

# 数据库系统事务一致性验证问题研究

## (工业软件产教联合主题交流会)

魏恒峰

hfwei@nju.edu.cn

2024 年 11 月 11 日



数据副本带来了**数据一致性**问题



(使用**形式化方法**) 解决**分布式系统**中的**数据一致性**问题



PostgreSQL



yugabyteDB



ORACLE®



TiDB



mongoDB®



Dgraph



TDSQL



fauna

# 它们正确吗?





- ▶ **数据库系统**实现正确了吗?
- ▶ **客户端程序**编写正确了吗?

► 数据库系统实现正确了吗?



数据库系统通过**事务**与**隔离级别** (Isolation Levels) 保障**数据的正确性**

声称实现了快照隔离 (Snapshot Isolation; SI)



Microsoft®  
SQL Server®



PostgreSQL



yugabyteDB



SQLite

ORACLE®



MEM  
GRAPH



TiDB



mongoDB®



Dgraph



TDSQL



fauna

但是, 仍存在违反快照隔离的数据异常 (Data Anomalies) 情况



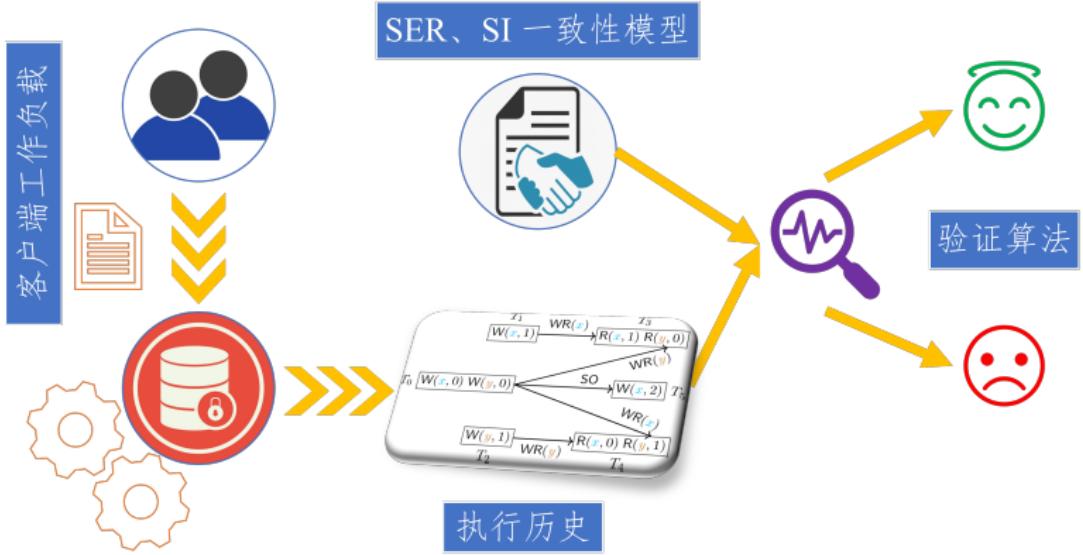
## Elle: Inferring Isolation Anomalies from Experimental Observations

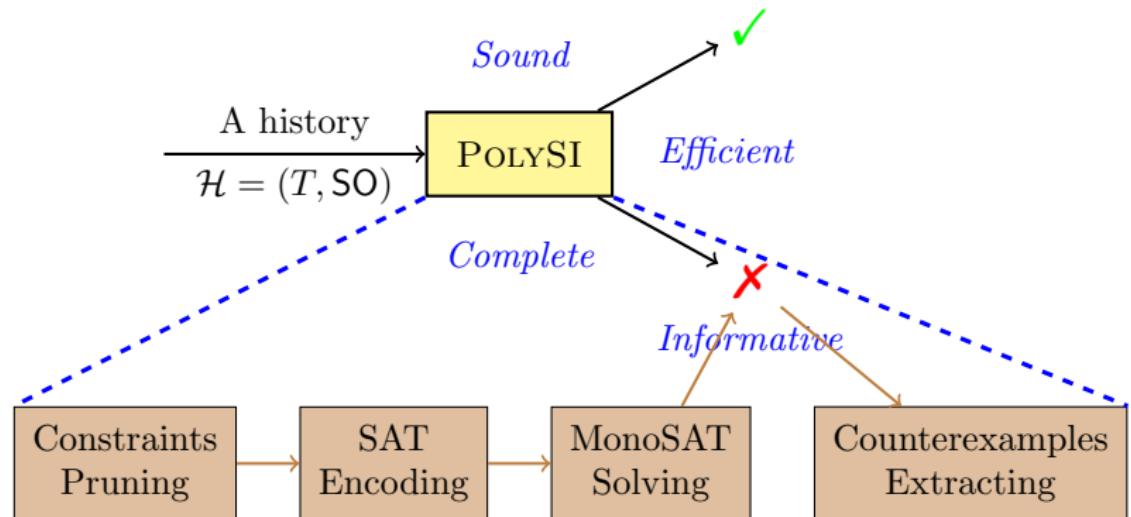
Kyle Kingsbury  
Jepsen  
aphyr@jepsen.io

Peter Alvaro  
UC Santa Cruz  
palvaro@ucsc.edu



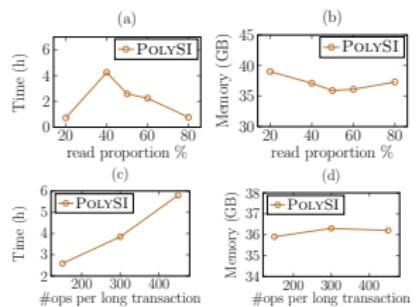
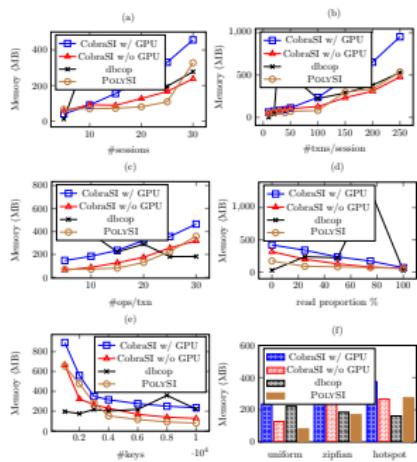
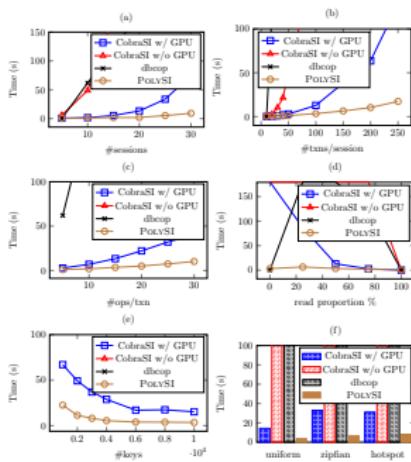
# 黑盒测试：数据库系统**执行历史**验证问题





[PolySI@VLDB'2023]

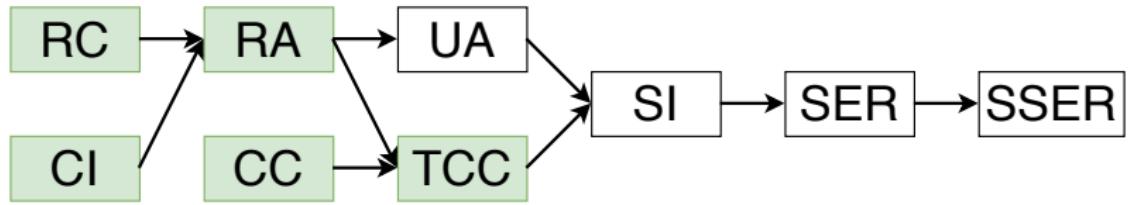




可扩展性

时间

内存



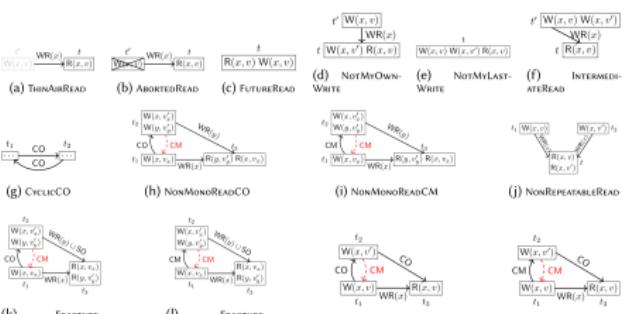
[Plume@OOPSLA'2024]

# 如何利用数据异常全面、准确地定义这些事务隔离级别

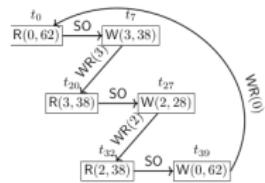
Isolation Level	Prohibited TAPs
Transactional Causal Consistency	All 14 TAPs
Read Atomicity	TAP-a to TAP-l (12 TAPs)
Read Committed	TAP-a to TAP-i (9 TAPs)
Cut Isolation	TAP-j

Table 1. The description and formalization of 14 TAPs (also visualized in Figure 3).

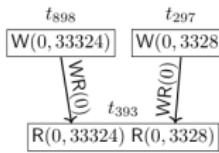
TAP	Description	TAP	Description
(a)	A transaction reads a value out of thin air. $\exists r \in R(T), \forall w \in W(T), \neg(w \xrightarrow{w} r).$	(b)	A transaction reads a value written by an aborted transaction. $\exists r \in R(T), \exists w \in W(T_0), w \xrightarrow{w} r.$
(c)	A transaction reads from a future write within the same transaction.	(d)	Transaction $t$ reads key $x$ from transaction $t' \neq t$ , but it has written to $x$ before this read. $\exists x \in K, \exists t, t' \neq t \in W_t(t), \exists r \in R_t(t), \exists w \in W_t(t).$ $\exists w' \in W_{t'}(t'), (w \xrightarrow{w(x)} r) \wedge w \xrightarrow{w} w'.$
(e)	Transaction $t$ reads key $x$ from a write $w$ itself, but $w$ is not the last write on $x$ in $t$ before this read. $\exists x \in K, \exists t, \exists w, w' \neq w \in W_t(t), \exists r \in R_t(t).$ $(w \xrightarrow{p_w} w' \xrightarrow{p_w} r \wedge w \xrightarrow{w(x)} r)$	(f)	Transaction $t$ reads key $x$ from a write $w$ in transaction $t' \neq t$ which writes $x$ more than once, but $w$ is not the last write on $x$ in $t'$ . $\exists x \in K, \exists t, \exists w, w' \neq w \in W_{t'}(t'), \exists r \in R_t(t),$ $\exists w'' \neq w' \in W_t(t), (w \xrightarrow{w(x)} r) \wedge w \xrightarrow{w} w''.$
(g)	The relation SO $\sqcap$ WR is cyclic. (SO $\sqcap$ WR) $\cap$ $I_T \neq \emptyset$ .	(h)	Transaction $t_1$ reads from $t_2$ and then reads $x \neq y$ from $t_2$ . Transaction $t_1$ also writes to $x$ but $t_1 \xrightarrow{CO} t_2$ . $\exists x, y \in K, \exists t_1, t_2 \neq t_1 \in W_{T_1}, \exists r_1 \in (R_{T_1} \cap RT_{T_2}) \setminus \{t_1, t_2\},$ $\exists w_x \in W_{t_1}(t_1), \exists w_y \in W_{t_2}(t_2), \exists r_x \in R_{t_1}(t_1), \exists r_y \in R_{t_2}(t_2),$ $(w_x \xrightarrow{w(x)} r_1 \wedge w_y \xrightarrow{w(y)} r_2 \wedge r_2 \xrightarrow{p_{R_{T_2}}} r_x \wedge t_1 \xrightarrow{CO} t_2)$
(i)	Transaction $t_1$ reads $y$ from $t_2$ and then reads $x \neq y$ from $t_1$ . Transaction $t_2$ also writes to $x$ but $t_2 \xrightarrow{CM} t_1$ . This is a general case of (h). $\exists x, y \in K, \exists t_1, t_2 \neq t_1 \in W_{T_1}, \exists r_1 \in (R_{T_1} \cap RT_{T_2}) \setminus \{t_1, t_2\},$ $\exists w_x \in W_{t_1}(t_1), \exists w_y \in W_{t_2}(t_2), \exists r_x \in R_{t_2}(t_2), \exists r_y \in R_{t_1}(t_1),$ $(w_x \xrightarrow{w(x)} r_2 \wedge w_y \xrightarrow{w(y)} r_1 \xrightarrow{p_{R_{T_1}}} r_x \wedge t_2 \xrightarrow{CM} t_1)$	(j)	A transaction reads from a key after other transactions more than once, but with different values. $\exists x \in K, \exists v, v' \in V, \exists t \in W_T, \exists t \neq t_1, t_2 \neq t \in W_T,$ $\exists r_1 \in R_{t_1}(t_1), \exists w_1 \in W_{t_1}(t_1),$ $\exists r_2 \in R_{t_2}(t_2), \exists w_2 \in W_{t_2}(t_2),$ $(t_1 \neq t_2, w_1 \xrightarrow{w(x)} r_1 \wedge w_2 \xrightarrow{w(x)} r_2)$
(k)	Transaction $t_1$ reads from $t_2$ and $y \neq x$ from $t_2$ . Transaction $t_2$ also writes to $x$ but $t_2 \xrightarrow{CO} t_1$ . $\exists x, y \in K, \exists t_1, t_2 \neq t_1 \in W_{T_1}, \exists r_1 \in (R_{T_1} \cap RT_{T_2}) \setminus \{t_1, t_2\},$ $\exists w_x \in W_{t_1}(t_1), \exists w_y \in W_{t_2}(t_2), \exists r_x \in R_{t_2}(t_2), \exists r_y \in R_{t_1}(t_1),$ $(w_x \xrightarrow{w(x)} r_2 \wedge w_y \xrightarrow{w(y)} r_1 \xrightarrow{p_{R_{T_1}}} r_x \wedge t_2 \xrightarrow{CO} t_1)$	(l)	Transaction $t_1$ reads from $t_2$ and $y \neq x$ from $t_2$ . Transaction $t_2$ also writes to $x$ but $t_2 \xrightarrow{CM} t_1$ . This is a general case of (i) and (k). $\exists x, y \in K, \exists t_1, t_2 \neq t_1 \in W_{T_1}, \exists r_1 \in (R_{T_1} \cap RT_{T_2}) \setminus \{t_1, t_2\},$ $\exists w_x \in W_{t_1}(t_1), \exists w_y \in W_{t_2}(t_2), \exists r_x \in R_{t_2}(t_2), \exists r_y \in R_{t_1}(t_1),$ $(w_x \xrightarrow{w(x)} r_2 \wedge w_y \xrightarrow{w(x)} r_1 \xrightarrow{p_{R_{T_1}}} r_x \wedge t_2 \xrightarrow{CM} t_1)$
(m)	Transaction $t_1$ reads $x$ from $t_1$ . There is a transaction $t_2$ that also writes to $x$ such that $t_1 \xrightarrow{CO} t_2 \xrightarrow{CO} t_1$ . $\exists x \in K, \exists t_1, t_2 \neq t_1 \in W_{T_1}, \exists r_1 \in RT_{T_2} \setminus \{t_1, t_2\},$ $(t_1 \xrightarrow{WR(x)} r_1 \wedge t_1 \xrightarrow{CO} t_2 \xrightarrow{CO} t_1)$	(n)	Transaction $t_1$ reads $x$ from $t_1$ . There is a transaction $t_2$ that also writes to $x$ such that $t_1 \xrightarrow{CO} t_2 \xrightarrow{CO} t_1$ . This is a general case of (l) and (m). $\exists x \in K, \exists t_1, t_2 \neq t_1 \in W_{T_1}, \exists r_1 \in RT_{T_2} \setminus \{t_1, t_2\},$ $(t_1 \xrightarrow{WR(x)} r_1 \wedge t_1 \xrightarrow{CO} t_2 \xrightarrow{CO} t_1)$



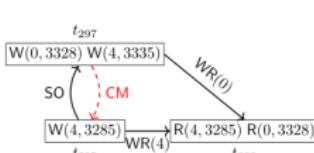




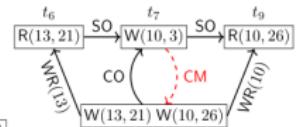
(a) TAP-g in AntidoteDB



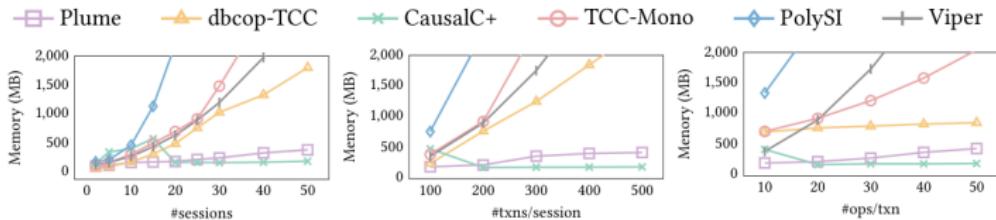
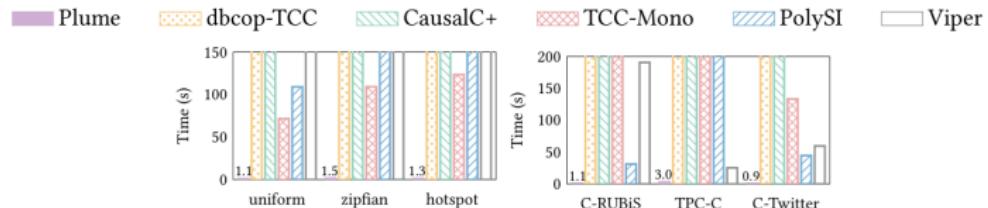
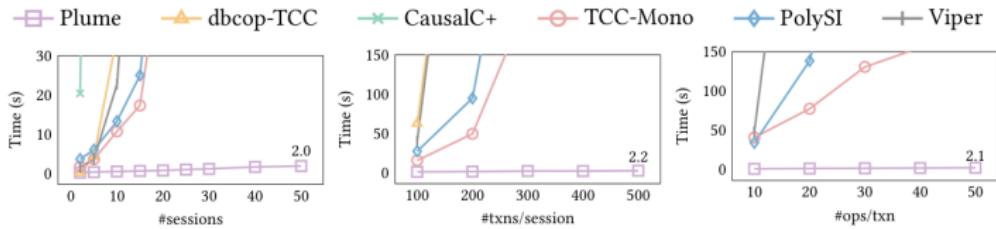
(b) TAP-j in MariaDB

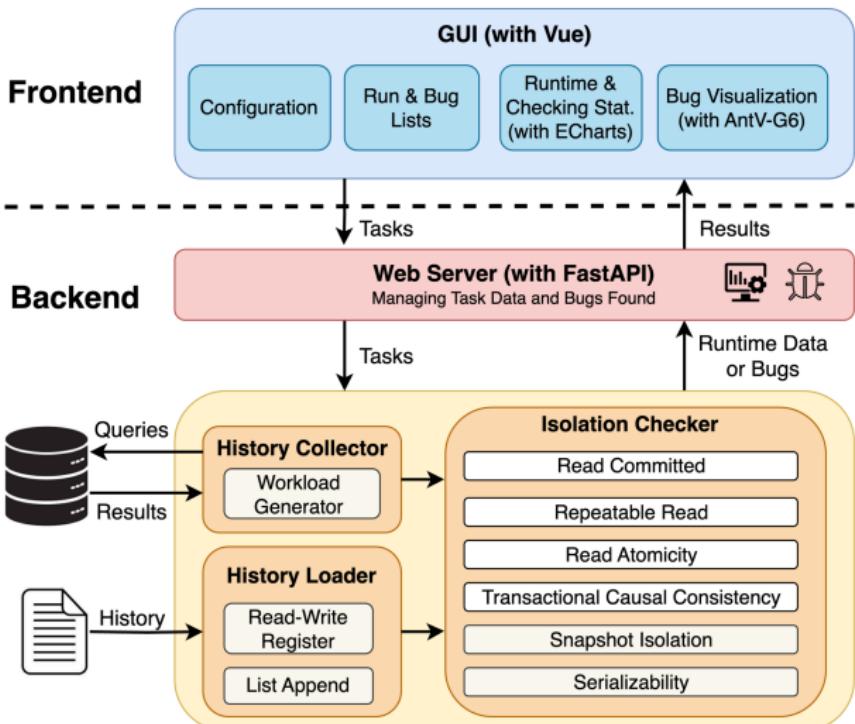


(c) TAP-k in MariaDB

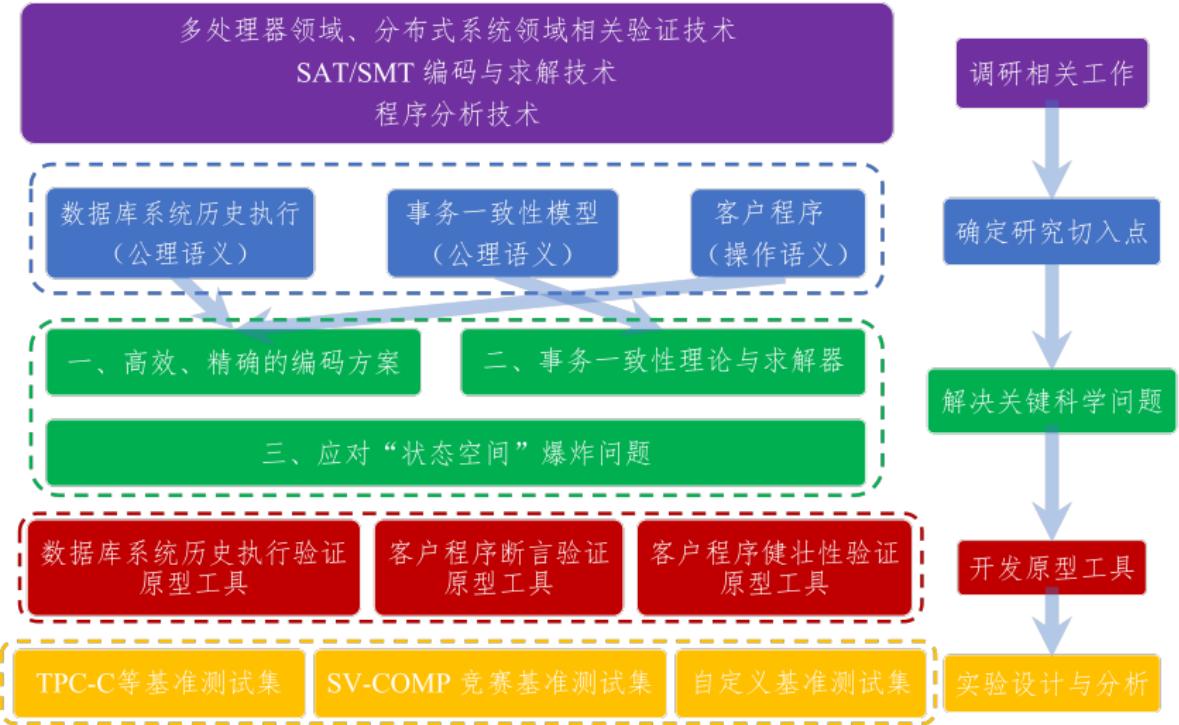


(d) TAP-m in YugabyteDB





[IsoVista@VLDB'2024 (Demo)]



# 混合 (Mixing) 隔离级别



# 谓词 (Predicates)

## Predicate Pushdown in SQL

Discover what goes  
under the hood



itsadityagupta



Hengfeng Wei (hfwei@nju.edu.cn)