



第十一届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2020

架构革新 高效可控



北京国际会议中心 | 2020/12/21-12/23

腾讯TDSQL 分布式事务处理

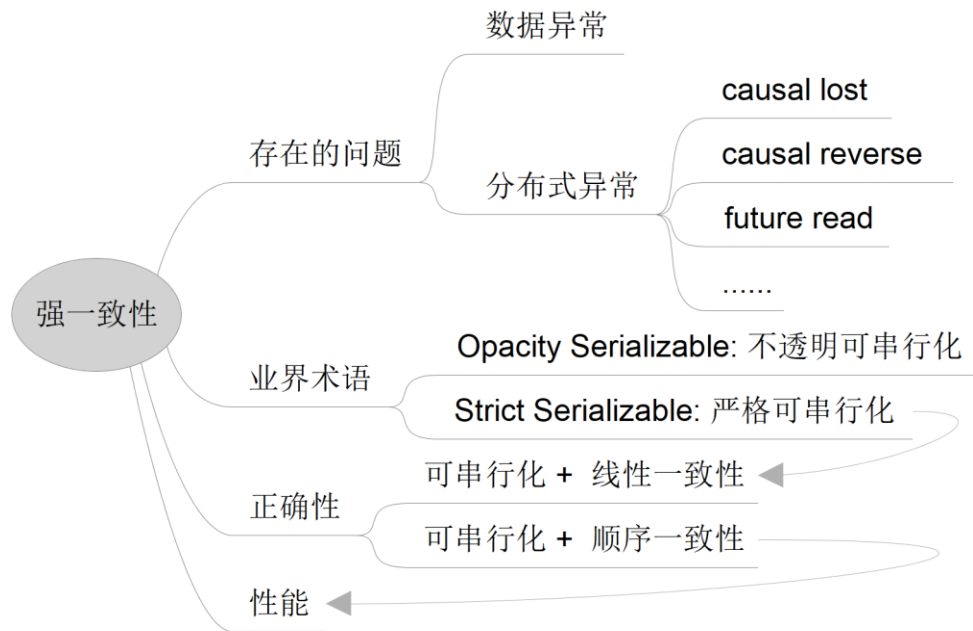
多级一致性技术（强一致性）

李海翔 @那海蓝蓝



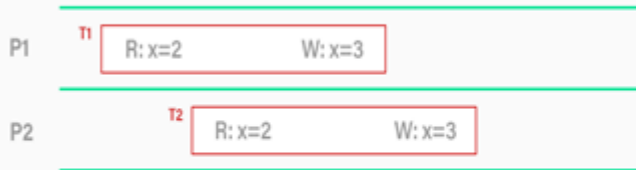
架构革新 11th
自主可控





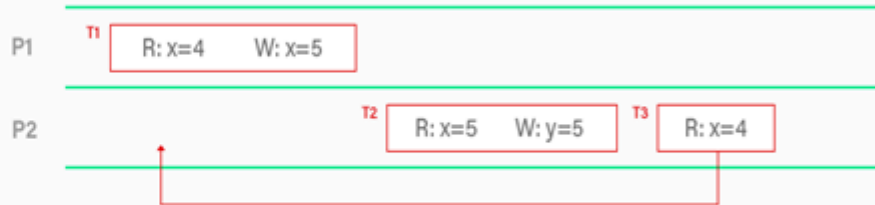


Consistency without isolation causes bugs



无事务保障：脏写数据异常

Isolation without consistency causes bugs



T3 time travels back in time. This is allowed because there is no consistency guarantee!

无一致性保障：违反因果序



分布式事务型数据库，
两种一致性缺一不可



数据准备

```
create table control ( deposit_no int not null ); insert into control values (1);
create table receipt ( receipt_no serial primary key, deposit_no int not null, payee text not null, amount money not null );
insert into receipt (deposit_no, payee, amount) values ((select deposit_no from control), 'Crosby', '100' );
insert into receipt (deposit_no, payee, amount) values ((select deposit_no from control), 'Stills', '200' );
insert into receipt (deposit_no, payee, amount) values ((select deposit_no from control), 'Nash', '300' );
```

理论分析时的并发情况:

T1	T2	T3
R1-x0		
W1-y1		
	R2-x0	
	W2-x1	
	C2	
		R3-y0
C1 ?		

事实并发情况:

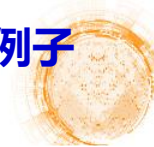
T1	T2 /T3
R1-x0	
W1-y1	
	R2-x0
	W2-x1
	C2
	R3-y0
C1 ?	



3TS

腾讯事务处理技术验证系统

@分布式一致性+事务一致性的问题— PostgreSQL 中一个神奇的例子



架构革新 ◎ 高效可控

第十一届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2020

P1	P2	
T1	T2	T3
<pre>begin; -- T1 insert into receipt (deposit_no, payee, amount) values ((select deposit_no from control), 'Young', '100'); select * from receipt; receipt_no deposit_no payee amount -----+-----+-----+----- 1 1 Crosby \$100.00 2 1 Stills \$200.00 3 1 Nash \$300.00 4 1 Young \$100.00 (4 rows)</pre>	<pre>begin; -- T2 select deposit_no from control; deposit_no ----- 1 update control set deposit_no = 2; commit;</pre>	<pre>begin; -- T3 select * from receipt where deposit_no = 1; receipt_no deposit_no payee amount -----+-----+-----+----- 1 1 Crosby \$100.00 2 1 Stills \$200.00 3 1 Nash \$300.00 (3 rows)</pre>
<pre>commit; ERROR: could not serialize access due to read/write dependencies among transactions DETAIL: Cancelled on identification as a pivot, during commit attempt. HINT: The transaction might succeed if retried.</pre>		

https://wiki.postgresql.org/wiki/SSI#Deposit_Report

2020



北京国际会议中心



2020/12/21-12/23





3TS

腾讯事务处理技术验证系统

@分布式一致性+事务一致性的问题— PostgreSQL 中一个神奇的例子

架构革新 © 高效可控

第十一届中国数据库技术大会
DATABASE TECHNOLOGY CONFERENCE CHINA 2020**可串行化:** T3-T1-T2**依赖图:** 无环**SSI:** 无环**Pg SSI:** 无环, 但有 “RW → T1 → RW”**可串行化:** T3-T1-T2 <-----> **因果序:** T2-T3**结论:**

单机系统的事务排序, 需要考虑因果序(可串行化的缺陷)

DTA: T3的lower >= 符合读已提交的数据的提交时间戳**多级一致性:** T3-T1-T2-T3, 故有环**Pg事实执行逻辑:**

T3 → T1 因为T1没有提交, T3根据MVCC读旧版本
 R1 (control) [deposit_no=1]
 R3 (receipt) → W1 (receipt) [deposit_no=1]

T1 → T2
 R1 (control) → R2 (control) W2 (control)

T2→T3(session order) 没有判断

T3 → RW → T1 → RW → T2
 发生了Stale read, 在事务型DB中, 是否合理?

按事务原理推理的执行逻辑:

T1 → T3 按real-time排定逻辑关系
 R1 (control) [deposit_no=1]
 W1 (receipt) [deposit_no=1] → R3 (receipt)

T1 → T2
 R1 (control) → R2 (control) W2 (control)

T2→T3(session order) 没有判断

无环

按事务+分布式一致性原理推理的执行逻辑:

T1 → T3 按real-time排定逻辑关系
 R1 (control) [deposit_no=1]
 W1 (receipt) [deposit_no=1] → R3 (receipt)

T1 → T2
 R1 (control) → R2 (control) W2 (control)

T2→T3(session order) 应该判断

有环

事务: 可串行化, 需要考虑因果一致性

https://wiki.postgresql.org/wiki/SSI#Deposit_ReportDTCC
2020

北京国际会议中心

2020/12/21-12/23

IT168.com

ChinaUnix

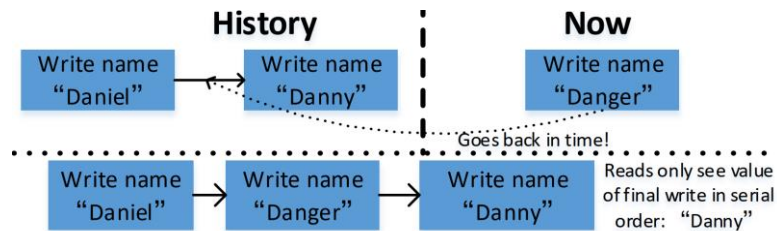
ITPUB



腾讯事务处理技术验证系统

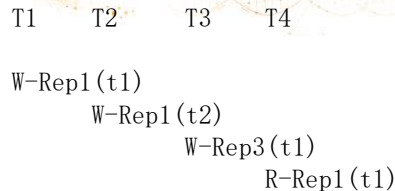
@分布式一致性+事务一致性的问题—业界观点：(异步)复制数据库

架构革新 ◎ 高效可控
第十一届中国数据库技术大会
DATABASE TECHNOLOGY CONFERENCE CHINA 2020

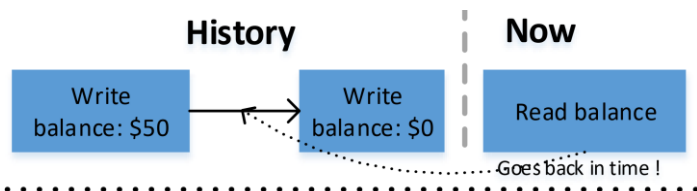


Immortal write

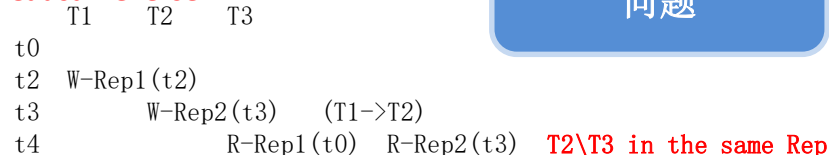
In multi-master asynchronously replicated database systems, where writes are allowed to occur at either replica, it is possible for conflicting writes to occur across the replicas.



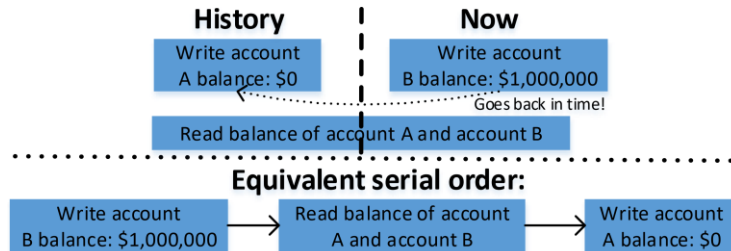
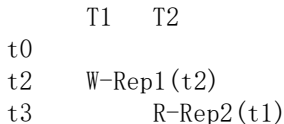
复制型架构中的问题



Causal reverse



Stale reads do not violate serializability. The system is simply time travelling the read transaction to a point in time in the equivalent serial order of transactions before the new writes to this data item occur. Therefore, asynchronously replicated database systems can allow stale reads without giving up its serializability (or even one-copy serializability) guarantee.



<http://dbmsmusings.blogspot.com/2019/08/an-explanation-of-difference-between.html>

DTCC
2020



北京国际会议中心

2020/12/21-12/23





总结:

1. 并发事务之间, 属于事务/数据的一致性问题
2. 非并发事务之间, 存在如下问题 (借用了Daniel的词, 但含义可能不同)

原因/异常	示例号	场景/现象	特点
causal lost (绑匪问题)	S1	P1 W1-x-t1 P2 R2-x-t0 不同session之间不存在显式的因果关系	因果关系是隐式无法知道的
causal reverse (Stale read (不能读到最新的数据, 读到一个旧数据))	S2	不同session之间存在显式的因果关系 P1 W1-x-t1 P2 R2-x-t1 R3-x-t0 (时间关系不确定导致因果反转)	因果关系是显式可知道的
	S3	或: P1 W1-x-t1 W1-y-t3 P2 R3-x-t0 (读已提交导致因果反转)	
	S4	P1 W1-x-t1 R1-x-t0 相同session之间存在显式的因果关系	
未来读读漂移	S5	P1 W1-x1-t1 C1 P2 R2-x1-t0 C2 (读后于自己在real-time上发生的数据, 无因果关系)	时钟漂移到未来, 无因果关系





太多的一致性，如何与事务一致性融合？

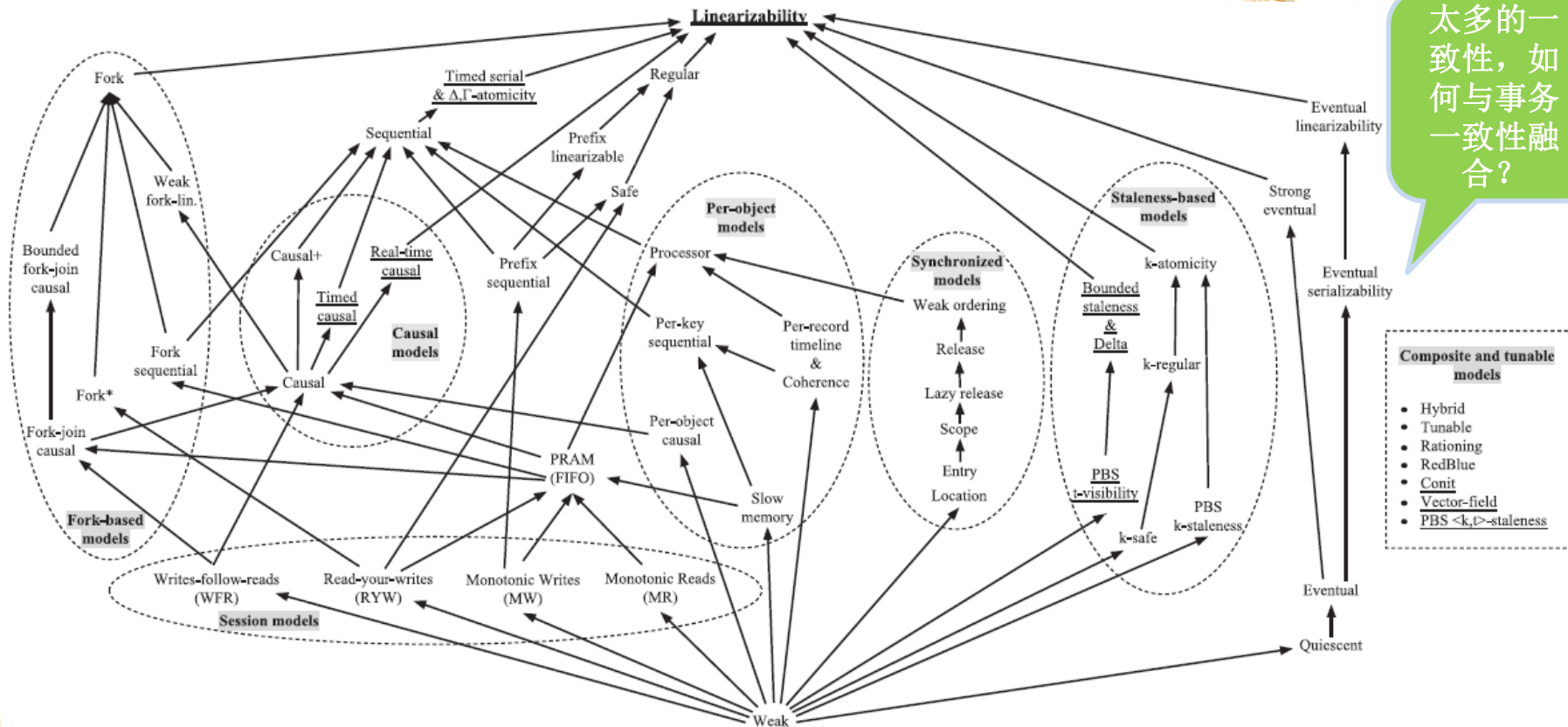




Table 1: *Ordering guarantees*. We use result visibility to define a range of coherence and consistency models. Checkmarks indicate constraints that are enforced in the model. We assume that programs, sessions and transactions are equivalent. Note that, because result visibility is an invariant, the constraints are additive and there is no need to mention executions. As a result, we can compactly represent a wide range of ordering guarantees using only result visibility to construct the columns of this table.

	Coherence (same object)				Consistency (different objects)				Session Guarantees				Isolation (transactional)						Real-time	
	Program ordering		Write Atomicity		Program ordering		Write Atomicity		Monotonic Views		Dependent Writes		Dependency cycle-free		Consistent Reads	Repeatable Reads	Non-inter. Reads	No Aborted Reads	$\alpha \prec_r \beta$	
	$\alpha(x) \prec_p \beta(x)$	$w_1(x) \prec w_2(x) \triangleleft \gamma$	$\alpha(x) \prec_p \beta(y)$	$w_1(x) \prec w_2(y) \triangleleft \gamma$	$\alpha(x) \prec_p \beta(y)$	$w_1(x) \prec w_2(y) \triangleleft \gamma$	$w \triangleleft r_1 \prec_p r_2$	$w_1 \triangleleft r \prec_p w_2 \triangleleft \gamma$	$w \triangleleft r_1$	$w_1 \triangleleft r_2$	$w_1 \triangleleft r$	$\alpha \in t_1 \triangleleft \beta \in t_2$	$w \in t_1 \triangleleft r \in t_2$	$w \in t_1 \triangleleft r \in t_2$	$w_1 \prec_p w_2 \in t_1$	$WS(t_1) \triangleleft RS(t_2)$	$\alpha \prec_r \beta$			
	$\alpha(x) \triangleleft \beta(x)$	$w_1(x) \triangleleft \gamma$	$\alpha(x) \triangleleft \beta(y)$	$w_1(x) \triangleleft \gamma$	$\alpha(x) \triangleleft \beta(y)$	$w_1(x) \triangleleft \gamma$	$w \triangleleft r_2$	$w_1 \triangleleft \gamma$	$S(t_1) \triangleleft S(t_2)$	$WS(t_1) \triangleleft RS(t_2)$	$w \triangleleft RS(t_2)$	$WS(t_1) \triangleleft RS(t_2)$	$WS(t_1) \triangleleft RS(t_2)$	$w \triangleleft RS(t_2)$	$w_2 \triangleleft r \in t_2$	$t_2 \in T_C$	$\alpha \triangleleft \beta$			
	$w \prec r \ w \prec w \ r \prec r \ r \prec w$	$\prec_p \triangleleft$	$w \prec r \ w \prec w \ r \prec r \ r \prec w$	$\prec_p \triangleleft$	$w \prec r \ w \prec w \ r \prec r \ r \prec w$	$\prec_p \triangleleft$			$WS \triangleleft WS \ RS \triangleleft WS \ WS \triangleleft RS$											
Strict Serializability	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓		
Serializability [7]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			
Snapshot Isolation [6]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			
Parallel SI [20]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			
Non-monotonic SI[4]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			
Monotonic Atomic Views [5]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓			
Repeatable Reads	✓	✓	✓	✓		✓	✓	✓	✓			✓	✓		✓	✓	✓			
Read Committed [14]	✓	✓	✓	✓		✓	✓	✓	✓			✓	✓		✓	✓	✓			
Read Uncommitted [14]	✓	✓	✓	✓		✓	✓	✓	✓			✓								
Linearizability [15]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓							✓		
Sequential Consistency [16]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓									
Total Store Order [1]	✓	✓	✓	✓	✓		✓	✓	✓	✓	✓									
Partial Store Order [1]	✓	✓	✓	✓	✓			✓	✓	✓	✓									
Weak Ordering [12]	✓	✓	✓	✓	✓				✓	✓	✓									
Processor Consistency [13]	✓	✓	✓	✓	✓		✓	✓	✓	✓	✓									
PRAM [17]	✓	✓	✓	✓	✓		✓	✓	✓	✓	✓									
Cache Consistency [13]					✓															
Read-Your-Writes [21]	✓					✓														
Monotonic Reads [21]		✓					✓			✓										
Monotonic Writes [21]		✓		✓			✓		✓											
Writes-Follow-Reads [21]			✓					✓			✓									



腾讯事务处理技术验证系统

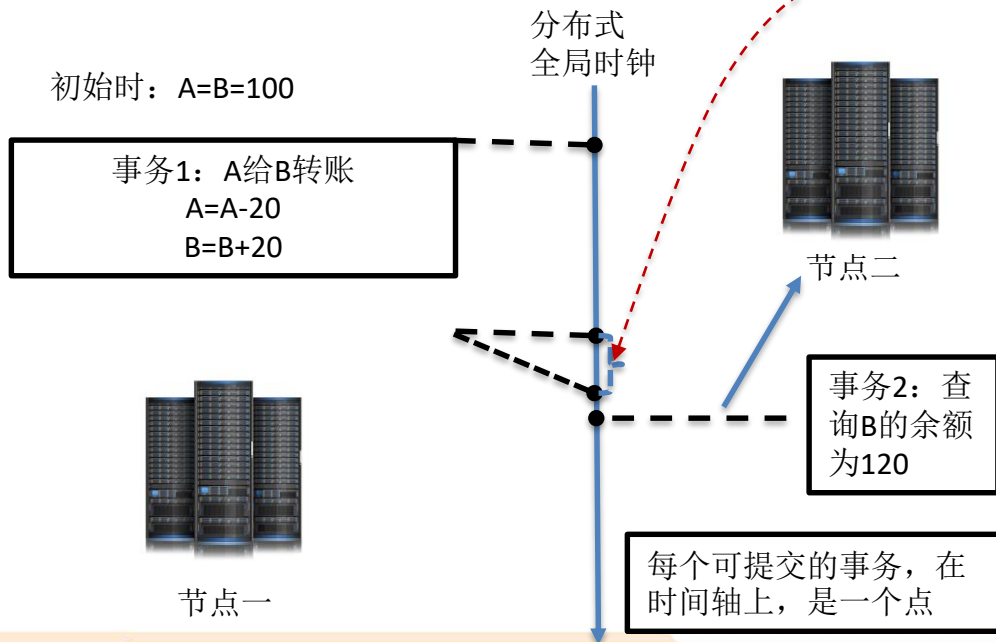
@分布式一致性+事务一致性的问题—业界典型方案



Spanner, 严格可串行化, 存在性能瓶颈

原子钟+TrueTime延迟提交 (Google Spanner)

分布式一致性+事务一致性 == 线性一致+可串行化



延迟提交时间 (TrueTime)

>
节点时钟差异的上限 (原子钟)

从 time master 同步时的网络延迟, 导致的误差 $\approx 1\text{ms}$

local clock 的漂移, 误差是时间的锯齿函数, 需校准

问题:

单分片写事务最大发生数:

$1s/2\epsilon = 125\text{个事务/每秒}$

ϵ 是延迟提交平均值, 约为 4ms



分布式一致性+事务一致性

1. **One-copy serializability** guarantees **sequential consistency** where no thread of execution can process more than one transaction.
2. **Strong session serializability** guarantees **sequential consistency** where each session corresponds to a separate thread of execution.
3. **Strong write serializability** corresponds to a **linearizability** guarantee for write transactions, but only **sequential consistency for read-only transactions**.
4. **Strong partition serializability** guarantees **linearizable consistency** within a partition, but only **sequential consistency** across partitions.
5. And **strict serializability** guarantees **linearizable consistency** at all times for all reads and writes (across non-concurrent transactions).

<http://dbmsmusings.blogspot.com/2019/08/an-explanation-of-difference-between.html>

强一致性/弱一致性

In a hybrid consistent system, **strong operations** are guaranteed to be seen **in some sequential order by all processes** (as in **sequential consistency**), while **weak operations** are designed to be fast, and they eventually become visible by all processes (much like in eventual consistency). Weak operations are only guaranteed to be ordered according to their interleaving with strong operations: if two operations belong to the same session and one of them is strong, then their **relative order** of invocation is respected and visible by all processes.

Similarly to hybrid consistency (see Section 3.8), **strong operations are ordered according to processor or sequential consistency**, whereas the ordering of weak operations is just restricted by the **relative ordering** with respect to the strong operations invoked by the same process.

Viotti P, Vukolić M. Consistency in Non-Transactional Distributed Storage Systems. ACM Comput. Surv., 2016, 49(1): 19:1-19:34.



- 线性一致性：任何2个写操作，串行执行；读、写可并行；写操作的输出值总是全序

$$\text{LINEARIZABILITY}(\mathcal{F}) \triangleq \text{SINGLEORDER} \wedge \text{REALTIME} \wedge \text{RVAL}(\mathcal{F}), \quad (7)$$

where

$$\text{SINGLEORDER} \triangleq \exists H' \subseteq \{op \in H : op.oval = \nabla\} : vis = ar \setminus (H' \times H) \quad (8)$$

and

$$\text{REALTIME} \triangleq rb \subseteq ar. \quad (9)$$

- Regular线性一致性：任何2个写操作，串行执行；读、写可并行；但获取到的值可能是最近并发写操作影响到的一个值
- Safe线性一致性：任何2个写操作，串行执行；读、写可并行；但获取到的值可能是副本中的任何一个值（某个副本可能没有收到最近的写操作的影响而改变值）

$$\text{REGULAR}(\mathcal{F}) \triangleq \text{SINGLEORDER} \wedge \text{REALTIMEWRITES} \wedge \text{RVAL}(\mathcal{F}) \quad (10)$$

$$\text{SAFE}(\mathcal{F}) \triangleq \text{SINGLEORDER} \wedge \text{REALTIMEWRITES} \wedge \text{SEQRVAL}(\mathcal{F}), \quad (11)$$

where

$$\text{REALTIMEWRITES} \triangleq rb|_{wr \rightarrow op} \subseteq ar \quad (12)$$

is a restriction of real-time ordering only for writes (preceding reads or other writes), and

$$\text{SEQRVAL}(\mathcal{F}) \triangleq \forall op \in H : \text{Concur}(op) = \emptyset \Rightarrow op.oval \in \mathcal{F}(op, \text{cxt}(A, op)), \quad (13)$$

which restricts the return value consistency only to read operations that are not concurrent with any write.

Viotti P, Vukolić M. Consistency in Non-Transactional Distributed Storage Systems. ACM Comput. Surv., 2016, 49(1): 19:1-19:34

理论上：

线性一致可以细分为多种级别

The difference between the three resides in the allowed set of return values for a read operation concurrent with a write

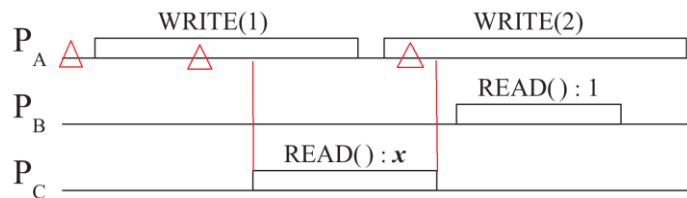


Fig. An execution exhibiting read-write concurrency (real-time flows from left to right). The register is initialized to 0. **Atomic (linearizable) semantics allow x to be 0 or 1.** Regular semantics allow x to be 0, 1, or 2. With safe semantics x may be any value



- 顺序一致性：对session内的写操作排序；不同session间的写操作，不用全局有序；任何读，都按照session排定的序读取（读操作读取的结果对于一个session写出的结果是稳定的）

$$\text{PRAM} \triangleq \text{so} \subset \text{vis}. \quad (20)$$

$$\text{SEQUENTIALCONSISTENCY}(\mathcal{F}) \triangleq \text{SINGLEORDER} \wedge \text{PRAM} \wedge \text{RVAL}(\mathcal{F}). \quad (21)$$

-so (session order) is a partial order defined as $\text{so} \triangleq \text{rb} \cap \text{ss}$.

Pipeline RAM (**PRAM** or FIFO) consistency [Lipton and Sandberg 1988] prescribes that all processes see write operations issued by a given process in the same order as they were invoked by that process

In a storage system implementing **sequential** consistency, all operations are serialized in the same order on all replicas, and the ordering of operations determined by each process is preserved.

Viotti P, Vukolić M. Consistency in Non-Transactional Distributed Storage Systems.
ACM Comput. Surv., 2016, 49(1): 19:1-19:34

In **sequential consistency**, all writes --- no matter which thread made the write, and no matter which data item was written to --- are **globally ordered**. Every single thread of execution must see the writes occurring in this order.

<https://dbmsmusings.blogspot.com/2019/07/overview-of-consistency-levels-in.html>
<http://dbmsmusings.blogspot.com/2019/08/an-explanation-of-difference-between.html>



S_{P_A} : W1 W2 W3 W5 W4 W7 W6 W8

S_{P_B} : W1 W3 W5 W7 W2 W4 W6 W8.

S_{Lin} : W1 W3 W2 W4 W5 W6 W8 W7.

An execution with processes issuing write operations on a shared object. Black spots are the chosen linearization points.

前2个，顺序一致；第三个，线性一致性

原始定义

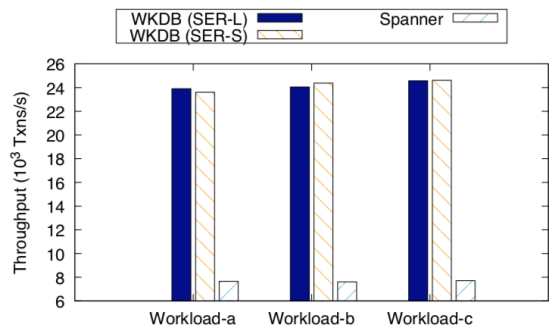
A high-speed processor may execute operations in a different order than is specified by the program. The correctness of the execution is guaranteed if the processor satisfies the following condition: **the result of an execution is the same as if the operations had been executed in the order specified by the program. A processor satisfying this condition will be called sequential.** Con-

Leslie Lamport. 1979. How to make a multiprocessor computer that correctly executes multiprocess programs. IEEE Transactions on Computers 28, 9 (1979), 690-691.

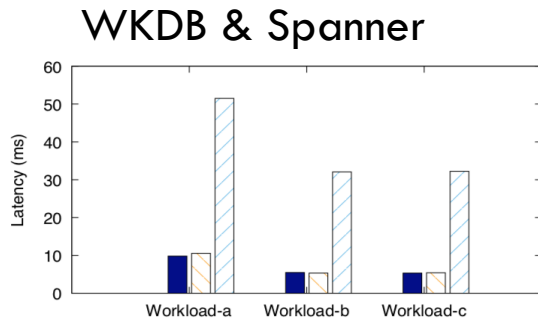
此处的“**globally ordered**”非对所有的源自不同session的写操作一同排序，而是对他们产生的结果使之在全局保持确定的顺序且符合so序



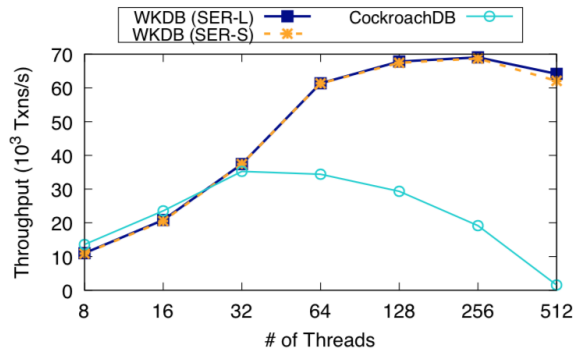
在YCSB基准测试集上，WKDB显著优于谷歌的Spanner及其开源实现CockroachDB。



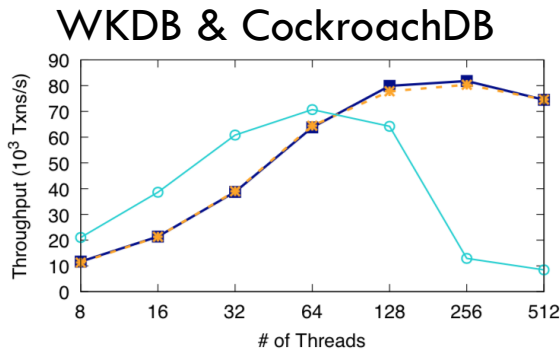
(a) Throughput



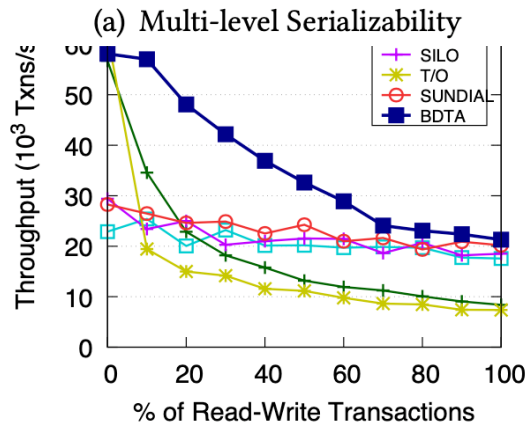
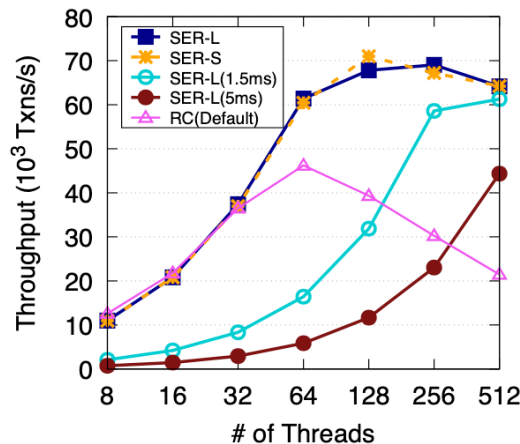
(b) Latency



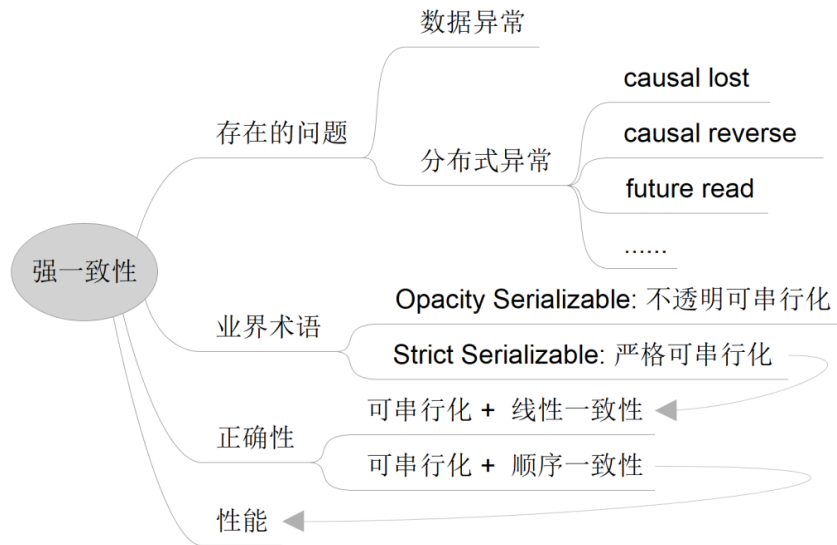
(a) YCSB-Workload-a



(b) YCSB-Workload-b



(b) Concurrency Control Evaluations



腾讯事务处理技术验证系统

- 1 第一个实现多种强一致性的工程系统
- 2 第一个高性能的强一致性工程系统
- 3 序幕拉开，真正的强一致高性能性系统将如雨后春笋.....



腾讯事务处理技术验证系统

Tencent Transaction Processing Testbed System

腾讯事务处理技术验证系统

<https://github.com/Tencent/3TS>



作为一个事务技术相关的研究框架，3TS致力于探索事务的本质问题：

1. 分布式数据库的事务处理技术，会受到哪些因素的影响（可用性？可靠性？安全性？一致性？扩展性？功能？性能？架构？新硬件？AI？.....）？
2. 分布式事务型数据库系统的评价、评测体系应该如何建立？
3. 世界上有多少种数据异常？数据异常的体系化研究方法如何建立？
4. 为什么并发访问控制算法会有很多种？各种并发访问控制算法之间，有没有的本质的关联关系？

THANKS

