

Transaction Scheduling: From Conflicts to Runtime Conflicts

YANG CAO, University of Edinburgh, UK

WENFEI FAN, Shenzhen Institute of Computing Sciences, China, University of Edinburgh, UK, and BDBC, Beihang University, China

WEIJIE OU, Shenzhen Institute of Computing Sciences, China

RUI XIE, Shenzhen Institute of Computing Sciences, China

WENYUE ZHAO, University of Edinburgh, UK

This paper studies how to improve the performance of main memory multicore OLTP systems for executing transactions with conflicts. A promising approach is to partition transaction workloads into mutually conflict-free clusters, and distribute the clusters to different cores for concurrent execution. We show that if transactions in each cluster are properly scheduled, transactions that are traditionally considered conflicting can be executed without conflicts at runtime. In light of this, we propose to schedule transactions and reduce runtime conflicts, instead of partitioning based on the conventional notion of conflicts. We formulate the transaction scheduling problem to minimize runtime conflicts, and show that the problem is NP-complete. This said, we develop an efficient scheduling algorithm to improve parallelism. Moreover, for transactions that are not packed in batches, we show that runtime conflict analysis also helps reduce conflict penalties, by proposing a proactive deferring method. Using standard and enhanced benchmarks, we show that on average our scheduling and proactive deferring methods improve the throughput of existing partitioners and concurrency control protocols by 131% and 109%, respectively, up to 294% and 152%.

CCS Concepts: • **Information systems** → **Database transaction processing**; **Database management system engines**.

Additional Key Words and Phrases: transaction scheduling, concurrency control, transaction processing

ACM Reference Format:

Yang Cao, Wenfei Fan, Weijie Ou, Rui Xie, and Wenyue Zhao. 2023. Transaction Scheduling: From Conflicts to Runtime Conflicts. *Proc. ACM Manag. Data* 1, 1, Article 26 (May 2023), 26 pages. <https://doi.org/10.1145/3588706>

1 INTRODUCTION

There has been increasing demand on high volume of concurrent transactions from *e.g.*, e-commerce, FinTech and cloud applications [30]. This and the growing dominance of multicore machines highlight the need for pursuing higher throughput and parallelism of multi-thread transaction processing [5, 33, 40, 58]. Due to contended operations that read and write the same data items, transaction execution has to be guarded against concurrency anomalies to uphold the desired isolation levels. There are mainly two types of approaches to maximizing concurrency while providing isolation guarantees: (a) partition-based approaches [15, 22, 33, 36, 40, 47] that use transaction partitioners to decide a transaction-to-thread assignment for a batch of transactions before their execution; and (b) concur-

Authors' addresses: Yang Cao, yang.cao@ed.ac.uk, University of Edinburgh, UK; Wenfei Fan, wenfei@inf.ed.ac.uk, Shenzhen Institute of Computing Sciences, China and University of Edinburgh, UK and BDBC, Beihang University, China; Weijie Ou, ouweijie@sics.ac.cn, Shenzhen Institute of Computing Sciences, China; Rui Xie, xierui@sics.ac.cn, Shenzhen Institute of Computing Sciences, China; Wenyue Zhao, wenyue.zhao@ed.ac.uk, University of Edinburgh, UK.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

2836-6573/2023/5-ART26 \$15.00

<https://doi.org/10.1145/3588706>

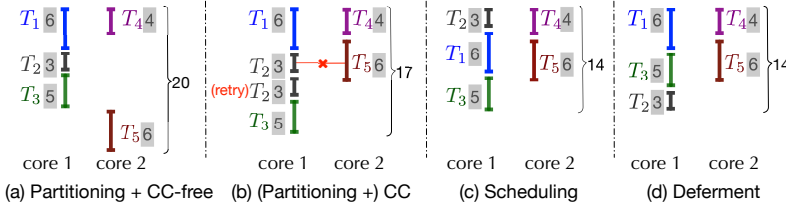


Fig. 1. Different transaction executions in Examples 1-2

rency control (CC) based methods [10, 26, 40, 51, 59] that, instead of searching for good transaction-to-thread assignments, focus on CC protocols to resolute contended transactions on-the-fly.

Transaction partitioning methods focus on a bundle of transactions that are *e.g.*, a set of transaction collected in a fixed duration [36] or workloads [22] for which the transaction logic is known in advance (*e.g.*, stored procedures and hard-coded templates), as targeted by deterministic databases [5, 38]. Such workloads allow the system to carry out analysis over the transactions, *i.e.*, transaction partitioning, to decide the best way of assigning transactions to threads. More specifically, given a workload \mathcal{W} over k threads, a transaction partitioning method computes a partitioning (P_1, P_2, \dots, P_k) of \mathcal{W} that assigns each P_i to a distinct thread i such that transactions in the same partition are executed serially, while concurrent execution is governed by CC for the desired isolation levels. Its objective is to minimize conflicts among transactions across different partitions to reduce CC cost, while maximizing load balance so as to minimize total parallel execution time.

Instead of searching for the best thread assignment for the transactions, CC aims to resolute conflicts between transactions that are concurrently executed, and do not require the access sets of transactions before execution. There are two types of CC cost: (a) CC overhead incurred on every transaction and (b) conflict penalty when a conflict between transactions happens at runtime, *e.g.*, block (waiting) time with locking-based CC [11, 18] or abort/retry penalty with optimistic concurrency control (OCC) protocols [6, 26]. In practice, CC protocols with higher overhead are more likely to reduce conflict penalties. Hence, the effectiveness of CC relies on delicate trade-offs between the overhead and the ability to reduce conflict penalties.

We refer to transactions targeted by the partitioning and CC approaches as *bundled* transactions and *unbundled* ones, respectively.

Runtime conflict. Conflicts are conventionally defined relative to the isolation level adopted by transaction systems. For serializability [21], transactions T and T' are in conflict if they access the same data item and at least one of them updates it. For snapshot isolation [8], T and T' are conflicting if they write to the same data item.

However, transactions that are conventionally considered in conflicts can be executed concurrently without conflicts at runtime as long as they are scheduled properly, as illustrated below.

Example 1: Consider a set \mathcal{W}_0 of transactions $\{T_1, T_2, T_3, T_4, T_5\}$:

- $T_1 = R[x_2]W[x_2]R[x_3]W[x_3]R[x_4]W[x_4]$,
- $T_2 = R[x_1]W[x_2]W[x_1]$,
- $T_3 = R[x_3]W[x_3]R[x_2]R[x_3]W[x_2]$,
- $T_4 = R[x_5]W[x_5]R[x_6]W[x_6]$, and
- $T_5 = R[x_1]W[x_1]R[x_5]W[x_5]R[x_1]W[x_1]$,

where $R[x_1]$ denotes a read of data item x_1 and $W[x_1]$ is a write to x_1 . Assume that the targeted

isolation level is serializability. Then T_1, T_2 and T_3 are in conflict; similarly for (T_2, T_5) and (T_4, T_5) .

Assume that we have two cores. Following [36], we partition \mathcal{W}_0 into $P_1 = \{T_1, T_2, T_3\}$ and $P_2 = \{T_4\}$, along with T_5 as a cross-partition transaction, as shown in Fig. 1(a). One can execute P_1 and P_2 concurrently without CC, followed by T_5 after both P_1 and P_2 are completed. Assuming that each of read and write operation takes a unit time; then the makespan of such execution (the concurrent execution of the transactions) of \mathcal{W}_0 is 20 time units. Note that the workloads P_1 and P_2 are not balanced, and the second core may idle for long before the cross-partition transaction T_5 can start.

Alternatively, by following [15, 22], one can start T_5 as soon as T_4 completes, as shown in Fig. 1(b); however, this requires CC since T_5 is in conflict with T_2 that is being concurrently executed at the other core. When using OCC protocols, this may cause a retry of T_2 and a makespan of 17 time units for \mathcal{W}_0 (ignoring CC overhead and taking into only the execution time of all transactions).

In contrast, if we schedule the transactions as two *queues* $Q_1 = \langle T_2, T_1, T_3 \rangle$ and $Q_2 = \langle T_4, T_5 \rangle$, i.e., imposing an execution order of transactions in the partitions, we can process \mathcal{W}_0 at two cores concurrently *without* CC as shown in Fig. 1(c), by executing Q_1 and Q_2 in their orders at the two cores independently. Although T_2 of Q_1 is in conflict with T_5 of Q_2 , their executions actually *do not overlap*. In fact, the executions are serializable even without CC. Note that the workloads are more balanced than the case with partitioning, and the makespan of such an execution of \mathcal{W} is 14 instead of 20. \square

This example suggests that we consider *runtime conflicts* instead for partition-based approaches, which provide a more accurate characterization of transaction executions and more opportunities to improve concurrency. Moreover, runtime conflicts could also help with CC-based approaches for unbundled transactions, as shown by the following example.

Example 2: Consider the execution of transactions in Example 1 following a CC-based approach: we execute T_1, T_2 and T_3 at core 1 and T_4, T_5 at core 2 in order, using OCC. Then T_2 retries due to conflict with T_5 , causing a total execution time of 17 (as shown Fig. 1(b)).

If prior to the start of T_2 , core 1 could find that T_5 is being executed at core 2, which will incur conflict with T_2 , then a better option is to “defer” T_2 and execute T_3 first instead. In this way, T_3, T_5 and T_2 can all commit without retry with a total time of 14 instead of 17 as shown in Fig. 1(d). That is, we can reduce runtime conflicts during execution for CC-based approaches, by proactively deferring transactions and changing their execution order on-the-fly. \square

Transaction scheduling. For bundled transactions that are handled by existing partitioners, we propose to schedule transactions so as to minimize runtime conflicts. Given a workload \mathcal{W} that has been partitioned over k cores, we generate k queues (Q_1, \dots, Q_k) and a residual set \mathcal{R} , such that Q_i ’s consist of transactions in \mathcal{W} without incurring runtime conflicts, and hence can even be executed without CC. As opposed to partitioning, transactions in Q_i are *scheduled*, by placing an ordering on the executions of their transactions. Transactions in \mathcal{R} are executed over all cores using CC.

Intuitively, transactions in Q_i and those in Q_j ($i \neq j$) have *no runtime conflict* if they are executed on schedule, although they may be conventionally considered in conflict; hence transactions in Q_i run *serially* at the i -th core, and *in parallel* with those in Q_j ($i \neq j$) without conflicts. Transactions in \mathcal{R} may inflict runtime conflict with other transactions in \mathcal{R} or Q_i , and hence are executed at all cores with CC. As shown in Example 1, when compared to transaction partitioning, scheduling allows more transactions to be put in Q_1, \dots, Q_k and executed concurrently without conflicts.

We show that transaction scheduling is intractable. This is not surprising since it is more intriguing than transaction partitioning, which is already NP-hard [36]. This said, we develop an efficient algorithm that is able to either refine a transaction partition into a schedule, or compute a schedule from scratch when a partition is not available. Scheduling targets bundled transactions

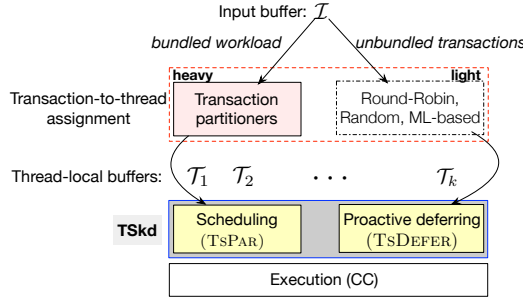


Fig. 2. A typical transaction system with TS_{KD}: TS_{KD} works as a plug-in for existing transaction systems that either use transaction partitioners to assign transactions to threads or directly rely on CC protocols. It sits between the transaction-to-thread assignment module and transaction execution engines to reduce runtime conflicts.

to which the partitioning methods are applicable. It not only reduces runtime conflicts but also balances workloads among the threads, improving the throughput. It is particularly effective for transactions with skewed costs or I/O latency, which existing partitioners do not handle very well.

Proactive transaction deferment. For unbundled transactions handled by CC-based approaches, we develop a lightweight method to reduce their CC cost by proactively deferring transactions on-the-fly that are on course to inflict runtime conflicts with other transactions being executed. Unlike partitioners that are often used for preprocessing batched transactions prior to their execution, CC is part of the execution phase and hence proactive deferment has to be at a low cost since otherwise the benefit of reduced runtime conflicts could be canceled out by the increased CC overhead.

To this end, we develop a lock-free structure that tracks the progresses of transaction execution of all the threads, which allows us to efficiently detect whether a runtime conflict would happen before executing a new transaction at runtime. The lock-free design avoids data racing overhead when threads update and look up execution progresses concurrently. In addition, it has a parametric complexity to trade the detection overhead for reduced conflict penalties and hence higher throughput.

Prototype and evaluation. As a proof of concept, we develop TS_{KD}, a lightweight tool as shown in Fig. 2 for improving existing partitioning-based or CC-based transaction systems via scheduling (module Ts_{PAR}) and proactive deferment (module Ts_{DEFER}).

Given a workload \mathcal{W} for partitioners, TS_{KD} (Ts_{PAR}) first learns rough runtime estimates of the transactions in \mathcal{W} via execution histories, partial dry-runs used by partitioners [5], or more sophisticated cost estimators [12, 44, 48, 53]. It then converts the partitioning of \mathcal{W} into a schedule, *i.e.*, k queues (Q_1, \dots, Q_k) and residual \mathcal{R} , to reduce runtime conflicts and improve load balance. The transactions in Q_1, \dots, Q_k can be concurrently executed as scheduled without CC if the estimates are accurate. A subtle issue arises when time estimates are not perfectly accurate; if so some transactions may end up inflicting runtime conflicts. To cope with this, TS_{KD} uses both CC and proactive transaction deferment to execute the scheduled queues and \mathcal{R} . This guarantees that TS_{KD} always executes transactions correctly, while still benefiting from reduced conflicts and balanced workload among the threads via scheduling.

For unbundled transactions that are executed directly using CC, TS_{KD} employs proactive transaction deferment (Ts_{DEFER}) to reduce conflict penalties. Ts_{DEFER} detects potential conflicts between concurrent transactions and it imposes no restrictions on the workload that CC schemes deal with.

To allow a fair and repeatable evaluation, we integrate TS_{KD} into DBx1000, a popular open-source testbed for in-memory transaction processing with built-in implementation of major CC

schemes [2, 58]. Using TPC-C and YCSB, we verify that TSKD improves popular partitioners and CC protocols by 131% and 109%, respectively, and reduces their retries by 45.3% and 45.7%, up to 294% and 152% for throughput and 61.1% and 63.9% for retry reduction. Moreover, TSKD remains effective for transactions with varying degrees of long-tail I/O latency, making existing partitioners and CC schemes robust against *e.g.*, databases hosted in disks. We observe that TSKD is less effective for extremely short-run transactions as conflict penalties play a much smaller role in their performance.

Contributions & organization. To summarize, in this work we develop techniques that improve both transaction partitioners and CC schemes by reducing runtime conflicts. More specifically:

- We formalize a notion of runtime conflicts and propose transaction scheduling to improve transaction partitioners by making conflicting transactions conflict-free at runtime (Section 2.2).
- We propose proactive transaction deferment for CC protocols to reduce runtime conflicts during execution time (Section 2.3).
- We develop TSKD, a tool that implements both methods. It serves as a plugin for transaction systems that use either transaction partitioners or CC protocols directly, providing delicate controls over and flexible trade-offs between transaction retry penalties and CC overheads (Section 3)
- We settle the complexity of transaction scheduling and develop an efficient algorithm (Section 4).
- We develop a lock-free structure that enables proactive transaction deferment (Section 5).
- We empirically verify that TSKD improves the throughput of both partitioners and CC schemes, and makes them robust against large skewed I/O latency and runtime (Section 6).

We remark that TSKD is not to replace existing CC schemes or transaction partitioners. Instead, it is positioned as a tool to expose opportunities to improve their performance in a non-intrusive way.

We discuss related work in Section 7 and future work in Section 8.

2 REDUCING RUNTIME CONFLICTS FOR TRANSACTIONS

In this section, we first review the basics of transaction processing (Section 2.1). We then present the notion of runtime conflicts, based on which we propose transaction scheduling (Section 2.2) and proactive transaction deferring (Section 2.3).

2.1 Transaction Preliminaries

A transaction is a sequence of database actions that are to be executed as atomic work units, including reads from and writes to a database. The tuples in the database that are read (resp. written) by a transaction T are referred to as the *read* (resp. *write*) set of T .

Conflicts. Two transactions T and T' are *in conflict w.r.t.* a CC protocol ρ if they contain contended operations under ρ . Contention is defined relative to the particular isolation level [8] upheld by ρ . For instance, if ρ enforces serializability, then T and T' are in conflict if they both access (read or write) the same data item x and at least one of them writes x . If ρ enforces snapshot isolation, then T and T' are in conflict if they both write the same data item. As an example, under serializability, T_2 and T_5 of Example 1 are in conflict; however, they do not conflict under snapshot isolation.

Transaction T is *conflict-free* with T' if they are not in conflict. Conflicting transactions require CC to ensure the correctness at the particular isolation level upheld by the CC protocol ρ .

Workload model. As shown in Fig. 2, we consider multi-thread transaction systems that either use transaction partitioners for bundled transactions or directly use CC for unbundled ones. Indeed, the majority of transaction systems fall into these two categories. They consist of (a) input buffer I that receives transactions to be processed; and (b) thread-local buffers $\mathcal{T}_1, \dots, \mathcal{T}_k$, where each

\mathcal{T}_i ($i \in [1, k]$) contains the transactions to be executed by thread i .

For *bundled* workloads \mathcal{W} , i.e., a set of transactions revealed to \mathcal{I} all at once prior to execution, the practice is to first decide an assignment for transactions in \mathcal{W} to the k thread-local buffers, by partitioning them based on their read/write sets [5, 15, 22, 33, 36, 40, 47]. All threads then execute assigned transactions in their local buffers concurrently. By spending more preprocessing time, partitioners reduce transaction execution time with fewer aborts and less CC cost.

For transactions that are coming *unbundled* in \mathcal{I} , they are periodically flushed to the thread-local buffers via much lighter method than transaction partitioning, e.g., round-robin, random or lightweight ML-based assignment methods [43]. All transactions are executed concurrently with CC. Hence, the impact of CC is greater on unbundled transactions when compared to bundled workloads.

Concurrency control (CC). When transactions in the local buffers are executed concurrently, conflicting transactions may cause anomalies that the desired isolation level prohibits, impairing the correctness of the execution results. To this end, transaction systems use CC protocols [10, 21, 26, 40, 51, 59] to prevent such anomalies. The cost of CC consists of (a) the overhead of the CC protocols, which is charged to every transaction even when it involves in no conflict, and (b) the cost of preventing anomalies when conflicts happen, e.g., transaction blocking (waiting) time for locking-based CC and abort/retry costs for optimistic CC methods.

Transaction partitioning. The partitioning methods assign bundled transactions (workload) to threads for concurrent execution by analyzing the conflicts between transactions. It is typically used as a preprocessing method and incurs computation cost [5].

Conceptually, the process that transaction partitioners [5, 15, 35] partition a workload \mathcal{W} can be modeled as a graph cut problem: by finding a minimum k -cut of the *conflict* graph G of \mathcal{W} , where vertices are transactions and an edge (T, T') indicates that T conflicts with T' , one breaks \mathcal{W} into k partitions, where the edge cut denotes conflicts between partitions. Each partition is executed serially by a dedicated thread and partitions are executed concurrently with CC.

One can also gathers all conflicting transactions into a single set \mathcal{R} such that the k partitions become mutually conflict-free. We refer to \mathcal{R} as a *residual* partition for \mathcal{W} and the remaining k sets of transactions as *CC-free* partitions. It has been shown that by executing CC-free partitions without CC first, and then executing \mathcal{R} with CC after all CC-free partitions complete, one can achieve higher throughput for high contention transactions [36].

In this setting, the objective of partitioning is to minimize the size of residual set (or k -cut) while balancing the partitions. For convenience, we represent a partitioning of \mathcal{W} simply as $(P_1, \dots, P_k, \mathcal{R})$, where P_i 's are CC-free partitions and \mathcal{R} is the residual.

2.2 From Transaction Partitioning to Scheduling

We next propose transaction scheduling based on runtime conflicts.

Execution time. Transactions in practice often have varying execution time. The *execution time* of a transaction T , denoted by $\text{time}(T)$, is the serial execution time of T by a single thread, from the start to the completion of T . It measures the duration (length) of T .

A variety of methods have been explored to estimate $\text{time}(T)$, from a brute-force one that counts reads and writes in T (e.g., Example 1), to advanced statistical and ML strategies [12, 44, 48, 53].

For a set P of transactions, $\text{time}(P)$ denotes the total execution time of transactions in P , assuming that they are executed serially.

Transaction schedule. A *schedule* of a transaction workload \mathcal{W} over k threads, denoted by $\mathcal{W}_f^<$, is a pair $(f, <)$, where

- f is a function that partitions \mathcal{W} into $k+1$ disjoint sets Q_1, \dots, Q_k and \mathcal{R}_s such that (a) $\bigcup_{i=1}^k Q_i \cup \mathcal{R}_s = \mathcal{W}$ and (b) T and T' are not in runtime-conflict for any $T \in Q_i, T' \in Q_j$ and $i \neq j$; and
- $<$ is an order relation on \mathcal{W} such that for $i \in [1, k]$ and each Q_i , $<$ is a total order on Q_i .

We refer to Q_i as a *conflict-free queue*.

Similar to conventional transaction partitioners, function f clusters transactions in \mathcal{W} by assigning them to threads. In contrast to partitioning, relation $<$ orders transactions assigned to each of the k queues. Moreover, the partition is guided by runtime conflicts (see below). As a consequence, the k conflict-free queues may include residual transactions for partitioning, and \mathcal{R}_s only contains transactions that incur runtime conflicts with those in the queues.

Given a schedule $(f, <)$ over k threads, each core serially executes transactions assigned to it by f , in the order specified by $<$.

Runtime conflict. For a transaction $T \in \mathcal{W}$ that is assigned to queue Q_i , denote by $\text{ts}(T)$ the *scheduled start time* of T by schedule $(f, <)$. It is defined as $\sum_{T' \in \text{pred}(T)} \text{time}(T')$, where $\text{pred}(T)$ is the set of transactions in \mathcal{W} that are assigned to the same Q_i as T by f and are ordered prior to T by $<$ (recall that $\text{time}(T)$ is the serial execution time of T). Similarly, $\text{tc}(T)$ denotes the *scheduled complete time* of T by $(f, <)$ and is defined as $\text{ts}(T) + \text{time}(T)$. The range $[\text{ts}(T), \text{tc}(T)]$ is called the *scheduled runtime* of T .

Two transactions T and T' are *in conflict at runtime* by schedule $(f, <)$ if they are in conflict and their scheduled runtime overlap. If they are not in conflict at runtime, T and T' are called *RC-free*.

Consider a conventional transaction partitioner that splits workload \mathcal{W} into a partition plan $f_c = (P_1, \dots, P_k, \mathcal{R})$. Let $(f, <)$ be a schedule of \mathcal{W} over k threads. Then we say $(f, <)$ *refines* partitioning f_c if P_i is a subset of Q_i partitioned by f for all $i \in [1, k]$.

Intuitively, the schedule is more tolerant than the partition by including conventionally conflicting transactions in the RC-free queues. As a consequence, \mathcal{R}_s is a subset of residual \mathcal{R} , and hence more transactions can be executed concurrently.

Example 3: Continuing with Example 1, the queues Q_1 and Q_2 given there form a schedule for the workload \mathcal{W}_0 . Moreover, the schedule refines the partitioning of \mathcal{W}_0 that consists of P_1 and P_2 and T_5 . As shown there, the scheduling reduces the makespan of the execution of \mathcal{W}_0 from 20 to 14, and improves the throughput. \square

2.3 Proactively Deferring Conflicting Transactions

For residual transactions or unbundled transactions that are directly executed using CC, we propose another method, called *proactive transaction deferment*, to further reduce conflict penalties for CC.

Recall that each thread i executes transactions in its local buffer \mathcal{T}_i using CC so that anomalies can be prevented when conflicting transactions are being executed concurrently. Proactive deferment works by detecting whether runtime conflicts would happen when a thread i is about to execute the next transaction T in \mathcal{T}_i after committing the current one; if so, it defers T and skips to the next transaction for execution. As will be seen in Section 5, the complexity of conflict detection is parameterized with knobs, to provide flexible trade-offs between the overhead of conflict detection and conflict penalties if runtime conflicts are not prevented in time.

For instance, as shown in Example 2, when core 1 detects that executing T_2 would inflict a conflict with T_5 being executed at core 2, it defers T_2 and skips to T_3 , reducing retry cost of T_2 .

We remark that proactive deferring is not a replacement for CC since it does not aim to detect and prevent *all* the runtime conflicts. Instead, it serves as a lightweight filter of transactions in the thread-local buffers before passing them to the transaction execution engine. It makes CC more

robust against transactions with varying complexities by virtue of its parametric cost.

3 A PROTOTYPE TRANSACTION SCHEDULER

In this section, we present an overview of TSKD, a lightweight tool for transactions to reduce runtime conflicts and CC penalties.

As shown in Fig. 2, TSKD works with existing transaction systems by serving as an intermediate layer between the transaction-to-thread assignment component and the execution engine. It consists of two components: (a) a *transaction scheduling* (TsPAR) module that schedules and optimizes the partitioning generated by an existing transaction partitioner, and (b) a *proactive transaction deferment* (TsDEFER) module that optimizes the thread-local buffers on-the-fly.

Scheduling (TsPAR module). Given a transaction workload \mathcal{W} and a partitioning of \mathcal{W} , TsPAR computes a schedule $(f, <)$ that turns the partitioning into k RC-free queues (Q_1, \dots, Q_k) and a residual \mathcal{R}_s . TsPAR then assigns transactions in Q_i to the local buffer \mathcal{T}_i of thread i in the order; all k threads can execute the assigned transactions in parallel without CC if time estimates are accurate.

After all threads finish their local transactions, transactions in \mathcal{R}_s are then executed with CC by all threads (with TsDEFER enabled; more below). By default, TsPAR assigns transactions to the thread by round-robin for its simplicity and load balance. One can also use other lightweight transaction-to-thread assignment methods supported by the underlying systems, e.g., random or ML-based [43].

To decide runtime conflicts for schedules, TSKD estimates the execution cost of the transactions in \mathcal{W} by following e.g., [12, 44, 48, 53]. TsPAR *does not* rely on the actual transaction execution time; instead, it is only sensitive to the relative length of transactions. Hence, any estimates that roughly preserve the relative costs of transactions suffice. By default, TsPAR uses execution histories to coarsely estimate costs; if T is instantiated with the same parameters as T' in the history, and if T and T' are based on the same template (e.g., stored procedure), then TsPAR uses the cost of T' as an estimate for T . When no T' has the same parameters as T , TsPAR picks a T' with parameters close to that of T as a coarse estimate. For cases when execution histories are not available, TsPAR adopts the partial dry-run approach used by e.g., deterministic databases [5] to generate the estimates. The idea is to execute (samples) of transactions partially such that no writes are physically executed during the dry-run; this has also been shown effective to deduce access sets of transactions for partitioning [5]. In the extreme case, TSKD uses the sizes of the access sets of the transactions as a fallback.

To cope with inaccurate or even missing estimates, TSKD uses CC and TsDEFER (more below) to guard the execution of RC-free transactions against potentially overlooked conflicts. This ensures that TSKD upholds the desired isolation level no matter how bad estimates can be. As will be shown in Section 6, TSKD still benefits from reduced conflicts and more balanced workload among the cores via scheduling, and yields higher throughput, especially for transactions with skewed runtime.

We will study the core problem underlying TsPAR in Section 4.

Proactive deferment (TsDEFER module). For unbundled transactions that are directly assigned to the thread-local buffers upon arrival or transactions of which runtime estimates are unavailable, TSKD uses TsDEFER to reduce conflict penalties by further reducing runtime conflicts. TsDEFER works by altering the ordering of transactions in the thread-local buffers during execution time, via carefully designed lightweight operations. It works as a dynamic filter of the transactions in the thread-local buffers and passes transactions that are not likely to cause runtime conflicts to the execution engine by deferring problematic transactions.

In a nutshell, TsDEFER checks, prior to passing a transaction T for execution (with CC) at thread i , whether T is in conflict with transactions that are (very likely) being or will be immediately executed by some other thread. If so, it defers the execution of T by moving T to the end of the transaction

queue in the local buffer for thread i ; otherwise it passes T to the transaction execution engine.

We will discuss the implementation of TsDEFER in Section 5.

Remarks. (1) TSKD aims to improve existing transaction systems, partition-based [5, 15, 17, 33, 35, 36, 39, 40, 47, 56, 61], CC-based [2, 6, 9, 11, 16, 18, 19, 26, 28, 32, 42, 46, 54, 55, 59, 60] or hybrid [45, 47]. It neither imposes any new restrictions on nor makes assumptions about how these systems work. This enables existing systems to deploy TSKD without changes to benefit from the reduced runtime conflicts by transaction scheduling and deferment, via reduced CC overhead and conflict penalties.

(2) TSKD also uses TsDEFER to reduce the execution cost of the unscheduled residual \mathcal{R}_s returned by TsPAR for bundled transactions.

(3) Neither TsPAR nor TsDEFER of TSKD is fixed to a specific isolation level such as serializability; instead, they work with arbitrary isolation levels that the underlying systems uphold, by observing conflicts according to the preferred isolation level.

Limitations. Because of its non-intrusive design, TSKD (TsPAR and TsDEFER) inherits some limitations of the systems it optimizes.

(1) *Access sets.* Transaction partitioners require that the access set of transactions is known upfront so that they could deduce conflicts between transactions. TsPAR inherits this prerequisite to schedule the partitioning of bundled transactions. As a result, no random client-driven transactions are expected for both partitioners and TsPAR; instead, they target transactions in the form of stored procedures or hard-coded templates, which are prevalent in *e.g.*, banking, e-commerce, business applications. Due to the importance of transaction partitioning, techniques have been developed to identify access sets in various scenarios [5, 38]. Moreover, for unbundled transactions where access sets are not known beforehand, we can fall back to CC-based approaches and use TsDEFER instead.

Another limitation inherited is that TSKD executes range queries directly with CC, since partitioners do not optimize range queries for which read/write-sets are not available. On the positive side, TsPAR inherits and reuses the conflict graphs constructed by transaction partitioners [15, 36] without the need for reconstruction, as will be seen in Section 4.

(2) *Application-specified dependencies.* Similar to CC-based approaches for unbundled transactions, TsDEFER operates on transactions visible during execution time only; hence, they do not have control on the global order of transaction execution to enforce transaction (*e.g.*, causal) dependencies implied by application logic. A possible way of mitigating this is to integrate consistency protocols with CC so that we have both isolation and causal consistency guarantees; TsDEFER could be extended for such cases by only deferring a transaction when a limited number of transactions depend on it.

Different from CC and TsDEFER, transaction partitioners and TsPAR can readily incorporate transaction dependencies by enforcing dependencies in partitions and during scheduling.

(3) *Generalization.* In principal TsPAR is not limited to the in-memory setting; it can be applied to shared-nothing distributed systems. In contrast, TsDEFER cannot be trivially generalized as it relies on lightweight probing operations to detect runtime conflict at execution time; such operations will incur too much overhead in the shared-nothing architecture due to network latency involved.

4 TRANSACTION SCHEDULING

In this section, we formulate the transaction scheduling problem, settle its complexity bound, and develop an efficient algorithm for it. The results serve as the foundation of the TsPAR module.

Problem & complexity. The key to TsPAR is to compute transaction schedules for transaction partitioners. Referred to as the *transaction scheduling problem*, this is abstracted as follows:

- INPUT: A transaction workload \mathcal{W} , number k of threads, and a partitioning of \mathcal{W} , i.e., $(P_1, \dots, P_k, \mathcal{R})$.
- OUTPUT: A transaction schedule $S = (f, <)$.
- OBJECTIVE: For the k RC-free queues (Q_1, \dots, Q_k) and residual set \mathcal{R}_s generated by S , to
 - (a) minimize the *makespan* of the k RC-free queues; and
 - (b) minimize the number of unscheduled transactions in \mathcal{R}_s .

Here (a) the makespan of the k RC-free queues is their concurrent execution time, calculated as the maximum of the serial execution time among all k queues. We (b) minimize the amount of work in the residual, and hence maximize work that is handled via RC-free queues. With both (a) and (b), we aim to find a schedule that minimizes the total execution time of all the transactions in workload \mathcal{W} and hence, improve the throughput for executing \mathcal{W} .

As opposed to transaction partitioning problems [15, 36, 40], this bi-criteria optimization problem is more challenging in that it considers runtime conflicts and execution time.

Theorem 1: *The transaction scheduling problem is NP-complete. It is already NP-hard to decide whether there exists a schedule such that*

- (a) *it has k non-empty RC-free queues; or*
- (b) *the makespan of the schedule is no larger than a given number, even when $k = 3$.*

The lower bounds hold even when all transactions take a unit time. □

Proof sketch: We give an NP algorithm that guesses a schedule for \mathcal{W} over k cores, and checks in PTIME whether the schedule satisfies the conditions. We prove the NP-hardness of (a) by reduction from the maximum independent set problem [34], and (b) by reduction from the bounded independent sets problem [29]. Both reductions use transactions that take a unit time only. □

Nonetheless, below we develop an efficient algorithm, denoted by TSgen, to compute transaction schedules for TsPAR. Algorithm TSgen works in two settings. Given a workload \mathcal{W} of transactions, (1) it can take as input a transaction partition plan and turns it into a schedule for \mathcal{W} , as shown in Algorithm 1. (2) It can also compute a schedule for \mathcal{W} starting from scratch. Below we first present TSgen for case (1), and then show how it handles case (2).

From partitions to schedules. Given a workload \mathcal{W} and a partition plan $f_c = (P_1, \dots, P_k, \mathcal{R})$ of \mathcal{W} , TSgen refines f_c into a schedule $(f, <)$ for \mathcal{W} . It preserves transactions of cluster P_i in queue Q_i , makes as many residual transactions in \mathcal{R} to be RC-free as possible by moving them from \mathcal{R} to Q_1, \dots, Q_k , and balances the workloads of the k RC-free queues Q_1, \dots, Q_k .

TSgen (Algorithm 1) starts with empty RC-free queues Q_i and empty set \mathcal{R}_s . It examines residual transactions in \mathcal{R} and decides whether to merge them into the RC-free queues of f and if so, how to do it. In the process, it schedules transactions in P_j 's and preserves their assignment by f_c , i.e., a transaction in P_j is added to queue Q_j at thread j . Along the way, TSgen generates RC-free queue P_j . The unscheduled residual transactions in \mathcal{R} of f_c remain in \mathcal{R}_s .

More specifically, TSgen works as follows. Initially, the set \mathcal{R}_s and RC-free queues Q_i are empty for all $i \in [1, k]$ (lines 1-2). TSgen iteratively expands Q_i by examining transactions in \mathcal{R} one by one, following an ordering \vec{R} of \mathcal{R} (lines 4-14; by default, TSgen picks a random ordering of \mathcal{R} as \vec{R}). For each transaction T_* in \mathcal{R} (line 5), it checks whether T_* can be merged into the input CC-free partition P_l with the smallest total execution time (line 6), in order to balance the workload of RC-free queues.

To do this, it first finds all transactions in P_j ($j \in [1, k], j \neq l$) that are in conflict with T_* , appends them to the corresponding RC-free queue Q_j , and removes them from P_j (lines 7-9). It then checks whether appending T_* to RC-free queue Q_l would cause runtime conflict with transactions that are

already in queue Q_j ($j \neq l$) via procedure `ckRCF` (omitted). If not, it appends T_* to Q_l , and adjusts the load (len_l) of thread l by including time(T_*) (lines 10-11). Otherwise, TSgen decides that T_* is a residual transaction that cannot be scheduled, and moves it into \mathcal{R}_s (line 12).

After all transactions in \mathcal{R} are examined and assigned to either the RC-free queues or \mathcal{R}_s , TSgen appends the remaining transactions in each partition P_i to the corresponding Q_i (lines 13-14). It returns the k RC-free queues (Q_1, \dots, Q_k), and the set \mathcal{R}_s of unscheduled residual transactions (line 15)

Example 4: Given the partition of \mathcal{W}_0 in Example 1, i.e., $P_1 = \{T_1, T_2, T_3\}$, $P_2 = \{T_4\}$ and residual $\mathcal{R} = \{T_5\}$, algorithm TSgen assigns T_5 of \mathcal{R} to P_2 , which turns to RC-free queue $Q_2 = \langle T_4, T_5 \rangle$, with $Q_1 = \langle T_2, T_1, T_3 \rangle$, exactly the same schedule as in Example 3. \square

To complete TSgen, we next present how it identifies transactions that conflict with T_* from each input CC-free partition P_i (lines 7-9). More specifically, to identify transactions in each input partition P_i that are in conflict with T_* , TSgen makes use of the conflict graph G_c of \mathcal{W} . Here G_c is an undirected graph in which (i) the nodes are the transactions of \mathcal{W} , and (ii) two transactions are connected by an edge if they are in conflict with each other. Note that variants of conflict graph have already been used in transaction partitioners (e.g., [15, 61]), and are found efficient to construct and effective in practice; TSgen re-uses their conflict graphs when looking up conflicts. TSgen identifies transactions that are in conflict with T_* by searching the neighbor nodes of T_* in G_c , and checking whether they are in the input CC-free partitions P_i (line 8). It appends such conflicting transactions in P_i to RC-free queue Q_i (line 9), and checks whether T_* incurs runtime conflicts with the expanded Q_i (`ckRCF`; line 10). It appends T_* to Q_l if there is no runtime conflict.

As will be seen in Section 6, TSgen improves the throughput by 131% on average, up to 294%; the overhead it adds to transaction partitioners is less than 5% of that of the partitioners it optimizes.

Scheduling without input partition. Algorithm TSgen is also able to compute a schedule for workload \mathcal{W} in the absence of a partition plan f_c . More specifically, given a transaction workload \mathcal{W} , we simply treat \mathcal{W} as the residual \mathcal{R} and runs TSgen with empty CC-free partitions, i.e., $P_1 = \dots = P_k = \emptyset$. TSgen works in the same way as how it refines a non-empty partitioning of \mathcal{W} .

Complexity. Algorithm TSgen can be implemented in $O(|\mathcal{W}| + (k-1)|\mathcal{R}|)$ -time, where $|\mathcal{W}|$ (resp. $|\mathcal{R}|$) is the number of transactions in \mathcal{W} (residual \mathcal{R}) and k is the number of threads, by reusing the conflict graph G_c from partitioners. That is, TSgen is linear in $|\mathcal{W}|$ when the number k of threads is a constant. Indeed, (a) each transaction in the CC-free partitions P_i is examined only once to decide its assignment. (b) Each transaction in \mathcal{R} is examined by `ckRCF` in $O(k)$ -time.

Remark. TSgen looks up conflicts between transactions from the conflict graph G_c of \mathcal{W} . However, it does not necessarily need to construct G_c during scheduling. Typically transaction partitioners already build G_c [15] or its variants [36] for partitioning in order to minimize cross-partition conflicts, and TSgen re-uses them from the partitioners. TSgen aims to strike a balance between the scheduling cost added to the partitioners and the quality of the schedules computed. It reduces the execution cost of residual queues by preserving conflict-free transactions of partitions P_i 's in the RC-free queues and moving residual transactions of \mathcal{R} to RC-free queues.

5 PROACTIVE TRANSACTION DEFERMENT

In this section, we present how TSDEFER reduces runtime conflicts at runtime for unscheduled residual \mathcal{R}_s from TSpar and for unbundled transactions that are directly assigned to threads.

TSDEFER acts directly on thread-local buffers and reduces runtime conflicts for CC by proactively deferring transactions that are likely to cause runtime conflicts with other transactions that are being executed. Unlike transaction partitioners and TSpar, TSDEFER inflicts extra overhead to the

ALGORITHM 1: Algorithm TSgen**Input:** Transaction workload \mathcal{W} , a partition $f_c = (P_1, \dots, P_k, \mathcal{R})$ of \mathcal{W} over k threads.**Output:** A transaction schedule $(f, <)$ that refines f_c for \mathcal{W} into two k queues (Q_1, \dots, Q_k) and a set \mathcal{R}_s .

```

1  $\mathcal{R}_s \leftarrow \emptyset$ ;
2 foreach  $i$  in  $[1, k]$  do  $\text{len}_i \leftarrow \sum_{T \in P_i} \text{time}(T)$ ;  $Q_i \leftarrow \emptyset$ ;
   /* denote by  $G_c$  the conflict graph of transactions in  $\mathcal{W}$  */
3 while  $\tilde{R} \neq \emptyset$  do
4    $T_* \leftarrow \tilde{R}.\text{pop}()$ ;
5    $l \leftarrow \arg \min_{i \in [1, k]} \text{len}_i$ ; // pick the RC-free queue with the least load
6   foreach  $i$  in  $[1, k]$  do
7      $S_i \leftarrow \{T \in P_i \mid (T, T_*) \text{ is an edge in } G_c\}$ ; //  $S_i$ : transactions in  $P_i$  that conflict with  $T_*$  in  $P_i$ 
8     append transactions in  $S_i$  to  $Q_i$ ;  $P_i \leftarrow P_i \setminus S_i$ ;
9   if  $\text{ckRCF}(Q_1, \dots, Q_{l-1}, Q_l \cup \{T_*\}, \dots, Q_k) = \text{true}$  then
10    append  $T_*$  to  $Q_l$ ;  $\text{len}_l \leftarrow \text{len}_l + \text{time}(T_*)$ ;
11  else  $\mathcal{R}_s \leftarrow \mathcal{R}_s \cup \{T_*\}$ ;
12 foreach  $i \in [1, k]$  do
13   if  $P_i \neq \emptyset$  then append transactions in  $P_i$  to  $Q_i$ ;
14 return  $Q_1, \dots, Q_k$  and  $\mathcal{R}_s$ ;

```

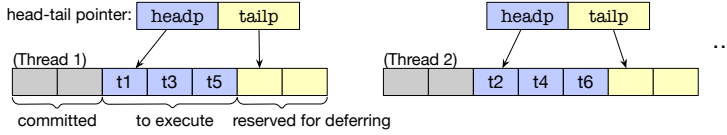


Fig. 3. Structure for transaction progress tracking: grey grids are committed transactions; blue grids represent transactions about to execute; yellow grids reserve slots for deferred transactions.

CC protocols for execution. Hence, it needs to be extremely lightweight while being effective since otherwise its overhead would cancel out any benefit from reduced conflicts.

This imposes challenges to the design of TSDEFER. In particular, tracking and sharing execution progress among threads are essentially contended operations, while locking is not an option due to its large overhead on execution. Moreover, a thread cannot afford to simply look at all concurrent transactions to decide runtime conflicts due to the overhead of reading their read/write sets.

To tackle the challenges, TSDEFER proposes two techniques.

- It uses a *lock-free* structure for all threads to keep track of their execution progress and look up the progress of other threads for conflict detection. This avoids the locking overhead and data races among threads on tracking and sharing progress.
- To suppress the overhead of checking conflicts among transactions on-the-fly, TSDEFER randomly probes data items of the read/write sets of active transaction at other threads within a limited times, and defers the transaction with a certain probability if conflicting items are found. The crux is that each data item probe takes only a constant time in a lock-free manner, independent of both the size of transactions and the number of threads.

Runtime progress tracking. As shown in Fig. 3, TSDEFER represents the transaction queue in each thread-local buffer as an array of transaction IDs; the buffer is shared globally across all threads. Each thread uses two pointers, *headp* and *tailp*, that point to the next transaction to

be executed and the end of the transaction list, respectively. Both the array and the pointers are writable by its own thread and are read-only for other threads.

TsDEFER implements three operations for each thread i :

- `regPos` to update the progress of thread i by moving `headp` to the next transaction ID once the current one commits; and
- `lookup` that, upon each invocation, returns in a *constant* time a data item in the write set of some *active transaction* of another thread, where a transaction is active for thread i if it (a) is at thread j ($j \neq i$) and (b) is pointed by `headp` of thread j .
- `defer` that moves a transaction to the end of the queue that maintains the deferred transactions of the current thread.

Proactive deferment. With these, TsDEFER works as follows.

- (1) When thread i is about to execute a new transaction T in its local buffer \mathcal{T}_i , it checks whether T would cause a runtime conflict with active transactions at other threads by invoking a bounded `#lookups` number of lookup operations.
- (2) Let d be the number of distinct data items retrieved from lookup. If `#lookups` – d is above a threshold (typically 1), TsDEFER decides that T is likely to cause a runtime conflict. In such cases, thread i defers T via the `defer` operation with a probability of `deferp%`.
- (3) If T is deferred, thread i then moves to the next transaction and updates its progress via `regPos`. In addition, it records the deferred T at `tailp`, and then moves `tailp` to the next slot.
- (4) If T is not deferred, thread i then goes on to execute T . If T subsequently commits successfully, then thread i updates its progress via `regPos`. If T aborts, thread i will retry T immediately until it commits. During the period, it does not update its progress.

Lock-free Implementation. We next describe the implementation of the progress tracking structure. The key idea is to keep the overhead low and controllable; the rationale is that TsDEFER does not replace CC, *i.e.*, it does not detect all runtime conflicts during execution time. In light of this, it implements both operations `regPos` and `lookup` with C++ atomic builtins [1] in a lock-free manner. Note that `lookup` may read slightly stale progress due to the lock-free design. However, we find that such staleness has negligible implication as TsDEFER is supposed to be lightweight and does not aim to identify all potential runtime conflicts. For long-run transactions of which conflicts are costly, one can compensate this by instructing `lookup` to check transactions that are further in the future *w.r.t.* the one it sees from `headp`, within bounded steps.

To ensure each lookup has a constant complexity, each thread maintains the predicted access sets of all other transactions. This does not have to be exact since some of the read/write items are not possibly known without execution; for such cases, TsDEFER only records the estimated access sets based on the constants and template of the transaction. For instance, the warehouse id (`w_id`), district id (`d_id`) and customer id (`c_id`) of an instantiated Payment transaction in TPC-C could largely determine its access sets; similarly, the access sets of YCSB transactions can also be accurately inferred by the instantiated parameters (*keys*) in most cases.

Each invocation of `lookup` then randomly takes a thread ID j and index p , via reservoir sampling; it then returns the p -th item in the write set of the transaction that is currently pointed by `headp` of thread j . Note that each invocation requires only one read to the global structure, to retrieve the transaction ID determined by j , say T_* . Then `lookup` simply reads the data item of T_* indexed by p in its local copy of the read/write set of T_* .

Example 5: Recall thread-local buffers $\mathcal{T}_1 = \langle T_1, T_2, T_3 \rangle$ and $\mathcal{T}_2 = \langle T_4, T_5 \rangle$ from Example 2. Assume $\text{deferp}\% = 100\%$. With TsDEFER, thread 1 can defer T_2 with a probability of 50% with 1 lookup (i.e., $\#lookups = 1$), and defer T_2 for certain with only 2 lookup (i.e., $\#lookups = 2$). Indeed, for thread 1 and T_2 , the only active transaction is T_5 , which has two data items, x_1 and x_5 . Hence one lookup has 50% of chance of returning x_1 , which witnesses a conflict with T_2 , while two lookup calls return x_1 for certain, triggering TsDEFER to defer T_2 . Note that reading T_5 alone would already cost 6 reads of data items, already a higher overhead than the deferment of T_2 . \square

Parameters and trade-off. Observe that the accuracy and overhead of TsDEFER is parameterized by the following two parameters:

- $\#lookups$: the number of lookup operations invoked; and
- $\text{deferp}\%$: the probability to deferring a candidate.

The overhead of TsDEFER is determined by the number of lookup operations, i.e., $\#lookups$. With larger $\#lookups$, however, more potentially contended operations in remote threads can be discovered; hence, the chance of *not* deferring a transaction that is in runtime conflict with other active transactions becomes lower. This helps us reduce the CC cost (abort/retry penalty) via reduced conflicts among transactions being executed. Probability $\text{deferp}\%$ allows TsDEFER to adapt to varying contention levels: for extremely high contention workloads, TsDEFER uses a relatively lower $\text{deferp}\%$ to avoid excessive number of transactions being deferred.

With these tunable parameters, TsDEFER can adapt to transaction workloads with varying characteristics, by trading deferring overhead for the accuracy of conflict reduction for CC. When conflict penalties (abort/retry cost) make a larger factor, e.g., for long-run transactions or transactions with heavy application-level abort penalties, one may prefer larger $\#lookups$ for a higher chance of identifying a runtime conflict, to avoid a costly abort/retry with a bit higher overhead. On the contrary, for light and simple transactions, e.g., key-access over a key-value table as YCSB transactions do, one may prefer smaller $\#lookups$ since abort/retry is cheap and overhead is more sensitive. In the extreme case, one can disable TsDEFER with $\#lookups = 0$, avoiding any overhead to CC at all.

Such flexible trade-off between overhead and conflict penalty reduction is in particular viable for transactions that have complex application logic and hence have varying costs of aborts/retries.

6 EXPERIMENTAL STUDY

Using benchmarks, we evaluated the effectiveness of TSKD in improving both partitioning-based systems for bundled transaction workloads (Section 6.2) and CC-based systems for unbundled transactions (Section 6.3), in particular for transactions with varying runtime skewness. We specify the experimental settings in Section 6.1 and summarize the evaluation results in Section 6.4.

6.1 Implementation and Experimental Setup

Systems. We have implemented a prototype of TSKD as described in Sections 3-5, and integrated it with DBx1000 [2, 58] as the transaction execution engine, which implements multiple CC protocols. Since DBx1000 directly initializes thread-local buffers with instantiated unbundled transactions, to evaluate the effectiveness of TSKD for partitioning-based systems, we extend it by also porting external transaction partitioners into DBx1000, making it a full testbed for both partitioning-based transaction processing and non-partitioning CC-based transaction processing methods.

This is done by initializing the transaction buffer for each thread in DBx1000 with a partitioning plan generated by the underlying partitioner, so that transactions are executed according to the partitioning. When TSKD is enabled for transaction partitioners, each thread local buffer receives the transaction queue generated by the partitioner and TSKD.

TSKD instances. We deployed five instances of TSKD, depending on whether it partitions transactions and what partitioner it employs:

- TSKD[S]: an instance of TSKD that targets systems with partitioning; it employs Strife [36], a partitioner that has been shown effective for highly contended transaction workloads [36].
- TSKD[C]: an instance of TSKD that uses the Schism [15] partitioner.
- TSKD[H]: TSKD with Horticulture [35] as the partitioner.
- TSKD[O]: TSKD without input partitions, *i.e.*, taking all transactions as the residual.
- TSKD[CC]: an instance of TSKD targeting unbundled transactions that are directly handled by DBx1000's default transaction-to-thread assignment with CC; it only uses TsDEFER.

Among the three partitioners, STRIFE generates partitioning of transactions with an explicit residual set, while SCHISM and HORTICULTURE do not. For the partitioning of SCHISM and HORTICULTURE, TSKD first extracts a residual set that contains all those transactions that are in conflict with some other transactions from another partition, and then carries out the scheduling as it does with STRIFE.

Baselines. We compared the performance of the scheduling-based TSKD[S], TSKD[C] and TSKD[H] with their partitioning counterparts STRIFE [36], SCHISM [15] and HORTICULTURE [35], respectively, for bundled transactions that are known to the partitioners prior to execution.

Both SCHISM [15] and STRIFE [36] are general partitioners that work with any given transaction workloads with known read/write sets. HORTICULTURE is hard-coded for TPC-C [50] and YCSB [13] workloads, and is not a full-fledged partitioner. We used an unofficial open source implementation of STRIFE from [4] and implemented SCHISM by following [15]. We extracted the implementation of HORTICULTURE from H-Store [3, 23], an OLTP system that implements HORTICULTURE.

To study the effectiveness of TSKD (TsDEFER) for unbundled transactions that are not known beforehand, we also compared TSKD[CC] with DBx1000's default configuration that executes transactions using CC without partitioning (denoted by DbCC).

Configuration. For all TSKD instances, transactions are executed with CC to guarantee correctness. This is a suboptimal implementation of TSKD and is in favor of the baselines since one can retain the lower cost of CC-free execution of the RC-free queues by enforcing the scheduled order via, *e.g.*, dependency tracking in [37, 38]. TsPAR uses the warm-up dry-run trails of DBx1000 as the source of histories to derive coarse cost estimates for transactions. The isolation level is set to serializability.

Benchmarks. We used two transaction benchmarks.

TPC-C [50]. We tested with full TPC-C transactions. Since the built-in TPC-C implementation in DBx1000 contains only the NewOrder and Payment transactions without insertion (update only), we extended it to cover full TPC-C. In particular, we enabled insertions in NewOrder and Payment and added OrderStatus, StockLevel and Delivery transactions to DBx1000 by following [7].

To evaluate the impact of contention levels, we also enabled TPC-C to vary the originally hard-coded percentage ($c\%$) of transactions that cross multiple warehouses. We set $c\%$ to 25% and $\#whn/\#core$ to 2 when varying the number of cores ($\#core$), to simulate workloads with high contention, where $\#whn$ is the number of warehouses.

YCSB [13]. We also tested with the YCSB benchmark. We used the built-in YCSB driver in DBx1000, which implements the YCSB core A workload [57]. We used a single table of 20M records, where each record is 128 bytes in size and is accessed by a unique key; each transaction accesses 16 records. Contention in YCSB is configured by a Zipfian distribution that controls data skewness; we varied the θ parameter of Zipfian from 0.7 to 0.9 ($\theta = 0.8$ by default).

Parameter	Definition	Range	Default
$c\%$	TPC-C cross-warehouse transaction(%)	[15%, 35%]	25%
$\#whn$	# of warehouses in TPC-C	[20, 60]	40
θ	Zipfian skewness parameter in YCSB	[0.7, 0.9]	0.8
$\#core$	# of threads for execution	[8, 32]	20
CC	CC protocols for execution	OCC, SILO, TICTOC	OCC
$\min T$	$\min T \cdot t_T$ = minimum transaction runtime lower bound, where t_T is the average transaction runtime	[1/8, 1]	1/2
p	$p \cdot \min T$ = maximum runtime lower bound	[32, 64]	48
θ_T	the skewness of the Zipfian distribution for runtime lower bounds	[0.7, 0.9]	0.8
l_{IO}	l_{IO} = max latency/min latency, where min latency is 5000 CPU cycles, about 1/6 (1/8) of the average TPC-C (YCSB) transaction runtime.	[0, 100]	50 ¹
θ_{IO}	the skewness parameter of the Zipfian distribution for I/O latency; larger θ_{IO} means <i>longer tail</i> latency	[0.8, 1.6]	1.2
$\#lookups$	# of lookup operations for TsPAR	[1, 5]	2
$\text{defer}p\%$	the probability to defer a candidate	[0.4, 0.8]	0.6

¹ by default we disabled I/O latency; we set $l_{IO} = 50$ for tests about I/O latency when varying θ_{IO} .

Table 1. Workload and system parameters: *blue/yellow* for TPC-C/YCSB parameters; *green* for system parameters; *red* for runtime skewness and I/O latency; and *gray* for TsDEFER parameters. When varying a parameter, we use the default for all the other parameters.

Extension with runtime skewness. Both TPC-C and YCSB transactions are short, *e.g.*, a YCSB transaction consists of key accesses in a key-value table. In practice transactions may have varying lengths (*i.e.*, runtime) and skewness in their distribution. To test more complicated transactions, we extend both TPC-C and YCSB by “lower bounding” the runtime of transactions: assuming that a transaction T is lower bounded by t_{\min} , if the actual execution time of T exceeds t_{\min} , then nothing changes and T commits as usual; otherwise, it delays its committing until the total runtime is t_{\min} .

More specifically, we extend TPC-C and YCSB with three parameters $\min T$, p and θ_T to emulate runtime skewness. Transactions are assigned with a minimum runtime *randomly* drawn from a range $[\min T \cdot t_T, p \cdot \min T \cdot t_T]$, following a Zipfian distribution with the skewness parameter θ_T . Here t_T is the average transaction runtime and $\min T (\leq 1)$ is a small coefficient such that $\min T \cdot t_T$ serves as the “unit” time of transaction execution, and integer p constrains the maximum “lower bound” runtime. These lower bounds follow a Zipfian distribution with skewness degree varied by θ_T .

By varying $\min T$, p and θ_T , we emulate different runtime patterns of the transactions. We set $\min T$ small, *e.g.*, as low as 1/8 so that the original short-run TPC-C and YCSB transactions are included as a subclass of the workloads for all the tests. These parameters improve the expressiveness of TPC-C and YCSB, not to restrict the benchmarks.

Extension with I/O latency. To evaluate the impact of I/O latency, we further extended DBx1000 with a new knob similar to [49], to add an artificial delay to simulate I/O latency at transaction commit time. The delays draw values from $[0, l_{IO} \cdot \min IO]$ for a Zipfian distribution with skewness parameter θ_{IO} , where (a) $\min IO$ is set to 5000 CPU cycles, about 1/6 (resp. 1/8) of the average TPC-C (resp. YCSB) transaction runtime, and (b) l_{IO} varies from $[0, 100]$. By varying l_{IO} and θ_{IO} , we get various patterns of I/O latency to DBx1000. In particular, larger l_{IO} means longer worst-case I/O latency and higher θ_{IO} indicates a “longer-tail” latency distribution.

Metrics. We measure the performance of systems by the following:

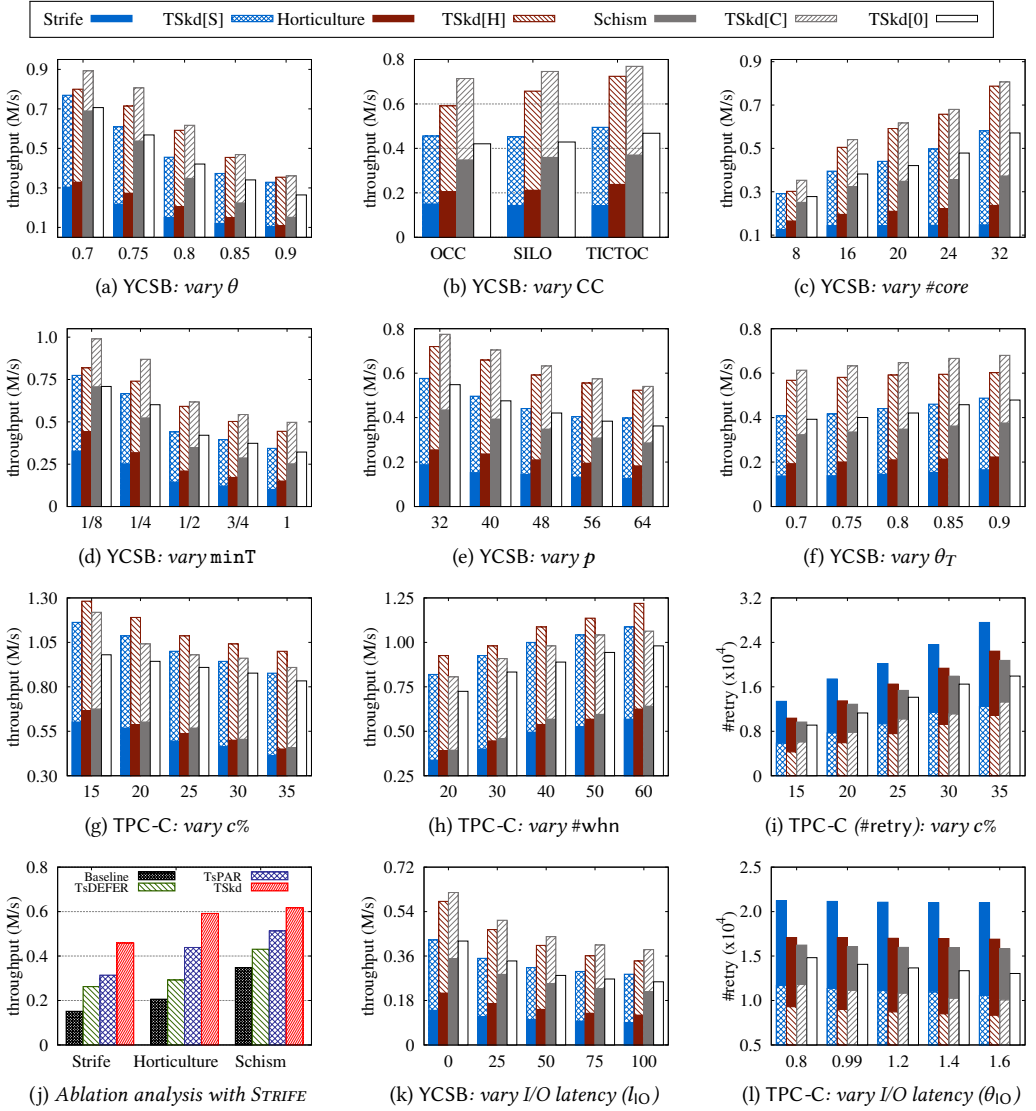


Fig. 4. Effectiveness of TSKD (TsPAR) on Partitioning-based Systems (Section 6.2)

- *throughput*: the number of transactions committed per second;
- *#retry*: the total number of retries per 100,000 transactions.

Experimental setup. The experiments were run on an AWS EC2 m5.8xlarge instance, with 32 vCPU and 128 GB of memory. Each experiment was run 3 times. The average is reported here.

All the parameters, including their configuration ranges and default settings, are summarized in Table 1. When varying a parameter, we use the default configuration of all other parameters. By default, each bundle consists of 10,000 transactions.

6.2 TSKD on Partitioning-based Systems

We first evaluated the effectiveness of transaction scheduling in improving the throughput and #retry of partition-based systems. We compared the performance of TSKD[S], TSKD[C] and TSKD[H] with STRIFE, SCHISM and HORTICULTURE, respectively, to find out improvement over transaction partitioners introduced by TSKD via scheduling. Varying each parameter, we tested the throughput and #retry of all methods on TPC-C, YCSB and their skewed extensions.

Throughput: scheduling vs. partitioning. We first compared the transaction throughput of schedules computed by TSKD[S], TSKD[C] and TSKD[H] with that of partitioning generated by STRIFE, SCHISM and HORTICULTURE, respectively, with varying parameters. The results over YCSB and TPC-C are shown in Figures 4a-4h, and tell us the following.

(1) TSKD consistently improves the throughput of STRIFE, SCHISM and HORTICULTURE. On average TSKD[S], TSKD[C] and TSKD[H] improve the throughput of STRIFE, SCHISM and HORTICULTURE by 211%, 78.9% and 184% over YCSB, respectively, up to 294%, 119% and 225%. Over TPC-C on average the throughput of TSKD[S], TSKD[C] and TSKD[H] is 133%, 75.2% and 101% higher than that of STRIFE, SCHISM and HORTICULTURE, respectively, up to 183%, 104% and 141%. The reason is two-fold.

(a) *Higher concurrency.* By scheduling (TsPAR), TSKD achieves better balanced load among the cores, yielding a higher level of concurrency. In contrast, existing partitioners can be sub-optimal in load balance, e.g., the average ratio of the largest partition over the smallest partition of STRIFE is 3.2 over YCSB, while this reduces to 1.2 after deploying TsPAR atop STRIFE. Moreover, by executing \mathcal{R}_s with CC and proactive transaction deferment (TsDEFER), TSKD further balances the workloads.

(b) *Reduced #retry.* Furthermore, by combining TsPAR and TsDEFER together, such higher level of concurrency does not incur larger CC costs due to the ability of TSKD to reduce runtime conflicts. Indeed, in contrast to the normal pattern that higher concurrency means larger #retry, the #retry of all systems even decreases with TSKD. For instance, the #retry of TSKD[S], TSKD[C] and TSKD[H] is consistently lower than that of STRIFE, SCHISM and HORTICULTURE, in all cases (e.g., Fig. 4i). On average, TSKD reduces 49.7%, 43.6% and 33.6% of the #retry of that of STRIFE, SCHISM and HORTICULTURE on YCSB, respectively, and 54.4%, 36.7% and 53.9% over TPC-C.

(2) The effectiveness of TSKD is robust w.r.t. CC protocols; it is even more evident with added cores as shown in Figures 4b-4c. For instance, over YCSB, TSKD[S] consistently achieves over 203% higher throughput than STRIFE with either OCC, SILO or TICTOC. The gap even increases with larger #core, e.g., TSKD[H] improves STRIFE by 133% with 8 cores, and 294% with 32 cores.

(3) Overall, TSKD is more effective for workloads with higher contention, e.g., the throughput improvement of TSKD[H] over HORTICULTURE increases from 143% to 225% when θ varies from 0.7 to 0.9 over YCSB, as shown in Fig. 4a; similarly, when $c\%$ increases from 15% to 35% for TPC-C, the throughput improvement of TSKD[C] over SCHISM increases from 80.5% to 98.2% as depicted in Fig. 4g.

(4) Without an input transaction partitioning, TSKD[0] still achieves throughput on average 85.8%, 65.1% and 69.7% higher than STRIFE, SCHISM and HORTICULTURE over TPC-C, respectively, and 184%, 29.1% and 101% higher over YCSB, as shown in Fig. 4. Compared to transaction partitioners, TSKD[0] achieves better load balancing via scheduling, by treating all transactions as residual; in addition, it has low retries due to reduced conflicts. This suggests that TSKD[0] is an option for non-partitioned workloads or transaction partitioners that do not produce a residual.

(5) To analyze the individual contribution of TsPAR and TsDEFER to the effectiveness of TSKD, we tested TSKD[S] against TsPAR[S] and TsDEFER[S] where TsPAR[S] is TSKD on STRIFE with TsDEFER disabled, and TsDEFER[S] is TSKD with only TsDEFER enabled to execute partitions of STRIFE; simi-

	system	scheduled pct. $s\%$	#retry/10k transactions	
			w/o TsDEFER	w/ TsDEFER
YCSB	TSKD[S]	64.5	1353	448
	TSKD[C]	52.7	2786	1029
	TSKD[H]	69.7	1614	611
TPC-C	TSKD[S]	34.6	542	255
	TSKD[C]	36.7	803	474
	TSKD[H]	20.8	416	253

Table 2. Accuracy of scheduling and effectiveness of TsDEFER

larly for SCHISM and HORTICULTURE. As shown in Fig. 4j, over YCSB, on average TSKD[S], TsPAR[S] and TsDEFER[S] improve the throughput of STRIFE by 203%, 108%, 73.2%, respectively; similarly for SCHISM and HORTICULTURE and over YCSB. This shows that for bundled workloads, TsPAR plays a bigger role in improving the performance. Moreover, TsPAR and TsDEFER perform the best when working together, e.g., even better than the sum of the improvement by the two separately on STRIFE.

Runtime skewness. TSKD improves the robustness of all systems for transactions with varying skewness, as shown in Figures 4d-4f over YCSB (results over TPC-C omitted due to space limit). In particular, the improvement in throughput by TSKD[S], TSKD[C] and TSKD[H] over STRIFE, SCHISM and HORTICULTURE, respectively, even increases when transactions become longer (*i.e.*, larger $\min T$) or more variable (*i.e.*, larger p and smaller θ_T). For example, on YCSB, TSKD[S] improves STRIFE in throughput by 136% with $\min T = 1/8$ while this increases to 240% with $\min T = 1$; similarly TSKD[H] improves HORTICULTURE on YCSB by 171% with p set to 32, while this increases to 186% with $p = 64$. This is because the quality of partitions degrades with transactions of varying lengths, while with TSKD they can adapt to the variance and skewness in transaction runtime. Moreover, longer transactions inflict larger conflict penalties, and TsPAR is more effective in reducing their runtime conflicts.

I/O latency. We next evaluated the impact of I/O latency. Varying l_{IO} and θ_{IO} , we tested the throughput of all methods; partial results over YCSB and TPC-C are shown in Figures 4k-4l. We find that TSKD still improves all three partitioners in all cases. Although the raw throughput degrades for all methods with larger l_{IO} or smaller θ_{IO} (*i.e.*, not so long-tail), the improvement brought by TSKD is relatively stable, e.g., consistently around 205% for STRIFE over YCSB; similarly for other partitioners and over TPC-C. This is because while latency reduces the accuracy of scheduling by TsPAR, it increases the cost of retries; and TSKD consistently reduces the retries of all partitioners (e.g., Fig. 4l) and balances the partitions in the presence of long-tail latency. Hence, the effectiveness of TSKD remains robust against I/O latency. This suggests that TSKD can also work with transactions involving I/O or network latency.

Overhead. We also tested the overhead that TsPAR adds to the transaction partitioners, measured as the ratio of the runtime of TsGen to the partitioning time of the partitioners, denoted by overheadR. For workloads consisting of 100,000 transactions, we find that the average overheadR of TsPAR over STRIFE and SCHISM over YCSB is 4.1% and 3.7%, respectively, and 4.6% and 4.4% over TPC-C. This verifies that TsPAR has moderate overhead.

Cost estimation. By default, TsPAR uses the warm-up process of DBx1000 to assess the cost of transactions when scheduling (recall Section 3). To assess the accuracy of scheduling with such coarse estimates, we measured (a) the average scheduled percentage ($s\%$) of residual transactions that were merged to RC-free queues; and (b) the #retry when executing the RC-free queues, with and without employing TsDEFER. The results are shown in Table 2.

We find that TSKD schedules a decent percentage of residual transactions, e.g., 30.7% and 62.3% over TPC-C and YCSB, respectively. Due to the coarse estimates used by TSKD, it is understandable

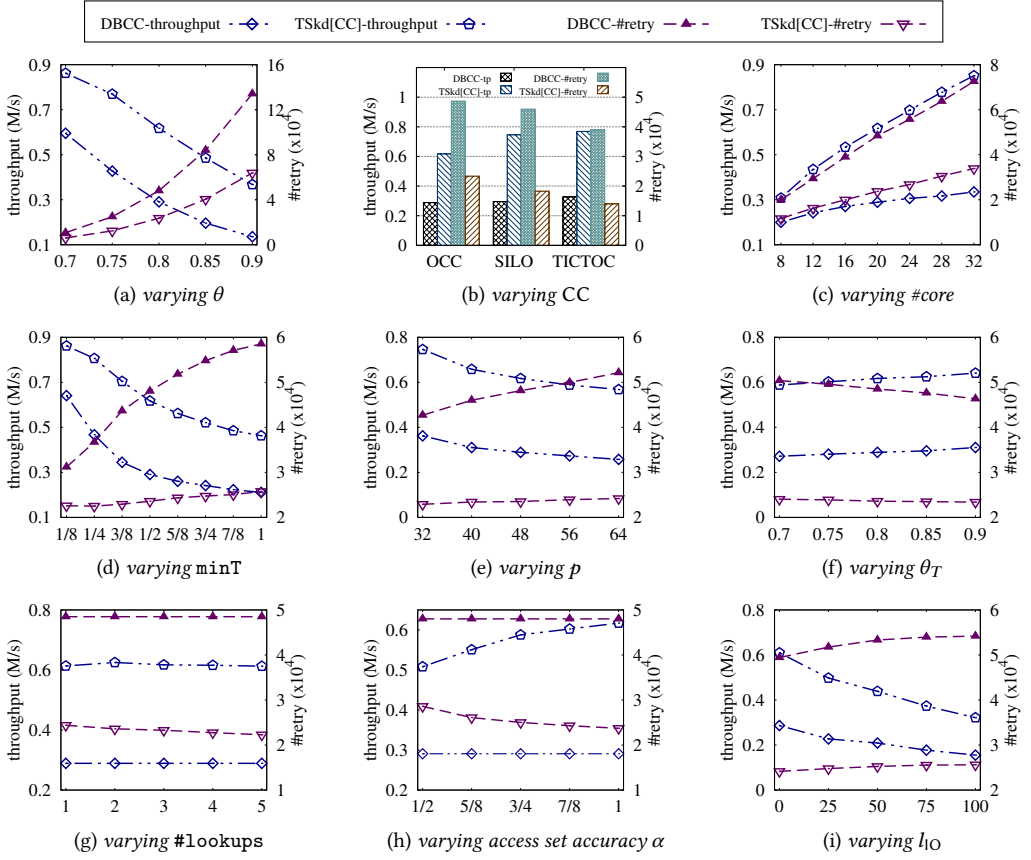


Fig. 5. Effectiveness of TSKD (TsDEFER) on CC-based Systems over YCSB (Section 6.3): X-axis for throughput and Y-axis for #retry.

that transactions in the RC-free queues incur conflicts and retries. Nonetheless, after employing TsDEFER, TSKD significantly reduces #retry for the RC-free transactions. For instance, TSKD[S], TSKD[C] and TSKD[H] reduce 66.9%, 63.1% and 62.2% of #retry for the RC-free transactions over YCSB, and by 53.0%, 41.0% and 39.2% over TPC-C, respectively. As shown earlier, this even makes the #retry of all systems much lower than without scheduling, although they have much higher level of concurrency and better load balancing with TSKD deployed.

Although TPC-C transactions are less likely to be scheduled than YCSB, the throughput improvement of TSKD over TPC-C is comparable to that of YCSB. This is because with even a moderate number of scheduled residual transactions, TSKD is already able to make TPC-C partitions balanced.

6.3 TSKD on Non-partitioning CC-based Systems

We also evaluated the effectiveness of TSKD (TsDEFER) for transaction systems that directly use CC without transaction partitioning. To do this, we compared the performance of TSKD[CC] with that of DBCC, with varying parameters and CC protocols. The results over YCSB are shown in Figure 5 (the results over TPC-C are similar and hence are omitted for saving space).

Contention. Varying θ from 0.7 to 0.9 for YCSB and $c\%$ from 15% to 35%, we tested the impact

of transaction contention level on the effectiveness of TsDEFER. We find that TsDEFER consistently improves DbCC in both throughput and #retry. As shown in Fig. 5a, on YCSB, on average the throughput of TSKD[CC] is 111% higher than that of DbCC, while the #retry of TSKD[CC] is 49.8% lower. By digging into the profiling statistics, we find that the reduction of retries by TsDEFER is largely consistent with that of mutex contention of the execution. In fact, TsDEFER reduces #contended_mutex [31], the total number of times that a mutex was contended (already locked when a lock request was made) in DbCC by 53.8% on average.

Moreover, the improvement is even more significant with larger θ . For instance, TsDEFER improves the throughput of DbCC by 44.8% and reduces the #retry by 43.9% when θ is 0.7, while these increase to 171% and 53.4%, respectively, with $\theta = 0.9$. The results over TPC-C are similar. This is because for workloads with higher contention, there exist more runtime conflicts, which could benefit from TsDEFER with a higher chance. This is also confirmed by the runtime statistics, *e.g.*, TSKD[CC] reduces #contended_mutex of DbCC by 44.5% with $\theta = 0.7$ and this increases to 56.4% with $\theta = 0.9$.

Runtime skewness. Varying parameters $\min T$, p and θ_T , we evaluated the effectiveness of TSKD in response to different runtime skewness. As shown in Figures 5d-5f, TsDEFER is particularly effective with longer (larger $\min T$) and more skewed and variable (larger p and lower θ_T) transactions. For instance, over YCSB when $\min T$ increases from 1/8 to 1, the throughput improvement of TSKD[CC] over DbCC increases from 34.5% to 119%. In contrast, when θ_T is 0.9, TSKD[CC] reduces #retry of DbCC by 49.5%, and this increases to 53.3% when $\theta_T = 0.7$; similarly over TPC-C.

I/O latency. Varying l_{IO} and θ_{IO} , we also evaluated the impact of I/O latency on the effectiveness of TsDEFER for unbundled transactions. We find that with larger l_{IO} or smaller θ_{IO} , the throughput of both DbCC and TSKD[CC] decreases as transactions are prolonged by the I/O stalls; however, TSKD[CC] remains robust in improving DbCC in both throughput and retries, as shown in Fig. 5i. This is because TsDEFER is insensitive to transaction lengths via runtime progress tracking; hence its effectiveness is robust against varying patterns of I/O latency.

Trade-offs. We next examined the trade-offs between the effectiveness and overhead of TsDEFER, by varying #lookups and deferp%. We find that larger #lookups naturally improves the effectiveness of TsDEFER in reducing runtime conflicts and retries, *e.g.*, over YCSB the #retry of DbCC is reduced by 49.7% with TSKD enabled when #lookups = 1, and this increases to 54.1% with #lookups = 5 (as shown in Fig. 5g). However, larger #lookups comes with higher overhead, which would work with longer transactions whose runtime and retry cost could cover the overhead of TsDEFER. For short-run TPC-C and YCSB transactions, we find that #lookups = 2 gives the best throughput, *e.g.*, TSKD[CC] improves TSKD by 116% over YCSB and 105% over TPC-C. Similarly, we find that higher deferp% gives better reduction on #retry for TSKD over DbCC.

CC and scalability. Varying CC protocols and the number of cores, we tested the impact of CC and the scalability of both TSKD[CC] and DbCC. The results over YCSB are depicted in Figures 5b-5c. We find that TSKD[CC] consistently improves the throughput and #retry of DbCC via TsDEFER, by 117% and 51.8%, respectively, on average. Moreover, the gap even widens with larger #core. TSKD works the best with TICTOC, yielding a throughput improvement of 152% and a reduction of #retry by 63.9% over YCSB. It has the least advantage with OCC (the default CC used in all tests), but still improves DbCC by 116% for throughput and 52.0% for #retry.

Impact of inaccurate access sets. Finally, we evaluated the impact of inaccurately determined transaction read/write sets on the effectiveness of TSKD[CC]. To do this, we restricted TSKD[CC] such that it could only use an α -fraction of the actual access sets, *i.e.*, the determined transaction read/write sets have an accuracy of α . Varying α from [0.5, 1], we tested the performance of

TSKD[CC] with DbCC. The results over YCSB are shown in Fig. 5h (the results over TPC-C are similar and thus omitted). We find that TSKD[CC] still improves the throughput of DbCC even when TSKD[CC] overlooked 50% of the actual access set, *i.e.*, when α is as low as 50% over YCSB. This is because TsDEFER only randomly probes access sets of concurrent transactions to detect potential conflicts and it needs only part of the access sets; hence, it still observes conflicts and improves DbCC even when it overlooked reads/writes moderately. Naturally, the effectiveness of TSKD[CC] improves when the access sets are more accurately determined with higher α .

6.4 Summary

From the experiments we find the following.

- (1) Transaction scheduling (TsPAR) is effective in improving the performance of partitioning-based transaction systems, *e.g.*, it improves the throughput of state-of-the-art transaction partitioners by 131% on average, up to 294% over TPC-C and YCSB benchmarks.
- (2) The benefit of TSKD (TsPAR) is even more evident when a larger number of cores are available. For instance, over YCSB, on average TSKD improves the throughput of partitioners by 75.3% when 8 cores are used; the improvement increases to 214% with 32 cores.
- (3) Proactive transaction deferment is effective for CC-based systems. On average it improves the throughput of DbCC by 109% and reduces its #retry by 45.7%, up to 152% and 63.9%, respectively.
- (4) With TsDEFER, TsPAR can handle inaccurate cost estimates of transactions while still improving the throughput of transaction partitioners.
- (5) TSKD is particularly effective for transactions with skewed or long runtime or I/O latency. It makes both partition-based and CC-based systems robust against runtime and I/O skewness.
- (6) TSKD works well with all CC protocols. Its improvement over transaction partitioners and DbCC (CC-based) is evident with each and every CC scheme tested.

7 RELATED WORK

We categorize the related work as follows.

Transaction partitioning. There has been a host of work on transaction partitioning for parallel OLTP systems [15, 35, 39, 56, 61]. The methods are designed for distributed systems under a shared-nothing architecture, where the database is partitioned across multiple computing nodes to minimize transactions that access data from multiple partitions since it requires costly distributed CC.

Closer to this work is the study of multi-core transaction processing [17, 33, 36, 40, 47]. Instead of assigning transactions to cores, [33] decomposes transactions into smaller read and write actions, and assigns such actions to threads to minimize contention. Adaptive concurrency control is proposed by [47] for changing workloads, by dynamically clustering data and adopting the optimistic CC (OCC) protocols for each cluster. [17] uses transaction batching and reordering throughout the life cycle of a transaction to reduce conflicts for OCC-based protocols. [40] develops Orthrus, a multi-core transaction system that separates transaction logic and concurrency control, and uses a set of dedicated cores for the latter; the deterministic approach of [41] pre-processes transaction workloads via partitioning. In particular, [36] proposes Strife, a transaction partitioner to dynamically cluster finer-grained batches of contended workloads that do not have a good static partition.

Different from the partitioning strategies, we propose to schedule transactions of possibly varying execution times; transaction scheduling not only partitions transactions but also sorts them in each partition. This is based on runtime conflicts, which characterize contention among transactions at a finer-grained granularity than traditional transaction conflicts, and allow conventionally conflicting

transactions to be executed without conflicts. We also develop an algorithm that can refine existing transaction partitioning and make it a sound schedule while minimizing runtime conflicts. Note that our concept of transaction scheduling is quite different from the transaction reordering of [17], in which batching and reordering take place among concurrent transactions in the validation stage of OCC-based protocols, to improve commit success rate. By contrast, our method is not restricted to OCC protocols; instead, we study partitioning-based transaction processing and aim to reduce runtime conflicts by scheduling. In addition, we develop a strategy that proactively detects and defers transactions that are highly likely to cause runtime conflicts on-the-fly, to reduce aborts and retries.

Transaction assignment for OLTP. Most OLTP systems use random or round-robin like strategies to assign unbundled transactions that are coming and processed on-the-fly [62]. Unlike partitioners, they do not have prior knowledge of the entire workload, and hence assign transactions one by one upon arrival. [43] proposes to use a lightweight ML model that predicts, for an incoming transaction, which thread it should be assigned to in order to have lower abort chance. Our work complements these methods by altering the assignment during execution, by tracking thread progress and proactively deferring transactions that are likely to cause conflicts.

Concurrency control. Concurrency control (CC) provides isolation guarantees for concurrent transaction processing. A variety of CC protocols have been proposed, which mostly fall in two classes: locking-based and timestamp based. Locking-based ones, such as two-phase locking [11, 18] and its variants, are pessimistic in that a transaction accesses a tuple only after it acquires a lock with the required permission. For high-contention transactions, timestamp-based protocols such as OCC [6, 26], multi-version concurrency control (MVCC) [10] and their variants [9, 16, 19, 28, 32, 42, 46, 54, 55, 59, 60] have been verified effective in reducing conflicts and blocking time. Hybrid approaches that combine the two strategies have also been explored [45, 47]. There has also been work on learned CC to further mix CC parameters and specialize to a given workload [52].

While these CC protocols provide controlled concurrent execution of conflicting transactions and ensure isolation guarantees, they do not specify how the transactions should be allocated to the threads and in what order they should be executed; both of these have significant impact on the performance of concurrent transaction execution. Our work aims to bridge the gap by minimizing runtime conflicts via transaction scheduling and proactive transaction deferment, as an alternative to conventional approaches that only cluster transactions or directly invoke CC protocols.

In-memory OLTP systems. A number of multi-core in-memory OLTP systems have been developed recently, e.g., TicTok [59], Cicada [28], Foedus [25], ERMIA [24], Silo [51] and Orthrus [40]. These systems focus on the design of CC protocols and architectural optimizations for OLTP workloads.

Different from these systems, we do not aim to build yet another full-fledged OLTP database system. Instead, we present TSKD as a lightweight tool that can be incorporated into these OLTP database systems to improve their throughput by reducing runtime conflict.

Deterministic approaches. Related is also the work on deterministic databases (see [5] for a survey), which execute bundled transaction workloads via a predetermined order that is typically decided via dependency analysis of the transactions [14, 20, 27, 37, 38, 40, 49]. They scale well in distributed environments, where the determined ordering ensures transactions to always read consistent data across multiple cache copies. In particular, QueCC [37] advocates a completely CC-free execution of contended transactions. It plans fragmented transaction pieces in multiple queues, which are aligned via their dependency graph that ensures serializability. Caracal [38] proposes a deterministic scheme that batches concurrency control operations with access to only write sets, and splits transactions into uncontended and contended pieces for higher concurrency.

Different from these, our work studies how to reduce runtime conflicts by proposing transaction

scheduling and proactive deferring. It neither imposes restrictions on how CC should work, *e.g.*, deterministically or nondeterministically, nor breaks down transactions. This said, this work and the previous work on deterministic databases complement each other. For instance, their techniques for determining read/write sets and analyzing workloads can be used by transaction partitioners for clustering and assigning transactions. Their dependency-based execution method can be adopted by transaction scheduling when runtime estimates are not reliable.

8 CONCLUSION

We have shown that by considering runtime conflicts, more transactions can be executed concurrently with low conflict penalties than adopting the conventional notion of conflicts. We have proposed transaction scheduling by placing an ordering on the execution of transactions, instead of transaction partitioning. We have shown that the scheduling problem is NP-complete. This said, we have developed an efficient algorithm to refine a partitioning into a schedule for bundled transaction. Moreover, for unbundled transactions targeted by CC protocols, we have proposed a proactive deferring strategy to reduce conflict penalties. Our experimental study has verified that transaction scheduling and deferment are promising in improving throughput and reducing conflict penalties.

One topic for future work is to develop ML models that decide TSKD parameters specialized for given workloads. Another topic is to combine transaction scheduling with other strategies, *e.g.*, transaction decomposition [33] and dynamic adjustment [47].

ACKNOWLEDGMENTS

Cao is supported by Royal Academy of Engineering Research Fellowship RF\201920\19\319 and Edinburgh-Huawei Joint Lab. Fan is supported in part by Royal Society Wolfson Research Merit Award WRM/R1/180014.

REFERENCES

- [1] [n. d.]. Atomic Builtins - Using the GNU Compiler Collection (GCC). https://gcc.gnu.org/onlinedocs/gcc/_005f_005fatomic-Builtins.html.
- [2] [n. d.]. DBx1000. <https://github.com/yxymit/DBx1000>.
- [3] [n. d.]. H-Store. <https://github.com/apavlo/h-store>.
- [4] [n. d.]. Strife. <https://github.com/zhr1201/STRIFE-public>.
- [5] Daniel J. Abadi and Jose M. Faleiro. 2018. An overview of deterministic database systems. *Commun. ACM* 61, 9 (2018), 78–88.
- [6] Rakesh Agrawal, Michael J. Carey, and Miron Livny. 1987. Concurrency Control Performance Modeling: Alternatives and Implications. *TODS* 12, 4 (1987), 609–654.
- [7] Tiemo Bang, Norman May, Ilia Petrov, and Carsten Binnig. 2022. The Full Story of 1000 Cores: An Examination of Concurrency Control on Real (ly) Large Multi-Socket Hardware—Measurements, Logs, Plots. *The VLDB Journal* (2022).
- [8] Hal Berenson, Philip A. Bernstein, Jim Gray, Jim Melton, Elizabeth J. O’Neil, and Patrick E. O’Neil. 1995. A Critique of ANSI SQL Isolation Levels. In *SIGMOD*. 1–10.
- [9] Philip A. Bernstein and Nathan Goodman. 1981. Concurrency Control in Distributed Database Systems. *ACM Comput. Surv.* 13, 2 (1981), 185–221.
- [10] Philip A. Bernstein and Nathan Goodman. 1983. Multiversion Concurrency Control - Theory and Algorithms. *TODS* 8, 4 (1983), 465–483.
- [11] Philip A. Bernstein, David W. Shipman, and Wing S. Wong. 1979. Formal Aspects of Serializability in Database Concurrency Control. *IEEE Trans. Software Eng.* 5, 3 (1979), 203–216.
- [12] Aleksey V. Burdakov, Victoria Proletarskaya, Andrey D. Ploutenko, Oleg Ermakov, and Uriy A. Grigorev. 2020. Predicting SQL Query Execution Time with a Cost Model for Spark Platform. In *IoTBDs*. 279–287.
- [13] Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. 2010. Benchmarking cloud serving systems with YCSB. In *SoCC*. 143–154.
- [14] James Cowling and Barbara Liskov. 2012. Granola: Low-Overhead Distributed Transaction Coordination. In *USENIX ATC*. 223–235.

- [15] Carlo Curino, Evan Jones, Yang Zhang, and Sam Madden. 2010. Schism: A Workload-Driven Approach to Database Replication and Partitioning. *PVLDB* 3, 1–2 (2010), 48–57.
- [16] Mohammad Dashti, Sachin Basil John, Amir Shaikhha, and Christoph Koch. 2016. Repairing Conflicts among MVCC Transactions. *CoRR* abs/1603.00542 (2016).
- [17] Bailu Ding, Lucja Kot, and Johannes Gehrke. 2018. Improving Optimistic Concurrency Control Through Transaction Batching and Operation Reordering. *PVLDB* 12, 2 (2018), 169–182.
- [18] Kapali P. Eswaran, Jim Gray, Raymond A. Lorie, and Irving L. Traiger. 1976. The Notions of Consistency and Predicate Locks in a Database System. *Commun. ACM* 19, 11 (1976), 624–633.
- [19] Jose M. Faleiro, Daniel Abadi, and Joseph M. Hellerstein. 2017. High Performance Transactions via Early Write Visibility. *PVLDB* 10, 5 (2017), 613–624.
- [20] Jose M. Faleiro and Daniel J. Abadi. 2015. Rethinking serializable multiversion concurrency control. *PVLDB* 8, 11 (2015), 1190–1201.
- [21] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. 2002. *Database systems - the complete book (international edition)*. Pearson Education.
- [22] Evan P. C. Jones, Daniel J. Abadi, and Samuel Madden. 2010. Low overhead concurrency control for partitioned main memory databases. In *SIGMOD*. 603–614.
- [23] Robert Kallman, Hideaki Kimura, Jonathan Natkins, Andrew Pavlo, Alex Rasin, Stanley B. Zdonik, Evan P. C. Jones, Samuel Madden, Michael Stonebraker, Yang Zhang, John Hugg, and Daniel J. Abadi. 2008. H-store: a high-performance, distributed main memory transaction processing system. *Proc. VLDB Endow* 1, 2 (2008), 1496–1499.
- [24] Kangnyeon Kim, Tianzheng Wang, Ryan Johnson, and Ippokratis Pandis. 2016. ERMIA: Fast Memory-Optimized Database System for Heterogeneous Workloads. In *SIGMOD*. 1675–1687.
- [25] Hideaki Kimura. 2015. FOEDUS: OLTP Engine for a Thousand Cores and NVRAM. In *SIGMOD*. 691–706.
- [26] H. T. Kung and John T. Robinson. 1981. On Optimistic Methods for Concurrency Control. *TODS* 6, 2 (1981), 213–226.
- [27] Jialin Li, Ellis Michael, and Dan RK Ports. 2017. Eris: Coordination-free consistent transactions using in-network concurrency control. In *SOSP*. 104–120.
- [28] Hyeontaek Lim, Michael Kaminsky, and David G. Andersen. 2017. Cicada: Dependably Fast Multi-Core In-Memory Transactions. In *SIGMOD*. 21–35.
- [29] Zbigniew Lonc. 1991. On complexity of some chain and antichain partition problems. In *International Workshop on Graph-Theoretic Concepts in Computer Science*. Springer, 97–104.
- [30] Microsoft. [n. d.]. In-Memory OLTP overview and usage scenarios. <https://docs.microsoft.com/en-us/sql/relational-databases/in-memory-oltp/overview-and-usage-scenarios?view=sql-server-ver16#usage-scenarios-for->.
- [31] Mutrace. [n. d.]. Measuring Lock Contention. <http://0pointer.de/blog/projects/mutrace.html>.
- [32] Thomas Neumann, Tobias Mühlbauer, and Alfons Kemper. 2015. Fast Serializable Multi-Version Concurrency Control for Main-Memory Database Systems. In *SIGMOD*. 677–689.
- [33] Ippokratis Pandis, Ryan Johnson, Nikos Hardavellas, and Anastasia Ailamaki. 2010. Data-Oriented Transaction Execution. *PVLDB* 3, 1 (2010), 928–939.
- [34] Christos H Papadimitriou. 1994. *Computational Complexity*. Addison-Wesley.
- [35] Andrew Pavlo, Carlo Curino, and Stanley B. Zdonik. 2012. Skew-aware automatic database partitioning in shared-nothing, parallel OLTP systems. In *SIGMOD*. 61–72.
- [36] Guna Prasaad, Alvin Cheung, and Dan Suciu. 2020. Handling Highly Contended OLTP Workloads Using Fast Dynamic Partitioning. In *SIGMOD*. 527–542.
- [37] Thamir M. Qadah and Mohammad Sadoghi. 2018. QueCC: A Queue-oriented, Control-free Concurrency Architecture. In *Middleware*. ACM, 13–25.
- [38] Dai Qin, Angela Demke Brown, and Ashvin Goel. 2021. Caracal: Contention Management with Deterministic Concurrency Control. In *SIGOPS*. 180–194.
- [39] Abdul Quamar, K. Ashwin Kumar, and Amol Deshpande. 2013. SWORD: Scalable Workload-Aware Data Placement for Transactional Workloads. In *EDBT*. 430–441.
- [40] Kun Ren, Jose M. Faleiro, and Daniel J. Abadi. 2016. Design Principles for Scaling Multi-core OLTP Under High Contention. In *SIGMOD*. 1583–1598.
- [41] Kun Ren, Alexander Thomson, and Daniel J. Abadi. 2014. An Evaluation of the Advantages and Disadvantages of Deterministic Database Systems. *PVLDB* 7, 10 (2014), 821–832.
- [42] Mohammad Sadoghi, Mustafa Canim, Bishwaranjan Bhattacharjee, Fabian Nagel, and Kenneth A. Ross. 2014. Reducing Database Locking Contention Through Multi-version Concurrency. *PVLDB* 7, 13 (2014), 1331–1342.
- [43] Yangjun Sheng, Anthony Tomic, Tieying Zhang, and Andrew Pavlo. 2019. Scheduling OLTP transactions via learned abort prediction. In *aiDM*. ACM, 1:1–1:8.

- [44] Rekha Singhal and Manoj K. Nambiar. 2016. Predicting SQL Query Execution Time for Large Data Volume. In *IDEAS*. 378–385.
- [45] Chunzhi Su, Natacha Crooks, Cong Ding, Lorenzo Alvisi, and Chao Xie. 2017. Bringing Modular Concurrency Control to the Next Level. In *SIGMOD*. 283–297.
- [46] Dixin Tang and Aaron J. Elmore. 2018. Toward Coordination-free and Reconfigurable Mixed Concurrency Control. In *USENIX ATC 18*. USENIX Association, 809–822.
- [47] Dixin Tang, Hao Jiang, and Aaron J. Elmore. 2017. Adaptive Concurrency Control: Despite the Looking Glass, One Concurrency Control Does Not Fit All. In *CIDR*.
- [48] Enrico Tedeschi, Tor-Arne S. Nordmo, Dag Johansen, and Hrule(a)vard D. Johansen. 2019. Predicting Transaction Latency with Deep Learning in Proof-of-Work Blockchains. In *BigData*. IEEE, 4223–4231.
- [49] Alexander Thomson, Thaddeus Diamond, Shu-Chun Weng, Kun Ren, Philip Shao, and Daniel J Abadi. 2012. Calvin: Fast distributed transactions for partitioned database systems. In *SIGMOD*. 1–12.
- [50] Transaction Processing Performance Council. [n. d.]. TPC-C Benchmark. <http://www.tpc.org/tpcc/>.
- [51] Stephen Tu, Wenting Zheng, Eddie Kohler, Barbara Liskov, and Samuel Madden. 2013. Speedy transactions in multicore in-memory databases. In *SIGOPS*. 18–32.
- [52] Jia-Chen Wang, Ding Ding, Huan Wang, Conrad Christensen, Zhaoguo Wang, Haibo Chen, and Jinyang Li. 2021. Polyjuice: High-Performance Transactions via Learned Concurrency Control. In *OSDI*. 198–216.
- [53] Gerhard Weikum, Arnd Christian Koenig, Achim Kraiss, and Markus Sinnwell. 1999. Towards self-tuning memory management for data servers. *IEEE Data Eng. Bull.* 22, 2 (1999), 3–11.
- [54] Yingjun Wu, Joy Arulraj, Jiexi Lin, Ran Xian, and Andrew Pavlo. 2017. An Empirical Evaluation of In-Memory Multi-Version Concurrency Control. *PVLDB* 10, 7 (2017), 781–792.
- [55] Yingjun Wu, Chee Yong Chan, and Kian-Lee Tan. 2016. Transaction Healing: Scaling Optimistic Concurrency Control on Multicores. In *SIGMOD*. ACM, 1689–1704.
- [56] Chao Xie, Chunzhi Su, Cody Little, Lorenzo Alvisi, Manos Kapritsos, and Yang Wang. 2015. High-performance ACID via modular concurrency control. In *SOSP*. 279–294.
- [57] Yahoo! [n. d.]. YCSB Core Workloads. <https://github.com/brianfrankcooper/YCSB/wiki/Core-Workloads>.
- [58] Xiangyao Yu, George Bezerra, Andrew Pavlo, Srinivas Devadas, and Michael Stonebraker. 2014. Staring into the Abyss: An Evaluation of Concurrency Control with One Thousand Cores. *PVLDB* 8, 3 (2014), 209–220.
- [59] Xiangyao Yu, Andrew Pavlo, Daniel Sánchez, and Srinivas Devadas. 2016. TicToc: Time Traveling Optimistic Concurrency Control. In *SIGMOD*. 1629–1642.
- [60] Yuan Yuan, Kaibo Wang, Rubao Lee, Xiaoning Ding, Jing Xing, Spyros Blanas, and Xiaodong Zhang. 2016. BCC: Reducing False Aborts in Optimistic Concurrency Control with Low Cost for In-Memory Databases. *PVLDB* 9, 6 (2016), 504–515.
- [61] Erfan Zamanian, Julian Shun, Carsten Binnig, and Tim Kraska. 2020. Chiller: Contention-Centric Transaction Execution and Data Partitioning for Modern Networks. In *SIGMOD*. 511–526.
- [62] Tieying Zhang, Anthony Tomasic, Yangjun Sheng, and Andrew Pavlo. 2018. Performance of OLTP via Intelligent Scheduling. In *ICDE*. 1288–1291.

Received April 2022; revised July 2022; accepted August 2022