

# Region-Linearizable Serializability: A Practical Consistency Model for Multi-Region Distributed Transactions

Ruijie Gong<sup>†</sup>  
University of Hong Kong  
Hong Kong SRA, China  
rjgong@cs.hku.hk

Haoze Song<sup>†</sup>  
University of Hong Kong  
Hong Kong SRA, China  
hzsong@cs.hku.hk

Xusheng Chen<sup>\*</sup>  
Huawei Cloud  
ShenZhen, China  
chenxusheng6@huawei.com

Sen Wang  
Huawei Technologies Co., Ltd.  
Hong Kong SRA, China  
wangsen31@huawei.com

Heming Cui  
University of Hong Kong  
Hong Kong SRA, China  
heming@cs.hku.hk

## ABSTRACT

Deploying databases across multiple regions has become the de facto choice for cloud-native applications that desire high availability, strong scalability, and efficient service localization. However, supporting serializable transactions in such databases presents a significant challenge. Coordinating cross-region transactions (CRTs) is inherently slow due to the extensive geographic distances. Even a few slow transactions can significantly degrade the performance of the entire system. Despite various proposals introduced to optimize the performance by either eliminating CRTs or minimizing the cost of cross-region coordination, CRTs remain crucial for general workloads, with their costs still being dozens of times higher than those of intra-region transactions (IRTs).

This paper contends that existing serializable consistency models are not well-designed for multi-region deployments. The root cause is the strong heterogeneity in deployments: certain transactions (e.g., CRTs) experience significantly higher latency in committing compared to others. To address this, we propose a new layered consistency model specifically tailored for multi-region deployments: region-linearizable serializability (RLS). Specifically, RLS ensures strict serializability (i.e., linearizability) for IRTs from the same region and provides regular serializability for CRTs.

To demonstrate the efficiency and applicability of our new consistency model (i.e., RLS), we design, implement, and evaluate variations of two representative database systems: Spanner and CockroachDB. Our evaluation demonstrates that these variations can significantly enhance performance (e.g., 4.5× throughput for Spanner on average) or provide better functionality without performance degradation (e.g., more consistency guarantees for CockroachDB).

## PVLDB Reference Format:

Ruijie Gong<sup>†</sup>, Haoze Song<sup>†</sup>, Xusheng Chen<sup>\*</sup>, Sen Wang, and Heming Cui.  
Region-Linearizable Serializability: A Practical Consistency Model for Multi-Region Distributed Transactions. PVLDB, 14(1): XXX-XXX, 2020.  
doi:XX.XX/XXX.XX

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights

## 1 INTRODUCTION

Today, cloud providers (e.g., AWS [6], Azure [7], and Huawei [12]) host computing infrastructures across multiple global locations. These infrastructures are typically organized into regions, each strategically designed to offer computing services in proximity to clients. As a result, multi-region deployment has emerged as a prominent choice for cloud-native applications aiming for stringent client-perceived latency, high availability, strong scalability, and effective service localization [9, 18, 40, 43, 50].

In pursuit of scalability and availability, these applications frequently partition and replicate their data storage across multiple servers (nodes). Each partition, known as a data shard, maintains a primary replica. The primary replica is consistently deployed in the region where the majority of access requests originate, ensuring optimal locality. Figure 1 provides an overview of our multi-region deployment model, which is fueled by the desire of global companies to not only build scalable applications but also control with fine granularity where data resides for good performance and meet the data governance policies (e.g., GDPR [1]).

However, supporting serializable ACID transactions in such applications presents a perpetual challenge. Coordinating a cross-region transaction (CRT) is always slow due to geographic distances. For example, a cross-region transaction incurs a  $\sim 50ms$  network delay from Hong Kong to Singapore, whereas the network delay within a region is less than  $2ms$  when powered by modern network technologies (e.g., dedicated inter-datacenter networks [6]).

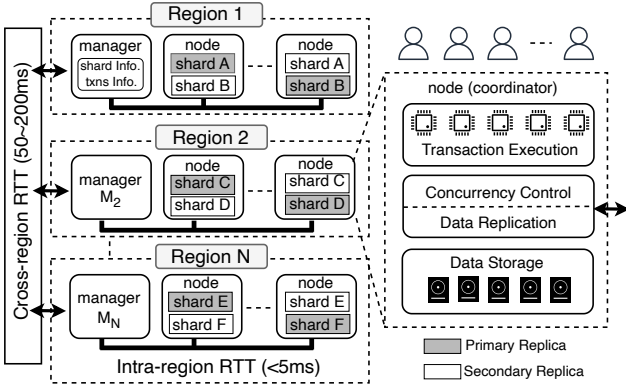
Several impressive works (e.g., [18, 33, 34, 38, 41]) have been proposed to enhance the performance of geo-distributed transactions. These strategies involve redesigning critical aspects of concurrency control protocols or eliminating CRTs by making certain assumptions about the workloads. However, we emphasize that CRTs remain essential for general workloads, such as those with limited prior knowledge of application semantics or those employing interactive transactions. Moreover, the performance disparity between CRTs and IRTs is still significant. For instance, Detock [34], a state-of-the-art geo-distributed transaction protocol meticulously

licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.  
doi:XX.XX/XXX.XX

<sup>\*</sup> Xusheng Chen is the corresponding author.

<sup>†</sup> Ruijie Gong and Haoze Song contribute equally.



**Figure 1:** This diagram presents a typical deployment for multi-region databases. The database is partitioned into multiple data shards spanning over multiple regions. Each shard includes a primary replica and several secondary replicas. Intra-region communications are much faster than inter-region communications.

optimized for CRTs, can execute and commit CRTs in a single cross-region network communication but still incurs an average latency of approximately  $\sim 100ms$  for CRTs and  $\sim 5ms$  for IRTs under default experimental setups. Furthermore, according to the real-world studies [10, 18, 34], as well as our experimental evaluations, even a few slow CRTs can entangle numerous IRTs, leading to deadlocks or aborts. This substantially degrades the overall database performance.

Diverging from these works, we address multi-region transactions by reevaluating and redesigning existing consistency models. Specifically, we contend that current serializable consistency models are inadequately designed for multi-region deployment. The primary issue arises from the inherent heterogeneity introduced by multi-region deployment in various transaction types. As data is closely associated with their respective home region through primary replicas, data access costs vary for transactions originating from distinct regions. We contend that a diverse range of consistency guarantees should be accessible for different transaction types, all while upholding strong consistency in requisite scenarios.

In reality, existing consistency models tend to be either overly strong or can be further enhanced without performance degradation. Aggressive models (e.g., strict serializability model [35]) abstract the entire system as a single node, necessitating heavy synchronizations between all computing servers. This approach leads to poor performance in multi-region deployments. Moreover, this model fails to preserve the advantages of near-client computing: even a local transaction has to be ordered with cross-region ones, which is at odds with the motivation of multi-region deployment. Other models provide weaker consistency (e.g., strong snapshot isolation [17]). While these models may suffice for certain application scenarios, the consistency guarantees for local transactions can be further enhanced without much performance overhead: the communication cost for coordinating local transactions can be cheap when using modern hardware and networks.

In this paper, we propose region-linearizable serializability (for short, RLS), the first consistency model to provide as strong as possible consistency for multi-region databases. RLS treats CRTs and IRTs differently. It ensures strict serializability (i.e., the strongest consistency guarantee) for IRTs from the same region and provides regular serializability for CRTs. Specifically, RLS ensures serializability and “no stale reads” property (a.k.a. regularity) for all transactions (i.e., both CRTs and IRTs) while preserving real-time order only for the transactions that at least access one same region. For a formal definition, we refer readers to §3.

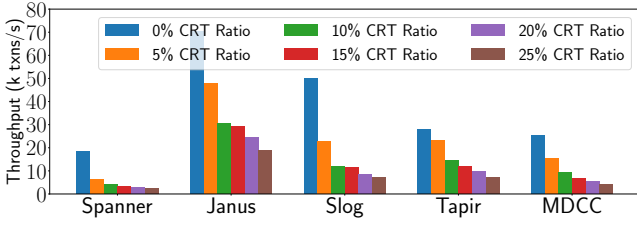
To demonstrate the efficiency and applicability of RLS, instead of creating a new transaction protocol from scratch, we design, implement, and evaluate variations of two database systems: Spanner and CockroachDB (for short, CRDB). We refer to these variants as Spanner-RLS and CRDB-RLS, respectively. We chose these two database systems because they complement each other in both the consistency model (strict serializability versus single-key linearizability) and the design of concurrency control protocols (two-phase locking versus timestamp ordering). These two proof-of-concept prototypes pave the way for efficient and practical optimization of distributed protocols when deployed across multiple regions based on a correct-by-construction approach.

Specifically, Spanner and Spanner-RLS follow the design of traditional pessimistic concurrency control: it orders transactions using locks and commits transactions using the two-phase commit. Our variation (Spanner-RLS) significantly reduces CRTs’ contention footprint (i.e., locking duration), thus allowing more parallel concurrency. As a result, Spanner-RLS attains  $1.16\times$  to  $89.01\times$  higher throughput than Spanner while having similar latency.

We present CRDB-RLS as a practical example to demonstrate that databases using weaker consistency models can also benefit from RLS by tightening their consistency guarantees. The original consistency model (i.e., single-key serializability) used by CRDB is considerably weaker than RLS. We implemented regional semantics into CRDB’s conflict detection protocol, requiring all conflicting IRTs to be ordered within the regions, even if they access different keys. Given that intra-region communication is significantly cheaper than cross-region communications, the performance overhead of CRDB-RLS can be ignored (e.g., less than  $\sim 15\%$  in our evaluation). Conversely, the stronger guarantee efficiently eliminates anomalies caused by violating the real-time order, thus simplifying application development [13, 30].

**Contributions.** In summary, this paper’s contributions stem from a fundamental insight that existing consistency models are inadequately designed for multi-region deployment. To our knowledge, this paper provides the first tailored consistency model for multi-region transactional processing. Our contributions are four-fold:

- We systematically analyze the multi-region deployment model and highlight the limitations of existing consistency models.
- We implemented two distinct system variations, Spanner-RLS and CRDB-RLS, to demonstrate the efficiency and applicability of our novel consistency model (RLS). Both of the two prototypes are built on open-source codebases, and the source code is available at [https://github.com/vldb24p771/spanner\\_rls](https://github.com/vldb24p771/spanner_rls) and [https://github.com/vldb24p771/crdb\\_rls](https://github.com/vldb24p771/crdb_rls), respectively.



**Figure 2: Impact of CRT ratio on the throughput of the state-of-the-art geo-distributed transaction systems.**

- We extensively evaluate these variations and present compelling results showcasing the substantial performance improvements achieved by RLS. Specifically, RLS enhances Spanner’s performance, achieving throughput improvements ranging from 1.16× to 89.01×, and provides more robust consistency guarantees for CRDB without significant performance degradation.
- Spanner-RLS and CRDB-RLS can serve as practical templates for the future adoption of RLS to other concurrency control protocols.

The rest of the paper is organized as follows. §2 discusses the system model of our multi-region databases deployment, the background of geo-distributed transaction processing, and the motivating applications. §3 details our new consistency model: RLS. §4 illustrates the application of RLS to Spanner. §5 delves into CRDB and CRDB-RLS. Finally, a discussion of related works is presented in §6, and §8 concludes the paper.

## 2 BACKGROUND AND MOTIVATION

This section provides background on multi-region databases and the issues of existing geo-distributed transaction protocols.

### 2.1 Multi-region Databases

Multi-region infrastructures motivate our design of RLS. Deploying a multi-region application has the following advantages.

- **[Low data access latency]** Multi-region deployment enables the placement of data in proximity to active client concentrations across different regions. Applications are designed to predominantly access a data shard from a single region, resorting to cross-region data access only when necessary. This approach facilitates applications in offering global data access with ultra-low latency.
- **[Meet data governance policy]** Privacy regulations impose multi-region deployments for global services with strict requirements on data residency. For instance, the General Data Protection Regulation (GDPR) forbids replicating European citizens’ data outside. Thus, multi-region deployments are becoming the de facto choice for multinational companies.
- **[Flexible failure model]** Multi-region deployment allows for data replication across regions, significantly enhancing availability by tolerating complete region failures. In practice, replicating data globally can be costly. Users have the flexibility to define distinct replication policies based on data criticality, such as replicating critical data across regions and keeping other data within the same region.

A typical multi-region deployment model is illustrated in [Figure 1](#). The database is partitioned into multiple data shards spanning over multiple regions. Each shard comprises a primary replica and several secondary replicas. Replicas can be configured to reside in cross-region or intra-region nodes (servers) based on replication policies. Each region is equipped with a centralized transaction manager responsible for globally consistent metadata (e.g., table schema, data placement policy, and globally unique transaction IDs). Nodes can communicate with each other over the network. Our assumptions include a partially synchronized network, and every message in the database is eventually delivered and processed.

### 2.2 Geo-distributed Transaction Protocols

Many influential works have been proposed to optimize the performance of geo-distributed transaction processing. We summarize the state-of-the-art systems in [Table 1](#). All of these systems support at least serializability as their consistency model or can be enhanced to serializability with minor modifications (e.g., MDCC).

Generally, these systems focus on optimizing the cost of coordinating transactions from the protocol scope. For instance, some works attempt to reduce either the number or the overhead of network round-trips in transaction coordination. Tapir [51] improves Spanner’s performance by integrating two-phase commit and consensus protocols into a single framework, eliminating redundant coordination and reducing the WAN round-trips. RedT [52] enhances system performance by further decreasing the network round-trips. RedT targets RDMA-capable networks for local communication and employs a pre-write-log mechanism to eliminate the synchronization of prepare messages (i.e., the first phase in two-phase commits) from the coordinator to primary replicas.

Calvin [42], Slog [38], Detock [34], Ocean Vista [18], Epaxos [32], MDCC [24], and Janus [33] follow the design of deterministic databases, which logically create a global log containing all transactions that have been input into the system. The system then ensures a concurrent execution schedule equivalent to processing all transactions serially in the order they appear in this log (i.e., a partial order). Consequently, after the transaction order is determined, the execution of both CRTs and IRTs can be local. The executors adhere to orders in the logs they receive.

However, all these proposals can not essentially prevent an IRT from being blocked by CRTs (as illustrated in [§3.2](#)). Briefly, Spanner and RedT employ two-phase locking for transaction ordering and commit transactions using two-phase commit. When an IRT conflicts with an ongoing CRT, the IRT has to be blocked for the required locks. Tapir uses a variation of OCC and enforces the IRT to abort in the validation phase, resulting in a high abort rate (as confirmed by other previous papers [10, 18]). Deterministic databases either order IRTs and CRTs together (e.g., Calvin, Janus, Epaxos, and Ocean Vista) or require an IRT to be blocked, waiting for the execution of CRTs that are scheduled ahead (e.g., Slog and Detock) for overly strong consistency guarantees.

Consequently, these systems can cause severe performance issues when CRTs occur in the database, and the contention between the IRTs and CRTs is relatively high. We experimentally studied the impact of CRT ratios on the five latest representative systems in [Figure 2](#). We used YCSB-T workloads with a Zipf parameter of 0.8.

System	Transaction Protocol	Consistency Models	Coordination Blocking
Spanner [13]	Read-write: 2PL + 2PC; Read-only: TSO	SS	Yes
Calvin [42]	Centralized coordinator	SS	Yes
Slog [38]	CRT: Centralized coordinator; IRT: Intra-region Sequencer	SS	Yes
Detock [34]	CRT: Dependency-graph; IRT: Intra-region Sequencer	SS	Yes
Janus [33]	Dependency-graph	SS	Yes
Epaxos [32]	Dependency-graph	SS	Yes
Ocean Vista [18]	TSO (Watermark)	SS	Yes
CRDB [40]	TSO (HLC)	SKL	Yes
RedT [52]	2PL + 2PC	Serial.	Yes
Tapir [51]	Variant of OCC	Serial.	No (by aborting IRTs)
MDCC [24]	Paxos	SI	Yes

**Table 1: This table summarizes the state-of-the-art geo-distributed transaction systems in the literature. These existing systems either block IRTs or enforce the IRTs to abort when the IRTs conflict with an ongoing CRT.**

For detailed evaluation setups, we refer readers to §4.2. From the experimental results, our key observation is that even a few CRTs can significantly degrade the whole system’s performance (e.g., up to 86% degradation with only 5% CRTs).

RLS addresses such issues by rethinking the limitations within existing consistency models and thus is fundamentally different from these proposals. RLS trades off consistency for performance with minimal intrusion (i.e., the consistency tradeoff in RLS is tightly necessary for addressing blocking issues). We regard RLS as orthogonal to these advanced geo-distributed transaction protocols. Consequently, new protocols may benefit from the methodology of RLS and the key optimizations of these advanced protocols.

### 3 REGION-LINEARIZABLE SERIALIZABILITY

In this section, we formally define our new consistency model: Region-linearizable Serializability (RLS), which is specially tailored for the database systems that are deployed in near-client computing facilities (e.g., Regions [6] in AWS). For comparisons with other serializable consistency models, we refer readers to §6.1.

#### 3.1 Definition of RLS

For clarity, we adopted the formalism from existing works [21]. Table 2 summarizes the notations used in our definition, which will be further illustrated later. Without loss of generality, we consider an OLTP service (either a relational database or a transactional key-value store) handling data objects identified by unique keys. We use  $\mathcal{K}$  to represent the global key spaces.  $\mathcal{K}$  is divided into multiple disjoint *shards* to facilitate transaction processing across many nodes (servers), which is common in multi-region deployed data-intensive applications.

**Shard Groups (e.g., grouped by regions).** Grouping semantics is one of the foremost distinctions of RLS when compared to other existing serializable consistency models. Specifically, RLS divides data shards into disjoint groups and ensures linearizability (i.e., the strongest consistency) for intra-group operations. For inter-group operations, RLS provides regional serializability.

Note that grouping and sharding construct a two-level division of the global key space ( $\mathcal{K}$ ): the OLTP service has many disjoint groups,

and each group contains a number of (not necessarily equivalent) disjoint shards. This division serves different purposes; sharding is for horizontally scaling the service to run on many servers, and the division policy is usually for reducing the ratio of cross-shard transactions to achieve high efficiency [34]. On the other hand, grouping is a consistency strategy where cross-group external ordering requirements are typically less important. Such semantics are usually related and directed by geo-distributed (multi-region) deployments.

We regard this two-level framework as essential because it efficiently bridges the division based on application semantics (e.g., the data items grouped by warehouse ID in TPC-C [14]) and division based on deployment topology (e.g., the data shards grouped by regions in geo-distributed deployments).

**Transactions and Operations.** Clients interact with the OLTP service through transactions. Each transaction comprises several single-key read or single-key write *operations*. Formally, each transaction  $T$  is a tuple  $(\Sigma_T, \xrightarrow{to})$ , where  $\Sigma_T$  is the set of operations in  $T$ , and  $\xrightarrow{to}$  is a total order on  $\Sigma_T$ . Each operation is either a read (denoted as  $o_1 = r(k_1, v_1)$ ) or a write (denoted as  $o_2 = w(k_2, v_2)$ ). We use  $\mathcal{R}_T = \{k | r(k, v) \in \Sigma_T\}$  to denote  $T$ ’s read set and  $\mathcal{W}_T = \{k | w(k, v) \in \Sigma_T\}$  as  $T$ ’s write set.

It should be noted that RLS, as a consistency model, does not essentially require the read and write set of each transaction to be determined upfront, which is a common but restrictive assumption in existing deterministic databases [24, 32, 34, 38, 42]. Essentially, RLS supports general transaction semantics (to be illustrated in our example system: Spanner-RLS, §4). In RLS, we consider a general transaction  $T_1 = \{r(x, n), w(n, v)\}$  that reads the value of key  $x$  as the key for the write operation, where  $\mathcal{W}_{T_1}$  can not be obtained before execution.

**Conflicts and Relevance.** We say two transactions conflict with each other if they access the same key, and at least one of the two accesses is “write” (which is known as read-write conflicts and write-write conflicts in other papers). We say a transaction  $T$  is relevant to shard group  $g$  if  $T$  accesses at least one key owned by  $g$ , and we use  $\mathcal{G}_T$  to represent the set of all groups relevant to  $T$ .



Ops	$o$	A database operation, e.g., read, write, insert, scan, etc.
	$r(k, v)$	Read the value $v$ using key $k$
	$w(k, v)$	Write value $v$ for key $k$
Txns	$T$	A transaction consists of operations ( $\Sigma_T$ ) with order ( $\xrightarrow{to}$ )
	$\mathcal{R}_T$	Read Set of Transaction $T$
	$\mathcal{W}_T$	Write Set of Transaction $T$
	$\mathcal{G}_T$	The set of all shard groups relevant to $T$
Data	$\mathcal{K}$	Global Key Space
	$g$	A shard group contains multiple shards
Order	$\mathcal{H}_i$	Transaction history on $node_i$ , $\mathcal{H}_i = (\mathcal{E}_i, po_i, \tau_i)$
	$\mathcal{H}$	Transaction history of the whole system, $\mathcal{H} = \bigcup \mathcal{H}_i$
	$\mathcal{S}$	A totally ordered serializable schedule for all transactions
	$\xrightarrow{rb}$	Real-time order imposed by runtime execution
	$\xrightarrow{so}$	Order for operations in $\mathcal{S}$
	$<_S$	Order for transactions in $\mathcal{S}$

**Table 2: Preliminaries and notations for RLS.**

Formally,

$$\mathcal{G}_T = \{g \mid \exists k : k \in g \wedge k \in (\mathcal{W}_T \cup \mathcal{R}_T)\}$$

**History and Equivalence.** A history of a data  $node_i$  ( $server_i$ ) is an associative triple  $\mathcal{H}_i = (\mathcal{E}_i, po_i, \tau_i)$ , where  $\mathcal{E}$  is a set of operations;  $po$  is a partial ordering on  $\mathcal{E}$  into processes; and  $\tau$  divides  $\mathcal{E}$  into transactions. We say two histories ( $\mathcal{H}_1$  and  $\mathcal{H}_2$ ) are equivalent if they have the same  $\mathcal{E}$ ,  $po$ , and  $\tau$ . Intuitively, two equivalent histories have the same sequence of operations for each client process and thus are indistinguishable inside the database.

**Real-time order.** An order of transactions is usually considered as a set of *return before* relations [25]. In our paper, we say a transaction  $T_1$  precedes another transaction  $T_2$  if  $T_1$  finishes (commits) before  $T_2$  starts (i.e., arrives at the database system), denoted as  $T_1 \xrightarrow{rb} T_2$ .

**Definition of RLS.** We then define RLS using the notations above. We say that an OLTP service ensures RLS, if for all execution histories,  $\mathcal{H} = \bigcup \mathcal{H}_i$ , are equivalent to a serial schedule  $\mathcal{S}$  and the following three properties hold for  $\mathcal{S}$ .

- **Serializability.** There exists serial schedule  $\mathcal{S}$  with total ordering  $so$  on  $\mathcal{E}$  such that ①  $\mathcal{S}$  is equivalent to  $\mathcal{H}$ ; and ② no two transactions overlap in  $so$ , i.e., either

$$o_1 \xrightarrow{so} o_2, \forall o_1 \in T_1, \forall o_2 \in T_2$$

or

$$o_2 \xrightarrow{so} o_1, \forall o_1 \in T_1, \forall o_2 \in T_2$$

Therefore, the property ② infers that  $\xrightarrow{so}$  defines a total order  $<_S$  among all transactions.

- **No Stale Reads.** Formally, for any two transactions  $T_1$  and  $T_2$

$$\mathcal{W}_T \cap (\mathcal{W}_T \cup \mathcal{R}_T) \neq \emptyset \wedge T_1 \xrightarrow{rb} T_2 \implies T_1 <_S T_2$$

- **Real-time Ordering inside all Shard Groups.** Formally,

$$\mathcal{G}_{T_1} \cap \mathcal{G}_{T_2} \neq \emptyset \wedge T_1 \xrightarrow{rb} T_2 \implies T_1 <_S T_2$$

### 3.2 Performance Issues in Strict Serializability

Strict serializability (SS), the most substantial consistency level for distributed databases, ensures that a replicated distributed database works as a single node that executes all client transactions serially. The serial order respects the real-time relations (i.e., the “return before” relation in §3.1) among all client transactions.

However, the strong guarantees of SS always come up with high-performance costs, especially when deployed in a multi-region environment. This has led both academia and industry to seek weaker consistency models. For example, numerous new consistency models were proposed in recent years (see §6.1), and almost all industrial systems do not provide SS by default.

### 3.3 Practical Implications

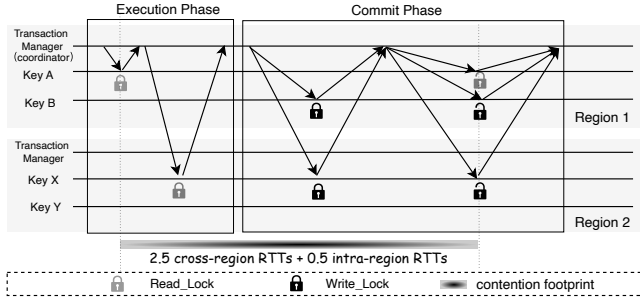
In essence, the guarantee of SS is considered excessive for many multi-region application scenarios. Specifically, ensuring a real-time relation between transactions pertains to external (out-of-band) causal relations among transactions. Since a system is unaware of external relations, SS regards *all* pairs of transactions without overlapping lifetime as potentially causally related and pertains to their ordering, albeit most transactions are independent.

RLS ensures the “no stale reads” property for all transactions, effectively preventing most application-level anomalies [44]. Additionally, RLS enforces real-time ordering among transactions accessing interleaved regions (i.e., conflict IRTs and CRTs), including transaction ordering requirements inferred by transitivity. Compared to SS, the only anomalies in RLS may arise from the potential disruption of real-time ordering among transactions happening independently within non-overlapping regions.

We argue that such anomalies do not compromise the correctness of multi-region databases for two primary reasons. First, multi-region databases optimally leverage data access locality to assign shards to regions (see §2.1). Typically, each region manages (e.g., being the leader of) shards containing data of nearby clients, making two transactions accessing non-overlapped regions causally unrelated. Thus, prioritizing their real-time order will not introduce application-level anomalies.

Second, the time window for breaking causal relations is narrow. RLS necessitates “no stale reads” for all transactions, whether intra-region or cross-region. To sever the causal relationship between two transactions, external communication must conclude faster than a transaction’s lifetime. Specifically, consider two transactions  $T_2, T_3$  accessing non-overlapped regions, where  $T_2 \xrightarrow{rb} T_3$ . If anomalies were present, it would imply the existence of another transaction  $T_1$  accessing both  $T_2$  and  $T_3$ ’s regions, leading to a final serial order of  $T_3 <_S T_1 <_S T_2$  (as depicted in Figure 11a). However, as RLS also mandates “no stale reads”,  $T_1$  must be concurrent with  $T_2$  and  $T_3$ , implying that the external causal relation must conclude within  $T_1$ ’s lifetime.

Therefore, RLS possesses the unique potential to significantly enhance the scalability and latency of multi-region databases while maintaining correctness and programmability. RLS stands out as the



**Figure 3:** This diagram shows how Spanner orders a CRT using 2PL and commits it using 2PC in a multi-region deployment. Replicas are removed for readability.

pioneering consistency model that takes into account real-world deployments and the inherent locality feature of data.

## 4 SPANNER AND SPANNER-RLS

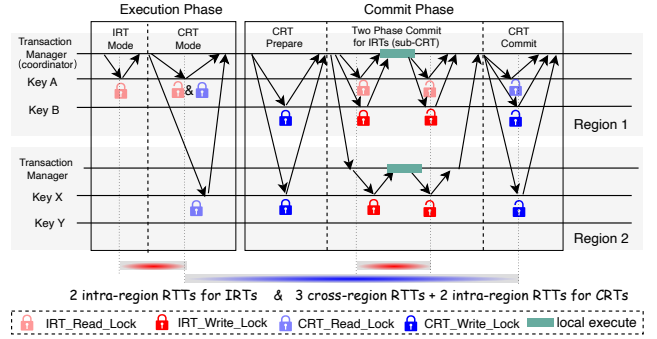
In this section, we present the design, implementation, and evaluation of Spanner and Spanner-RLS. We regard Spanner-RLS as a specific instance of using our new consistency model (RLS) to advance the performance of existing strictly serializable databases.

### 4.1 Protocols and Implementations

**Spanner Background.** Google’s Spanner provides strict serializability (a.k.a. external consistency) for read-write transactions by coordinating them using two-phase locking (2PL) and then committing the transactions using two-phase commit (2PC).

Figure 3 presents an example, where a transaction  $T$  reads the keys  $A$  and  $X$ , and then updates the keys  $B$  and  $X$ . The keys  $A$  and  $B$  are located at different data nodes in *region 1*, and the key  $X$  is located at a data node in *region 2*. To commence,  $T$  is forwarded to the transaction manager that  $T$  firstly accesses, acting as the coordinator and assigning a globally unique TID to  $T$ . In this example, the transaction manager of *region 1* serves as the coordinator since  $T$  reads the key  $A$  in the first access. Subsequently, the coordinator sequentially executes all the transaction operations. During execution, it acquires read locks for each read operation and buffers write operations in temporary memory. After buffering all writes, the coordinator obtains exclusive locks for all write keys (i.e., the keys  $B$  and  $X$ ) and installs the writes if all the required locks are acquired. Following this, all the locks (both read and exclusive locks) are released immediately. As such,  $T$  is successfully executed and committed.

We are now prepared to introduce the performance issues in Spanner. Spanner coordinates intra-region transactions (IRTs) and cross-region transactions (CRTs) similarly, where a read lock held by a CRT will block all writes from both CRTs and IRTs, and an exclusive lock held by a CRT will prevent all reads and writes, correspondingly. Consequently, a CRT’s contention footprint (i.e., the lock duration) is extremely large. More than two cross-region network round trips will block all conflict transactions. In our example, it takes 2.5 cross-region network round trips and 0.5 intra-region round trips. Even worse, such blockings can rapidly accumulate



**Figure 4:** This diagram shows how Spanner-RLS orders a CRT using a variant of 2PL and commits it using 2PC.

through transitive relations. For instance, considering another CRT  $T'$  that accesses the key  $B$  and  $Y$ ,  $T'$  can successfully acquire the exclusive lock on the key  $Y$  while having to wait for the lock on the key  $B$ . Consequently, all other IRTs and CRTs that access the key  $Y$  have to compete with  $T'$  for ownership of the lock on the key  $Y$ , enlarging the affected key space.

**Spanner-RLS.** Following the methodology of our new consistency model (RLS), we treat the IRTs and CRTs differently in the variation of Spanner, termed Spanner-RLS. To achieve this, we distinguish the locks acquired by IRTs and CRTs. The two types of locks order transactions independently. A CRT lock does not block IRTs, and vice versa. In our design, CRT locks only provide the functionality for reservation and maintain the partial order between CRTs.

Algorithm 1 shows the pseudocode of Spanner-RLS, and we highlight the regional semantics in blue. Figure 4 illustrates how Spanner-RLS executes and commits the same transaction  $T$ . Without loss of generality, we assume that the read and write sets of the transaction are unknown to the transaction manager. Therefore, Spanner-RLS can support general transactions without prior knowledge.  $T$  is firstly executed as an IRT and acquire IRT\_Read\_Lock for the key  $A$  during the execution.  $T$  switches to the CRT mode when it attempts to perform remote reads (i.e., reads the key  $X$  in *region 2*). Before that, it releases the acquired IRT\_Read\_Lock for the key  $A$  and updates the lock type to CRT\_Read\_Lock. If the key  $A$  is already exclusively locked by other CRTs,  $T$  aborts and directly retries using CRT mode. Otherwise,  $T$  successfully enters the CRT mode and employs CRT\_Read\_Lock for the remaining reads (e.g., reads the key  $X$ ). Since all the changes are handled by intra-region communication, it will not incur much overhead. Then, following the transaction logic,  $T$  computes its write set and proceeds to the commit phase.

In the commit phase,  $T$  acquire CRT\_Write\_Lock for the key  $B$  and  $X$ . When all CRT\_Write\_Lock is successfully acquired (i.e., the order between  $T$  and other CRTs has been determined), the coordinator notifies all transaction managers of the region that  $T$  accessed. Each transaction manager commits  $T$  independently using IRT mode. As we already allow  $T$  to hold the read locks for all read keys (i.e., the keys  $A$  and  $X$ ), the read keys of  $T$  can not be modified by any other CRTs. In case of any IRTs that have modified the read keys of  $T$ , we re-execute it locally. The tricky is that even

### Algorithm 1: Algorithm of Spanner-RLS

```

1 function Execution phase:
2   read_set & write_set  $\leftarrow \emptyset$ 
3   txnType  $\leftarrow$  IRT  $\triangleright$  Start a new transaction as IRT.
4   touchedRegions  $\leftarrow \emptyset$   $\triangleright$  Regions involved in the transaction.
5    $\triangleright$  Execute the transaction commands, which triggers the events:
6   Event read(key)
7     value = find_record(key)
8     read_set.append(key)
9     touchedRegions  $\leftarrow$  touchedRegions  $\cup$  key.region
10    if |touchedRegions|  $\geq 2$  then
11      txnType  $\leftarrow$  CRT
12      Release_IRT_Read_Lock(k) for all  $k \in$  read_set
13      CRT_Read_Lock(k) for all  $k \in$  read_set
14    if txnType == IRT then
15      IRT_Read_Lock(key)
16    else
17      CRT_Read_Lock(key)
18   Event write(key, value)  $\triangleright$  Writes are only buffered
19     write_set.append(<key, value>)
20     Execute Line 9 ~ 13
21 function Commit phase:
22   if txnType == IRT then
23     IRT_Write_Lock(k) for all  $k \in$  write_set
24     wait for all ACKs from storage  $\triangleright$  Abort if fail
25     Commit(txn)
26     Release_IRT_Read_Lock(k) for all  $k \in$  read_set
27     Release_IRT_Write_Lock(k) for all  $k \in$  write_set
28   else
29     CRT_Write_Lock(k) for all  $k \in$  write_set
30     wait for all ACKs from storage  $\triangleright$  Abort if fail
31     Send Commit to txn managers in  $r, r \in$  touchedRegions
32      $\triangleright$  Each transaction manager commits the transaction as IRT
33     wait for all ACKs from the txn managers  $\triangleright$  Abort if fail
34     Commit(txn)
35     Release_CRT_Read_Lock(k) for all  $k \in$  read_set
36     Release_CRT_Write_Lock(k) for all  $k \in$  write_set
37

```

if the re-execution depends on remote reads, we can defer the IRT lock acquisition of the execution until the remote reads have been finished since we now have obtained the read and write set of the transaction. One exception is that the transaction  $T$ 's read and write set may differ during re-execution, or  $T$  has cycle dependency in the transaction logic (e.g., the execution in *region 1* depends on the reads in *region 2*, the execution in *region 2* depends on the reads in *region 1*, and both the read keys of *region 1* and *region 2* has been changed). In such cases, we can easily revert from Spanner-RLS to Spanner by using IRT locks directly for the CRT.

We then analyze the tradeoff in Spanner-RLS. By differentiating CRTs and IRTs locks, Spanner-RLS eliminates both the “commit blocking” and the “coordination blocking” for IRTs. In our example, the contention footprint for conflict IRTs is reduced to two intra-region network round-trip communication. On the other hand, CRTs may incur slightly more communication costs between the coordinator and the transaction managers. However, in practical

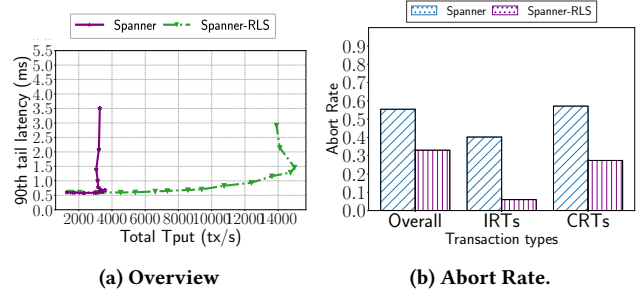


Figure 5: Overall performance and abort rate of Spanner and Spanner-RLS on YCSB-T (default setting) using NO\_WAIT.

workloads, IRTs are always the dominators and are more critical and sensitive to performance. In fact, if CRTs constitute the majority of the workloads, the performance degradation induced by the heterogeneous network is less of a problem. Hence, Spanner-RLS can fall back to the classic Spanner, which is switchable at runtime. **Read-only Transactions.** Leveraging error-bounded timing service (e.g., TrueTime API), Spanner can execute read-only transactions in a single network round trip. Using its TrueTime API, Spanner assigns a commit timestamp to each transaction, guaranteed to be between the transaction’s real start and end times. Therefore, when using the TrueTime API for read-only transactions, they can safely read from the replicas without coordination. Spanner-RLS follows this design for enhanced performance. In our evaluation, we emulated TrueTime error as 10 ms, which is used in the previous paper [21] and matches the p99.9 value observed in practice.

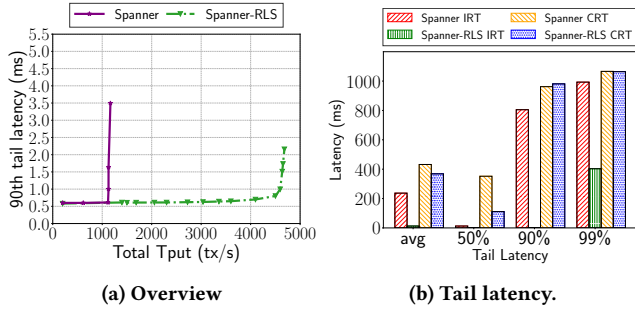
## 4.2 Evaluation and Discussion

**4.2.1 Experimental Setups.** We implemented Spanner-RLS in C++ utilizing the third-party implementation [2] since the original version of Spanner is not open-sourced. We employed libevent for message passing between processes on distinct nodes and between threads in the same process. Transactions were implemented as stored procedures containing read and write operations over a set of keys.

**Cluster Setups.** All experiments were conducted on our cluster comprising 10 machines, each with a 2.60GHz Intel E5-2690 CPU with 24 cores, 40Gbps NIC, and 64GB memory. We executed each data node (shard) in a docker container and utilized tc [22] to regulate the RTT among nodes. The server ran on the Ubuntu 18.04 operating system.

**Deployments.** To simulate a multi-region deployment, we abstracted each server as an individual region. Subsequently, We set the cross-region round-trip latency as 50ms using tc, aligning with the real-world statistics [8]. We partitioned the database into 300 data shards, with each region containing 30 data shards and a replication level of 3, alongside transaction managers. By default, we utilized 40 clients per region to attain the peak throughput for both Spanner and Spanner-RLS.

**Workloads.** We employed the standard YCSB-T benchmark for our evaluation. We generated a total of 3,000,000 keys, distributing 100,000 keys per shard. Each transaction had 10 operations, encompassing 5 read operations and 5 read-modify-write operations. By



**Figure 6: Overall performance and latency of Spanner and Spanner-RLS on YCSB-T (default setting) using WAIT\_DIE.**

default, we tuned the percentage of CRTs to be 10% and varied the amount of contention in the system by choosing keys according to a Zipf distribution with a Zipf coefficient = 0.75 (medium and high contention). Each of our experiments lasted 3 minutes, with the first 30s and the last 30s excluded from results to avoid performance fluctuations during start-up and cool-down.

**Deadlock Mechanisms.** We considered two different deadlock mechanisms in our evaluation since these mechanisms impart different scopes to how RLS benefits Spanner.

- **NO\_WAIT.** When using this deadlock mechanism, if a lock request is denied, the database will immediately abort the requesting transaction, and the client will retire the transaction.
- **WAIT\_DIE.** Unlike NO\_WAIT, WAIT\_DIE allows a transaction to wait for the requested lock if the transaction is older than the one that holds the lock. Otherwise, the transaction is forced to abort and will be retired by the client.

We do not consider other deadlock mechanisms since they are either subsumed by the two mechanisms or will incur significant overhead in a multi-region deployment. For instance, deadlock detection necessitates a centralized deadlock detector for cycle detection, which can be expensive due to cross-region communication.

**4.2.2 Performance Overview.** We first evaluated the performance under the default setting. For an apple-to-apple comparison, we refrained from using prior knowledge of read and write sets in our experiments, even though the read and compose set of YCSB-T can be revealed before execution. As shown in Figure 5a and Figure 6a, Spanner-RLS significantly outperformed Spanner on YCSB-T. In particular, Spanner-RLS achieved 3.95 $\times$  and 4.27 $\times$  higher peak throughput when utilizing NO\_WAIT and WAIT\_DIE, respectively. Spanner-RLS’s 90th resembled that of Spanner, essentially representing the IRTs latency. We observed that employing the NO\_WAIT mechanism resulted in substantially higher throughput than using WAIT\_DIE, given the default workload’s write-intensive nature with medium contention.

To comprehend how our new design contributes to performance improvement, we gathered data on the abort rate for NO\_WAIT and tail latency for WAIT\_DIE when both Spanner and Spanner-RLS achieved the peak throughput. Figure 5b and Figure 6b illustrate the results.

Regarding NO\_WAIT, Spanner-RLS can efficiently reduce the abort rate for both IRTs and CRTs. The overall abort rate decreased from

56% to 33% (i.e., 41.1% reduction). In particular, Spanner-RLS achieved a more significant reduction for IRTs (from 40% to 6%) due to the “non-blocking” property in IRT coordination and commitment. It’s worth mentioning that the 6% abort rate was only caused by the contention among IRTs. Meanwhile, the abort rate of CRTs also saw a reduction. However, compared to IRTs, CRTs still exhibited a much higher abort rate (i.e., 27.42%) due to the larger contention footprint.

For WAIT\_DIE, Spanner-RLS achieved a significantly lower average latency for IRTs, while the average latency of CRTs was roughly the same as in Spanner. This is attributed to the fact that in Spanner-RLS, an IRT will never be blocked by CRTs. The results on 50th and 90th latency support this assertion. Both Spanner and Spanner-RLS exhibited low 50th latency, while the 90th latency of Spanner and Spanner-RLS was 805ms and 1.4ms, respectively. The 99th latency of Spanner-RLS increased due to the queuing effect in the software stack.

Next, we delve into understanding how Spanner-RLS and Spanner are affected by various workload parameters. These experiments were conducted using YCSB’s APIs as they offer flexibility in configuration.

**4.2.3 Impact of Concurrency.** We first compare the performance of Spanner-RLS and Spanner under various concurrencies. As illustrated in Figure 7a and Figure 8a, Spanner’s throughput reached saturation rapidly as the number of clients increased. Consequently, the peak throughput of Spanner was 3911 and 1181 transactions per second using NO\_WAIT and WAIT\_DIE, respectively. In contrast, Spanner-RLS could serve more clients and achieve a substantially higher peak throughput.

**4.2.4 Impact of CRT Ratio.** We studied the impact of the CRT ratio by fine-tuning the workload generation. As shown in Figure 7b and Figure 8b, when CRTs were enabled, Spanner experienced severe performance degradation (e.g., throughput dropping from 18526 transactions per second to 6184 transactions per second when the CRT ratio increased from 0% to 5%), aligning with our discussion in §1 and §2.2. In contrast, Spanner-RLS’s performance degraded slightly, attributed to the elimination of cross-region costs for IRTs. In scenarios where all transactions were IRTs (i.e., a special case in our experiments), Spanner-RLS demonstrated slightly lower throughput than Spanner due to the cost for extra steps in concurrency control (i.e., checking transaction types even if all transactions are IRTs). With a continuous increase in CRT ratios, the throughput of both Spanner and Spanner-RLS decreased due to the cross-region communication cost. In practice, the CRT ratio of workloads should not be too high since the cost of CRT itself is still relatively high compared to IRTs. Real-world workloads show good data locality under multi-region deployment (§2.1), facilitating low-latency data access.

**4.2.5 Impact of Cross-Region RTT.** Next, we studied the impact of the cross-region network delays, a critical factor affecting the overall cost of CRTs. Larger cross-region network delays generally result in longer transaction coordination and commit times for CRTs. The results, illustrated in Figure 7c and Figure 8c, clearly indicate that Spanner-RLS outperforms Spanner regardless of the cross-region network delays. In fact, Spanner-RLS demonstrates



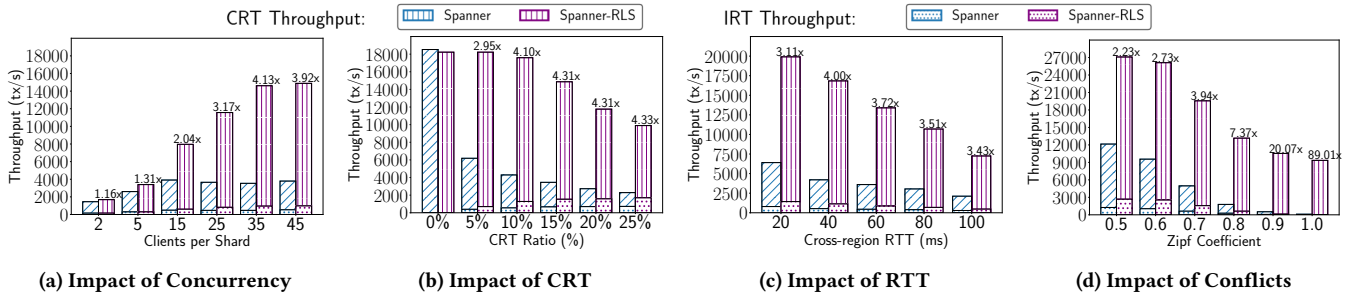


Figure 7: Performance of Spanner and Spanner-RLS on YCSB-T with different experimental parameters using NO\_WAIT.

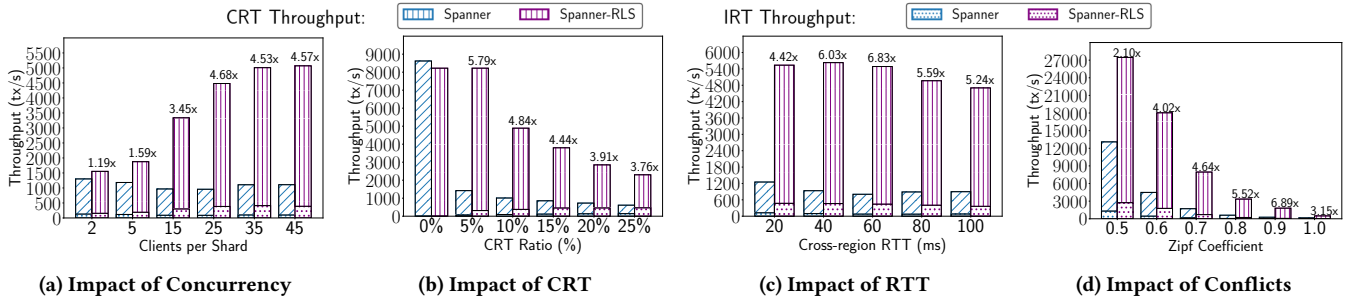


Figure 8: Performance of Spanner and Spanner-RLS on YCSB-T with different experimental parameters using WAIT\_DIE.

more when the network delays are moderate (e.g., 40s and 60s for a cross-region network round trip). In addition, the network delay amplifies the advantages of Spanner-RLS while also affecting Spanner-RLS’s CRTs, causing a drop in throughput from 1429 transactions per second to 476 transactions per second as the network delay increased from 20 seconds to 100 seconds using NO\_WAIT.

**4.2.6 Impact of Contention.** In our final experiment, we compared the performance under various contention by adjusting the skewness of Zipf distribution in YCSB-T while keeping other parameters consistent with the default settings. The results depicted in Figure 7d and Figure 8d consistently show that the Spanner-RLS’s throughput outperforms Spanner’s across all levels of contention. This improvement stems from Spanner-RLS efficiently reducing the contention footprint by ordering IRTs and CRTs independently. As discussed in §2.2, the contention footprint of IRTs eliminates both “coordination blocking” and “commit blocking”, which is extremely expensive in multi-region deployments. As expected by our analysis, Spanner-RLS gained larger margins under high contention. This is because, under high contention, IRTs in Spanner have more chance to be blocked by CRTs, leading to poor performance. It should be noted that the performance of Spanner-RLS also degraded due to the cost of acquiring locks, and NO\_WAIT consistently outperformed WAIT\_DIE since WAIT\_DIE suffers more from lock thrashing and timestamp allocation when the contention is higher.

### 4.3 Takeaways

By adhering to the principles of RLS, developers can significantly enhance the performance of Spanner. Spanner-RLS serves as a practical example, illustrating how RLS can assist multi-region databases

in achieving an optimal balance between consistency and performance. Further enhancements in the performance of Spanner-RLS can be achieved by implementing advanced optimizations (e.g., a pre-write-log mechanism in RedT [52]), which is orthogonal to our paper. We believe our study on Spanner-RLS holds the potential to guide future research by encouraging researchers and developers to consider multi-level two-phase locking and two-phase commit. Even though the design of Spanner-RLS may not be directly applicable to other strictly serializable concurrency protocols due to potential differences in transaction coordination mechanisms, the fundamental concept of RLS remains applicable.

## 5 CRDB AND CRDB-RLS

In this section, we present the design, implementation, and evaluation of Cockroachdb (for short, CRDB) and CRDB-RLS. We show that RLS has the potential to evolve the consistency model of existing databases (i.e., providing tighter and stronger consistency guarantees) without sacrificing performance.

### 5.1 Protocols and Implementations

**CRDB Background.** CRDB [40] is an open-source production-grade database system that began as an external Spanner clone. Same as Spanner, CRDB aims to build a resilient geo-distributed SQL Database with serializable ACID transactions.

Overall, CRDB provides single-key linearizability (i.e., no stale reads for each key) by supporting multi-version timestamp ordering (MVTO). The transaction manager nodes in CRDB are the special nodes for interacting with clients, assigning timestamps to transactions, and driving the coordination of transactions. CRDB assumes

**Algorithm 2: Algorithm of CRDB-RLS Coordinator**

```

1 function CRDB-RLS Coordinator:
2   inflightOps  $\leftarrow \emptyset$   $\triangleright$  Ongoing operations.
3   touchedRegions  $\leftarrow \emptyset$   $\triangleright$  Regions involved in the transaction.
4   txnTS  $\leftarrow \text{now}()$   $\triangleright$  Timestamp of the transaction.
5   for op  $\leftarrow$  KV operation received from SQL layer do
6     if op.commit then
7       op.deps  $\leftarrow$  inflightOps
8       send (commit, txnTS) to transaction managers
9       wait for all ACKs
10    else
11      r  $\leftarrow$  op.key.region
12      if r  $\notin$  touchedRegions then
13        txnTS  $\leftarrow \max(\text{txnTS}, \text{GetFinishedTs}(r))$ 
14        VerifyReads(txnTS)
15        touchedRegions  $\leftarrow$  touchedRegions  $\cup \{r\}$ 
16      op.deps  $\leftarrow \{x \in \text{inflightOps} \mid x.\text{key} = \text{op.key}\}$ 
17      inflightOps  $\leftarrow (\text{inflightOps} - \text{op.deps}) \cup \{\text{op}\}$ 
18      resp  $\leftarrow$  send(op, keyLeader(op.key))
19      txnTS  $\leftarrow \max(\text{txnTS}, \text{GetFinishedTs}(r))$ 
20      VerifyReads(txnTS)

```

a maximum clock offset among transaction managers (i.e., using 500ms by default), which is critical to its correctness.

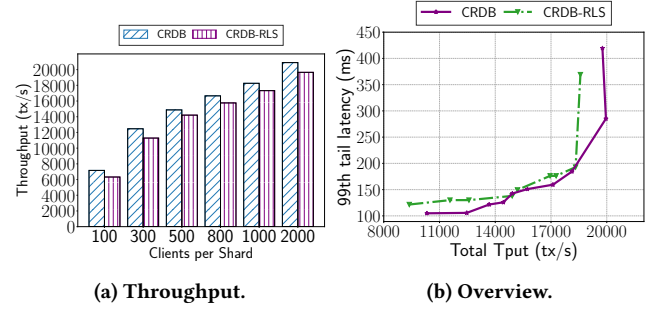
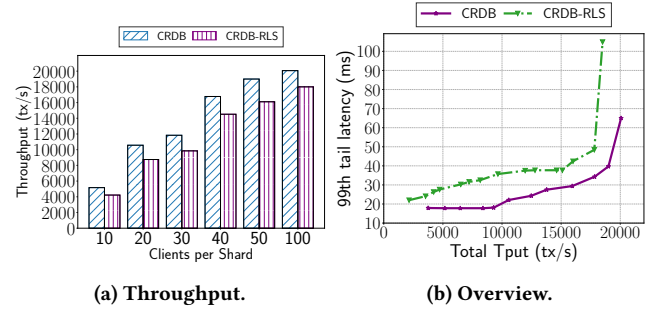
CRDB’s consistency model (i.e., single-key linearizability) is strictly weaker than RLS as CRDB ensures only a subset of RLS’s guarantees: CRDB does not preserve the real-time ordering for any pair of non-conflicting transactions, while RLS provides real-time order among IRTs in the same region. We refer readers to §6.1 for detailed comparisons between RLS and SKL.

CRDB makes such a design choice because, without the two-layered design of RLS, the developers could only choose the consistency model from two extremes on the spectrum: either enforcing all real-time constraints among non-conflicting transactions (i.e., strict serializability) or enforcing none of them. Since implementing strict serializability across regions can easily overwhelm the benefits of data locality and lead to terrible performance, CRDB had to select the latter choice. Consequentially, CRDB cannot even ensure the causal relation of two transactions from the same client when they access different keys in the same region.

CRDB performs its reads and writes at its commit timestamp and relies heavily on multi-version concurrency control (MVCC) to process concurrent requests. When a transaction is detected as in conflict with other transactions, CRDB adjusts the commit timestamp of the transaction to ensure single-key linearizability. Since conflict detection is only conducted in the critical granularity, CRDB does not provide any consistency guarantees when two transactions are not in conflict with each other.

RLS provides a better design point in the spectrum. To demonstrate the pros and cons of RLS, we re-design the protocol of CRDB by incorporating multi-region semantics into the part of conflict detection and name the delivery as CRDB-RLS.

**CRDB-RLS.** Algorithm 2 shows the pseudocode of CRDB-RLS, and we highlight the multi-region semantics in blue. In order to

**Figure 9: Performance Comparison on Micro Workload.****Figure 10: Performance Comparison on YCSB-T (skewed).**

achieve RLS, for any two transactions  $T_1$  and  $T_2$  that access overlapped regions, if  $T_1$  have finished, CRDB-RLS must ensure that  $T_2$ ’s timestamp is larger than  $T_1$ ’s. CRDB achieves such guarantees by comparing with  $T_1$ ’s write timestamp when  $T_2$ ’s read arrives. Specifically, if  $T_2.ts > T_1.ts$ ,  $T_2$  must see  $T_1$ ’s write; otherwise, if  $T_2.ts < T_1.ts$ ,  $T_1$  may still finish before  $T_2$  starts due to clock skewness, but the skewness should have an assumed bound (i.e., 500ms in the codebase of CRDB). Consequentially, if  $T_1.ts - T_2.ts < bound$ , CRDB cannot determine the order between  $T_1$  and  $T_2$ . In such a case, CRDB enforces  $T_2$  to abort and then lets it retry automatically.

Our observation of CRDB-RLS is that, to know whether a transaction  $T_1$  may have finished before  $T_2$  starts, instead of directly comparing the timestamps, a more intuitive method should be using active inquiry. CRDB uses transaction managers to maintain the status of each transaction, and one can get the status of each transaction from the managers. In its original design, CRDB does not adopt this active inquiry design because it can be inefficient and non-scalable to let all transactions contact a single node in a geo-distributed deployment. However, RLS’s region-based approach enabled CRDB to record transactions’ state in a per-region manner (Algorithm 2, Line 13) and let each transaction inquire only relevant regions’ managers (Algorithm 2, Line 14), achieving both stronger consistency and high efficiency.

Same as CRDB, CRDB-RLS allocates timestamps using hybrid-logical clocks (HLC), where physical time is based on a node’s coarsely-synchronized system clock, and logical time is based on Lamport’s clocks. Such a mechanism benefits the design of CRDB-RLS as valid transactions can now adjust their timestamps based on logical time without re-obtaining new timestamps frequently.

**Implementations.** We implemented CRDB-RLS using the open-source codebase of CRDB with version v20.3 from the official sites. We modified the logic of obtaining a valid timestamp by recoding all used timestamps inside a region using sets. Our implementation is orthogonal to the optimizations introduced by its origin paper [40] (e.g., write pipelining, parallel commits, and follower read).

## 5.2 Evaluation and Discussion

We used the same hardware and cluster setups as Spanner. For the deployments, we used 10 containers spanning uniformly over the machines. Each container runs as a CRDB node. We used the multi-region SQL for Table partition and data replication.

**5.2.1 Performance on Micro Workload.** As the configuration file of CRDB is complex and always critical to the performance, we calibrated our results by running the most basic built-in workloads provided by the CRDB codebase (i.e., a transaction reads and writes to three keys spread uniformly across the cluster). The results are shown in Figure 9. The peak throughput of CRDB was  $\sim 20k$  tps, which aligns with the results shown on the official sites. Therefore, we believe our experimental results are representative.

In addition to CRDB, we also test the performance of our implemented variation: CRDB-RLS. CRDB-RLS’s peak throughput ( $\sim 18.6k$  tps) is slightly lower (7%) than CRDB, while the latency is roughly the same. The slight performance degradation is caused by the higher cost of obtaining valid timestamps for IRTs.

**5.2.2 Performance on YCSB-T.** We further compared the performance of CRDB and CRDB-RLS on the default YCSB-T (§4.2) using the calibrated configurations. The results are shown in Figure 10. Overall, CRDB-RLS achieved  $0.87\times$  to  $0.91\times$  throughput compared to CRDB with various concurrencies. The peak throughput of CRDB-RLS was 12% lower than CRDB, and the 99th tail latency of CRDB-RLS was  $1.2\times$  to  $1.6\times$  higher than CRDB.

CRDB-RLS incurred a more server performance drop on the skewed YCSB-T workloads since, compared to CRDB, CRDB-RLS may expand the contention footprint by involving more ongoing timestamps (i.e., line 16, Algorithm2). However, the overhead is still marginal, and the performance degradation is smaller than  $\sim 15\%$ .

## 5.3 Takeaways.

CRDB-RLS achieved similar performance as CRDB while providing strong consistency guarantees on real-time orders. This is because CRDB-RLS can efficiently capture ongoing transactions inside a region due to the fast networks. Recording ongoing transactions (known as transaction tables) is a common approach for processing transactions on a single machine or in a smaller cluster due to its simplicity and easy extension. We reused such an approach but with a more lightweight tracking method for ordering IRTs.

## 6 RELATED WORKS

Transaction processing represents a well-explored area of research, with a plethora of influential works. We will provide an overview of related works in this section.

## 6.1 Proximal Consistency Models

Figure 12 compares RLS to its proximal consistency models. We describe three of them in detail. All of them are serializable.

**Regular Sequential Serializability (RSS)** complements RLS, as they focus on different aspects of distributed databases. RSS is primarily tailored for read-only transactions, permitting two read-only transactions to observe partial results of a committed read-write transaction in arbitrary orders (see our example in Figure 11b). This behavior essentially violates the real-time ordering among read-only transactions. On the contrary, RLS focuses on data locality within multi-region deployments, allowing a database system to relax the real-time ordering among transactions that access non-interleaved regions (refer to Figure 11a). As illustrated in Figure 11b, RSS allows transaction  $T_2$  to read the writes made by a concurrent transaction  $T_1$ , while  $T_3$  following  $T_2$  reads a version preceding  $T_1$ . Consequently, the real-time order between  $T_2$  and  $T_3$  is disrupted: the real-time order between  $T_2$  and  $T_3$  is  $T_2 \rightarrow T_3$ ; the serializable order enforced by RSS is  $T_3 \rightarrow T_1 \rightarrow T_2$ , which implies  $T_3 \rightarrow T_2$ . This execution is not allowed by RLS since RLS ensures strict serializability (i.e., real-time order) for IRTs within the same region.

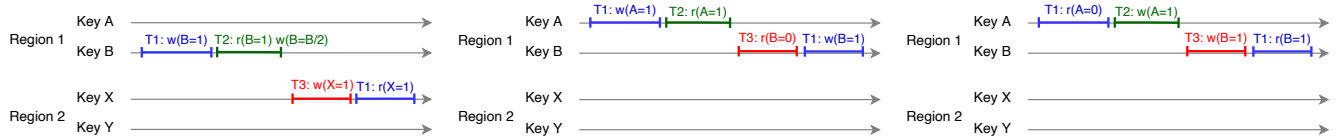
**Single-key Linearizability (SKL)** was initially proposed by CockroachDB [40]. SKL stands as a strictly weaker variant of RLS. Like RLS, SKL guarantees serializability and “no stale-reads”. However, unlike RLS, SKL does not preserve real-time orders between non-conflicting transactions. For instance, the execution illustrated in Figure 11c is permissible by SKL since there are “no stale reads” for each accessed key. However, it violates the real-time ordering between  $T_2$  and  $T_3$  in the same region, as the serial order is  $T_3 \rightarrow T_1 \rightarrow T_2$ , a violation not allowed by RLS. While weaker SKL may suffice for certain applications, other applications might necessitate stronger guarantees. Moreover, preventing consistency anomalies can significantly streamline application development.

**Process-ordered Serializability (PoS)** complements RLS. In a real deployment scenario, RLS can be stronger than PoS by associating each client with its nearby region. Specifically, PoS tracks the causal relations of each client and ensures the system preserves the ordering within each client’s requests. If each client is associated with a region (e.g., sending requests to nodes within its region), RLS can prove to be strictly stronger than PoS since it guarantees real-time ordering for each client.

## 6.2 Transaction Priority

Compared to strict serializability (SS), one of the pivotal innovations of RLS is scheduling IRTs ahead of CRTs until the CRTs’ order is established. Therefore, RLS operates under the assumption that IRTs have a higher priority than CRTs until CRTs have been ordered, after which both IRTs and CRTs are given equal priority for execution.

In this context, multi-region database developers can leverage existing transaction priority protocols to transition from SS to RLS, enhancing performance. For instance, Polaris [48] represents a transaction priority protocol rooted in a variant of OCC. Polaris embeds priority-related conflict detection within each record and permits priority preemption during runtime. Furthermore, Polaris avoids global operations, mitigating substantial overhead in a multi-region deployment. Consequently, achieving RLS should involve drawing inspiration from OCC-like concurrency control protocols.



(a) Allowed by RLS but disallowed by SS, RSS. (b) Allowed by RSS but disallowed by RLS. (c) Allowed by CRDB but disallowed by RLS

Figure 11: Comparison of RLS with proximal levels of consistency models.

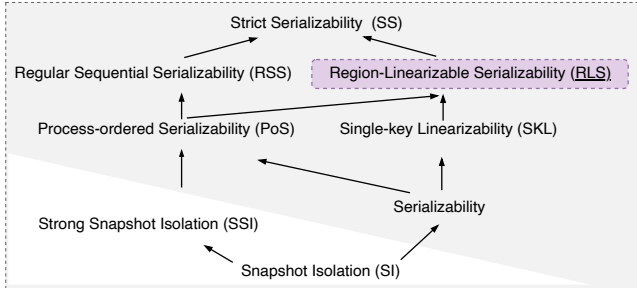


Figure 12: This diagram shows how RLS compared to its proximal consistency models: SS [35], RSS [21], SKL [40], PoS [16, 29], SSI [19], serializability [44], and SI [5]. We highlight all those serializable consistency models in grey.

### 6.3 Mixed Consistency Models

Several prior works [20, 23, 24, 27, 28, 31, 45, 47] have been on manipulating weakly and strongly consistent transactions within a single database. For instance, MixT [31] advocates that consistency is a property of information. It proposes a new embedded language enabling users to configure the consistency guarantees for each operation manually. Similarly, Red-Blue consistency [27] allows strongly and causally consistent operations to co-exist in a single system in the transaction’s granularity and depends on application semantics to make “consistency choices”. AutoGR [45] automatically analyzes and identifies the minimal set of the required consistency guarantees based on applications using the Z3 theorem prover. However, it relies on the application codes as inputs and only provides serializable guarantees without real-time order.

In contrast, RLS is geared towards multi-region deployments, directly integrating network semantics into the consistency model. Consequently, RLS does not depend on prior application knowledge to manually set distinct consistency guarantees for different transactions or operations. Moreover, RLS provides serializability (i.e., the strongest isolation levels) for all transactions while tailoring real-time properties to achieve heightened performance.

## 7 FUTURE WORKS

The multi-region transactions set three clear objectives: high throughput, low client-perceived latency (especially for IRTs), and as strong as possible consistency. Earning all three requires careful designs to strike a balance between the three individual objectives. Our exploration of RLS opens up a new design space for achieving such goals. We draw some lessons worth further research:

- **Multi-Layered Consistency Model.** Modern network exhibits a multi-layered structure [46, 52]. For instance, Cloud providers (e.g., Huawei [12]) are diligent in using CXL- and RDMA-based networks inside a data center, a dedicated network between data centers inside a region, and public networks across the regions. In this work, we have explored the consistency model for multi-region deployment by treating IRTs and CRTs differently. However, a more fine-grained design may still be desirable to fit consistency guarantees into the stack of the network tightly. For instance, in-network ordering technologies [11, 26, 32, 37], which leverage the properties of the network for order, are proposed for in-data center deployment. How to combine such technologies with geo-distributed transactions is still an open problem. A multi-layered consistency model may work as the glue to bridge the design of in-network transaction processing technologies with geo-distributed transaction processing technologies.
- **Data Locality and Partial Replication.** Even though RLS provides practical mitigation for achieving high throughput and low latency for both CRTs and IRTs, the cost of remote reads cannot be fundamentally removed. Therefore, an efficient data partition and replication policy are still critical for real-world usage: a good partition policy [3, 4, 15, 36, 49] can vastly reduce the ratio of remote reads, and a cautious replication policy [9, 39, 40] can balance the overhead of data synchronization and remote access.

## 8 CONCLUSION

This work uncovered fundamental bottlenecks in strictly serializable concurrency control algorithms used in multi-region deployments. As the consistency model inherently enforces these bottlenecks, many proposals turn to lift the restrictions by adopting weaker consistency. However, we found that all existing consistency models are inadequately designed for multi-region deployments: they are either overly stringent or have room for improvement without incurring significant performance penalties.

In response, we propose Region-Linearizable Serializability (RLS), the first consistency model meticulously tailored for multi-region deployment. Following the RLS methodology, we design, implement, and evaluate two practical system variations based on open-sourced codebases: Spanner-RLS and CRDB-RLS. The code of our stereotypes is available at [https://github.com/vldb24p771/spanner\\_rls](https://github.com/vldb24p771/spanner_rls) and [https://github.com/vldb24p771/crdb\\_rls](https://github.com/vldb24p771/crdb_rls), respectively.

Our evaluation results demonstrate that RLS can significantly enhance the performance of Spanner (i.e., from 1.16× to 89.01× higher throughput) and further strengthen the consistency guarantees of CRDB without significant performance drop (i.e., < 15%).



## REFERENCES

- [1] [n.d.]. General Data Protection Regulation (GDPR) – Official Legal Text. <https://gdpr-info.eu/>. (Accessed on 10/01/2021).
- [2] [n.d.]. Github: UWSysLab/tapir. <https://github.com/UWSysLab/tapir>.
- [3] Michael Abebe, Brad Glasbergen, and Khuzaima Daudjee. 2020. DynaMast: Adaptive dynamic mastering for replicated systems. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. IEEE, 1381–1392.
- [4] Michael Abebe, Brad Glasbergen, and Khuzaima Daudjee. 2020. MorphoSys: automatic physical design metamorphosis for distributed database systems. *Proceedings of the VLDB Endowment* 13, 13 (2020), 3573–3587.
- [5] Todd Anderson, Yuri Breitbart, Henry F Korth, and Avishai Wool. 1998. Replication, consistency, and practicality: are these mutually exclusive?. In *Proceedings of the 1998 ACM SIGMOD international conference on Management of data*. 484–495.
- [6] AWS. [n.d.]. Regions, Availability Zones, and Local Zones. <https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/using-regions-availability-zones.html>.
- [7] Azure. [n.d.]. Regions and availability zones. <https://docs.microsoft.com/en-us/azure/availability-zones/az-overview>.
- [8] Microsoft Azure. [n.d.]. Azure network round-trip latency statistics. <https://docs.microsoft.com/en-us/azure/networking/azure-network-latency>.
- [9] Xusheng Chen, Haoze Song, Jianyu Jiang, Chaoyi Ruan, Cheng Li, Sen Wang, Gong Zhang, Reynold Cheng, and Heming Cui. 2021. Achieving low tail-latency and high scalability for serializable transactions in edge computing. In *Proceedings of the Sixteenth European Conference on Computer Systems*. 210–227.
- [10] Xusheng Chen, Haoze Song, Jianyu Jiang, Chaoyi Ruan, Cheng Li, Sen Wang, Gong Zhang, Reynold Cheng, and Heming Cui. 2021. Achieving low tail-latency and high scalability for serializable transactions in edge computing. In *Proceedings of the Sixteenth European Conference on Computer Systems*. 210–227.
- [11] Inho Choi, Ellis Michael, Yunfan Li, Dan RK Ports, and Jialin Li. 2023. Hydra: {Serialization-Free} Network Ordering for Strongly Consistent Distributed Applications. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI '23)*. 293–320.
- [12] Huawei Cloud. [n.d.]. HUAWEI CLOUD Regions and Service Endpoints. <https://developer.huaweicloud.com/intl/en-us/endpoint>.
- [13] James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, J. J. Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson Hsieh, Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David Mwaura, David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Szymaniak, Christopher Taylor, Ruth Wang, and Dale Woodford. 2012. Spanner: Google's Globally-distributed Database. In *Proceedings of the Tenth Symposium on Operating Systems Design and Implementation (OSDI '12)*.
- [14] THE TRANSACTION PROCESSING COUNCIL. 2014. TPC-C. <http://www.tpc.org/tpcc/>.
- [15] Carlo Curino, Evan Philip Charles Jones, Yang Zhang, and Samuel R Madden. 2010. Schism: a workload-driven approach to database replication and partitioning. (2010).
- [16] Khuzaima Daudjee and Kenneth Salem. 2004. Lazy database replication with ordering guarantees. In *Proceedings. 20th International Conference on Data Engineering*. IEEE, 424–435.
- [17] Khuzaima Daudjee and Kenneth Salem. 2006. Lazy database replication with snapshot isolation. In *Proceedings of the 32nd international conference on Very large data bases*. 715–726.
- [18] Hua Fan and Wojciech Golab. 2019. Ocean vista: gossip-based visibility control for speedy geo-distributed transactions. *Proceedings of the VLDB Endowment* 12 (2019), 1471–1484.
- [19] Alan Fekete, Dimitrios Liarokapis, Elizabeth O'Neil, Patrick O'Neil, and Dennis Shasha. 2005. Making snapshot isolation serializable. *ACM Transactions on Database Systems (TODS)* 30, 2 (2005), 492–528.
- [20] Lei Gao, Mike Dahlin, Amol Nayate, Jiandan Zheng, and Arun Iyengar. 2003. Application specific data replication for edge services. In *Proceedings of the 12th international conference on World Wide Web*. 449–460.
- [21] Jeffrey Helt, Matthew Burke, Amit Levy, and Wyatt Lloyd. 2021. Regular Sequential Serializability and Regular Sequential Consistency. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*. 163–179.
- [22] Bert Hubert. [n.d.]. tc(8), Linux manual page. <https://man7.org/linux/man-pages/man8/tc.8.html>.
- [23] Tim Kraska, Martin Hentschel, Gustavo Alonso, and Donald Kossmann. 2009. Consistency rationing in the cloud: Pay only when it matters. *Proceedings of the VLDB Endowment* 2, 1 (2009), 253–264.
- [24] Tim Kraska, Gene Pang, Michael J Franklin, Samuel Madden, and Alan Fekete. 2013. MDCC: Multi-data center consistency. In *Proceedings of the 8th ACM European Conference on Computer Systems*. 113–126.
- [25] Leslie Lamport. 2019. Time, clocks, and the ordering of events in a distributed system. In *Concurrency: the Works of Leslie Lamport*. 179–196.
- [26] Bojie Li, Gefei Zuo, Wei Bai, and Lintao Zhang. 2021. 1pize: Scalable total order communication in data center networks. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*. 78–92.
- [27] Cheng Li, Daniel Porto, Allen Clement, Johannes Gehrke, Nuno Preguica, and Rodrigo Rodrigues. 2012. Making {Geo-Replicated} systems fast as possible, consistent when necessary. In *10th USENIX Symposium on Operating Systems Design and Implementation (OSDI '12)*. 265–278.
- [28] Cheng Li, Nuno Preguica, and Rodrigo Rodrigues. 2018. Fine-grained consistency for geo-replicated systems. In *2018 USENIX Annual Technical Conference (USENIX ATC '18)*. 359–372.
- [29] Haonan Lu, Christopher Hodsdon, Khiem Ngo, Shuai Mu, and Wyatt Lloyd. 2016. The {SNOW} Theorem and {Latency-Optimal} {Read-Only} Transactions. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI '16)*. 135–150.
- [30] Haonan Lu, Shuai Mu, Siddhartha Sen, and Wyatt Lloyd. 2023. NCC: Natural Concurrency Control for Strictly Serializable Datastores by Avoiding the Timestamp-Inversion Pitfall. *arXiv preprint arXiv:2305.14270* (2023).
- [31] Mae Milano and Andrew C Myers. 2018. MixT: A language for mixing consistency in geodistributed transactions. *ACM SIGPLAN Notices* 53, 4 (2018), 226–241.
- [32] Iulian Moraru, David G. Andersen, and Michael Kaminsky. 2013. There is More Consensus in Egalitarian Parliaments. In *Proceedings of the 13th ACM Symposium on Operating Systems Principles (SOSP '91)*.
- [33] Shuai Mu, Lamont Nelson, Wyatt Lloyd, and Jinyang Li. 2016. Consolidating concurrency control and consensus for commits under conflicts. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI '16)*. 517–532.
- [34] Cuong DT Nguyen, Johann K Miller, and Daniel J Abadi. 2023. Detock: High Performance Multi-region Transactions at Scale. *Proceedings of the ACM on Management of Data* 1, 2 (2023), 1–27.
- [35] Christos H Papadimitriou. 1979. The serializability of concurrent database updates. *Journal of the ACM (JACM)* 26, 4 (1979), 631–653.
- [36] Andrew Pavlo, Carlo Curino, and Stanley Zdonik. 2012. Skew-aware automatic database partitioning in shared-nothing, parallel OLTP systems. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*. 61–72.
- [37] Ji Qi, Xusheng Chen, Yunpeng Jiang, Jianyu Jiang, Tianxiang Shen, Shixiong Zhao, Sen Wang, Gong Zhang, Li Chen, Man Ho Au, and Heming Cui. 2021. BIDL: A High-throughput, Low-latency Permissioned Blockchain Framework for Datacenter Networks. In *The 28th ACM Symposium on Operating Systems Principles*.
- [38] Kun Ren, Dennis Li, and Daniel J Abadi. 2019. SLOG: serializable, low-latency, geo-replicated transactions. *Proceedings of the VLDB Endowment* 12 (2019), 1747–1761.
- [39] Nicolas Schiper, Pierre Sutra, and Fernando Pedone. 2010. P-store: Genuine partial replication in wide area networks. In *2010 29th IEEE Symposium on Reliable Distributed Systems*. IEEE, 214–224.
- [40] Rebecca Taft, Irfan Sharif, Andrei Matei, Nathan VanBenschoten, Jordan Lewis, Tobias Gieger, Kai Niemi, Andy Woods, Anne Birzin, Raphael Poss, et al. 2020. Cockroachdb: The resilient geo-distributed sql database. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 1493–1509.
- [41] 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 *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data (Scottsdale, Arizona, USA) (SIGMOD '12)*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/2213836.2213838>
- [42] Alexander Thomson, Thaddeus Diamond, Shu-Chun Weng, Kun Ren, Philip Shao, and Daniel J. Abadi. 2014. Fast Distributed Transactions and Strongly Consistent Replication for OLTP Database Systems. In *SIGMOD '12: Proceedings of the 2012 ACM SIGMOD international conference on Management of data*.
- [43] Nathan VanBenschoten, Arul Ajmani, Marcus Gartner, Andrei Matei, Aayush Shah, Irfan Sharif, Alexander Shraer, Adam Storm, Rebecca Taft, Oliver Tan, et al. 2022. Enabling the next generation of multi-region applications with cockroachdb. In *Proceedings of the 2022 International Conference on Management of Data*. 2312–2325.
- [44] Paolo Viotti and Marko Vukolić. 2016. Consistency in non-transactional distributed storage systems. *ACM Computing Surveys (CSUR)* 49, 1 (2016), 1–34.
- [45] Jiawei Wang, Cheng Li, Kai Ma, Jingze Huo, Feng Yan, Xinyu Feng, and Yinlong Xu. 2021. AUTOGR: automated geo-replication with fast system performance and preserved application semantics. *Proceedings of the VLDB Endowment* 14, 9 (2021), 1517–1530.
- [46] Kevin C Webb, Alex C Snoeren, and Kenneth Yocum. 2011. Topology switching for data center networks. In *Workshop on Hot Topics in Management of Internet, Cloud, and Enterprise Networks and Services (Hot-ICE '11)*.
- [47] Yingyi Yang, Yi You, and Bochuan Gu. 2017. A Hierarchical Framework with Consistency Trade-off Strategies for Big Data Management. In *2017 IEEE International Conference on Computational Science and Engineering (CSE) and IEEE International Conference on Embedded and Ubiquitous Computing (EUC)*, Vol. 1. IEEE, 183–190.
- [48] Chenhao Ye, Wuh-Chwen Hwang, Keren Chen, and Xiangyao Yu. 2023. Polarix: Enabling Transaction Priority in Optimistic Concurrency Control. *Proceedings of the ACM on Management of Data* 1, 1 (2023), 1–24.

- [49] Erfan Zamanian, Carsten Binnig, and Abdallah Salama. 2015. Locality-aware partitioning in parallel database systems. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. 17–30.
- [50] Irene Zhang, Naveen Kr Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan RK Ports. 2018. Building consistent transactions with inconsistent replication. *ACM Transactions on Computer Systems (TOCS)* 35, 4 (2018), 1–37.
- [51] Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan R. K. Ports. 2018. Building Consistent Transactions with Inconsistent Replication. *ACM Trans. Comput. Syst.* 35, 4 (Dec. 2018), 1–37. <https://doi.org/10.1145/3269981>
- [52] Qian Zhang, Jingyao Li, Hongyao Zhao, Quanqing Xu, Wei Lu, Jinliang Xiao, Fusheng Han, Chuanhui Yang, and Xiaoyong Du. 2023. Efficient Distributed Transaction Processing in Heterogeneous Networks. *Proceedings of the VLDB Endowment* 16, 6 (2023), 1372–1385.