtual Event, August 16–20, 2021, Proceedings, Part III* 41. Springer, 718–747.
- [19] Lioba Heimbach and Roger Wattenhofer. 2022. Eliminating sandwich attacks with the help of game theory. In *Proceedings of the 2022 ACM on Asia Conference on Computer and Communications Security*. 153–167.
- [20] Mahimna Kelkar, Kushal Babel, Philip Daian, James Austgen, Vitalik Buterin, and Ari Juels. 2023. Complete Knowledge: Preventing Encumbrance of Cryptographic Secrets. *Cryptology ePrint Archive* (2023).
- [21] Mahimna Kelkar, Soubhik Deb, and Sreeram Kannan. 2021. Order-Fair Consensus in the Permissionless Setting. *IACR Cryptol. ePrint Arch.* 2021 (2021), 139.
- [22] Mahimna Kelkar, Soubhik Deb, Sishan Long, Ari Juels, and Sreeram Kannan. 2021. Themis: Fast, Strong Order-Fairness in Byzantine Consensus. *Cryptology ePrint Archive* (2021).
- [23] Mahimna Kelkar, Fan Zhang, Steven Goldfeder, and Ari Juels. 2020. Order-fairness for byzantine consensus. In *Annual International Cryptology Conference*. Springer, 451–480.
- [24] Dahlia Malkhi and Pawel Szalachowski. 2022. Maximal Extractable Value (MEV) Protection on a DAG. *arXiv preprint arXiv:2208.00940* (2022).
- [25] Kaihua Qin, Liyi Zhou, and Arthur Gervais. 2021. Quantifying Blockchain Extractable Value: How dark is the forest? *arXiv preprint arXiv:2101.05511* (2021).
- [26] Ronald L. Rivest, Adi Shamir, and Yael Tauman. 2001. How to Leak a Secret. In *Advances in Cryptology — ASIACRYPT 2001*, Colin Boyd (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 552–565.
- [27] Ronald L Rivest, Adi Shamir, and David A Wagner. 1996. Time-lock puzzles and timed-release crypto. (1996).
- [28] Sikka inc. 2022. Sikka Projects. <https://sikka.tech/projects/>.
- [29] Nik Unger and Ian Goldberg. 2015. Deniable key exchanges for secure messaging. In *Proceedings of the 22nd acm sigsac conference on computer and communications security*. 1211–1223.
- [30] Ricardo Vieitez Parra et al. 2018. The Impact of Attestation on Deniable Communications. (2018).
- [31] Sarisht Wadhwa, Luca Zanolini, Francesco D’Amato, Aditya Asgaonkar, Chengrui Fang, Fan Zhang, and Kartik Nayak. 2023. Data Independent Order Policy Enforcement: Limitations and Solutions. *Cryptology ePrint Archive*, Paper 2023/868. <https://doi.org/10.1145/3658644.3670367> <https://eprint.iacr.org/2023/868>.
- [32] Matheus Venturynne Xavier Ferreira and David C. Parkes. 2023. Credible Decentralized Exchange Design via Verifiable Sequencing Rules. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing (STOC 2023)*. Association for Computing Machinery, New York, NY, USA, 723–736. <https://doi.org/10.1145/3564246.3585233>
- [33] Sen Yang, Fan Zhang, Ken Huang, Xi Chen, Youwei Yang, and Feng Zhu. 2022. SoK: MEV countermeasures: Theory and practice. *arXiv preprint arXiv:2212.05111* (2022).
- [34] Liyi Zhou, Kaihua Qin, Christof Ferreira Torres, Duc V Le, and Arthur Gervais. 2021. High-frequency trading on decentralized on-chain exchanges. In *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 428–445.