# Cerberus: Minimalistic Multi-shard Byzantine-resilient Transaction Processing

Jelle Hellings*      Daniel P. Hughes†
Joshua Primero†      Mohammad Sadoghi*

*Exploratory Systems Lab, Department of Computer Science
University of California, Davis, CA, 95616-8562, USA
†Radix DLT Ltd, Argyle Works, 29-31 Euston Road, London, NW1 2SD

## Abstract

To enable high-performance and scalable blockchains, we need to step away from traditional consensus-based fully-replicated designs. One direction is to explore the usage of *sharding* in which we partition the managed dataset over many shards that—independently—operate as blockchains. Sharding requires an efficient fault-tolerant primitive for the ordering and execution of multi-shard transactions, however.

In this work, we seek to design such a primitive suitable for distributed ledger networks with high transaction throughput. To do so, we propose Cerberus, a set of minimalistic primitives for processing single-shard and multi-shard UTXO-like transactions. Cerberus aims at maximizing parallel processing at shards while minimizing coordination within and between shards. First, we propose *Core*-Cerberus, that uses strict environmental requirements to enable simple yet powerful multi-shard transaction processing. In our intended UTXO-environment, Core-Cerberus will operate *perfectly* with respect to all transactions proposed and approved by well-behaved clients, but does not provide any guarantees for other transactions.

To also support more general-purpose environments, we propose *two* generalizations of Core-Cerberus: we propose *Optimistic*-Cerberus, a protocol that does not require any additional coordination phases in the well-behaved optimistic case, while requiring intricate coordination when recovering from attacks; and we propose *Pessimistic*-Cerberus, a protocol that adds sufficient coordination to the well-behaved case of Core-Cerberus, allowing it to operate in a general-purpose fault-tolerant environments without significant costs to recover from attacks. Finally, we compare the three protocols, showing their potential scalability and high transaction throughput in practical environments.

## 1  Introduction

The advent of blockchain applications and technology has rejuvenated interest of companies, governments, and developers in resilient distributed fully-replicated systems and the distributed ledger technology (DLT) that powers them. Indeed, in the last decade we have seen a surge of interest in reimagining systems and build them using DLT networks. Examples can be found in the financial and banking sector [15, 36, 47], IoT [41], health care [28, 37], supply chain tracking, advertising, and in databases [3, 5, 23, 44, 45]. This wide interest is easily explained, as blockchains promise to improve resilience, while enabling the federated management of data by many participants.

To illustrate this, we look at the financial sector. Current traditional banking infrastructure is often rigid, slow, and creates substantial frictional costs. It is estimated that the yearly cost of transactional friction alone is $71 billion [8] in the financial sector, creating a strong desire for alternatives. This sector is a perfect match for DLT, as it enables systems that manage digital assets
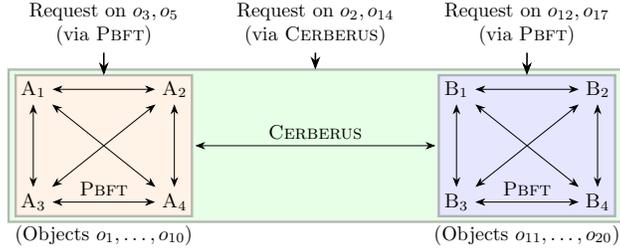
Figure 1: A *sharded* design in which two resilient blockchains each hold only a part of the data. Local decisions within a cluster are made via *traditional* PBFT *consensus*, whereas multi-shard transactions are processed via CERBERUS (proposed in this work).

and financial transactions in more flexible, fast, and open federated infrastructures that eliminate the friction caused by individual private databases maintained by banks and financial services providers. Consequently, it is expected that a large part of the financial sector will move towards DLT [18].

At the core of DLT is the *replicated state* maintained by the network in the form of a ledger of transactions. In traditional blockchains, this ledger is fully replicated among all participants using consensus protocols [14, 35, 41, 43]. For many practical use-cases, one can choose to use either permissionless consensus solutions that are operated via economic self-incentivization through cryptocurrencies (e.g., Nakamoto consensus [42, 51]), or permissioned consensus solutions that require vetted participation (e.g, PBFT [16]). Unfortunately, the design of consensus protocols utilized by todays DLT networks are severely limited in their ability to provide the *high transaction throughput* that is needed to address practical needs, e.g., in the financial and banking sector.

On the one hand, we see that permissionless solutions can easily scale to thousands of participants, but are severely limited in their transaction processing throughput. E.g., in Ethereum, a popular public permissionless DLT platform, the rapid growth of decentralized finance applications [12] has caused its network fees to rise precipitously as participants bid for limited network capacity [7], while Bitcoin can only process a few transactions per second [47]. On the other hand, permissioned solutions can reach much higher throughputs, but still lack scalability as their performance is bound by the speed of individual participants.

In this paper, we focus on a fundamental solution to significantly increase the throughput of DLT that may apply to either permissionless or permissioned networks. While this paper primarily discuss this solution through the lens of permissioned networks, similar techniques apply to permissionless DLT with the necessary extensions for these kinds of networks, such as self-incentivization, Sybil attack protection, and tolerance of validator set churn. These kinds of permissionless networks are the focus of Radix, and their impetus for their original creation of the CERBERUS concept that this paper will discuss.

A direction one can take to improve on the limited throughput of a DLT network, is to incorporate *sharding* in their design: instead of operating a single fully-replicated consensus-based DLT network, one can partition the data in the DLT network among several *shards* that each have the potential to operate mostly-independent on their data, while only requiring cooperation between shards to process transactions that affect data on several shards. In such a sharded design, transactions that only affect objects within a single shard can be processed via normal consensus (e.g., PBFT). Transactions that affect objects within several shards require additional coordination, however. The choice of protocol for such multi-shard transaction processing determines greatly the scalability benefits of sharding and the overhead costs incurred by sharding. We have sketched a basic sharded design in Figure 1.

To provide multi-shard transaction processing with high throughput in practical environments with a large number of shards, including permissionless networks, Radix proposed CERBERUS—a

technique for performing multi-shard transactions. In this paper, we propose and analyze a family of multi-shard transaction processing protocol variants using the original CERBERUS concept. To be able to adapt to the needs of specific use-cases, we propose three variants of CERBERUS: Core-CERBERUS, Optimistic-CERBERUS, and Pessimistic-CERBERUS.[1]

First, we propose Core-CERBERUS (CCERBERUS), a design specialized for processing *UTXO-like transactions*. CCERBERUS is a simplified variant of CERBERUS, and uses the strict environmental assumptions on UTXO-transactions to its advantage to yield a *minimalistic* design that does as little work as possible per involved shard. Even with this minimalistic design, CCERBERUS will operate *perfectly* with respect to all transactions proposed and approved by well-behaved clients (although it may fail to process transactions originating from malicious clients).

Next, to also support more general-purpose environments, we propose *two* generalizations of CCERBERUS, namely Optimistic-CERBERUS and Pessimistic-CERBERUS , that each deal with the strict environmental assumptions of CCERBERUS, while preserving the minimalistic design of CCERBERUS. In the design of Optimistic-CERBERUS (OCERBERUS), we assume that malicious behavior is rare and we optimize the normal-case operations. We do so by keeping the normal-case operations as minimalistic as possible. In specific, compared to CCERBERUS, OCERBERUS does not require any additional coordination phases in the well-behaved optimistic case, while still being able to lift the environmental assumptions of CCERBERUS. In doing so, OCERBERUS does require intricate coordination when recovering from attacks. In the design of Pessimistic-CERBERUS, we assume that malicious behavior is common and we add sufficient coordination to the normal-case operations of CCERBERUS to enable a simpler and localized recovery path, allowing PCERBERUS to recover from attacks at lower cost, at the expense of increased complexity in normal-case operation. Both variants we believe may be productive directions for consideration for different network deployment situations depending on desired trade-offs.

To show the strengths of each of the CERBERUS protocols, we show that CERBERUS can provide serializable transaction execution for UTXO-like transactions. Furthermore, we show that each of the protocol variants have excellent scalability in practice, even when exclusively dealing with multi-shard workloads.

**Organization**  First, in Section 2, we present the terminology and notation used throughout this paper. Then, in Section 3, we specify the *correctness criteria* by which we evaluate our CERBERUS multi-shard transaction processing protocols. Then, in Sections 4, 5, and 6, we present the three variants of CERBERUS, namely Core-CERBERUS (CCERBERUS), Optimistic-CERBERUS (OCERBERUS), and Pessimistic-CERBERUS (PCERBERUS). In Section 7, we further analyze the practical strengths, properties, and performance of CERBERUS. Then, in Section 8, we discuss related work, while we conclude on our findings in Section 9.

## 2    Preliminaries

Before we proceed with our detailed presentation of CERBERUS, we first introduce the system model, the sharding model, the data model, the transaction model, and the relevant terminology and notation used throughout this paper.

**Sharded fault-tolerant systems**  If $S$ is a set of replicas, then we write $\mathcal{G}(S)$ to denote the non-faulty *good replicas* in $S$ that always operate as intended, and we write $\mathcal{F}(S) = S \setminus \mathcal{G}(S)$ to denote the remaining replicas in $S$ that are *faulty* and can act *Byzantine*, deviate from the intended operations, or even operate in coordinated malicious manners. We write $\mathbf{n}_S = |S|$, $\mathbf{g}_S = |\mathcal{G}(S)|$, and

---

[1]The ideas underlying CERBERUS was outlined in an earlier whitepaper of our Radix team available at `https://www.radixdlt.com/wp-content/uploads/2020/03/Cerberus-Whitepaper-v1.0.pdf`.

$\mathbf{f}_S = |S \setminus \mathcal{G}(S)| = \mathbf{n}_S - \mathbf{g}_S$ to denote the number of replicas in $S$, good replicas in $S$, and faulty replicas in $S$, respectively.

Let $\mathfrak{R}$ be a set of replicas. In a *sharded fault-tolerant system* over $\mathfrak{R}$, the replicas are partitioned into sets $\mathtt{shards}(\mathfrak{R}) = \{\mathcal{S}_0, \ldots, \mathcal{S}_\mathbf{z}\}$ such that the replicas in $\mathcal{S}_i$, $0 \le i \le \mathbf{z}$, operate as an independent Byzantine fault-tolerant system. As each $\mathcal{S}_i$ operates as an independent fault-tolerant system, we require $\mathbf{n}_{\mathcal{S}_i} > 3\mathbf{f}_{\mathcal{S}_i}$, a minimal requirement to enable Byzantine fault-tolerance in an asynchronous environment [20, 21]. We assume that every shard $\mathcal{S} \in \mathtt{shards}(\mathfrak{R})$ has a unique identifier $\mathtt{id}(\mathcal{S})$.

We assume *asynchronous communication*: messages can get lost, arrive with arbitrary delays, and in arbitrary order. Consequently, it is impossible to distinguish between, on the one hand, a replica that is malicious and does not send out messages, and, on the other hand, a replica that does send out proposals that get lost in the network. As such, CERBERUS can only provide *progress* in periods of *reliable bounded-delay communication* during which all messages sent by good replicas will arrive at their destination within some maximum delay [25, 27]. Further, we assume that communication is *authenticated*: on receipt of a message $m$ from replica $\textsc{r} \in \mathfrak{R}$, one can determine that $\textsc{r}$ did sent $m$ if $\textsc{r} \in \mathcal{G}(\mathfrak{R})$. Hence, faulty replicas are able to impersonate each other, but are not able to impersonate good replicas. To provide authenticated communication under practical assumptions, we can rely on cryptographic primitives such as message authentication codes, digital signatures, or threshold signatures [38, 48].

**Assumption 2.1.** Let $\mathtt{shards}(\mathfrak{R})$ be a sharded fault-tolerant system. We assume *coordinating adversaries* that can—at will—choose and control any replica $\textsc{r} \in \mathcal{S}$ in any shard $\mathcal{S} \in \mathtt{shards}(\mathfrak{R})$ as long as, for each shard $\mathcal{S}'$, the adversaries only control up to $\mathbf{f}_{\mathcal{S}'}$ replicas in $\mathcal{S}'$.

**Object-dataset model** We use the *object-dataset model* in which data is modeled as a collection of *objects*. Each object $o$ has a unique *identifier* $\mathtt{id}(o)$ and a unique *owner* $\mathtt{owner}(o)$. In the following, we assume that all owners are *clients* of the system that manages these objects. The only operations that one can perform on an object are *construction* and *destruction*. An object cannot be recreated, as the attempted recreation of an object $o$ will result in a new object $o'$ with a distinct identifier $(\mathtt{id}(o) \ne \mathtt{id}(o'))$.

**Object-dataset transactions** Changes to object-dataset data are made via transactions requested by clients. We write $\langle \tau \rangle_c$ to denote a transaction $\tau$ requested by a client $c$. We assume that all transactions are *UTXO-like transactions*: a transaction $\tau$ first produces resources by destructing a set of *input objects* and then consumes these resources in the construction of a set of *output objects*. We do not rely on the exact rules regarding the production and consumption of resources, as they are highly application-specific. Given a transaction $\tau$, we write $\mathtt{Inputs}(\tau)$ and $\mathtt{Outputs}(\tau)$ to denote the input objects and output objects of $\tau$, respectively, and we write $\mathtt{Objects}(\tau) = \mathtt{Inputs}(\tau) \cup \mathtt{Outputs}(\tau)$.

**Assumption 2.2.** Given a transaction $\tau$, we assume that one can determine $\mathtt{Inputs}(\tau)$ and $\mathtt{Outputs}(\tau)$ a-priori. Furthermore, we assume that every transaction has inputs. Hence, $|\mathtt{Inputs}(\tau)| \ge 1$.

Owners of objects $o$ can *express their support* for transactions $\tau$ that have $o$ as their input. To provide this functionality, we can rely on cryptographic primitives such as digital signatures [38].

**Assumption 2.3.** If an owner is well-behaved, then an expression of support cannot be forged or provided by any other party. Furthermore, a well-behaved owner of $o$ will only express its support for *a single* transaction $\tau$ with $o \in \mathtt{Inputs}(\tau)$, as only one transaction can consume the object $o$, and the owner will only do so after the construction of $o$.

**Multi-shard transactions**  Let $o$ be an object. We assume that there is a well-defined function $\texttt{shard}(o)$ that maps object $o$ to the single shard $\mathcal{S} \in \texttt{shards}(\mathfrak{R})$ that is responsible for maintaining $o$. Given a transaction $\tau$, we write

$$\texttt{shards}(\tau) = \{\texttt{shard}(o) \mid o \in \texttt{Objects}(\tau)\}$$

to denote the shards that are affected by $\tau$. We say that $\tau$ is a *single-shard transaction* if $|\texttt{shards}(\tau)| = 1$ and is a *multi-shard transaction* otherwise. We assume

**Assumption 2.4.** Let $D(\mathcal{S})$ be the dataset maintained by shard $\mathcal{S}$. We have $o \in D(\mathcal{S})$ only if $\texttt{shard}(o) = \mathcal{S}$.

# 3  Correctness of multi-shard transaction processing

Before we introduce CERBERUS, we put forward the correctness requirements we want to maintain in a multi-shard transaction system in which each shard is itself a set of replicas operated as a Byzantine fault-tolerant system. We say that a shard $\mathcal{S}$ performs an action if every good replica in $\mathcal{G}(\mathcal{S})$ performs that action. Hence, any processing decision or execution step performed by $\mathcal{S}$ requires the usage of a *consensus protocol* to coordinate the replicas in $\mathcal{S}$:

**Fault-tolerant primitives**  At the core of resilient systems are *consensus protocols* [14, 16, 40, 41] that coordinate the operations of individual replicas in the system, e.g., a Byzantine fault-tolerant system driven by PBFT [16] or HOTSTUFF [53], or a crash fault-tolerant system driven by PAXOS [40]. As these systems are fully-replicated, each replica holds exactly the same data, which is determined by the *sequence of transactions*—the journal—agreed upon via consensus:

**Definition 3.1.** A *consensus protocol* coordinate decision making among the replicas of a resilient cluster $\mathcal{S}$ by providing a reliable ordered replication of *decisions*. To do so, consensus protocols provide the following guarantees:

1. If good replica $\text{R} \in \mathcal{S}$ makes a $\rho$-th decision, then all good replicas $\text{R}' \in \mathcal{S}$ will make a $\rho$-th decision (whenever communication becomes reliable).

2. If good replicas $\text{R}, \text{Q} \in \mathcal{S}$ make $\rho$-th decisions, then they make the same decisions.

3. Whenever a good replica learns that a decision $D$ needs to be made, then it can force consensus on $D$.

Let $\tau$ be a transaction processed by a sharded fault-tolerant system. Processing of $\tau$ does not imply execution: the transaction could be invalid (e.g., the owners of affected objects did not express their support) or the transaction could have inputs that no longer exists. We say that the system *commits* to $\tau$ if it decides to apply the modifications prescribed by $\tau$, and we say that the system *aborts* $\tau$ if it decides to not do so. Using this terminology, we put forward the following requirements for any sharded fault-tolerant system:

R1. *Validity.* The system must only processes transaction $\tau$ if, for every input object $o \in \texttt{Inputs}(\tau)$ with a well-behaved owner $\texttt{owner}(o)$, the owner $\texttt{owner}(o)$ supports the transaction.

R2. *Shard-involvement.* The shard $\mathcal{S}$ only processes transaction $\tau$ if $\mathcal{S} \in \texttt{shards}(\tau)$.

R3. *Shard-applicability.* Let $D(\mathcal{S})$ be the dataset maintained by shard $\mathcal{S}$ at time $t$. The shards $\texttt{shards}(\tau)$ only commit to execution of transaction $\tau$ at $t$ if $\tau$ consumes only existing objects. Hence, $\texttt{Inputs}(\tau) \subseteq \bigcup\{D(\mathcal{S}) \mid \mathcal{S} \in \texttt{shards}(\tau)\}$.

R4. *Cross-shard-consistency.* If shard $\mathcal{S}$ commits (aborts) transaction $\tau$, then all shards $\mathcal{S}' \in$ shards($\tau$) eventually commit (abort) $\tau$.

R5. *Service.* If client $c$ is well-behaved and wants to request a valid transaction $\tau$, then the sharded system will eventually *process* $\langle \tau \rangle_c$. If $\tau$ is shard-applicable, then the sharded system will eventually *execute* $\langle \tau \rangle_c$.

R6. *Confirmation.* If the system processes $\langle \tau \rangle_c$ and $c$ is well-behaved, then $c$ will eventually learn whether $\tau$ is committed or aborted.

We notice that shard-involvement is a *local requirement*, as individual shards can determine whether they need to process a given transaction. In the same sense, shard-applicability and cross-shard-consistency are *global* requirements, as assuring these requirements requires coordination between the shards affected by a transaction.

# 4 Core-Cerberus: simple yet efficient transaction processing

The core idea of CERBERUS is to minimize the coordination necessary for multi-shard ordering and execution of transactions. To do so, CERBERUS combines the semantics of transactions in the object-dataset model with the minimal coordination required to assure shard-applicability and cross-shard consistency. This combination results in the following high-level three-step approach towards processing any transaction $\tau$:

1. *Local inputs.* First, every affected shard $\mathcal{S} \in$ shards($\tau$) locally determines whether it has all inputs from $\mathcal{S}$ that are necessary to process $\tau$.

2. *Cross-shard exchange.* Then, every affected shard $\mathcal{S} \in$ shards($\tau$) exchanges these inputs to all other shards in shards($\tau$), thereby pledging to use their local inputs in the execution of $\tau$.

3. *Decide outcome* Finally, every affected shard $\mathcal{S} \in$ shards($\tau$) decides to commit $\tau$ if all affected shards were able to provide all local inputs necessary for execution of $\tau$.

Next, we describe how these three high-level steps are incorporated by CERBERUS into normal consensus steps at each shards. Let shard $\mathcal{S} \in$ shards($\mathfrak{R}$) receive client request $\langle \tau \rangle_c$. The good replicas in $\mathcal{S}$ will first determine whether $\tau$ is valid and applicable. If $\tau$ is not valid or $\mathcal{S} \notin$ shards($\tau$), then the good replicas discard $\tau$. Otherwise, if $\tau$ is valid and $\mathcal{S} \in$ shards($\tau$), then the good replicas utilize *consensus* to force the primary $\mathcal{P}(\mathcal{S})$ to propose in some consensus round $\rho$ the message $m(\mathcal{S}, \tau)_\rho = (\langle \tau \rangle_c, I(\mathcal{S}, \tau), D(\mathcal{S}, \tau))$, in which $I(\mathcal{S}, \tau) = \{o \in \text{Inputs}(\tau) \mid \mathcal{S} = \text{shard}(o)\}$ is the set of objects maintained by $\mathcal{S}$ that are input to $\tau$ and $D(\mathcal{S}, \tau) \subseteq I(\mathcal{S}, \tau)$ is the set of currently-available inputs at $\mathcal{S}$. Only if $I(\mathcal{S}, \tau) = D(\mathcal{S}, \tau)$ will $\mathcal{S}$ pledge to use the local inputs $I(\mathcal{S}, \tau)$ in the execution of $\tau$.

The acceptance of $m(\mathcal{S}, \tau)_\rho$ in round $\rho$ by all good replicas completes the *local inputs* step. Next, during execution of $\tau$, the *cross-shard exchange* and *decide outcome* steps are performed. First, the *cross-shard exchange* step. In this step, $\mathcal{S}$ broadcasts $m(\mathcal{S}, \tau)_\rho$ to all other shards in shards($\tau$). To assure that the broadcast arrives, we rely on a reliable primitive for *cross-shard exchange*, e.g., via an efficient cluster-sending protocol [29, 32]. Then, the replicas in $\mathcal{S}$ wait until they receive messages $m(\mathcal{S}', \tau)_{\rho'} = (\langle \tau \rangle_c, I(\mathcal{S}', \tau), D(\mathcal{S}', \tau))$ from all other shards $\mathcal{S}' \in$ shards($\tau$).

After cross-shard exchange comes the final *decide outcome* step. After $\mathcal{S}$ receives $m(\mathcal{S}', \tau)_{\rho'}$ from all shards $\mathcal{S}' \in$ shards($\tau$), it decides to *commit* whenever $I(\mathcal{S}', \tau) = D(\mathcal{S}', \tau)$ for all $\mathcal{S}' \in$ shards($\tau$). Otherwise, it decides *abort*. If $\mathcal{S}$ decides commit, then all good replicas in $\mathcal{S}$ destruct all objects in $D(\mathcal{S}, \tau)$ and construct all objects $o \in$ Outputs($\tau$) with $\mathcal{S} = \text{shard}(o)$. Finally, each good replica informs $c$ of the outcome of execution. If $c$ receives, from every shard $\mathcal{S}'' \in$ shards($\tau$), identical
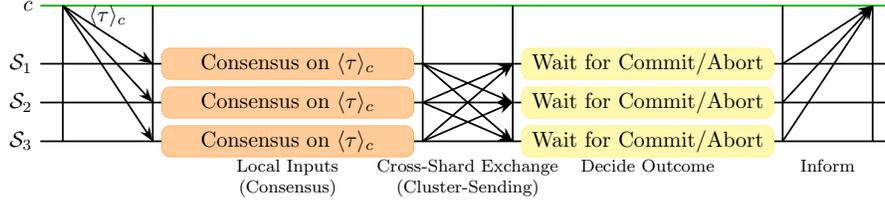
Figure 2: The message flow of CCERBERUS for a 3-shard client request $\langle\tau\rangle_c$ that is committed.

outcomes from $\mathbf{g}_{\mathcal{S}''} - \mathbf{f}_{\mathcal{S}''}$ distinct replicas in $\mathcal{S}''$, then it considers $\tau$ to be successfully executed. In Figure 2, we sketched the working of CCERBERUS.

The *cross-shard exchange* step of CCERBERUS at $\mathcal{S}$ involves waiting for other shards $\mathcal{S}'$. Hence, there is the danger of deadlocks if the other shards $\mathcal{S}'$ never perform the cross-shard exchange step:

**Example 4.1.** *Consider distinct transactions $\langle\tau_1\rangle_{c_1}$ and $\langle\tau_2\rangle_{c_2}$ that both affect objects $o_1$ in $\mathcal{S}_1$ and $o_2$ in $\mathcal{S}_2$. Hence, we have $\mathtt{Inputs}(\tau_1) = \mathtt{Inputs}(\tau_2) = \{o_1, o_2\}$ and with $\mathbf{shard}(o_1) = \mathcal{S}_1$ and $\mathbf{shard}(o_2) = \mathcal{S}_2$. We assume that $\mathcal{S}_1$ processes $\tau_1$ first and $\mathcal{S}_2$ processes $\tau_2$ first. Shard $\mathcal{S}_1$ will start by sending $m(\mathcal{S}, \tau_1)_{\rho_1} = (\langle\tau_1\rangle_{c_1}, \{o_1\}, \{o_1\})$ to $\mathcal{S}_2$. Next, $\mathcal{S}_1$ will wait, during which it receives $\tau_2$. At the same time, $\mathcal{S}_2$ follows similar steps for $\tau_2$ and sends $m(\mathcal{S}, \tau_2)_{\rho_2} = (\langle\tau_2\rangle_{c_2}, \{o_2\}, \{o_2\})$ to $\mathcal{S}_1$. As a result, $\mathcal{S}_1$ is waiting for information on $\tau_1$ from $\mathcal{S}_2$, while $\mathcal{S}_2$ is waiting for information on $\tau_2$ from $\mathcal{S}_1$.*

To assure that the above example does not lead to a deadlock, we employ two techniques.

1. *Internal propagation.* To deal with situations in which some shards $\mathcal{S} \in \mathtt{shards}(\tau)$ did not receive $\langle\tau\rangle_c$ (e.g., due to network failure or due to a faulty client that fails to send $\langle\tau\rangle_c$ to some shards), we allow each shard to learn $\tau$ from any other shard. In specific, any shard $\mathcal{S} \in \mathtt{shards}(\tau)$ will start consensus on $\langle\tau\rangle_c$ after receiving *cross-shard exchange* related to $\langle\tau\rangle_c$.

2. *Concurrent resolution.* To deal with concurrent transactions, as in Example 4.1, we allow each shard to accept and execute transactions for different rounds concurrently. To assure that such concurrent execution does not lead to inconsistent state updates, each replica implements the following *first-pledge* and *ordered-commit* rules. Let $\tau$ be a transaction with $\mathcal{S} \in \mathtt{shards}(\tau)$ and $\mathrm{R} \in \mathcal{S}$. The *first-pledge* rule states that $\mathcal{S}$ pledges $o$, constructed in round $\rho$, to transaction $\tau$ only if $\tau$ is the first transaction proposed after round $\rho$ with $o \in \mathtt{Inputs}(\tau)$. The *ordered-commit* rule states that $\mathcal{S}$ can abort $\tau$ in any order, but will only commit $\tau$ that is accepted in round $\rho$ after previous rounds finished execution.

Next, we apply the above two techniques to the situation of Example 4.1:

**Example 4.2.** *While $\mathcal{S}_1$ is waiting for $\tau_1$, it received cross-shard exchange related to $\langle\tau_2\rangle_{c_2}$. Hence, in a future round $\rho_1' > \rho_1$, it can propose and accept $\langle\tau_2\rangle_{c_2}$. By the first-pledge rule, $\mathcal{S}_1$ already pledged $o_1$ to the execution of $\tau_1$. Hence, it cannot pledge any objects to the execution of $\tau_2$. Consequently, $\mathcal{S}_1$ will eventually be able to send $m(\mathcal{S}_1, \tau_2)_{\rho_1'} = (\langle\tau_2\rangle_{c_2}, \{o_1\}, \emptyset)$ to $\mathcal{S}_2$. Likewise, $\mathcal{S}_2$ will eventually be able to send $m(\mathcal{S}_2, \tau_1)_{\rho_2'} = (\langle\tau_1\rangle_{c_1}, \{o_2\}, \emptyset)$ to $\mathcal{S}_1$. Consequently, both shards decide abort on $\tau_1$ and $\tau_2$, which they can do in any order due to the ordered-commit rule.*

Finally, we notice that abort decisions at shard $\mathcal{S}$ on a transaction $\tau$ can often be made without waiting for all shards $\mathcal{S}' \in \mathtt{shards}(\tau)$. Shard $\mathcal{S}$ can decide abort after it detects $I(\mathcal{S}, \tau) \neq D(\mathcal{S}, \tau)$ or after it receives the first message $(\langle\tau\rangle_c, I(\mathcal{S}'', \tau), D(\mathcal{S}'', \tau))$ with $I(\mathcal{S}'', \tau) \neq D(\mathcal{S}'', \tau)$, $\mathcal{S}'' \in \mathtt{shards}(\tau)$. For efficiency, we allow $\mathcal{S}$ to abort in these cases.

**Theorem 4.3.** *If, for all shards $\mathcal{S}^*$, $\mathbf{g}_{\mathcal{S}^*} > 2\mathbf{f}_{\mathcal{S}^*}$, and Assumptions 2.1, 2.2, 2.3, and 2.4 hold, then Core-*CERBERUS *satisfies Requirements R1–R6 with respect to all transactions that are not requested by malicious clients and that do not involve objects with malicious owners.*

*Proof.* Let $\tau$ be a transaction. As good replicas in $\mathcal{S}$ discard $\tau$ if it is invalid or if $\mathcal{S} \notin \texttt{shards}(\tau)$, CCERBERUS provides *validity* and *shard-involvement*. Next, *shard-applicability* follow directly from the decide outcome step.

If a shard $\mathcal{S}$ commits or aborts transaction $\tau$, then it must have completed the decide outcome and cross-shard exchange steps. As $\mathcal{S}$ completed cross-shard exchange, all shards $\mathcal{S}' \in \texttt{shards}(\tau)$ must have exchanged the necessary information to $\mathcal{S}$. By relying on cluster-sending for cross-shard exchange, $\mathcal{S}'$ requires cooperation of all good replicas in $\mathcal{S}'$ to exchange the necessary information to $\mathcal{S}$. Hence, we have the guarantee that these good replicas will also perform cross-shard exchange to any other shard $\mathcal{S}'' \in \texttt{shards}(\tau)$. As such, every shard $\mathcal{S}'' \in \texttt{shards}(\tau)$ will receive the same information as $\mathcal{S}$, complete cross-shard exchange, and make the same decision during the decide outcome step, providing *cross-shard consistency*.

Finally, due to internal propagation and concurrent resolution, every valid transaction $\tau$ will be processed by CCERBERUS as soon as it is send to any shard $\mathcal{S} \in \texttt{shards}(\tau)$. Hence, every shard in $\texttt{shards}(\tau)$ will perform the necessary steps to eventually inform the client. As all good replicas $\text{R} \in \mathcal{S}$, $\mathcal{S} \in \texttt{shards}(\tau)$, will inform the client of the outcome for $\tau$, the majority of these inform-messages come from good replicas, enabling the client to reliably derive the true outcome. Hence, CCERBERUS provides *service* and *confirmation*. $\qquad\square$

Notice that in the object-dataset model in which we operate, each object can be constructed once and destructed once. Hence, each object $o$ can be part of at-most two committed transactions: the first of which will construct $o$ as an output, and the second of which has $o$ as an input and will consume and destruct $o$. This is independent of any other operations on other objects. As such these two transactions *cannot* happen concurrently. Consequently, we only have concurrent transactions on $o$ if the owner $\texttt{owner}(o)$ expresses support for several transactions that have $o$ as an input. By Assumption 2.3, the owner $\texttt{owner}(o)$ must be malicious in that case. As such, transactions of well-behaved clients and owners will *never abort*.

In the design of CCERBERUS, we take advantage of this observation that *aborts* are always due to malicious behavior by clients and owners of objects: to minimize coordination while allowing graceful resolution of concurrent transactions, we do not undo any pledges of objects to the execution of any transactions. This implies that objects that are involved in malicious transactions can get lost for future usage, while not affecting any transactions of well-behaved clients.

# 5 Optimistic-Cerberus: robust transaction processing

In the previous section, we introduced CCERBERUS, a minimalistic and efficient multi-shard transaction processing protocol that relies on practical properties of UTXO-like transactions. Although the design of CCERBERUS is simple yet effective, we see two shortcomings that limits its use cases. First, CCERBERUS operates under the assumption that any issues arising from concurrent transactions is due to malicious behavior of clients. As such, CCERBERUS chooses to lock out objects affected by such malicious behavior for any future usage. Second, CCERBERUS requires consecutive consensus and cluster-sending steps, which increases its transaction processing latencies. Next, we investigate how to deal with these weaknesses of CCERBERUS *without giving up* on the minimalistic nature of CCERBERUS.

To do so, we propose Optimistic-CERBERUS (OCERBERUS), which is optimized for the *optimistic* case in which we have no concurrent transactions, while providing a recovery path that can recover from concurrent transactions without locking out objects. At the core of OCERBERUS is assuring that any issues due to malicious behavior, e.g., concurrent transactions, are *detected* in such a

way that individual replicas can start a recovery process. At the same time, we want to minimize transaction processing latencies. To bridge between these two objectives, we integrate detection and cross-shard coordination within a single consensus round that runs at each affected shard.

Let $\langle\tau\rangle_c$ be a multi-shard transaction, let $\mathcal{S} \in \texttt{shards}(\tau)$ be a shard with primary $\mathcal{P}(\mathcal{S})$, and let $m(\mathcal{S},\tau)_{v,\rho}$ be the round-$\rho$ proposal of $\mathcal{P}(\mathcal{S})$ of view $v$ of $\mathcal{S}$. To enable detection of concurrent transactions, OCERBERUS modifies the consensus-steps of the underlying consensus protocol by applying the following high-level idea:

> A replica $\textsc{r} \in \mathcal{S}$, $\mathcal{S} \in \texttt{shards}(\tau)$, only accepts proposal $m(\mathcal{S},\tau)_{v,\rho} = (\langle\tau\rangle_c, I(\mathcal{S},\tau), D(\mathcal{S},\tau))$ for some transaction $\tau$ if it gets confirmation that replicas in each other shard $\mathcal{S}' \in \texttt{shards}(\tau)$ are also accepting proposals for $\tau$. Otherwise, replica $\textsc{r}$ detects failure.

Next, we illustrate how to integrate the above idea in the three-phase design of PBFT, thereby turning PBFT into a multi-shard aware consensus protocol:

1. *Global preprepare.* Primary $\mathcal{P}(\mathcal{S})$ must send $m(\mathcal{S},\tau)_{v,\rho}$ to all replicas $\textsc{r}' \in \mathcal{S}'$, $\mathcal{S}' \in \texttt{shards}(\tau)$. Replica $\textsc{r} \in \mathcal{S}$ only finishes the global preprepare phase after it receives a *global preprepare certificate* consisting of a set $M = \{m(\mathcal{S}'',\tau)_{v'',\rho''} \mid \mathcal{S}'' \in \texttt{shards}(\tau)\}$ of preprepare messages from all primaries of shards affected by $\tau$.

2. *Global prepare.* After $\textsc{r} \in \mathcal{S}$, $\mathcal{S} \in \texttt{shards}(\tau)$, finishes the global preprepare phase, it sends prepare messages for $M$ to all other replicas in $\textsc{r}' \in \mathcal{S}'$, $\mathcal{S}' \in \texttt{shards}(\tau)$. Replica $\textsc{r} \in \mathcal{S}$ only finishes the global prepare phase for $M$ after, for every shard $\mathcal{S}' \in \texttt{shards}(\tau)$, it receives a *local prepare certificate* consisting of a set $P(\mathcal{S}')$ of prepare messages for $M$ from $\mathbf{g}_{\mathcal{S}'}$ distinct replicas in $\mathcal{S}'$. We call the set $\{P(\mathcal{S}'') \mid \mathcal{S}'' \in \texttt{shards}(\tau)\}$ a *global prepare certificate*.

3. *Global commit.* After replica $\textsc{r} \in \mathcal{S}$, $\mathcal{S} \in \texttt{shards}(\tau)$, finishes the global prepare phase, it sends commit messages for $M$ to all other replicas in $\textsc{r}' \in \mathcal{S}'$, $\mathcal{S}' \in \texttt{shards}(\tau)$. Replica $\textsc{r} \in \mathcal{S}$ only finishes the global commit phase for $M$ after, for every shard $\mathcal{S}' \in \texttt{shards}(\tau)$, it receives a *local commit certificate* consisting of a set $C(\mathcal{S}')$ of commit messages for $M$ from $\mathbf{g}_{\mathcal{S}'}$ distinct replicas in $\mathcal{S}'$. We call the set $\{P(\mathcal{S}'') \mid \mathcal{S}'' \in \texttt{shards}(\tau)\}$ a *global commit certificate*.

To minimize inter-shard communication, one can utilize threshold signatures and cluster-sending to carry over local prepare and commit certificates between shards via a few constant-sized messages. The above three-phase *global-*PBFT protocol already takes care of the *local input* and *cross-shard exchange* steps. Indeed, a replica $\textsc{r} \in \mathcal{S}$ that finishes the global commit phase has accepted global preprepare certificate $M$, which contains all information of other shards to proceed with execution. At the same time, $\textsc{r}$ also has confirmation that $M$ is prepared by a majority of all good replicas in each shard $\mathcal{S}' \in \texttt{shards}(\tau)$ (which will eventually be followed by execution of $\tau$ within $\mathcal{S}'$). With these ingredients in place, only the *decide outcome* step remains.

The decide outcome step at shard $\mathcal{S}$ is entirely determined by the global preprepare certificate $M$. Shard $\mathcal{S}$ decides to *commit* whenever $I(\mathcal{S}',\tau) = D(\mathcal{S}',\tau)$ for all $(\langle\tau\rangle_c, I(\mathcal{S}',\tau), D(\mathcal{S}',\tau)) \in M$. Otherwise, it decides *abort*. If $\mathcal{S}$ decides commit, then all good replicas in $\mathcal{S}$ destruct all objects in $D(\mathcal{S},\tau)$ and construct all objects $o \in \texttt{Outputs}(\tau)$ with $\mathcal{S} = \texttt{shard}(o)$. Finally, each good replica informs $c$ of the outcome of execution. If $c$ receives, from every shard $\mathcal{S}' \in \texttt{shards}(\tau)$, identical outcomes from $\mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}}$ distinct replicas in $\mathcal{S}'$, then it considers $\tau$ to be successfully executed. In Figure 3, we sketched the working of OCERBERUS.

Due to the similarity between OCERBERUS and CCERBERUS, it is straightforward to use the details of Theorem 4.3 to prove that OCERBERUS provides *validity*, *shard-involvement*, and *shard-applicability*. Next, we will focus on how OCERBERUS provides *cross-shard-consistency*. As a first step, we illustrate the ways in which the normal-case of OCERBERUS can fail (e.g., due to malicious behavior of clients, failing replicas, or unreliable communication).
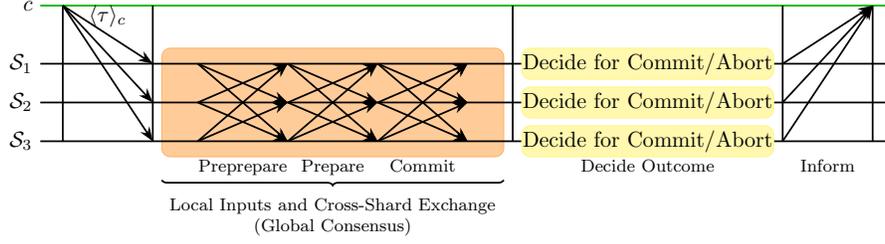
Figure 3: The message flow of OCERBERUS for a 3-shard client request $\langle\tau\rangle_c$ that is committed.

**Example 5.1.** *Consider a transaction $\tau$ proposed by client $c$ and affecting shard $\mathcal{S} \in \textbf{\textit{shards}}(\tau)$. First, we consider the case in which $\mathcal{P}(\mathcal{S})$ is malicious and tries to set up a coordinated attack. To have maximum control over the steps of* OCERBERUS*, the primary sends the message $m(\mathcal{S}, \tau)_{v,\rho}$ to only $\mathbf{g}_{\mathcal{S}''} - \mathbf{f}_{\mathcal{S}''}$ good replicas in each shard $\mathcal{S}'' \in \textbf{\textit{shards}}(\tau)$. By doing so, $\mathcal{P}(\mathcal{S})$ can coordinate the faulty replicas in each shard to assure failure of any phase at any replica $\text{R}' \in \mathcal{S}'$, $\mathcal{S}' \in \tau$:*

1. *To prevent $\text{R}'$ from finishing the global preprepare phase (and start the global prepare phase) for an $M$ with $m(\mathcal{S}', \tau)_{v',\rho'} \in M$, $\mathcal{P}(\mathcal{S})$ simply does not send $m(\mathcal{S}, \tau)_{v,\rho}$ to $\text{R}'$.*

2. *To prevent $\text{R}'$ from finishing the global prepare phase (and start the global commit phase) for $M$, $\mathcal{P}(\mathcal{S})$ instructs the faulty replicas in $\mathcal{F}(\mathcal{S})$ to not send prepare messages for $M$ to $\text{R}'$. Hence, $\text{R}'$ will receive at-most $\mathbf{g}_{\mathcal{S}} - \mathbf{f}_{\mathcal{S}}$ prepare messages for $M$ from replicas in $\mathcal{S}$, assuring that it will not receive a local prepare certificate $P(\mathcal{S})$ and will not finish the global prepare phase for $M$.*

3. *Likewise, to prevent $\text{R}'$ from finishing the global commit phase (and start execution) for $M$, $\mathcal{P}(\mathcal{S})$ instructs the faulty replicas in $\mathcal{F}(\mathcal{S})$ to not send commit messages to $\text{R}'$. Hence, $\text{R}'$ will receive at-most $\mathbf{g}_{\mathcal{S}} - \mathbf{f}_{\mathcal{S}}$ commit messages for $M$ from replicas in $\mathcal{S}$, assuring that it will not receive a local commit certificate $C(\mathcal{S})$ and will not finish the global commit phase for $M$.*

*None of the above attacks can be attributed to faulty behavior of $\mathcal{P}(\mathcal{S})$. First, unreliable communication can result in the same outcomes for $\text{R}'$. Furthermore, even if communication is reliable and $\mathcal{P}(\mathcal{S})$ is good, we can see the same outcomes:*

1. *The client $c$ can be malicious and not send $\tau$ to $\mathcal{S}$. At the same time, all other primaries $\mathcal{P}(\mathcal{S}'')$ of shards $\mathcal{S}'' \in \textbf{\textit{shards}}(\tau)$ can be malicious and not send anything to $\mathcal{S}$ either. In this case, $\mathcal{P}(\mathcal{S})$ will never be able to send any message $m(\mathcal{S}, \tau)_{v,\rho}$ to $\text{R}'$, as no replica in $\mathcal{S}$ is aware of $\tau$.*

2. *If any primary $\mathcal{P}(\mathcal{S}'')$ of $\mathcal{S}'' \in \textbf{\textit{shards}}(\tau)$ is malicious, then it can prevent some replicas in $\mathcal{S}$ from starting the global prepare phase, thereby preventing these replicas to send prepare messages to $\text{R}'$. If $\mathcal{P}(\mathcal{S}'')$ prevents sufficient replicas in $\mathcal{S}$ from starting the global prepare phase, $\text{R}'$ will be unable to finish the global prepare phase.*

3. *Likewise, any malicious primary $\mathcal{P}(\mathcal{S}'')$ of $\mathcal{S}'' \in \textbf{\textit{shards}}(\tau)$ can prevent replicas in $\mathcal{S}$ from starting the global commit phase, thereby assuring that $\text{R}'$ will be unable to finish the global commit phase.*

To deal with malicious behavior, OCERBERUS needs a robust recovery mechanism. We cannot simply build that mechanism on top of traditional view-change approaches: these traditional view-change approaches require that one can identify a single source of failure (when communication is reliable), namely the current primary. As Example 5.1 already showed, this property does not hold for OCERBERUS. To remedy this, the recovery mechanisms of OCERBERUS has components that

10

1: **event** R $\in \mathcal{S}$ is unable to finish round $\rho$ of view $v$ **do**
2:  **if** R finished in round $\rho$ the global prepare phase for $M$,
           but is unable to finish the global commit phase **then**
3:    Let $P$ be the global prepare certificate of R for $M$.
4:    **if** R has a local commit certificate $C(\mathcal{S}'')$ for $M$ **then**
5:      **for** $\mathcal{S}' \in \mathtt{shards}(\tau)$ **do**
6:        **if** R did not yet receive a local commit certificate $C(\mathcal{S}')$ **then**
7:          Broadcast $\langle \mathtt{VCGlobalSCR} : M, P, C(\mathcal{S}'') \rangle$ to all replicas in $\mathcal{S}'$.
8:      **else**
9:        Detect the need for local state recovery of round $\rho$ of view $v$ (Figure 5).
10:   **else**
11:     Detect the need for local state recovery of round $\rho$ of view $v$ (Figure 5).
12:   (Eventually repeat this event if R remains unable to finish round $\rho$.)

13: **event** R$'$ $\in \mathcal{S}'$ receives a message $\langle \mathtt{VCGlobalSCR} : M, P, C(\mathcal{S}'') \rangle$ from R $\in \mathcal{S}$ **do**
14:   **if** R$'$ did not reach the global commit phase for $M$ **then**
15:     Use $M$, $P$, and $C(\mathcal{S}'')$ to reach the global commit phase for $M$.
16:   **else**
17:     Send a commit message for $M$ to R.

Figure 4: The view-change *global short-cut recovery path* that determines whether R already has the assurance that the current transaction will be committed. If this is the case, then R requests only the missing information to proceed with execution. Otherwise, R requires at-least local recovery (Figure 5).

perform *local view-change* and that perform *global state recovery*. The pseudo-code for this recovery protocol can be found in Figure 4. Next, we describe the working of this recovery protocol in detail. Let R $\in \mathcal{S}$ be a replica that determines that it cannot finish a round $\rho$ of view $v$.

First, R determines whether it already has a *guarantee* on which transaction it has to execute in round $\rho$. This is the case when the following conditions are met: R finished the global prepare phase for $M$ with $m(\mathcal{S}, \tau)_{v,\rho} \in M$ and has received a local commit certificate $C(\mathcal{S}'')$ for $M$ from some shard $\mathcal{S}'' \in \mathtt{shards}(\tau)$. In this case, R can simply request all missing local commit certificates directly, as $C(\mathcal{S}'')$ can be used to prove to any involved replica R$'$ $\in \mathcal{S}'$, $\mathcal{S}' \in \mathtt{shards}(\tau)$, that R$'$ also needs to commit to $M$. To request such missing commit certificates of $\mathcal{S}'$, replica R sends out $\mathtt{VCGlobalSCR}$ messages to all replicas in $\mathcal{S}'$ (Line 7 of Figure 4). Any replica R$'$ that receives such a VCGLOBALSCR message can use the information in that message to reach the global commit phase for $M$ and, hence, provide R with the requested commit messages (Line 13 of Figure 4).

If R does not have a *guarantee* itself on which transaction it has to execute in round $\rho$, then it needs to determine whether any other replica (either in its own shard or in any other shard) has already received and acted upon such a guarantee. To initiate such local and global state recovery, R simply detects the current view as faulty. To do so, R broadcasts a $\mathtt{VCRecoveryRQ}$ message to all other replicas in $\mathcal{S}$ that contains all information R collected on round $\rho$ in view $v$ (Line 4 of Figure 5). Other replicas Q $\in \mathcal{S}$ that already have *guarantees* for round $\rho$ can help R by providing all missing information (Line 6 of Figure 5). On receipt of this information, R can proceed with the round (Line 7 of Figure 5). If no replicas can provide the missing information, then eventually all good replicas will detect the need for local recovery, this either by themselves (Line 1 of Figure 5) or after receiving $\mathtt{VCRecoveryRQ}$ messages of at-least $\mathbf{f}_{\mathcal{S}} + 1$ distinct replicas in $\mathcal{S}$, of which at-least a single replica must be good (Line 10 of Figure 5).

Finally, if a replica R receives $\mathbf{g}_{\mathcal{S}}$ $\mathtt{VCRecoveryRQ}$ messages, then it has the guarantee that at least $\mathbf{g}_{\mathcal{S}} - \mathbf{f}_{\mathcal{S}} \geq \mathbf{f}_{\mathcal{S}} + 1$ of these messages come from good replicas in $\mathcal{S}$. Hence, due to Line 10 of Figure 5,

1: **event** R $\in \mathcal{S}$ detects the need for local state recovery of round $\rho$ of view $v$ **do**
2:     Let $M$ be the latest global preprepare certificate accepted for round $\rho$ by R (if any).
3:     Let $S$ be $M$ and any prepare and commit certificates for $M$ collected by R.
4:     Broadcast $\langle \texttt{VCRecoveryRQ} : v, \rho, S \rangle$.

5: **event** Q $\in \mathcal{S}$ receives messages $\langle \texttt{VCRecoveryRQ} : v, \rho, S \rangle$ of R $\in \mathcal{S}$ and Q has

      1. started the global prepare phase for $M$ with $m(\mathcal{S}, \tau)_{w,\rho} \in M$;

      2. a global prepare certificate for $M$;

      3. a local commit certificate $C(\mathcal{S}'')$ for $M$

   **do**
6:     Send $\langle \texttt{VCLocalSCR} : M, P, C(\mathcal{S}'') \rangle$ to R $\in \mathcal{S}$.

7: **event** R $\in \mathcal{S}$ receives a message $\langle \texttt{VCLocalSCR} : M, P, C(\mathcal{S}'') \rangle$ from Q $\in \mathcal{S}$ **do**
8:     **if** R did not reach the global commit phase for $M$ **then**
9:       Use $M$, $P$, and $C$ to reach the global commit phase for $M$.

10: **event** R $\in \mathcal{S}$ receives messages $\langle \texttt{VCRecoveryRQ} : v_i, \rho, S_i \rangle$, $1 \le i \le \mathbf{f}_{\mathcal{S}} + 1$,
       from distinct replicas in $\mathcal{S}$ **do**
11:     R detects the need for local state recovery of round $\rho$ of view $\min\{v_i \mid 1 \le i \le \mathbf{f}_{\mathcal{S}} + 1\}$.

12: **event** R $\in \mathcal{S}$ receives messages $\langle \texttt{VCRecoveryRQ} : v, \rho, S_i \rangle$, $1 \le i \le \mathbf{g}_{\mathcal{S}}$,
       from distinct replicas in $\mathcal{S}$ **do**
13:     **if** $\texttt{id}(\text{R}) \ne (v + 1) \bmod \mathbf{n}_{\mathcal{S}}$ **then**
14:       (R awaits the $\texttt{NewView}$ message of the new primary, Line 15 of Figure 6.)
15:     **else**
16:       Broadcast $\langle \texttt{NewView} : \langle \texttt{VCRecoveryRQ} : v, \rho, S_i \rangle \mid 1 \le i \le \mathbf{g}_{\mathcal{S}} \rangle$ to all replicas in $\mathcal{S}$.
17:       **if** there exists a $S_i$ that contains global preprepare certificate $M$,
         but no $S_j$ contains a local commit certificate for $M$ **then**
18:         R initiates global state recovery of round $\rho$ (Line 1 of Figure 6).

Figure 5: The view-change *local short-cut recovery path* that determines whether some Q can provide R with the assurance that the current transaction will be committed. If this is the case, then R only needs this assurance, otherwise $\mathcal{S}$ requires a new view (Figure 6).

all $\mathbf{g}_{\mathcal{S}}$ good replicas in $\mathcal{S}$ will send $\texttt{VCRecoveryRQ}$, and, when communication is reliable, also receive these messages. Consequently, at this point, R can start the new view by electing a new primary and awaiting the $\texttt{NewView}$ proposal of this new primary (Line 12 of Figure 5). If R is the new primary, then it starts the new view by proposing a $\texttt{NewView}$. As other shards *could* have already made final decisions depending on local prepare or commit certificates of $\mathcal{S}$ for round $\rho$, we need to assure that such certificates are not invalidated. To figure out whether such final decisions have been made, the new primary will query other shards $\mathcal{S}'$ for their state whenever the $\texttt{NewView}$ message contains global preprepare certificates for transactions $\tau$, $\mathcal{S}' \in \texttt{shards}(\tau)$, but not a local commit certificate to *guarantee* execution of $\tau$ (Line 17 of Figure 5).

    The new-view process has three stages. First, the new primary P proposes the new-view via a $\texttt{NewView}$ message (Line 12 of Figure 5). If necessary, the new primary P also requests the relevant global state from any relevant shard (Line 1 of Figure 6). The replicas in other shards will respond to this request with their local state (Line 9 of Figure 6). The new primary collects these responses

1: **event** P $\in \mathcal{S}$ initiates global state recovery of round $\rho$ using $\langle \texttt{NewView}: V \rangle$ **do**
2:     Let $T$ be the transactions with global preprepare certificates for round $\rho$ of $\mathcal{S}$ in $V$.
3:     Let $S$ be the shards affected by transactions in $T$.
4:     Broadcast $\langle \texttt{VCGlobalStateRQ}: v, \rho, V \rangle$ to all replicas in $\mathcal{S}' \in S$.
5:     **for** $\mathcal{S}' \in S$ **do**
6:         Wait for VCGlobalStateRQ messages for $V$ from $\mathbf{g}_{\mathcal{S}'}$ distinct replicas in $\mathcal{S}'$.
7:         Let $W(\mathcal{S}')$ be the set of received VCGLOBALSTATERQ messages.
8:     Broadcast $\langle \texttt{NewViewGlobal}: V, \{W(\mathcal{S}') \mid \mathcal{S}' \in S\} \rangle$ to all replicas in $\mathcal{S}$.

9: **event** R$' \in \mathcal{S}'$ receives message $\langle \texttt{VCGlobalStateRQ}: v, \rho, V \rangle$ from P $\in \mathcal{S}$ **do**
10:     **if** R$'$ has a global preprepare certificate $M$ with $m(\mathcal{S}, \tau)_{w,\rho} \in M$
            and reached the global commit phase for $M$ **then**
11:         Let $P$ be the global prepare certificate for $M$.
12:         Send $\langle \texttt{VCGlobalStateR}: v, \rho, V, M, P \rangle$ to P.
13:     **else**
14:         Send $\langle \texttt{VCGlobalStateR}: v, \rho, V \rangle$ to P.

15: **event** R $\in \mathcal{S}$ receives a valid $\langle \texttt{NewView}: V \rangle$ message from replica P **do**
16:     **if** there exists a $\langle \texttt{VCRecoveryRQ}: v_i, \rho, S_i \rangle \in V$ that contains
            a global preprepare certificate $M$ with $m(\mathcal{S}, \tau)_{w,\rho} \in M$,
            a global prepare certificate $P$ for $M$, and a local commit certificate $C(\mathcal{S}'')$ for $M$ **then**
17:         Use $M$, $P$, and $C$ to reach the global commit phase for $M$.
18:     **else if** there exists a $\langle \texttt{VCRecoveryRQ}: v_i, \rho, S_i \rangle \in V$ that contains
            a global preprepare certificate $M$,
            but no $\langle \texttt{VCRecoveryRQ}: v_j, \rho, S_j \rangle \in V$ contains a local commit certificate for $M$ **then**
19:         R detects the need for global state recovery of round $\rho$ (Line 22 of Figure 6).
20:     **else**
21:         (P must propose for round $\rho$.)

22: **event** R $\in \mathcal{S}$ receives a valid $\langle \texttt{NewViewGlobal}: V, W \rangle$ from P $\in \mathcal{S}$ **do**
23:     **if** any message in $W$ is of the form $\langle \texttt{VCGlobalStateR}: v, \rho, V, M, P \rangle$ **then**
24:         Select $\langle \texttt{VCGlobalStateR}: v, \rho, V, M, P \rangle \in W$ with latest view $w$, $m(\mathcal{S}, \tau)_{w,\rho} \in M$.
25:         Use $M$ and $P$ to reach the global commit phase for $M$.
26:     **else**
27:         (P must propose for round $\rho$.)

Figure 6: The view-change *new-view recovery path* that recovers the state of the previous view based on a `NewView` proposal of the new primary. As part of the new-view recovery path, the new primary can construct a global new-view that contains the necessary information from other shards to reconstruct the local state.

and sends them to all replicas in $\mathcal{S}$ via a `NewViewGlobal` message.

Then, after P sends the `NewView` message to R $\in \mathcal{S}$, R determines whether the `NewView` message contains sufficient information to recover round $\rho$ (Line 16 of Figure 6), contains sufficient information to wait for any relevant global state (Line 18 of Figure 6), or to determine that the new primary must propose for round $\rho$ (Line 21 of Figure 6). If R determines it needs to wait for any relevant global state, then R will wait for this state to arrive via a `NewViewGlobal` message. Based on the received global state, R determines to recover round $\rho$ (Line 23 of Figure 6), or determines that the new primary must propose for round $\rho$ (Line 26 of Figure 6).

Next, we shall prove the correctness of the view-change protocol outlined in Figures 4, 5, and 6.

First, using a standard quorum argument, we prove that in a single round of a single view of $\mathcal{S}$, only a single global preprepare message affecting $\mathcal{S}$ can get committed by any other affected shards:

**Lemma 5.2.** *Let $\tau_1$ and $\tau_2$ be transactions with $\mathcal{S} \in (\mathtt{shards}(\tau_1) \cap \mathtt{shards}(\tau_2))$. If $\mathbf{g}_\mathcal{S} > 2\mathbf{f}_\mathcal{S}$ and there exists shards $\mathcal{S}_i \in \mathtt{shards}(\tau_i)$, $i \in \{1, 2\}$, such that good replicas $\mathrm{R}_i \in \mathcal{G}(\mathcal{S}_i)$ reached the global commit phase for global preprepare certificate $M_i$ with $m(\mathcal{S}, \tau_i)_{v,\rho} \in M_i$, then $\tau_1 = \tau_2$.*

*Proof.* We prove this property using contradiction. We assume $\tau_1 \neq \tau_2$. Let $P_i(\mathcal{S})$ be the local prepare certificate provided by $\mathcal{S}$ for $M_i$ and used by $\mathrm{R}_i$ to reach the global commit phase, let $S_i \subseteq \mathcal{S}$ be the $\mathbf{g}_\mathcal{S}$ replicas in $\mathcal{S}$ that provided the prepare messages in $P_i(\mathcal{S})$, and let $T_i = S_i \setminus \mathcal{F}(\mathcal{S})$ be the good replicas in $S_i$. By construction, we have $|T_i| \geq \mathbf{g}_\mathcal{S} - \mathbf{f}_\mathcal{S}$. As all replicas in $T_1 \cup T_2$ are good, they will only send out a single prepare message per round $\rho$ of view $v$. Hence, if $\tau_1 \neq \tau_2$, then $T_1 \cap T_2 = \emptyset$, and we must have $2(\mathbf{g}_\mathcal{S} - \mathbf{f}_\mathcal{S}) \leq |T_1 \cup T_2|$. As all replicas in $T_1 \cup T_2$ are good, we also have $|T_1 \cup T_2| \leq \mathbf{g}_\mathcal{S}$. Hence, $2(\mathbf{g}_\mathcal{S} - \mathbf{f}_\mathcal{S}) \leq \mathbf{g}_\mathcal{S}$, which simplifies to $\mathbf{g}_\mathcal{S} \leq 2\mathbf{f}_\mathcal{S}$, a contradiction. Hence, we conclude $\tau_1 = \tau_2$. □

Next, we use Lemma 5.2 to prove that any global preprepare certificate that *could* have been accepted by any good affected replica is preserved by OCERBERUS:

**Proposition 5.3.** *Let $\tau$ be a transaction and $m(\mathcal{S}, \tau)_{v,\rho}$ be a preprepare message. If, for all shards $\mathcal{S}^*$, $\mathbf{g}_{\mathcal{S}^*} > 2\mathbf{f}_{\mathcal{S}^*}$, and there exists a shard $\mathcal{S}' \in \mathtt{shards}(\tau)$ such that $\mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}'}$ good replicas in $\mathcal{S}'$ reached the global commit phase for $M$ with $m(\mathcal{S}, \tau)_{v,\rho} \in M$, then every successful future view of $\mathcal{S}$ will recover $M$ and assure that the good replicas in $\mathcal{S}$ reach the commit phase for $M$.*

*Proof.* Let $v^* \leq v$ be the first view in which a global prepare certificate $M^*$ with $m(\mathcal{S}, \tau^*)_{v^*,\rho} \in M^*$ satisfied the premise of this proposition. Using induction on the number of views after the first view $v^*$, we will prove the following two properties on $M^*$:

1. every good replica that participates in view $w$, $v^* < w$, will recover $M^*$ upon entering view $w$ and reach the commit phase for $M^*$; and

2. no replica will be able to construct a local prepare certificate of $\mathcal{S}$ for any global preprepare certificate $M^\dagger \neq M^*$ with $m(\mathcal{S}, \tau^\dagger)_{w,\rho} \in M^\dagger$, $v^* < w$.

The base case is view $v^* + 1$. Let $S' \subseteq \mathcal{G}(\mathcal{S}')$ be the set of $\mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}'}$ good replicas in $\mathcal{S}'$ that reached the global commit phase for $M^*$. Each replica $\mathrm{R}' \in S'$ has a local prepare certificate $P(\mathcal{S})$ consisting of $\mathbf{g}_\mathcal{S}$ prepare messages for $M^*$ provided by replicas in $\mathcal{S}$. We write $S(\mathrm{R}') \subseteq \mathcal{G}(\mathcal{S})$ to denote the at-least $\mathbf{g}_\mathcal{S} - \mathbf{f}_\mathcal{S}$ good replicas in $\mathcal{S}$ that provided such a prepare message to $\mathrm{R}'$.

Consider any valid new-view proposal $\langle \mathtt{NewView} : V \rangle$ for view $v^* + 1$. If the conditions of Line 16 of Figure 6 hold for some global preprepare certificate $M^\dagger$ with $m(\mathcal{S}, \tau^\ddagger)_{w,\rho} \in M^\ddagger$, then we recover $M^\ddagger$. As there is a local commit certificate for $M^\ddagger$ in this case, the premise of this proposition holds on $M^\ddagger$. As $v^*$ is the first view in which the premise of this proposition hold, we can use Lemma 5.2 to conclude that $w = v^*$, $M^\ddagger = M^*$, and, hence, that the base case holds if the conditions of Line 16 of Figure 6 hold. Next, we assume that the conditions of Line 16 of Figure 6 do not hold, in which case $M^*$ can only be recovered via global state recovery. As the first step in global state recovery is proving that the condition of Line 18 of Figure 6 holds. Let $T \subseteq \mathcal{G}(\mathcal{S})$ be the set of at-least $\mathbf{g}_\mathcal{S} - \mathbf{f}_\mathcal{S}$ good replicas in $\mathcal{S}$ whose $\mathtt{VCRecoveryRQ}$ message is in $V$ and let $\mathrm{R}' \in S'$. We have $|S(\mathrm{R}')| \geq \mathbf{g}_\mathcal{S} - \mathbf{f}_\mathcal{S}$ and $|T| \geq \mathbf{g}_\mathcal{S} - \mathbf{f}_\mathcal{S}$. Hence, by a standard quorum argument, we conclude $S(\mathrm{R}') \cap T \neq \emptyset$. Let $\mathrm{Q} \in (S(\mathrm{R}') \cap T)$. As $\mathrm{Q}$ is good and send prepare messages for $M^*$, it must have reached the global prepare phase for $M^*$. Consequently, the condition of Line 18 of Figure 6 holds and to complete the proof, we only need to prove that any well-formed $\mathtt{NewViewGlobal}$ message will recover $M^*$.

Let $\langle \mathtt{NewViewGlobal} : V, W \rangle$ be any valid global new-view proposal for view $v^* + 1$. As $\mathrm{Q}$ reached the global prepare phase for $M^*$, any valid global new-view proposal must include messages from $\mathcal{S}' \in \mathtt{shards}(\tau)$. Let $U' \subseteq \mathcal{S}'$ be the replicas in $\mathcal{S}'$ of whom messages $\mathtt{VCGlobalStateR}$ are included

in $W$. Let $V' = U' \setminus \mathcal{F}(\mathcal{S}')$. We have $|S'| \geq \mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}'}$ and $|V'| \geq \mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}'}$. Hence, by a standard quorum argument, we conclude $S' \cap V' \neq \emptyset'$. Let $\mathbf{Q}' \in (S' \cap V')$. As $\mathbf{Q}'$ reached the global commit phase for $M^*$, it will meet the conditions of Line 25 of Figure 6 and provide both $M^*$ and a global prepare certificate for $M^*$. Let $M^{\ddagger}$ be any other global preprepare certificate in $W$ accompanied by a global prepare certificate. Due to Line 24 of Figure 6, the global preprepare certificate for the newest view of $\mathcal{S}$ will be recovered. As $v^*$ is the newest view of $\mathcal{S}$, $M^{\ddagger}$ will only prevent recovery of $M^*$ if it is also a global preprepare certificate for view $v^*$ of $\mathcal{S}$. In this case, Lemma 5.2 guarantees that $M^{\ddagger} = M^*$. Hence, any replica $\mathbf{R}$ will recover $M^*$ upon receiving $\langle \texttt{NewViewGlobal} : V, W \rangle$.

Now assume that the induction hypothesis holds for all views $j$, $v^* < j \leq i$. We will prove that the induction hypothesis holds for view $i + 1$. Consider any valid new-view proposal $\langle \texttt{NewView} : V \rangle$ for view $i+1$ and let $M^{\ddagger}$ with $m(\mathcal{S}, \tau^{\ddagger})_{w,\rho} \in M^{\ddagger}$ be any global preprepare certificate that is recovered due to the new-view proposal $\langle \texttt{NewView} : V \rangle$. Hence, $M^{\ddagger}$ is recovered via either Line 17 of Figure 6 or Line 25 of Figure 6. In both cases, there must exist a global prepare certificate $P$ for $M^{\ddagger}$. As $\langle \texttt{NewView} : V \rangle$ is valid, we must have $w \leq i$. Hence, we can apply the second property of the induction hypothesis to conclude that $w \leq v^*$. If $w = v^*$, then we can use Lemma 5.2 to conclude that $M^{\ddagger} = M^*$. Hence, to complete the proof, we must show that $w = v^*$. First, the case in which $M^{\ddagger}$ is recovered via Line 17 of Figure 6. Due to the existence of a global commit certificate $C$ for $M^{\ddagger}$, $M^{\ddagger}$ satisfies the premise of this proposition. By assumption, $v^*$ is the first view for which the premise of this proposition holds. Hence, $w \geq v^*$, in which case we conclude $M^{\ddagger} = M^*$. Last, the case in which $M^{\ddagger}$ is recovered via Line 25 of Figure 6. In this case, $M^{\ddagger}$ is recovered via some message $\langle \texttt{NewViewGlobal} : V, W \rangle$. Analogous to the proof for the base case, $V$ will contain a message $\texttt{VCRecoveryRQ}$ from some replica $\mathbf{Q} \in S(\mathbf{R}')$. Due to Line 2 of Figure 5, $\mathbf{Q}$ will provide information on $M^*$. Consequently, a prepare certificate for $M^*$ will be obtained via global state recovery, and we also conclude $M^{\ddagger} = M^*$. □

Lemma 5.2 and Proposition 5.3 are technical properties that assures that no transaction that could-be-committed by any replica will ever get lost by the system. Next, we bootstrap these technical properties to prove that all good replicas can always recover such could-be-committed transactions.

**Proposition 5.4.** Let $\tau$ be a transaction and $m(\mathcal{S}, \tau)_{v,\rho}$ be a preprepare message. If, for all shards $\mathcal{S}^*$, $\mathbf{g}_{\mathcal{S}^*} > 2\mathbf{f}_{\mathcal{S}^*}$, and there exists a shard $\mathcal{S}' \in \texttt{shards}(\tau)$ such that $\mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}'}$ good replicas in $\mathcal{S}'$ reached the global commit phase for $M$ with $m(\mathcal{S}, \tau)_{v,\rho} \in M$, then every good replica in $\mathcal{S}$ will accept $M$ whenever communication becomes reliable.

*Proof.* Let $\mathbf{R} \in \mathcal{S}$ be a good replica that is unable to accept $M$. At some point, communication becomes reliable, after which $\mathbf{R}$ will eventually trigger Line 1 of Figure 4. We have the following cases:

1. If $\mathbf{R}$ meets the conditions of Line 4 of Figure 4, then $\mathbf{R}$ has a local commit certificate $C(\mathcal{S}'')$, $\mathcal{S}'' \in \texttt{shards}(\tau)$. This local commit certificate certifies that at least $\mathbf{g}_{\mathcal{S}''} - \mathbf{f}_{\mathcal{S}''}$ good replicas in $\mathcal{S}''$ finished the global prepare phase for $M$. Hence, the conditions for Proposition 5.3 are met for $M$ and, hence, any shard in $\texttt{shards}(\tau)$ will maintain or recover $M$. Replica $\mathbf{R}$ can use $C(\mathcal{S}'')$ to prove this situation to other replicas, forcing them to commit to $M$, and provide any commit messages $\mathbf{R}$ is missing (Line 13 of Figure 4).

2. If $\mathbf{R}$ does not meet the conditions of Line 4 of Figure 4, but some other good replica $\mathbf{Q} \in \mathcal{S}$ does, then $\mathbf{Q}$ can provide all missing information to $\mathbf{R}$ (Line 6 of Figure 5). Next, $\mathbf{R}$ uses this information (Line 7 of Figure 5), after which it meets the conditions of Line 4 of Figure 4.

3. Otherwise, if the above two cases do not hold, then all $\mathbf{g}_{\mathcal{S}}$ good replicas in $\mathcal{S}$ are unable to finish the commit phase. Hence, they perform a view-change. Due to Proposition 5.3, this view-change will succeed and put every replica in $\mathcal{S}$ into the commit phase for $M$. As all good

15

replicas in $\mathcal{S}$ are in the commit phase, each good replica in $\mathcal{S}$ will be able to make a local commit certificate $C(\mathcal{S})$ for $M$, after which they meet the conditions of Line 4 of Figure 4. $\square$

Finally, we use Proposition 5.4 to prove *cross-shard-consistency*.

**Theorem 5.5.** *Optimistic*-CERBERUS *maintains cross-shard-consistency.*

*Proof.* Assume a single good replica R $\in \mathcal{S}$ commits or aborts a transaction $\tau$. Hence, it accepted some global preprepare certificate $M$ with $m(\mathcal{S}, \tau)_{v,\rho} \in M$. Consequently, R has local commit certificates $C(\mathcal{S}')$ for $M$ of every $\mathcal{S}' \in \mathtt{shards}(\tau)$. Hence, at least $\mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}'}$ good replicas in $\mathcal{S}'$ reached the global commit phase for $M$, and we can apply Proposition 5.4 to conclude that any good replica R$'' \in \mathcal{S}''$, $\mathcal{S}'' \in \mathtt{shards}(\tau)$ will accept $M$. As R$''$ bases its commit or abort decision for $\tau$ on the same global prepare certificate $M$ as R, they will both make the same decision, completing the proof. $\square$

As already argued, it is straightforward to use the details of Theorem 4.3 to prove that OCER-BERUS provides *validity*, *shard-involvement*, and *shard-applicability*. Via Theorem 5.5, we proved *cross-shard-consistency*. We cannot prove *service* and *confirmation*, however. The reason for this is simple: even though OCERBERUS can detect and recover from accidental faulty behavior and accidental concurrent transactions, OCERBERUS is not designed to gracefully handle targeted attacks.

**Example 5.6.** *Recall the situation of Example 4.1. Next, we illustrate how* OCERBERUS *deals with these concurrent transactions. We again consider distinct transactions* $\langle\tau_1\rangle_{c_1}$ *and* $\langle\tau_2\rangle_{c_2}$ *with* $\mathtt{Inputs}(\tau_1) = \mathtt{Inputs}(\tau_2) = \{o_1, o_2\}$ *and with* $\mathtt{shard}(o_1) = \mathcal{S}_1$ *and* $\mathtt{shard}(o_2) = \mathcal{S}_2$. *We assume that* $\mathcal{S}_1$ *processes* $\tau_1$ *first and* $\mathcal{S}_2$ *processes* $\tau_2$ *first.*

*The primary* $\mathcal{P}(\mathcal{S}_1)$ *will propose* $\tau_1$ *by prepreparing* $m(\mathcal{S}_1, \tau_1)_{v_1, \rho_1}$. *In doing so,* $\mathcal{P}(\mathcal{S}_1)$ *sends* $m(\mathcal{S}_1, \tau_1)_{v_1, \rho_1}$ *to all replicas* $\mathcal{S}_1 \cup \mathcal{S}_2$. *Next, the replicas in* $\mathcal{S}_1$ *will wait for a message* $m(\mathcal{S}_2, \tau_1)_{v_2', \rho_2'}$ *from* $\mathcal{S}_2$. *At the same time,* $\mathcal{P}(\mathcal{S}_2)$ *already proposed* $\tau_2$ *at the same time by sending out* $m(\mathcal{S}_2, \tau_2)_{v_2, \rho_2}$, *and the replicas in* $\mathcal{S}_2$ *will wait for a message* $m(\mathcal{S}_1, \tau_2)_{v_1', \rho_1'}$. *Hence, the replicas in* $\mathcal{S}_1$ *will never receive* $m(\mathcal{S}_2, \tau_1)_{v_2', \rho_2'}$ *and the replicas in* $\mathcal{S}_2$ *will never receive* $m(\mathcal{S}_1, \tau_2)_{v_1', \rho_1'}$. *Consequently, no replica will finish global preprepare, the consensus round will fail for all replicas, and all good replicas will initiate a view-change. As no replica reached the global prepare phase, transactions* $\tau_1$ *and* $\tau_2$ *do not need to be recovered during the view-change. After the view-changes, both* $\mathcal{S}_1$ *and* $\mathcal{S}_2$ *can process other transactions (or retry* $\tau_1$ *or* $\tau_2$), *but if they both process* $\tau_1$ *and* $\tau_2$ *again, the system will again initiate a view-change.*

As said before, OCERBERUS is optimistic in the sense that it is optimized for the situation in which faulty behavior (including concurrent transactions) is rare. Still, in all cases, OCERBERUS maintains cross-shard consistency, however. Moreover, in the optimistic case in which shards have good primaries and no concurrent transactions exist, progress is guaranteed whenever communication is reliable:

**Proposition 5.7.** *If, for all shards* $\mathcal{S}^*$, $\mathbf{g}_{\mathcal{S}^*} > 2\mathbf{f}_{\mathcal{S}^*}$, *and Assumptions 2.1, 2.2, 2.3, and 2.4 hold, then Optimistic*-CERBERUS *satisfies Requirements R1–R6 in the optimistic case.*

If the optimistic assumption does not hold, then this can result in coordinated attempts to prevent OCERBERUS from making progress. At the core of such attacks is the ability for malicious clients and malicious primaries to corrupt the operations of shards coordinated by good primaries, as already shown in Example 5.1. To reduce the impact of targeted attacks, one can opt to make primary election non-deterministic, e.g., by using shard-specific distributed coins to elect new primaries in individual shards [11, 13].

As a final note, we remark that we have presented OCERBERUS with a per-round checkpoint and recovery method. In this simplified design, the recovery path only has to recover at-most a single round. Our approach can easily be generalized to a more typical multi-round checkpoint and recovery method, however. Furthermore, we believe that the way in which OCERBERUS extends PBFT can easily be generalized to other consensus protocols, e.g., HOTSTUFF.
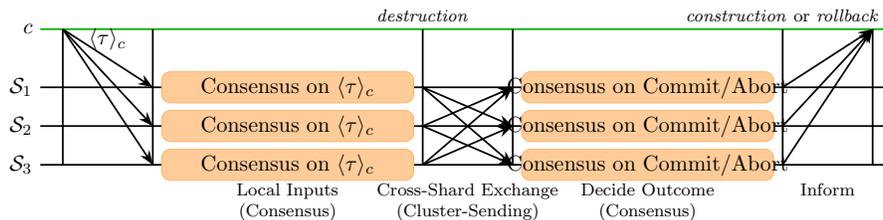
Figure 7: The message flow of PCERBERUS for a 3-shard client request $\langle\tau\rangle_c$ that is committed.

# 6   Pessimistic-Cerberus: transaction processing under attack

In the previous section, we introduced OCERBERUS, a general-purpose minimalistic and efficient multi-shard transaction processing protocol. OCERBERUS is designed with the assumption that malicious behavior is rare, due to which it can minimize coordination in the normal-case while requiring intricate coordination when recovering from attacks. As an alternative to the optimistic approach of OCERBERUS, we can apply a *pessimistic* approach to CCERBERUS to gracefully recover from concurrent transactions that is geared towards minimizing the influence of malicious behavior altogether. Next, we explore such a pessimistic design via *Pessimistic*-CERBERUS (PCERBERUS).

The design of PCERBERUS builds upon the design of CCERBERUS by adding additional coordination to the cross-shard exchange and decide outcome steps. As in CCERBERUS, the acceptance of $m(\mathcal{S}, \tau)_\rho$ in round $\rho$ by all good replicas completes the *local inputs* step. Before cross-shard exchange, the replicas in $\mathcal{S}$ destruct the objects in $D(\mathcal{S}, \tau)$, thereby fully pledging these objects to $\tau$ until the commit or abort decision. Then, $\mathcal{S}$ performs cross-shard exchange by broadcasting $m(\mathcal{S}, \tau)_\rho$ to all other shards in $\mathtt{shards}(\tau)$, while the replicas in $\mathcal{S}$ wait until they receive messages $m(\mathcal{S}', \tau)_{\rho'} = (\langle\tau\rangle_c, I(\mathcal{S}', \tau), D(\mathcal{S}', \tau))$ from all other shards $\mathcal{S}' \in \mathtt{shards}(\tau)$.

After cross-shard exchange comes the final *decide outcome* step. After $\mathcal{S}$ receives $m(\mathcal{S}', \tau)_{\rho'}$ from all shards $\mathcal{S}' \in \mathtt{shards}(\tau)$, the replicas force a *second consensus step* that determines the round $\rho^*$ at which $\mathcal{S}$ decides *commit* (whenever $I(\mathcal{S}', \tau) = D(\mathcal{S}', \tau)$ for all $\mathcal{S}' \in \mathtt{shards}(\tau)$) or *abort*. If $\mathcal{S}$ decides commit, then, in round $\rho^*$, all good replicas in $\mathcal{S}$ construct all objects $o \in \mathtt{Outputs}(\tau)$ with $\mathcal{S} = \mathtt{shard}(o)$. If $\mathcal{S}$ decides abort, then, in round $\rho^*$, all good replicas in $\mathcal{S}$ reconstruct all objects in $D(\mathcal{S}, \tau)$ (rollback). Finally, each good replica informs $c$ of the outcome of execution. If $c$ receives, from every shard $\mathcal{S}' \in \mathtt{shards}(\tau)$, identical outcomes from $\mathbf{g}_{\mathcal{S}'} - \mathbf{f}_{\mathcal{S}}$ distinct replicas in $\mathcal{S}'$, then it considers $\tau$ to be successfully executed. In Figure 7, we sketched the working of PCERBERUS.

We notice that processing a multi-shard transaction via PCERBERUS requires *two* consensus steps per shard. In some cases, we can eliminate the second step, however. First, if $\tau$ is a multi-shard transaction with $\mathcal{S} \in \mathtt{shards}(\tau)$ and the replicas in $\mathcal{S}$ accept $(\langle\tau\rangle_c, I(\mathcal{S}, \tau), D(\mathcal{S}, \tau))$ with $I(\mathcal{S}, \tau) \neq D(\mathcal{S}, \tau)$, then the replicas can immediately abort whenever they accept $(\langle\tau\rangle_c, I(\mathcal{S}, \tau), D(\mathcal{S}, \tau))$. Second, if $\tau$ is a single-shard transaction with $\mathtt{shards}(\tau) = \{\mathcal{S}\}$, then the replicas in $\mathcal{S}$ can immediately decide commit or abort whenever they accept $(\langle\tau\rangle_c, I(\mathcal{S}, \tau), D(\mathcal{S}, \tau))$. Hence, in both cases, processing of $\tau$ at $\mathcal{S}$ only requires a single consensus step at $\mathcal{S}$.

Next, we illustrate how PCERBERUS deals with concurrent transactions.

**Example 6.1.** *Recall the situation of Example 4.1. Next, we illustrate how* PCERBERUS *deals with these concurrent transactions. We again consider distinct transactions $\langle\tau_1\rangle_{c_1}$ and $\langle\tau_2\rangle_{c_2}$ with $\mathtt{Inputs}(\tau_1) = \mathtt{Inputs}(\tau_2) = \{o_1, o_2\}$ and with $\mathtt{shard}(o_1) = \mathcal{S}_1$ and $\mathtt{shard}(o_2) = \mathcal{S}_2$. We assume that $\mathcal{S}_1$ processes $\tau_1$ first and $\mathcal{S}_2$ processes $\tau_2$ first.*

*Shard $\mathcal{S}_1$ will start by destructing $o_1$ and sends $(\langle\tau_1\rangle_{c_1}, \{o_1\}, \{o_1\})$ to $\mathcal{S}_2$. Next, $\mathcal{S}_1$ will wait, during which it receives $\tau_2$. At the same time, $\mathcal{S}_2$ follows similar steps for $\tau_2$ and sends $(\langle\tau_2\rangle_{c_2}, \{o_2\}, \{o_2\})$ to $\mathcal{S}_1$. While $\mathcal{S}_1$ is waiting for information on $\tau_1$ from $\mathcal{S}_2$, it receives $\tau_2$ and starts processing of $\tau_2$.*

17

*Shard $\mathcal{S}_1$ directly determines that $o_1$ does no longer exist. Hence, it sends $(\langle\tau_2\rangle_{c_2},\{o_1\},\emptyset\})$ to $\mathcal{S}_2$. Likewise, $\mathcal{S}_2$ will start processing of $\tau_1$, sending $(\langle\tau_1\rangle_{c_1},\{o_2\},\emptyset)$ to $\mathcal{S}_1$ as a result.*

*After the above exchange, both $\mathcal{S}_1$ and $\mathcal{S}_2$ conclude that transactions $\tau_1$ and $\tau_2$ must be aborted, which they eventually both do, after which $o_1$ is restored in $\mathcal{S}_1$ and $o_2$ is restored in $\mathcal{S}_2$.*

We notice that this situation leads to *both* transactions being aborted. Furthermore, we see that even though transactions get aborted, individual replicas can all determine whether their shard performed the necessary steps and, hence, whether their primary operated correctly. Next, we prove the correctness of PCerberus:

**Theorem 6.2.** *If, for all shards $\mathcal{S}^*$, $\mathbf{g}_{\mathcal{S}^*} > 2\mathbf{f}_{\mathcal{S}^*}$, and Assumptions 2.1, 2.2, 2.3, and 2.4 hold, then Pessimistic-*Cerberus* satisfies Requirements R1–R6.*

*Proof.* Let $\tau$ be a transaction. As good replicas in $\mathcal{S}$ discard $\tau$ if it is invalid or if $\mathcal{S} \notin \mathtt{shards}(\tau)$, PCerberus provides *validity* and *shard-involvement*. Next, *shard-applicability* follow directly from the decide outcome step.

If a shard $\mathcal{S}$ commits or aborts transaction $\tau$, then it must have completed the decide outcome and cross-shard exchange steps. As $\mathcal{S}$ completed cross-shard exchange, all shards $\mathcal{S}' \in \mathtt{shards}(\tau)$ must have exchanged the necessary information to $\mathcal{S}$. By relying on cluster-sending for cross-shard exchange, $\mathcal{S}'$ requires cooperation of all good replicas in $\mathcal{S}'$ to exchange the necessary information to $\mathcal{S}$. Hence, we have the guarantee that these good replicas will also perform cross-shard exchange to any other shard $\mathcal{S}'' \in \mathtt{shards}(\tau)$. Hence, every shard $\mathcal{S}'' \in \mathtt{shards}(\tau)$ will receive the same information as $\mathcal{S}$, complete cross-shard exchange, and make the same decision during the decide outcome step, providing *cross-shard consistency.*

A client can force service on a transaction $\tau$ by choosing a shard $\mathcal{S} \in \mathtt{shards}(\tau)$ and sending $\tau$ to all good replicas in $\mathcal{G}(\mathcal{S})$. By doing so, the normal mechanisms of consensus can be used by the good replicas in $\mathcal{G}(\mathcal{S})$ to force acceptance on $\tau$ in $\mathcal{S}$ and, hence, bootstrapping acceptance on $\tau$ in all shards $\mathcal{S}' \in \mathcal{G}(\mathcal{S})$. Due to cross-shard consistency, every shard in $\mathtt{shards}(\tau)$ will perform the necessary steps to eventually inform the client. As all good replicas $R \in \mathcal{S}$, $\mathcal{S} \in \mathtt{shards}(\tau)$, will inform the client of the outcome for $\tau$, the majority of these inform-messages come from good replicas, enabling the client to reliably derive the true outcome. Hence, CCerberus provides *service* and *confirmation*. $\square$

# 7 The strengths of Cerberus

Having introduced the three variants of Cerberus in Sections 4, 5, and 6, we will now analyze the strengths and performance characteristics of each of the variants. First, we will show that Cerberus provides serializable execution [6, 9]. Second, we look at the ability of Cerberus to maximize per-shard throughput by supporting out-of-order processing. Finally, we compare the costs, the attainable performance, and the scalability of the three protocols.

## 7.1 The ordering of transactions in Cerberus

The data model utilized by CCerberus, OCerberus, and PCerberus guarantees that any object $o$ can only be involved in at-most *two* committed transactions: one that *constructs* $o$ and another one that *destructs* $o$. Assume the existence of such transactions $\tau_1$ and $\tau_2$ with $o \in \mathtt{Outputs}(\tau_1)$ and $o \in \mathtt{Inputs}(\tau_2)$. Due to *cross-shard-consistency* (Requirement R4), the shard $\mathtt{shard}(o)$ will have to execute both $\tau_1$ and $\tau_2$. Moreover, due to *shard-applicability* (Requirement R3), the shard $\mathtt{shard}(o)$ will execute $\tau_1$ strictly before $\tau_2$. Now consider the relation

$$\prec := \{(\tau,\tau') \mid (\text{the system committed to } \tau \text{ and } \tau') \wedge (\mathtt{Outputs}(\tau) \cap \mathtt{Inputs}(\tau') \neq \emptyset)\}.$$

Obviously, we have $\prec(\tau_1, \tau_2)$. Next, we will prove that all committed transactions are executed in a *serializable* ordering. As a first step, we prove the following:

**Lemma 7.1.** *If we interpret transactions as nodes and $\prec$ as an edge relation, then the resulting graph is* acyclic.

*Proof.* The proof is by contradiction. Let $G$ be the graph-interpretation of $\prec$. We assume that graph $G$ is cyclic. Hence, there exists transactions $\tau_0, \ldots, \tau_{m-1}$ such that $\prec(\tau_i, \tau_{i+1})$, $0 \le i < m-1$, and $\prec(\tau_{m-1}, \tau_0)$. By the definition of $\prec$, we can choose objects $o_i$, $0 \le i < m$, with $o_i \in (\mathtt{Outputs}(\tau_i) \cap \mathtt{Inputs}(\tau_{(i+1) \bmod m}))$. Due to *cross-shard-consistency* (Requirement R4), the shard $\mathtt{shard}(o_i)$, $0 \le i < m$, executed transactions $\tau_i$ and $\tau_{(i+1) \bmod m}$.

Consider $o_i$, $0 \le i < m$, and let $t_i$ be the time at which shard $\mathtt{shard}(o_i)$ executed $\tau_i$ and constructed $o_i$. Due to *shard-applicability* (Requirement R3), we know that shard $\mathtt{shard}(o_i)$ executed $\tau_{(i+1) \bmod m}$ strictly after $t_i$. Moreover, also shard $\mathtt{shard}(o_{(i+1) \bmod m})$ must have executed $\tau_{(i+1) \bmod m}$ strictly after $t_i$ and we derive $t_i < t_{(i+1) \bmod m}$. Hence, we must have $t_0 < t_1 < \cdots < t_{m-1} < t_0$, a contradiction. Consequently, $G$ must be acyclic. □

To derive a serializable execution order for all committed transactions, we simply construct a directed acyclic graph in which transactions are nodes and $\prec$ is the edge relation. Next, we *topologically sort* the graph to derive the searched-for ordering. Hence, we conclude:

**Theorem 7.2.** *A sharded fault-tolerant system that uses the object-dataset data model, processes UTXO-like transactions, and satisfies Requirements R1-R5 commits transactions in a serializable order.*

We notice that CERBERUS only provides serializability for *committed* transactions. As we have seen in Example 6.1, concurrent transactions are not executed in a serializable order, as they are aborted. It is this flexibility in dealing with aborted transactions that allows all variants of CERBERUS to operate with minimal and fully-decentralized coordination between shards; while still providing strong isolation for all committed transactions.

## 7.2   Out-of-order processing in Cerberus

In normal consensus-based systems, the *latency* for a single consensus decision is ultimately determined by the message delay $\delta$. E.g., with the three-phase design of PBFT, it will take at least $3\delta$ before a transaction that arrives at the primary is executed by all replicas. To minimize the influence of message delay on *throughput*, some consensus-based systems support *out-of-order decision making* in which the primary is allowed to maximize bandwidth usage by continuously proposing transactions for future rounds (while previous rounds are processed by the replicas). To illustrate this, one can look at fine-tuned implementations of PBFT running at replicas that have sufficient memory buffers available [16, 30]. In this setting, replicas can work on several consensus rounds at the same time by allowing the primary to propose for rounds within a window of rounds.

As the goal of CERBERUS is to maximize performance—both in terms of latency (OCERBERUS) and in terms of throughput—we have designed CERBERUS to support out-of-order processing (if provided by the underlying consensus protocol, in the case of CCERBERUS and PCERBERUS). The only limitation to these out-of-order processing capabilities are with respect to transactions affecting a shared object: such transactions must be proposed strictly in-order, as otherwise the set of pledged inputs cannot be correctly determined by the good replicas. This is not a limitation for the normal-case operations, however, as such concurrent transactions only happen due to malicious behavior.

## 7.3   A comparison of the three Cerberus variants

Finally, we compare the practical costs of the three CERBERUS multi-shard transaction processing protocols. First, in Figure 8, we provide a high-level comparison of the costs of each of the protocols

| Protocol name | Normal-case complexity Consensus | Exchange | Phases | Concurrent Transactions | View-changes |
|---|---|---|---|---|---|
| CCERBERUS | $s$ | 1 | 4 | Objects pledged | Single-shard |
| OCERBERUS | $s$ | 3 | 3 | View-change & Abort | Multi-shard |
| PCERBERUS | $2s$ | 1 | 7 | Normal-case Abort | Single-shard |

Figure 8: Comparison of the three CERBERUS protocols for processing a transaction that affects $s$ shards. We compare the normal-case complexity, how they deal with concurrent transactions (due to malicious clients), and how they deal with malicious primaries.

to process a single transaction $\tau$ that affects $s = |\mathtt{shards}(\tau)|$ distinct shards. For the normal-case behavior, we compare the complexity in the number of *consensus steps* per shard and the number of *cross-shard exchange* steps between shards (which together determine the maximum throughput), and the number of *consecutive communication phases* (which determines the minimum latency).

Next, we compare how the three protocols deal with malicious behavior by clients and by replicas. If no clients behave malicious, then all transactions will *commit*. In all three protocols, malicious behavior by clients can lead to the existence of concurrent transactions that affect the same object. Upon detection of such concurrent transactions, all three protocols will *abort*. The consequences of such an abort are different in the three protocols.

In CCERBERUS, objects affected by aborted transactions remain pledged and cannot be reused. In practice, this loss of objects can provide an incentive for clients to not behave malicious, but does limit the usability of CCERBERUS in non-incentivized environments. OCERBERUS is optimized with the assumption that conflicting concurrent transactions are rare. When conflicts occur, they can lead to the failure of a global consensus round, which can lead to a view-change in one or more affected shards (even if none of the primaries is faulty). Finally, PCERBERUS deals with concurrent transactions by aborting them via the normal-case of the protocol. To be able to do so, PCERBERUS does require additional consensus steps, however.

The three CERBERUS protocols are resilient against malicious replicas: only malicious primaries can affect the normal-case operations of these protocols. If malicious primaries behave sufficiently malicious to affect the normal-case operations, their behavior is detected, and the primary is replaced. In both CCERBERUS and PCERBERUS, dealing with a malicious primary in a shard can be done completely in isolation of all other shards. In OCERBERUS, which is optimized with the assumption that failures are rare, the failure of a primary while processing a transaction $\tau$ can lead to view-changes in all shards affected by $\tau$.

Finally, we illustrate the performance of CERBERUS. To do so, we have modeled the maximum throughput of each of these protocols in an environment where each shard has seven replicas (of which two can be faulty) and each replica has a bandwidth of $100 \, \mathrm{Mbit \, s^{-1}}$. We have chosen to optimize CCERBERUS, OCERBERUS, and PCERBERUS to minimize *processing latencies* over minimizing bandwidth usage (e.g., we do not batch requests and the cross-shard exchange steps do not utilize threshold signatures; with these techniques in place we can boost throughput by a constant factor at the cost of the per-transaction processing latency). In Figure 9, we have visualized the maximum attainable throughput for each of the protocols as function of the number of shards. In Figure 10, we have visualized the number of per-shard steps performed by the system (for CCERBERUS and OCERBERUS, this is equivalent to the number of per-shard consensus steps, for PCERBERUS this is half the number of per-shard consensus steps). As one can see from these figures, all three protocols have excellent scalability: increasing the number of shards will increase the overall throughput of the system. Sharding does come with clear overheads, however, increasing the number of shards also increases the number of shards affected by each transaction, thereby increasing the overall number of consensus steps. This is especially true for very large transactions that affect many objects (that can affect many distinct shards).
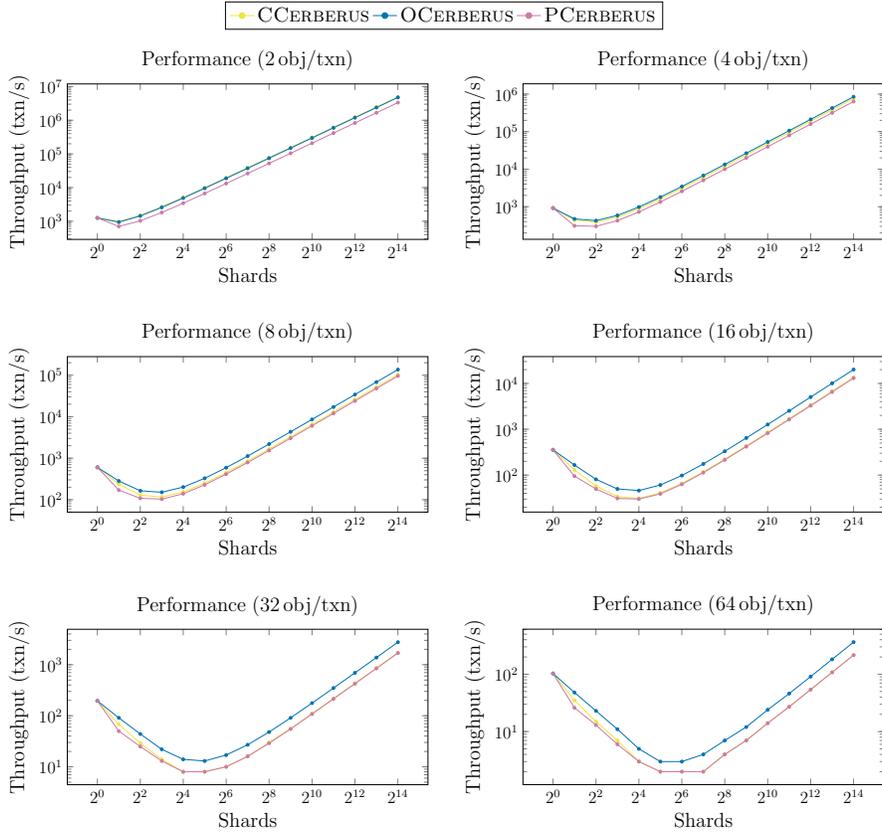
Figure 9: Throughput of the three CERBERUS protocols as a function of the number of shards.
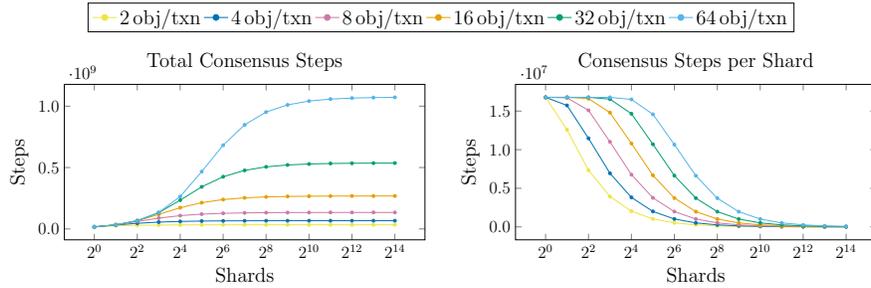


Figure 10: Amount of work, in terms of consensus steps, for the shards involved in processing the transactions.

# 8 Related Work

Distributed systems are typically employed to either increase reliability (e.g., via consensus-based fault-tolerance) or to increase performance (e.g., via sharding). Consequently, there is abundant literature on such distributed systems, distributed databases, and sharding (e.g., [46, 49, 50]) and on consensus-based fault-tolerant systems (e.g., [10, 14, 19, 31, 49]). Next, we shall focus on the few works that deal with sharding in fault-tolerant systems.

Several recent system papers have proposed specialized systems that combine sharding with consensus-based resilient systems. Examples include AHL [17], CAPER [3], CHAINSPACE [1], and

SHARPER [4], which all use sharding for data management and transaction processing. Systems such as AHL and CAPER are designed with single-shard workloads in mind, as they rely on centralized orderers to order and process multi-shard transactions, whereas systems such as CHAINSPACE and SHARPER are closer to the decentralized design of CERBERUS. In specific, CHAINSPACE uses a consensus-based commit protocol that performs three consecutive consensus and cross-shard exchange steps that resemble the two-step approach of PCERBERUS (although the details of the recovery path are rather different). In comparison, CERBERUS greatly improves on the design of CHAINSPACE by reducing the number of consecutive consensus steps necessary to process transactions and by introducing out-of-order transaction processing capabilities. Finally, SHARPER integrates global consensus steps in a consensus protocol in a similar manner as OCERBERUS. Their focus is mainly on a crash-tolerant PAXOS protocol, however, and they do not fully explorer the details of a full Byzantine fault-tolerant recovery path.

A few fully-replicate consensus-based systems utilize sharding at the level of consensus decision making, this to improve consensus throughput [2, 22, 26, 29]. In these systems, only a small subset of all replicas, those in a single shard, participate in the consensus on any given transaction, thereby reducing the costs to replicate this transaction without improving storage and processing scalability. Finally, the recently-proposed *delayed-replication algorithm* aims at improving scalability of resilient systems by separating fault-tolerant data storage from specialized data processing tasks [33], the latter of which can be distributed over many participants.

Recently, there has also been promising work on sharding and techniques supporting sharding for permissionless blockchains. Examples include techniques to enable sidechains, blockchain relays, and atomic swaps [23, 24, 34, 39, 52], which each enable various forms of cooperation between blockchains (including simple cross-chain communication and cross-chain transaction coordination). Unfortunately, these permissionless techniques are several orders of magnitudes slower than comparable techniques for traditional fault-tolerant systems, making them incomparable with the design of CERBERUS discussed in this work.

# 9    Conclusion

In this paper, we introduced Core-CERBERUS, Optimistic-CERBERUS, and Pessimistic-CERBERUS, three fully distributed approaches towards multi-shard fault-tolerant transaction processing. The design of these approaches is geared towards processing UTXO-like transactions in sharded distributed ledger networks with minimal cost, while maximizing performance. By using the properties of UTXO-like transactions to our advantage, both Core-CERBERUS and Optimistic-CERBERUS are optimized for cases with fewer expected malicious behaviors, in which case they are able to provide serializable transaction processing with only a single consensus step per affected shard, whereas Pessimistic-CERBERUS is optimized to efficiently deal with a broad-range of malicious behavior at the cost of a second consensus step during normal operations.

The core ideas of CERBERUS are not tied to any particular underlying consensus protocol. In this work, we have chosen to build CERBERUS on top of PBFT, as our experience shows that well-tuned implementations that use *out-of-order processing* of this protocol can outperform most other protocols in raw throughput [30]. Combining other consensus protocols with CERBERUS will result in other trade-offs between maximum throughput, per-transaction latency, bandwidth usage, and (for protocols that do not support out-of-order processing) vulnerability to message delays. Applying the ideas of CERBERUS fully onto other consensus protocols in a fully fine-tuned manner remains open, however. E.g., we are very interested in seeing whether incorporating CERBERUS into the more-resilient four-phase design of HOTSTUFF can sharply reduce the need for multi-shard view-changes in OCERBERUS (at the cost of higher per-transaction latency).

# References

[1] Mustafa Al-Bassam, Alberto Sonnino, Shehar Bano, Dave Hrycyszyn, and George Danezis. Chainspace: A sharded smart contracts platform, 2017. URL: `http://arxiv.org/abs/1708.03778`.

[2] Yair Amir, Claudiu Danilov, Danny Dolev, Jonathan Kirsch, John Lane, Cristina Nita-Rotaru, Josh Olsen, and David Zage. Steward: Scaling byzantine fault-tolerant replication to wide area networks. *IEEE Transactions on Dependable and Secure Computing*, 7(1):80–93, 2010. `doi:10.1109/TDSC.2008.53`.

[3] Mohammad Javad Amiri, Divyakant Agrawal, and Amr El Abbadi. CAPER: A cross-application permissioned blockchain. *Proc. VLDB Endow.*, 12(11):1385–1398, 2019. `doi:10.14778/3342263.3342275`.

[4] Mohammad Javad Amiri, Divyakant Agrawal, and Amr El Abbadi. SharPer: Sharding permissioned blockchains over network clusters, 2019. URL: `https://arxiv.org/abs/1910.00765v1`.

[5] Elli Androulaki, Artem Barger, Vita Bortnikov, Christian Cachin, Konstantinos Christidis, Angelo De Caro, David Enyeart, Christopher Ferris, Gennady Laventman, Yacov Manevich, Srinivasan Muralidharan, Chet Murthy, Binh Nguyen, Manish Sethi, Gari Singh, Keith Smith, Alessandro Sorniotti, Chrysoula Stathakopoulou, Marko Vukolić, Sharon Weed Cocco, and Jason Yellick. Hyperledger Fabric: A distributed operating system for permissioned blockchains. In *Proceedings of the Thirteenth EuroSys Conference*, pages 30:1–30:15. ACM, 2018. `doi:10.1145/3190508.3190538`.

[6] Vijayalakshmi Atluri, Elisa Bertino, and Sushil Jajodia. A theoretical formulation for degrees of isolation in databases. *Inform. Software Tech.*, 39(1):47–53, 1997. `doi:10.1016/0950-5849(96)01109-3`.

[7] Paddy Baker and Omkar Godbole. Ethereum fees soaring to 2-year high: Coin metrics. *CoinDesk*, 2020.

[8] Guillaume Bazot. Financial intermediation cost, rents, and productivity: An international comparison. Technical report, European Historical Economics Society, 2018.

[9] Hal Berenson, Phil Bernstein, Jim Gray, Jim Melton, Elizabeth ONeil, and Patrick ONeil. A critique of ANSI SQL isolation levels. *SIGMOD Rec.*, 24(2):1–10, 1995. `doi:10.1145/568271.223785`.

[10] Christian Berger and Hans P. Reiser. Scaling byzantine consensus: A broad analysis. In *Proceedings of the 2nd Workshop on Scalable and Resilient Infrastructures for Distributed Ledgers*, pages 13–18. ACM, 2018. `doi:10.1145/3284764.3284767`.

[11] Gabi Bracha and Ophir Rachman. Randomized consensus in expected $\mathcal{O}((n^2 \log n))$ operations. In *Distributed Algorithms*, pages 143–150. Springer Berlin Heidelberg, 1992. `doi:10.1007/BFb0022443`.

[12] Christopher Brookins. DeFi boom has saved bitcoin from plummeting. *Forbes*, 2020.

[13] Christian Cachin, Klaus Kursawe, Frank Petzold, and Victor Shoup. Secure and efficient asynchronous broadcast protocols. In *Advances in Cryptology — CRYPTO 2001*, pages 524–541. Springer Berlin Heidelberg, 2001. `doi:10.1007/3-540-44647-8_31`.

[14] Christian Cachin and Marko Vukolic. Blockchain consensus protocols in the wild (keynote talk). In *31st International Symposium on Distributed Computing*, volume 91, pages 1:1–1:16. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2017. `doi:10.4230/LIPIcs.DISC.2017.1`.

[15] Michael Casey, Jonah Crane, Gary Gensler, Simon Johnson, and Neha Narula. The impact of blockchain technology on finance: A catalyst for change. Technical report, International Center for Monetary and Banking Studies, 2018. URL: `https://www.cimb.ch/uploads/1/1/5/4/115414161/geneva21_1.pdf`.

[16] Miguel Castro and Barbara Liskov. Practical byzantine fault tolerance and proactive recovery. *ACM Trans. Comput. Syst.*, 20(4):398–461, 2002. `doi:10.1145/571637.571640`.

[17] Hung Dang, Tien Tuan Anh Dinh, Dumitrel Loghin, Ee-Chien Chang, Qian Lin, and Beng Chin Ooi. Towards scaling blockchain systems via sharding. In *Proceedings of the 2019 International Conference on Management of Data*, pages 123–140. ACM, 2019. `doi:10.1145/3299869.3319889`.

[18] Nikhilesh De. CFTC chair: 'a large part' of financial system could end up in blockchain format. *CoinDesk*, 2020.

[19] Tien Tuan Anh Dinh, Rui Liu, Meihui Zhang, Gang Chen, Beng Chin Ooi, and Ji Wang. Untangling blockchain: A data processing view of blockchain systems. *IEEE Trans. Knowl. Data Eng.*, 30(7):1366–1385, 2018. `doi:10.1109/TKDE.2017.2781227`.

[20] D. Dolev. Unanimity in an unknown and unreliable environment. In *22nd Annual Symposium on Foundations of Computer Science*, pages 159–168. IEEE, 1981. `doi:10.1109/SFCS.1981.53`.

[21] Danny Dolev. The byzantine generals strike again. *J. Algorithms*, 3(1):14–30, 1982. `doi:10.1016/0196-6774(82)90004-9`.

[22] Michael Eischer and Tobias Distler. Scalable byzantine fault-tolerant state-machine replication on heterogeneous servers. *Computing*, 101:97–118, 2019. `doi:10.1007/s00607-018-0652-3`.

[23] Muhammad El-Hindi, Carsten Binnig, Arvind Arasu, Donald Kossmann, and Ravi Ramamurthy. BlockchainDB: A shared database on blockchains. *Proc. VLDB Endow.*, 12(11):1597–1609, 2019. `doi:10.14778/3342263.3342636`.

[24] Ethereum Foundation. BTC Relay: A bridge between the bitcoin blockchain & ethereum smart contracts, 2017. URL: `http://btcrelay.org`.

[25] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. *J. ACM*, 32(2):374–382, 1985. `doi:10.1145/3149.214121`.

[26] Yossi Gilad, Rotem Hemo, Silvio Micali, Georgios Vlachos, and Nickolai Zeldovich. Algorand: Scaling byzantine agreements for cryptocurrencies. In *Proceedings of the 26th Symposium on Operating Systems Principles*, SOSP, pages 51–68. ACM, 2017. `doi:10.1145/3132747.3132757`.

[27] Seth Gilbert and Nancy Lynch. Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services. *SIGACT News*, 33(2):51–59, 2002. `doi:10.1145/564585.564601`.

[28] William J. Gordon and Christian Catalini. Blockchain technology for healthcare: Facilitating the transition to patient-driven interoperability. *Computational and Structural Biotechnology Journal*, 16:224–230, 2018. `doi:10.1016/j.csbj.2018.06.003`.

[29] Suyash Gupta, Sajjad Rahnama, Jelle Hellings, and Mohammad Sadoghi. ResilientDB: Global scale resilient blockchain fabric. *Proc. VLDB Endow.*, 13(6):868–883, 2020. `doi:10.14778/3380750.3380757`.

[30] Suyash Gupta, Sajjad Rahnama, and Mohammad Sadoghi. Permissioned blockchain through the looking glass: Architectural and implementation lessons learned. In *40th International Conference on Distributed Computing Systems*. IEEE, 2020.

[31] Suyash Gupta and Mohammad Sadoghi. *Blockchain Transaction Processing*, pages 1–11. Springer International Publishing, 2018. `doi:10.1007/978-3-319-63962-8_333-1`.

[32] Jelle Hellings and Mohammad Sadoghi. Brief announcement: The fault-tolerant cluster-sending problem. In *33rd International Symposium on Distributed Computing (DISC 2019)*, pages 45:1–45:3. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2019. `doi:10.4230/LIPIcs.DISC.2019.45`.

[33] Jelle Hellings and Mohammad Sadoghi. Coordination-free byzantine replication with minimal communication costs. In *23rd International Conference on Database Theory (ICDT 2020)*, volume 155, pages 17:1–17:20. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2020. `doi:10.4230/LIPIcs.ICDT.2020.17`.

[34] Maurice Herlihy. Atomic cross-chain swaps. In *Proceedings of the 2018 ACM Symposium on Principles of Distributed Computing*, pages 245–254. ACM, 2018. `doi:10.1145/3212734.3212736`.

[35] Maurice Herlihy. Blockchains from a distributed computing perspective. *Commun. ACM*, 62(2):78–85, 2019. `doi:10.1145/3209623`.

[36] Matt Higginson, Johannes-Tobias Lorenz, Bjrn Mnstermann, and Peter Braad Olesen. The promise of blockchain. Technical report, McKinsey&Company, 2017.

[37] Maged N. Kamel Boulos, James T. Wilson, and Kevin A. Clauson. Geospatial blockchain: promises, challenges, and scenarios in health and healthcare. *International Journal of Health Geographics*, 17(1):1211–1220, 2018. `doi:10.1186/s12942-018-0144-x`.

[38] Jonathan Katz and Yehuda Lindell. *Introduction to Modern Cryptography*. Chapman and Hall/CRC, 2nd edition, 2014.

[39] Jae Kwon and Ethan Buchman. Cosmos whitepaper: A network of distributed ledgers, 2019. URL: `https://cosmos.network/cosmos-whitepaper.pdf`.

[40] Leslie Lamport. Paxos made simple. *ACM SIGACT News*, 32(4):51–58, 2001. Distributed Computing Column 5. `doi:10.1145/568425.568433`.

[41] Laphou Lao, Zecheng Li, Songlin Hou, Bin Xiao, Songtao Guo, and Yuanyuan Yang. A survey of iot applications in blockchain systems: Architecture, consensus, and traffic modeling. *ACM Comput. Surv.*, 53(1), 2020. `doi:10.1145/3372136`.

[42] Satoshi Nakamoto. Bitcoin: A peer-to-peer electronic cash system, 2009. URL: `https://bitcoin.org/bitcoin.pdf`.

[43] Arvind Narayanan and Jeremy Clark. Bitcoin's academic pedigree. *Commun. ACM*, 60(12):36–45, 2017. `doi:10.1145/3132259`.

[44] Senthil Nathan, Chander Govindarajan, Adarsh Saraf, Manish Sethi, and Praveen Jayachandran. Blockchain meets database: Design and implementation of a blockchain relational database. *Proc. VLDB Endow.*, 12(11):1539–1552, 2019. `doi:10.14778/3342263.3342632`.

[45] Faisal Nawab and Mohammad Sadoghi. Blockplane: A global-scale byzantizing middleware. In *35th International Conference on Data Engineering (ICDE)*, pages 124–135. IEEE, 2019. `doi:10.1109/ICDE.2019.00020`.

[46] M. Tamer Özsu and Patrick Valduriez. *Principles of Distributed Database Systems*. Springer, 2020. `doi:10.1007/978-3-030-26253-2`.

[47] Michael Pisa and Matt Juden. Blockchain and economic development: Hype vs. reality. Technical report, Center for Global Development, 2017.

[48] Victor Shoup. Practical threshold signatures. In *Advances in Cryptology — EUROCRYPT 2000*, pages 207–220. Springer Berlin Heidelberg, 2000. `doi:10.1007/3-540-45539-6_15`.

[49] Gerard Tel. *Introduction to Distributed Algorithms*. Cambridge University Press, 2nd edition, 2001.

[50] Maarten van Steen and Andrew S. Tanenbaum. *Distributed Systems*. Maarten van Steen, 3th edition, 2017. URL: `https://www.distributed-systems.net/`.

[51] Gavin Wood. Ethereum: a secure decentralised generalised transaction ledger, 2016. EIP-150 revision. URL: `https://gavwood.com/paper.pdf`.

[52] Gavin Wood. Polkadot: vision for a heterogeneous multi-chain framework, 2016. URL: `https://polkadot.network/PolkaDotPaper.pdf`.

[53] Maofan Yin, Dahlia Malkhi, Michael K. Reiter, Guy Golan Gueta, and Ittai Abraham. HotStuff: BFT consensus with linearity and responsiveness. In *Proceedings of the ACM Symposium on Principles of Distributed Computing*, pages 347–356. ACM, 2019. `doi:10.1145/3293611.3331591`.