

Replay Attacks and Defenses Against Cross-shard Consensus in Sharded Distributed Ledgers

Alberto Sonnino, Shehar Bano, Mustafa Al-Bassam, George Danezis
University College London & chainspace.io

Abstract—We present a family of replay attacks against sharded distributed ledgers targeting cross-shard consensus protocols, such as the recently proposed Chainspace and Omniledger. They allow an attacker, with network access only, to double-spend or lock resources with minimal efforts. The attacker can act independently without colluding with any nodes, and succeed even if all nodes are honest; most of the attacks can also exhibit themselves as faults under periods of asynchrony. These attacks are effective against both shard-led and client-led cross-shard consensus approaches. We present Byzcuit—a new cross-shard consensus protocol that is immune to those attacks. We implement a prototype of Byzcuit and evaluate it on a real cloud-based testbed, showing that our defenses impact performance minimally, and overall performance surpasses previous works.

Index Terms—Distributed Ledgers, Sharding, Attacks

1. Introduction

Sharding is one of the key approaches to address blockchain scalability issues [2], and a growing number of systems are implementing sharded blockchains [1], [2], [6], [8], [9], [12]. The key idea is to create groups (or shards) of nodes that handle only a subset of all transactions and system state, relying on classical Byzantine Fault Tolerance (BFT) protocols for reaching *intra-shard consensus*. These systems achieve optimal performance and scalability because: (i) non-conflicting transactions can be processed in parallel by multiple shards; and (ii) the system can scale up by adding new shards. This separation of transaction handling across shards is not perfectly ‘clean’—a transaction might rely on data managed by multiple shards, requiring an additional step of *cross-shard consensus* across the concerned shards. An atomic commit protocol (such as the two-phase commit protocol [7]) typically runs across all the concerned shards to ensure the transaction is accepted by all or none of those shards.

We present the first replay attacks on cross-shard consensus in sharded blockchains. An attacker can launch these attacks with minimal effort, without subverting any nodes, and assuming a weakly synchronous network (and in some cases, without relying on any network assumption)—even when the byzantine safety assumptions are satisfied. These attacks compromise key system properties of safety and liveness, effectively enabling the attacker to double-spend coins (or any other objects managed by the blockchain) and create coins out of thin air. Our attacks apply to the two main approaches to achieve cross-shard consensus [2]: (i) shard-led protocols that only involve the concerned shards, and require no

external entity for coordination (Section 4); and (ii) client-led protocols that are coordinated by the client (Section 5).

We concretely sketch the replay attacks in the context of two representative systems: Chainspace [1] as an example of shard-led protocols; and Omniledger [8] as an example of client-led protocols. Not only those systems were recently presented at top security conferences, but they form the basis of numerous start-ups and open-source projects such as chainspace.io¹ and Harmony². For each of the two cross-shard consensus approaches, Appendix A describes how an attacker can actively stage the attack by eliciting from the system the messages to replay (in contrast to passively observing the network traffic, and waiting to detect and record the target messages). We also discuss the feasibility of these attacks and their real-world impact; and we responsibly disclosed them to the concerned companies. We implement and open-source³ a demo of the replay attacks described in this paper (Section 3).

The replay attacks we present are generic and apply to other systems that are based on similar models, like RapidChain [12]. Ethereum’s cross-shard ‘yanking’ proposal [4] also faces similar challenges; Section 9 describes their cross-shard consensus protocol and compares their current proposal to mitigate cross-shard replay attacks with this work. We note that account-based blockchains like Ethereum defend against transactions replay using account sequence numbers, in an entirely different context; *i.e.* each account holds a monotonically increasing counter to prevent attackers from re-submitting old transactions to the network. On the other hand, this work focuses on attacks due to replaying messages in cross-shard atomic commit protocols. Based on our detailed analysis of replay attacks, we develop a defense strategy (Section 6).

Drawing insights from our analysis of performance trade-offs and replay attack vulnerabilities in existing shard-led and client-led cross-shard consensus protocols, we present a hybrid system, Byzcuit (Section 7). It combines useful features from both these design approaches to achieve high performance and scalability, and leverages our proposed defense to achieve resilience against replay attacks. Byzcuit employs a Transaction Manager to coordinate cross-shard communication, reducing its cost to $O(n)$ communication, between n shards, in the absence of faults. We implement a prototype of Byzcuit in Java as a fork of the Chainspace code [1], and release it as

1. <https://chainspace.io>

2. <https://harmony.one>

3. <https://github.com/sheharbano/byzcuit/tree/replay-attacks>

an open-source project⁴. We evaluate Byzcuit on a real cloud-based testbed under varying transaction loads and show that Byzcuit has a client-perceived latency of less than a second, even for a system load of 1000 transactions per second (tps). Byzcuit’s transaction throughput scales linearly with the number of shards by 250–300 tps for each shard added, handling up to 1550 tps with 10 shards—which is about 8 times higher than Chainspace with a similar setup. We quantify the overhead of our replay defenses and find that as expected those reduce the throughput by 20–250 tps.

Contributions. We make the following key contributions: we (i) develop the first replay attacks against shard-led and client-led cross-shard consensus protocols, and illustrate their impact on important academic and implemented designs; and (ii) present defenses; (iii) design a hybrid, new system Byzcuit with improved performance trade-offs, and which integrates our proposed defense to achieve resilience against the replay attacks; and (iv), we implement a prototype of Byzcuit and evaluate its performance and scalability on a real distributed set of nodes and under varying transaction loads, and illustrate how it is superior to previous proposals.

2. Background and Related Work

We present background and related work on cross-shard consensus protocols.

Sharded blockchains. Earlier systems like Bitcoin [11] probabilistically elect a single node which can extend the blockchain. However, such systems assume synchrony, have no finality (*i.e.*, forks can exist and be eventually accepted) and low performance (*i.e.*, high latency and low throughput). Consequently, there has been a shift to committee-based designs [2] where a group of nodes collectively extends the blockchain typically *via* classical byzantine fault tolerance (BFT) consensus protocols such as PBFT [5]. While these systems offer better performance, single-committee consensus is not scalable—as every node handles every transaction, adding more nodes to the committee decreases throughput due to the increased communication overhead. This motivated the design of *sharded* systems, where multiple committees handle a subset of all the transactions—allowing parallel execution of transactions. Every committee has its own blockchain and set of objects (or unspent transaction outputs, UTXO) that they manage. Committees run an *intra-shard consensus protocol* (*e.g.*, PBFT) within themselves, and extend the blockchain in parallel.

Cross-shard consensus. In sharded systems, some transactions may operate on objects handled by different shards, effectively requiring the relevant shards to additionally run a *cross-shard consensus protocol* to enable agreement across the shards. If any of the shards relevant to the transaction rejects it, all the other shards should likewise reject the transaction to ensure atomicity.

The typical choice for implementing cross-shard consensus is the two-phase atomic commit protocol [7]. This protocol has two phases which are run by a *coordinator*. In the first *voting* phase, the nodes tentatively write changes

locally, lock resources, and report their status to the coordinator. If the coordinator does not receive status message from a node (*e.g.*, because the node crashed or the status message was lost), it assumes that the node’s local write failed and sends a rollback message to all the nodes to ensure any local changes are reversed, and locks released. If the coordinator receives status messages from all the nodes, it initiates the second *commit* phase and sends a commit message to all the nodes so they can permanently write the changes and unlock resources. In the context of sharded blockchains, the atomic commit protocol operates on shards (which make the local changes associated with the voting phase *via* an intra-shard consensus protocol like PBFT), rather than individual nodes. A further consideration is who assumes the role of the coordinator.

Related Work. Replay attacks in general have seen extensive study in the security literature. This is the first paper that presents replay attacks on cross-shard consensus protocols. Traditionally, the most stringent threat models assumed by consensus protocols involve byzantine adversaries who are able to control or subvert consensus nodes and cause them to behave arbitrarily. Repurposing those protocols to open permissionless networks (*e.g.*, blockchains) opens up new attack avenues such as replay attacks as shown in this paper. There are currently two key approaches to cross-shard consensus [2]. The first approach involves *client-led protocols* (such as Omniledger [8] and RSCoin [6]), where the client acts as a coordinator. These protocols assume that clients are incentivized to proceed to the unlock phase. While such incentives may exist in a cryptocurrency application where an unresponsive client loses its own coins if the inputs are permanently locked, these do not however hold for a general-purpose platform where transaction inputs may have shared ownership. The second approach involves *shard-led protocols* (such as Chainspace [1], Rapidchain [12] and Elastico [9]), where shards collectively assume the role of a coordinator. All the concerned shards collaboratively run the protocol between them. This is achieved by making the entire shard act as a ‘resource manager’ for the transactions it handles. We describe our replay attacks in the context of two representative systems: Chainspace [1] as an example of shard-led protocols (Section 4); and Omniledger [8] as an example of client-led protocols (Section 5). We provide a more detailed description of these systems in the relevant sections.

3. Attack Overview

Sections 4 and 5 discuss replay attacks on both shard-led and client-led cross-shard consensus protocols, respectively. We present a high-level description of these attacks, the threat model, demo attack implementation, and the notation used in this paper.

Replay Attacks on Cross-Shard Consensus. The attacker records a target shard’s responses to the atomic commit protocol, and replays them during another instance of the protocol. We present (i) attacks against the first phase (*voting*), and (ii) attacks against the second phase (*commit*) of the atomic commit protocol.

To attack the first phase (*voting*) of the atomic commit protocol, the attacker replaces messages generated by

4. <https://github.com/sheharbano/byzcuit>

the target shard by replaying pre-recorded messages. In practice, the attacker does not *replace* those messages—it achieves a similar result by making its replayed messages arrive at the coordinator faster (racing the target shard’s original message), exploiting the fact that the coordinator makes progress based on the first message it receives. Replaying messages in this fashion enables the attacker to compromise the system safety (by creating inconsistent state on the shards) and/or liveness (by causing valid transactions to be rejected).

To attack the second phase (*commit*) of the atomic commit protocol, the attacker simply replays prerecorded messages to target shards, and compromises consistency. The attacker can replay those messages at any time of its choice, and does not rely on any racing condition as in the previous case.

Threat Model. The attacker can successfully launch the described attacks without colluding with any shard nodes, and under the BFT honest majority safety assumption for nodes within shards (*i.e.*, the attacks are effective even if *all* nodes are honest). We assume an attacker that can observe and record messages generated by shards; this can be achieved by (i) monitoring the network, or (ii) reading the blockchain (which is more practical). The attacker can be an external observer that passively collects the target messages at the level of the network, or it can act as a client and actively interact with the system to elicit the target messages. The attacks against the first phase of the atomic commit protocol (Sections 4.3 and 5.3) assume a weakly synchronous network in which an attacker may delay messages and race target shards by replaying pre-recorded messages. The attacks against the second phase of the atomic commit protocol (Section 4.4 and 5.4) do not make any such assumptions on the underlying network.

Attack Implementation. We implemented a demo of the replay attacks against Chainspace [1], as an example of systems that implement shard-led cross-shard consensus protocol, in Java.⁵ We are open-sourcing the demo of our attacks⁶, and a document describing a step-by-step tutorial to execute the attacks⁷. The demo shows, in the context of a simple payment application that supports account creation and coin transfer, how an attacker can use the replay attacks described in this paper to create coins out of thin air. Note that the attacks do not rely on any strict timing assumptions—the same attacker could control the accounts of both payer and payee, as well as the client.

Notation. Operations on the blockchain are specified as *transactions*. A transaction defines some transformation on the blockchain state, and has input and output *objects* (such as UTXO entries). An object is some data managed by the blockchain, such as a bank account, a specific coin, or a hotel room. For example, $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ represents a transaction with two inputs, x_1 managed by *shard 1* and x_2 managed by *shard 2*; and three outputs, y_1 managed by *shard 1*, y_2 managed by *shard 2*, and y_3 managed by *shard 3*. We call the shards that manage the input objects *input shards*, and the shards that manage

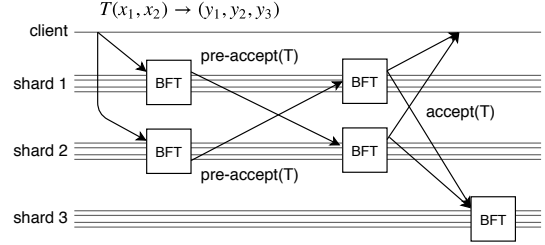


Figure 1: An example execution of S-BAC for a valid transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ with two inputs (x_1 and x_2 , both are active) and three outputs (y_1, y_2, y_3), where the final decision is $\text{accept}(T)$.

the output objects *output shards*. It is possible for a shard to be both the input and output shard. Objects can be in two states: *active* (on unspent) objects are available for being processed by a transaction; and *inactive* (or spent) objects cannot be processed by any transaction. Additionally, some systems also associate *locked* state with objects that are currently being processed by a transaction to protect against manipulation by other concurrent transactions involving those objects. The attacks we describe in this paper generalize to transactions with k inputs and k' outputs managed by an arbitrary number of shards.

4. Shard-led Cross-Shard Consensus

In shard-led cross-shard consensus protocols, the shards collectively take on the role of the coordinator in the atomic commit protocol. We describe replay attacks on shard-led cross-shard consensus protocols. To make the discussion concrete, we illustrate these attacks in the context of Chainspace [1] (Section 4.1), though we note that these attacks can be generalized to other similar systems. We discuss how the attacker can record shard messages to replay in future attacks (Section 4.2). In Sections 4.3 and 4.4, we describe replay attacks on the first and second phase of the cross-shard consensus protocol, and discuss the real-world impact of these attacks (Section 4.5).

4.1. Chainspace Overview

Chainspace uses a shard-led cross-shard consensus protocol called S-BAC. The client submits a transaction to the input shards. Each shard internally runs a BFT protocol to tentatively decide whether to accept or abort the transaction locally, and broadcasts its local decision ($\text{pre-accept}(T)$ or $\text{pre-abort}(T)$) to other relevant shards. Figure 2 shows the state machine representing the life cycle of objects in Chainspace. A shard generates $\text{pre-abort}(T)$ if the transaction fails local checks (*e.g.*, if any of the input objects are ‘inactive’ or ‘locked’). If a shard generates $\text{pre-accept}(T)$, it changes the state of the input objects to ‘locked’. This is the first step of S-BAC, and is equivalent to the voting phase in the two-phase atomic commit protocol (Section 2).

Each shard collects responses from other relevant shards, and commits the transaction if all shards respond with $\text{pre-accept}(T)$, or aborts the transaction otherwise. This is the second step of S-BAC, and is equivalent to the commit phase in the two-phase atomic commit protocol

5. Attacks against systems with client-led cross-shard consensus such as Omniledger [8] can be similarly implemented.

6. <https://github.com/sheharbano/byzcuit/tree/replay-attacks>

7. <https://github.com/sheharbano/byzcuit/blob/master/docs/Chainspace-Replay-Attack-Demo.pdf>

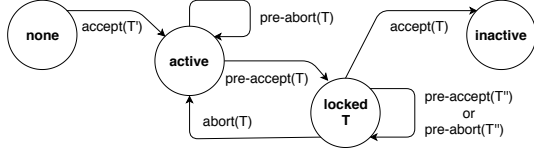


Figure 2: State machine representing the life cycle of Chainspace objects. An object becomes ‘active’ as a result of a previous successful transaction. The object state changes to ‘locked’ if a shard locally emits $\text{pre-accept}(T)$ in the first phase of the cross-shard consensus protocol for a transaction T . A locked object cannot be processed by other transactions T'' . If the second phase of the protocol results in $\text{accept}(T)$, the object becomes ‘inactive’; alternatively, if the result is $\text{abort}(T)$ the object becomes ‘active’ again and is available for being processed by other transactions.

(Section 2). The shards communicate this decision to the client as well as the output shards by sending them the $\text{accept}(T)$ or $\text{abort}(T)$ messages. If the shard’s decision is $\text{accept}(T)$, it changes the input object state to ‘inactive’. If the shard’s decision is $\text{abort}(T)$, it changes the input object state to ‘active’ (effectively unlocking it). Upon receiving $\text{accept}(T)$, the client concludes that the transaction was committed, and the output shards create the output objects (with the state ‘active’) of the transaction.

Figure 1 shows an example execution of S-BAC for a valid transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ with two inputs (x_1 and x_2 , both are active) and three outputs (y_1, y_2, y_3), where the final decision is $\text{accept}(T)$. The client submits T to shard 1 and shard 2 . Upon receiving T , both shard 1 and shard 2 confirm that the transaction is to commit, and emit $\text{pre-accept}(T)$ at the end of the first phase of S-BAC. Each shard receives $\text{pre-accept}(T)$ from the other shard, and emits $\text{accept}(T)$ at the end of the second phase of S-BAC. As a result, the input objects x_1 and x_2 become inactive, and the output shards respectively create objects y_1, y_2 , and y_3 .

4.2. Message Recording

Prior to the replay attacks, the attacker records responses generated by shards. The attacker can record shard responses in the first phase of S-BAC (*i.e.*, $\text{pre-accept}(T)$ or $\text{pre-abort}(T)$), enabling the family of attacks described in Section 4.3. The attacker can also record shard responses in the second phase of S-BAC (*i.e.*, $\text{accept}(T)$ or $\text{abort}(T)$), enabling the family of attacks described in Section 4.4. In the general case, the attacker passively collects the messages either by sniffing the network on protocol executions, or by downloading the blockchain and selecting the messages to replay⁸. Section A.1 shows how the attacker can act as client to actively elicit the messages necessary for the attacks—this empowers the attacker to actively orchestrate the attacks.

4.3. Attacks on the First Phase of S-BAC

We present replay attacks on the first phase of S-BAC by taking the example of a transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ as described in Section 3. These attacks easily generalize to transactions with k inputs and k' outputs

8. Since those messages need to be recorded on chain for verification, just using transport layer encryption between nodes is not effective.

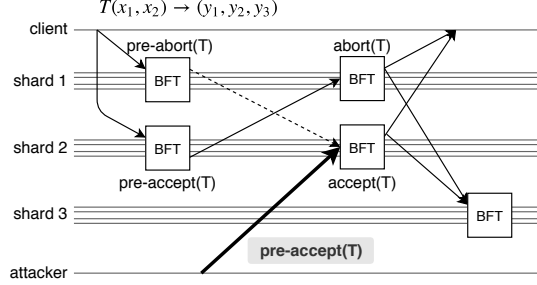


Figure 3: Illustration of the replay attack depicted in row 6 of Table 1. The attacker replays to shard 2 a prerecorded $\text{pre-accept}(T)$ message (shown as a bold line) from shard 1 , which precludes shard 1 ’s $\text{pre-abort}(T)$ message (shown as a dotted line).

managed by an arbitrary number of shards. The replay attacks work in two steps; (i) the attacker records $\text{pre-accept}(T)$ or $\text{pre-abort}(T)$ messages (as described in Section 4.2 and Section A.1); and (ii) then replays those messages during a new instance of the protocol.

Table 1 shows the replay attacks that the attacker can launch, for all possible combinations of messages emitted by shard 1 and shard 2 in the first phase of S-BAC. The caption includes details about how to interpret this table. All attacks exploit the parallel composition of multiple S-BAC instances, and insufficient binding of messages to its S-BAC instance. We describe row 6 of Table 1, to help readers interpret rest of the table on their own. In the correct execution (row 5), shard 1 and shard 2 emit $\text{pre-abort}(T)$ (because x_1 is not active) and $\text{pre-accept}(T)$ in the first phase, respectively. In the second phase, both shards emit $\text{abort}(T)$ and the protocol terminates. Figure 3 illustrates the replay attack corresponding to row 6 of Table 1. The attacker races shard 1 by sending to shard 2 the prerecorded $\text{pre-accept}(T)$ message from shard 1 . As a result, shard 2 emits $\text{accept}(T)$, inactivates object x_2 and creates object y_2 . This leads to inconsistent state across the shards. In a correct execution: (i) if T is accepted all its inputs (x_1 and x_2) should become inactive, and all the outputs (y_1, y_2, y_3) should be created; and (ii) if T is aborted, all its inputs (x_1 and x_2) should become active again, and none of the outputs (y_1, y_2, y_3) should be created. However, here we have an incorrect termination of S-BAC: at the end of the protocol x_1 could be active and x_2 is inactive; y_1 is not created, y_2 and y_3 are created.

Table 1 shows that through careful selection of the messages to replay from different S-BAC instances, the attacks can be effective against any shard. All the attacks (except row 4) compromise consistency; the attacker can trick the input shards to inactivate arbitrary objects, and trick the output shards into creating new objects in violation of the protocol. The attack depicted in row 4 only affects availability.

4.4. Attacks on the Second Phase of S-BAC

We present replay attacks on the second phase of S-BAC. The attacker prerecords $\text{accept}(T)$ messages as described in Section 4.2 and Section A.1. Table 2 shows replay attacks for all possible combinations of messages emitted by shard 1 and shard 2 in the second phase. Since the attacks we describe in this section assume that

| Phase 1 of S-BAC | | | Phase 2 of S-BAC | | |
|------------------|------------------------------------|------------------------------------|--|--|-------------------------------|
| | Shard 1 (potential victim) | Shard 2 (potential victim) | Shard 1 (potential victim) | Shard 2 (potential victim) | Shard 3 (potential victim) |
| 1 | pre-accept(T) lock x_1 | pre-accept(T) lock x_2 | accept(T) create y_1 ; inactivate x_1 | accept(T) create y_2 ; inactivate x_2 | - create y_3 |
| 2 | \triangleright pre-abort(T) | | accept(T) create y_1 ; inactivate x_1 | abort(T) unlock x_2 | - create y_3 |
| 3 | | \triangleright pre-abort(T) | abort(T) unlock x_1 | accept(T) create y_2 ; inactivate x_2 | - create y_3 |
| 4 | \triangleright pre-abort(T) | \triangleright pre-abort(T) | abort(T) unlock x_1 | abort(T) unlock x_2 | - |
| 5 | pre-abort(T) - | pre-accept(T) lock x_2 | abort(T) - | abort(T) unlock x_2 | - |
| 6 | \triangleright pre-accept(T) | | abort(T) - | accept(T) create y_2 ; inactivate x_2 | - create y_3 |
| 7 | pre-accept(T) lock x_1 | pre-abort(T) - | abort(T) unlock x_1 | abort(T) - | - |
| 8 | | \triangleright pre-accept(T) | accept(T) create y_1 ; inactivate x_1 | abort(T) - | - create y_3 |
| 9 | pre-abort(T) - | pre-abort(T) - | abort(T) - | abort(T) - | - |

TABLE 1: List of replay attacks against the first phase of S-BAC for all possible executions of the transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ as described in Section 3. The highlighted rows indicate correct executions of S-BAC (*i.e.*, without the attacker), and the other rows indicate incorrect executions due to the replay attacks. In multirows, the top sub-rows show the protocol messages emitted by shards, and the bottom sub-rows indicate local shard actions as a result of emitting those messages. For example, (column 3, row 2) means that *shard* 1 emits **accept(T)** (top sub-row), and creates a new object y_1 and inactivates x_1 (bottom sub-row). The first two columns indicate the messages emitted by each shard at the end the first phase of S-BAC. The attacker races shards at the end of the first phase of S-BAC by replaying prerecorded messages, marked with the symbol \triangleright in the first two columns of Table 1. For example \triangleright pre-abort(T) at (column 1, row 2) means that the attacker sends to other relevant shards (in this case *shard* 2) a prerecorded **pre-abort(T)** message impersonating *shard* 1 that races the original **pre-accept(T)** (column 1, row 1) emitted by *shard* 1. The last three columns indicate the messages emitted at the end of the second phase of S-BAC.

the first phase of S-BAC concluded correctly (*i.e.*, all the relevant shards unanimously decide to accept or reject a transaction), both the shards generate **abort(T)** (row 1) or **accept(T)** (row 5). The caption includes details about how to interpret this table. We describe row 6 of Table 2, to help readers interpret rest of the table on their own. In the correct execution (row 5), both shards emit **abort(T)** and no output objects are created. In the attack in row 6, the attacker replays a prerecorded **accept(T)** from *shard* 1 to all the relevant shards (in this case *shard* 3). Upon receiving this message, *shard* 3 (incorrectly) creates y_3 .

The potential victims of replay attacks corresponding to the second phase of S-BAC are the shards that *only* act as output shards (*i.e.*, do not simultaneously act as input shards). The attacker can replay **accept(T)** multiple times tricking *shard* 3 into creating y_3 multiple times. These attacks are possible because shards do not keep records of inactive objects (following the UTXO model) for scalability reasons⁹, and because *shard* 3 takes part in only the second phase of S-BAC. The attacker can double-spend y_3 repeatedly by replaying a single prerecorded message multiple times, and spending the object (*i.e.* purging it from *shard* 3's UTXO) before each replay.

Contrarily to the attacks against the first phase of S-BAC (Section 4.3), these attacks do not rely on any racing conditions; there is no need to race any honest messages.

9. Requiring shards to remember the full history of inactive objects would increase their memory requirements monotonically over time, reaching at some point memory limits preventing further operations. Thus this is a poor mitigation for the attacks presented.

4.5. Real-world Impact

The real-world impact and attacker incentives to conduct these attacks depends on the nature and implementation of the smart contract handling the target objects. We discuss the impact of these attacks in the context of two common smart contract applications, which are also described in the Chainspace paper [1]. To take a concrete example, we illustrate the attack depicted in row 3 of Table 1, but similar results can be obtained with the other attacks described in Table 1 and Table 2.

One of the most common blockchain application is to manage cryptocurrency (or coins) and enable payments for processing transactions, implemented by the CSCoin smart contract in Chainspace. Lets suppose object x_1 (handled by *shard* 1) represents Alice's account, and object x_2 (handled by *shard* 2) represents Bob's account. To transfer v coins to Bob, Alice submits a transaction $T(x_1, x_2) \rightarrow (y_1, y_2)$, where y_1 and y_2 respectively represent the new account objects of Alice and Bob, with updated account balances. By executing the attack described in row 3 of Table 1, an attacker can trick *shard* 1 to abort the transaction and unlock x_1 (thus reestablishing Alice's account balance as it was prior to the coin transfer), and *shard* 2 to accept the transaction and create y_2 (thus adding v coins to Bob's account). This attack effectively allows any attacker to double-spend coins on the ledger; and shows how to create v coins out of thin air.

Another common blockchain use case is a platform for decision making (or electronic petitions), implemented by the SVote smart contract in Chainspace. Upon ini-

| Phase 2 of S-BAC | | | |
|------------------|--|--|-------------------------------|
| | Shard 1 | Shard 2 | Shard 3 (potential victim) |
| 1 | accept(T) create y_1 ; inactivate x_1 | accept(T) create y_2 ; inactivate x_2 | - |
| 2 | \triangleright accept(T) | | create y_3 |
| 3 | | \triangleright accept(T) | create y_3 |
| 4 | \triangleright accept(T) | \triangleright accept(T) | create y_3 |
| 5 | abort(T) (unlock x_1) | abort(T) (unlock x_2) | - |
| 6 | \triangleright accept(T) | | create y_3 |
| 7 | | \triangleright accept(T) | create y_3 |
| 8 | \triangleright accept(T) | \triangleright accept(T) | create y_3 |

TABLE 2: List of replay attacks against the second phase of S-BAC for all possible executions of the transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ as described in Section 3. The highlighted rows indicate correct executions of S-BAC (*i.e.*, without the attacker), and the other rows indicate incorrect executions due to the replay attacks. In multirows, the top sub-rows show the protocol messages emitted by shards, and the bottom sub-rows indicate local shard actions as a result of emitting those messages. For example, (column 1, row 1) means that *shard* 1 emits $\text{accept}(T)$ (top sub-row), and creates a new object y_1 and inactivates x_1 (bottom sub-row). The first two columns indicate the messages emitted by each shard at the end the second phase of S-BAC, and the last column shows the effect of these messages on the output *shard* 3. Replayed messages are marked with the symbol \triangleright . For example $\triangleright \text{accept}(T)$ at (column 1, row 2) means that the attacker sends to other relevant shards (in this case *shard* 3) a prerecorded $\text{accept}(T)$ message impersonating *shard* 1.

tialization, the *SVote* contract creates two objects: (i) x_1 representing the tally’s public key, a list of all voters’ public keys, and the tally’s signature on these; and (ii) x_2 representing a vote object at the initial stage of the election (all candidates having a score of zero) along with a zero-knowledge proof asserting the correctness of the initial stage. To vote, clients submit a transaction $T(x_1, x_2) \rightarrow (y_1, y_2)$, where y_1 and y_2 are respectively the updated voting list (*i.e.*, the voting list without the client’s public key), and the election stage updated with the client’s vote. By executing the attack described by row 3 of Table 1, an attacker can trick *shard* 1 to abort the transaction and thus not update the voting list, and *shard* 2 to accept the transaction and thus update the election stage. This allows clients to vote multiple times during an election while remaining undetected (due to the privacy-preserving properties of the *SVote* smart contract).

5. Client-led Cross-shard Consensus

We describe replay attacks on client-led cross-shard consensus protocols. We illustrate these attacks in the context of Omniledger [8] (Section 5.1) to make the discussion concrete. However, we note that these attacks can be generalized to other similar systems. We discuss how the attacker can record shard messages to replay in future attacks (Section 5.2). We describe replay attacks on the first (Section 5.3) and second (Section 5.4) phase of the protocol. Finally, we discuss the real-world impact of these attacks (Section 5.5).

5.1. Omniledger Overview

Omniledger uses a client-led cross-shard consensus protocol called Atomix. The client submits the transaction T to the input shards. Each shard runs a BFT protocol locally to decide whether to accept or reject the transaction, and communicates its response ($\text{pre-accept}(T)$ or $\text{pre-abort}(T)$) to the client.¹⁰ A shard emits $\text{pre-abort}(T)$ if

10. For consistency and clarity, we use the terminology used in Section 4. In Omniledger, $\text{pre-accept}(T)$ is actually a *proof-of-accept* and $\text{pre-abort}(T)$ is a *proof-of-abort* [8].

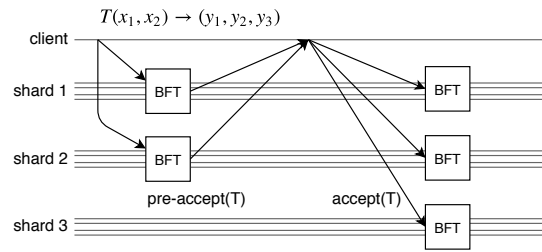


Figure 4: An example execution of Atomix for a valid transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ with two inputs (x_1 and x_2 , both are active) and three outputs (y_1, y_2, y_3), where the final decision is $\text{accept}(T)$.

the transaction fails local checks. Alternatively, if a shard emits $\text{pre-accept}(T)$, it inactivates the input objects it manages. This is the first phase of Atomix, and is similar to the voting phase in the two-phase atomic commit protocol (Section 2), but differs in that the protocol proceeds optimistically. The write changes made by the input shards in the first phase of Atomix are considered permanent (*i.e.*, there is no ‘locked’ object state), unless the client requests the input shards to revert their changes in the second phase. After the client has collected $\text{pre-accept}(T)$ from all input shards, it submits $\text{accept}(T)$ message (containing proof of the $\text{pre-accept}(T)$ messages) to the output shards which create the output objects. Alternatively, if any of the input shards emits $\text{pre-abort}(T)$, the client sends $\text{abort}(T)$ (containing proof of $\text{pre-abort}(T)$) to the relevant input shards which make the input objects active again. This is the second phase of Atomix, and is similar to the commit phase in the two-phase atomic commit protocol (Section 2).

Figure 4 shows the execution of Atomix for a valid transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$, with two active inputs (x_1 managed by *shard* 1, and x_2 managed by *shard* 2) and producing three outputs (y_1, y_2, y_3) managed by *shard* 1, *shard* 2 and *shard* 3, respectively. The client sends T to the input shards, both of which reply with $\text{pre-accept}(T)$ and make the input objects x_1 and x_2 inactive. The client sends $\text{accept}(T)$ to the output shards which respectively create objects y_1, y_2 , and y_3 .

| | Phase 1 of Atomix | | | Phase 2 of Atomix | | |
|----|---------------------------------------|---------------------------------------|--------------------|-------------------------------|-------------------------------|-------------------------------|
| | Shard 1 (potential victim) | Shard 2 (potential victim) | Client (victim) | Shard 1 (potential victim) | Shard 2 (potential victim) | Shard 3 (potential victim) |
| 1 | pre-accept(T) inactivate x_1 | pre-accept(T) inactivate x_2 | accept(T) | - create y_1 | - create y_2 | - create y_3 |
| 2 | ▷ pre-abort(T) | | abort(T) | - re-activate x_1 | - re-activate x_2 | - |
| 3 | | ▷ pre-abort(T) | abort(T) | - re-activate x_1 | - re-activate x_2 | - |
| 4 | ▷ pre-abort(T) | ▷ pre-abort(T) | abort(T) | - re-activate x_1 | - re-activate x_2 | - |
| 5 | pre-abort(T) - | pre-accept(T) inactivate x_2 | abort(T) | - - | - re-activate x_2 | - |
| 6 | ▷ pre-accept(T) | | accept(T) | - create y_1 | - create y_2 | - create y_3 |
| 7 | pre-accept(T) inactivate x_1 | pre-abort(T) - | abort(T) | - re-activate x_1 | - - | - |
| 8 | | ▷ pre-accept(T) | accept(T) | - create y_1 | - create y_2 | - create y_3 |
| 9 | pre-abort(T) - | pre-abort(T) - | abort(T) | - - | - - | - |
| 10 | ▷ pre-accept(T) | ▷ pre-accept(T) | accept(T) | - create y_1 | - create y_2 | - create y_3 |

TABLE 3: List of replay attacks against the first phase of Atomix for all possible executions of a transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ as described in Section 3. The highlighted rows indicate correct executions of Atomix (*i.e.*, without the attacker), and the other rows indicate incorrect executions due to the replay attacks. In multirows, the top sub-rows show the protocol messages emitted by shards, and the bottom sub-rows indicate local shard actions as a result of emitting those messages. For example, (column 1, row 1) means that *shard* 1 emits **pre-accept(T)** (top sub-row), and inactivates x_1 (bottom sub-row). The first two columns indicate the messages emitted by each shard at the end the first phase of Atomix. Replayed messages are marked with the symbol ▷, for example ▷**pre-abort(T)** at (column 1, row 2) means that the attacker sends to the client a prerecorded **pre-abort(T)** message impersonating *shard* 1 that races the original **pre-accept(T)** (column 1, row 1) emitted by *shard* 1. The third column indicates the messages sent by the client to the relevant shards, and the last three columns indicate the local actions performed by shards at the end of Atomix.

5.2. Message Recording

Before launching the attacks, the attacker first records the target shard responses. The attacker can record shard responses in the first phase of Atomix (*i.e.*, **pre-accept(T)** or **pre-abort(T)**), enabling the attacks described in Section 5.3. The attacker can also record shard responses in the second phase of Atomix (*i.e.*, **accept(T)** or **abort(T)**), enabling the attacks described in Section 5.4. In the general case, the attacker passively collects the messages to replay, for example by protocol executions on the network, or by downloading the blockchain and selecting the appropriate messages. Section A.2 shows how the attacker can act as a client to actively elicit and record the target messages to later use in the replay attacks.

5.3. Attacks on the First Phase of Atomix

We present replay attacks on the first phase of Atomix by taking the example of a transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ as described in Section 3. These attacks easily generalize to transactions with k inputs and k' outputs managed by an arbitrary number of shards. The replay attacks work in two steps: (i) the attacker observes the traffic and records **pre-accept(T)** or **pre-abort(T)** messages as described in Section 5.2; and (ii) then replay those messages.

Table 3 shows the replay attacks that the attacker can launch, for all possible combinations of responses generated by *shard* 1 and *shard* 2 in the first phase

of Atomix. The caption includes details about how to interpret this table. We describe row 6 of Table 3, to help readers interpret rest of the table on their own. In the correct execution (row 5), *shard* 1 emits **pre-abort(T)**, and *shard* 2 emits **pre-accept(T)** and inactivates the input objects x_2 . Upon receiving these messages, the client sends **abort(T)** to the output shards *shard* 1, *shard* 2 and *shard* 3, and *shard* 2 re-activates x_2 ; and the protocol terminates. In the attack illustrated in row 6 of Table 3, the attacker races *shard* 1 by sending to the client the prerecorded **pre-accept(T)** message from *shard* 1. As a result, the client sends **accept(T)** message to the output shards *shard* 1, *shard* 2 and *shard* 3, which respectively create the output objects y_1 , y_2 , and y_3 . As a result, the system ends up in an inconsistent state because the output objects (y_1 , y_2 , y_3) have been created, while the input object (x_1) was not active—this results in a double-spend of the input object x_1 .

Table 3 shows that through careful selection of the messages to replay, the attacks can be effective against any shard. The attacks illustrated in row 2, row 3, and row 4 only affect availability, while the other attacks compromise consistency (*i.e.*, the attacker can trick the input shards to reactivate arbitrary objects, and trick the output shards into creating new objects in violation of the protocol). The potential victims of these attacks include the client (*e.g.*, when the attacker replays the shard messages to it in the first phase of Atomix) and any input or output shards.

5.4. Attacks on the Second Phase of Atomix

We present replay attacks on the second phase of Atomix. The attacker prerecords `accept(T)` and `abort(T)` messages as described in Section 5.2 and Section A.2.

Table 4 shows replay attacks corresponding to the messages emitted by the client in the second phase—*i.e.*, `accept(T)` in row 1, or `abort(T)` in row 3. The caption includes details about how to interpret this table. The `abort(T)` message at (column 1, row 2) means that the attacker sends a prerecorded `abort(T)` message to the input shards (*shard* 1 and *shard* 2) impersonating the client. Upon receiving this message, *shard* 1 and *shard* 2 (incorrectly) re-activate x_1 and x_2 , respectively. Furthermore, all output shards create the output objects when the correct `accept(T)` message emitted by the client (row 1, column 1) reaches them. This results in inconsistent state, as the output objects are created, but the input objects are not consumed and are reactivated by the `abort(T)` message replayed by the adversary. The potential victims of `abort(T)` replay attack are the input shards.

Similarly, `accept(T)` at (row 4, column 1) means that the attacker sends a prerecorded `accept(T)` message to the output shards (*shard* 1, *shard* 2 and *shard* 3) impersonating the client. Upon receiving this message, the output shards (incorrectly) create y_1 , y_2 and y_3 . Furthermore, the input shards (*shard* 1 and *shard* 2) reactivate x_1 and x_2 upon receiving the correct `abort(T)` message emitted by the client (row 3, column 1). This creates inconsistent state: the input objects have not been consumed and have been reactivated by the `abort(T)` message emitted by the client, but the output objects have been created due to the `accept(T)` message replayed by the attacker. The potential victims of `accept(T)` replay attack are the output shards.

These attacks are possible because output shards create objects directly upon receiving `accept(T)`; they do not check if the objects have been previously invalidated because shards do not keep records of inactive objects (per the UTXO model) for scalability reasons.¹¹ The attacker can double-spend the output objects repeatedly from a single prerecorded message by replaying it multiple times, and spending the object (and effectively purging it from the output shards' UTXO) before each replay.

Similar to the attacks against the second phase of S-BAC (Section 4.4), these attacks do not exploit any racing condition and can be mounted by an adversary at a leisurely pace.

5.5. Real-world Impact

Contrarily to Chainspace, Omniledger does not support smart contracts and only handles a cryptocurrency. The attacks described in Sections 5.3 and 5.4 allow an attacker to: (i) double-spend the coins of any user, by reactivating spent coins (*e.g.*, the attacker may execute the attack depicted by row 2 of Table 4 to re-activate the objects x_1 and x_2 after the transfer is complete); and (ii) create coins out of thin air by replaying the message to create coins (*e.g.*, an attacker may execute the attack

11. Verifying that objects have not been previously invalidated implies either keep a forever-growing list of invalidated objects, or download and check the shard's entire blockchain.

depicted by row 4 of Table 4 to create multiple times object y_3 , by purging it from the UTXO list of *shard* 3 prior to each instance of the attack).

If the attacker colludes with the client, it can trigger the prerecorded messages needed for the attacks as described in Section 5.2. Alternatively, the attacker can passively observe the network and collect the target messages to replay. Similar results can be obtained using the attacks described in Table 3. Note that since transaction are recorded on the blockchain, these attacks can be detected retrospectively. This can lead to the attacker being exposed, or the attacker can inculpate innocent users (the attacker can replay messages of any user).

6. Defenses Against Replay Attacks

We identify two issues that lead to the replay attacks described in Section 4 and Section 5, and discuss how to fix those:

- First, the input shards do not have a way to know that particular protocol messages received correspond to a specific instance (or session) of the protocol. This gap in the input shards' knowledge enables an attacker to replay, mix and match, old messages leading to attacks. To address this limitation, we associate a session identifier with each transaction, which has to be crafted carefully to not degrade the performance of the protocols significantly—such as, for example, by requiring nodes to store state linearly in the number of past transactions.
- Second, in some cases the output shards are only involved in the second phase of the protocol, and therefore have no knowledge of the transaction context (to determine freshness) that is available to the input shards. This limitation can be addressed by ensuring that all shards—input and output—witnesses the entire protocol execution, rather than just a subset of protocol execution phases.

Note that the two mitigation techniques described above must be used together, as part of a single defense strategy against replay attacks.

7. The Byzcuit Protocol

We showed that both S-BAC (Sections 4.3 and 4.4) and Atomix (Sections 5.3 and 5.4) are vulnerable to replay attacks that can compromise system liveness and safety. Atomix is the simpler protocol of the two, and using the client to coordinate cross-shard communication can reduce the cost to $O(n)$ in the number of shards (by aggregating shard messages). However, an unresponsive or malicious client can permanently lock input objects by never initiating the second phase of the protocol, requiring additional design considerations (*e.g.*, a new entity that periodically unlocks input objects for transactions on which no progress has been made). On the other hand, S-BAC runs the protocol among the shards, without relying on client coordination. But this comes at the cost of increased cross-shard communication: all input shards communicate with all other input shards, which leads to communication complexity of $O(n^2)$ in the number of input shards.

Motivated by these insights, we present Byzcuit—a cross-shard atomic commit protocol (based on S-BAC)

| Phase 2 of Atomix | | | |
|-------------------------------------|-------------------------------|-------------------------------|-------------------------------|
| Client | Shard 1 (potential victim) | Shard 2 (potential victim) | Shard 3 (potential victim) |
| 1 $\text{accept}(T)$ | - create y_1 | - create y_2 | - create y_3 |
| 2 $\triangleright \text{abort}(T)$ | - re-activate x_1 | - re-activate x_2 | - |
| 3 $\text{abort}(T)$ | - re-activate x_1 | - re-activate x_2 | - |
| 4 $\triangleright \text{accept}(T)$ | - create y_1 | - create y_2 | - create y_3 |

TABLE 4: List of replay attacks against the second phase of Atomix for all possible executions of the transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ as described in Section 3. The highlighted rows indicate correct executions of Atomix (*i.e.*, without the attacker), and the other rows indicate incorrect executions due to the replay attacks. In multirows, the top sub-rows show the protocol messages emitted by shards, and the bottom sub-rows indicate local shard actions. Note that we use the multirow format for consistency reasons; in this table the first column indicates the messages emitted by the client at the beginning of the second phase of Atomix, and the last two column shows the effect of these messages on the relevant shards. Replayed messages are marked with the symbol \triangleright . For example, $\triangleright \text{abort}(T)$ at (column 1, row 2) means that the attacker sends a prerecorded $\text{abort}(T)$ message to the input shards impersonating the client.

that integrates design features from Atomix—and offers better performance and security against replay attacks. Byzcuit allocates a Transaction Manager (TM) to coordinate cross-shard communication, reducing its cost to $O(n)$ in the happy case¹²; alternatively Byzcuit also has a fallback mode in case the TM fails, similar to Atomix and traditional two phase commit protocols.

Byzcuit achieves resilience against the replay attacks described in Section 4 and Section 5, by leveraging the defense proposed in Section 6.

7.1. Byzcuit Protocol Design

We describe how Byzcuit integrates the defense presented in Section 5. To map particular protocol messages to a specific protocol instance (or session), Byzcuit associates a session identifier with each transaction. To ensure that all the relevant (input and output) shards witness all phases of the protocol execution, Byzcuit leverages the notion of *dummy objects*: each shard creates a fixed number of dummy objects upon configuration; if a shard only serves as an output shard for a transaction (and therefore will only be involved in the second phase of the protocol), Byzcuit forces it to be involved in the first phase of the protocol by implicitly including a dummy object managed by the output shard in the transaction inputs, which will create a new dummy object upon completion. As a result, the output shard also becomes an input shard (because of the inclusion of its dummy object in the transaction inputs) and witnesses the entire protocol execution, rather than just the second phase.

Byzcuit Protocol Execution. We illustrate Byzcuit taking the example of a transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ with two input objects, x_1 managed by *shard* 1 and x_2 managed by *shard* 2; and three outputs, y_1 managed by *shard* 1, y_2 managed by *shard* 2, and y_3 managed by *shard* 3.

Figure 5 illustrates the Byzcuit protocol; the client first sends the transaction to all input and output shards. Note that this is different than other protocols like S-BAC and Atomix, where the transaction is only sent to the input shards. As mentioned previously, to achieve

resilience against replay attacks, Byzcuit forces a shard that is *only* involved in creating the output objects to also become an input shard (and witness the transactional context by participating in the first phase of the protocol) by implicitly consuming one of its dummy inputs (which creates a new dummy object upon completion). Byzcuit associates a sequence number s_{x_i} to each object and dummy object (when the object is created $s_{x_i} = 0$). The sequence number is intrinsically linked to the object: when clients query shards to obtain an object x_i , they also receive the associated sequence number s_{x_i} .

When submitting the transaction T , the client also sends along a transaction sequence number $s_T = \max\{s_{x_1}, s_{x_2}, s_{d_3}\}$, where the transaction sequence number s_T is the maximum of the sequence numbers s_{x_i} of each input object x_i and dummy objects d_i (●).

Upon receiving a new pair (T, s_T) , each shard saves (T, s_T) in a local cache memory—the transaction sequence number s_T acts as session identifier associated with the transaction T . Each shard internally verifies that the transaction passes local checks, and that s_T is equal to (or bigger than) the sequence numbers of the objects they manage (*i.e.*, *shard* 1 checks $s_T \geq s_{x_1}$, *shard* 2 checks $s_T \geq s_{x_2}$, *shard* 3 checks $s_T \geq s_{d_3}$). The shards send their local decision to the TM: $\text{pre-accept}(T, s_T)$ for local accept (and the shard locks the objects it manages), or $\text{pre-abort}(T, s_T)$ for local abort.

After receiving all the messages corresponding to the first phase of Byzcuit from the concerned shards, the TM sends a suitable message to the shards ($\text{accept}(T, s_T)$ if all the shards respond with $\text{pre-accept}(T, s_T)$, or $\text{abort}(T, s_T)$ otherwise). Upon receiving $\text{accept}(T, s_T)$ or $\text{abort}(T, s_T)$ from the TM, shards first verify that they previously cached the pair (T, s_T) associated with the message; otherwise they ignore it (●).

The $\text{accept}(T, s_T)$ or $\text{abort}(T, s_T)$ messages sent by the TM provide enough evidence to the shards to verify whether s_T is correctly computed; *i.e.* shards verify that s_T is at least the maximum of the sequence numbers of each input and dummy object by inspecting the transaction T signed by each shard. If $\text{accept}(T, s_T)$ has a correct s_T , the shards inactivate the input objects and create the output objects (y_1, y_2, y_3) , and *shard* 3 creates a new dummy object d_3 ; otherwise, they update the sequence numbers

12. The communication complexity can be reduced to $O(n)$ in the number of shards by aggregating shard messages as described by Omniledger.

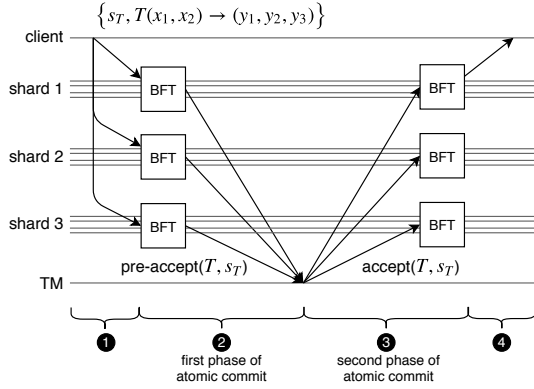


Figure 5: An example execution of Byzcuit for a valid transaction $T(x_1, x_2) \rightarrow (y_1, y_2, y_3)$ with two input objects (x_1 and x_2 , both are active), and three outputs (y_1, y_2, y_3), where the final decision is $\text{accept}(T, s_T)$.

of each input object (s_{x_1}, s_{x_2}) and dummy object d_3 to ($s_T + 1$), i.e. shards locally update $s_{x_1} \leftarrow (s_T + 1)$ and $s_{x_2} \leftarrow (s_T + 1)$, and $s_{d_3} \leftarrow (s_T + 1)$. Shards delete (T, s_T) from their local cache (③).

Since we assume that shards are honest—inline with the threat model of the systems discussed—it suffices if only one shard notifies the client of the protocol outcome; we may set any arbitrary rule to decide which shard notifies the client (e.g., the shard handling the first input object) (④). Figure 6 shows the finite state machine describing the life cycle of Byzcuit objects.

Transaction Manager. The Transaction Manager (TM) coordinates cross-shard communication in Byzcuit. We now discuss who might play the role of the TM, and argue that Byzcuit guarantees liveness even if the TM is faulty (byzantine) or crashes.

Keeping with the overall design goal of decentralization, we envision that a designated shard will act as the TM. If the shard is honest, the TM is live—and therefore progress is always made. The input shards contact in turn each node of the TM shard until they reach one honest node. The TM shard may have up to f dishonest nodes; therefore, the client or the input shards need to send messages to at least $f + 1$ nodes of the TM shard to ensure that it is received by at least one honest node¹³. Thus, as soon as the first honest node receives the message, the protocol progresses.

If the TM is the client or any centralized party, it may act arbitrarily—but this does not stall the protocol because anyone can make the protocol progress by taking over at any time the role of the TM. This is possible because the TM does not act on the basis of any secrets, therefore anyone else can take over and complete the protocols. This “anyone” may be an honest node in a shard that wants to finally unlock an object (e.g., upon a timeout); or other clients that wish to use a locked object; or it may be an external service that has a job to periodically close open Byzcuit instances. Byzcuit ensures

13. Clients may take a statistical view of availability. Given that fewer than $2/3$ of nodes in a shard are dishonest, sending the transaction to ρ nodes will fail to reach an honest node with probability only $(1/3)^\rho$. Clients may send messages sequentially to nodes, and only continue if they do not observe progress within some timeout to further reduce costs.

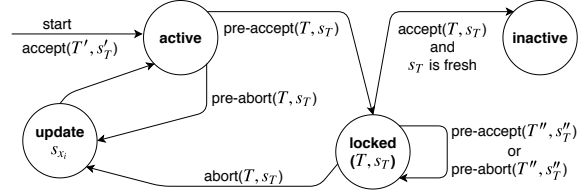


Figure 6: State machine representing the life cycle of objects in Byzcuit. Objects are initially ‘active’. Upon receiving a transaction that passes local checks, a shard changes its input objects’ state to ‘locked’ (objects are locked for a given transaction T and transaction sequence number s_T) and emits $\text{pre-accept}(T, s_T)$; otherwise it updates the sequence number of every object it manages and emits $\text{abort}(T, s_T)$. Once a shard locks input objects for a given (T, s_T) , any $\text{accept}(T, s_T)$ and $\text{abort}(T, s_T)$ with malformed transaction sequence numbers are ignored, and do not cause any transition (not included in the figure). Any incoming transaction T' that requires processing ‘locked’ input object(s) is aborted. Upon receiving $\text{accept}(T, s_T)$ with a well formed s_T , a shard makes its input objects ‘inactive’ and creates the output objects. Alternatively, upon receiving $\text{abort}(T, s_T)$ shards unlock their inputs and updates the corresponding sequence numbers.

such parties may attempt to make progress asynchronously and concurrently safely. As a result, Byzcuit guarantees liveness as long as there is at least one honest entity in the system, willing to act as TM and drive the protocol.

Handling Sequence Number Overflow. An attacker can try to exhaust the possible sequence numbers to make them overflow. The attacker submits a pair (T, s_T) such that the sequence number s_T is just below the system overflow value; the sequence numbers associate with the inputs would be again prone to the attacks described in Section 4.3¹⁴. To mitigate this issue, shards define a *clone* procedure allowing to update any of their objects to an unchanged version of themselves (i.e. it creates a fresh copy of the object). This clone procedure effectively creates a new object with serial number $s'_x = 0$. When shards detect that the serial number of one of their objects approaches the overflow value, they execute internally this clone procedure. The attacker may exploit this mechanism to DoS the system, forcing shards to constantly update their objects; as a result, the target objects are not available to users. DoS countermeasures are out of scope, and are typically addressed by introducing transaction fees.

7.2. Security against Replay Attacks

We argue that Byzcuit is resilient to replay attacks. We recall the Honest Shard assumption from Chainspace and Omniledger under which Byzcuit operates, and assume that messages are authenticated as in traditional BFT protocols.

Assumption 1. (Honest Shard [1]) *The adversary may create arbitrary smart contracts, and input arbitrary transactions into Byzcuit, however they are bound to only control up to f faulty nodes in any shard. As a result, and to ensure the correctness and liveness properties of Byzantine consensus, each shard must have a size of at least $3f + 1$ nodes. (From Chainspace [1].)*

14. Note that this overflow vulnerability is common to every system relying on nonces chosen by the users, like Byzantine Quorum Systems [10].

Any message emitted by shards comes with at least $f + 1$ signatures from nodes. Assuming honest shards, the attacker can forge at most f signatures, which is not enough to impersonate a shard. We use the Lemma below to prove the security of Byzcuit.

Lemma 1. *Under Honest Shard assumption, no attacker can obtain prerecorded messages containing a fresh transaction sequence number s_T .*

Proof. The core idea protecting Byzcuit from these replay attacks is that the attacker can only obtain prerecorded messages associated with old transaction sequence numbers s_T . The transaction sequence number s_T is fresh only if it is at least equal the maximum of the sequence number of all input and dummy objects of the transaction T . Shards update every input and dummy object sequence number upon aborting transactions in such a way that sequence numbers only increase. That is, after emitting $\text{pre-accept}(T, s_T)$ or $\text{pre-abort}(T, s_T)$, either the sequence number of all input and dummy objects of T are updated to a value bigger than s_T (in case of $\text{pre-abort}(T, s_T)$), or the objects are inactivated which prevents any successive transaction to use them as input (in case of $\text{pre-accept}(T, s_T)$). It is therefore impossible for the adversary to hold a prerecorded message for a fresh s_T since the only prerecorded messages that the adversary can obtain contain sequence numbers smaller than s_T . \square

Security of the first phase of Byzcuit. An attacker may try to replay $\text{pre-accept}(T, s_T)$ and $\text{pre-abort}(T, s_T)$ during the first phase of the protocol, similarly to the attacks described in Sections 4.3 and 5.3; the TM then aggregates these messages into either $\text{accept}(T, s_T)$ or $\text{abort}(T, s_T)$, and forwards them to the shards during the second phase of the protocol.

Theorem 1 shows that Byzcuit detects that they originate from replayed messages and ignores them. Intuitively, the transaction sequence number s_T acts as a monotonically increasing session identifier associated with the transaction T ; the attacker cannot obtain prerecorded messages containing a fresh s_T . Byzcuit shards can then distinguish replayed messages (i.e., messages with old s_T) from the messages of the current instance of the protocol (i.e., messages with fresh s_T).

Theorem 1. *Under Honest Shard assumption, Byzcuit ignores $\text{accept}(T, s_T)$ and $\text{abort}(T, s_T)$ messages issued from replayed $\text{pre-accept}(T, s_T)$ and $\text{pre-abort}(T, s_T)$.*

Proof. Figure 6 shows that once Byzcuit locks objects for a particular pair (T, s_T) , the protocol can only progress toward $\text{accept}(T, s_T)$ or $\text{abort}(T, s_T)$; i.e. shards can either accept or abort the transaction T . The attacker aims to trick one or more shards to incorrectly accept or abort T by injecting prerecorded messages during the first phase of Byzcuit; we show that the attacker fails in every scenario.

Suppose that a transaction T should abort (the TM outputs $\text{abort}(T, s_T)$), but the attacker tries to trick some shards to accept the transaction. Figure 6 shows that the attacker can only succeed the attack if they gather $\text{accept}(T, s_T)$ containing a fresh transaction sequence number s_T . Lemma 1 states that no attacker can obtain prerecorded messages over a fresh transaction sequence number s_T ; therefore the only messages available to the

adversary at this point of the protocol are (at most) $n - 1$ $\text{pre-accept}(T, s_T)$ and (at most) n $\text{abort}(T, s_T)$, where n is the number of concerned shards. This is not enough to form an $\text{accept}(T, s_T)$ message with a fresh transaction sequence number s_T (which is composed of n $\text{pre-accept}(T, s_T)$ messages); therefore the attacker cannot trick any shard to accept the transaction.

Suppose that a transaction T should be accepted (the TM outputs $\text{accept}(T, s_T)$ with a fresh s_T), but the attacker tries to trick some shards to abort the transaction. Figure 6 shows that Byzcuit does not require a fresh transaction sequence number s_T to abort transactions (the freshness of s_T is only enforced upon accepting a transaction); but shards locked the input and dummy objects of the transaction for the pair (T, s_T) (with fresh s_T), so the attacker needs to gather $\text{abort}(T, s_T)$ containing the same transaction sequence number s_T locked by shards. Lemma 1 shows that the attacker cannot obtain prerecorded messages over fresh s_T ; therefore the only messages available to the adversary containing the (fresh) s_T locked by shards at this point of the protocol are n $\text{pre-accept}(T, s_T)$. This is not enough to form an $\text{abort}(T, s_T)$ message (which is composed of at least one $\text{pre-abort}(T, s_T)$); therefore the attacker cannot trick any shard to abort the transaction. \square

Security of the second phase of Byzcuit. An attacker may try to replay $\text{accept}(T, s_T)$ and $\text{abort}(T, s_T)$ messages during the second phase of the protocol, similarly to the attacks described in Sections 4.4 and 5.4.

Theorem 2 shows that Byzcuit ignores those replayed messages. Intuitively, these attacks target shards acting only as output shards (and not also as input shards) and exploit the fact that they are only involved in the second phase of the protocol, and therefore have no knowledge of the transaction context (to determine freshness) that is available to the input shards. Byzcuit is resilient to these replay attacks as it is designed in such a way that there are no shards that act only as output shards; all output shards are forced to also become input shards, by introducing dummy objects if they do not manage any input objects; this prevents the attacks by removing the attack target.

Theorem 2. *Under Honest Shard assumption, Byzcuit ignores replayed $\text{accept}(T, s_T)$ and $\text{abort}(T, s_T)$ messages.*

Proof. Figure 6 shows that shards only act upon $\text{accept}(T, s_T)$ and $\text{abort}(T, s_T)$ messages if they have the pair (T, s_T) saved in their local cache¹⁵. Shards save a pair (T, s_T) in their local cache upon emitting $\text{pre-accept}(T, s_T)$ or $\text{pre-abort}(T, s_T)$, and delete it at the end of the protocol; therefore the only attack windows where the adversary can replay $\text{accept}(T, s_T)$ and $\text{abort}(T, s_T)$ messages is while the transaction T (associated with s_T) is being processed by the second phase of the protocol. This forces the attacker to operate under the same conditions as Theorem 1. \square

Appendix B shows that Byzcuit guarantees liveness, consistency and validity, similarly to S-BAC.

15. Contrarily to S-BAC and Atomix, all Byzcuit shards have the pair (T, s_T) in their local cache after as they all participate to the first phase of the protocol.

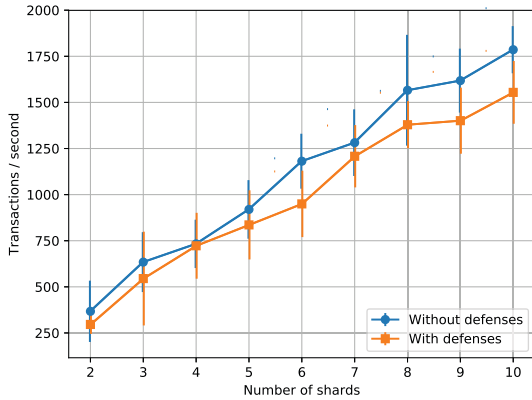


Figure 7: The effect of the number of shards on throughput. Each transaction has 2 input objects and 5 output objects, both chosen randomly from shards.

8. Implementation & Evaluation

We implement a prototype of Byzcuit (Section 7) in Java and evaluate its performance and scalability. To analyze the overhead introduced by our replay attack defenses (*i.e.*, with message sequence numbers and dummy objects), we compare Byzcuit with replay defenses (`byzcuit`) with the baseline of Byzcuit without any replay attack defenses (`byzcuit-baseline`).

Our implementation of Byzcuit is a fork of the Chainspace code [1], and is released as an open-source project¹⁶. For BFT consensus, we use the BFT-SMART [3] Java library (based on PBFT [5]), which is one of the very few maintained open source BFT libraries. End users run a client to communicate with Byzcuit nodes, which sends transactions according to the BFT-SMART protocol. The Byzcuit client also acts as the Transaction Manager (TM) and is responsible for driving the cross-shard consensus.

We evaluate the performance and scalability of our Byzcuit implementation through deployments on Amazon EC2 containers. We also compare Byzcuit with Chainspace to measure performance improvements, by running our evaluations in a similar setup as Chainspace. We launch up to 96 instances for shard nodes and 96 instances for clients on *t2.medium* virtual machines, each containing 8 GB of RAM on 2 virtual CPUs and running GNU/Linux Debian 8.1. We use 4 nodes per shard. Each measured data point corresponds to 10 runs represented by error bars. The error bars in Figure 7 and Figure 8 show the average and standard deviation, and the error bars in Figure 9 show the median and the 75th and 25th percentiles.

Throughput and Scalability. Figure 7 shows the throughput of Byzcuit (the number of transactions processed per second, tps) corresponding to an increasing number of shards. Each transaction has 2 input objects and 5 output objects, chosen randomly from shards. We test transactions with 5 output objects for a fair evaluation of Byzcuit’s replay defenses by triggering the creation of dummy objects (*i.e.*, a large number of output objects and a small number of input objects implies a higher probability of output-only shards getting selected, triggering the

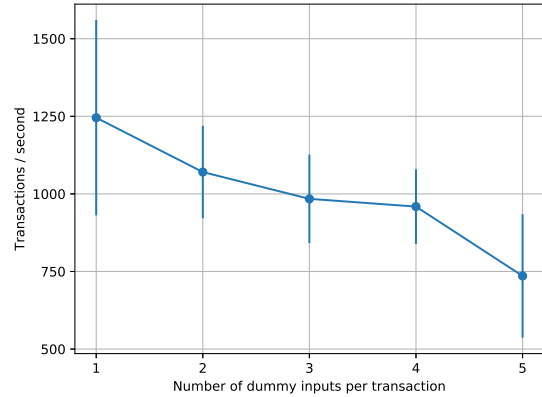


Figure 8: Decrease of Byzcuit throughput with the number of dummy objects. Each transaction has 1 input object, and up to 5 dummy objects randomly selected from unique non-input shards. 6 shards are used.

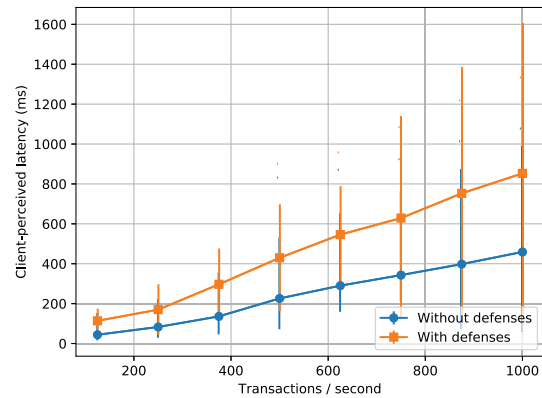


Figure 9: Client-perceived latency vs. system load (number of transactions received per second by Byzcuit), for 6 shards with 2 inputs and 5 outputs per transaction (both chosen randomly from shards).

creation of dummy objects). We find that `byzcuit` has a throughput of 260 tps for 2 shards, and linearly scales with the addition of more shards achieving up to 1550 tps for 10 shards. As expected, the throughput of `byzcuit` is lower than `byzcuit-baseline` by a somewhat constant factor ranging from 20–200 tps, but still increases linearly. This is expected because the creation of dummy objects in `byzcuit` leads to a higher number of shards processing the same transaction compared to `byzcuit-baseline`, leading to lower concurrency and lower throughput.

Another interesting observation is that the design and implementation optimizations in `byzcuit` lead to significantly higher throughput than Chainspace, even though the former has lower concurrency due to the dummy objects. For similar experimental setup and for 2–10 shards, Chainspace achieves 70–180 tps, while `byzcuit` achieves 260–1550 tps. This is due to the improved design of the cross-shard consensus protocol (Section 7), which results in communication complexity of $O(n)$ in contrast to Chainspace’s $O(n^2)$ (where n is the number of input shards). Another reason for `byzcuit`’s significant throughput improvement is that unlike Chainspace, all interactions between the Transaction Manager and the shards are asynchronous. This eliminates the blocking condition in Chainspace where a shard cannot commit a

16. <https://github.com/sheharbano/byzcuit>

transaction in the second phase of the cross-shard consensus protocol, until it receives messages from all concerned shards corresponding to the first phase.

The Effect of Dummy Objects on Throughput. We previously observed that dummy objects reduce the throughput of byzcuit with respect to byzcuit-baseline. Figure 8 shows the extent of throughput degradation due to dummy objects. We submit specially crafted transactions to 6 shards, such that each transaction has 1 input object, and we vary the number of dummy objects from 1–5 selected from unique shards, resulting in a corresponding decrease in concurrency because as many shards end up processing the transaction. For example, 2 dummy objects means that 3 shards process the transaction (1 input shard, and 2 more shards corresponding to the dummy objects). As expected, the throughput decreases by 20–250 tps with the addition of each dummy object, and reaches 750 tps when all 6 shards handle all transactions.

Client-perceived Latency. Figure 9 shows the client-perceived latency—the time from when a client submits a transaction, until it receives a decision from Byzcuit about whether the transaction has been committed—under varying system loads (expressed as transactions submitted to Byzcuit per second). We submit a total of 1200 transactions at 200–1000 transactions per second to Byzcuit with 6 shards. Each transaction has 2 inputs objects and 5 output objects, both chosen randomly from shards. When the system is experiencing a load of up to 1000 tps, clients hear back about their transactions in less than a second on average, even with our replay attack defenses.

9. Comparison with Mutex-based Cross-shard Consensus Protocols

Mutex-based schemes for cross-shard transactions, such as Ethereum’s cross-shard “yanking” proposal [4], find a way to avoid complex cross-shard coordination for transactions that involve objects managed by different shards. The key idea is to require all objects that a transaction reads or writes to be in the same shard (*i.e.*, all locks for a transaction are local to the shard). Cross-shard transactions are enabled by transferring the concerned objects between shards, effectively giving shards a lock on those objects. When *shard* 1 transfers an object to *shard* 2, *shard* 1 includes a transfer “receipt” in its blockchain. A client can then send to *shard* 2 a Merkle proof of this receipt being included in *shard* 1’s blockchain, which makes the object active in *shard* 2.

Mutex-based schemes also need to consider replay attacks. Clients can claim the same receipt multiple times, unless shards store information about previously claimed receipts. Naïvely, shards have to store information about all previously claimed receipts permanently. However, two intermediate options with trade-offs have been proposed [4]:

- Shards only store information about receipts for l blocks; so clients can only claim receipts within l blocks, and objects are permanently lost if not claimed within l blocks.
- Shards only store information about receipts for l blocks, and include the root of a Merkle tree of claimed receipts in their blockchain every l blocks.

If a receipt is not claimed within l blocks, the client must provide one Merkle proof every l blocks that have passed to show that the receipt has not been previously claimed, in order to claim it. The longer the receipt was not claimed, the greater the number of proofs that are needed to claim a receipt. These proofs need to be also stored on-chain to allow other nodes to validate them.

Byzcuit forgoes the need for shards to store information about old state (such as inactive objects or old receipts) as shards only need to know the set of active objects they manage, and does not impose a trade-off between the amount of information about old state that needs to be stored and the cost of recovering old state that was held up in an incomplete cross-shard transaction (*i.e.*, an unclaimed receipt).

10. Conclusion

We presented the first replay attacks against cross-shard consensus protocols in sharded distributed ledgers. These attacks affect both shard-driven and client-driven consensus protocols, and allow attackers to double-spend or lock objects with minimal efforts. The attacker can act independently without colluding with any nodes, and succeed even if all nodes are honest; most of the attacks work without making any assumptions on the underlying network. While addressing these attacks seems like an implementation detail, their many variants illustrate that a fundamental re-think of cross-shard commit protocols is required to protect against them.

We developed Byzcuit, a new cross-shard consensus protocol merging features from shard-led and client-led consensus protocols, and withstanding replay attacks. Byzcuit can be seen as unifying Atomix (from Omniledger) and S-BAC (from Chainspace), into an $O(n)$ protocol, that is efficient and secure. We implemented a prototype of Byzcuit and evaluated it on a real cloud-based testbed, showing that it is more efficient than Chainspace, and on par with Omniledger performance. The resulting protocol is a drop-in replacement for either, and can be adopted to immunize systems based on those designs.

Acknowledgements

At the time of this work, George Danezis, Shehar Bano and Alberto Sonnino were supported in part by EP-SRC Grant EP/N028104/1 and the EU H2020 DECODE project under grant agreement number 732546 as well as chainspace.io. Mustafa Al-Bassam is supported by The Alan Turing Institute. We thank Eleftherios Kokoris-Kogias for helpful suggestions on early manuscripts. We appreciate the valuable feedback we received from our shepherd and the anonymous reviewers.

References

- [1] AL-BASSAM, M., SONNINO, A., BANO, S., HRYCYSZYN, D., AND DANEZIS, G. Chainspace: A Sharded Smart Contracts Platform. In *Proceedings of the Network and Distributed System Security Symposium* (2018).

- [2] BANO, S., SONNINO, A., AL-BASSAM, M., AZOUVI, S., MCCORRY, P., MEIKLEJOHN, S., AND DANEZIS, G. SoK: Consensus in the Age of Blockchains. In *Proceedings of the ACM Conference on Advances in Financial Technologies* (2019).
- [3] BESSANI, A., SOUSA, J. A., AND ALCHIERI, E. E. P. State Machine Replication for the Masses with BFT-SMART. In *Proceedings of the IEEE/IFIP International Conference on Dependable Systems and Networks* (2014).
- [4] BUTERIN, V. Cross-Shard Contract Yanking. <https://ethresear.ch/t/cross-shard-contract-yanking/1450>, 2018.
- [5] CASTRO, M., AND LISKOV, B. Practical Byzantine Fault Tolerance. In *Proceedings of the Symposium on Operating Systems Design and Implementation* (1999).
- [6] DANEZIS, G., AND MEIKLEJOHN, S. Centrally Banked Cryptocurrencies. In *Proceedings of the Network and Distributed System Security Symposium* (2016).
- [7] GRAY, J. Notes on data base operating systems. In *Operating Systems, An Advanced Course* (1978).
- [8] KOKORIS KOGIAS, E., JOVANOVIĆ, P. S., GASSER, L., GAILLY, N., SYTA, E., AND FORD, B. A. Omniledger: A secure, scalable, decentralized ledger via sharding. In *Proceedings of the IEEE Symposium on Security and Privacy* (2018).
- [9] LUU, L., NARAYANAN, V., ZHENG, C., BAWEJA, K., GILBERT, S., AND SAXENA, P. A secure sharding protocol for open blockchains. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security* (2016).
- [10] MALKHI, D., AND REITER, M. Byzantine quorum systems. *Distributed computing* 11, 4 (1998), 203–213.
- [11] NAKAMOTO, S. Bitcoin: A Peer-to-Peer Electronic Cash System. <https://bitcoin.org/bitcoin.pdf>, 2008.
- [12] ZAMANI, M., MOVAHEDI, M., AND RAYKOVA, M. Rapidchain: Scaling blockchain via full sharding. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security* (2018).

Appendix A.

Eliciting Messages to Replay

This appendix shows how the attacker can act as (or collude with) a client to actively elicit and record the target messages to later use in the replay attacks. This empowers the attacker to actively orchestrate the attacks.

We describe how the attacker can trigger target messages in the context of an example, without loss of generality. Lets assume that *shard* 1 manages objects x_1 ('active') and object \widetilde{x}_1 ('inactive' or non-existent), and *shard* 2 manages object x_2 ('active'); \widetilde{x}_* means any inactive object on the shard, and y_* means any output object (*i.e.*, their details do not matter).

A.1. Shard-led Cross-Shard Consensus

We show how the attacker can act as (or collude with) a client to actively elicit and record the target messages, in the context of shard-led cross-shard consensus protocols as illustrated by Section 4. To elicit $\text{pre-accept}(T)$ for a transaction $T(x_1, x_2) \rightarrow (y_*)$ (the output y_* is not relevant here) from *shard* 1, the key consideration is to closely precede the transaction with another transaction T' that: (i) locks the inputs managed by at least one other shard (in this case x_2 on *shard* 2); and (ii) to ensure that the preceding transaction T' gets ultimately aborted, and x_2 becomes active again. The steps look as follows:

- The attacker submits $T'(x_2, \widetilde{x}_*) \rightarrow (y_*)$ to *shard* 2. This locks x_2 .

- The attacker quickly follows up by submitting $T(x_1, x_2) \rightarrow (y_*)$ to *shard* 1 and *shard* 2. *Shard* 1 generates $\text{pre-accept}(T)$, which is the target message that the attacker records. *Shard* 2 generates $\text{pre-abort}(T)$ because x_2 is locked by T' . Consequently, in the second phase of S-BAC, both *shard* 1 and *shard* 2 end up aborting T .
- T' is eventually aborted, making x_2 active again.

To elicit $\text{pre-abort}(T)$ for a transaction $T(x_1, x_2) \rightarrow (y_*)$ (the output y_* is not relevant here) from *shard* 1, the key consideration is to closely precede the transaction with another transaction T' that locks the input managed by the shard (in this case x_1 on *shard* 1). The steps look as follows:

- The attacker submits $T'(x_1, \widetilde{x}_*) \rightarrow (y_*)$ to *shard* 1. This locks x_1 .
- The attacker quickly follows up by submitting $T(x_1, x_2) \rightarrow (y_*)$ to *shard* 1 and *shard* 2. *Shard* 1 generates $\text{pre-abort}(T)$ because x_1 is locked by T' , which is the target message that the attacker records. *Shard* 2 generates $\text{pre-accept}(T)$. Consequently, in the second phase of S-BAC, both *shard* 1 and *shard* 2 end up aborting T .
- T' is eventually aborted, making x_1 active again.

To elicit $\text{accept}(T)$ used by the attacks described in Section 4.4, the attacker simply submits transaction T and observes and records its successful execution. The attacker has no incentive to record $\text{abort}(T)$ messages as these are ignored by shards (see Table 2).

A.2. Client-led Cross-Shard Consensus

We show how the attacker can act as (or collude with) a client to actively elicit and record the target messages, in the context of client-led cross-shard consensus protocols as illustrated by Section 5.

To elicit $\text{pre-accept}(T)$ from *shard* 1 for a transaction $T(x_1, x_2) \rightarrow (y_*)$ (the output y_* is not relevant here) from *shard* 1, the key consideration is to closely precede the transaction with another transaction that: (i) temporarily spends the inputs managed by at least one other shard (in this case x_2 on *shard* 2); and (ii) to ensure that the preceding transaction is ultimately aborted so that x_2 becomes active again. The steps look as follows:

- The attacker submits $T'(x_2, \widetilde{x}_*) \rightarrow (y_*)$ to *shard* 2, where \widetilde{x}_* is managed by a different shard. *Shard* 2 emits $\text{pre-accept}(T')$ and marks x_2 as inactive.
- The attacker follows up by submitting $T(x_1, x_2) \rightarrow (y_*)$ to *shard* 1 and *shard* 2. *Shard* 1 generates $\text{pre-accept}(T)$, which is the target message that the attacker records. *Shard* 2 generates $\text{pre-abort}(T)$ because x_2 is inactive.
- The attacker submits $\text{abort}(T)$ to *shard* 1 to reactivate x_1 , and sends $\text{abort}(T')$ to *shard* 2 to reactivate x_2 .

For the attacks described in Section 5.4, the attacker needs to elicit $\text{abort}(T)$ and $\text{accept}(T)$ from the target shards. For the former, the attacker can follow the steps described previously to elicit $\text{pre-accept}(T)$ and $\text{pre-abort}(T)$. To elicit $\text{accept}(T)$, the attacker simply submits transaction T and observes and records its successful execution.

Appendix B.

Byzcuit Security & Correctness

We show that Byzcuit guarantees liveness, consistency, and validity similarly to S-BAC.

Theorem 3. (*Liveness [1]*) *Under Honest Shards assumption, a transaction T that is proposed to at least one honest concerned node, eventually results in either being committed or aborted, namely all parties deciding $\text{accept}(T, s_T)$ or $\text{abort}(T, s_T)$. (From Chainspace [1].)*

Proof. We rely on the liveness properties of the byzantine agreement (shards with only f nodes eventually reach consensus on a sequence), and the broadcast from nodes of shards to all other nodes of shards, channelled through the Transaction Manager. Assuming T has been given to an honest node, it will be sequenced withing an honest shard BFT sequence, and thus a $\text{pre-accept}(T, s_T)$ or $\text{pre-abort}(T, s_T)$ will be sent from the $2f + 1$ honest nodes of this shard, aggregated into $\text{accept}(T, s_T)$ or $\text{abort}(T, s_T)$, and sent to the $f + 1$ nodes of the other concerned shards. Upon receiving these messages the honest nodes from other shards will process the transaction within their shards, and the BFT will eventually sequence it. Thus the user will eventually receive a decision from at least $f + 1$ nodes of a shard. \square

Theorem 4. (*Consistency [1]*) *Under Honest Shards assumption, no two conflicting transactions, namely transactions sharing the same input will be committed. Furthermore, a sequential executions for all transactions exists. (From Chainspace [1].)*

Proof. A transaction is accepted only if some nodes receive $\text{accept}(T, s_T)$, which presupposes all shards have provided enough evidence to conclude $\text{pre-accept}(T, s_T)$ for each of them. Two conflicting transaction, sharing an input, must share a shard of at least $3f + 1$ concerned nodes for the common object—with at most f of them being malicious. Without loss of generality upon receiving the $\text{pre-accept}(T, s_T)$ message for the first transaction, this shard will sequence it, and the honest nodes will emit messages for all—and will lock this object until the two phase protocol concludes. Any subsequent attempt to $\text{pre-accept}(T, s_T)$ for a conflicting T' will result in a $\text{pre-abort}(T, s_T)$ and cannot yield a accept , if all other shards are honest majority too. After completion of the first $\text{accept}(T, s_T)$ the shard removes the object from the active set, and thus subsequent T' would also lead to $\text{pre-abort}(T, s_T)$. Thus there is no path in the chain of possible interleavings of the executions of two conflicting transactions that leads to them both being committed. \square

Theorem 5. (*Validity [1]*) *Under Honest Shards assumption, a transaction may only be accepted if it is valid according to the smart contract (or application) logic. (From Chainspace [1].)*

Proof. A transaction is committed only if some nodes conclude that $\text{accept}(T, s_T)$, which presupposes all shards have provided enough evidence to conclude $\text{pre-accept}(T, s_T)$ for each of them. The concerned nodes include at least one shard per input object for the transaction; for any contract logic represented in the transaction,

at least one of those shards will be managing object from that contract. Each shard checks the validity rules for the objects they manage (ensuring they are active) and the contracts those objects are part of (ensuring the transaction is valid with respect to the contract logic) in order to $\text{pre-accept}(T, s_T)$. Thus if all shards say $\text{pre-accept}(T, s_T)$ to conclude that $\text{accept}(T, s_T)$, all object have been checked as active, and all the contract calls within the transaction have been checked by at least one shard—whose decision is honest due to at most f faulty nodes. If even a single object is inactive or locked, or a single trace for a contract fails to check, then the honest nodes in the shard will emit $\text{pre-abort}(T, s_T)$, and the final decision will be $\text{abort}(T, s_T)$. \square