

Lion: Minimizing Distributed Transactions through Adaptive Replica Provision (Extended Version)

Qiushi Zheng[†], Zhanhao Zhao[†], Wei Lu[†], Chang Yao[‡], Yuxing Chen[§], Anqun Pan[§], Xiaoyong Du[†]

[†] Renmin University of China [‡] Zhejiang University [§] Tencent Inc.

[†]{zhengqiushi, zhanhaozhao, lu-wei, duyong}@ruc.edu.cn

[‡]changy@zju.edu.cn [§]{axingguchen, aaronpan}@tencent.com

Abstract—Distributed transaction processing often involves multiple rounds of cross-node communications, and therefore, tends to be slow. To improve performance, existing approaches convert distributed transactions into single-node transactions by either migrating co-accessed partitions onto the same nodes or establishing a super node housing replicas of the entire database. However, migration-based methods might cause transactions to be blocked due to waiting for data migration, while the super node can become a bottleneck.

In this paper, we present **Lion**, a novel transaction processing protocol that utilizes partition-based replication to reduce the occurrence of distributed transactions. Inspired by the fact that modern distributed databases horizontally partition data, with each partition having multiple replicas, **Lion** aims to assign a node with one replica from each partition involved in a given transaction’s read or write operations. To ensure such a node is available, we propose an adaptive replica provision mechanism, enhanced with an LSTM-based workload prediction algorithm, to determine the appropriate node for locating replicas of co-accessed partitions. The adaptation of replica placement is conducted preemptively and asynchronously, thereby minimizing its impact on performance. By employing this adaptive replica placement strategy, we ensure that the majority of transactions can be efficiently processed on a single node without additional overhead. Only a small fraction of transactions will need to be treated as regular distributed transactions when such a node is unavailable. Consequently, **Lion** effectively minimizes distributed transactions, while avoiding any disruption caused by data migration or the creation of a super node. We conduct extensive experiments to compare **Lion** against various transaction processing protocols. The results show that **Lion** achieves up to 2.7x higher throughput and 76.4% better scalability against these state-of-the-art approaches.

Index Terms—Transaction Processing, Replica Provision

I. INTRODUCTION

Modern distributed databases, with representative examples including Spanner [1], CockroachDB [2], TiDB [3], are essential to today’s large-scale online applications. These databases horizontally partition data across several nodes to improve scalability. However, data partition comes at the cost of requiring distributed transaction processing, where a transaction must be executed and committed through rounds of communication with multiple nodes. As the involved network overhead is non-negligible, it is commonly accepted that distributed transactions are slow [4].

To boost performance, it is prevalent to execute transactions on a single node as many as possible to avoid distributed transaction processing. Thus far, existing approaches, either

based on data migration [5]–[9] or full-replication [10], [11], are sub-optimal. Data migration improves the partition quality by transferring the required data for a transaction to a specific node. This synchronous migration may lead to substantial performance degradation, as transactions relying on those partitions will get blocked until the migration completes. Worse still, when distinct transactions on different nodes require the same data, a “ping-pong” [12] problem can emerge, with the data continuously migrating between nodes back and forth to satisfy transaction demands. Full-replication-based methods, on the other hand, establish an additional node containing replicas of the entire database. Instead of involving data migration, distributed transactions can be converted into single-node transactions by executing on this “super node”. However, this node can potentially become a performance bottleneck. More importantly, due to the large data volume, establishing such a node may be infeasible. Therefore, designing a transaction processing protocol to efficiently minimize distributed transactions still remains an open problem.

In real-world distributed databases [1]–[3], it is fundamental that each partition has multiple replicas for high availability. Given the fact that a node can host several replicas of different partitions, it is possible that one node houses all the replicas needed for a given transaction’s read/write operations. Therefore, we consider reducing distributed transactions by executing them as single-node transactions on such nodes. However, using this idea to improve transaction performance requires addressing two major challenges: First, it is not trivial to ensure a node meeting the transaction’s specific requirements always exists. A sophisticated replica placement mechanism is needed to strategically locate replicas on appropriate nodes, thus avoiding the creation of a super node. Second, optimizing replica placement without causing transactions to be delayed or blocked is not straightforward. We must prepare a node hosting all the necessary replicas in advance of the transaction’s execution to avoid data migration overheads.

In this paper, we present **Lion**, an efficient transaction processing protocol that employs partition-based replication to minimize distributed transactions. We propose an adaptive replica provision mechanism to optimize the replica placement, ensuring that transactions can find a node with all the required replicas. Specifically, we first employ a graph-based workload analysis algorithm to identify partitions that

are frequently accessed together and then adjust the replica placement to ensure the co-accessed partitions are allocated onto an appropriate node. The replica placement adaptation is guided by a tailor-designed cost model for minimized distributed transactions and load balancing. Further, we introduce a workload prediction algorithm based on Long Short Term Memory (LSTM) to facilitate non-intrusive replica adjustment. By analyzing these predicated workloads, we anticipate partitions that are likely to be co-accessed in the future and pre-allocate their replicas to the appropriate nodes in advance. Therefore, we guarantee the replica adjustment is performed asynchronously with transaction processing. During the adjustment, we add a secondary replica to the corresponding node in the background, without interrupting the execution of transactions on the primary replica. We utilize the remastering technique [13] to prompt the secondary replica to the primary replica when necessary. Because replication protocols [14], [15] are adopted to guarantee the majority (or even all) of the replicas for a partition are consistent with each other, this remastering process is generally lightweight and free of data migration overheads.

Based on the replica placement, we execute transactions on a single node as many as possible. Due to the constraint that each partition typically has one primary replica and several secondary replicas, but only the primary replica handles write requests, we process transactions efficiently and correctly as follows. 1) If a transaction finds a node housing primary replicas of all relevant data, it can directly execute on that node as a single-node transaction. 2) In situations where the node lacks the necessary primary replicas but contains secondary ones, the transaction can still execute on that node after remastering these replicas to be primary. 3) Otherwise, transactions that fail to find such a node are treated as regular distributed transactions. We utilize the proposed cost model to ensure most transactions can be directly executed without remastering and distributed processing. Further, *Lion* supports two kinds of transaction processing schemes, namely standard (ad-hoc) execution and batch execution, ensuring its general applicability to real-world distributed databases.

In summary, we make the following contributions.

- We introduce *Lion*, a new transaction processing protocol that minimizes distributed transactions. *Lion* is generally applicable to modern distributed databases that leverage partition-based replication.
- We propose an adaptive replica provision mechanism that uses a graph-based workload analysis algorithm to determine a proper replica placement and a cost-model-based strategy to minimize distributed transactions as well as achieve load balance.
- We design an LSTM-based workload prediction algorithm to ensure the replica placement adjustment is non-intrusive while accommodating workload changes with sustainable performance.
- We conduct extensive evaluations on two popular benchmarks, namely YCSB and TPC-C, and compare *Lion* against various existing works for optimizing distributed

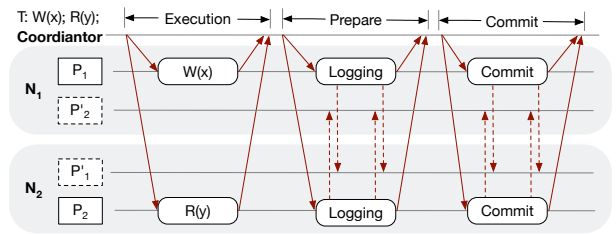


Fig. 1: Standard distributed transaction processing

transactions. The results show that *Lion* is efficient, and outperforms the state-of-the-art approaches by up to 2.7x.

The remainder of the paper is structured as follows. The next section provides relevant background and presents the problem statement. Section III overviews the architecture of *Lion*. Section IV details the adaptive replica provision mechanism, and elaborates on the prediction-based optimization. Section V describes the system implementation, and Section VI presents the experimental results. Section VII discusses the related work, and Section VIII concludes.

II. PRELIMINARIES

In this section, we first describe distributed transaction processing and existing approaches to minimize distributed transactions. We then present the problem statement.

A. Distributed Transaction Processing

In modern distributed database systems, data is typically divided into several partitions. Each of these partitions has one primary replica and may have multiple secondary replicas. For example, as shown in Figure 1, there are two nodes N_1 and N_2 and two partitions P_1 and P_2 . The primary replica of P_1 is located on N_1 while its secondary replica is on N_2 . Conversely, P_2 has its primary replica on N_2 and its secondary on N_1 . When the context is clear, we use P_1 to denote the primary replica of partition P_1 , and use P'_1 to represent the secondary replica of partition P_1 .

Distributed transaction processing is essential to manipulate data on multiple nodes while ensuring ACID properties and availability. The standard approach [1] for processing a distributed transaction T consists of three phases: 1) execution phase, 2) prepare phase, and 3) commit phase. As shown in Figure 1, suppose the transaction T involves a write operation on the data item x , denoted as $W(x)$, and a read operation on the data item y , denoted as $R(y)$. In the execution phase, the coordinator of transaction T first distributes these read/write requests to the corresponding primary replicas for local execution. For example, T sends the write operation $W(x)$ to P_1 , while sending the read operation $R(y)$ to P_2 , respectively. These primary replicas (P_1 and P_2) are then regarded as participants.

Once all read/write operations of T are successfully executed, the coordinator employs the two-phase commit (2PC) to commit/abort T . First, the prepare phase begins, during which the coordinator sends a prepare message to all participants. Upon receiving this message, each participant votes for

whether to commit T and replicates its prepare log, containing this vote and the data item to be written (if applicable), to the corresponding secondary replicas. For example, supposing P_1 decides to commit T , it replicates the commit vote with the data item x to P'_1 . If all participants agree to commit T , the coordinator then moves T to the commit phase. During this phase, the commit decision is sent to these participants and replicated to all involved secondary replicas.

According to this process, executing and committing a distributed transaction requires multiple (≥ 5) round trips of blocking communication, which limits transaction throughput [4]. In this paper, we therefore focus on enabling most transactions can be executed as single-node transactions to avoid the costly network overheads.

B. Distributed Transactions Minimization

We now analyze existing approaches that minimize distributed transactions to clarify the design space of `Lion`.

1) *The data migration technique*: Data migration is commonly used to co-locate correlated partitions, through either a “pull” or “push” method [16]–[19]. However, during the migration, transactions accessing these partitions may abort or get blocked, resulting in service downtime and performance degradation.

Several earlier works are designed as offline tools for repartitioning the physical layout of databases. For example, Schism [5] and Sword [6] identify co-accessed partitions based on historical data and initiate the partition replacement before the database restarts. However, the downtime is unbearable and the layout becomes stale when the workload changes. Plus, they do not account for the placement of secondary replicas, leading to unnecessary migrations.

Leap [8] pursues adaptivity to dynamical workload through an aggressive migration strategy at the transaction level. This approach allows each node to migrate remote data to the local node for each operation before execution. However, this approach may lead to two potential issues: the “ping-pong” problem and the load imbalance problem since all partitions will be migrated to the same node facing the skewed workload.

In contrast, Clay [7] and Hermes [12] adopt a sophisticated migration strategy by analyzing more transactions. Clay periodically monitors transaction execution on each node and devises a migration plan upon detecting load imbalances. Its approach involves transferring partitions from overloaded nodes to others to resolve the imbalance. However, as the primary focus is load balancing through repartitioning, Clay can not eliminate all distributed transactions. Because it sometimes misinterprets the overloaded node, running single-node transactions, as having a similar load to nodes with fewer distributed transactions and fails to start the repartitioning. Hermes combines migration with deterministic protocols. Similar to Leap, it migrates data on the fly with distributed transaction processing. But it collects transactions in batches and reorders the batch sequence to keep transactions accessing the same partitions together. This approach reduces the “ping-pong” effect since the following transactions may reuse the

TABLE I: Comparison of `Lion` with existing approaches

	Key Designs	Dynamic Adaptivity	Migration Efficiency	Load Balancing	Execution Constraints
2PC	Distributed Transactions	N/A	N/A	✗	N/A
Schism [5]	Offline Repartitioning	✗	✗	✗	N/A
Leap [8]	Aggressive Migration	✓	✗	✗	N/A
Clay [7]	Periodical Migration	✓	✗	✓	N/A
Hermes [12]	Deterministic Migration	✓	✗	✓	In batches
Star [10]	Full Replication	N/A	✓	✗	In batches
<code>Lion</code>	Adaptive Replication	✓	✓	✓	N/A

migration caused by earlier ones. However, migration is still inevitable and severely interferes with transaction processing.

2) *The full-replication technique*: The full-replication technique ensures that at least one node contains complete replicas of all partitions. It is usually equipped with data remastering techniques [13], [20] to eliminate distributed transactions. Compared with data migration, remastering is a lightweight technique since it only transfers from the primary replica to another secondary one, instead of the entire partition copy. Star [10] introduces a “super node” with full replication, routing all distributed transactions to that node. The system periodically remasters primary replicas to the “super node” to process them as single-node ones. However, as the cross-partition ratio increases, the “super node” becomes overloaded and causes the bottleneck. DynaMast [11] assumes that all nodes have full replicas and uses dynamic mastering to ensure that each node has an appropriate proportion of mastership, which is not within the scope of partition-based replication methods. However, it is unfeasible to have certain nodes store all data replicas. As the number of nodes increases, the synchronization costs between all replicas can become unsustainable.

3) *Summary*: In Table I, we compare existing approaches with `Lion` across four dimensions. All methods, excluding 2PC, aim to enhance transaction processing efficiency through data migration and full replication. However, Schism and Star lack adaptability to dynamic workloads. Schism relies on an offline strategy, while Star cannot dynamically modify its replica placement. Leap, Clay, and Hermes suffer performance degradation during migration. Leap and Star experience bottlenecks due to inadequate load balancing considerations. Additionally, Star and Hermes may impose batch or deterministic constraints on transaction execution. In contrast, `Lion` considers improving across all dimensions by its adaptive replica provision.

C. Problem Definition

We now present a formal definition of the problem addressed by `Lion`. The database comprises multiple nodes denoted as $N = \{N_1, \dots, N_n\}$. Each node’s storage includes a collection of partitions represented as $P = \{P_1, \dots, P_m\}$. To ensure high availability, each partition must contain a minimum of k replicas, distributed in a default round-robin fashion. Given a batch of transactions $B = \{T_1, \dots, T_b\}$, our objective is to determine a new replica placement P' . Let ϵ represent the percentage of load imbalance permissible within the system, and let θ denote the average load across all nodes in P' multiplied by $1 + \epsilon$. Further, we outline three requirements in

the following formula. First, we aim to minimize distributed transactions. Second, the replica rearrangement cost from P to P' should be minimized. Third, the load should be balanced under P' .

$$\begin{aligned} \text{minimize } C_e(B, P') = \sum_{i=1}^B f_c(n_i, T_i), C_p(P, P') \quad (1) \\ \text{s.t. } \forall x \in N, f_b(x) < \theta \end{aligned}$$

Given a transaction T_i routed to the node n_i , $f_c(n_i, T_i)$ represents the execution cost of T_i on n_i , mainly including the remote access and commit cost. $C_e(B, P')$ represents the cost summation for B under P' . $C_p(P, P')$ signifies the replica rearrangement cost, including migration and remastering costs. We impose an upper limit of θ on the load $f_b(x)$ for each node to maintain load balance.

III. SYSTEM OVERVIEW

`Lion` is applicable for modern distributed databases [1]–[3] that maintain multiple replicas of data partitions. For illustration purposes, we overview `Lion` based on the share-nothing architecture consisting of multiple monolithic nodes. Each node is responsible for storing data replicas and processing incoming transactions. Figure 2 shows the system overview of `Lion`. To minimize distributed transactions, `Lion` introduces two system components: 1) **planner**, a specific kind of node, including a workload analyzer and a plan generator, to determine an optimal replica placement plan; 2) **adaptor**, a component within each node to adjust the replica placement based on the generated plan.

We dynamically adjust replica provision according to the workload through the following steps:

- *Workload analysis*: After collecting B recently received transactions, we invoke the workload analyzer to identify co-accessed partitions. In addition to analyzing the past B transactions, our approach includes K predicted future transactions, which are generated by the proposed workload prediction technique. We construct a graph based on the read/write sets of these $B + K$ transactions. In the graph, we represent partitions as vertices and transactions as edges connecting partitions that are co-accessed. We utilize a clustering algorithm to precisely group these partitions into clumps. Each clump consists of partitions and indicates that the partitions within it should be placed at the same node. This method represents both current and anticipated future access patterns of transactions. The detailed workload analysis and prediction techniques will be presented in Section IV-A and Section IV-C.
- *Plan generation*: We then employ the plan generator to determine the optimal replica placement plan. We design a cost model to assign each clump to a specific node, taking into account factors such as replica rearrangement cost and load balancing. By considering multiple replicas, this model ensures less adjustment overhead and better load balancing than the prior methods designed for the single replica setting. We shall detail the plan generation in Section IV-B.

- *Asynchronous adjustment*: Lastly, the adaptor asynchronously adjusts the replicas according to the established plan. In particular, the adaptor will add or remove secondary replicas by invoking the replica manipulation functions [2], [3] inherent in the replication-based database.

After pre-allocating replicas according to the plan, `Lion` dynamically adjusts the position of primary replicas and converts the distributed transactions into single-node ones. This is achieved by the following two steps. Firstly, we dispatch the transaction T to the node with well-prepared replicas. To accomplish this, we introduce a set of transaction routers, each of which is equipped with a cost model identical to the planner’s. The router will dispatch T to a node with maximum requisite replicas, where the execution cost is the lowest. Secondly, we utilize the remastering technique to finalize the conversion. On the executor, T is processed through the execution, prepare, and commit phases. During the execution phase, operations are directly executed on the node if it possesses the necessary primary replicas. `Lion` remasters before executing that operation if it has a secondary replica. If all operations can be executed on a single node, the transaction can be directly committed, omitting the prepare phase. Otherwise, `Lion` processes T as a standard distributed transaction with 2PC. The process of remastering is performed as follows, which is widely used in [2], [3], [21], [22]. First, a secondary replica is selected as the candidate based on our proposed replica rearrangement algorithm and new read/write operations will get blocked on the primary replica. Then the lagging logs will be synchronized from the leader to the target secondary replica, ensuring its state is consistent with that of the primary replica. After that, a leader election starts to prompt the candidate to be the new primary replica, who then continues to perform operations. During remastering, there are potential risks of data inconsistency and split-brain problems [23]. We address these following the TiDB’s approach [3]. This approach utilizes log synchronization to maintain consistency between the new primary replica and the old one during remastering. Further, to prevent split-brain issues, it blocks new operations during remastering to ensure that only one primary replica provides service at any given time.

The replica provisioning ensures that co-access partitions will be placed on one node and transactions accessing the same partitions are deliberately routed to the same node, which reduces “ping-pong” remastering across the nodes. In scenarios where a remastering conflict emerges, one transaction completes the remastering successfully while others resort to committing as distributed transactions. For example, transactions T_1 and T_2 are routed to nodes N_1 and N_2 , respectively, and both attempt to remaster the overlapping replicas simultaneously. Assuming the success of T_1 and the failure of T_2 , subsequent transactions resembling T_2 , which access the same partitions, will endeavor to route to N_1 whenever possible. Otherwise, they will execute through 2PC.

Example 1: As shown in Figure 2, there are three nodes N_1 , N_2 , and N_3 , and three partitions P_1 , P_2 , and P_3 . Suppose x , y , and z is stored in P_1 , P_2 , and P_3 . The primary replica of

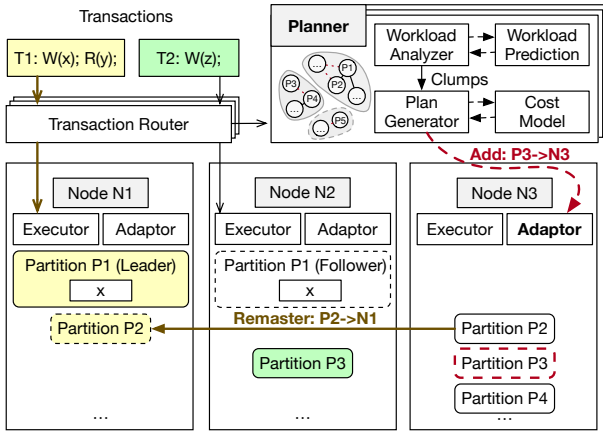


Fig. 2: System overview of `Lion`

P_1 , P_2 , and P_3 is located on N_1 , N_3 , and N_2 , respectively. Let us consider two transactions, T_1 and T_2 . T_1 contains a write $W(x)$ and a read $R(y)$, while T_2 involves a write $W(z)$. For T_1 , the router sends T_1 to N_1 as it has the primary replica of P_1 . However, the replica of P_2 is on N_3 . Therefore, after T_1 executes $W(x)$, its read is executed after the adaptor on N_1 remasters P_2 to N_1 . In contrast, T_2 is routed and executed as a single-node transaction on N_2 without remastering since all the required primary replica exists in N_2 . Consider a transaction T_3 writes to partitions P_3 and P_4 . Suppose the primary replicas of P_3 and P_4 are on N_2 and N_3 , respectively. T_3 can not be converted into a single-node transaction through remastering since neither of the nodes has all replicas for it. In this case, T_3 would be executed as a distributed transaction. After collecting and analyzing transactions periodically, `Lion` decides to add a secondary replica of P_3 on N_3 to co-locate partitions P_3 and P_4 . As a result, when a subsequent transaction T_4 arrives that reads P_3 and P_4 , it can be processed as a single-node transaction on N_3 . The details of that adaptive replica mechanism are further explained in Section IV-B (Example 2) and Section IV-C (Example 3).

IV. THE DESIGN OF `LION`

In this section, we present the design of `Lion` in detail.

A. Workload Analysis Technique

To determine which partitions should be placed on the same node, we employ a graph-based algorithm within the workload analyzer. This algorithm treats a batch of transactions as a graph and produces a set of clumps, where each clump represents co-accessed partitions. The process of this algorithm is depicted in Figure 3, comprising two stages: 1) graph construction, defining transactions as a graph, and 2) clump generation, extracting co-access information from the graph.

Graph Construction. Within `Lion`, we retain the partition IDs accessed by each transaction, alongside its *TxnMeta* information, such as *TxnID* in the transaction context. The involved partitions of a given query are determined after the SQL parsing procedure and will be further pruned by the

query optimization. We record these results as a new variable *TxnParts* in the *TxnMeta*.

Based on that, we begin by modeling the workload’s access patterns into a heat graph, denoted as an indirect weighted graph $G(V, E)$, with vertices in set V and edges in set E . We accumulate the weight of both vertices and edges to represent the access frequency of the partition and the co-access possibility between partitions. Each partition accessed by a transaction is treated as a vertex, labeled as v , with its weight denoted as $w(v)$. Edges, referred to as e , connect pairs of vertices representing partitions accessed by the same transaction. Notably, the weights of edges connecting multiple nodes denoted as e_c ($e_c = (u, v) | u \in N_i, v \in N_j, i \neq j$), are considerably greater than the weights of edges within a single node, termed e_s ($e_s = (u, v) | u \in N_i, v \in N_j, i = j$). This emphasizes the higher priority given to e_c when considering edge weights. Moreover, we employ a priority queue *hVertices* to record the most frequently accessed vertices, aiming to expedite the subsequent clustering process. In Figure 3a, we provide an illustrative example illustrating the transactions batch collected by the planner and the resultant construction of the $G(V, E)$ graph.

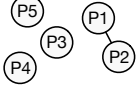
Clump Generation. After the graph $G(V, E)$ is constructed, a clustering algorithm will be triggered to identify partitions that are frequently accessed by the same transactions as clumps. Each clump c contains a set of involved vertices ($c.pids$), the weighted sum of the vertices ($c.w$), the destination to place the clump ($c.n$), etc.

The clustering algorithm selects a vertex and expands on its neighbors until all co-accessed vertices are included. It starts with the hottest unused vertex v from *hVertices* as a seed and evaluates neighboring vertices v_{adj} based on their connection weight $w(e)$ against the threshold α . If this weight surpasses α , indicating high co-access or different nodes, these vertices are grouped in a clump. Conversely, vertices with weaker correlations or independent access are placed into separate clumps. As the clump expands, $c.w$ will get updated with each inclusion of new vertices, which facilitates load balancing for later clump reallocation. Once all related neighbors have been explored, the algorithm concludes the search for the current clump, selects a new seed and proceeds to the next clump. This process continues until all vertices in *hVertices* have been visited. The generated clumps serve as input for the replica rearrangement strategy algorithm (in Section IV-B). Figure 3b illustrates the creation of four clumps, representing co-located partition sets. C_1 comprises partitions P_1 and P_2 , with a weight of 4 (each vertex weighs 2). Meanwhile, C_2 , C_3 , and C_4 contain one partition each.

B. Replica Rearrangement Strategy

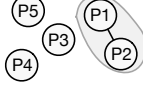
Once the clumps are identified, `Lion` assigns them across nodes to co-locate replicas for the partitions within a clump. Nevertheless, an imprudent clump placement can lead to a significant amount of partition movement, consuming system resources and potentially interfering with current transaction execution. To mitigate these issues, we introduce a replica

T-ID	P-ID	T-ID	P-ID
T1	P1,P2	T5	P5
T2	P3	T6	P4
T3	P4	T7	P5
T4	P1,P2		



(a) Graph construction

C-ID	P-ID	W
C1	P1,P2	4
C2	P3	1
C3	P4	2
C4	P5	2



(b) Clumps generation

Fig. 3: Example for workload analysis technique

rearrangement strategy. This strategy aims to distribute clumps with minimal movement cost while ensuring load balancing.

1) *Rearrangement Objective*: We first define the problem addressed by this algorithm and refine the objective as discussed in Section II-C. Given a collection of clumps C and the current replica placement P , this algorithm is designed to determine a reconfiguration plan RP for the clump placement. To facilitate this, we employ a `ReconfigurationPlan` structure. The RP constitutes a mapping between clumps and nodes, where each entry $\langle c, n \rangle$ denotes that clump c will be assigned to the node n . Applying the RP will adjust the replica placement and result in a new placement P' . During this process, the assignment of C will incur additional operational costs. Each partition v within clump c should be relocated to node n through a series of operations, involving data migration and remastering. The cost fluctuates across different scenarios, as outlined below:

- **Case 1.** If node n hosts the primary replica of v , i.e. $N_p(v, p) = n$, no additional cost is incurred. $N_p(v, p)$ represents the node that hosts the primary replica of the given partition v .
- **Case 2.** If node n has a secondary replica of v , i.e. $n \in N_s(v, p)$, remastering can be used to effect the conversion at a cost of w_r . $N_s(v, p)$ represents the set of nodes with the secondary replica v .
- **Case 3.** Otherwise, if node n lacks any replica of partition v , data migration becomes inevitable, incurring a cost of w_m .

The RP should adhere to the objectives in Equation 1. More specifically, the second objective is extended by Equation 2. Here, $f_o(n_i, c_i)$ represents the cost for assigning c_i to n_i .

$$\text{minimize } C_p(P, P') = \sum_{\{n_i, c_i\} \in RP_i} f_o(n_i, c_i) \quad (2)$$

2) *Cost Evaluation*: Subsequently, we delineate the process of evaluating cost and identifying the destination for each clump. Our approach involves a cost model to assess the cost for each clump c across all nodes and heuristically selects the node n with the lowest cost as the destination. Let us start with the following three scenarios for cost evaluation: If n hosts all primary partitions, the cost model prefers to place c on it with no additional cost. In case n lacks all primary replicas but possesses all secondary ones, it is still an ideal choice since expensive migration can be avoided. The placement only introduces the lightweight remastering cost for the secondary

replicas in the later process. However, if n doesn't possess all required replicas, data migration becomes necessary to add replication for the lacking partitions.

We summarize these scenarios in Equation 3 to calculate the cost of placing c on n .

$$f_o(n, c) = w_r * \sum_{v \in V_c} cnt_r(v, n) + w_m * \sum_{v \in V_c} cnt_m(v, n) \quad (3)$$

where V_c represents the partitions that c possesses. $\sum_{v \in V_c} cnt_r(v, n)$ and $\sum_{v \in V_c} cnt_m(v, n)$ stands for the number of secondary and the lacking partitions on n respectively. The calculation of these two kinds of partition is defined in Equation 4. Note that we also track the normalized access frequency of replicas for each partition as $f(v, n)$. A higher $f(v, N_p(v, p))$ indicates that the remastering cost would be more substantial. When the current primary replica is actively accessed, it might disrupt the ongoing transactions or encounter blocking until the execution is completed.

$$cnt_r(v, n) = \begin{cases} 1 + \log_2(f(v, N_p(v, p)) + 1), & n \in N_s(v, p) \\ 0, & \text{else} \end{cases}$$

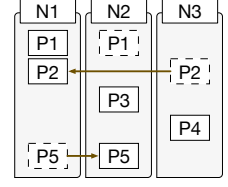
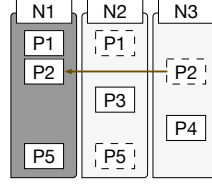
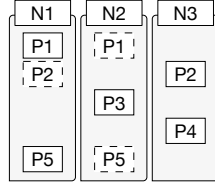
$$cnt_m(v, n) = \begin{cases} 1, & n \notin N_p(v, p) \cup N_s(v, p) \\ 0, & \text{else} \end{cases} \quad (4)$$

We allow users to set a maximum replica limit for each partition according to their requirements. Note that this replica number is not a theoretical limit of `Lion`. Upon exceeding this limit after adding new replicas, we remove one replica from the replica group. We opt to remove the secondary replica with the lowest $f(v, n)$ by designating it with a `delete_flag`. Subsequently, replica synchronization ceases to update the data in the flagged one. Different from `Clay` and `Hermes`, `Lion` benefits two-fold by considering the placement of multiple replicas: first, it minimizes migration expenses, thereby minimizing disruption to the current system execution; second, it avoids full replication, thereby lessening the synchronization overhead for consistency.

3) *Rearrangement Algorithm*: We finally deploy the replica rearrangement algorithm within the plan generator to achieve the aforementioned objective. This algorithm, operating on a set of clumps C and the existing replica placement P , facilitates the determination of the RP . The algorithm operates in two distinct steps:

Clump dispatching. The plan generator initially assigns the clump to a destination node with a minimal cost. It iterates through every clump c_i , and employs the function `FindDstNode()` to select the suitable target node n . This function utilizes the cost model, as expressed in Equation 3, to evaluate costs across all nodes, recording the interim costs in m_c . Subsequently, it identifies the destination with the lowest cost for c_i (lines 5), and the updated c_i is added to the reconfiguration plan RP . To monitor load balance, the balance factor b_i is updated once the plan generator determines the destination for each clump. This factor incorporates the weight c_w of each clump (lines 7). Additionally, a priority queue q is employed to log the clumps assigned to individual

C-ID	P-ID	W
C1	P1,P2	4
C2	P3	1
C3	P4	2
C4	P5	2



(a) The input clumps

(b) The original replica layout

(c) The clump dispatching phase

(d) The load fine-tuning phase

Fig. 4: Example for the replica rearrangement algorithm

Algorithm 1: Replica Rearrangement Algorithm

```

1 Function Rearrangement ( $C, P$ ):
2    $N \leftarrow$  the number of nodes in  $P$ . ;  $rp \leftarrow []$ ;  $m_c \leftarrow []$ 
3    $b_i \leftarrow 0 \forall i = 1, \dots, N$ ;  $q_i \leftarrow [] \forall i = 1, \dots, N$ 
4   for  $c_i \in C$  do
5      $c_i.n \leftarrow \text{FindDstNode}(c_i, P, m_c)$ 
6      $RP.push(c_i)$ ;  $q[c_i.n].push(c_i)$ 
7      $\text{UpdateBalance}(c_i, b)$ 
8    $avg \leftarrow \frac{\text{LoadSum}(C)}{N}$ ;  $is\_done \leftarrow false$ 
9   while  $! \text{CheckBalance}(avg, b)$  and  $! is\_done$  do
10     $step \leftarrow A$ 
11     $oN, iN \leftarrow \text{FindOINodes}(b, avg)$ 
12    if  $|oN| = 0$  or  $|iN| = 0$  then
13      break
14    while  $! \text{CheckBalance}(avg, b)$  and  $step > 0$  do
15       $idx, i_n, is\_find \leftarrow$ 
16         $\text{PickClump}(q, oN, iN, m_c, b)$ 
17      if  $! is\_find$  then
18        break
19       $RP[idx].n \leftarrow i_n$ 
20       $\text{UpdateBalance}(RP[idx], b)$ 
21      if  $|oN| = 0$  or  $|iN| = 0$  then
22         $step \leftarrow 0$ 
23      else
24         $step \leftarrow step - 1$ 
25    if  $step = A$  then
26       $is\_done \leftarrow true$ 
27 return  $RP$ 

```

nodes, sorted by their respective weights in ascending order (lines 6). These data structures facilitate efficient fine-tuning in subsequent steps if necessary.

Load fine-tuning. Subsequently, the plan generator adjusts the current RP to maintain load balance among nodes if necessary. The function $\text{CheckBalance}()$ calculates the balance variance and compares it against a predefined threshold θ (line 9). If the variance is below θ , the algorithm concludes. Otherwise, suggesting a potential load imbalance caused by RP , a fine-tuning mechanism is activated to address this issue (lines 9-25). The basic idea is to reassign clumps from overloaded nodes to idle ones to bridge the load disparity with additional operation costs. Firstly, the function $\text{FindOINodes}()$ identifies overloaded nodes oN and idle nodes iN based on the b_i (line 11). If no overloaded nodes are detected, the loop terminates. Otherwise, the function

$\text{PickOneClump}()$ is employed to transfer a clump from oN to iN (line 15). This function initially selects an overloaded node o_n from oN and calculates the required clump size sz based on the load gap between o_n and iN . Then it searches the $q[o_n]$ to find a clump less than sz . If succeeds, it returns the idle node i_n from iN with the lowest cost. Otherwise, the function retries on other overloaded nodes. If all retries fail, the loop exits (line 25). Upon successfully identifying a qualified clump, it updates its destination, along with the balance factors and oN and iN (lines 18-19). A variable $step$ is defined to save recalculations of $\text{FindOINodes}()$. The replica rearrangement algorithm concludes its iterations when the load balance falls within θ .

Example 2: To illustrate, consider the clumps and replica placement in Figure 4a and 4b following the example discussed in Section III. For the sake of simplicity in the following examples, we assume that all replicas have approximately the same access frequency. Upon executing the first step outlined in Algorithm 1, all clumps are initially dispatched to nodes with minimized costs, as shown in Figure 4c. For instance, C_1 opts for N_1 as its destination. Since the costs for it to N_1 , N_2 , and N_3 are w_r , $w_m + w_r$, and w_m , the cost at N_1 is lowest. Similarly, C_2 , C_3 , and C_4 choose N_2 , N_3 , and N_1 without costs. However, this allocation leads to load imbalance, where N_1 is overloaded with a weight of 6, while other nodes maintain weights of 1 and 2. To rectify this, the second step is triggered to transfer a clump from N_1 to others. The algorithm selects C_4 and designates N_2 as the new destination. Given that N_2 is an idle node and has a secondary replica, the reassignment of C_4 to N_2 incurs an additional cost of w_r . The final replica placement depicted in Figure 4d demonstrates a balanced load with an operation cost of $2 * w_r$.

C. Workload Predication

While Lion adapts to dynamic workloads through the planner, the time-consuming process of replica rearrangement still poses a formidable challenge. As the workload changes, all transactions bear the burden of 2PC overhead until replication becomes fully prepared. To expedite this process, we combine Lion with a workload prediction mechanism. It is aimed at forecasting workload patterns and proactively adding replication for co-accessed partitions, termed pre-replication. Unlike Hermes [12], which relies on foreknowledge of future read/write operations and suits specific deterministic database

scenarios, our prediction mechanism imposes no such constraints and is suitable for non-deterministic systems.

1) *How to predict co-accessed partitions:* The basic idea is to convert discrete transaction access into a time-series analysis problem. The prediction technique can be divided into the following three phases:

Template Identification. The arrival rate history [24], [25] is a crucial metric used to characterize the access pattern of transaction queries, formulated in Equation 5. The *ar* of a query signifies its access frequency changes over time. It can be visualized as a curve with a sampling interval i , where the x-axis represents the sampling time point t and the y-axis denotes the sum of query frequency $f(n)$ within i . However, maintaining the *ar* information for every query can be costly. To mitigate this, we establish a partition-based rule for labeling transactions based on their access patterns. Transactions accessing the same partitions receive the same label, forming identical templates. Once these templates are identified, we track the arrival rate history of each template instead of individual queries. In Figure 5a, we identify five templates under the partition-based rule, corresponding to the example in Section III. Of these, four templates (P_1P_2 , P_3 , P_4 , and P_5) were active before the timestamp t_1 , while the other two templates (P_3P_4 and P_5P_6) became dominant thereafter.

$$ar(t, i) = \sum_{n=t}^{t+i} f(n) \quad (5)$$

Workload Classification. Two templates are deemed similar if their arrival rates increase and decrease simultaneously, a similarity evaluated by computing the cosine distance between their *ar* values. To improve prediction efficiency, templates with a calculated distance below a predefined threshold β are merged into the same workload class. Consequently, subsequent predictions will be performed for the merged workloads rather than individual templates. Each workload comprehensively stores information for all its constituent templates, including partition IDs and their associated access frequencies. During pre-replication initiation, reservoir sampling [26] assists in identifying partitions from the workload that are highly likely to appear soon. In Figure 5b, we demonstrate the consolidation of five templates into two distinct workloads. Specifically, P_1P_2 , P_3 , P_4 , and P_5 constitute workload W_1 , while W_2 encompasses P_3P_4 and P_5P_6 . Timestamp t_1 serves as the boundary between these two workloads.

Time-series Prediction. We utilize an LSTM model to forecast future workload trends based on the historical *ar* of each workload. Compared to LSTM, traditional methods like linear regression and traditional RNNs struggle to effectively capture long-term dependencies and handle non-linear dynamics within sequences, thus exhibiting limitations in handling complex time series patterns [24], [25], [27], such as the co-accessed partitions focused in our paper. We utilize a lightweight LSTM model in `Lion`, which does not necessarily require a GPU for training because the training latency is acceptable even on a CPU. We train the model periodically based on the logs that record the partitions accessed by

transactions in a given time period. In particular, we build the query arrival rates (as formulated in Equation 5) based on the raw log. When the mean squared error (MSE) between predicted and actual results falls below a predefined threshold, we retrain the model to maintain the model accuracy. We acknowledge that more orthogonal optimizations can be explored such as using GPU to expedite model training [28] and incorporating superior algorithms [29], [30] to improve the prediction accuracy. As we primarily focus on leveraging replicas to minimize distributed transactions, we directly use a standard LSTM for workload prediction without specific improvement or fine-tuning. We defer these optimizations to our future work.

For the upcoming workload with high anticipated arrival rates, we select templates based on their access frequencies. Subsequently, the co-accessed partitions within the template are integrated into the graph $G(V, E)$. `Lion` employs a replica rearrangement algorithm that creates a unified plan based on both historical and predicted workloads, rather than producing separate adjustment plans for each type of workload. We model historical and predicted workloads as a single graph for that purpose. To fine-tune the influence of predicted workloads on our planning process, we employ a parameter w_p that determines their weight within the graph. This parameter serves as a weighted coefficient for incorporating predicted information into the graph. A weight of 0 signifies that the prediction algorithm is inactive. By default, w_p is set to 1. For instance, in Figure 5b, assuming the current timestamp is t_2 and the active workload is W_1 . After forecasting, it indicates that the arrival rate of W_2 will surpass W_1 . Consequently, we sample the template P_3P_4 from W_2 and incorporate it into G . The predicted co-accessed partitions are included as additional edge weights, visualized by the red dashed line connecting P_3 and P_4 in Figure 5c.

2) *When to trigger pre-replication:* We’ve devised a workload variation metric wv to assess the variance between the present workload and its anticipated future state, as defined in Equation 6. After periodic evaluations, the pre-replication will be triggered when wv surpasses a predefined threshold γ .

$$wv(t, h) = \sqrt{\frac{1}{n} \sum_{k=1}^n (a_k(t+h, \delta) - a_k(t, \delta))^2} \quad (6)$$

Here, h delineates the prediction horizon, depicting how far into the future a model can forecast. $ar_k(t, i)$ represents the arrival rate of workload W_k at timestamp t within sampling interval i . The divergence from timestamp t to $t+h$ is expressed as $wv(t, h)$ across all potential workloads. When $wv(t, h) > \gamma$, a significant impending workload shift will occur soon, prompting the initiation of pre-replication for forthcoming transactions. For example, in Figure 5b, the wv reaches its peak at timestamp t_2 , signifying a drastic alteration in workload expected at the future timestamp t_2+h . This indicates an imminent and significant shift in workload dynamics.

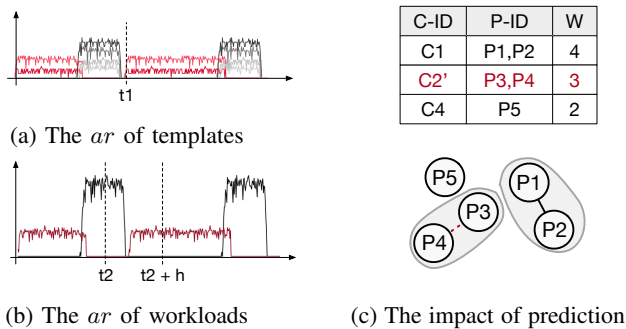


Fig. 5: Example for workload prediction

Example 3: Finally, we recap the Example 1 and explain the impact of the prediction mechanism on replica placement rearrangement in Section III. In the context of the aforementioned transaction batch and replica placement, the prediction mechanism anticipates that P_3 and P_4 will be co-accessed soon, corresponding to the transaction T_3 that writes to these two partitions. Consequently, it merges C_2 and C_3 with a collective weight of 3, as shown in Figure 5c. During the replica arrangement process, the new C'_2 is relocated to N_3 to maintain load balance, incurring an additional replication cost of w_m . That’s why the plan generator instructs N_3 to execute $Add : P3 \rightarrow N3$ instead of N_2 in Section III.

D. Batch Optimization

To further reduce the cost of remastering, we’ve introduced an asynchronous remastering optimization for the batch version of `Lion`. The idea is to remaster the partition asynchronously for the whole batch before starts transaction processing, which allows overlapping of network delays and reduces the remastering overhead significantly.

In batch processing, transactions routed to the executors are buffered as a batch and wait to get executed until the global batch epoch is incremented, either by reaching a specified batch size (default 10k) or after a set time window. This optimization adds a remastering phase before execution. Upon dispatching a transaction to the executor’s buffer queue on the local node, the system checks if it’s feasible to convert the transaction into a single-node one, as discussed in Section III. If viable, the executor sends a remastering request asynchronously; otherwise, it bypasses it directly. Unlike the standard execution mode, the executor doesn’t stall for the remastering response but proceeds to subsequent incoming transactions. The transaction index within the batch will be piggybacked along with the remastering messages to help the executor locate the transaction context in the buffer. Once the batch epoch is incremented, the executor proceeds to the next execution phase only after ensuring acknowledgment of all remastering requests. To achieve this, barriers are implemented using network round trips between the remastering and execution phases. As the execution begins, the executor can locally commit those transactions that have successfully undergone conversion. For instance, let’s consider the scenario

where transactions T_1 and T_2 are received successively within a time window. Suppose both T_1 and T_2 meet the conversion conditions. When T_1 arrives, the executor initiates the remastering messages. If T_2 arrives while T_1 ’s remastering is ongoing, the executor asynchronously triggers the remastering for T_2 , without waiting for the acknowledgment of T_1 .

V. IMPLEMENTATION

Our implementation, built upon the C++ codebase of `Star` [10], comprises two node types: the distributor node and the executor node. Every node is equipped with a global router table, which stores the locations of master and secondary replicas through a hash map. The distributor node will generate transactions and route the transaction to the appropriate executor node. The executor nodes are responsible for partitioned data storage and transaction processing. We run multiple threads of three types on these nodes: messenger, worker, and adaptor. A messenger thread sends and receives network messages through TCP sockets. A worker thread is for handling read/write requests. We’ve implemented an adaptor thread that manages data migration and replication with several `MHandler` functions that can possibly be invoked. We can use the `MigReqHandler()` to migrate partitions from one node to another and use `AddRepReqHandler()` to add a replication for the given partition on a certain destination node. `Lion` operates independently as threads on distributor nodes. It collects transactions and adds them into a `Clumps` structure to get coarse-grained clumps. After that, `RearrangeFunc()` will be invoked to generate the rearrangement plan (RP), which is a list and each entry stands for the target partition ID and its destination. The RP will be routed to the adaptor on the executor nodes and adjust the replica layout by calling the corresponding `MHandler` functions. `Lion` adopts an epoch-based group commit mechanism [31] to reduce the cost of replication synchronization. A global epoch is incremented at 10 millisecond intervals or when reaching a 10k batch size. The committed transactions within an epoch are buffered and asynchronously dispatched to all replicas through `AsyncReplicationReqHandler()`. Each transaction’s outcome remains invisible until the epoch ends and all nodes collectively agree to commit transactions from that specific epoch. Note that distributed databases, such as `TiDB` [32] and `Oceanbase` [33] typically offer APIs for replica remastering and replica management, e.g., adding a replica. For example, `transfer-leader` API provided in `TiDB` can enable replica remastering without the necessity of a leader failure. While these databases are indeed built upon consensus protocols, their remastering processes generally follow the approach we described in Section III. Therefore, `Lion` can be integrated into these real-world distributed databases by calling the APIs they provide.

VI. PERFORMANCE EVALUATION

In this section, we evaluate the performance of `Lion` and compare it against state-of-the-art systems. After introducing

TABLE II: Settings of ablation experiments

Variant	Partitioning Strategy	Workload Prediction	Batch Optimization
2PC		✗	✗
Lion (S)	Schism	✗	✗
Lion (R)	Replica Rearrangement	✗	✗
Lion (SW)	Schism	✓	✗
Lion (RW)	Replica Rearrangement	✓	✗
Lion (RB)	Replica Rearrangement	✗	✓
Lion	Replica Rearrangement	✓	✓

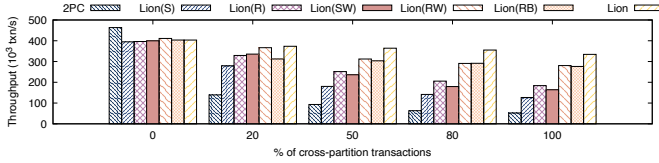


Fig. 6: Performance of ablation experiments

the experimental setup, we evaluate `Lion` in a range of settings to demonstrate its superior performance.

A. Experiment Setup

We run our experiments on a cluster architecture of multiple nodes. The setup comprises 1 distributor node equipped with 36 Intel(R) Xeon(R) Gold 5220 CPUs and 196 GB of DRAM. Additionally, we employ 10 executor nodes, each featuring 8 Intel Core Processor (Skylake) CPUs with 32 GB of DRAM. Experiments are conducted using 4 executor nodes by default. Each partition is initially configured to have 2 replicas. We set the max replica number to 4 due to the constraints on the number of nodes used in default experiments. We use a lightweight LSTM encoder with 2 layers and 20 hidden units for time series prediction and training forecasting models based on the preceding ten-period historical data logs. Iperf3 indicates a network throughput of approximately 937 Mbits/sec between each node. We run 8 worker threads on each executor node, yielding a total of 64 threads. Each node has 2 threads for network communication.

1) *Benchmarks*: All experiments are conducted over the following two benchmarks.

YCSB. The Yahoo! Cloud Serving Benchmark (YCSB) is a simplified transactional workload specifically designed to facilitate performance comparisons across various database and key-value systems [34]. In our evaluation, each data node maintains 24 million data items, resulting in a storage space of 200MB per data node. A parameter called `skew_factor` is used to control the distribution of the accessed data items. The `skew_factor` is set to 0.8 under a skewed workload, resulting in a high load imbalance where 80% transaction tends to access the partitions in the one node. Under the uniform workload, the `skew_factor` is set to 0. The cross-partitioned transactions always access two partitions.

TPC-C. The TPC-C benchmark [35] stands as the industry standard for evaluating OLTP databases. Its dataset comprises 9 relations, and each warehouse is equipped with 100MB of data. By default, we allocate 24 warehouses per node in our experiments. Specifically focusing on NewOrder transactions, the benchmark emulates customers submitting orders to their local district within a warehouse. We simulate scenarios where the same customer makes purchases from different warehouses over time.

2) *Baselines*: To ensure an apples-to-apples comparison, we implemented existing approaches in the same framework as `Lion`. In our experiments, we use the OCC as the concurrency control algorithm, and compare `Lion` with both standard execution and batch execution approaches.

a. Standard execution approaches.

2PC. A classic distributed protocol based on OCC [36]. The processing of the distributed transaction always undergoes the execute, prepare, and commit phases.

Leap. An aggressive transaction management approach. Before executing the operations, it always migrates the master replica from the remote node to the local. When all operations are executable, it commits directly and skips the prepare phase.

Clay. An online partitioning approach. The repartitioning starts when it detects the load imbalance among nodes. Then it generates a partition reconfiguration based on the co-access frequency and adjusts the partitions through data migration. To better compare the cleverness of the reconfiguration, we implement the asynchronous replication and remastering for Clay as `Lion`.

b. Batch execution approaches.

Star. An asymmetric replication approach with a two-phase switching algorithm. It ensures one node has all the partitions. The transactions will be collected in batches. The distributed transactions in the batch will be routed to that node as the single-node one and get committed without 2PC.

Calvin. A classic distributed deterministic approach. It executes the same transaction batch on each replica to avoid 2PC. It requires the declaration of the read/write set before transaction execution. It uses a lock manager to obtain locks for each transaction in the fixed order and the transaction will not be executed until all locks are acquired.

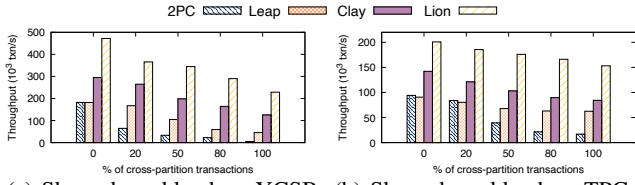
Hermes. A deterministic approach equipped with data migration. It migrates the partition in demand before the lock manager starts to get the locks. It utilizes a prescient transaction routing algorithm to mitigate the “ping-pong” effect while achieving load balance.

Aria. A distributed deterministic approach. It introduces an optimistic write reservation technique to execute the transactions without coordination and without prior knowledge of a transaction’s read/write set.

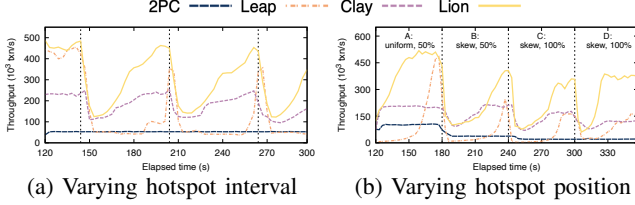
Lotus [37]. Another distributed epoch-based approach. It is implemented with granule locks to enhance concurrency and introduces batch execution/commit for overlapping computation, communication, and asynchronous replication.

B. Optimization Analysis

We first evaluate the effectiveness of three optimizations in `Lion`: the replica rearrangement algorithm in Section IV-B, the workload prediction mechanism in Section IV-C, and the batch optimization in Section IV-D. We conduct comprehensive ablation studies to evaluate the individual contributions of these optimizations in `Lion`. The default workload is uniformed YCSB with 100% distributed transactions. We outline all the variants of `Lion` with different optimizations in Table II, and plot the throughput in Figure 6. As observed, each



(a) Skewed workload on YCSB (b) Skewed workload on TPC-C
Fig. 7: Impact of varying cross-partition ratios (non-batch)



(a) Varying hotspot interval (b) Varying hotspot position
Fig. 8: Impact of dynamic workloads (non-batch)

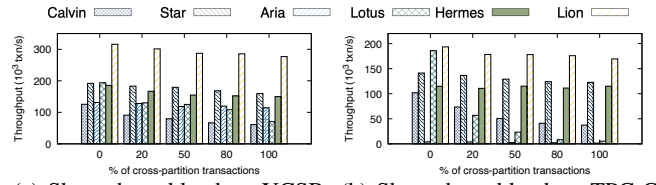
optimization improves transaction performance. First, compared to 2PC, `Lion(R)` which represents `Lion` with only our proposed replica rearrangement strategy, demonstrates up to $3.5\times$ performance improvement. Second, by additionally leveraging workload prediction in `Lion(R)`, `Lion(RW)` shows a performance increase of up to 52.3% over `Lion(R)`. Third, with the further employment of batch optimization, `Lion` achieves up to 20% higher throughput than `Lion(RW)`.

We further compare the proposed partitioning strategy with `Schism` [5]. We implement an alternation of `Lion` using `Schism` as the partitioning strategy, denoted as `Lion(S)`. As depicted in Figure 6, `Lion` outperforms `Lion(S)` by up to $1.7\times$. This can be attributed to the fact that, unlike `Schism`, our partitioning strategy additionally considers replications and takes future transactions into account. We further compare `Lion(S)` with `Lion(R)`, where `Lion(R)` represents `Lion` with only our proposed replica rearrangement strategy, to evaluate the partitioning strategy effectiveness exclusively. As observed in Figure 6, the replica rearrangement strategy outperforms `Schism` by up to 31.1% , primarily because `Schism` does not account for the placement of secondary replicas, leading to unnecessary migrations. We also examine the effectiveness of the prediction mechanism by integrating it with `Schism`, denoted as `Lion(SW)`. As shown in Figure 6, `Lion(SW)` outperforms `Lion(S)` by up to 29.4% , because of the reduced migration cost facilitated by predictions.

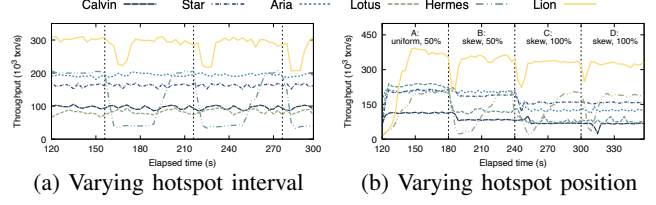
C. Overall Performance

We now compare `Lion` with other standard execution approaches using both YCSB and TPC-C workloads.

1) *Workloads with varying the cross-partition ratio:* We first measure the throughput by increasing the percentage of cross-partition transactions under skewed workloads. The experiments are conducted on YCSB and TPC-C with an 80% skew_factor. We set the default remastering delay to 3000 microseconds to simulate the remastering overhead in real-world scenarios. The results in Figure 7 demonstrate that `Lion` achieves up to $1.9\times$ higher throughput than other approaches. `Lion` is designed to eliminate distributed transactions and uses the replica rearrangement algorithm to dissipate



(a) Skewed workload on YCSB (b) Skewed workload on TPC-C
Fig. 9: Impact of varying cross-partition ratios (batch)



(a) Varying hotspot interval (b) Varying hotspot position
Fig. 10: Impact of dynamic workloads (batch)

imbalances. 2PC costs expensive network coordination to process distributed transactions. `Leap` and `Clay` exhibit better performance since they can adapt to the workload by migration. However, `Leap`'s aggressive strategy does not consider load balance. `Clay` struggles to eliminate all distributed transactions since it perceives the overloaded node running single-node transactions as having an equal load to nodes with fewer distributed transactions.

2) *Dynamic workloads with a changing hotspot:* We evaluate `Lion`'s performance under two dynamically changing workload scenarios. In each scenario, the workload cycles through multiple periods, alternating every 60 seconds. Each period features distinct access patterns, with transactions accessing non-overlapping partitions to create unique hotspots. In the varying hotspot interval scenario, we create three custom queries with a uniform access pattern. The partition ID intervals within each query are fixed in one period but shift among different periods. In the varying hotspot position scenario, we design a combined workload to mimic changes in the most frequently accessed keys in the Zipfian distribution. This workload consists of four periods (A, B, C, D), encompassing uniform access with a 50% cross-ratio, skew with a 50% cross-ratio, skew with a 100% cross-ratio, and skew with a 100% cross-ratio with distribution shift via partition ID offsets.

We then evaluate `Lion`'s throughput fluctuates over time within dynamic workloads. As shown in Figure 8, `Lion` can adapt to new workloads faster as well as maintain higher stable throughput, due to its prediction mechanism and wise partitioning strategy. In contrast, 2PC's performance remains consistently low since it fails to adapt to workloads in Figure 8a. `Leap`'s aggressive migration disrupts transactions and prolongs jitter, evident in skew scenarios B, C, and D (as shown in Figure 8b). `Clay` can not eradicate all distributed transactions when facing cross-partitioned and skewed workloads (as shown in Figure 8a and C, D in Figure 8b).

D. Comparison with Batch Execution Approaches

We next compare `Lion` with other batch execution approaches, including `Star`, `Lotus`, and other deterministic methods. We deploy a single-threaded lock manager for all de-

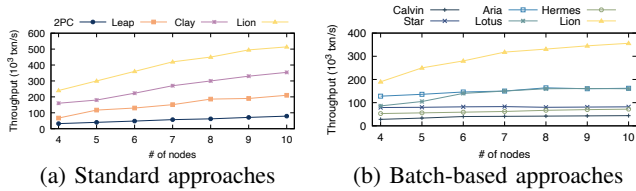


Fig. 11: Scalability

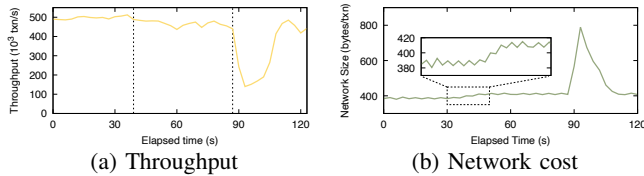


Fig. 12: The analysis of migration process

terministic methods to grant locks to multiple executors following the deterministic order. To make a fair comparison, we implemented a batch-processing version of `Lion` which is introduced in Section IV-D. The batch size is set to be 10k transactions.

1) *Workloads with varying the cross-partition ratio:* We measure the throughput following the same way in VI-C1.

The results shown in Figure 9 show that `Lion` has up to 1.7 \times higher than the next-best approach. `Lion` distributes unrelated co-located partitions evenly among nodes, which enables all workers in each node can independently process transactions with high parallelism. `Calvin` avoids the 2PC through deterministic execution but still suffers from the remote read caused by distributed transactions. As illustrated in Figure 9a and Figure 9b, the performance of `Star` and `Hermes` remains stable when varying the cross-ratio. Because they eliminate distributed transactions either through the full replication or the migration. But their throughput is limited by the bottleneck of a “super node” or a single lock manager. `Aria` and `Lotus` perform well in low cross-ratio scenarios due to their specific epoch-based transaction processing. However, their performance decreases significantly as the cross-ratio increases, because they require a costly commit protocol for distributed transactions and lack optimizations for load balancing. Moreover, their epoch-based schema exacerbates contention. For instance, `Lotus` maintains locks until the end of an epoch, leading to transaction aborts and re-executions, as evident in Figure 14a.

2) *Dynamic workloads with a changing hotspot:* Following the workload introduced in Section VI-C2, we assess `Lion`’s performance using deterministic and batch-based approaches. Notably, `Lion` can adapt to the new workload with 4 \times faster than its non-batch version, experiencing a mere 26.6% throughput degradation in Figure 10a. That can be attributed to its batch optimization with asynchronous remastering. `Hermes` shows its superior performance against other baselines. Because it eliminates distributed transactions through migration in a deterministic pre-defined order for each replica group. However, it experiences a severe performance jitter when adapting to new workloads due to its deterministic migration. Subsequent transactions will be unable to proceed until the preceding distributed transactions complete migration, which

leads to a staggering 62.9% even lower performance compared to `Lotus`(as shown in Figure 10a and the B-C switching boundary in Figure 10b). Others suffer from poor performance since they fail to adapt to the workload. Additionally, they overlook the load-balancing problem when facing skew workloads (as depicted in scenarios B, C, and D in Figure 10b).

E. Scalability

We now study the scalability of each approach, varying the number of executor nodes from 4 to 10 under the same workload (100% cross-partition with uniform access pattern) detailed in Sections VI-D and VI-D. Observing the results as presented in Figure 11, `Lion` achieves approximately 2 \times higher throughput with 10 nodes compared to the scenario involving 4 nodes and has up to 76.4% better scalability than other approaches. This superiority stems from `Lion`’s replica rearrangement strategy, considering distributed transaction elimination and load balancing. This approach ensures optimized performance across each node, and this cumulative advantage becomes more pronounced with an increase in the number of nodes. We note that the throughput of all non-deterministic approaches scales almost linearly as the number of nodes increases. However, `Star` demonstrates poor scalability due to a bottleneck arising from the super node. Deterministic approaches reach their limit post an increase in the number of nodes to 7, mainly resulting from the predefined order established by the sequencer and lock manager, which becomes less conducive in larger node clusters.

F. Migration and remastering analysis

We now evaluate `Lion`’s adaptability to new workloads in depth. Our evaluation begins by analyzing the overhead of our proposed adaptive replica provision mechanism, employing the workload detailed in Section VI-C2, with anticipated changes at the 90s, shown as the right dashed line in Figure 12a. The planner, detecting an impending workload shift, initiates rearranging the replica arrangement by adding replications at the 30s, highlighted by the dashed line in the left segment of Figure 12a. This introduces additional synchronization costs to ensure consistency for these newly-added replicas, which elevates the network cost per transaction execution from 380 to nearly 420 bytes, as illustrated in Figure 12b. However, with our replication arrangement strategy and group commit optimization, we cap the throughput decrease within 5%. A notable surge in network cost emerges around the 90s, peaking at 700 bytes per transaction (Figure 12b) due to remastering requests for single-node conversion. With the pre-replication mechanism, `Lion` mitigates expensive data transmissions and demonstrates better adaptability to dynamic workloads, resulting in a 4x increase in efficiency as shown in Figure 13a. Furthermore, Figure 13b showcases that batch processing, coupled with asynchronous remastering, experiences minimal impact from latency induced by the remastering process.

G. Latency breakdown

We finally analyze the latency of `Lion` and the breakdown of individual phases in comparison to deterministic

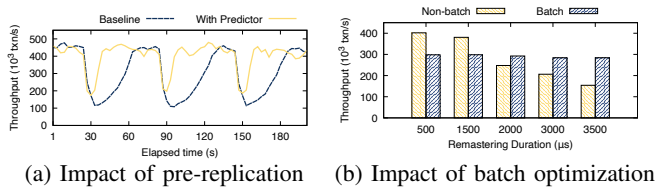


Fig. 13: The analysis of optimization

approaches. Figure 14a and Figure 14b illustrate that `Lion` exhibits latency at the 95th percentile which is 48% lower than `Hermes`. This stable latency is attributed to the group commit optimization and it dedicates 35% of the time to replication synchronization while other deterministic approaches execute the same transaction batch on each replica to maintain consistency. `Calvin` consistently presents high latency across all percentiles due to the necessity of remote reads during distributed transaction execution, consuming over 90% of the execution time. `Aria` employs an optimistic write reservation technique, fostering highly parallel execution without coordination. To reduce the abort ratio, it designs a reordering mechanism that costs an additional 20% latency. `Lotus` introduces nearly zero scheduling time due to its epoch-based execution strategy with asynchronous commit and replication. However, `Aria` and `Lotus` lead to occasional transaction aborts during successive batch processing rounds, resulting in high latency, especially at the 95th percentile. `Hermes` demonstrates adaptability to workloads through partition migration, yet it still requires 19% of the time for scheduling by a single lock manager, affecting executor concurrency.

VII. RELATED WORKS

Substantial efforts have been devoted to optimizing distributed transaction performance by reducing network overheads has gained increasing attention in recent years. For these works [38]–[42], they focus on minimizing network round trips by unifying 2PC and replica synchronization in a single framework. In these approaches, the coordinator simultaneously engages both primary and secondary replicas for voting on whether to commit a transaction. Consequently, they reduce network round trips by removing the need for sequential cross-node coordination and replica synchronization in 2PC-Paxos. However, they commit a transaction only after achieving a majority consensus among all involved replicas. This increased consensus across numerous nodes could potentially intensify the cross-node coordination cost.

As opposed to optimizing the number of network round-trips, quite a few studies [12], [43]–[47] explores deterministic execution, where each transaction is decomposed into multiple sub-transactions that are executed individually on different nodes. Because the equivalent serializable schedule of sub-transactions in each individual node follows the same pre-determined order, expensive coordination among nodes, e.g., 2PC, can be eliminated. Furthermore, lots of works are exploring the benefits of the group-based processing schema [31], [37], [48]. However, these methods require prior knowledge in advance and cannot handle interactive transactions.

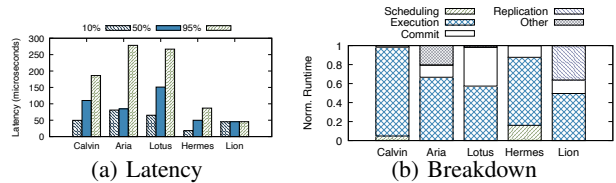


Fig. 14: The analysis of average latency

Another line of research utilizes modern hardware such as RDMA for enhancing distributed transaction processing efficiency. Several studies [49]–[51] offer valuable insights and architectural guidelines to mitigate network bottlenecks. Further, other works [52]–[54] re-implement concurrency control algorithms to be RDMA-compatible, employing efficient one-sided and two-sided verbs. However, these approaches often necessitate special network hardware support and substantial system modifications to leverage this technology.

VIII. CONCLUSION

In this paper, we present `Lion`, a general transaction processing protocol that minimizes distributed transactions by employing replication in today’s distributed databases. With an adaptive replica provision mechanism, `Lion` asynchronously adjusts the replica placement based on the workload, ensuring most transactions can execute on a single node containing all the necessary data replicas. Further, we introduce a workload prediction technique to ensure the replica adjustment can be proactive, which maintains sustainable performance even as the workload dynamically changes. We conduct extensive experiments to compare `Lion` against various transaction processing protocols. The results show that `Lion` achieves up to 2.7x higher throughput and 76.4% better scalability against these state-of-the-art approaches.

REFERENCES

- [1] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, W. C. Hsieh, S. Kanthak, E. Kogan, H. Li, A. Lloyd, S. Melnik, D. Mwaura, D. Nagle, S. Quinlan, R. Rao, L. Rong, Y. Saito, M. Szymaniak, C. Taylor, R. Wang, and D. Woodford, “Spanner: Google’s globally-distributed database,” in *OSDI*. USENIX Association, 2012, pp. 251–264.
- [2] R. Taft, I. Sharif, A. Matei, N. VanBenschoten, J. Lewis, T. Grieger, K. Niemi, A. Woods, A. Birzin, R. Poss, P. Bardea, A. Ranade, B. Darnell, B. Gruneir, J. Jaffray, L. Zhang, and P. Mattis, “Cockroachdb: The resilient geo-distributed SQL database,” in *SIGMOD Conference*. ACM, 2020, pp. 1493–1509.
- [3] D. Huang, Q. Liu, Q. Cui, Z. Fang, X. Ma, F. Xu, L. Shen, L. Tang, Y. Zhou, M. Huang, W. Wei, C. Liu, J. Zhang, J. Li, X. Wu, L. Song, R. Sun, S. Yu, L. Zhao, N. Cameron, L. Pei, and X. Tang, “Tidb: A raft-based HTAP database,” *Proc. VLDB Endow.*, vol. 13, no. 12, pp. 3072–3084, 2020.
- [4] R. Harding, D. V. Aken, A. Pavlo, and M. Stonebraker, “An evaluation of distributed concurrency control,” *Proc. VLDB Endow.*, vol. 10, no. 5, pp. 553–564, 2017.
- [5] C. Curino, Y. Zhang, E. P. C. Jones, and S. Madden, “Schism: a workload-driven approach to database replication and partitioning,” *Proc. VLDB Endow.*, vol. 3, no. 1, pp. 48–57, 2010.
- [6] A. Quamar, K. A. Kumar, and A. Deshpande, “SWORD: scalable workload-aware data placement for transactional workloads,” in *EDBT*. ACM, 2013, pp. 430–441.
- [7] M. Serafini, R. Taft, A. J. Elmore, A. Pavlo, A. Aboulmaga, and M. Stonebraker, “Clay: Fine-grained adaptive partitioning for general database schemas,” *Proc. VLDB Endow.*, vol. 10, no. 4, pp. 445–456, 2016.

- [8] Q. Lin, P. Chang, G. Chen, B. C. Ooi, K. Tan, and Z. Wang, "Towards a non-2pc transaction management in distributed database systems," in *SIGMOD Conference*. ACM, 2016, pp. 1659–1674.
- [9] M. Abebe, B. Glasbergen, and K. Daudjee, "Morphosys: Automatic physical design metamorphosis for distributed database systems," *Proc. VLDB Endow.*, vol. 13, no. 13, pp. 3573–3587, 2020.
- [10] Y. Lu, X. Yu, and S. Madden, "STAR: scaling transactions through asymmetric replication," *Proc. VLDB Endow.*, vol. 12, no. 11, pp. 1316–1329, 2019.
- [11] M. Abebe, B. Glasbergen, and K. Daudjee, "Dynamast: Adaptive dynamic mastering for replicated systems," in *ICDE*. IEEE, 2020, pp. 1381–1392.
- [12] Y. Lin, C. Tsai, T. Lin, Y. Chang, and S. Wu, "Don't look back, look into the future: Prescient data partitioning and migration for deterministic database systems," in *SIGMOD Conference*. ACM, 2021, pp. 1156–1168.
- [13] X. Wei, S. Shen, R. Chen, and H. Chen, "Replication-driven live reconfiguration for fast distributed transaction processing," in *USENIX Annual Technical Conference*. USENIX Association, 2017, pp. 335–347.
- [14] D. Ongaro and J. K. Ousterhout, "In search of an understandable consensus algorithm," in *USENIX Annual Technical Conference*. USENIX Association, 2014, pp. 305–319.
- [15] L. Lamport, "Paxos made simple, fast, and byzantine," in *OPODIS*, ser. Studia Informatica Universalis, vol. 3. Suger, Saint-Denis, rue Catulienne, France, 2002, pp. 7–9.
- [16] A. J. Elmore, V. Arora, R. Taft, A. Pavlo, D. Agrawal, and A. E. Abbadi, "Squall: Fine-grained live reconfiguration for partitioned main memory databases," in *SIGMOD Conference*. ACM, 2015, pp. 299–313.
- [17] A. J. Elmore, S. Das, D. Agrawal, and A. E. Abbadi, "Zephyr: live migration in shared nothing databases for elastic cloud platforms," in *SIGMOD Conference*. ACM, 2011, pp. 301–312.
- [18] U. Cubukcu, O. Erdogan, S. Pathak, S. Sannakkayala, and M. Slot, "Citrus: Distributed postgresql for data-intensive applications," in *SIGMOD Conference*. ACM, 2021, pp. 2490–2502.
- [19] S. Das, S. Nishimura, D. Agrawal, and A. E. Abbadi, "Albatross: Lightweight elasticity in shared storage databases for the cloud using live data migration," *Proc. VLDB Endow.*, vol. 4, no. 8, pp. 494–505, 2011.
- [20] S. Shen, X. Wei, R. Chen, H. Chen, and B. Zang, "Drtm+b: Replication-driven live reconfiguration for fast and general distributed transaction processing," *IEEE Trans. Parallel Distributed Syst.*, vol. 33, no. 10, pp. 2628–2643, 2022.
- [21] Z. Yang, C. Yang, F. Han, M. Zhuang, B. Yang, Z. Yang, X. Cheng, Y. Zhao, W. Shi, H. Xi, H. Yu, B. Liu, Y. Pan, B. Yin, J. Chen, and Q. Xu, "Oceanbase: A 707 million tpmc distributed relational database system," *Proc. VLDB Endow.*, vol. 15, no. 12, pp. 3385–3397, 2022. [Online]. Available: <https://www.vldb.org/pvldb/vol15/p3385-xu.pdf>
- [22] M. Wu, X. Yi, H. Yu, Y. Liu, and Y. Wang, "Nebula graph: An open source distributed graph database," *CoRR*, vol. abs/2206.07278, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2206.07278>
- [23] D. Ongaro, "Consensus: bridging theory and practice," Ph.D. dissertation, Stanford University, USA, 2014.
- [24] L. Ma, D. V. Aken, A. Hefny, G. Mezerhane, A. Pavlo, and G. J. Gordon, "Query-based workload forecasting for self-driving database management systems," in *SIGMOD Conference*. ACM, 2018, pp. 631–645.
- [25] Y. Gao, X. Huang, X. Zhou, X. Gao, G. Li, and G. Chen, "Dbaugur: An adversarial-based trend forecasting system for diversified workloads," in *ICDE*. IEEE, 2023, pp. 27–39.
- [26] J. S. Vitter, "Random sampling with a reservoir," *ACM Trans. Math. Softw.*, vol. 11, no. 1, pp. 37–57, 1985.
- [27] B. Boots, G. J. Gordon, and A. Gretton, "Hilbert space embeddings of predictive state representations," in *UAI*. AUAI Press, 2013.
- [28] R. Lee, M. Zhou, C. Li, S. Hu, J. Teng, D. Li, and X. Zhang, "The art of balance: A rateupdb experience of building a CPU/GPU hybrid database product," *Proc. VLDB Endow.*, vol. 14, no. 12, pp. 2999–3013, 2021.
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, E. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [30] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [31] Y. Lu, X. Yu, L. Cao, and S. Madden, "Epoch-based commit and replication in distributed OLTP databases," *Proc. VLDB Endow.*, vol. 14, no. 5, pp. 743–756, 2021.
- [32] PingCAP, "Replica management in tidb: Pd control user guide," <https://docs.pingcap.com/tidb/stable/pd-control>.
- [33] Oceanbase, "Modify locality," <https://en.oceanbase.com/docs/common-oceanbase-database-1000000001105992>.
- [34] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears, "Benchmarking cloud serving systems with YCSB," in *SoCC*. ACM, 2010, pp. 143–154.
- [35] R. O. Nambiar, N. Wakou, A. Masland, P. Thawley, M. Lanken, F. Carman, and M. Majdalany, "Shaping the landscape of industry standard benchmarks: Contributions of the transaction processing performance council (TPC)," in *TPCTC*, ser. Lecture Notes in Computer Science, vol. 7144. Springer, 2011, pp. 1–9.
- [36] B. Vandiver, H. Balakrishnan, B. Liskov, and S. Madden, "Tolerating byzantine faults in transaction processing systems using commit barrier scheduling," in *SOSP*. ACM, 2007, pp. 59–72.
- [37] X. Zhou, X. Yu, G. Graefe, and M. Stonebraker, "Lotus: Scalable multi-partition transactions on single-threaded partitioned databases," *Proc. VLDB Endow.*, vol. 15, no. 11, pp. 2939–2952, 2022.
- [38] J. W. Stamos and F. Cristian, "A low-cost atomic commit protocol," in *SRDS*. IEEE Computer Society, 1990, pp. 66–75.
- [39] I. Zhang, N. K. Sharma, A. Szekeres, A. Krishnamurthy, and D. R. K. Ports, "Building consistent transactions with inconsistent replication," in *SOSP*. ACM, 2015, pp. 263–278.
- [40] T. Kraska, G. Pang, M. J. Franklin, S. Madden, and A. D. Fekete, "MDCC: multi-data center consistency," in *EuroSys*. ACM, 2013, pp. 113–126.
- [41] S. Maiyya, F. Nawab, D. Agrawal, and A. E. Abbadi, "Unifying consensus and atomic commitment for effective cloud data management," *Proc. VLDB Endow.*, vol. 12, no. 5, pp. 611–623, 2019.
- [42] Q. Zhang, J. Li, H. Zhao, Q. Xu, W. Lu, J. Xiao, F. Han, C. Yang, and X. Du, "Efficient distributed transaction processing in heterogeneous networks," *Proc. VLDB Endow.*, vol. 16, no. 6, pp. 1372–1385, 2023.
- [43] J. M. Faleiro and D. J. Abadi, "Rethinking serializable multiversion concurrency control," *Proc. VLDB Endow.*, vol. 8, no. 11, pp. 1190–1201, 2015.
- [44] A. Thomson, T. Diamond, S. Weng, K. Ren, P. Shao, and D. J. Abadi, "Calvin: fast distributed transactions for partitioned database systems," in *SIGMOD Conference*. ACM, 2012, pp. 1–12.
- [45] Y. Lu, X. Yu, L. Cao, and S. Madden, "Aria: A fast and practical deterministic OLTP database," *Proc. VLDB Endow.*, vol. 13, no. 11, pp. 2047–2060, 2020.
- [46] T. Qadah, S. Gupta, and M. Sadoghi, "Q-store: Distributed, multi-partition transactions via queue-oriented execution and communication," in *EDBT*. OpenProceedings.org, 2020, pp. 73–84.
- [47] D. Qin, A. D. Brown, and A. Goel, "Caracal: Contention management with deterministic concurrency control," in *SOSP*. ACM, 2021, pp. 180–194.
- [48] Z. Lai, H. Fan, W. Zhou, Z. Ma, X. Peng, F. Li, and E. Lo, "Knock out 2pc with practicality intact: a high-performance and general distributed transaction protocol," in *ICDE*. IEEE, 2023, pp. 2317–2331.
- [49] E. Zamanian, C. Binnig, T. Kraska, and T. Harris, "The end of a myth: Distributed transactions can scale," *CoRR*, vol. abs/1607.00655, 2016.
- [50] C. Binnig, A. Crotty, A. Galakatos, T. Kraska, and E. Zamanian, "The end of slow networks: It's time for a redesign," *Proc. VLDB Endow.*, vol. 9, no. 7, pp. 528–539, 2016.
- [51] A. Kalia, M. Kaminsky, and D. G. Andersen, "Design guidelines for high performance RDMA systems," *login Usenix Mag.*, vol. 41, no. 3, 2016.
- [52] X. Wei, Z. Dong, R. Chen, and H. Chen, "Deconstructing rdma-enabled distributed transactions: Hybrid is better!" in *OSDI*. USENIX Association, 2018, pp. 233–251.
- [53] D. Y. Yoon, M. Chowdhury, and B. Mozafari, "Distributed lock management with RDMA: decentralization without starvation," in *SIGMOD Conference*. ACM, 2018, pp. 1571–1586.
- [54] H. Chen, R. Chen, X. Wei, J. Shi, Y. Chen, Z. Wang, B. Zang, and H. Guan, "Fast in-memory transaction processing using RDMA and HTM," *ACM Trans. Comput. Syst.*, vol. 35, no. 1, pp. 3:1–3:37, 2017.