



# RCBench: an RDMA-enabled transaction framework for analyzing concurrency control algorithms

Hongyao Zhao<sup>1</sup> · Jingyao Li<sup>1</sup> · Wei Lu<sup>1</sup> · Qian Zhang<sup>1</sup> · Wanqing Yang<sup>1</sup> · Jiajia Zhong<sup>1</sup> · Meihui Zhang<sup>2</sup> · Haixiang Li<sup>3</sup> · Xiaoyong Du<sup>1</sup> · Anqun Pan<sup>3</sup>

Received: 21 October 2022 / Revised: 30 April 2023 / Accepted: 20 October 2023 / Published online: 14 December 2023  
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2023

## Abstract

Distributed transaction processing over the TCP/IP network suffers from the *weak transaction scalability* problem, i.e., its performance drops significantly when the number of involved data nodes per transaction increases. Although quite a few of works over the high-performance RDMA-capable network are proposed, they mainly focus on accelerating distributed transaction processing, rather than solving the weak transaction scalability problem. In this paper, we propose *RCBench*, an RDMA-enabled transaction framework, which serves as a unified evaluation tool for assessing the transaction scalability of various concurrency control algorithms. The usability and advancement of RCBench primarily come from the proposed concurrency control primitives, which facilitate the convenient implementation of RDMA-enabled concurrency control algorithms. Various optimization principles are proposed to ensure that concurrency control algorithms in RCBench can fully benefit from the advantages offered by RDMA-capable networks. We conduct extensive experiments to evaluate the scalability of mainstream concurrency control algorithms. The results show that by exploiting the capabilities of RDMA, concurrency control algorithms in RCBench can obtain 42X performance improvement, and transaction scalability can be achieved in RCBench.

**Keywords** Concurrency control · Distributed transaction · Transaction scalability · RDMA

## 1 Introduction

The capability to support distributed transaction processing in database systems is indispensable for many mission-critical applications, such as e-banking and e-commerce. However, it is generally believed that distributed transaction processing over TCP/IP networks cannot scale [66]. That is, the increasing number of involved data nodes per transaction makes the system performance drop significantly. This phenomenon is referred to as *weak transaction scalability*. To show this phenomenon, we implement a two-phase locking concurrency control algorithm (a.k.a. 2PL) on distributed framework Deneva [29] open-sourced by MIT, and conduct an evaluation over YCSB benchmark. We plot the throughput by varying the number of accessed data nodes per transaction in Fig. 1. The throughput of distributed transaction

✉ Wei Lu  
lu-wei@ruc.edu.cn

Hongyao Zhao  
hongyaozhao@ruc.edu.cn

Jingyao Li  
li-jingyao@ruc.edu.cn

Qian Zhang  
zhangqianzq@ruc.edu.cn

Wanqing Yang  
wanqingyang@ruc.edu.cn

Jiajia Zhong  
zhongjiajia@ruc.edu.cn

Meihui Zhang  
meihui\_zhang@bit.edu.cn

Haixiang Li  
blueseali@tencent.com

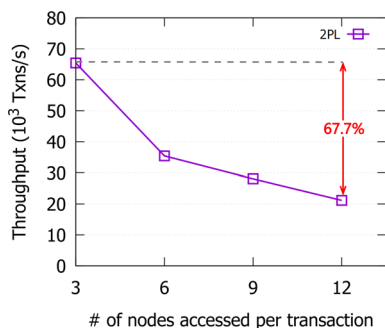
Xiaoyong Du  
duyong@ruc.edu.cn

Anqun Pan  
aaronpan@tencent.com

<sup>1</sup> Renmin University of China, Beijing, China

<sup>2</sup> Beijing Institute of Technology, Beijing, China

<sup>3</sup> Tencent Inc., Shenzhen, China



**Fig. 1** The throughput drops sharply when the number of accessed data nodes per transaction increases

processing decreases by a factor of 67.7% when the number of accessed data nodes per transaction increases from 2 to 5.

The reasons that cause the weak transaction scalability are two-fold. On the one hand, communications from the coordinator to the participants increase heavily as more accessed data nodes per transaction are involved, making the system performance drop; on the other hand, coordinating more nodes in each transaction would introduce additional overhead, bottlenecked by the slowest one in the distributed transaction processing. Interestingly, to the best of our knowledge, none of the existing works directly target solving weak transaction scalability. Instead, most of them target accelerating the performance of distributed transaction processing. For these works, they are divided into three categories. Approaches of the first category attempt to either eliminate distributed transactions or minimize their number as much as possible. For these proposals, various partitioning approaches and partition placement schemes [1, 2, 16, 21, 39, 44, 47, 50, 56, 67] are carefully designed so that partitions involved in the same transaction locate in the same data node. Yet, static partitioning works only if the optimal data placement is known a priori and never changes, while dynamic partitioning often suffers from an expensive data migration overhead. As opposed to the first category, approaches of the second category directly optimize the execution logic of distributed transactions. To do this, variants of 2PC are carefully designed to reduce the number of communications. For instance, Early Prepare [49] eliminates the prepare phase of 2PC, and hence, one round-trip from the coordinator to the participants is reduced; besides, deterministic concurrency control algorithms [24–26, 36, 37, 41, 43, 53] can completely eliminate 2PC, and hence, two round-trips from the coordinator to the participants are reduced; despite these optimizations, the master-worker architectures in these works are still bottlenecked by the high network latency and coordination overhead, which hinder the transaction scalability. To alleviate the network and coordination overhead in the second category, approaches [5, 14, 20, 60, 61, 64] of the third category adopt the shared-memory archi-

ture, in which nodes are connected via RDMA-capable networks, i.e., high-performance Infiniband networks with the remote direct memory access (a.k.a. RDMA) capability. As reported in [9], RDMA-capable networks provide relatively comparable latency to main memory but take significantly lower latency than TCP/IP networks. Thus, under the shared-memory architecture, a worker is scheduled to execute transactions entirely without the coordination of the coordinators by reading/writing remote data items directly via low-latency RDMA verbs. This idea is similar to that in centralized systems, as opposed to dividing each distributed transaction into multiple sub-transactions in the master-worker architectures.

In this paper, we aim to build a unified transaction framework over RDMA-capable networks with two requirements. Requirement R1: this framework must take full advantage of RDMA-capable networks. Requirement R1 helps verify whether re-implementations of concurrency control algorithms are transaction-scalable when fully utilizing the advantages of RDMA-capable networks. Requirement R2: it is generic to re-implement mainstream concurrency control algorithms in this framework. Requirement R2 facilitates achieving convenient re-implementations and a fair comparison among them. Thus far, although quite a few works are proposed for RDMA-capable networks, they do not satisfy either requirement R1 or requirement R2 or both. For instance, the extensions of 2PL [5, 14, 61, 64] and optimistic concurrency control (OCC) [20, 60] mainly focuses on locking/unlocking data items and validation, respectively. However, these works lack generality (requirement R2) because applying any of them individually is not enough to re-implement other concurrency control algorithms, e.g., such as choosing proper versions of data items in multi-version concurrency control algorithms (MVCC); additionally, some implementations, e.g., DrTM-H [60], do not take full advantages of RDMA-capable networks, which may potentially affect the evaluation of transaction scalability.

To satisfy the above two requirements, we propose a unified transaction framework called RCBench. To address requirement R1, we make a careful redesign of the data access methods for concurrency control algorithms completely using one-sided RDMA verbs, which take full advantage of RDMA-capable networks but come with the limitation of requiring prior knowledge of a data item's address before accessing it. To solve this problem, in RCBench, we carefully design a key-to-address index, with which, to access a data item, we first obtain its remote address based on its key, and then access the data item based on the remote address using one-sided RDMA verbs. Besides, existing optimizations for RDMA-capable networks, including coroutine, doorbell batching, outstanding requests, and passive ack are integrated into RCBench.

To satisfy requirement R2, we collect the metadata for each concurrency control algorithm and abstract six primitives. Every concurrency control algorithm can be conveniently re-implemented by invoking the primitives on its metadata, without directly touching the RDMA programming but enjoying all advantages of RDMA-capable networks. Because a concurrency control algorithm could be re-implemented using a different set of primitives, or the same set of primitives but with varying numbers of RDMA verb invocations, we propose five optimization principles for the re-implementations that aim to achieve transaction scalability by minimizing the number of RDMA verb invocations. Note that, our proposed primitives differ from those designed in FaRM [19] which aim to speed up message communication, while ours are used for manipulating remote metadata or data. Following these principles, we employ the primitives to re-implement multiple mainstream state-of-the-art concurrency control algorithms, including (1) widely-used protocols such as 2PL [27]: No-Wait [7], Wait-Die [45], Wound-Wait [45], T/O [6, 7, 45], and MVCC [62]; (2) modern protocols such as Silo [57], Maat [28], and Cicada [35]; (3) a deterministic protocol Calvin [53].

To answer the questions that whether or not the re-implementations are transaction-scalable, and which re-implementations achieve the best performance, we conduct comprehensive experiments over the widely used benchmarks. We report our findings below.

- In RCBench, it is convenient to re-implement the concurrency control algorithms using our proposed six primitives that enjoy all benefits of RDMA-capable networks. Interestingly, some results show that our re-implementations can even achieve a better performance than some customized implementations under the same settings.
- All re-implementations, except for Calvin, demonstrate significant performance improvement (ranging from 18X to 42X) against their counterparts over TCP/IP networks. The degree of improvement closely relies on the number of primitive calls. Particularly, Silo is reported to perform the best in most cases.
- Among all optimizations, the coroutine brings the greatest performance improvement (1.7X to 2.5X). Besides, the selection of appropriate lock types is also important. For example, using a single type of exclusive lock instead of exclusive/shared locks in 2PL can significantly improve performance in moderate or low contention scenarios.
- Our experimental results show that transaction scalability can be achieved. We must emphasize this finding is rather important because, in this way, it is unnecessary to do expensive data migration across data nodes in order to eliminate distributed transactions.

## 2 Preliminaries

In this section, we first review the distributed transaction processing and then discuss the background of RDMA techniques.

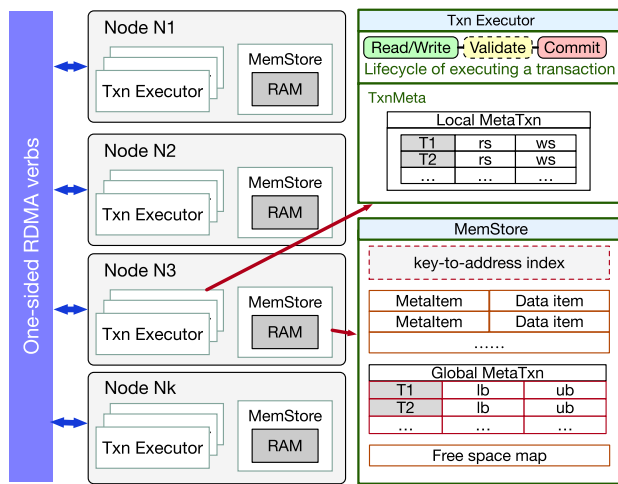
### 2.1 Distributed transaction processing

In shared-nothing database systems, data are horizontally partitioned, and data nodes are responsible for storing and accessing the partitions that are assigned to them. These systems adopt multi-coordinator system architecture to do distributed transaction processing, such as Percolator [40]. Specifically, each coordinator individually (1) accepts the transactions, (2) breaks the transaction into several sub-transactions and distributes sub-transactions to the appropriate data nodes, also known as participants, for execution, and (3) issues 2PC to coordinate the commit/abort of the transaction. In step (2), each data node is equipped with a database instance that manages the concurrent execution of the sub-transactions in this node. In step (3), once all sub-transactions decide to commit, the coordinator coordinates the commit of the transaction; otherwise, it coordinates the abort of the transaction.

### 2.2 RDMA

RDMA is a conceptual extension of direct memory access (DMA) technology and has become an increasingly popular technique to accelerate the performance of a system. It is capable of providing the following three properties. (1) **Zero-copy** property. Applications can perform data transfers without the involvement of the network software stack. Data are sent and received directly to the buffers without being copied between the network layers. (2) **Kernel bypass** property. Applications can perform data transfers directly from user space without kernel involvement. (3) **No CPU involvement** property. Applications can access remote memory without consuming any CPU cycles in the remote machine. Although RDMA is supported on both Ethernet and InfiniBand networks that provide the same common user APIs for programming, RDMA over InfiniBand networks provides much higher bandwidth and lower latency. Hence, in this paper, we focus on RDMA over InfiniBand networks.

It provides two categories of verbs for programming: (1) one-sided verbs, including READ, WRITE, WRITE With Immediate, and two atomic operations: FETCH-And-ADD (a.k.a. FAA) as well as COMPARE-And-SWAP (a.k.a. CAS), and (2) two-sided verbs, including SEND and RECEIVE. Programming using one-sided verbs enjoys all three proper-



**Fig. 2** An overview of RCBench

ties of RDMA. For example, the atomic RDMA CAS allows a machine to do compare-and-swap in the remote machine atomically without any intervention of the remote CPU. Nevertheless, programming using two-sided verbs can only have zero-copy and kernel-passing properties. That is, two-sided verbs still require CPU involvement of the remote machine. Therefore, a one-sided verb always seems more favorable than a two-sided verb.

Although one-sided verbs have all benefits of RDMA, in programming, we must know the memory address of the data to be read/written in the remote node a priori. That is, we may issue multiple one-sided verbs to obtain the memory address in the remote node and issue another verb to do the remote read/write using the remote address. On the contrary, we can read/write a remote data item by issuing a single two-sided verb without knowing its remote memory address. As reported in Sect. 7.3, we observe that 10–20 one-sided verb invocations are roughly equivalent to a single two-sided verb invocation, and system throughput decreases as the number of one-sided verb invocations increases. Therefore, in order to fully leverage the advantages of RDMA-capable networks, RCBench completely utilizes one-sided verb invocation to access remote data while reducing the number of one-sided verb invocations as much as possible.

### 3 Overview of RCBench

Figure 2 presents an overview of RCBench, which processes distributed transactions with the full benefits of RDMA-capable networks. We have released its source code on Github.<sup>1</sup> RCBench adopts the shared-memory architecture

**Table 1** Symbols and their meanings

Symbol	Meaning
$X$	A data item
$X^d$	The data value of $X$
$X.PK$	The primary key of $X$
$X^m$	The ItemMeta of $X$
$T$	A transaction
$T.Tid$	ID of $T$
$T^l(T^g)$	Local (global) metadata of $T$
$T^l.rs$	The read set of $T$
$T^l.ws$	The write set of $T$
$T^g.lock$	The lock used to prevent concurrent writes on $T^g$
$X^m.lock$	The lock used to prevent concurrent writes on $X^d$ and $X^m$

[9], in which data items are horizontally partitioned using a hash-based method. Each partition is assigned to one node and maintained in its main memory, and all operations on remote data items are executed via one-sided RDMA verbs. In RCBench, multiple clients generate transactions and send them to different servers in a round-robin fashion [29]. Each server acts as both a data node equipped with *MemStore* and a compute node equipped with multiple *Txn executors*. For ease of illustration, Table 1 summarizes the notations used throughout the paper.

#### 3.1 MemStore

MemStore is a pre-allocated, RDMA-registered memory that manages data items and provides auxiliary metadata to facilitate Txn executors in performing concurrency control. To facilitate the access of data items and metadata, we construct a key-to-address index on each node, where each index entry records the primary key and the main memory address of the corresponding data item residing in MemStore. Txn executors traverses the remote index to fetch data items using one-sided RDMA verbs. We will elaborate on the access of data items in Sect. 4.1.

We organize MemStore into several memory blocks in fixed-length, and maintain a bitmap *free\_space\_map* to reflect the memory usage. In *free\_space\_map*, we use 1 bit to indicate the status of a memory block: 1 means the block is used, while 0 means it is free. When inserting or deleting a data item, we update the corresponding bit in *free\_space\_map* atomically, which will be further discussed in Sect. 4.2.

<sup>1</sup> <https://github.com/dbiir/RCBench>.

Each data item  $X$  is stored in a fixed-length memory block, containing its corresponding metadata  $X^m$  and data value  $X^d$ . This allows remote Txn executors to access both  $X^m$  and  $X^d$  with a single one-sided RDMA verb. For metadata, we design the following two categories of metadata for performing concurrency control.

- **MetaItem**  $X^m$ .  $X^m$  represents the metadata corresponding to each data item  $X$ . For example, in 2PL,  $X^m$  includes the ID of a running transaction that has ever written  $X$ , or the ID list of running transactions that have ever read  $X$ , serving as the exclusive/shared lock identifiers. In the real implementation,  $X^m$  is set as a fixed-length data structure. By doing this, our access method guarantees that  $X^m$  and  $T^s$  can be fetched using a single one-sided RDMA verb if we have the address of  $X$ . Suppose in 2PL, we store the ID list of transactions in  $X^m$ . Because we cannot know how many transactions read each data item  $X$  a priori, we create  $X^m$  with a very large list to store them, which leads to a prohibitively expensive space overhead. As a compromise, we set a proper size of the list. When a transaction  $T$  attempts to acquire a shared lock on  $X$ , and the ID list of  $X^m$  is full, we simply abort  $T$  to ensure correctness.
- **MetaTxn**  $T^l/T^s$ . MetaTxn denotes metadata maintained by a transaction  $T$ , including local metadata  $T^l$  and global metadata  $T^s$ . Local metadata  $T^l$  can only be manipulated by the local Txn executors, so we store  $T^l$  in the local memory of Txn executors. In contrast, global metadata  $T^s$  can be accessed by remote Txn executors. For example, in 2PL,  $T^l$  maintains the transaction's read/write sets, and global MetaTxn  $T^s$  contains the transaction status, including running, aborted, or committed. We manage  $T^s$  similarly to data items and allocate fixed-length memory for each  $T^s$ , which will be presented in detail in Sect. 4.3.

### 3.2 Txn executor

Every Txn executor in our framework works like its counterpart in the centralized system except that the former is enriched with the capability to directly access the memory of a remote node. After receiving a transaction  $T$ , a Txn executor processes  $T$  through two or three phases: 1) read/write phase, 2) validation phase (if any), and 3) commit phase. In the read/write phase, the Txn executor checks the condition, e.g., acquire a lock using 2PL, to read a data item  $X$ ,

or/and possibly write  $X$ . Before this, it is necessary to fetch the address of  $X$  in the local/remote node through the key-to-address index. In the validation phase (if any), the Txn executor examines the conflicts between  $T$  and other concurrent transactions and determines whether or not  $T$  can commit or abort. If  $T$  aborts, any modifications by  $T$  need to be rolled back; otherwise, any modifications by  $T$  need to be persisted in MemStore.

Although it is general to abstract the execution of each transaction through two or three phases, concurrency control algorithms could have different logic to ensure transaction correctness. Obviously, individual re-implementations of them are inefficient for development, and more importantly, we attempt to make a fair comparison over these re-implementations. For these reasons, we first collect the metadata for each concurrency control algorithm and abstract six primitives. Every concurrency control algorithm can be conveniently re-implemented by invoking the primitives on its metadata, without directly touching the RDMA programming, but enjoying all advantages of RDMA-capable networks. Because a concurrency control algorithm can be re-implemented with a different set of primitives, or with the same set of primitives but having a significantly different number of RDMA verb invocations, aiming to achieve transaction scalability, we propose five optimization principles for the re-implementations that are able to reduce the number of RDMA verb invocations.

To further improve performance, we provide two transaction processing modes for Txn executor, including the thread-to-transaction mode and coroutine-to-transaction mode. In the first mode, a thread is created for one executor to execute transactions sequentially; in the second mode, a thread is composed of multiple coroutines, each of which is created for one executor to execute transactions sequentially. By using finer-grained scheduling in the second mode, when a coroutine is blocked, the transaction execution can be switched to another coroutine of the same thread. The scheduling overhead of coroutines is significantly smaller than that of threads, and hence, coroutine-to-transaction mode can achieve better performance.

## 4 The access method

In this section, we present the access and maintenance method of remote data items and global metadata completely using one-sided RDMA verbs.

#### 4.1 The access to remote data items

As mentioned in Sect. 3, we develop an RDMA-friendly key-to-address hash index  $IDX$  to manage data items on each node. In general, reading or writing a data item is accomplished using two one-sided RDMA verbs. When accessing a data item, given its primary key, we first retrieve the corresponding index entry using a one-sided RDMA verb and then obtain the data item based on the address stored in the index entry with another one-sided RDMA verb. However, additional one-sided RDMA verbs for fetching an index entry are required when hash collisions happen, resulting in increased network overhead. To address this issue, we adopt the empirical methods from [38] and implement the hash index  $IDX$  as an  $n$ -way Cuckoo hash table. With this hash scheme, we apply  $n$  orthogonal hash functions, denoted as  $hash_k$  ( $k = 1, 2, \dots, n$ ), to assign  $n$  locations for a primary key, indicating that every key is either at one of  $n$  possible locations or absent. By default, we set  $n = 3$ , referring to the optimal setting indicated in [38]. By iteratively examining these  $n$  possible locations, we guarantee that an index entry can be located using  $\leq n$  single one-side verbs. We elaborate on how to fetch an index entry through  $IDX$  in detail:

In our design, we store an index entry as a triple  $\hat{X}$  of  $\langle X.addr, X.size, X.PK \rangle$  occupying 24 bytes, where  $X.addr$ ,  $X.size$ , and  $X.PK$  are the memory address, size, and primary key of the data item  $X$ , respectively. Besides, each index entry maintains a *lock*, occupying 8 bytes, to prevent concurrent manipulation on  $IDX$ . By so doing, we formulate the function *getRtItemAddr* to fetch the appropriate index entry  $\hat{X}$  of a given key  $PK$  below: (1) we calculate a location using  $hash_k(X.PK)$  and issue an RDMA READ to fetch its corresponding  $\hat{X}$ ; (2) we then return  $\hat{X}$  if  $\hat{X}.X.PK = PK$  and  $\hat{X}$  is not locked by another transaction. If  $\hat{X}$  has been locked by another transaction, we wait until the lock is released. If  $\hat{X}.X.PK \neq PK$ , we set  $k = k + 1$  and re-execute (1)(2) until either finding a correct  $\hat{X}$  with  $\hat{X}.X.PK = PK$ , or retrying until  $k = n - 1$ .

To support range queries, we adopt the approach proposed in [30]. Specifically, we horizontally partition index entries into several hash indexes based on the prefix of the stored primary key. When querying data items with primary key ranges from  $PK_i$  to  $PK_j$ , we traverse hash indexes that exhibit intersection with the requested range to obtain required index entries. To prevent phantom read anomalies, we maintain an extra lock meta  $IDX.lock$  for each hash index. Before/After the range query, shared locks are acquired/released on involved hash indexes.

#### 4.2 The manipulation of remote data items

We elaborate on the manipulations of remote data items in terms of insert and delete. Based on the design of  $IDX$ ,

inserting a new index entry  $\langle X.addr, X.size, PK \rangle$  of  $X$  requires the following procedures. (1) First,  $n$  RDMA CAS invocations are issued to concurrently lock  $n$  potential locations with addresses that are calculated by  $hash_k(X.PK)$  for  $k = 1, 2, \dots, n$ . In case of any failure, additional RDMA CAS are required to retry locking until success. (2) Then, we issue  $n$  concurrent RDMA READ invocations to obtain these locations, and examine them locally. (3) If any free location is available, an RDMA WRITE invocation is used to write the new index entry into the corresponding location. (4) Otherwise, insertion with preemption is executed by invoking an RDMA WRITE to kick out and replace one of the existing locations with the new one. (5) Finally, we release all prior locks with  $n$  RDMA CAS invocations. Note that, if preemption exists in this procedure, similar steps are repeated to insert the kicked-out entry. To prevent successive kicks, the hash table will be resized to accommodate more data if the number of kicks reaches a pre-defined limit [38].

Before insertion of the index entry, a free memory space needs to be found for the data item  $X$  in the remote MemStore. To do this, we issue an RDMA READ to fetch *free\_space\_map*, identify one of the available memory spaces based on it, and invoke an RDMA CAS to reset the corresponding bit to 1 before writing  $X$ .

When deleting a data item  $X$ , we (1) issue  $n$  RDMA READ to obtain index entries with addresses calculated by  $hash_k(X.PK)$  for  $k = 1, 2, \dots, n$ . (2) use an RDMA CAS to lock the target entry  $\hat{X}$ , and invoke another RDMA READ afterward to confirm that  $\hat{X}$  has not been modified by any other transaction. (3) clean  $\hat{X}$  from  $IDX$  with an RDMA WRITE, and release the memory space of  $X$  by issuing an RDMA CAS to update the corresponding bit in *free\_space\_map* from 1 to 0.

In contrast to the shared locks for range queries in Sect. 4.1, we acquire/release exclusive locks on hash indexes before/after insert or delete operations to prevent the phantom read. We will present the implementation of shared and exclusive locks in Sect. 6.2.1.

#### 4.3 The access to global metadata

We manage the global metadata  $T^g$  in the same manner as that of data items, which occupies a fixed-length  $\Theta$  of memory space in MemStore. We use an  $n$ -way cuckoo hash table to store the index entries of global MetaTxn  $T^g$ . With this hash table, the access, insert and delete procedures for a given  $T^g$  are basically the same as that of data items described in Sects. 4.1 and 4.2. Each index entry is stored as a tuple  $\hat{T}$  of  $\langle T^g.addr, T.Tid \rangle$ , where  $T.Tid$  is taken as an integer. Besides, each index entry maintains a *lock*, occupying 8 bytes, to prevent concurrent manipulations. For reference, we name the function used to obtain the appropriate index

entry  $\hat{T}$  of a given  $Tid$  as  $getRtTgAddr$ , and omit its implementation as it is self-explanatory.

### 5 Concurrency control primitives

In this section, we first formulate the operation logic of concurrency control algorithms in centralized database systems, then abstract six primitives that take full benefits of RDMA-capable networks, and finally extend the operations based on the primitives to facilitate the re-implementations of concurrency control algorithms.

#### 5.1 Operation abstraction in centralized systems

A transaction can be modeled as a sequence of read/write operations, ended with a commit/abort operation. In the centralized database system, upon any read/write from/to data item  $X^d$ , the concurrency control algorithms need to examine whether the transaction has the qualification to read/write  $X^d$  by acquiring a lock on  $X^d$  in 2PL, or examine the timestamp on  $X^d$  in T/O, or others. To verify whether a transaction can commit or needs to abort, the concurrency control algorithms might need to do the validation, by either checking the conflict between its read set and the write set of concurrent transactions, or adjusting its timestamp interval based on its concurrent transactions, or others. To do the commit/abort, the concurrency control algorithms might need to release the locks, and persist the write set of the transaction. For reference, we list the operations that are necessary for some classic concurrency control algorithms in Table 2.

The execution of all operations in Table 2 can be formulated into three steps, shown in Algorithm 1: (1) fetch data item  $X$  (line 1), (2) perform the operation logic based on  $X$  (line 2), and (3) update  $X$  if necessary (line 3). We take a lock request on a data item  $X$  in lock-based algorithms for an example to illustrate these three steps. To acquire a lock on  $X$ , we (1) first fetch  $X^m$ , (2) then examine whether the lock meta  $X^m.lock$  in  $X^m$  has been modified by other transactions; (3) and finally update  $X^m$  by setting  $X^m.lock$  to indicate that  $X$  has been locked by it if  $X^m.lock$  has not been modified by other transactions.

---

**Algorithm 1:** Operation logic abstraction in centralized database systems

---

- 1 Fetch  $X^d$  and/or  $X^m$ ;
- 2 Perform operation logic based on  $X^d$  and/or  $X^m$ ;
- 3 Update  $X^d$  and/or  $X^m$  if necessary.

---

**Table 2** Operations of classic algorithms. Each class of algorithms includes operations involving marking with a ✓

Phase	Operations	Lock-based algorithms	OCC algorithms	Timestamp-based algorithms	MVCC-based algorithms
Read/Write	Read	✓	✓	✓	✓
	Write	✓	✓	✓	✓
	Lock	✓			
Validation	Timestamp-Examination			✓	
	Version-Retrieval				✓
	Read/Write-Set-Validation		✓		
Commit	Timestamp-Interval-Adjust		✓		
	Lock		✓		
	Persist-Data	✓	✓	✓	✓
	Roll-Back-Modification	✓	✓	✓	✓
	Unlock	✓	✓		

## 5.2 Primitive abstraction in RDMA-capable networks

To extend the operations of concurrency control algorithms in RDMA-capable networks, and make the underlying RDMA programming transparent to developers, we abstract six primitives, which are  $Read^D$ ,  $Write^D$ ,  $Atomic^D$  for data items and  $Read^T$ ,  $Write^T$ ,  $Atomic^T$  for transaction's meta items.

- $Read^D$  (Primitive 1) is used to read remote  $X$ . In each node, we maintain the address of  $IDX$  of every data node. By taking  $PK$  as input,  $Read^D$  fetches  $\hat{X}$  by issuing  $getRtItemAddr$  function (line 1). If  $\hat{X} = NULL$ ,  $Read^D$  fails and returns  $NULL$  since  $\hat{X}$  does not exist (line 2); otherwise,  $Read^D$  issues another RDMA READ to read and return  $X$  via  $\hat{X}$  (line 3).

---

### Primitive 1: $Read^D(PK)$

---

```

1  $\hat{X} \leftarrow getRtItemAddr(PK)$ ;
2 if  $\hat{X} = NULL$  then return  $NULL$ ;
3 return  $RDMA\_READ(\hat{X}.X.addr, \hat{X}.X.size)$ 

```

---

- $Write^D$  (Primitive 2) is designed to overwrite remote  $X$ . Similarly,  $Write^D$  issues  $getRtItemAddr$  function to fetch  $\hat{X}$  (lines 1–2). If  $\hat{X}$  exists, then  $Write^D$  issues another RDMA WRITE that writes the new value  $newV$  to  $X$  via  $\hat{X}$  (line 3).

---

### Primitive 2: $Write^D(PK, newV)$

---

```

1  $\hat{X} \leftarrow getRtItemAddr(PK)$ ;
2 if  $\hat{X} = NULL$  then return  $false$ ;
3  $RDMA\_WRITE(\hat{X}.X.addr, newV, \hat{X}.X.size)$ ;
4 return  $true$ ;

```

---

- $Atomic^D$  (Primitive 3) is designed to conditionally update an item (specified by input parameter  $meta$ ) of  $X^m$ , e.g., the lock of  $X$  ( $X^m.lock$  with  $meta = MT\_LOCK$ ) or the read timestamp of  $X$  ( $X^m.rts$  with  $meta = MT\_RTS$ ), with the atomicity guarantee. Given  $PK$ ,  $Atomic^D$  issues  $getRtItemAddr$  to fetch  $\hat{X}$  (line 1). If  $\hat{X}$  exists,  $Atomic^D$  calculates the remote address  $X.meta\_addr$  of  $meta$  by adding  $meta$ 's offset in  $X^m$  to  $\hat{X}.X.addr$  (line 3). and issues RDMA CAS to conditionally update  $meta$  by new value  $newV$  with atomicity guarantee (line 4).

- $Read^T$  (Primitive 4) is designed to read remote  $T^g$ . Taking  $Tid$  as the input,  $Read^T$  issues  $getRtTgAddr$  function to fetch  $\hat{T}$  (line 1). If  $\hat{T} = NULL$ ,  $Read^D$  fails and returns  $NULL$  since  $\hat{T}$  does not exist (line 2); otherwise,  $Read^T$  issues an RDMA READ to read  $T^g$  via  $\hat{T}.T^g.addr$  and  $\Theta$ , the length of data structure  $T^g$ , and return  $T^g$  (line 3).

---

### Primitive 3: $Atomic^D(PK, meta, oldV, newV)$

---

```

1  $\hat{X} \leftarrow getRtItemAddr(PK)$ ;
2 if  $\hat{X} = NULL$  then return  $false$ ;
3  $X.meta\_addr \leftarrow \hat{X}.X.addr + meta.offset$ ;
4  $t \leftarrow RDMA\_CAS(X.meta\_addr, oldV, newV)$ ;
5 return  $t = oldV$ ;

```

---



---

### Primitive 4: $Read^T(Tid)$

---

```

1  $\hat{T} \leftarrow getRtTgAddr(Tid)$ ;
2 if  $\hat{T} = NULL$  then return  $NULL$ ;
3 return  $RDMA\_READ(\hat{T}.T^g.addr, \Theta)$ 

```

---

- $Write^T$  (Primitive 5) is designed to write remote  $T^g$ .  $Write^T$  first issues  $getRtTgAddr$  function to fetch  $\hat{T}$  (lines 1–2), and issues an RDMA WRITE to overwrite  $T^g$  with  $newV$  (line 3).

---

### Primitive 5: $Write^T(Tid, newV)$

---

```

1  $\hat{T} \leftarrow getRtTgAddr(Tid)$ ;
2 if  $\hat{T} = NULL$  then return  $false$ ;
3  $RDMA\_WRITE(\hat{T}.T^g.addr, newV, \Theta)$ ;
4 return  $true$ ;

```

---

- $Atomic^T$  (Primitive 6) is designed to conditionally update an item (specified by the input parameter  $meta$ ) of  $T^g$ , e.g.,  $T^g.lock$ ,  $T^g.st$  with atomicity guarantee.  $Atomic^T$  fetches  $\hat{T}$  by issuing function  $getRtTgAddr$  (lines 1–2), calculates  $T^g.meta\_addr$  of  $meta$  by adding  $meta$ 's offset in  $T^g$  to  $\hat{T}.T^g.addr$  locally (line 3) and issues an RDMA CAS to conditionally update  $meta$  by new value  $newV$  with atomicity guarantee (line 4).

---

### Primitive 6: $Atomic^T(Tid, meta, oldV, newV)$

---

```

1  $\hat{T} \leftarrow getRtTgAddr(Tid)$ ;
2 if  $\hat{T} = NULL$  then return  $false$ ;
3  $T^g.meta\_addr \leftarrow \hat{T}.T^g.addr + meta.offset$ ;
4  $t \leftarrow RDMA\_CAS(T^g.meta\_addr, oldV, newV)$ ;
5 return  $t = oldV$ 

```

---

## 5.3 Operation extension in RDMA-capable networks

With the above six primitives that take full benefits of RDMA-capable networks, we extend the operations of concurrency control algorithms in centralized database systems to RDMA-capable networks. As for the access to remote data items, we propose RDMA-BasicD shown in Algorithm 2 as



the extension. RDMA-BasicD takes the primary key  $PK$  of a data item and current transaction  $T$  as the input. It (1) first issues  $Read^D$  to read remote data items (line 4), (2) then performs local logic based on  $X$  (line 7), and (3) finally modifies  $X$  using  $Write^D$  (line 9). For example, to acquire a lock on remote data item  $X$  in lock-based algorithms, we first use  $Read^D$  to fetch  $X^m$ , then locally examine whether the lock meta  $X^m.lock$  has been modified by other transactions, and finally update  $X^m$  with  $X^m.lock$  using  $Write^D$  to indicate that  $X$  has been locked if  $X^m.lock$  has not been modified by other transactions. In our implementation, we acquire a latch to ensure the atomicity of the three steps using  $Atomic^D$ . As shown in Algorithm 2, we issue an  $Atomic^D$  to acquire the latch of data item  $X$  (lines 2,3), and set  $X^m.latch$  to 0 to release this latch (lines 8,9). As for the access to remote global transaction metadata  $T^g$ , we further propose RDMA-BasicT shown in Algorithm 2 as the extension. We omit the details of RDMA-BasicT which follow the same logic as RDMA-BasicD.

Overall, RDMA-BasicD and RDMA-BasicT can be used to implement all remote operations of mainstream concurrency control algorithms. We take lock-based algorithms as an example, their operations, listed in Table 2, can be implemented using the same three steps outlined in the lock acquisition logic. The implementation of 2PL using RDMA-BasicD and RDMA-BasicT is described in detail in our technical report.<sup>2</sup> In conclusion, RDMA-BasicD can be used to implement timestamp-examination, version-retrieval, and read/write-set validation operations, and RDMA-BasicT can be used to implement timestamp-interval-adjust operations, covering all the operations mentioned in Table 2.

## 6 Design principles and Re-implementations of concurrency control algorithms

### 6.1 Optimization principles

As discussed, remote operations of various algorithms can be implemented with RDMA-BasicD and RDMA-BasicT conveniently. However, since RDMA-BasicD and RDMA-BasicT are general operations feasible for implementing all the concurrency control algorithms, further optimizations are required to fit the characteristics of different algorithms. We introduce five optimization principles for RDMA-BasicD and RDMA-BasicT, which either reduce the number of one-sided verb invocations, or eliminate explicit latch acquisition. By deeply customizing the re-implementations of concurrency control algorithms using these principles, we fully leverage the advantages of RDMA-capable networks to achieve transaction scalability.

<sup>2</sup> <https://github.com/dbiir/RCBench/blob/master/RCBench.pdf>.

### Algorithm 2: Operation extension in RDMA-capable networks

```

1 Function RDMA-BasicD( $PK, T$ ):
2    $r \leftarrow Atomic^D(PK, MT\_LATCH, 0, T.Tid)$ ;
3   if  $\neg r$  then Abort  $T$ ;
4    $X \leftarrow Read^D(PK)$ ;
5   if  $X = NULL$  then
6      $\lfloor Atomic^D(PK, MT\_LATCH, T.Tid, 0); Abort\ T$ ;
7   Perform local logic based on  $X$ ;
8    $\hat{X}^m.latch \leftarrow 0$ ;
9    $Write^D(PK, X)$ ;

10 Function RDMA-BasicT( $Tid_i, T$ ):
11   $r \leftarrow Atomic^T(Tid_i, MT\_LATCH, 0, T.Tid)$ ;
12  if  $\neg r$  then Abort  $T$ ;
13   $T_i^g \leftarrow Read^T(Tid_i)$ ;
14  if  $T_i^g = NULL$  then
15     $\lfloor Atomic^T(Tid_i, MT\_LATCH, T.Tid, 0); Abort\ T$ ;
16  Perform local logic based on  $T_i^g$ ;
17   $\hat{X}^m.latch \leftarrow 0$ ;
18   $Write^T(Tid_i, T_i^g)$ ;

```

First, as discussed in Sect. 2.2, **reducing the number of one-sided verb invocations can improve performance.**

As shown in Algorithm 2, both BasicD and BasicT require three primitives to complete a remote operation. However, depending on the characteristics of the algorithm, we find that not all three primitives are necessary. For example, the lock acquisition operation in lock-based algorithms can be implemented using a single  $Atomic^D$ . By eliminating primitives in RDMA-BasicD based on the requirements of various algorithms, we have designed four optimization principles: One-Cas, One-Write, One-Read, which uses only one  $Atomic^D$ ,  $Write^D$ ,  $Read^D$  in RDMA-BasicD; and Read-Cas, which uses a  $Read^D$  and a  $Atomic^D$  in RDMA-BasicD. These four principles are also applicable to RDMA-BasicT. We provide an example of One-Cas, Read-Cas, and One-Write in Sect. 6.2.1, and an example of One-Read in Sect. 6.2.2, respectively.

- **One-Cas.** If an operation only modifies metadata with a maximum of 8 bytes and knows the original and new values in advance, we can use a single  $Atomic^D$  ( $Atomic^T$ ) to modify this metadata.
- **One-Write.** If a transaction  $T$  has acquired a lock or latch for the target metadata,  $T$  can directly use a  $Write^D$  ( $Write^T$ ) to write back the metadata modified.
- **One-Read.** If there is no need to modify the remote metadata, a single  $Read^D$  ( $Read^T$ ) would be enough to read it.
- **Read-Cas.** If an operation only modifies the metadata with a maximum of 8 bytes, but its new value needs to be calculated against the old value, we can issue one

$Read^D(Read^T)$  to read back the metadata and another  $Atomic^D(Atomic^T)$  to update it.

Second, **explicit latch acquisition may reduce concurrency when accessing data items, incurring higher abort rate**. For example, the read operation in the T/O algorithm can be performed without the explicit latch to improve concurrency, as described in Sect. 6.2.3. To eliminate latch acquisition, we design the Double-Read principle as follows.

- **Double-Read.** Double-Read is an alternative for atomicity guarantee other than explicit latches. Specifically, it issues two  $Read^D(Read^T)$  before and after a potential modification of the target metadata using an  $Atomic^D(Atomic^T)$ . If the content re-read is different from the first one in an unexpected way, a concurrent modification may have occurred. Note that the Double-Read principle does not have the inherent capability to prevent the ABA problem. Thus, when employing this principle, users must ensure that the data information modified using  $Atomic^D(Atomic^T)$  does not encounter the ABA issue.

## 6.2 Re-implementations

Based on our proposed primitives and optimization principles, we manage to re-implement mainstream algorithms by minimizing the number of one-sided verbs and reducing explicit latch acquisition. To illustrate the detailed re-implementation methods, we give examples of No-Wait, Silo, and T/O algorithms in this section, and place the complete descriptions of other re-implementations, including Wait-Die, Wound-Wait, MVCC, MaaT, and Cicada, in the technical report. The operations using the proposed principles are highlighted in pink.

### 6.2.1 No-wait

**No-Wait** [7] is a variant of 2PL concurrency control algorithm. For any data item  $X$ , it always tries to acquire a certain type of lock on  $X$  before any read or write on it and aborts immediately in case of locking failures to avoid deadlocks. Specifically, in the read/write phase of No Wait, for each read of  $X$ , a shared lock is acquired on  $X$ , while for each write of  $X$ , an exclusive lock of  $X$  is acquired. In the Commit phase, we update the values of the data items that need to be modified, and release all the locks.

To boost No-Wait using RDMA, it is necessary to re-implement its logic that (1) acquire remote exclusive/shared locks, (2) perform remote reads/writes, and (3) release remote exclusive/shared locks. This involves creating a new lock metadata,  $X^m.lock$ , in the metadata of each data item  $X^m$ . The  $X^m.lock$  is a 64-bit value, where the least significant bit

is used to store the  $lock\_type$  (0 for shared lock or SL, and 1 for exclusive lock or EL), and the remaining 63 bits store the number of shared locks held on  $X$  (i.e.,  $num\_of\_locks$ ). By doing this, we can acquire an exclusive lock on  $X$  by updating  $X^m.lock$  from 0 to  $EL$ , satisfying the requirements of the One-Cas. And we can acquire a shared lock by adding 1 into  $X^m.lock.num\_of\_locks$ , satisfying the requirements of the Read-Cas. As transaction  $T$  has already acquired locks for data items when it enters the commit phase, the One-Write principle is applicable for the commit phase. Besides,  $MetaTxn T^l$  maintains the read set  $T^l.rs$  and write set  $T^l.ws$  of transaction  $T$ . Algorithm 3 exhibits the key functions of the re-implementation RDMA-No-Wait.

Function  $AcqIMExcLock$  and  $RlsIMExcLock$  are used to acquire/release an exclusive lock on  $X$ , which apply One-Cas principle. We acquire/release the exclusive lock on  $X$  by modifying the remote  $X^m.lock$  from 0/ $EL$  to  $EL$ /0 through only one  $Atomic^D$  (lines 2, 5).

Function  $AcqIMShLock$  uses the Read-Cas optimization principle to acquire a shared lock. By taking its primary key  $PK$  as the input, we first issue a  $Read^D$  to obtain  $X^m.lock$  (line 7) and check whether there is an exclusive lock on  $X$  locally (line 8). If there is, it fails to acquire the lock; otherwise, we make a local copy  $newL$  of  $X^m.lock$ , and update  $newL$  by recording a new shared lock on  $X$  (lines 11–12); we then issue an  $Atomic^D$  to update  $X^m.lock$  with  $newL$  to declare that a new shared lock on  $X$  is granted (line 13). Note that,  $Atomic^D$  compares the remote  $X^m.lock$  with local  $X^m.lock$  and modifies remote  $X^m.lock$  to  $newL$  only if they are the same; that is, if any changes to  $X$  have been made between  $ReadDI$  (line 7) and  $Atomic^D$  (line 13), the remote write would fail, and the returned value  $r$  is set to be false (line 14). Following the reverse logic of  $AcqIMShLock$ , function  $RlsIMShLock$  also uses this principle to release remote shared locks (lines 15–20).

Function  $Commit$  is used to commit transaction  $T$  and uses the One-Write principle. For each data item read before, we sequentially release the shared locks that  $T$  has acquired (lines 30–31) using  $RlsIMShLock$ . For each data item  $X$  in the write set  $T^l.ws$ , we update  $X^d$  with its new value, set  $x^m.lock$  to 0 locally, and issue a  $Write^D$  to overwrite the remote  $X$  (line 33).

Functions  $Read$  and  $Write$  are used to do remote reads and writes, respectively. For  $Read$  or  $Write$ , we first try to acquire the lock on  $X$  with  $X.PK = PK$ , then read the remote  $X$  if the lock acquisition success, and add  $X$  to the local read/write set  $T^l.rs/T^l.ws$ . Note, for  $Write$ , we need to update  $X^d$  locally before adding it into the local write set  $T^l.ws$  (line 27).

**Discussion.** As shown in Algorithm 3, it requires at least two primitive calls ( $Read^D$  and  $Atomic^D$ ) to acquire/release a shared lock while it requires only one primitive call ( $Atomic^D$ ) to acquire/release an exclusive lock. Based on

**Algorithm 3: RDMA-No-Wait**


---

```

1 Function AcqIMExcLock (PK) :
2    $r \leftarrow \text{Atomic}^D(PK, MT\_LOCK, 0, EL)$ ;
3   return  $r$ ;
4 Function RlsIMExcLock (PK) :
5    $\text{Atomic}^D(PK, MT\_LOCK, EL, 0)$ ;
6 Function AcqIMShLock (PK) :
7    $X \leftarrow \text{Read}^D(PK)$ ;
8   if  $X^m.lock.lock\_type = EL$  then
9     return false;
10  else
11     $newL \leftarrow X^m.lock$ ;
12     $newL.num\_of\_locks++$ ;
13     $r \leftarrow \text{Atomic}^D(PK, MT\_LOCK, X^m.lock, newL)$ ;
14    return  $r$ ;
15 Function RlsIMShLock (PK) :
16    $X \leftarrow \text{Read}^D(PK)$ ;
17    $newL \leftarrow X^m.lock$ ;
18    $newL.num\_of\_locks--$ ;
19    $r \leftarrow \text{Atomic}^D(PK, MT\_LOCK, X^m.lock, newL)$ ;
20   if  $\neg r$  then goto line 11;
21 Function Read (PK, T) :
22   if  $\neg \text{AcqIMShLock}(PK)$  then Abort  $T$ ;
23    $X \leftarrow \text{Read}^D(PK)$ ;
24    $T^l.rs \leftarrow \{X\} \cup T^l.rs$ ;
25 Function Write (PK, newV, T) :
26   if  $\neg \text{AcqIMExcLock}(PK)$  then Abort  $T$ ;
27    $X \leftarrow \text{Read}^D(PK)$ ;  $X^d \leftarrow newV$ ;
28    $T^l.ws \leftarrow \{X\} \cup T^l.ws$ ;
29 Function Commit (T) :
30   foreach  $X \in T^l.rs$  do
31      $\text{RlsIMShLock}(X.PK)$ ;
32   foreach  $X \in T^l.ws$  do
33      $X^m.lock \leftarrow 0$ ;  $\text{Write}^D(X.PK, X)$ ;

```

---

this observation, for the low-conflict application scenarios where reads/writes on the same data item rarely occur, using a single exclusive lock instead of exclusive/shared locks could potentially bring performance benefits by reducing extra remote primitive calls. To verify the benefits of using the single exclusive locking mechanism, we make an extensive experimental evaluation in Sect. 7.6, and the result shows that single exclusive locking can outperform exclusive/shared locking by a factor of 1.3X in the low-conflict application scenarios. *For this reason, in this paper, we adopt the single exclusive locking mechanism instead of the exclusive/shared locking mechanism by default.* In the real implementation, we collectively use *AcqIMExcLock/RlsIMExcLock* instead of *AcqIMShLock / RlsIMShLock* in Algorithm 3.

**6.2.2 Silo**

**Silo** [57] is a classic optimistic concurrency control algorithm. In Silo, each transaction  $T$  is scheduled to execute through three phases: read/write phase, validation phase and commit/abort phase. In the read/write phase, for each read of  $X$ , we store  $X$  to read set  $T^l.rs$ ; for each write of  $X$ , we store  $X$  to write set  $T^l.ws$ . In the validation phase,  $\forall X \in T^l.ws$ , it needs to acquire an exclusive lock on  $X$  and check whether  $X$  has been modified by other transactions. For  $\forall X \in T^l.rs$ , it checks whether  $X$  has been modified as well. If  $T$  cannot acquire all the locks successfully or if there exists a data item that has been modified by other transactions, we abort  $T$  and release all the locks held by  $T$ ; otherwise, we commit  $T$ , and  $\forall X \in T^l.ws$ , we accordingly update  $X^d$  and  $X^m.wts$ .

To boost Silo using RDMA, we re-implement its logic that (1) do remote reads/writes, (2) acquire remote locks and validate data items in the write set, and (3) validate data items in the read set. For the logic of (1) and (3), they satisfy the requirement of the One-Read principles as they do not modify metadata. In logic (2), Silo uses only the exclusive lock mechanism, satisfying the requirement of the One-Cas principle. To optimize the commit phase in the Silo algorithm, we can use the One-Write principle, as in the No-Wait algorithm.

To discuss the implementation of Silo in detail, we first design the metadata  $X^m$  in Silo, including two fields: (1) lock metadata  $X^m.lockb$ , recording the ID of a transaction that is currently granted with an exclusive lock on  $X$ ; and (2)  $X^m.wts$ , the maximum commit timestamp of transactions that have ever written  $X$ . Besides  $T^l.rs$  and  $T^l.ws$ ,  $T^l$  includes the commit timestamp ( $T^l.cts$ ) of transactions  $T$ . Then we discussed the key functions of the re-implementation called RDMA-Silo in Algorithm 4.

Functions *Read* and *Write* are used to perform the remote read/write on data item  $X$  by using One-Read principle. For *Read*, we read remote  $X$  by using a  $\text{Read}^D$ , and add it to the local read set  $T^l.rs$  (line 2). For *Write*, we read remote  $X$  by using a  $\text{Read}^D$ , update  $X^d$  locally, and add  $X$  to local write set  $T^l.ws$  (lines 4–5).

Function *Validation* is used to perform the operation in the validation phase. First, we obtain a commit timestamp of  $T$  locally. Then,  $\forall X \in T^l.ws$ , we try to acquire an exclusive lock on  $X$  following the One-Cas principle (line 9). If the lock acquisition fails,  $T$  aborts; otherwise, we continue to validate  $T$ . Subsequently,  $\forall X \in T^l.ws \cup T^l.rs$ , we examine whether  $X$  has been modified by other transactions based on a One-Read principle to read  $X$  again by using  $\text{Read}^D$  (line 10–11). If  $X$  has been locked or modified, we abort  $T$  (lines 12–13).

Upon commit of  $T$ ,  $\forall X \in T^l.ws$ , we make a single remote write by issuing  $\text{Write}^D$  to release the lock on  $X$ , as well as

**Algorithm 4: RDMA-Silo**


---

```

1 Function Read( $PK, T$ ):
2    $X \leftarrow \text{Read}^D(PK)$ ;  $T^l.rs \leftarrow \{X\} \cup T^l.rs$ ;
3 Function Write( $PK, T, newV$ ):
4    $X \leftarrow \text{Read}^D(PK)$ ;  $X^d \leftarrow newV$ ;
5    $T^l.ws \leftarrow \{X\} \cup T^l.ws$ ;
6 Function Validation( $T$ ):
7    $T^l.cts \leftarrow$  get the current timestamp;
8   foreach  $X \in T^l.ws$  do
9     if  $\neg \text{Atomic}^D(PK, MT\_LOCKB, 0, T.Tid)$  then
10      Abort  $T$ ;
11   foreach  $X \in T^l.rs \cup T^l.ws$  do
12      $\bar{X} \leftarrow \text{Read}^D(X.PK)$ ;
13     if  $\bar{X}^m.lockb \neq 0$  and  $\bar{X}^m.lockb \neq T.Tid$  then Abort  $T$ ;
14     if  $\bar{X}^m.wts \neq X^m.wts$  then Abort  $T$ ;
14 Function Commit( $T$ ):
15   foreach  $X \in T^l.ws$  do
16      $X^m.wts \leftarrow T^l.cts$ ;
17      $X^m.lockb \leftarrow 0$ ;
18      $\text{Write}^D(X.PK, X)$ ;

```

---

update  $X^d$  (lines 15–17). Upon abort of  $T$ , we only need to release all the locks on  $\forall X \in T^l.ws$ .

**6.2.3 T/O**

**T/O** orders transactions based on their beginning timestamps ( $T^l.bts$ ). If the execution order of the transactions does not match their beginning timestamp order, one of them needs to abort. To compare the timestamp, each data item maintains the maximum beginning timestamp of the transactions that have ever read/ written  $X$  ( $X^m.rts/X^m.wts$ ). For any read on data item  $X$  of transaction  $T$ , if there does not exist any conflicts, we update  $X^m.rts$  to  $\max\{X^m.rts, T^l.bts\}$ ; for any write on  $X$  of  $T$ , if there does not exist any conflicts, we update  $X^m.wts$  to  $\max\{X^m.wts, T^l.bts\}$ . Upon a conflict of  $T$  with some other transaction  $\bar{T}$  over  $X$ , we abort  $T$  if  $T^l.bts < \bar{T}^l.bts$ , meaning that  $T$  is supposed to be ordered before  $\bar{T}$  but  $T$  reads/writes  $X$  that  $\bar{T}$  has ever written<sup>3</sup>; otherwise,  $T$  must wait to commit/abort after  $\bar{T}$  commits/aborts to guarantee correctness.

To boost T/O using RDMA, it is necessary to re-implement its logic that performs (1) remote reads/writes on data item  $X$ , (2) remote update of  $X^m.rts$  or  $X^m.wts$ , and (3) wait-commit or cascading abort. In logic (2) of RDMA-T/O, when updating  $X^m.rts$ , only 8 bytes of metadata are modified, and the new value of  $X^m.rts$  needs to be calculated against the old value, satisfying the requirement of the Read-Cas principle. Unfortunately, while a remote transaction is modi-

**Algorithm 5: RDMA-T/O**


---

```

1 Function UpdateRTS( $PK, X, T$ ):
2   if  $X^m.rts < T^l.bts$  then
3      $\text{Atomic}^D(PK, MT\_RTS, X^m.rts, T^l.bts)$ ;
4      $X^m.rts \leftarrow T^l.bts$ ;
5     if  $X \neq \text{Read}^D(PK)$  then return false;
6   return true;
7 Function Read( $PK, T$ ):
8    $X \leftarrow \text{Read}^D(PK)$ ;
9   if  $X^m.wts > T^l.bts$  or  $X^m.lock \neq 0$  then
10    Abort  $T$ ;
11   if  $\neg \text{UpdateRTS}(X, PK)$  then Abort  $T$ ;
12    $T^l.bL \leftarrow X^m.wL \cup T^l.bL$ ;
13    $T^l.rs \leftarrow \{X\} \cup T^l.rs$ ; return  $X$ ;
14 Function Write( $PK, T, newV$ ):
15   if  $\neg \text{Atomic}^D(PK, MT\_LATCH, 0, T.Tid)$  then
16    Abort  $T$ ;
17    $X \leftarrow \text{Read}^D(PK)$ ;
18   if  $\max\{X^m.rts, X^m.wts\} > T^l.bts$  then
19      $\text{Atomic}^D(PK, MT\_LATCH, T.Tid, 0)$ ;
20     Abort  $T$ ;
21    $T^l.ws \leftarrow \{X\} \cup T^l.ws$ ;
22    $X^m.wts \leftarrow T^l.bts$ ;  $X^d \leftarrow newV$ ;
23    $T^l.bL \leftarrow X^m.wL \cup T^l.bL$ ;
24    $X^m.wL \leftarrow X^m.wL \cup \{T.Tid\}$ ;
25    $X^m.latch \leftarrow 0$ ;  $\text{Write}^D(PK, X)$ ;
26 Function CascadingAbortCheck( $T$ ):
27   foreach  $\bar{T}id \in T^l.bL$  do
28     do  $\bar{T} \leftarrow \text{Read}^I(\bar{T}id)$ ;
29     while  $\bar{T}^s.st = RN$ ;
30     if  $\bar{T}^s.st = AB$  then
31        $T^s.st \leftarrow AB$ ; return false;
32    $T^s.st \leftarrow CM$ ; return true;
33 Function Commit( $T$ ):
34   if  $\neg \text{CascadingAbortCheck}(T)$  then Abort  $T$ ;
35   foreach  $X \in T^l.ws$  do
36     while  $\neg \text{Atomic}^D(X.PK, MT\_LATCH, 0, T.Tid)$ ;
37      $\bar{X} \leftarrow \text{Read}^D(X.PK)$ ;
38      $\bar{X}^m.wL \leftarrow \bar{X}^m.wL - T.Tid$ ;
39      $\bar{X}^m.latch \leftarrow 0$ ;  $\text{Write}^D(PK, \bar{X})$ ;

```

---

fyng  $X^m.rts$ , another write transaction may simultaneously modify other metadata  $X^m.wts$ . Using the Read-Cas principle alone cannot ensure the atomicity of this modification. To address this issue, we adopt the Double-Read principle to ensure atomicity. This principle avoids the need for an explicit latch and ensures the atomicity of the read operation. Furthermore, since logic (3) does not modify the remote MetaTxn metadata, we can use the One-Read principle to implement it.

The metadata  $X^m$  in T/O maintains four fields: (1)  $X^m.latch$ , the latch of  $X$  to prevent concurrent modification, (2)  $X^m.rts$ /(3)  $X^m.wts$ , the maximum beginning timestamp

<sup>3</sup> Specifically, for any read of  $T$ , we abort  $T$  if  $T^l.bts < X^m.wts$ , and for any write of  $T$ , we abort  $T$  if  $T^l.bts < X^m.rts$  or  $T^l.bts < X^m.wts$ .

of the transactions that have ever read/written  $X$ , and (4) a list  $X^m.wL$ , each item of which is the ID of an uncommitted transaction that has ever written  $X$ . Besides  $T^l.rs$ ,  $T^l.ws$ ,  $T^l$  maintains the beginning timestamp  $T^l.bts$  of  $T$ , and an extra list  $T^l.bL$ , recording transactions that are ordered before  $T$ .  $T^s$  additionally maintains a transaction status  $T^s.st$ , which is *running* (*RN*), *committed* (*CM*) or *aborted* (*AB*). Key functions of the re-implementation called RDMA-T/O are shown in Algorithm 5.

Function *Read* is used to do remote read using the Double-Read principle. We first read  $X$  by issuing *Read<sup>D</sup>* (line 8) and check whether  $X$  is readable by  $T$  (line 9). If  $X$  is not readable by  $T$ , we abort  $T$  (lines 10); otherwise, we try to do a remote update on  $X^m.rts$  using  $T^l.bts$ , and perform another remote read on  $X$  to check whether the remote update is successful (lines 11, 2–6). If the update fails or the results of the two *Read<sup>D</sup>* operations are different, we abort  $T$ ; otherwise, we perform a local update on  $T^l.bL$  (the dependent transactions of  $T$ ) and store  $X$  in the read set  $T^l.rs$  (lines 12–14).

Function *Write* is used to do remote write. We first use *Atomic<sup>D</sup>* to acquire latch, then issue *Read<sup>D</sup>* to read back  $X$  and check whether  $X$  is writable by comparing  $\max\{X^m.rts, X^m.wts\}$  with  $T^l.bts$  (lines 15–20). If  $\max\{X^m.rts, X^m.wts\} > T^l.bts$ , meaning  $X$  is not writable by  $T$ , then we release the latch and abort  $T$  (lines 18–20); otherwise, it means that  $X$  is writable by  $T$ . We then store the original  $X$  to the write set  $T^l.ws$  for restoring  $X$  in case that  $T$  aborts (line 21). Subsequently, we do a local update on  $X^m.wts$ ,  $X^d$ ,  $T^l.bL$ ,  $X^m.wL$  (the running transactions with writes on  $X$ ),  $X^m.latch$ , and finally we apply the local updates on remote  $X$  by issuing *Write<sup>D</sup>* (lines 22–25).

Upon committing transaction  $T$ , it is necessary to execute function *CascadingAbortCheck* to check the status of each dependent transaction, maintained in  $T^l.bL$ , of  $T$  (lines 27–31). If any of them aborts, we would make a cascading abort of  $T$  (lines 30–31); if all of them commit, we first set the status  $T^s.st$  of  $T$  to *CM* (line 32), use a repeated invocation of *Atomic<sup>D</sup>* to acquire the latch of  $X$  (line 36), and do a remote update of  $\bar{X}^m$  by removing  $T.Tid$  from  $\bar{X}^m.wL$  for each data item  $\bar{X}$  that  $T$  has ever written (lines 37–39). Upon abort of transaction  $T$ , it is necessary to restore its modifications maintained in  $T^l.ws$  on each data item  $X$  that has ever been written. Note, if a dependent transaction of  $T$  aborts, and restores  $X$  that  $T$  has ever written, in this case,  $T$  cannot restore  $X$ ; otherwise, the value of  $X$  would be restored incorrectly. For example, suppose there exists a data item  $X$  with  $X^d = 1$ . Transaction  $T_1$  first updates  $X^d$  to 2, and transaction  $T_2$  then updates  $X^d$  to 3. Subsequently,  $T_1$  aborts, and restore  $X^d$  to 1. Because the abort of  $T_1$  causes a cascading abort of  $T_2$ ,  $T_2$  aborts. In this case, if  $T_2$  does a restore of  $X$  which changes  $X^d$  to 2, then  $X$  would be set in an incorrect value.

**Discussion.** One potential limitation of RDMA-T/O is that cascading aborts could waste CPU cycles. To eliminate cas-

cading aborts, one possible solution is to postpone the writes of each transaction  $T$  until the commit of  $T$ . In this case, upon any read or write of  $X$  from  $T$ , if there exists another uncommitted transaction  $\bar{T}$  with a smaller timestamp that writes  $X$ ,  $T$  must wait until  $\bar{T}$  commits. To notify the transactions that wait for  $\bar{T}$ , for each data item  $X$ , it is necessary to additionally maintain two lists,  $X^m.prL$  and  $X^m.pwL$  that record the pending transactions with reads and writes on  $X$ , respectively. After  $\bar{T}$  commits, we sequentially check each  $X$  that  $\bar{T}$  has written, and the transaction in  $X^m.prL \cup X^m.pwL$  with the smallest timestamp is scheduled to execute, while the other transactions still need to wait. Although the above solution can eliminate cascading aborts, it instead incurs other potential overheads, e.g., the prohibitive maintenance overhead of  $X^m.prL$  and  $X^m.pwL$  for each data item  $X$ , as well as a large amount of CPU idle time. For comparison, we follow the work proposed by P.A. Bernstein and N. Goodman [7] to adjust RDMA-T/O without cascading abort, and report the experimental evaluation over the two implementations in Sect. 7.6.

## 6.2.4 Deterministic algorithms

**Calvin** accepts and executes transactions in batches. In each batch, it determines the order of transactions in a first-come-first-serve manner. For each transaction  $T$ , if there does not exist any transaction that conflicts with  $T$  and is ordered before  $T$ ,  $T$  is scheduled to execute. Note, in Calvin, no transaction is aborted because of conflict. To be specific, Calvin is designed with three components to do concurrency control.

- **Sequencer** collects the transactions in batches. For each batch, it determines the order of transactions, breaks every transaction into several sub-transactions, and distributes sub-transactions to schedulers based on the data items they access.
- **Scheduler** schedules sub-transactions to execute in a predefined order. All sub-transactions attempt to acquire locks on data items to be read/written. When different sub-transactions apply for locks on the same data item, the scheduler grants the locks to sub-transactions in a predetermined order. If all locks are acquired, we schedule this sub-transaction to execute.
- **Transaction executor** executes sub-transactions. For each sub-transaction  $Ts$ , it performs local reads. For another sub-transaction that belongs to the same transaction with  $Ts$ , its writes may rely on reads of  $Ts$ , the transaction executor then synchronizes the reads to the other sub-transaction if necessary. Finally, the executor performs local writes and releases the acquired locks and commits.

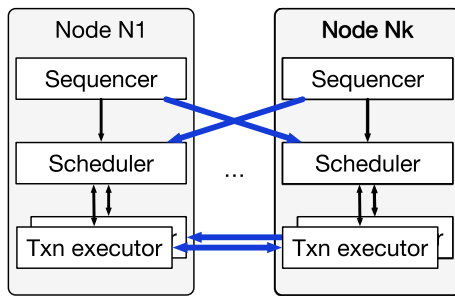


Fig. 3 Calvin overview

There are two opportunities to optimize Calvin using RDMA: 1) distributing sub-transactions from sequencers to schedulers, and 2) synchronizing reads among executors. For illustration purposes, we show these two opportunities by the blue lines in Fig. 3. For opportunity one, we implement a circular buffer following FaRM [19] to replace the original TCP/IP network message transmission mechanism. And for opportunity two, we design a fixed-length read set buffer stored in  $Ts^s$  for each sub-transaction  $Ts$ . To do this, when a sub-transaction  $Ts_1$  on Node  $N1$  synchronizes its read set to another sub-transaction  $Ts_2$  on the remote node  $N2$ , we remotely write local read set into  $Ts_2^s.N2.rs$  by issuing a  $Write^D$ . After  $Ts_2$  gathers all read sets of remote nodes,  $Ts_2$  can continue to perform the execution.

**Discussion.** However, as reported in [29], network overhead is comparatively trivial and not the bottleneck of Calvin. Therefore, the above two optimizations cannot help effectively improve the system performance, we argue using RDMA to optimize Calvin cannot bring obvious benefits and verify this observation in the experiment section. The same reason in Calvin is also applicable for other deterministic algorithms such as Q-Store [41], LADS [63] and QueCC [42] whose network usage is quite limited.

### 6.3 Optimizations

To further optimize the performance of concurrency control algorithms, we implement the following four optimizations in RCBench.

- **Coroutine.** One thread can have multiple coroutines, each of which executes transactions sequentially. Coroutines of the same thread are switched in a round-robin manner upon one coroutine is blocked by waiting for the results of RDMA verbs: Upon sending an RDMA verb from one coroutine, another coroutine of the same thread is switched to execute transactions. By using coroutines in the Boost C++ library with low context switch overhead (about 20 ns), we can reduce the CPU idle time, and hence improve the throughput of each algorithm.

- **Doorbell Batching (a.k.a. DB).** Usually, a verb takes one Memory Mapping I/Os (MMIOs). Instead, DB encapsulates multiple verbs into a batch and calls a single MMIO to send the beginning address of the batch. This address serves as a ringing doorbell to notify RNIC to fetch the batch through one or more DMAs. In this way, expensive MMIOs are replaced by a low-cost CPU and bandwidth-efficient DMA, therefore improving the performance of each algorithm. Given a batch of verbs, DB works only if there is no dependency among these verbs. For example, we cannot batch  $Read^D$  with  $Atomic^D$  in Function  $AcqIMShLock$  (lines 7,13 in Algorithm 3) because the inputs of  $Atomic^D$  primitive rely on the result of  $Read^D$  primitive. For this reason, we sequentially check the primitives evoked by the concurrency algorithms and batch continuous primitives without dependency using DB. For example, we batch continuous  $Atomic^D$  and  $Read^D$  at lines 3,5 in Algorithm 5 in terms of the same node to eliminate expensive MMIOs.
- **Outstanding Requests (a.k.a. OR).** RDMA NIC (RNIC) executes one-sided verbs sequentially, meaning that a verb starts to be sent until the results of previous verbs return. To improve the degree of concurrency, OR optimizes the mechanism of message communication by starting to send the verb upon the accomplishment of sending the previous one. In our framework, based on DB, we further optimize the batches to be sent to different nodes using OR.
- **Passive Ack (a.k.a. PA).** In our framework, RNIC processes RDMA verbs in a first-come-first-serve manner with a reliable connection. Therefore, given a batch of verbs, as long as we acknowledge the result of the last verb, we can guarantee that the results of previous verbs have also been returned. By doing this, redundant acknowledgment for previous verbs can be eliminated and hence save the bandwidth of RNIC. In our framework, based on DB, we further optimize multiple batches to be sent using PA.

## 7 Experiment

### 7.1 Setup

We use two popular OLTP benchmarks, YCSB [15] and TPCC [55], to evaluate our re-implementations of concurrency control algorithms.

**YCSB [15]** is a comprehensive benchmark that simulates large-scale Internet applications. Its dataset contains a single 10-column relation, in which each tuple occupies 1KB. The table is horizontally partitioned, and each node is set to have a fixed number of 10 million records, resulting in 10GB of data per node. Each transaction of the workloads

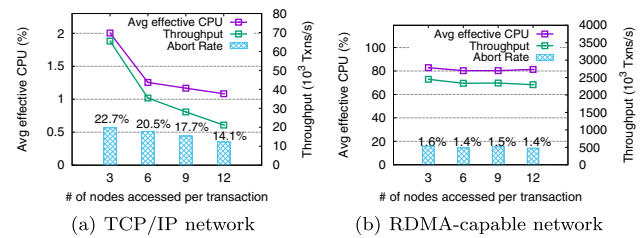
is set to have a fixed number of 10 read/write operations that access data items following the Zipfian distribution. YCSB provides adjustable parameters to simulate workloads with diverse characteristics. The skew factor parameter determines the degree of contention where the access of records follows the Zipfian distribution. The write-ratio parameter controls the ratio of write operations in transactions. Setting the write ratio to 0.2 means that there are 80% reads and 20% writes among all transactions. In this scenario, there might be approximately 10% ( $0.8^{10} \approx 0.1$ ) read-only transactions in the system, while the remaining 90% could be mixed read/write transactions. Additionally, YCSB offers the ( $\theta$ ) parameter, which enables controlling the number of data nodes that will be accessed per transaction. **By default, we set both the write ratio and the skew factor to 0.2, and set  $\theta$  to 2.**

**TPCC** [55] is another OLTP benchmark that simulates warehouse ordering applications. Its dataset contains 9 relations, and each warehouse is set to have 100MB of data. By default, we set the number of warehouses per node to 32. TPCC contains 5 types of transactions, among which Payment, New-order, and Delivery are read-write transactions, Stock-level and Order-status are read-only transactions. As Delivery, Stock-level and Order-status only involve local operations, similar to previous works [29, 65], we focus on Payment and New-order only to evaluate the transaction scalability of distributed transactions.

We evaluate our system on 4 machines of an RDMA-capable EDR cluster. Each node is equipped with one Intel(R) Xeon(R) Gold 5220 CPU @ 2.20GHz (18 cores  $\times$  2 HT) processor, 128GB RAM, and one ConnectX-5 EDR 100Gb/s InfiniBand MT27800. By default, each machine is assigned one RCBench node and one client. Each RCBench node is configured to have 4 threads to receive the transactions and send the response from/to clients, and 24 threads to execute transactions, each thread is set to have 8 coroutines. Each client is configured to have 4 threads to generate new transactions, and 4 threads to send the transactions and receive the response to/from RCBench nodes. For each experiment, we initiate a 30-second warm-up period followed by the collection of results for the subsequent 60 s.

## 7.2 Effective CPU utilization rate

In order to evaluate the efficiency of RCBench, we propose a metric in this section to assess the extent to which our designs take full advantage of RDMA. Specifically, by noticing that the introduction of RDMA significantly reduces the network overhead, it is widely recognized that the bottleneck has shifted from network to CPU in RDMA-capable distributed clusters. The CPU utilization rate appears to be a reasonable metric in this context. However, the existing formulation of CPU utilization rate generally incorporates



**Fig. 4** The throughput follows the same trend of  $\overline{U}_{CPU}^e$

irrelevant overheads, such as the thread/coroutine scheduling costs, etc. As a result, they do not perform as a precise indicator. To address this issue, we propose a new metric called “effective CPU utilization rate ( $\overline{U}_{CPU}^e$ )” to measure the actual utilization of the CPU. The formula is shown below.

$$\overline{U}_{CPU}^e = \frac{\text{Effective time to execute committed transactions}}{\text{Total time to execute transactions}}$$

$\overline{U}_{CPU}^e$  is defined as the proportion of total CPU execution time spent on the execution of read/write and concurrent control operations of committed transactions. Specifically, time spent on network communication, thread/coroutine coordination, and idle time is excluded from the calculation of  $\overline{U}_{CPU}^e$ . By tracing, breaking down, and averaging the time spent on each operation in threads, we can measure  $\overline{U}_{CPU}^e$  and use it to assess the effectiveness of our re-implementations. The higher  $\overline{U}_{CPU}^e$  is, the better RDMA is leveraged for efficient transaction processing. After calculating  $\overline{U}_{CPU}^e$  for each node, we determine the average  $\overline{U}_{CPU}^e$  of  $\mathcal{N}$  nodes, denoted as  $\overline{U}_{CPU}^e$ , for the entire distributed system. Formally,  $\overline{U}_{CPU}^e$  is calculated as follows:

$$\overline{U}_{CPU}^e = \frac{\sum_{i=1}^{\mathcal{N}} (\overline{U}_{CPU}^e \text{ of node } i)}{\mathcal{N}}$$

$\overline{U}_{CPU}^e$  can be interpreted as the effective working time per unit of time to execute committed transactions, which is similar to throughput quantified by the number of committed transactions per unit of time. Figure 4a shows the throughput,  $\overline{U}_{CPU}^e$  and abort rate of 2PL algorithm in conventional shared-nothing TCP/IP architecture. As the number of accessed data nodes per transaction increases from 3 to 12, the  $\overline{U}_{CPU}^e$  follows the same trend as that of the throughput, verifying it as a reliable metric of our re-implementation efficiency. Moreover, the abort rate decreases as the throughput decreases, indicating that the drop in performance is not attributable to aborts. For comparison, we plot the results of RCBench under the same setup in Fig. 4b, which displays a much stabler throughput,  $\overline{U}_{CPU}^e$ , and abort rate.

Note that the CPU utilization metric proposed by Binnig et al. [9, 66] is somewhat similar to our proposed  $\overline{U}_{CPU}^e$ . However, they mainly focus on CPU cycles consumed by

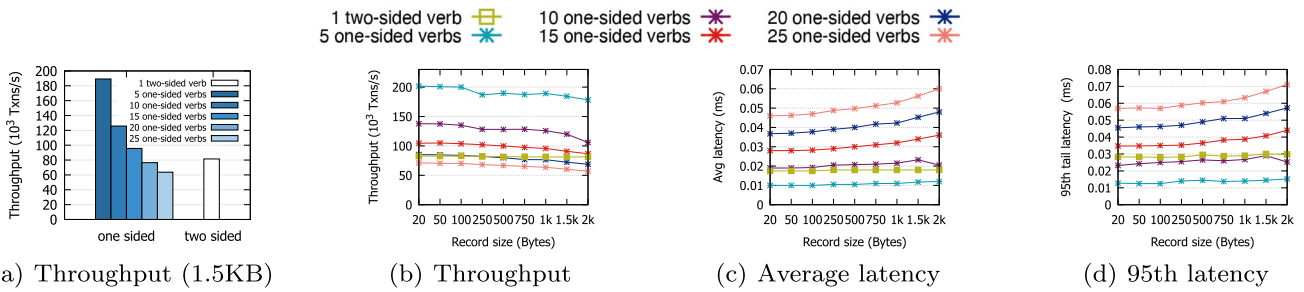


Fig. 5 A performance comparison between RDMA verbs

Table 3 A comparison of network types for YCSB workload

Algorithm	TCP baseline			+two-sided			+one-sided			
	TPS(k)	$\overline{U}_{CPU}^e$	AR	TPS(k)	$\overline{U}_{CPU}^e$	AR	TPS(k)	$\overline{U}_{CPU}^e$	AR	$\sigma$
No-Wait	26.74	1.6 %	20.8 %	85.36	6.2 %	15.4 %	<b>991.07(37.1X)</b>	15.2 %	0.0 %	23.5
Wait-Die	31.03	2.1 %	13.4 %	85.64	7.2 %	7.8 %	<b>807.78(26.0X)</b>	13.7 %	0.8 %	30.2
Wound-Wait	27.98	1.9 %	6.9 %	85.02	6.5 %	10.3 %	<b>714.57(25.5X)</b>	15.6 %	0.0 %	31.2
T/O	32.10	2.3 %	1.4 %	85.66	6.9 %	1.4 %	<b>853.63(26.6X)</b>	17.5 %	0.5 %	25.4
MVCC	32.34	2.3 %	0.9 %	84.84	7.2 %	1.2 %	<b>939.12(29.0X)</b>	17.9 %	0.3 %	22.8
Silo	25.01	1.8 %	17.9 %	69.57	5.5 %	17.5 %	<b>1071.59(42.8X)</b>	18.8 %	0.3 %	17.7
MaaT	25.58	3.2 %	7.2 %	49.40	10.2 %	10.0 %	<b>460.68(18.0X)</b>	14.3 %	0.0 %	51.3
Cicada	27.11	2.1 %	3.5 %	82.62	9.5 %	2.4 %	<b>720.66(26.6X)</b>	15.4 %	0.3 %	32.2
Calvin	249.62	2.0 %	0.0 %	<b>261.26</b>	2.0 %	0.0 %	207.78(0.8X)	2.0 %	0.0 %	

network communications itself. Instead, it is considered ineffective and excluded from  $\overline{U}_{CPU}^e$ . Harding [29] discusses the breakdown of distributed transaction executions in detail. However, they do not identify ineffective works, and abstract no concept concerning CPU utilization rate, which is different from our idea of  $\overline{U}_{CPU}^e$ .

### 7.3 Comparison between RDMA verbs

We investigate the efficiency difference between one-sided and two-sided verb invocations. We assume that executing a transaction requires either one two-sided verb or 5, 10, 15, 20, or 25 one-sided verbs, with each verb accessing a fixed length of data, such as 1500 bytes. This is because some algorithms, like Silo, require reading the metadata of the same data items multiple times to check if it has been modified by other transactions. To ensure fairness, both one-sided and two-sided verbs are called sequentially, that is, another one/two-sided verb cannot be called until the result of the previous message is returned.

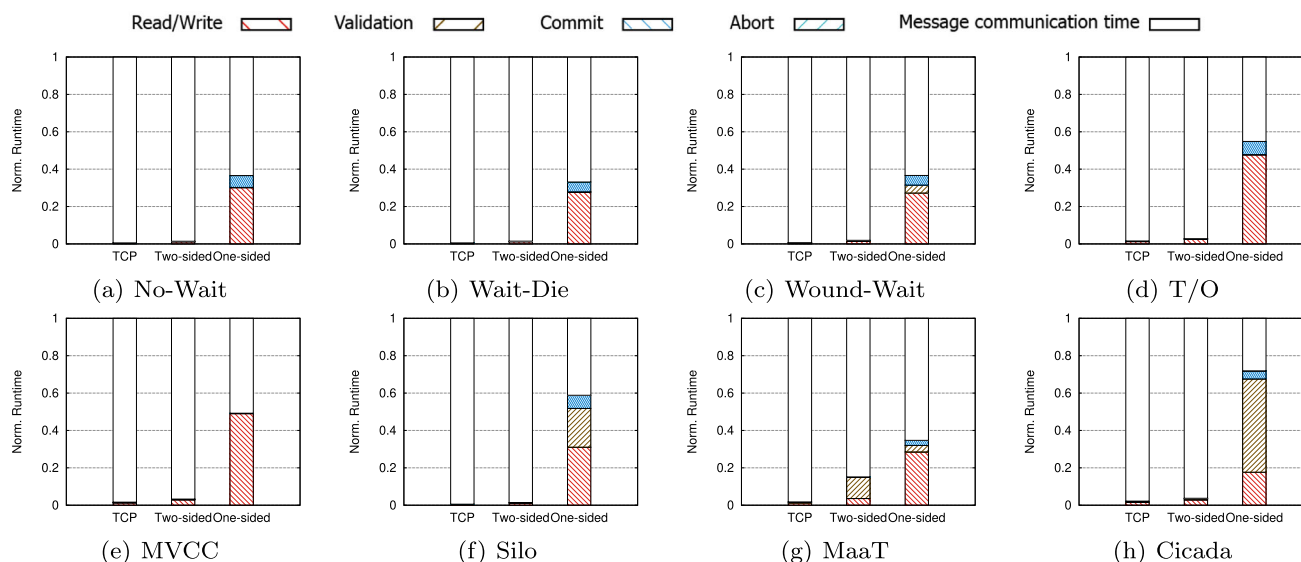
As shown in Fig. 5a, we plotted the throughput of transactions that required either one two-sided verb or 5, 10, 15, 20, or 25 one-sided verbs. The graph shows that when the data access size is fixed at 1500 bytes, the throughput of one two-sided verb transaction is equivalent to using 20 one-sided verbs transaction. To further evaluate the performance of

one-sided and two-sided verbs, we plotted the latency and tail (95th) latency of the transactions in Fig. 5c, d, respectively. It can be observed that the latency and tail latency of one two-sided verb transaction (yellow lines) are equivalent to that of using 10 one-sided verbs transaction. Even when we adjust the data access size from 20 to 2000, as shown in Fig. 5b, c, d, the difference in throughput and latency between using one-sided verbs and two-sided verbs remains consistent. Note that the latency of two-sided verbs we have observed in our evaluation is considerably higher than those reported in [9, 38]. This disparity arises because the two works primarily concentrated on assessing the efficiency of one-sided and two-sided verbs. In contrast, our experiments involving two-sided verbs take into account additional factors, including the overhead associated with CPU scheduling on remote machines. For instance, after receiving a message, a remote machine must place the message in a queue and allocate a Txn executor to execute the message. Indeed, the choice of such a setup is justified because two-sided verbs inherently involve the CPU of the remote machine.

### 7.4 Effect of RDMA networks

We study the effect of RDMA networks on YCSB in terms of throughput (TPS),  $\overline{U}_{CPU}^e$ , and abort rate (AR). For illustration, the re-implementations over TCP/IP Ethernet networks, two-





**Fig. 6** Time breakdown under different networks

sided RDMA networks, and one-sided RDMA networks are referred to as *TCP/IP*, *two-sided*, *one-sided*. Note, *two-sided* is simply implemented by replacing the network invocations in Deneva [29] with two-sided verb invocations.

The results are reported in Table 3. We use bold to highlight the best throughput of each algorithm over the three networks. As we can see, for *TCP/IP*,  $\overline{U}_{\text{CPU}}^e$  is rather low, ranging only from 1.6%–3.2%, leading to a poor throughput; for two-sided, compared with *TCP/IP*, except Calvin,  $\overline{U}_{\text{CPU}}^e$  improves by at least 3.0X, causing an improvement of throughput ranging from 1.9X to 3.2X; compared with *two-sided*, *one-sided* further improves  $\overline{U}_{\text{CPU}}^e$  and throughput ranging from 1.4X to 3.4X and 8.4X to 11.6X, respectively. This is because, in our framework, we completely eliminate the expensive 2PC overhead and enjoy all benefits of RDMA networks. Compared with *TCP/IP*, *one-sided* achieves a significant improvement of throughput ranging from 18.0X to 42.8X except for Calvin.

To precisely analyze the reason why *one-sided* outperforms *TCP/IP* and *two-sided*, as shown in Fig. 6, we evaluate the time breakdown for the eight algorithms except for Calvin. We break down the execution time of a transaction into five parts, including the read/write operations, the validation operations, the commit operations, the abort operations, and the message communication time in all three phases. As we can see, for each algorithm, *TCP/IP* and *two-sided* consume almost all of the time to perform message communication, while *one-sided* uses 28% (Cicada)—67% (Wait-Die) of the time to perform remote one-sided verbs. While the two-sided verb benefits from the high-performance InfiniBand network, its message communication time is still relatively long due to the additional overhead incurred by the communication mechanism. This includes operations such

as message sending and receiving, thread/coroutine scheduling before executing the message, and the creation of the transaction context. These operations cannot be avoided by two-sided verbs and contribute to the overall communication overhead. Thus, *one-sided* achieves such a high speedup over the *TCP/IP* and *two-sided* variants. In Table 3, we can observe that the *one-sided* variants have a much lower abort rate (AR) than the *TCP/IP* and *two-sided* variants. This is because *one-sided* variants have lower message communication time overhead, lower transaction processing time, and lower concurrency conflicts with other transactions, making them more efficient and less prone to conflicts and aborts.

In addition, for the throughput of different algorithms under *one-sided*, Silo performs the best. The reason is that the throughput is closely related to the averaged number ( $\sigma$ ) of primitive invocations per transaction. For reference, we report  $\sigma$  for each algorithm in Table 3, showing Silo takes the smallest  $\sigma$ . Often, a smaller  $\sigma$  leads to a higher throughput. However, the complexity of concurrency control algorithms may slightly affect the performance. For example, compared with No-Wait, MVCC takes a slightly smaller  $\sigma$  but has a lower throughput due to its complexity of traversing versions. Note that in the current YCSB workload, each transaction consists of 10 operations and accesses two nodes, which means about half of the operations (approximately 5 operations) may be remote operations. In Table 3,  $\sigma$  represents the total number of primitive invocations for approximately 5 remote read/write, 1 validation, and 1 commit operation of a transaction. For *two-sided*, the same workload requires 7 rounds of two-sided network communication (5 remote read/write + 1 validation + 1 commit operations), which can be translated to about 70–140 one-sided verb invocations, as described in Sect. 2.2 and 7.3. Moreover,  $\sigma$  is not

**Table 4** A comparison of one-sided and two-sided RDMA verbs combinations

Algorithm	two-sided	r/w	co				all	
No-Wait	85.36	359.20	113.22				991.07	
Wait-Die	85.64	387.28	139.39				807.78	
Wound-Wait	85.02	384.48	110.39				714.57	
T/O	85.66	106.98	96.14				853.63	
MVCC	84.84	236.08	103.25				939.12	
	two-sided	r/w	va	co	r/w+va	va+co	r/w+co	all
Silo	69.57	235.06	199.38	196.74	143.25	139.96	379.65	1071.59
MaaT	49.40	66.84	100.60	71.76	228.83	141.79	151.63	460.68
Cicada	82.62	163.87	126.83	126.04	345.63	169.23	215.82	720.66

**Table 5** A comparison of one-sided implementation with different optimization

Algorithm	one-sided		+ db&pa&or		+ coroutine		+ cor+db&pa&or	
	TPS(k)	$\overline{U}_{CPU}^e$	TPS(k)	$\overline{U}_{CPU}^e$	TPS(k)	$\overline{U}_{CPU}^e$	TPS(k)	$\overline{U}_{CPU}^e$
No-Wait	991.1	15 %	1091.0	21 %	<b>2274.3</b>	78 %	2145.6	58 %
Wait-Die	807.8	14 %	928.9	20 %	<b>1909.7</b>	78 %	1749.6	60 %
Wound-Wait	714.5	16 %	907.7	19 %	<b>1803.1</b>	75 %	1771.7	60 %
T/O	853.7	17 %	887.7	18 %	1663.3	50 %	<b>1706.3</b>	50 %
MVCC	939.1	17 %	970.0	17 %	1839.6	48 %	<b>1865.0</b>	51 %
Silo	1071.5	19 %	1184.3	21 %	2383.4	73 %	<b>2385.0</b>	75 %
MaaT	460.7	14 %	531.9	15 %	<b>796.7</b>	42 %	791.7	37 %
Cicada	720.7	15 %	759.9	15 %	1492.3	54 %	<b>1506.0</b>	55 %

expected to increase significantly in the TPCC workload since a New-order transaction has 5–15 remote read/write operations, and a Payment transaction has only 2 remote read/write operations. As opposed to the other concurrency control algorithms, Calvin takes a similar throughput under both TCP/IP and RDMA networks. This is because Calvin is bottlenecked by its scheduler rather than the networks, while RDMA is only used to alleviate the bottleneck by the networks. In the rest of the paper, we focus on one-sided and do not report the results for Calvin.

## 7.5 Comparison with hybrid variants

Following the hybrid approach in DrTM+H [60], we study the throughput under the combination of one-sided and two-sided RDMA verbs, as shown in Table 4. Because the lock-based algorithms and the timestamp-based algorithms do not require the validation phase, we implement two extra scenarios for these algorithms, only using one-sided verbs in the read/write phase (r/w in Table 4) or the commit phase (co in Table 4). For optimistic algorithms, we re-implement six scenarios for them, using one-sided verbs in (1) the read/write phase (r/w), (2) the validation phase (va), (3) the commit phase (co), (4) the read/write and validation phase (r/w+va), (5) the validation and commit phase (va+co), and (6) the read/write and commit phase (r/w+co).

As we can see, by introducing one-sided verbs at some phases, the throughput of most algorithms can be improved but is still not comparable to that using one-sided verbs at all phases (all in Table 4). The reason is that the hybrid variants still rely on two-sided verbs, which are subject to the overhead of remote node participation and scheduling, thus affecting performance. Therefore, in RCBench, one-sided verbs are more suitable for promoting each algorithm.

## 7.6 Effect of various optimizations

### 7.6.1 General optimizations

We investigate the effect of applying optimizations including coroutine, OR, DB, and PA. Table 5 reports the throughput and  $\overline{U}_{CPU}^e$  by combining different optimizations. We use bold to highlight the best throughput of each algorithm after optimization. As we can see, by introducing DB, PA, and OR,  $\overline{U}_{CPU}^e$  and the throughput are improved slightly ranging from 1X to 1.47X and 1.03X to 1.27X, respectively; interestingly, by introducing coroutine only, both  $\overline{U}_{CPU}^e$  and throughput are improved greatly, indicating that coroutine is a more important optimization factor than the others; adding DB, PA, and OR with coroutine can only improve  $\overline{U}_{CPU}^e$  and throughput of some algorithms, such as T/O and Cicada, by at most 2%, while cannot bring benefits for the others. This is because the

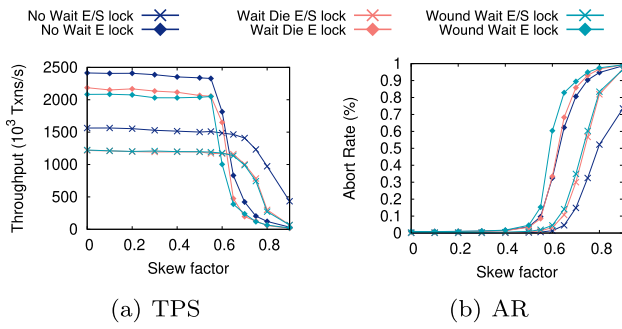


Fig. 7 Effect of different lock types

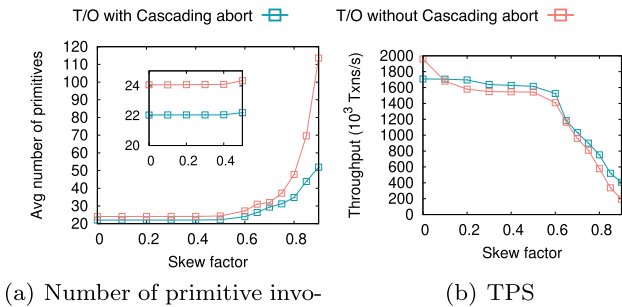


Fig. 8 Effect with or without cascading abort

coroutine reduces the idle time of the CPU, which dominates the benefit brought by DB, PA, and OR. Again, Silo performs the best, followed by No-Wait and MVCC.

7.6.2 Variants of 2PL

We re-implement and compare variants of 2PL under two types of locking mechanisms: exclusive/shared locking (abbreviated as *E/S lock*) and exclusive locking (abbreviated as *E lock*). The results are reported in Fig. 7. When the skew factor ranged from 0 to 0.6, variants under *E lock* outperformed those under *E/S lock*. This is because *E lock* used fewer one-sided verb invocations than *E/S lock* to acquire locks. In a low contention scenario, the benefit of improving concurrency was outweighed by the overhead of conflict examination. However, when the skew factor exceeded 0.7, the abort rate of *E lock* increased more quickly than that of *E/S lock*. This indicated that the benefit of improving concurrency became comparable to, or even outweighed, the overhead of conflict examination. In this scenario, *E/S lock* performed better than *E lock*.

7.6.3 Variants of T/O

We plot the result of T/O with or without enabling cascading abort in Fig. 8. As we can see, T/O with or without cascading abort performs similarly, and in most cases, as opposed to the phenomenon in the centralized environment, T/O with

cascading abort performs slightly better. This is because, as shown in Fig. 8a, T/O without cascading abort consumes extra one-sided verb invocations to manipulate two pending lists ( $X^m.prL$  and  $X^m.pwL$ ), which is comparable to that brought by the cascading abort.

7.7 Effect of contention levels

We evaluate the performance by varying the skew factor from 0 to 0.95. The throughput and  $\overline{U}_{CPU}^e$  under the low write-ratio (write-ratio = 0.2) are reported in Fig. 10b and 10a, respectively. As we can see, the throughput and  $\overline{U}_{CPU}^e$  of all algorithms remain relatively stable when the skew factor varies from 0 to 0.6. However, when the skew factor is greater than 0.65, the performance of all algorithms drops sharply. The optimal choice of the algorithm also tends to change with the skew factor. While Silo excels under the low and moderate contention (i.e., skew factor less than 0.8) due to the same reason mentioned in Sect. 7.6, Cicada performs the best under high contention, mostly attributable to the multi-version mechanism it uses to enable concurrent read and write operations. Figure 10d, c reports the throughput and  $\overline{U}_{CPU}^e$  under the high write-ratio (write-ratio = 0.8), respectively. The results follow a similar trend to those under low write ratios.

7.8 Effect of write ratios

We evaluate the performance by varying write ratios from 0 to 1.0. Figure 9b and 9a reports the throughput and  $\overline{U}_{CPU}^e$  under low contention (skew factor = 0.2), respectively. It can be observed that with increasing the write ratio, the throughput of Silo, Cicada, MVCC, and T/O drops slightly. Because each read or write operation takes the same number of primitive calls and acquires exclusive lock in No-Wait, Wait-Die, Wound-Wait, and MaaT, the throughput of them remains stable. We then report the results under high contention (skew factor = 0.8) in Fig. 9d, c, respectively. Due to the same reason, the throughput and  $\overline{U}_{CPU}^e$  of No-Wait, Wait-Die, Wound-Wait, and MaaT remain stable. However, the throughput and  $\overline{U}_{CPU}^e$  of Silo, Cicada, MVCC, and T/O drop significantly with increasing the write ratio due to the high contention level. Note that due to the multi-version mechanism, Cicada performs the best when the write ratio is larger than 0.2.

7.9 Scalability

Scalability is evaluated from two perspectives. First, we measure transaction scalability by varying the number ( $\theta$ ) of data nodes to be accessed per transaction while leaving the total number ( $\mathcal{N}$ ) of data nodes in the system to be constant. In contrast, system scalability is measured with a fixed  $\theta$  and

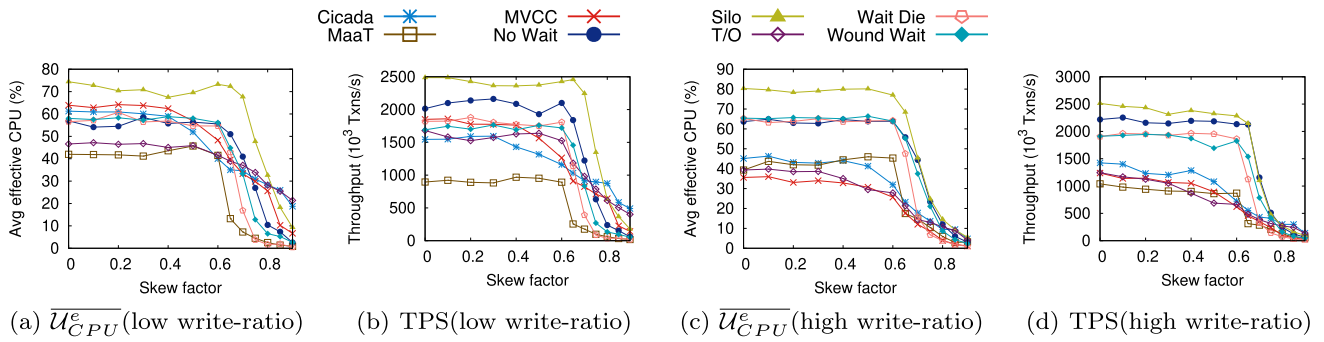


Fig. 9 The comparison of different contention levels

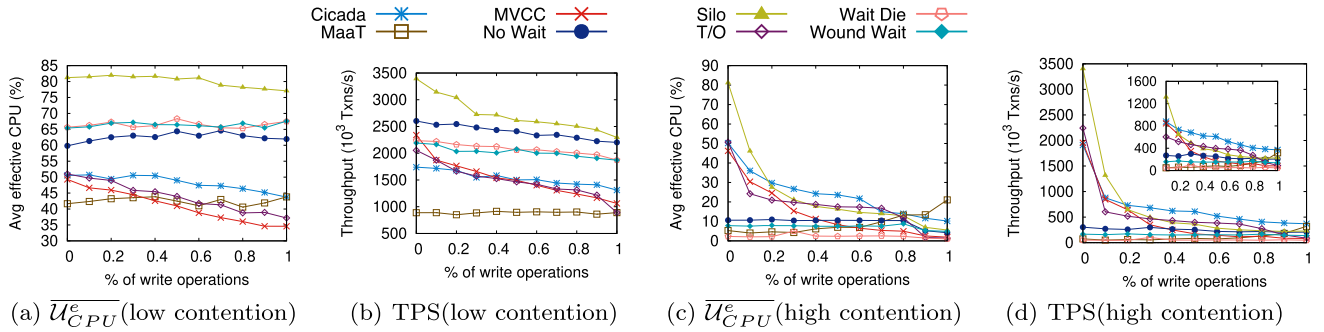


Fig. 10 The comparison of different write-ratios

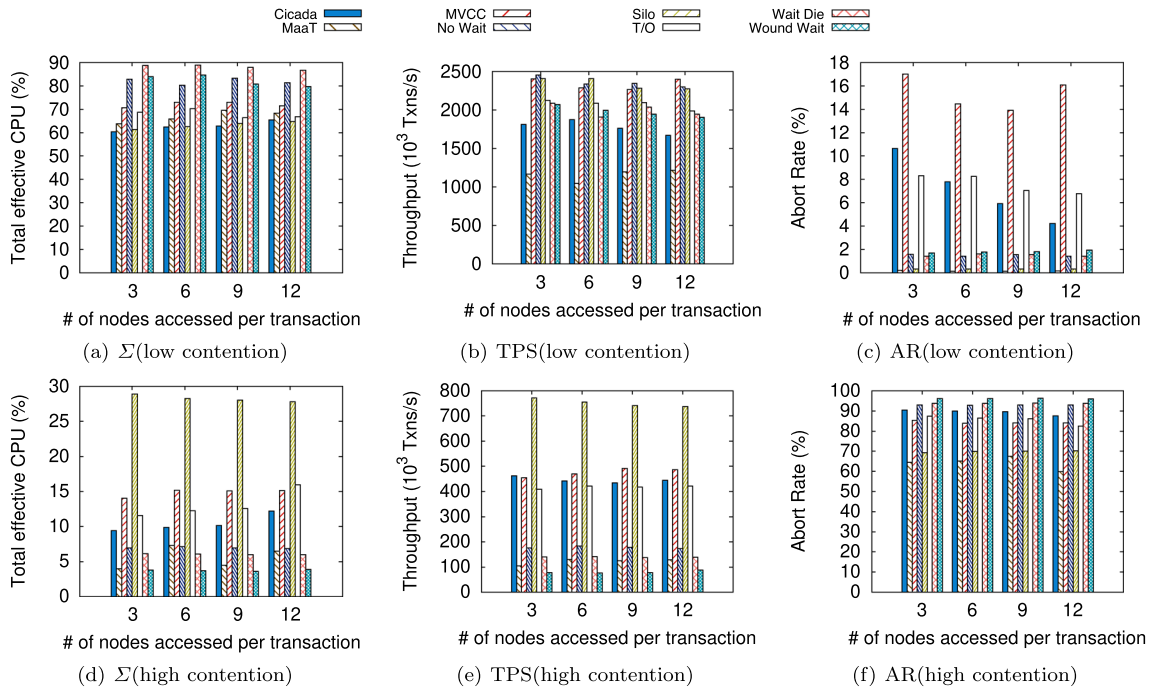


Fig. 11 Transaction scalability

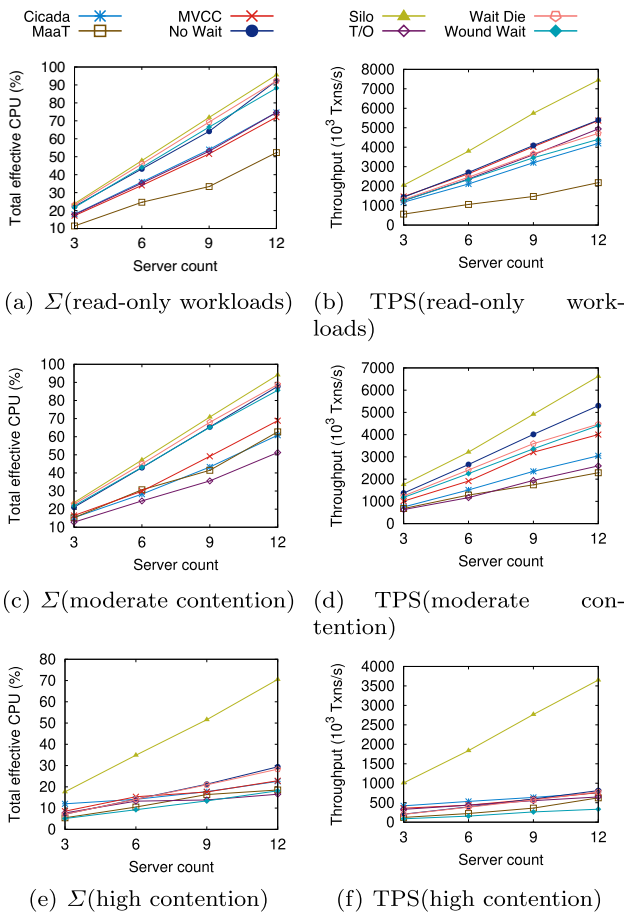


Fig. 12 System scalability under YCSB

varied  $\mathcal{N}$ . To fix  $\theta$ , we follow the implementation in Deneva, where the client knows the storage locations of all data items based on their hash value and generates transactions with fixed  $\theta$  based on the storage locations of all data items.

### 7.9.1 Transaction scalability

We test the transaction scalability by varying  $\theta$  from 3 to 12, and fixing  $\mathcal{N}$  to 12. Each transaction is composed of 12

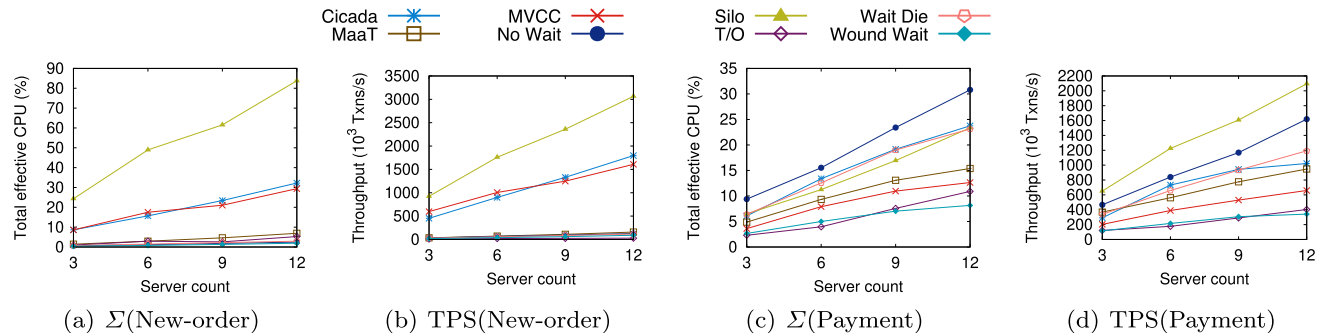


Fig. 13 System scalability under TPCC

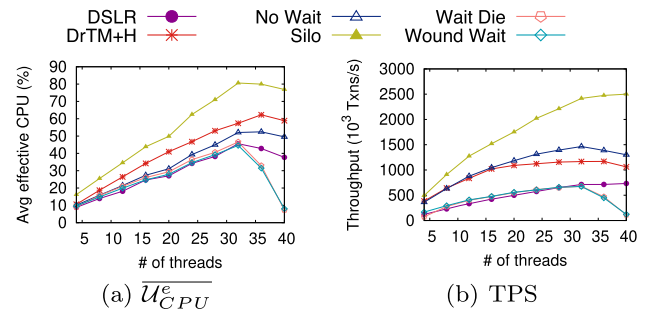


Fig. 14 Comparison with RDMA-based algorithms [60, 64]

read/write operations, with the percentage of remote operations fixed to 11/12. We report the throughput, total CPU utilization rate ( $\Sigma = \mathcal{N} \times \overline{U}_{CPU}^e$ ), and abort rate (AR) under low contention (skew factor = 0.2) in Fig. 11b, a, c, respectively. The throughput for different algorithms ranges from 1.0M TPS to 2.4M TPS, and for all algorithms,  $\Sigma$  as well as the throughput drops rather slightly. We report the results under high contention (skew factor = 0.8) in Fig. 11e, d, f, respectively. For the same reason,  $\Sigma$  and the throughput follow the same trend as those under low contention. This finding verifies that our re-implementations of algorithms using RDMA networks demonstrate arguably transaction scalability.

Note the throughput in this experiment cannot be directly compared with those of other experiments, for two reasons: first, the number of YCSB operations is set to 12 in this experiment, while it is set to 10 in other experiments; second, while the transactions in this experiment involve 11 remote operations, the transactions in the other experiments typically involve around 5 remote operations.

### 7.9.2 System scalability

We evaluate the system scalability by varying  $\mathcal{N}$  from 3 to 12 and fixing  $\theta$  to 2. We report  $\Sigma$  under read-only workload, moderate contention (skew factor = 0.5) workload and high contention (skew factor = 0.8) workload in Fig. 12a, c,

e, respectively. As we can see, when  $\mathcal{N}$  varies,  $\Sigma$  increases linearly. Besides, for the same  $\mathcal{N}$ ,  $\Sigma$  drops slightly when the contention increases. We also plot the throughput in Fig. 12b, d, f, respectively. The results show that the throughput increase linearly, which follows the same trend with  $\Sigma$ . We plot the results over TPCC in Fig. 13a–d. As we can see, the result follows a similar trend to that over YCSB.

### 7.10 Comparison with other RDMA-based algorithms

To show the effectiveness of our re-implementations, we compare Silo against DrTM+H [60], another RDMA-based re-implementation of Silo, and compare No Wait, Wait Die, and Wound Wait against DSLR [64], another RDMA-based re-implementation of 2PL variants. We report the comparison in Fig. 14. It can be observed that our Silo outperforms DrTM+H in terms of both throughput and  $\overline{\mathcal{U}}_{\text{CPU}}^e$  significantly. Besides, there is an obvious trend of growing returns when the number of threads increases. This is because, in DrTM+H, the remote accesses in the validation phase and the commit phase are all implemented using two-sided verbs, while in our work, we optimize all the remote accesses completely using one-sided verbs, thus improving  $\overline{\mathcal{U}}_{\text{CPU}}^e$  and the throughput significantly. In addition, as we can see, our No-Wait under *E/S lock* mechanism performs better than DSLR, and our Wait-Die and Wound-Wait take a similar throughput as DSLR. This is because DSLR consumes a similar number of one-sided verb invocations as Wait-Die and Wound-Wait, while No-Wait consumes the smallest one-sided verb invocations. However, the throughput and  $\overline{\mathcal{U}}_{\text{CPU}}^e$  of Wait-Die and Wound-Wait drop significantly when the number of threads is greater than 32. This is because Wait-Die and Wound-Wait maintain a fixed size of  $X^m \cdot pL$ , which causes extra aborts of transactions under the high contention levels. Additionally, we evaluated the transaction scalability of the algorithms in RCbench and compared them with the DrTM+H system, as shown in Fig. 15. We observed that although both the algorithms re-implemented in RCbench and DrTM+H can achieve transactional scalability, RCbench ensures up to 1.1X better transaction scalability than DrTM+H. Moreover, the throughput of concurrency control algorithms in RCbench, such as Silo and No-Wait, is higher than that in DrTM+H.

### 7.11 Summary

After conducting extensive evaluations, we summarize the major experimental findings below.

1. Transaction scalability can be arguably achieved (Sect. 7.9.1).

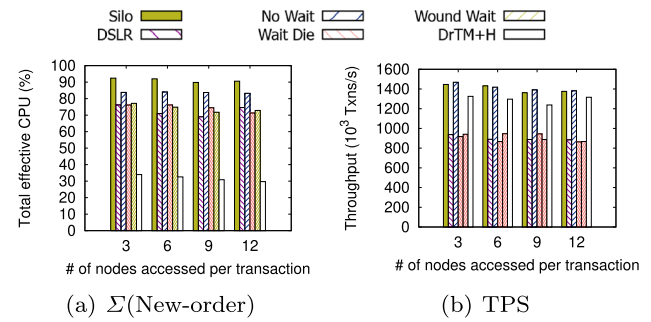


Fig. 15 Transaction scalability compared with RDMA-based algorithms

2. Concurrency control algorithms can achieve system scalability. (Sect. 7.9.2).
3. For the same algorithm,  $\overline{\mathcal{U}}_{\text{CPU}}^e$  is the dominant factor to scale out distributed transaction processing (Sect. 1 and 7.9).
4. It is practical and convenient to re-implement concurrency control algorithms using our proposed primitives, and even performs better than some customized implementations in many cases (Sect. 7.10).
5. Optimizing Calvin using RDMA cannot bring obvious benefits (Sects. 6.2.4 and 7.4).  
For other algorithms, RDMA networks bring significant performance improvement, and the degree of improvement closely relies on the number of primitive invocations and metadata complexity. Based on these criteria, Silo is reported to be the best in most cases (Sects. 7.4, 7.7 and 7.8).
6. Among all optimization techniques, coroutine brings the maximal performance improvement (1.7X to 2.5X). For 2PL variants, single exclusive locking is preferred in moderate or low contention scenarios, while exclusive/shared locking is preferred in high contention scenarios (Sects. 7.6.1 and 7.6.2).
7. Optimizations in the centralized environment may not work well in RDMA networks due to a prohibitive number of primitive invocations, e.g., RDMA-T/O (Sect. 7.6.3).

## 8 Related work

To the best of our knowledge, this is the first to comprehensively evaluate the transaction scalability of concurrency control algorithms using RDMA-based primitives. Our study is related to previous works on 1) transaction scalability evaluations of concurrency control algorithms over TCP/IP networks and 2) optimizations for them using RDMA networks.

We have witnessed a wide spectrum of concurrency control algorithms to guarantee the serializable isolation level. Because there does not exist a single algorithm that can perform the best in all scenarios, figuring out the best algorithm for a specific application scenario naturally becomes an important problem. A few works [3, 8, 10, 11, 52, 54] make theoretical analyses over the benefits of different kinds of algorithms. More works are proposed to put concurrency control algorithms in the same centralized framework to do evaluations [13, 22, 31–33, 51, 62, 65]. Recently, there is an increasing interest in evaluating algorithms for distributed transaction processing [11, 12, 29]. The results show that these algorithms cannot achieve transaction scalability due to the limitations of slow network and coordination overhead. There are also several works [4, 16–18, 21, 23, 34, 39, 46, 48, 53, 59, 68] that focus on improving the system scalability. The intuitive idea is to transform distributed transactions into local transactions by carefully designing locality-aware partitioning approaches. Yet, static partitioning works only if the optimal data placement is known a priori, while dynamic partitioning often suffers from an expensive data migration overhead.

Using RDMA networks to optimize concurrency control algorithms has become a hotspot in both academia and industry. Most of them are tailored for some particular concurrency control algorithms. Dragojevi et al. [19, 20] implement a distributed in-memory database called FaRM and leverage one-sided verbs to optimize Silo. DrTM+H [60] also re-implements Silo and additionally proposes optimizations based on hybrid one-sided and two-sided verbs. NAM-DB [9] optimizes the snapshot isolation algorithm, which can only achieve the snapshot isolation level. To achieve the serializable isolation level in NAM-DB, a complete redesign of the concurrency control algorithm using RDMA is required. A few works [5, 14, 61, 64] concentrate on optimizations of 2PL algorithms. More related to our work, Wang et al. [58] builds a unified framework to re-implement and evaluate the concurrency control algorithms using RDMA. Yet, this framework lacks generality, making each algorithm implemented from scratch independently. Moreover, the evaluation in this framework mainly focuses on the performance with a fixed number the data nodes accessed per transaction. Different from the above works, we focus on transaction scalability under the serializable isolation level. We first explore the dominant factors to scale out distributed transaction processing and propose a general framework RCBench using one-sided RDMA verbs only, enjoying all benefits of RDMA networks. To enable the implementation of various algorithms, we introduce six primitives and five optimization principles. RCBench utilizes the TCP/IP implementation in Deneva [29] to compare TCP/IP algorithm variants with one-sided algorithm variants. Compared to Deneva, RCBench is built on a shared-memory

architecture on top of the RDMA one-sided network, whereas Deneva is a shared-nothing system entirely implemented on the TCP/IP network. In comparison with NAM-DB [9], both RCBench and NAM-DB are built on the RDMA one-sided network. However, RCBench adopts a shared-memory architecture, while NAM-DB uses a compute-storage separation architecture known as network-attached memory. Furthermore, based on RCBench, we provide an additional six primitives and five principles that enable developers to use one-sided RDMA verbs to enhance different algorithms.

## 9 Conclusions

In this paper, we investigate the problem of whether it is scalable to process distributed transactions using RDMA networks under serializable isolation level. We observe that  $\mathcal{U}_{\text{CPU}}^e$  is the dominant factor to scale out distributed transactions. To improve  $\mathcal{U}_{\text{CPU}}^e$ , we first propose a framework with six abstracted primitives using one-sided verbs. We then re-implement state-of-the-art concurrency control algorithms with various optimizations simply based on the primitives. Our implementations can enjoy all benefits of RDMA networks. Finally, we conduct a comprehensive experimental study of the transaction scalability of our implementations. We list a few findings that have not yet been reported elsewhere. We are confident that these findings provide a potential guideline to develop highly scalable distributed databases.

**Acknowledgements** The paper is supported by the National Natural Science Foundation of China under Grant No. 61972403, and the paper is supported by Public Computing Cloud, Renmin University of China.

## References

1. Abebe, M., Glasbergen, B., Daudjee, K.: Dynamast: adaptive dynamic mastering for replicated systems. In: 36th IEEE international conference on data engineering, ICDE 2020, Dallas, TX, USA, April 20–24, 2020, pp. 1381–1392. IEEE (2020). <https://doi.org/10.1109/ICDE48307.2020.00123>
2. Abebe, M., Glasbergen, B., Daudjee, K.: Morphosys: automatic physical design metamorphosis for distributed database systems. *Proc. VLDB Endow.* **13**(13):3573–3587 (2020). <https://doi.org/10.14778/3424573.3424578>
3. Agrawal, R., Carey, M.J., Livny, M.: Concurrency control performance modeling: alternatives and implications. *ACM Trans. Database Syst.* **12**(4), 609–654 (1987)
4. Agrawal, S., Narasayya, V.R., Yang, B.: Integrating vertical and horizontal partitioning into automated physical database design. In: SIGMOD Conference, pp. 359–370. ACM (2004)
5. Barthels, C., Müller, I., Taranov, K., Alonso, G., Hoeffler, T.: Strong consistency is not hard to get: two-phase locking and two-phase commit on thousands of cores. *Proc. VLDB Endow.* **12**(13), 2325–2338 (2019)
6. Bernstein, P.A., Goodman, N.: Timestamp-based algorithms for concurrency control in distributed database systems. In: VLDB, pp. 285–300. IEEE Computer Society (1980)

7. Bernstein, P.A., Goodman, N.: Concurrency control in distributed database systems. *ACM Comput. Surv.* **13**(2), 185–221 (1981)
8. Bhide, A., Stonebraker, M.: A performance comparison of two architectures for fast transaction processing. In: *International Conference on Data Engineering* (1988)
9. Binnig, C., Crotty, A., Galakatos, A., Kraska, T., Zamanian, E.: The end of slow networks: it's time for a redesign. *Proc. VLDB Endow.* **9**(7), 528–539 (2016)
10. Carey, M.J.: An abstract model of database concurrency control algorithms. In: *SIGMOD Conference*, pp. 97–107. ACM Press (1983)
11. Carey, M.J., Livny, M.: Distributed concurrency control performance: a study of algorithms, distribution, and replication. In: *14th International Conference on Very Large Data Bases* (1988)
12. Carey, M.J., Livny, M.: Parallelism and concurrency control performance in distributed database machines. In: *SIGMOD Conference*, pp. 122–133. ACM Press (1989)
13. Carey, M.J., Muhanna, W.A.: The performance of multiversion concurrency control algorithms. *ACM Trans. Comput. Syst.* **4**(4), 338–378 (1986)
14. Chen, Y., Wei, X., Shi, J., Chen, R., Chen, H.: Fast and general distributed transactions using RDMA and HTM. In: *EuroSys*, pp. 26:1–26:17. ACM (2016)
15. Cooper, B.F., Silberstein, A., Tam, E., Ramakrishnan, R., Sears, R.: Benchmarking cloud serving systems with YCSB. In: *SoCC*, pp. 143–154. ACM (2010)
16. Curino, C., Zhang, Y., Jones, E.P.C., Madden, S.: Schism: a workload-driven approach to database replication and partitioning. *Proc. VLDB Endow.* **3**(1), 48–57 (2010)
17. Das, S., Nishimura, S., Agrawal, D., Abbadi, A.E.: Albatross: lightweight elasticity in shared storage databases for the cloud using live data migration. *Proc. VLDB Endow.* **4**(8), 494–505 (2011)
18. Dashti, M., John, S.B., Shaikhha, A., Koch, C.: Transaction repair for multi-version concurrency control. In: *SIGMOD Conference*, pp. 235–250. ACM (2017)
19. Dragojevic, A., Narayanan, D., Castro, M., Hodson, O.: Farm: fast remote memory. In: *NSDI*, pp. 401–414. USENIX Association (2014)
20. Dragojevic, A., Narayanan, D., Nightingale, E.B., Renzelmann, M., Shamis, A., Badam, A., Castro, M.: No compromises: distributed transactions with consistency, availability, and performance. In: *SOSP*, pp. 54–70. ACM (2015)
21. Elmore, A.J., Arora, V., Taft, R., Pavlo, A., Agrawal, D., Abbadi, A.E.: Squall: Fine-grained live reconfiguration for partitioned main memory databases. In: *SIGMOD Conference*, pp. 299–313. ACM (2015)
22. Elmore, A.J., Das, S., Agrawal, D., Abbadi, A.E.: Zephyr: live migration in shared nothing databases for elastic cloud platforms. In: *SIGMOD Conference*, pp. 301–312. ACM (2011)
23. Elmore, A.J., Das, S., Agrawal, D., El Abbadi, A.: Towards an elastic and autonomic multitenant database. In: *Proc. of NetDB Workshop*. sn (2011)
24. Faleiro, J.M., Abadi, D.J.: Rethinking serializable multiversion concurrency control. *Proc. VLDB Endow.* **8**(11), 1190–1201 (2015)
25. Faleiro, J.M., Abadi, D.J., Hellerstein, J.M.: High performance transactions via early write visibility. *Proc. VLDB Endow.* **10**(5), 613–624 (2017). <https://doi.org/10.14778/3055540.3055553>
26. Faleiro, J.M., Thomson, A., Abadi, D.J.: Lazy evaluation of transactions in database systems. In: *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data, SIGMOD '14*, p. 15–26. Association for Computing Machinery, New York (2014). <https://doi.org/10.1145/2588555.2610529>
27. Gray, J., Lorie, R.A., Putzolu, G.R., Traiger, I.L.: Granularity of locks and degrees of consistency in a shared data base. In: *Readings in database systems* (3rd ed.) (1976)
28. Härder, T.: Observations on optimistic concurrency control schemes. *Inf. Syst.* **9**(2), 111–120 (1984)
29. Harding, R., Aken, D.V., Pavlo, A., Stonebraker, M.: An evaluation of distributed concurrency control. *PVLDB* **10**(5), 553–564 (2017)
30. Higuchi, K., Tsuji, T.: A linear hashing enabling efficient retrieval for range queries. In: *2009 IEEE International Conference on Systems, Man and Cybernetics*, pp. 4557–4562 (2009). <https://doi.org/10.1109/ICSMC.2009.5346783>
31. Huang, J., Stankovic, J.A., Ramamritham, K., Towsley, D.F.: Experimental evaluation of real-time optimistic concurrency control schemes. In: *VLDB*, pp. 35–46. Morgan Kaufmann (1991)
32. Huang, Y., Qian, W., Kohler, E., Liskov, B., Shrira, L.: Opportunities for optimism in contended main-memory multicore transactions. *Proc. VLDB Endow.* **13**(5), 629–642 (2020)
33. Jipping, M.J., Ford, R.: Predicting performance of concurrency control designs. In: *SIGMETRICS*, pp. 132–142. ACM (1987)
34. Jones, E.P.C.: Fault-tolerant distributed transactions for partitioned OLTP databases. Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, MA, USA (2012)
35. Lim, H., Kaminsky, M., Andersen, D.G.: Cicada: dependably fast multi-core in-memory transactions. In: *SIGMOD Conference*, pp. 21–35. ACM (2017)
36. Lin, Y.S., Tsai, C., Lin, T.Y., Chang, Y.S., Wu, S.H.: Don't look back, look into the future: prescient data partitioning and migration for deterministic database systems. In: *Proceedings of the 2021 International Conference on Management of Data, SIGMOD'21*, pp. 1156–1168. Association for Computing Machinery, New York (2021). <https://doi.org/10.1145/3448016.3452827>
37. Lu, Y., Yu, X., Cao, L., Madden, S.: Aria: A fast and practical deterministic oltp database. *Proc. VLDB Endow.* **13**(12), 2047–2060 (2020). <https://doi.org/10.14778/3407790.3407808>
38. Mitchell, C., Geng, Y., Li, J.: Using one-sided RDMA reads to build a fast, cpu-efficient key-value store. In: *USENIX Annual Technical Conference*, pp. 103–114. USENIX Association (2013)
39. Pavlo, A., Curino, C., Zdonik, S.B.: Skew-aware automatic database partitioning in shared-nothing, parallel OLTP systems. In: *SIGMOD Conference*, pp. 61–72. ACM (2012)
40. Peng, D., Dabek, F.: Large-scale incremental processing using distributed transactions and notifications. In: *OSDI*, pp. 251–264. USENIX Association (2010)
41. Qadah, T., Gupta, S., Sadoghi, M.: Q-store: Distributed, multi-partition transactions via queue-oriented execution and communication. In: *EDBT*, pp. 73–84. OpenProceedings.org (2020)
42. Qadah, T.M., Sadoghi, M.: Quecc: a queue-oriented, control-free concurrency architecture. In: *Middleware*, pp. 13–25. ACM (2018)
43. Qin, D., Brown, A.D., Goel, A.: Caracal: contention management with deterministic concurrency control. In: *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles, SOSP'21*, pp. 180–194. Association for Computing Machinery, New York (2021). <https://doi.org/10.1145/3477132.3483591>
44. Quamar, A., Kumar, K.A., Deshpande, A.: SWORD: scalable workload-aware data placement for transactional workloads. In: *EDBT*, pp. 430–441. ACM (2013)
45. Rosenkrantz, D.J., Stearns, R.E., II, P.M.L.: System level concurrency control for distributed database systems. *ACM Trans. Database Syst.* **3**(2), 178–198 (1978)
46. Schiller, O., Cipriani, N., Mitschang, B.: Prorea: live database migration for multi-tenant RDBMS with snapshot isolation. In: *EDBT*, pp. 53–64. ACM (2013)
47. Serafini, M., Taft, R., Elmore, A.J., Pavlo, A., Aboulmaga, A., Stonebraker, M.: Clay: Fine-grained adaptive partitioning for general database schemas. *Proc. VLDB Endow.* **10**(4), 445–456 (2016). <https://doi.org/10.14778/3025111.3025125>



48. Shute, J., Vingralek, R., Samwel, B., Handy, B., Whipkey, C., Rollins, E., Oancea, M., Littlefield, K., Menestrina, D., Ellner, S., Cieslewicz, J., Rae, I., Stancescu, T., Apte, H.: F1: a distributed SQL database that scales. *Proc. VLDB Endow.* **6**(11), 1068–1079 (2013)
49. Stamos, J., Cristian, F.: A low-cost atomic commit protocol. In: *Proceedings 9th Symposium on Reliable Distributed Systems*, pp. 66–75 (1990). <https://doi.org/10.1109/RELDIS.1990.93952>
50. Taft, R., Mansour, E., Serafini, M., Duggan, J., Elmore, A.J., Aboulnaga, A., Pavlo, A., Stonebraker, M.: E-store: Fine-grained elastic partitioning for distributed transaction processing systems. *Proc. VLDB Endow.* **8**(3), 245–256 (2014). <https://doi.org/10.14778/2735508.2735514>
51. Tanabe, T., Hoshino, T., Kawashima, H., Tatebe, O.: An analysis of concurrency control protocols for in-memory databases with cbench. *Proc. VLDB Endow.* **13**, 3531–3544 (2020)
52. Thomasian, A.: Concurrency control: methods, performance, and analysis. *ACM Comput. Surv.* **30**(1), 70–119 (1998)
53. Thomson, A., Diamond, T., Weng, S., Ren, K., Shao, P., Abadi, D.J.: Calvin: fast distributed transactions for partitioned database systems. In: *SIGMOD Conference*, pp. 1–12. ACM (2012)
54. Thuraisingham, B., Ko, H.: Concurrency control in trusted database management systems: a survey. *SIGMOD Rec.* **22**(4), 52–59 (1993)
55. TPC-C: <http://www.tpc.org/tpcc/> (1988)
56. Tran, K.Q., Naughton, J.F., Sundarmurthy, B., Tsirogiannis, D.: Jecb: a join-extension, code-based approach to oltp data partitioning. In: *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data, SIGMOD'14*, pp. 39–50. Association for Computing Machinery, New York (2014). <https://doi.org/10.1145/2588555.2610532>
57. Tu, S., Zheng, W., Kohler, E., Liskov, B., Madden, S.: Speedy transactions in multicore in-memory databases. In: *SOSP*, pp. 18–32. ACM (2013)
58. Wang, C., Qian, X.: Rdma-enabled concurrency control protocols for transactions in the cloud era. *IEEE Trans. Cloud Comput.* **PP**, 1–1 (2021)
59. Wang, T., Kimura, H.: Mostly-optimistic concurrency control for highly contended dynamic workloads on a thousand cores. *Proc. VLDB Endow.* **10**(2), 49–60 (2016)
60. Wei, X., Dong, Z., Chen, R., Chen, H.: Deconstructing rdma-enabled distributed transactions: hybrid is better! In: *OSDI*, pp. 233–251. USENIX Association (2018)
61. Wei, X., Shi, J., Chen, Y., Chen, R., Chen, H.: Fast in-memory transaction processing using RDMA and HTM. In: *SOSP*, pp. 87–104. ACM (2015)
62. Wu, Y., Arulraj, J., Lin, J., Xian, R., Pavlo, A.: An empirical evaluation of in-memory multi-version concurrency control. *Proc. VLDB Endow.* **10**(7), 781–792 (2017)
63. Yao, C., Agrawal, D., Chen, G., Lin, Q., Ooi, B.C., Wong, W., Zhang, M.: Exploiting single-threaded model in multi-core in-memory systems. *IEEE Trans. Knowl. Data Eng.* **28**(10), 2635–2650 (2016)
64. Yoon, D.Y., Chowdhury, M., Mozafari, B.: Distributed lock management with RDMA: decentralization without starvation. In: *SIGMOD Conference*, pp. 1571–1586. ACM (2018)
65. Yu, X., Bezerra, G., Pavlo, A., Devadas, S., Stonebraker, M.: Starving into the abyss: an evaluation of concurrency control with one thousand cores. *Proc. VLDB Endow.* **8**(3) (2014)
66. Zamanian, E., Binnig, C., Harris, T., Kraska, T.: The end of a myth: Distributed transactions can scale. *Proc. VLDB Endow.* **10**(6), 685–696 (2017). <https://doi.org/10.14778/3055330.3055335>
67. Zamanian, E., Binnig, C., Salama, A.: Locality-aware partitioning in parallel database systems. In: *SIGMOD Conference*, pp. 17–30. ACM (2015)
68. Zhao, Z.: Efficiently supporting adaptive multi-level serializability models in distributed database systems. In: *SIGMOD Conference*, pp. 2908–2910. ACM (2021)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.