

HYDRA: Breaking the Global Ordering Barrier in Multi-BFT Consensus

Hanzheng Lyu*, Shaokang Xie[†], Jianyu Niu[‡]✉, Mohammad Sadoghi[†],
Yinqian Zhang[§], Cong Wang[‡], Ivan Beschastnikh[¶], Chen Feng*

University of British Columbia (*Okanagan Campus, [¶]Vancouver Campus), [†]University of California, Davis,

[‡]City University of Hong Kong, [§]Southern University of Science and Technology

*{lyuhanzheng@gmail.com, chen.feng@ubc.ca}, [†]{skxie,msadoghi}@ucdavis.edu

[‡]{njianyu@gmail.com, congwang@cityu.edu.hk}, [§]yinqianz@acm.org, [¶]bestchai@cs.ubc.ca

Abstract—Multi-Byzantine Fault Tolerant (Multi-BFT) consensus, which runs multiple BFT instances in parallel, has recently emerged as a promising approach to overcome the leader bottleneck in classical BFT protocols. However, existing designs rely on a global ordering layer to serialize blocks across instances, an intuitive yet costly mechanism that constrains scalability, amplifies failure propagation, and complicates deployment. In this paper, we challenge this conventional wisdom. We present HYDRA, the first Multi-BFT consensus framework that eliminates global ordering altogether. HYDRA introduces an object-centric execution model that partitions transactions by their accessed objects, enabling concurrent yet deterministic execution across instances. To ensure consistency, HYDRA combines lightweight lock-based coordination with a deadlock resolution mechanism, achieving both scalability and correctness. We implement HYDRA and evaluate it on up to 128 replicas in both LAN and WAN environments. Experimental results show HYDRA outperforms several state-of-the-art Multi-BFT protocols in the presence of a straggler. These results demonstrate strong consistency and high performance by removing global ordering, opening a new direction toward scalable Multi-BFT consensus design.

Index Terms—Multi-BFT consensus, Blockchain, Byzantine fault tolerance, Leader bottleneck, Partial ordering.

I. INTRODUCTION

Byzantine Fault Tolerant (BFT) consensus is a cornerstone of modern decentralized systems, powering diverse applications such as blockchains [1]–[4], decentralized finance (DeFi) [5], and decentralized storage [6]. Most classical BFT protocols, including PBFT [7] and Zyzzyva [8], adopt a leader-based architecture, where a designated replica (also called the leader) proposes transactions and coordinates with others to reach agreement. However, as system scale grows, this traditional leader-based design becomes a key performance bottleneck [9]–[14]: the leader’s coordination workload increases linearly with the number of replicas, making it the dominant factor limiting throughput and latency.

To address this scalability bottleneck, Multi-BFT consensus has emerged as a promising direction [10]–[13], [15]. Multi-BFT runs multiple leader-based consensus instances (*e.g.*, PBFT) in parallel, as illustrated in Fig. 1. Client transactions are partitioned into disjoint buckets, each handled by a separate instance, enabling concurrent agreement among multiple leaders. Each instance commits its own blocks locally,

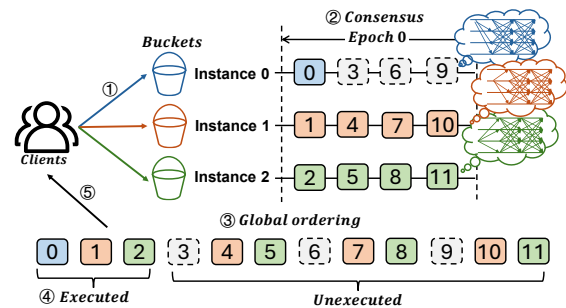


Figure 1: Multi-BFT consensus paradigm. Global ordering is the core of ordering blocks across instances.

which are then merged through a global ordering phase to form a single, totally ordered ledger. This design significantly increases throughput by better utilizing available bandwidth and computational capacity.

Despite its benefits, the global ordering layer—used to serialize blocks across instances—remains a major scalability barrier. While intuitively necessary for consistency, it introduces significant latency, particularly in the presence of slow instances. In practice, a single slow or faulty instance can stall the merging process, delaying confirmations for otherwise fast replicas. For example, in ISS [12], a state-of-the-art Multi-BFT protocol, the global ordering phase alone accounts for up to 92.8% of total transaction latency with just one straggler among 16 replicas. Although recent optimizations such as Ladon [16] and Orthrus [17] mitigate straggler effects, global ordering still dominates latency, contributing over 70% of total delay in Orthrus.

Global ordering also introduces execution inefficiency. Once blocks are serialized, replicas typically execute transactions sequentially according to the global log. Yet, many workloads contain semantically independent transactions that operate on disjoint objects. Hence, these computation workloads can be safely executed in parallel. This creates a double serialization barrier: 1) the ordering phase serializes instance outputs; and then 2) sequential execution re-serializes independent transactions. As a result, current Multi-BFT systems fail to exploit inherent transaction-level parallelism.

These limitations motivate a central question: *Can Multi-BFT consensus preserve safety and deterministic execution*

✉ Corresponding author: Jianyu Niu.

without enforcing a costly global order?

In this paper, we propose HYDRA, the first Multi-BFT consensus framework that eliminates the need for global ordering. Instead of enforcing a total transaction order and executing afterward, HYDRA unifies ordering and execution by leveraging transaction dependencies to expose concurrency while maintaining consistency. The key insight is that transactions operate on objects, which are independently updatable state items (*e.g.*, accounts or contract variables). As long as operations on the same object are executed sequentially, the correctness of the overall system state is preserved. Unlike sharding, where replicas maintain only a subset of objects and must coordinate across partitions, each replica in HYDRA retains a global view of all objects, enabling local coordination without expensive cross-instance protocols.

To realize this model, HYDRA partitions the global object space into disjoint groups, each managed by a dedicated instance responsible for ordering and committing transactions on its assigned objects. Transactions are classified into two types: intra-instance transactions, which affect objects within a single group, and cross-instance transactions, which span multiple groups. Replicas execute instances concurrently: intra-instance transactions are executed locally, whereas cross-instance transactions are carried out by coordinating their sub-operations across relevant instances. This object partitioning introduces a new atomicity challenge: sub-operations distributed across instances must either all commit or all abort for consistency.

HYDRA addresses this using a locking-based coordination mechanism: a transaction is prohibited from accessing objects involved in another in-flight transaction until the latter completes. Each transaction must acquire locks on all relevant objects before execution, guaranteeing that no intermediate state is exposed and that transaction outcomes remain atomic. However, locking introduces the potential for deadlocks. Different instances may acquire locks in different orders for the same cross-instance transaction set. Such inconsistent lock acquisition orders can create cyclic waiting dependencies, preventing involved transactions from making progress.

To resolve this, HYDRA implements a deadlock detection and resolution mechanism. When a transaction is blocked during lock acquisition, the system expands a distributed deadlock group by tracing dependency edges from all involved instances. If a dependency cycle is detected, HYDRA deterministically aborts one or more transactions in the group—ensuring that all instances resolve the conflict consistently and the system continues to make progress.

We build an end-to-end prototype of HYDRA in Go [18] and conduct extensive experiments on AWS to evaluate its performance. We compare HYDRA with the state-of-the-art Multi-BFT systems, including ISS [12], RCC [13], Mir-BFT [15], DQBFT [19], Ladon [16] and Orthrus [17].

Contributions. The main contributions are as follows.

- We propose HYDRA, the first Multi-BFT consensus framework that eliminates the need for global ordering. HYDRA significantly improves throughput and reduces end-to-end

latency, especially in the presence of slow instances.

- We design a concurrent execution model for Byzantine environments, allowing replicas to execute transactions locally while preserving deterministic consistency. This model leverages lock-based execution and distributed deadlock detection to enable high-throughput parallel processing without violating correctness.
- We build a prototype of HYDRA and evaluate its performance in both WAN and LAN environments. Experimental results show substantial throughput improvement, with up to 9.0× higher throughput in WAN deployments and 7.4× higher throughput in LAN settings, compared to existing Multi-BFT protocols in the presence of a straggler.

II. MOTIVATIONS

A. Anatomy of Multi-BFT Consensus

A Multi-BFT consensus system runs multiple BFT instances in parallel to improve throughput. Each transaction follows a typical five-stage process, as shown in Fig. 1: ① A client submits its transaction to replicas; ② Replicas reach instance-level agreement to commit transactions using a BFT protocol; ③ Replicas order committed transactions across instances into a global ledger; ④ Replicas execute the ordered transactions; ⑤ Replicas return results to the client.

The end-to-end latency and throughput are jointly determined by the following components: the *transmission phase* (① and ⑤), which is dominated by network delay; the *consensus phase* (②), which depends on the underlying BFT protocol; and the *global ordering phase* (③), which is determined by the global ordering algorithm; and finally, the *execution phase* (④), whose performance depends on the system’s execution model, *i.e.*, whether transactions are executed serially, speculatively, or in parallel based on their dependencies.

Among these components, the first two are well-studied and highly optimized in modern BFT systems, and such optimizations are largely orthogonal to the parallelization benefits brought by Multi-BFT designs. Thus, the major bottlenecks and opportunities for improvement lie in: (i) reducing or eliminating the global ordering cost, and (ii) enabling highly parallel execution. However, in practice, these two components fundamentally limit the scalability of Multi-BFT systems. We now examine these scalability challenges.

B. Analyzing Scalability Challenges

Global ordering limits scalability. While parallel instances can improve consensus throughput, the system must still serialize their results through a global ordering phase to ensure consistent execution. This stage becomes the dominant bottleneck as it couples the progress of all instances: a single delayed or crashed instance can block the entire system from advancing. Specifically, when one instance lags behind others, the system cannot finalize subsequent transactions from faster instances because their assigned global indices remain unfilled.

To understand the above limitation, Fig. 1 illustrates a case, in which instance 0 significantly lags behind the other

Table I: Latency breakdown (seconds) under one straggler among 16 replicas in a WAN environment.

Protocol	Transmission	Consensus	Ordering	Execution
ISS [12]	0.177	2.502	34.513	0.551
Orthrus [17]	0.176	2.553	7.430	0.558

instances: while the others have already produced four blocks, instance 0 has only produced one. This delay creates gaps at positions 3, 6, and 9 in the global log. As a result, the execution stalls after block 2, and subsequent blocks (e.g., blocks 4, 5, 7, 8, 10, and 11) cannot be executed until the missing blocks arrive. This dependency chain stalls the global log and causes the end-to-end latency to grow dramatically, even though the majority of instances continue to make local progress.

Table I shows experimental results to evaluate the impact of the global ordering phase, which exceeds 90% of the total end-to-end latency in ISS, which is a state-of-the-art pre-determined global ordering Multi-BFT consensus. Although recent works such as Ladon [16] and Orthrus [17] adopt dynamic or hybrid ordering to mitigate blocking, they still retain a global ordering barrier. Consequently, Orthrus continues to spend approximately 70% of its total latency in global ordering under the same conditions. Prior studies [12], [17] also obtain similar findings that under certain conditions, the global ordering stage can dominate overall latency.

Execution parallelism left untapped. Beyond ordering, existing Multi-BFT systems also underutilize execution parallelism. After transactions are globally ordered, replicas typically execute them sequentially to preserve determinism. This design overlooks the fact that many transactions are independent and can safely execute in parallel. Consequently, Multi-BFT consensus faces two layers of serialization: one at the ordering stage and another at the execution stage. This *double serialization barrier* fundamentally limits scalability, even when the consensus layer scales linearly with the number of instances.

We now simply analyze the theoretical execution time under sequential and parallel models. Let t denote the execution time of a single transaction. In the sequential model, execution time grows linearly with the number of transactions N , i.e., $T_{\text{seq}} = N \cdot t$. In contrast, parallel execution allows multiple transactions to be processed concurrently. With k parallel execution units, the theoretical execution time becomes $T_{\text{par}} = \lceil N/k \rceil \cdot t$. Thus, increasing the degree of parallelism effectively mitigates the execution bottleneck.

Summary. This work is motivated by the observation that ordering and execution are inherently connected: the way transactions are ordered determines how they can be executed. By rethinking this relationship, we can design a system that parallelizes both dimensions simultaneously by removing the global ordering in Multi-BFT consensus.

III. SYSTEM MODEL AND GOALS

A. System Model

We consider a system composed of $n = 3f + 1$ replicas, collectively denoted as the set \mathcal{N} , responsible for processing

transactions from a group of clients. We assume a subset of up to f replicas as *Byzantine*, represented as \mathcal{F} . Byzantine replicas may behave arbitrarily. The remaining replicas (in $\mathcal{N} \setminus \mathcal{F}$) are considered honest and strictly follow the protocol. We assume a single, computationally bounded adversary that controls all Byzantine replicas and cannot break cryptographic primitives to falsify messages from honest replicas (with a negligible probability). Each replica maintains a public/private key pair for signing and verifying messages.

Network model. We assume each pair of honest replicas is connected by an authenticated and reliable communication link. We adopt the partial synchrony model proposed by Dwork *et al.* [20], commonly used in BFT consensus [7], [21]. There is an established bound, denoted as Δ , and an undefined Global Stabilization Time (GST). After the GST point, the delivery of any message transmitted between two honest replicas within the Δ limit is guaranteed. That is, the system behaves *synchronously* after the GST.

B. Data Model

Objects. We adopt an object-centric design [22]–[28], where objects serve as the fundamental data units processed within the system. These objects are persistent, similar to accounts, and are formally represented as a tuple: $o = (\text{key}, \text{value})$, where *key* is a unique identifier, and *value* represents the current state of the object, which can be updated through transactions. For example, in Ethereum, each account can be treated as a distinct object, where the account address is the key and its balance is the value.

To ensure consistency and prevent conflicts in concurrent execution, objects can be explicitly locked and unlocked during transaction processing. When a transaction intends to modify an object, it must first acquire a lock on the object, preventing other transactions from modifying it simultaneously. Once the transaction is confirmed, the lock is released, allowing subsequent transactions to access the object.

Transactions. A transaction is structured as a Directed Acyclic Graph (DAG), formally defined as: $tx = (\text{id}, V, E)$, where *id* is a unique identifier for the transaction. The set V consists of vertices, where each vertex represents an operation on an object. Formally, a vertex v is expressed as: $v = (o, p, c, \sigma)$, where o is the object being modified, p denotes the operation applied to o , c represents the conditions required to execute p , and σ is the cryptographic signature of the owner of o . The set E contains directed edges, where an edge $e_{i,j} \in E$ signifies a dependency constraint, indicating that vertex v_i must be processed before vertex v_j can be executed.

Block. A block is defined as a tuple $b = (\text{txs}, \text{ins}, \text{sn}, \sigma)$, where *txs* denotes a batch of transactions, and *ins* specifies the instance processing the block. The sequence number *sn* is the index of a block within its instance. Finally, σ represents the cryptographic signature on b , guaranteeing both authenticity and integrity.

C. Preliminaries

Sequenced broadcast (SB). SB is a variant of Byzantine total order broadcast that provides a consistent total ordering of messages among replicas. In SB, a designated leader *broadcasts* messages that are associated with monotonically increasing sequence numbers, while all replicas collaborate to *deliver* these messages in a consistent global order. It employs a failure detector to identify and handle faulty or silent leaders. SB guarantees two properties:

- *Termination.* All honest replicas eventually deliver exactly one message for each sequence number.
- *Agreement.* The delivered messages are identical across all honest replicas

We adopt SB as a black-box abstraction for our consensus instances, where it accepts client transactions as input and produces a totally ordered stream of delivered transactions that are consistent among all honest replicas.

D. System Goals

We consider a Multi-BFT system composed of m BFT instances, indexed from 0 to $m-1$. Transactions are generated by clients and forwarded to replicas for processing, serving as the system’s input. Each instance operates in sequential rounds of a SB protocol, wherein a designated leader broadcasts a block of transactions and coordinates with backup replicas to deliver it. Transactions are executed after being partially ordered within an SB instance. A transaction is considered confirmed once it has been executed, regardless of whether the execution is successful or unsuccessful. The system state S is represented as a tuple, where each element corresponds to the maximum sequence number sn of an SB instance. Formally, the Multi-BFT system state is defined as:

$$S = (sn_0, sn_1, \dots, sn_{m-1}) \quad (1)$$

where sn_i denotes the maximum sequence number for the SB instance indexed by i . The system must ensure two fundamental properties:

- **Safety.** If two honest replicas reach the same state S , they must have identical values for every object in S .
- **Liveness.** If a transaction tx is received by at least one honest replica, then the client will eventually receive a response for tx .

IV. DESIGN RATIONALE

We present key insights of removing global ordering in Multi-BFT consensus, followed by challenges and corresponding solutions introduced by this new architectural design.

A. Key Insights

Traditional Multi-BFT systems enforce a global order on all transactions, although correctness only requires that all operations on the same object be processed in a consistent sequential order across replicas. Transactions accessing disjoint objects do not need global ordering and can safely execute in parallel.

This insight motivates an object-centric redesign of Multi-BFT architecture: by maintaining per-object ordering, HYDRA preserves deterministic execution while eliminating global ordering, thereby unlocking significantly higher concurrency.

A concrete example. To illustrate the idea, consider a transaction where user A transfers 10 tokens to user B . This transaction consists of two operations: $A.value-10$ and $B.value+10$, each targeting a distinct object. In conventional Multi-BFT systems, both operations are encapsulated in a single transaction that must be totally ordered with all others to ensure deterministic replay. In contrast, HYDRA allows these operations to be processed without global ordering. If both objects A and B reside in the same instance, the transaction follows that instance’s local order; if they are placed in different instances, each operation is ordered locally within its own instance.

B. Challenges and Solutions

Decomposing a transaction across multiple instances introduces two key challenges.

Challenge 1: Atomicity across instances. Suppose instance 0 executes $B.value+10$ while instance 1 later rejects $A.value-10$ due to insufficient balance. The resulting state is inconsistent— B gains tokens that were never deducted from A . A straightforward remedy would be to roll back $B.value+10$ when $A.value-10$ fails, but this is only safe if B ’s updated balance has not been used in subsequent transactions. Otherwise, reverting $B.value+10$ could invalidate dependent operations, leading to cascading inconsistencies.

To prevent such inconsistencies, HYDRA adopts a lock-based execution model inspired by classical concurrency control mechanisms such as Two-Phase Locking (2PL) [29]–[31]. Each transaction must acquire locks on all of its objects before execution. Only after obtaining all locks is the transaction executed. After that, all acquired locks in the transaction are released. This ensures that no transaction is ever partially executed: either all operations complete atomically or none take effect.

Challenge 2: Conflicting ordering among instances. Different instances may impose different orders on the same set of transactions, creating deadlocks. For example, consider two transactions: $tx_1 : A \rightarrow B$ and $tx_2 : B \rightarrow A$. If instance 0 locks A for tx_1 while instance 1 locks B for tx_2 , both transactions will wait indefinitely for the other’s lock. To resolve deadlocks, HYDRA integrates a deadlock detection and resolution mechanism [30]: once a transaction is confirmed but blocked on lock acquisition, the system searches for cycles in the wait-for graph and deterministically aborts one or more transactions to break the cycle.

These two mechanisms jointly ensure that all operations on the same object are guaranteed to be executed in a consistent and sequential order across replicas. It ensures that no transaction is ever partially executed—each transaction either completes all its intended operations atomically or is entirely

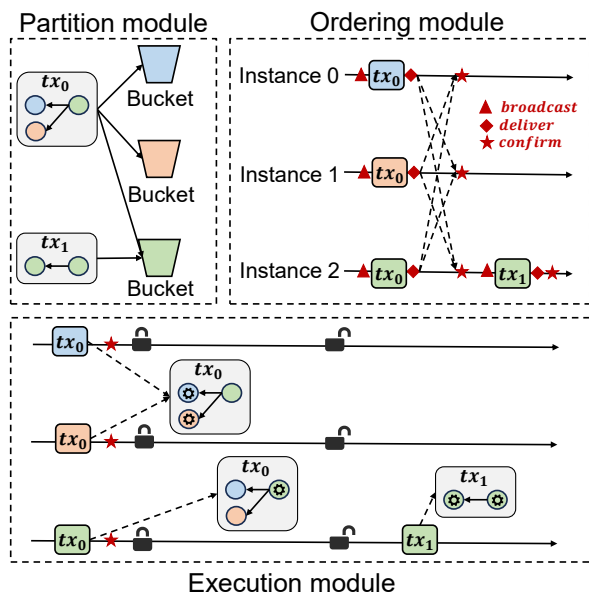


Figure 2: Overview of HYDRA. HYDRA contains three key steps: transaction partition, transaction ordering, and transaction execution.

aborted without leaving side effects. Importantly, this design allows the system to make progress with high concurrency.

V. SYSTEM DESIGN

We present HYDRA, which contains three key steps: transaction partition (Sec. V-B), transaction ordering (Sec. V-C), and transaction execution (Sec. V-D). Except for these transaction processes, HYDRA also performs checkpointing and garbage collection (Sec. V-E) to ensure efficient resource management and state consistency.

A. Architectural Overview

Fig. 2 presents an architectural overview of HYDRA, which operates multiple instances of the SB protocol in parallel. Transactions are assigned to buckets based on the objects they contain, and each bucket is mapped to a specific instance. HYDRA confirms transactions through the following phases:

- **Transaction partition:** Each transaction is assigned to the corresponding instances based on the objects it involves. If all objects within a transaction belong to the same instance, it is classified as an *intra-instance transaction*; otherwise, it is a *cross-instance transaction*.
- **Transaction ordering:** Replicas execute SB protocols within their respective instances to order transactions. A transaction is considered confirmed once all its corresponding instances have successfully delivered it.
- **Transaction execution:** Replicas execute transactions by respecting both the order (sequence numbers) and the dependencies (represented by the DAG).

A concrete example. Fig. 2 illustrates both cross-instance and intra-instance transactions. Consider a transaction tx_0 that involves three objects distributed across three different

instances. The objects in the first layer are assigned to instance 0 and instance 1, while the object in the second layer belongs to instance 2. The execution proceeds as follows: all three objects are proposed and reach consensus concurrently within their respective instances. Once all involved objects have been delivered, tx_0 is considered confirmed.

The replica then acquires locks on all related objects, ensuring atomicity. Execution begins with the first-layer objects in instance 0 and instance 1, followed by the second-layer object in instance 2, respecting the topological order of the transaction’s internal DAG. In contrast, transaction tx_1 is an intra-instance transaction that involves two objects, both assigned to instance 2. Since all related objects are processed within the same instance, tx_1 is confirmed immediately upon delivery in instance 2. Once confirmed, the transaction is executed directly according to its internal DAG, without requiring coordination across instances.

B. Transaction Partition

Algorithm 1 shows the detailed process to partition transactions. Upon receiving a transaction tx (Line 11), replica r assigns tx to the appropriate buckets based on the objects it involves. The process begins with verifying the transaction’s format and checking the owner’s signature for authenticity (Line 12). Once validated, replica r invokes the assign function to determine the corresponding bucket(s) for tx (Line 14). The design of the assign function follows an off-the-shelf strategy [15], [32], and we defer a detailed discussion of alternative assignment policies to Section VIII. Specifically, the replica obtains the appropriate bucket index for each object in the transaction. The transaction is then appended to the corresponding bucket (Line 15). Each bucket operates as an append-only list for backup replicas but supports both push and pull operations for the leader. To prevent duplication, if tx is already present in a bucket, it will not be added again.

C. Transaction Ordering

Each replica maintains a *log* for each SB instance, which facilitates the partial ordering of transactions within that instance. The log comprises multiple entries, each capable of storing a batch of transactions. The SB protocol defines two key events: $\langle \text{sb-broadcast} | b \rangle$ and $\langle \text{sb-deliver} | b \rangle$. The event $\langle \text{sb-broadcast} | b \rangle$ occurs when a block b is broadcasted within the SB, while $\langle \text{sb-deliver} | b \rangle$ denotes the successful ordering and delivery of block b .

Broadcast transactions. As shown in Algorithm 1, upon initializing HYDRA (Line 1), if replica r is the leader of $instance_i$, it enters a loop where it creates and broadcasts a block b for each sequence number sn in the current epoch (Lines 3-4). In each iteration, the leader pulls a specified number of the oldest transactions from the instance’s bucket (Line 5). If there are insufficient transactions to meet this threshold, the leader waits for a timeout. If the threshold is still not met after the timeout, it pulls all available transactions. It then assigns the instance index and sequence number to the block, signs it, and broadcasts it within the SB instance (Lines

Algorithm 1 HYDRA main algorithm on replica r

```
1: upon initialize system
2:   for  $i \in [0, m - 1]$  do  $\triangleright m$  is the number of instances
3:     if isLeader( $i, r$ ) then  $\triangleright r$  is instance $_i$ 's leader
4:       for  $sn \in \{0, 1, 2, \dots\}$  do
5:          $b.txs \leftarrow \text{pullValidTx}(\text{bucket}_i)$ 
6:          $b.ins \leftarrow i, b.sn \leftarrow sn, b.\sigma \leftarrow \text{sign}(b, r)$ 
7:         trigger  $\langle \text{sb-broadcast} | b \rangle$ 
8:       end for
9:     end if
10:  end for
11:   $\triangleright$  transaction partition
12: upon receive  $tx$ 
13:   if validateTx( $tx$ ) then  $\triangleright$  check format and signature
14:     for all  $v \in tx.V$  do
15:        $i \leftarrow \text{assign}(v.o)$   $\triangleright$  map object to an instance
16:        $\text{push}(tx, \text{bucket}_i)$ 
17:     end for
18:    $\triangleright$  transaction ordering
19: upon event  $\langle \text{sb-deliver} | b \rangle$ 
20:    $\text{log}[b.ins][b.sn] \leftarrow b$ 
21:   for all  $tx \in b$  do
22:     for all  $v \in tx.V$  where  $\text{assign}(v.o) = b.ins$  do
23:        $\text{delivered}(v) \leftarrow \text{true}$ 
24:     end for
25:     if  $\forall w \in tx.V : \text{delivered}(w)$  then
26:        $\text{confirm}(tx)$ 
27:     end if
28:   end for
```

6-7). All replicas participate in the ordering and delivery of the block. If r is not the leader, it acts as a backup, assisting in the ordering and delivery process within the SB instance.

Deliver transactions. As shown in Algorithm 1, upon delivering a block b from an SB instance (Line 18), the replica orders the block by appending it to the log of the instance indexed $b.ins$ at sequence number $b.sn$ (Line 19). Then, for each transaction tx contained in b , the replica iterates over all vertices v in tx . If the object $v.o$ is assigned to the current instance $b.ins$, the replica marks v as delivered (Lines 21-23). Once all vertices in tx are delivered, the transaction tx is confirmed (Lines 24-26).

Failure detector. In HYDRA, a failure detection module is integrated into the SB protocol to ensure liveness in the presence of faulty leaders. This mechanism enables replicas to replace a misbehaving or crashed leader for each SB instance, and has been widely adopted in prior BFT systems [21], [33]–[35]. For example, in PBFT, when replicas suspect the leader of Byzantine behavior, they initiate a view change to elect a new one. Similarly, in HYDRA, when the leader l_i of instance instance_i fails at sequence number sn , all honest replicas detect the failure, reach agreement on the state of instance_i and the new leader l'_i , and then restart the instance

Algorithm 2 Execute transactions on replica r

```
1: upon firstPending( $\text{log}[i]$ )  $\neq \perp$ 
2:    $\triangleright \text{log}[i]$  denotes the ordered transaction log of instance $_i$ 
3:    $tx \leftarrow \text{firstPending}(\text{log}[i])$ 
4:   for all  $v \in tx.V$  do
5:     if  $\text{assign}(v.o) = i$  then
6:        $\text{lock}(v.o)$ 
7:     end if
8:   end for
9:   if  $\forall v \in tx.V : \text{locked}(v.o)$  then
10:    for all  $v \in \text{toposort}(tx.V)$  do
11:       $\text{execute}(tx)$   $\triangleright$  execute vertices by dependency
12:    end for
13:    if  $\exists v \in tx.V : \text{failed}(v)$  then
14:       $\text{rollback}(tx)$   $\triangleright$  undo all executed vertices in  $tx$ 
15:       $\text{replyToClient}(tx, \text{FAILURE})$ 
16:    else
17:       $\text{replyToClient}(tx, \text{SUCCESS})$ 
18:    end if
19:    for all  $v \in tx.V$  do
20:       $\text{unlock}(v)$ 
21:    end for
22:  end if
```

from sequence number sn under the new leader.

D. Transaction Execution

HYDRA's execution model is inspired by well-established ideas in distributed databases. The requirement that a transaction acquires all necessary object locks before execution closely echoes 2PL [29]–[31], which ensures atomicity and isolation by preventing other transactions from accessing locked resources. Similarly, the deadlock resolution mechanism parallels classical deadlock detection techniques that rely on identifying cycles in wait-for graphs [30], [36].

As shown in Algorithm 2, when a replica identifies the first pending transaction tx in $\text{log}[i]$ (Lines 1-2), it iterates over all vertices $v \in tx.V$. For each vertex v , if its associated object $v.o$ is assigned to the current instance, the replica attempts to acquire a lock on $v.o$ (Lines 3-6). Once all required objects are successfully locked (Line 8), the replica executes the transaction according to the topological order of the DAG (Lines 9-11). If any vertex fails during execution, the replica rolls back all successfully executed vertices and replies to the client with a failure (Lines 12-14). Otherwise, if all vertices are executed successfully, the replica replies with a success (Line 16). Finally, the replica releases all locks held by tx (Lines 18-20), enabling subsequent transactions to access those objects.

Note that the algorithm does not explicitly handle the case where tx appears in only a subset of its target instances. In such cases, if tx remains unconfirmed at the end of an epoch, the replica will abort tx and release any locks it may have acquired, and a failure response is sent back to the client.

Deadlock resolution. Deadlocks may occur in this execution model, as transactions can appear in different orders across

Algorithm 3 Deadlock resolution on replica r

```
1: upon confirmed( $tx$ )  $\wedge$  isInterTx( $tx$ )
2:    $D \leftarrow$  ExpandDeadlockGroup( $tx$ )
3:    $victimList \leftarrow$  SelectVictims( $D$ )
4:   for all  $tx \in victimList$  do
5:     abort( $tx$ )
6:     reAdd2Buckets( $tx$ )
7:   end for
8:    $\triangleright$  search transactions forming a deadlock with  $tx$ 
9:   function ExpandDeadlockGroup( $tx$ )
10:     $D \leftarrow \{tx\}$ 
11:    repeat
12:       $D_{prev} \leftarrow D$ 
13:      for all  $t \in D_{prev}$  do
14:        for all  $t' \in$  findOrderingConflicts( $t$ ) do
15:           $D \leftarrow D \cup \{t'\}$ 
16:        end for
17:      end for
18:      until  $D = D_{prev}$ 
19:      return  $D$ 
20:    end function
21:     $\triangleright$  select victim transactions in  $D$  to eliminate deadlocks
22:    function SelectVictims( $D$ )
23:       $victims \leftarrow \emptyset$ 
24:      while hasDeadlock( $D$ ) do
25:         $victim \leftarrow$  transaction in  $D$  with smallest index
26:         $victims \leftarrow victims \cup \{victim\}$ 
27:         $D \leftarrow D \setminus \{victim\}$ 
28:      end while
29:      return  $victims$ 
30:    end function
31:     $\triangleright$  find transactions ordered inconsistently with  $tx$ 
32:    function findOrderingConflicts( $tx$ )
33:       $C \leftarrow \emptyset$ 
34:      for all  $(i, j) \in$  instances( $tx$ ),  $i \neq j$  do
35:        for all  $tx' \in$  prefix $[i](tx)$  do
36:          if  $tx' \in$  suffix $[j](tx)$  then
37:             $\triangleright tx'$  is ordered after  $tx$  in  $j$ 
38:             $C \leftarrow C \cup \{tx'\}$ 
39:          else if  $tx' \notin$  prefix $[j](tx) \wedge j \in$  instances( $tx'$ )
40:             $\triangleright tx'$  pending in  $j$ , will be ordered after  $tx$  in  $j$ 
41:             $C \leftarrow C \cup \{tx'\}$ 
42:          end if
43:        end for
44:      end for
45:      return  $C$ 
46:    end function
```

multiple instances, leading to conflicting execution dependencies. This is a well-known issue in concurrent systems, especially when each instance executes transactions independently based on local orderings. To address this, we adopt a deterministic deadlock resolution mechanism, as shown in Algorithm 3.

When a cross-instance transaction tx is confirmed (Line 1), *i.e.*, it has been delivered by all the instances it is assigned to, the replica initiates deadlock detection by invoking `expandDeadlockGroup` (Line 2). This procedure recursively constructs a set of transactions that may be involved in a cycle with tx due to inconsistent instance-level orderings (Lines 13-15). Specifically, the replica locates the position of tx in the log of the corresponding instance (Line 31), and scans backwards to identify earlier pending transactions (Line 32). If such a transaction tx' is found to appear after tx in another instance to which both are assigned (Line 33), or if tx' has not yet been ordered relative to tx in another instance (Line 35), the two are considered to have an inconsistent ordering. The replica then adds tx' to the conflict set C (Line 36). This ensures that both explicit and potential future inconsistencies are captured, even when tx' has not yet appeared in all relevant instance logs. Among these, it selects transactions that also access other objects in tx and appear after tx in those corresponding logs—indicating a conflicting relative order (Line 13). These transactions are added to the deadlock group (Line 14), and the process repeats until no new transactions are included (Line 17).

Once the deadlock group D is constructed, the replica calls `selectVictims` (Line 3) to determine which transactions should be aborted to break the cycle. This is done by iteratively removing the transaction with the smallest index (the index can be the hash of the transaction) from the group until the remaining transactions have no deadlock (Lines 22-26). The selected victims are then aborted and re-added to buckets again (Lines 4-7). This process ensures that deadlocks are resolved deterministically and consistently across replicas, without requiring additional consensus rounds.

E. Checkpoint and Garbage Collection

HYDRA employs a checkpoint protocol at the end of each epoch to facilitate state pruning and garbage collection. Once an epoch completes, each replica broadcasts a checkpoint message containing a signed digest summarizing the blocks it has processed during that epoch. Upon collecting at least $2f + 1$ matching checkpoint messages, a replica establishes a stable checkpoint. This checkpoint allows replicas to safely discard transactions, whether successfully executed or aborted, from the completed epoch to reduce storage overhead.

VI. CORRECTNESS AND ROBUSTNESS ANALYSIS

A. Safety and Liveness Analysis

We provide the proof sketch for safety and liveness (defined in Sec. III-D) here, while deferring the full proofs to Appendix A [37] due to space constraints.

Safety. Safety in HYDRA follows from two key invariants. First, at the same delivery frontier, all honest replicas observe identical delivered transactions from each instance. This is guaranteed by the Agreement property of the underlying SB protocol. Second, deadlock detection and abort decisions are deterministic functions of these delivered transactions. Given

the same set of delivered instance logs, all honest replicas compute the same deadlock group by recursively expanding conflicts. Abort decisions are then derived from this deadlock group using a deterministic victim selection rule. As a result, all honest replicas abort the same set of transactions and commit the same remaining ones.

Liveness. Liveness requires that every transaction submitted by a correct client eventually receives a response. In HYDRA, this is ensured by the Termination property of the SB protocol and the epoch-based execution model. Specifically, any transaction proposed by an honest leader is eventually delivered by the SB protocol. Transactions that cannot be committed due to conflicts or deadlocks are deterministically aborted and reinserted for execution in subsequent epochs. The epoch mechanism bounds lock holding time and prevents transactions from being blocked indefinitely.

B. Impact of Byzantine Behaviors

We analyze how HYDRA behaves under Byzantine behaviors that deliberately target its locking mechanisms and cross-instance coordination.

Lock manipulation by Byzantine leaders. A Byzantine leader could manipulate locking behavior from two aspects.

(i) *Delaying locked transactions.* A Byzantine leader may intentionally delay proposing a cross-instance transaction in one instance after it has already been proposed and has acquired locks in another instance. For example, for a transaction spanning instances ins_1 and ins_2 , a Byzantine leader in ins_2 may withhold or delay its proposal, causing locks acquired in ins_1 to be held for an extended period. This behavior increases blocking for other transactions accessing the same objects and degrades throughput.

(ii) *Creating deadlocks.* A Byzantine leader may also attempt to induce deadlocks by proposing conflicting orderings across instances. For instance, two cross-instance transactions tx_1 and tx_2 may be ordered as $tx_1 \rightarrow tx_2$ in ins_1 , while a Byzantine leader in ins_2 intentionally proposes them in the reverse order $tx_2 \rightarrow tx_1$, creating a deadlock.

Adversarial transaction injection by malicious clients. Malicious clients may deliberately generate multiple transactions involving the same set of objects and submit them to different instances in adversarial orders, intentionally creating deadlocks across instances and amplifying the likelihood of cross-instance deadlocks.

Impact analysis of Byzantine behaviors. The above adversarial behaviors target the liveness of the system rather than safety. Byzantine leaders or malicious clients may delay locked transactions or induce cross-instance deadlocks, thereby degrading throughput and increasing latency. However, the impact of such adversarial behaviors is bounded. First, deadlock resolution ensures that any induced deadlock is consistently identified by honest replicas. Once detected, the deadlock resolution protocol deterministically aborts a bounded set of transactions, preventing indefinite blocking. Second, the

epoch-based execution model enforces an upper bound on lock holding time: transactions that fail to make progress are aborted and retried in subsequent epochs, eliminating the possibility of persistent lock monopolization by Byzantine behavior. As a result, while adversaries can temporarily reduce system efficiency, they cannot cause unbounded delays, inconsistent execution, or divergence among honest replicas.

VII. PERFORMANCE EVALUATION

In this section, we evaluate the performance of HYDRA and compare it against other Multi-BFT protocols: Ladon [16], Orthrus [17], ISS [12], RCC [13], Mir [15], and DQBFT [19]. We implemented HYDRA¹ in Go [18], using the PBFT consensus protocol [33] as SB instances. Our evaluation aims to answer the following research questions:

- **Q1:** How does HYDRA perform compared to other Multi-BFT protocols? (Sec. VII-B)
- **Q2:** How does HYDRA perform in varying proportions of cross-instance transactions? (Sec. VII-C)
- **Q3:** How does HYDRA perform under faults? (Sec. VII-D)
- **Q4:** How does HYDRA’s memory usage scale with the number of instances? (Sec. VII-E)

A. Experimental Setup

We deploy our protocols on AWS EC2 instances (c5a.2xlarge) with one instance per node. Each instance is equipped with 8vCPUs, 16GB RAM, and runs Ubuntu Linux 22.04. Experiments are conducted in both LAN and WAN environments. In the LAN setting, machines communicate over a private network with 1 Gbps bandwidth. For the WAN setup, machines are distributed across four Amazon cloud regions (France, North America, Australia, and Japan) with both public and private network interfaces limited to 1Gbps. We use NTP for clock synchronization.

Each replica acts as the leader for one instance and as a backup for the others, *i.e.*, $m = n$. To maximize throughput, we use a large batch size of 4096 transactions, with each transaction carrying a 500-byte payload. We evaluate system performance under two network conditions: one with uniform performance across all instances, and another with a straggler scenario where one instance operates at one-tenth the speed of the others. Each experiment is repeated five times, and we report the median of the results.

We evaluate two performance metrics: (1) throughput, defined as the number of transactions successfully responded to clients per second, and (2) latency, measured as the average end-to-end delay from the moment clients submit transactions until they receive $f+1$ replicas’ replies. We report the peak throughput in kilo-transactions per second (ktps) before reaching saturation and the corresponding latency in seconds (s).

B. Performance Comparison

Fig. 3 and Fig. 4 compare the throughput and latency of HYDRA and other Multi-BFT protocols with one straggler

¹The source code is available at <https://github.com/ShaoKangXie/hydra.git>

and without stragglers. We evaluate the throughput and latency with 8–128 replicas in both LAN and WAN environments. We only report results with a single straggler, since performance is primarily limited by the slowest replica and adding more stragglers does not significantly change throughput or latency [16]. For clarity, the compared protocols are categorized into three groups: *pre-determined ordering schemes*, including ISS [12], Mir-BFT [15] and RCC [13], which enforce a fixed global transaction order; *dynamic ordering schemes*, including DQBFT [19] and Ladon [16], which adapt ordering based on instance progress to mitigate straggler effects; and *hybrid ordering scheme* represented by Orthrus [17], which partially relaxes global ordering by leveraging fast-path confirmation for certain transactions. Following prior studies [32], we set the ratio of cross-instance transactions to 12%. Each cross-instance transaction spans three instances and accesses three objects, reflecting a representative coordination pattern.

Performance in WAN. Fig. 3 compares HYDRA and other Multi-BFT protocols in WAN, with and without stragglers. With a straggler present, HYDRA significantly outperforms pre-determined ordering schemes: with 128 replicas, throughput improves by approximately 9.0 \times . This is because in pre-determined ordering schemes, a slow instance delays the entire system, amplifying straggler impact. Dynamic (DQBFT, Ladon) and hybrid (Orthrus) ordering schemes alleviate this coupling but still require a global ordering phase. For DQBFT, throughput decreases with more replicas because its single ordering instance becomes a bottleneck as the system scales. Consequently, HYDRA continues to achieve superior performance. With 8 replicas, it achieves 47% and 38% higher throughput than Ladon and Orthrus, respectively, and with 128 replicas, HYDRA surpasses DQBFT by 63%.

Without stragglers, HYDRA achieves comparable throughput to ISS, RCC, Ladon, and Orthrus, while consistently maintaining the lowest or near-lowest latency. This demonstrates that eliminating global ordering and incorporating deadlock resolution introduces negligible additional overhead, and scalability remains robust across different replica configurations.

Performance in LAN. Fig. 4 compares HYDRA and other Multi-BFT protocols in LAN, with and without stragglers. The overall trends are consistent with the WAN results: HYDRA maintains higher throughput and lower latency than the compared Multi-BFT protocols with one straggler. With 32 replicas, HYDRA achieves a 7.4 \times throughput improvement and a 70% latency reduction compared to the pre-determined ordering schemes. These designs remain bottlenecked by fixed global ordering. Besides, HYDRA continues to outperform dynamic and hybrid ordering schemes by eliminating global ordering. When no stragglers are present, HYDRA delivers throughput comparable to ISS, RCC, Ladon, and Orthrus, while sustaining consistently low latency.

Latency breakdown. To better understand the source of performance differences, we present a detailed latency breakdown of ISS, Orthrus and HYDRA under one straggler with 16 replicas in WAN. We select ISS and Orthrus because

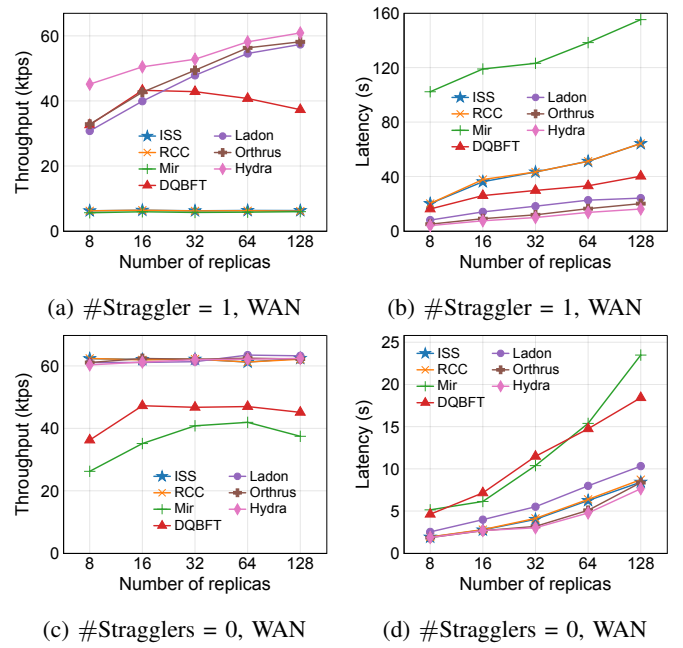


Figure 3: Throughput and latency of HYDRA, ISS, RCC, Mir, DQBFT, Ladon, and Orthrus in WAN.

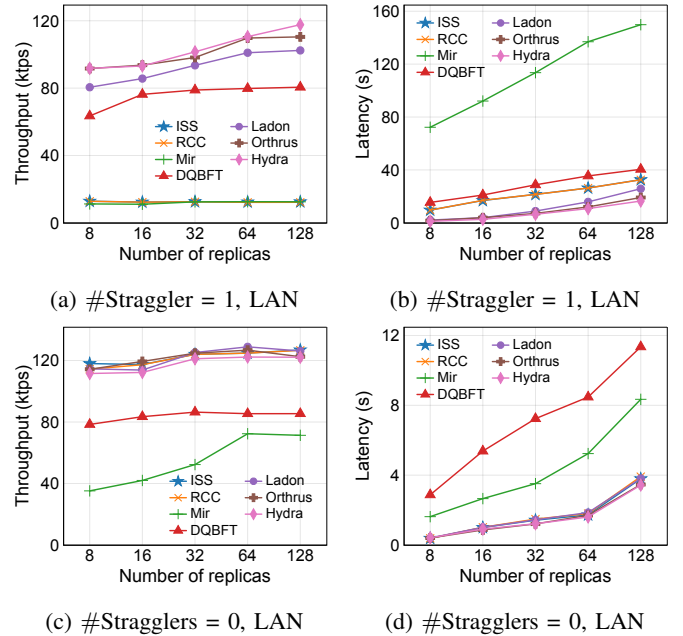


Figure 4: Throughput and latency of HYDRA, ISS, RCC, Mir, DQBFT, Ladon, and Orthrus in LAN.

ISS represents state-of-the-art Multi-BFT protocols with pre-determined global ordering, while Orthrus illustrates dynamic and hybrid ordering mechanisms that mitigate straggler effects. The total latency is divided into four stages: transmission, consensus, global ordering, and execution.

As shown in Fig. 5, the total latency of ISS is dominated by the global ordering phase, which accounts for 91.5% of its total latency. Although Orthrus reduces this cost by 78.5%,

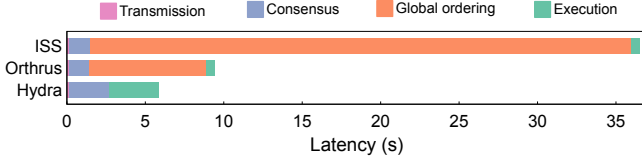


Figure 5: Breakdown of latency in ISS, Orthrus, and HYDRA.

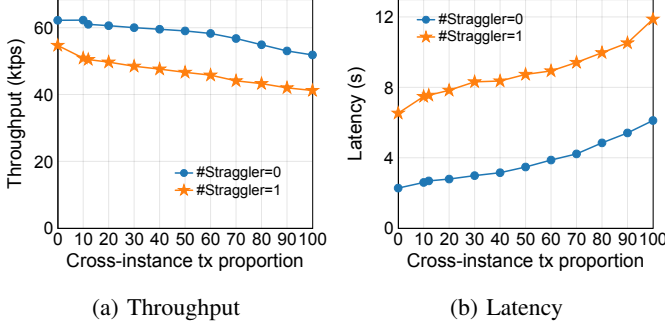


Figure 6: Throughput and latency of HYDRA under different cross-instance transactions proportions in WAN.

it still constitutes 69.3% of total latency. In contrast, HYDRA eliminates global ordering entirely, resulting in significantly lower overall latency. Specifically, HYDRA reduces end-to-end latency to 7.73s, outperforming ISS and Orthrus by 79.5% and 27.9%, respectively.

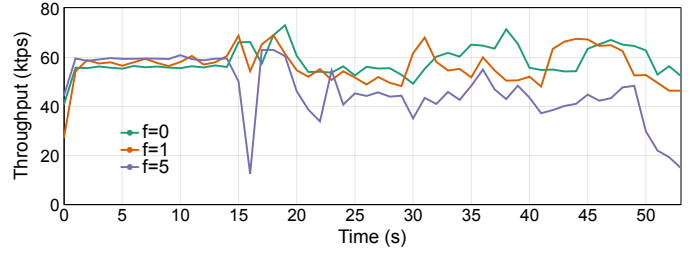
Among the remaining components, the consensus latency of HYDRA is slightly higher than in ISS and Orthrus. This is because we account for the time between receiving transactions and packing them into blocks as part of the consensus phase. Furthermore, since some transactions are assigned to multiple instances in HYDRA, the preprocessing overhead for these transactions is slightly larger. Notably, HYDRA incurs a slightly higher execution cost due to lock acquisition and deadlock resolution. However, this overhead is far outweighed by the savings from removing global ordering. Overall, these results highlight the core advantage of HYDRA’s design: by removing global ordering and enabling concurrent, lock-based execution, it substantially improves responsiveness without compromising consistency.

C. Impact of Cross-Instance Transactions

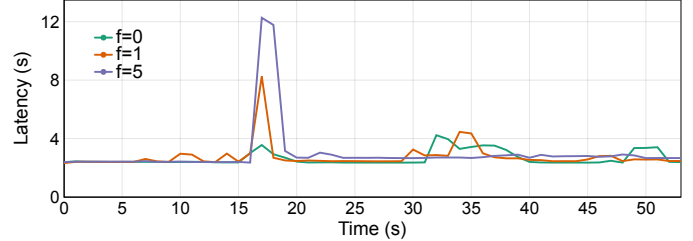
Fig. 6 shows the throughput and latency of HYDRA with different proportions of cross-instance transactions with 16 replicas in WAN. We only report WAN results here, as the LAN experiments exhibit the same overall trends.

As the proportion of cross-instance transactions increases, HYDRA experiences performance degradation, indicating that coordination across instances introduces additional lock contention and higher chances of deadlock resolution.

With one straggler, when the proportion grows from 0% to 100%, throughput decreases by 24.6% and latency increases by 81.7%. Without stragglers, the trend is similar: throughput decreases by 16% and latency increases by 1.7 \times . In both cases, the changes are relatively gradual, showing that the cost of



(a) Throughput average (over 1s intervals) over time.



(b) Latency average (over 1s intervals) over time.

Figure 7: Throughput and latency average of HYDRA over time with 0, 1, 5 faults. The faults occur at 15 seconds.

handling cross-instance dependencies grows moderately and does not cause severe performance degradation.

D. Performance under Faults

Crash faults. We evaluate the performance of HYDRA under crash faults in a WAN setting with 16 replicas, considering three scenarios with 0, 1, and 5 faulty replicas. Faults are injected at approximately 15s, and detected within a few seconds, triggering a view change after which performance stabilizes.

Fig. 7 shows the average throughput and latency over time. With one faulty replica, latency exhibits a short spike, peaking at 8.2s, and then stabilizes. In steady state (after view change), average latency increases modestly to 2.7s (about 8.7% higher than the pre-fault value), while throughput declines slightly by around 3-4%. This is because most instances continue processing normally, and only the instance led by the faulty replica temporarily stalls.

With five faulty replicas, the effect is more pronounced. Latency briefly spikes to 12s then settles to a higher steady-state level of 2.8s, roughly 15% higher than in the pre-fault state. Throughput also decreases by about 75% immediately after injection, and stabilizes at a steady-state throughput approximately 30% lower than normal. Even so, unaffected instances maintain stable progress and HYDRA fully recovers performance once the view change completes.

Overall, these results demonstrate that failures do not cascade across instances, and HYDRA remains resilient with bounded performance degradation even under multiple faults.

Byzantine faults. Table II summarizes system performance under different adversarial scenarios, including high-contention workloads generated by malicious clients (MC),

Table II: Throughput and latency under different adversarial scenarios. BL: Byzantine Leader, MC: Malicious Clients.

Metric / Scenario	Failure-free	MC	BL	BL + MC
Throughput (ktps)	61.50	57.01	60.13	52.86
Latency (s)	2.69	3.26	3.02	3.60

Table III: Memory usage and utilization under different numbers of instances. Each instance: 8 vCPUs and 16 GB RAM.

# Instances	8	16	32	64
Memory Usage (MB)	11046	10648	10549	10560
Utilization	67.4%	65.0%	64.4%	64.5%

Byzantine leader behavior (BL), and the combination of both (BL+MC). In the BL scenario, we inject a single Byzantine leader among 16 replicas, which deliberately delay locked transactions until near timeout, and proposes adversarial transaction orderings to induce cross-instance deadlocks. In the MC scenario, malicious clients generate high-contention workloads by repeatedly submitting transactions involving the same set of 20 objects to different instances in adversarial orders, intentionally inducing cross-instance deadlocks and increasing coordination overhead. Each scenario captures realistic adversarial strategies that may increase lock contention or coordination delays, as discussed in Section VI-B.

As the results show, the impact of adversarial behaviors is limited. Malicious clients submitting high-contention transactions reduce throughput by 7.3% and increase latency by 21.2%. Byzantine leaders attempting to manipulate locks or propose conflicting orders cause a smaller throughput reduction of 2.2% and a latency increase of 12.3%. When both adversarial behaviors are combined (BL+MC), throughput decreases by 14.0% and latency increases by 33.8%. These results indicate that HYDRA’s locking and deadlock handling mechanisms effectively contain the effects of adversarial manipulation, restricting the impact to performance degradation.

E. Memory Usage

We measure the memory usage of HYDRA while scaling the number of consensus instances from 8 to 64. Table III reports the memory usage and utilization per instance. We observe that memory utilization remains stable as the number of instances increases. This indicates that HYDRA exhibits near-constant memory scaling with respect to the number of instances, and that fully replicated state does not become a scalability bottleneck as the system scales horizontally.

This behavior is attributed to HYDRA’s epoch-based processing. At the end of each epoch, replicas perform checkpointing, truncate committed logs, and apply garbage collection to reclaim memory. As a result, the in-memory state maintained by each replica is bounded by the most recent checkpoint and does not grow unbounded over time, ensuring sustainable memory usage for long-running deployments.

VIII. DISCUSSIONS

Comparison with DAG-based protocols. Both Multi-BFT

and DAG-based protocols aim to improve throughput by exploiting parallelism. In Multi-BFT, each block in a round references a single block from the previous round, whereas DAG-based protocols typically reference at least $2f + 1$ blocks from the previous round. DAG-based systems require a global ordering to ensure consistency, which can become a performance bottleneck in the presence of slow replicas or stragglers. Moreover, DAGs may suffer from transaction duplication due to multiple block proposers each round.

Comparison with sharding protocols. Sharding-based protocols increase parallelism by partitioning both replicas and transactions across multiple shards. However, this design introduces additional coordination overhead for cross-shard transactions, as operations spanning multiple shards require multi-shard atomicity protocols. HYDRA’s core advantage lies in its communication-free cross-instance execution model, which fundamentally differentiates it from existing approaches. This point is further elaborated in the Appendix C of [37].

Distributed transaction. The key of distributed transaction processing is to ensure atomicity across partitioned data, typically through coordination protocols such as 2PC [38], [39]. While effective, these protocols incur at least two rounds of communication and may block if the coordinator fails or if some shards are slow to respond. In contrast, HYDRA operates in a fully replicated model where every replica maintains the complete object state. Cross-instance transactions are executed atomically through local locking, avoiding explicit inter-instance coordination and communication overhead while preserving deterministic consistency.

Transaction partition strategies. HYDRA does not assume a specific transaction partitioning algorithm, but instead is designed to operate correctly and efficiently under a wide range of existing partitioning strategies. For example, a simple and widely adopted approach assigns objects to instances by hashing their keys, *i.e.*, $insIndex = \text{Hash}(key) \bmod m$ [12], [39], [40]. A transaction is then routed to the instances responsible for the objects it accesses. This simple strategy achieves good load balance in expectation, but it may increase the frequency of cross-instance transactions when objects that frequently interact are mapped to different instances.

To mitigate this effect, HYDRA can naturally benefit from prior work such as TxAllo [32], which reduces the frequency of cross-instance transactions while maintaining a balanced load across instances. By co-locating frequently interacting objects when possible, such approaches lower cross-instance coordination overhead and deadlock probability, thereby improving overall throughput.

Importantly, these partitioning strategies are orthogonal to the core design of HYDRA. The system’s correctness guarantees hold under any deterministic partitioning scheme agreed upon by all replicas.

IX. RELATED WORK

Multi-BFT consensus. The foundation of Byzantine fault-tolerant consensus was established by Castro and Liskov in

PBFT [7], which inspired a long line of practical protocols such as Zyzzyva [8], Tendermint [41], and HotStuff [21]. These leader-based protocols streamline the process of reaching agreement but inevitably concentrate decision power on a single replica, resulting in a leader bottleneck that constrains throughput and scalability [42]–[46].

To alleviate this limitation, Multi-BFT consensus architectures were introduced, allowing replicas to run several consensus instances concurrently and thereby removing the dependence on a single leader [12], [13], [15], [16], [47]. Among early designs, Mir-BFT [15] pioneered this direction by executing multiple instances in parallel and establishing a global order of blocks based on pre-assigned indices. However, because Mir-BFT triggers an epoch change whenever any instance leader behaves incorrectly, the system is vulnerable to Byzantine disruptions. ISS [12] refined this by introducing no-op deliveries to mitigate unnecessary epoch changes, while RCC [13] similarly relies on a static global order of outputs. In both cases, a slow or faulty instance can delay progression for all instances, creating performance bottlenecks.

Other variants attempt to optimize the global ordering process. For example, DQBFT [19] dedicates a separate BFT instance to globally order outputs from other instances, simplifying coordination but making the system more prone to targeted attacks [48]. Ladon [16] introduces monotonic ranking to dynamically determine cross-instance ordering, mitigating straggler-induced delays. Building upon this, Orthrus [17] further introduces a fast path for independent transactions, enabling them to bypass the global ordering phase, while still maintaining global consistency for dependent ones. Although these techniques alleviate synchronization overheads, the global ordering phase still constitutes a large fraction of overall system latency.

Unlike the systems above, TELL [49] shifts optimization to the execution layer rather than to the consensus protocol. It employs a State Hash Table (SHT) to track read/write dependencies, enabling concurrent execution both across instances and within blocks. Conflicts are handled by selective re-execution and periodic merging of instance states at the epoch boundary. While this strategy reduces execution wait time, its improvement in end-to-end latency remains limited due to the overhead of reprocessing conflicting transactions.

Partial ordering design. A complementary research direction focuses on partially ordered execution to relax the global ordering requirement, particularly within payment-oriented systems. CryptoConcurrency [50] dynamically detects overspending by checking account balances, enabling concurrent transactions without full consensus. Astro [51] maintains per-client logs to prevent double spending, while ABC [52] allows validators to process transactions in parallel without global coordination, thereby enhancing efficiency. FastPay [53] leverages payment semantics to minimize shared state between accounts, facilitating asynchronous execution with high concurrency. Flash [54] bypasses reliable broadcast through a DAG-based, partially ordered structure. A non-sequential

model for monetary transfers is proposed in [55], based on reliable broadcast abstraction. Pastro [56] further defines a partially ordered transaction set that determines active participants and stake distribution, offering flexibility for applications that do not depend on a total order. While these approaches achieve efficient execution for payment systems, they cannot generally support complex smart contract semantics.

Hybrid ordering design. Recent advances in blockchain systems have explored hybrid ordering mechanisms that selectively apply global ordering only when necessary [28], [57]. Sui Lutris [28] introduces a dual-path model, where single-owner transactions follow a lightweight fast path, while shared-object transactions are ordered through a consensus path. This design reduces latency for independent operations but depends on client-side orchestration, which requires gathering additional certificate signatures and performing periodic checkpointing to ensure finality. Mysticeti [57] generalizes the same concept using a DAG-based consensus protocol, offering a tighter integration between the fast and consensus paths. However, Mysticeti still separates execution and finality: locally executed transactions may later be reverted if not sufficiently confirmed at epoch boundaries, complicating persistence across epochs. Orthrus [17] extends this hybrid ordering paradigm into a Multi-BFT framework while retaining a standard BFT core. It introduces a fast path where conflict-free transactions can be finalized immediately upon instance commitment, thus eliminating deferred confirmation and preventing rollback between epochs. Thunderbolt [58] introduces a hybrid ordering mechanism that deterministically orders cross-shard transactions while preserving partial concurrency for single-shard execution, achieving both high scalability and consistency across shards.

X. CONCLUSION

We presented HYDRA, a new Multi-BFT consensus architecture that eliminates global ordering while preserving safety and liveness. By adopting an object-centric execution model and enforcing only per-object ordering, HYDRA decouples the progress of parallel BFT instances and unlocks significantly higher concurrency. Our locking mechanism and deterministic deadlock handling further ensure that cross-instance execution remains correct without additional coordination rounds. We implemented HYDRA and evaluated it under both WAN and LAN environments. The results demonstrate that HYDRA sustains high throughput and low latency, outperforming state-of-the-art Multi-BFT protocols.

ACKNOWLEDGEMENT

This work is supported in part by NSFC under Grants 62302204; in part by the NSERC Discovery Grants RGPIN-2023-04962, RGPIN-2020-05203, RGPIN-2025-06826; in part by Fundamental and Interdisciplinary Disciplines Breakthrough Plan of the Ministry of Education of China under Grants JYB2025XDXM114; in part by Hong Kong Research Grants Council (RGC) under Grants R1012-21 and RFS2122-1S04; and in part by NSF Award Number 2245373.

XI. AI-GENERATED CONTENT ACKNOWLEDGMENT

Portions of this paper were assisted by the use of the GPT-5 model from OpenAI. Specifically, AI assistance was used for language refinement. All conceptual contributions, algorithmic designs, experimental methodologies, and scientific claims are solely made by the authors. The authors reviewed and verified the correctness of all AI-assisted content.

REFERENCES

- [1] Mohammad M. Jalalzai, Chen Feng, Costas Busch, Golden G. Richard, and Jianyu Niu. The Hermes BFT for blockchains. *TDSC*, 2022.
- [2] Marko Vukolić. The quest for scalable blockchain Fabric: Proof-of-Work vs. BFT Replication. 2016.
- [3] Dahlia Malkhi. Blockchain in the lens of BFT. USENIX Association, 2018.
- [4] Suyash Gupta, Jelle Hellings, Sajjad Rahnama, and Mohammad Sadoghi. Proof-of-execution: Reaching consensus through fault-tolerant speculation. 2021.
- [5] Xin Wang, Sisi Duan, James Clavin, and Haibin Zhang. Bft in blockchains: From protocols to use cases. *CSUR*, 2022.
- [6] Zhiqin Zhu, Guanqiu Qi, Mingyao Zheng, Jian Sun, and Yi Chai. Blockchain based consensus checking in decentralized cloud storage. *SMPT*, 2020.
- [7] Miguel Castro and Barbara Liskov. Practical Byzantine fault tolerance. In *OSDI*, 1999.
- [8] Ramakrishna Kotla, Lorenzo Alvisi, Mike Dahlin, Allen Clement, and Edmund Wong. Zyzzyva: speculative Byzantine fault tolerance. *SIGOPS*, 2007.
- [9] Fangyu Gai, Ali Farahbakhsh, Jianyu Niu, Chen Feng, Ivan Beschastnikh, and Hao Duan. Dissecting the performance of chained-BFT. In *ICDCS*, 2021.
- [10] Chrysoula Stathakopoulou, Tudor David, and Marko Vukolic. Mir-BFT: High-throughput BFT for blockchains. *arXiv preprint arXiv:1906.05552*, 2019.
- [11] Zeta Avarikioti, Lioba Heimbach, Roland Schmid, Laurent Vanbever, Roger Wattenhofer, and Patrick Wintermeyer. FnF-BFT: Exploring performance limits of BFT protocols. *arXiv preprint arXiv:2009.02235*, 2020.
- [12] Chrysoula Stathakopoulou, Matej Pavlovic, and Marko Vukolić. State machine replication scalability made simple. In *EuroSys*, 2022.
- [13] Suyash Gupta, Jelle Hellings, and Mohammad Sadoghi. RCC: Resilient concurrent consensus for high-throughput secure transaction processing. In *ICDE*, 2021.
- [14] Mohammad Javad Amiri, Chenyuan Wu, Divyakant Agrawal, Amr El Abbadi, Boon Thau Loo, and Mohammad Sadoghi. The bedrock of byzantine fault tolerance: A unified platform for {BFT} protocols analysis, implementation, and experimentation. In *NSDI*, 2024.
- [15] Chrysoula Stathakopoulou, Tudor David, and Marko Vukolic. Mir-bft: Scalable and robust BFT for decentralized networks. *JSys*, 2022.
- [16] Hanzheng Lyu, Shaokang Xie, Jianyu Niu, Chen Feng, Yinqian Zhang, and Ivan Beschastnikh. Ladon: High-Performance Multi-BFT Consensus via Dynamic Global Ordering. In *EuroSys*, 2025.
- [17] Hanzheng Lyu, Shaokang Xie, Jianyu Niu, Ivan Beschastnikh, Yinqian Zhang, Mohammad Sadoghi, and Chen Feng. Orthrus: Accelerating multi-bft consensus through concurrent partial ordering of transactions. In *ICDE*, 2025.
- [18] Golang. <https://go.dev/>.
- [19] Balaji Arun and Binoy Ravindran. Scalable byzantine fault tolerance via partial decentralization. 2022.
- [20] Cynthia Dwork, Nancy Lynch, and Larry Stockmeyer. Consensus in the presence of partial synchrony. *Journal of the ACM (JACM)*, 1988.
- [21] Maofan Yin, Dahlia Malkhi, Michael K. Reiter, Guy Golan Gueta, and Ittai Abraham. HotStuff: BFT consensus with linearity and responsiveness. *PODC*, 2019.
- [22] Andrew D Birrell and Bruce Jay Nelson. Implementing remote procedure calls. *ACM Transactions on Computer Systems (TOCS)*, 2(1):39–59, 1984.
- [23] CORPORATE Open Software Foundation. *Introduction to OSF DCE (rev. 1.0)*. Prentice-Hall, Inc., 1992.
- [24] Ann Wollrath, Roger Riggs, and Jim Waldo. A Distributed Object Model for the Java System. *Computing Systems*, 9:265–290, 1996.
- [25] Frank E Redmond. *Dcom: Microsoft Distributed Component Object Model with Cdrom*. IDG Books Worldwide, Inc., 1997.
- [26] Krzysztof Ostrowski, Ken Birman, Danny Dolev, and Jong Hoon Ahn. Programming with live distributed objects. In *European Conference on Object-Oriented Programming*, pages 463–489, 2008.
- [27] William E Weihl. Commutativity-based Concurrency Control for Abstract Data Types. *IEEE Trans Comput*, 1988.
- [28] Same Blackshear, Andrey Chursin, George Danezis, Anastasios Kichidis, Lefteris Kokoris-Kogias, Xun Li, Mark Logan, Ashok Menon, Todd Nowacki, Alberto Sonnino, et al. Sui Lutris: A Blockchain Combining Broadcast and Consensus. In *CCS*, 2024.
- [29] Kapali P. Eswaran, Jim N Gray, Raymond A. Lorie, and Irving L. Traiger. The notions of consistency and predicate locks in a database system. *Communications of the ACM*, 19(11):624–633, 1976.
- [30] Philip A Bernstein, Vassos Hadzilacos, Nathan Goodman, et al. *Concurrency control and recovery in database systems*, volume 370. Addison-wesley Reading, 1987.
- [31] Jim Gray and Andreas Reuter. *Transaction processing: concepts and techniques*. 1992.
- [32] Yuanzhe Zhang, Shirui Pan, and Jiangshan Yu. TxAllo: Dynamic Transaction Allocation in Sharded Blockchain Systems. In *IEEE ICDE*, 2023.
- [33] Miguel Castro and Barbara Liskov. Practical Byzantine fault tolerance and proactive recovery. *TOCS*, 2002.
- [34] Yehonatan Buchnik and Roy Friedman. Fireledger: A high throughput blockchain consensus protocol. *VLDB*, 2020.
- [35] Guy Golan Gueta, Ittai Abraham, Shelly Grossman, Dahlia Malkhi, Benny Pinkas, Michael Reiter, Dragos-Adrian Seredinschi, Orr Tamir, and Alin Tomescu. SBFT: a scalable and decentralized trust infrastructure. In *DSN*, 2019.
- [36] Nima Kaveh and Wolfgang Emmerich. Deadlock detection in distribution object systems. *SIGSOFT Softw. Eng. Notes*, pages 44–51, 2001.
- [37] Hanzheng Lyu, Shaokang Xie, Jianyu Niu, Mohammad Sadoghi, Yinqian Zhang, Cong Wang, Ivan Beschastnikh, and Chen Feng. Hydra: Breaking the global ordering barrier in multi-bft consensus. *arXiv preprint arXiv:2511.05843*, 2025.
- [38] Mustafa Al-Bassam, Alberto Sonnino, Shehar Bano, Dave Hryczyn, and George Danezis. Chainspace: A Sharded Smart Contracts Platform. In *NDSS*, 2018.
- [39] E. Kokoris-Kogias, P. Jovanovic, L. Gasser, N. Gailly, E. Syta, and B. Ford. Omniledger: A secure, scale-out, decentralized ledger via sharding. In *2018 IEEE Symposium on Security and Privacy (SP)*, pages 583–598, May 2018.
- [40] M Zamani, R Jurdak, Y Liu, B O’Flynn, and P Dutta. Rapidchain: A fast and scalable blockchain protocol. In *2018 IEEE International Conference on Pervasive Computing and Communications (PerCom)*, pages 1–10. IEEE, 2018.
- [41] Ethan Buchman. Tendermint: Byzantine fault tolerance in the age of blockchains. *M.Sc. Thesis, University of Guelph, Canada*, Jun 2016.
- [42] Salem Alqahtani and Murat Demirbas. Bottlenecks in blockchain consensus protocols. In *COINS*, 2021.
- [43] Aleksey Charapko, Ailidani Ailijiang, and Murat Demirbas. Pigpaxos: Devouring the communication bottlenecks in distributed consensus. In *SIGMOD*, 2021.
- [44] Fangyu Gai, Jianyu Niu, Ivan Beschastnikh, Chen Feng, and Sheng Wang. Scaling Blockchain Consensus via a Robust Shared Mempool. In *ICDE*, 2023.
- [45] Dakai Kang, Suyash Gupta, Dahlia Malkhi, and Mohammad Sadoghi. Hotstuff-1: Linear consensus with one-phase speculation. *Proceedings of the ACM on Management of Data*, 2025.
- [46] Suyash Gupta, Sajjad Rahnama, Shubham Pandey, Natacha Crooks, and Mohammad Sadoghi. Dissecting bft consensus: In trusted components we trust! In *EuroSys*, 2023.
- [47] Dakai Kang, Sajjad Rahnama, Jelle Hellings, and Mohammad Sadoghi. Spotless: Concurrent rotational consensus made practical through rapid view synchronization. In *ICDE*. IEEE, 2024.
- [48] Ernesto Estrada. Network robustness to targeted attacks. the interplay of expansibility and degree distribution. *EPJ B*, 2006.
- [49] Xing Tong, Zheming Ye, Zhao Zhang, Cheqing Jin, and Aoying Zhou. TELL: Efficient Transaction Execution Protocol Towards Leaderless Consensus. In *ICDE*, 2024.
- [50] Andrei Tonkikh, Pavel Ponomarev, Petr Kuznetsov, and Yvonne-Anne Pignolet. Cryptoconcurrency: (Almost) Consensusless Asset Transfer with Shared Accounts. In *ACM CCS*, pages 1556–1570, 2023.

- [51] Daniel Collins, Rachid Guerraoui, Jovan Komatovic, Petr Kuznetsov, Matteo Monti, Matej Pavlovic, Yvonne-Anne Pignolet, Dragos-Adrian Seredinschi, Andrei Tonkikh, and Athanasios Xygkis. Online Payments by Merely Broadcasting Messages. In *DSN*, 2020.
- [52] Jakub Sliwinski and Roger Wattenhofer. Asynchronous Proof-of-Stake. In *International Symposium on Stabilizing, Safety, and Security of Distributed Systems*, pages 194–208, 2021.
- [53] Mathieu Baudet, George Danezis, and Alberto Sonnino. Fastpay: High-Performance Byzantine Fault Tolerant Settlement. In *Proceedings of the 2nd ACM Conference on Advances in Financial Technologies*, pages 163–177, 2020.
- [54] Andrew Lewis-Pye, Oded Naor, and Ehud Shapiro. Flash: An Asynchronous Payment System with Good-Case Linear Communication Complexity. In *arXiv preprint arXiv:2305.03567*, 2023.
- [55] Alex Auvolat, Davide Frey, Michel Raynal, and François Taïani. Money Transfer Made Simple: a Specification, a Generic Algorithm, and its Proof. In *arXiv preprint arXiv:2006.12276*, 2020.
- [56] Petr Kuznetsov, Yvonne-Anne Pignolet, Pavel Ponomarev, and Andrei Tonkikh. Permissionless and Asynchronous Asset Transfer. *Distributed Computing*, 36(3):349–371, 2023.
- [57] Kushal Babel, Andrey Chursin, George Danezis, Anastasios Kichidis, Lefteris Kokoris-Kogias, Arun Koshy, Alberto Sonnino, and Mingwei Tian. MYSTICETI: Reaching the Latency Limits with Uncertified DAGs. In *NDSS*, 2025.
- [58] Junchao Chen, Alberto Sonnino, Lefteris Kokoris-Kogias, and Mohammad Sadoghi. Thunderbolt: Concurrent smart contract execution with nonblocking reconfiguration for sharded DAGs. In *EDBT*, 2026.