# Verifying Transactional Consistency of MongoDB
### (submitted to VLDB'2022)

Hengfeng Wei

hfwei@nju.edu.cn

December 26, 2021

MongoDB 的三种经典部署架构

MongoDB 事务的三阶段发展过程

*A Fundamental Question:*

*What transactional consistency guarantee do MongoDB transactions in each deployment provide?*

# 挑战一：MongoDB 官方规约不清楚, SI 有多种变体



Figure 3: Back-to-Back Transactions with and without Speculative Snapshot Isolation

挑战二：MongoDB 缺少精简的事务协议描述, 更没有严格证明

挑战三: SI 检测问题是 NP-complete 问题, 复杂度高

THEOREM 3.2. *For any criterion $C \in \{$PREFIX CONSISTENCY, SNAPSHOT ISOLATION, SERIALIZABILITY$\}$ the problem of checking whether a given history satisfies $C$ is NP-complete.*

贡献一: 使用 (VIS, AR) 框架, 为多种 SI 变体提供形式化规约

贡献二: 为 MongoDB 事务一致性协议提供精简的伪代码描述

# 贡献三: 证明 WIREDTIGER、REPLICASET、SHARDEDCLUSTER 事务协议分别满足 STRONGSI、REALTIMESI、SESSIONSI 变体

贡献四: 设计并评估了多项式时间 SI 变体白盒检测算法

1. 事务 $T : (E, \mathsf{po})$
   - ▶ $\mathsf{po}$ : Program Order
   - ▶ $start(T)$ : 事务开始时间
   - ▶ $commit(T)$ : 事务提交时间

2. 历史 $\mathcal{H} : (\mathbb{T}, \mathsf{SO})$
   - ▶ $\mathbb{T}$ : 已提交事务集合
   - ▶ $\mathsf{SO}$ : Session Order

3. 执行 $\mathcal{A} : (\mathcal{H}, \mathsf{VIS}, \mathsf{AR})$
   - ▶ $\mathsf{VIS}$ : 可见性 (Visibility) 偏序关系
   - ▶ $\mathsf{AR}$ : 仲裁 (Arbitration) 全序关系
   - ▶ $\mathsf{VIS} \subseteq \mathsf{AR}$

一个事务一致性模型可定义为一组一致性公理的集合 $\Phi$。

历史 $\mathcal{H}$ 满足事务一致性模型 $\Phi$, 如果存在 VIS 与 AR 使得

$$\exists \text{VIS}, \text{AR}.\ (\mathcal{H}, \text{VIS}, \text{AR}) \models \Phi。$$

$$\forall (E, \mathsf{po}) \in \mathcal{H}. \; \forall e \in \mathsf{Event}. \; \forall key, val. \; \big(\mathsf{op}(e) = \mathsf{read}(key, val) \wedge \{f \mid (\mathsf{op}(f) = \_(key, \_) \wedge f \xrightarrow{\mathsf{po}} e\} \neq \emptyset\big)$$
$$\implies \mathsf{op}(\max_{\mathsf{po}}\{f \mid \mathsf{op}(f) = \_(key, \_) \wedge f \xrightarrow{\mathsf{po}} e\}) = \_(key, val) \hspace{2cm} (\textsc{Int})$$

$$\forall T \in \mathcal{H}. \; \forall key, val. \; T \vdash \mathsf{read}(key, val) \implies \max_{\textsc{ar}}(\textsc{vis}^{-1}(T) \cap \mathsf{WriteTx}_{key}) \vdash \mathsf{write}(key, val) \hspace{1cm} (\textsc{Ext})$$

| | | |
|---|---|---|
| $\textsc{so} \subseteq \textsc{vis}$ (Session) | $\textsc{ar} ; \textsc{vis} \subseteq \textsc{vis}$ | (Prefix) |
| $\textsc{rb} \subseteq \textsc{vis}$ (ReturnBefore) | $\textsc{cb} \subseteq \textsc{ar}$ | (CommitBefore) |
| $\textsc{vis} \subseteq \textsc{rb}$ (RealtimeSnapshot) | $\forall S, T \in \mathcal{H}. \; S \bowtie T \implies (S \xrightarrow{\textsc{vis}} T \vee T \xrightarrow{\textsc{vis}} S)$ | (NoConflict) |

$$SI = \text{INT} \land \text{EXT} \land \text{PREFIX} \land \text{NOCONFLICT}$$

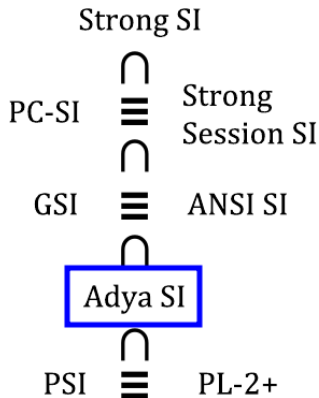$$\textsc{SessionSI} = \textsc{SI} \land \textsc{Session}$$

$$\text{SessionSI} = \text{SI} \wedge \text{Session}$$

$$\text{RealtimeSI} = \text{SI} \wedge \text{ReturnBefore} \wedge \text{CommitBefore}$$

$$\text{SessionSI} = \text{SI} \land \text{Session}$$

$$\text{RealtimeSI} = \text{SI} \land \text{ReturnBefore} \land \text{CommitBefore}$$

$$\text{GSI} = \text{SI} \land \text{RealtimeSI} \land \text{CommitBefore}$$

$$\textsc{SessionSI} = \textsc{SI} \land \textsc{Session}$$

$$\textsc{RealtimeSI} = \textsc{SI} \land \textsc{ReturnBefore} \land \textsc{CommitBefore}$$

$$\textsc{GSI} = \textsc{SI} \land \textsc{RealtimeSI} \land \textsc{CommitBefore}$$

$$\textsc{StrongSI} = \textsc{GSI} \land \textsc{ReturnBefore}$$

- ANSI-SI
- SI
- GSI
- STRONGSI
- STRONGSESSIONSI
- PSI
- WRITESI
- NMSI
- PCSI

Strong SI

$\cap$

PC-SI $\equiv$ Strong Session SI

$\cap$

GSI $\equiv$ ANSI SI

$\cap$

Adya SI

$\cap$

PSI $\equiv$ PL-2+

# Conclusion

Hengfeng Wei (hfwei@nju.edu.cn)