



2019

05

08-10

北京新云南皇冠假日酒店

数据风云 十年变迁

DTCC

第十届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2019



+

○

○

○

基于MGR的读写强一致性数据库

滴滴出行 数据库团队



SPEAKER

INTRODUCE

田佳伟 数据库开发专家

曾就职于百度、Thomson Reuters 等公司，现就职于滴滴出行，负责数据库产品的架构及研发工作；多年中间件和数据库开发经验，研究方向主要为分布式数据库系统的理论与实践。



目录

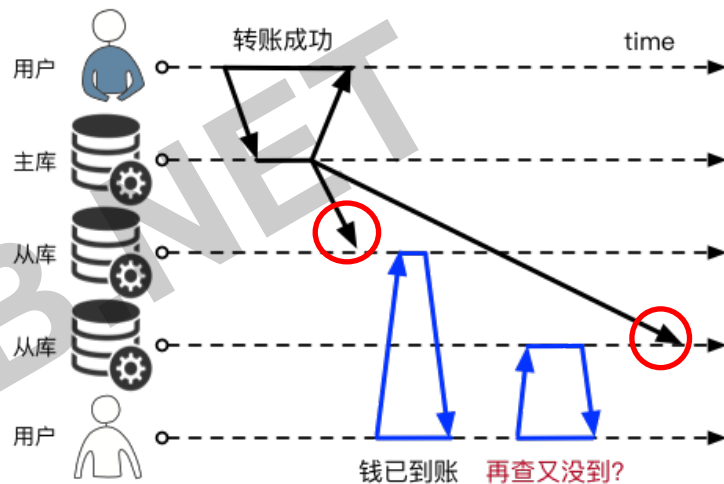
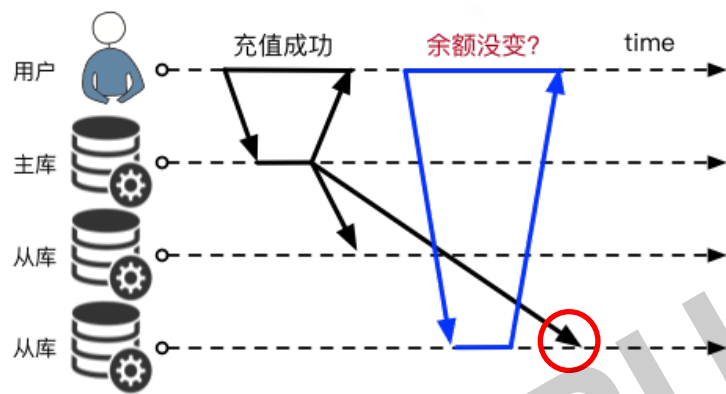
Contents

- 现实问题
- 写一致性
- 读一致性
- 规划和展望

PART ONE

What's the reality ?

现实中的交易场景

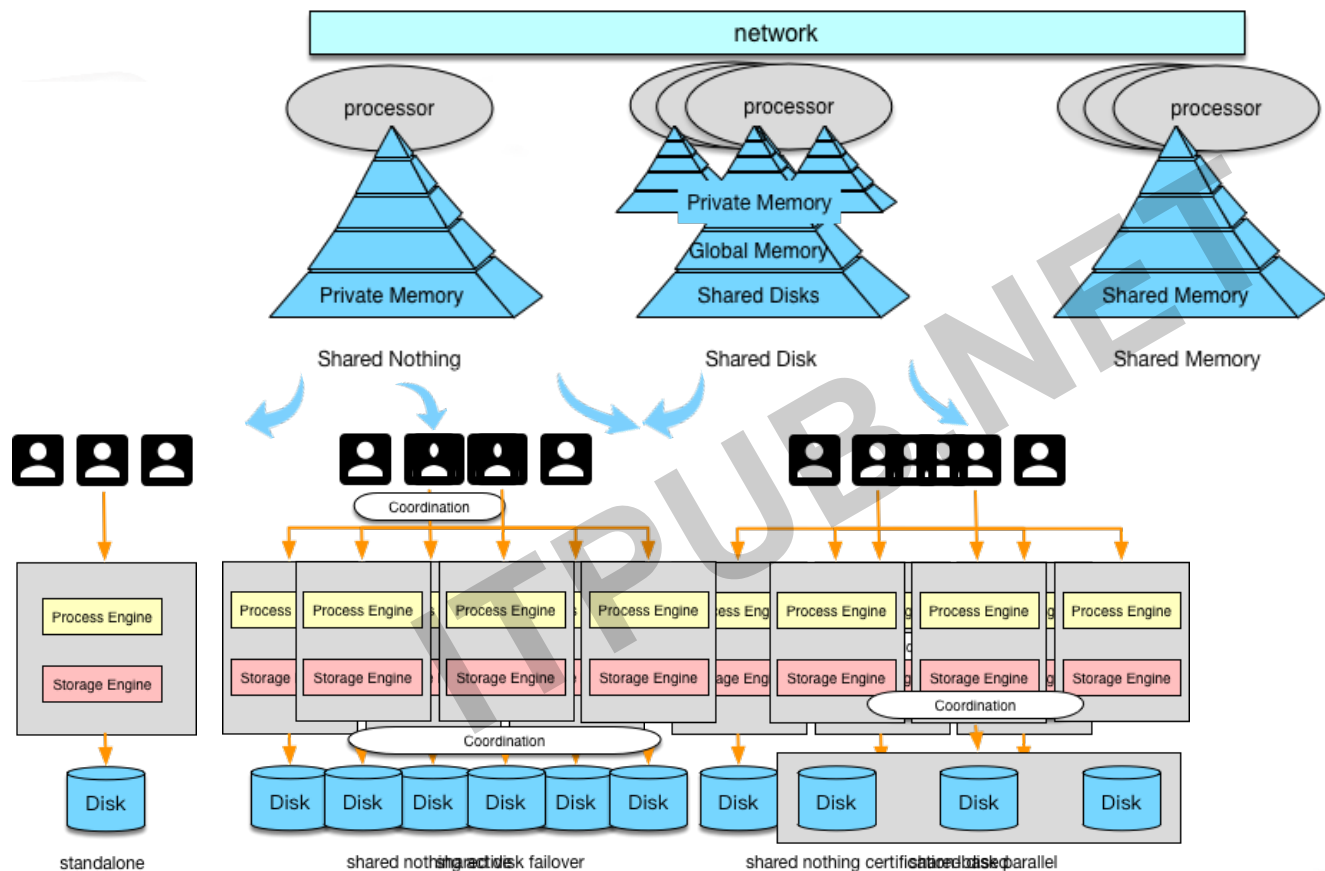


2

PART TWO

Write Consistency

ITPUB.NET



分布式数据库的核心能力

01. Scalability

弹性的扩展能力

05. Availability

保证99.99%的可用性

02. Fault-Tolerant

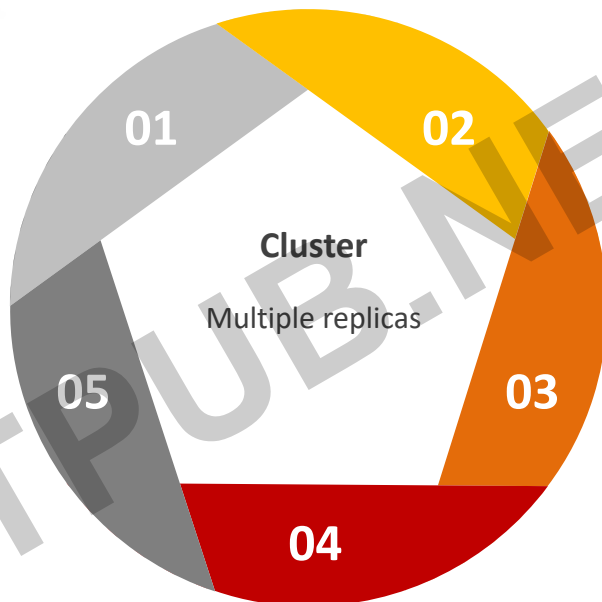
可以容忍实例、机器、存储和网络的异常

03. Write Consistency

数据一旦写入成功，
不会丢失

04. Read Consistency

保证数据读取的线性一致性

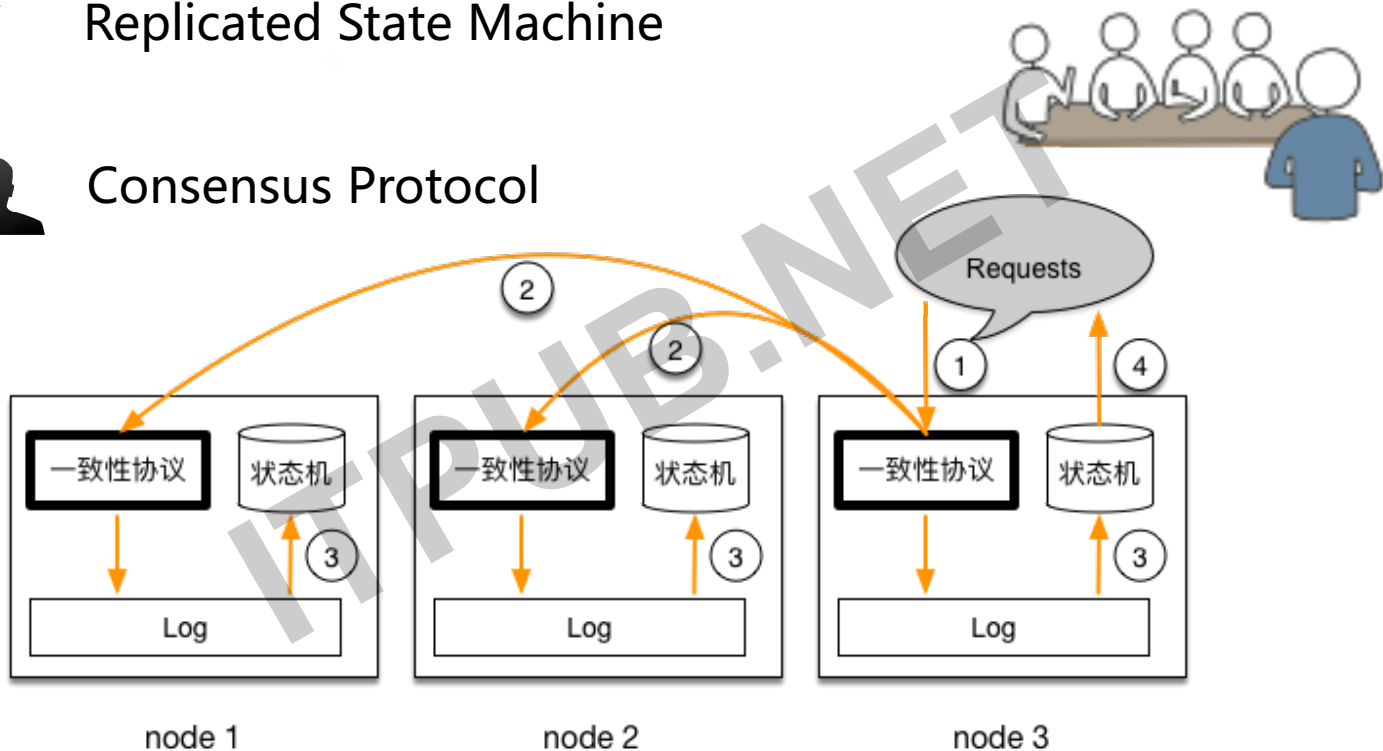




Replicated State Machine



Consensus Protocol



为什么选择MySQL Group Replication

01. Consensus Protocol

Mencious (Paxos variant)
保证线性一致性
全局同序处理

02. Stability

官方支持
社区生态完善

03. Multiple Primary

多点写入
提升可用性

04. Focus

读取的线性一致性
完善一致性语义

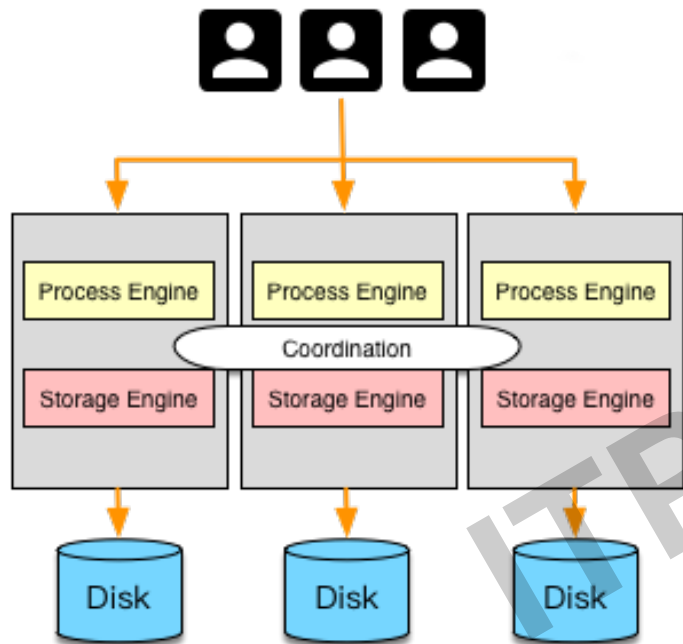
MGR
MySQL Group
Replication

3

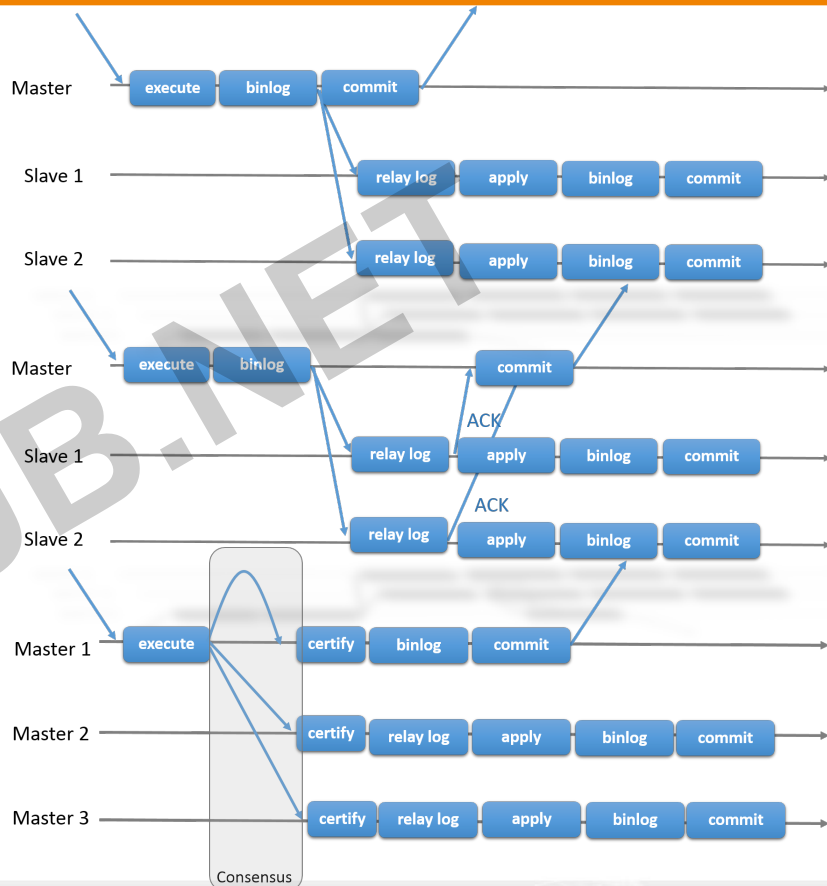
PART THREE

Read Consistency

MGR的架构和复制方式



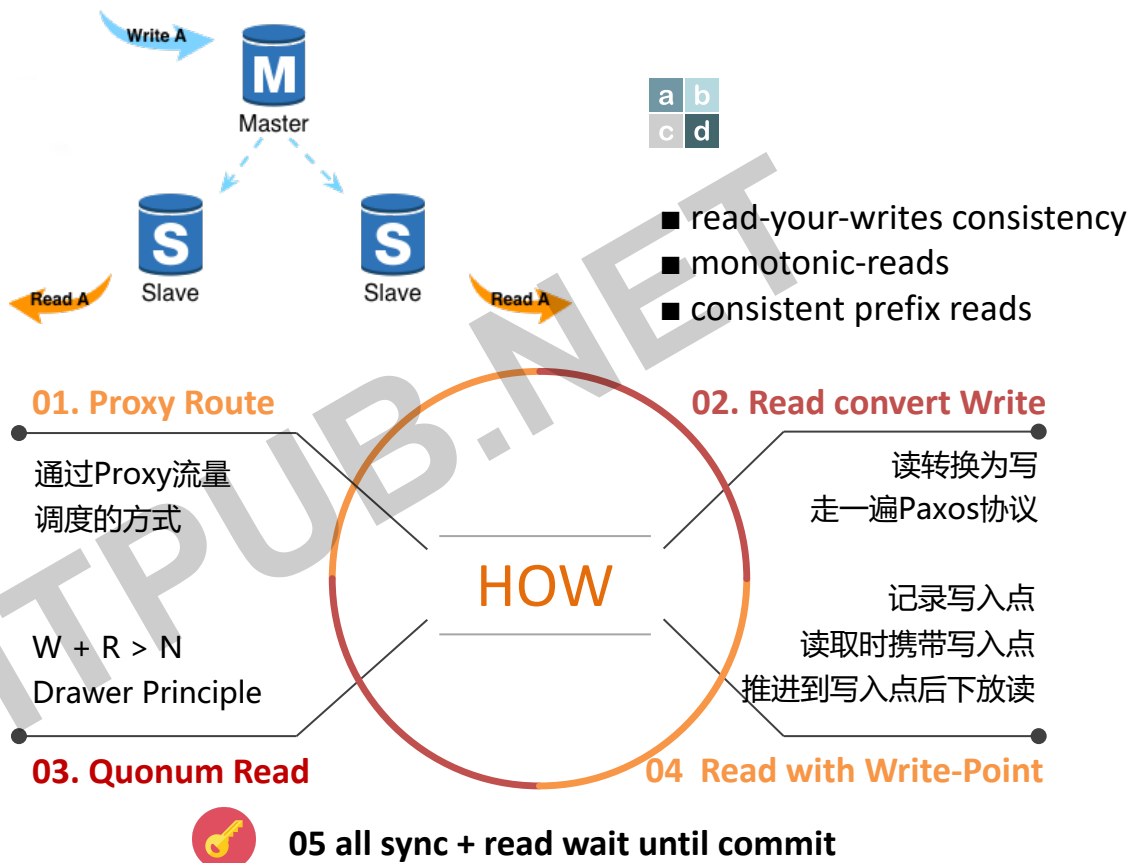
shared nothing certification-based



01

选择方案

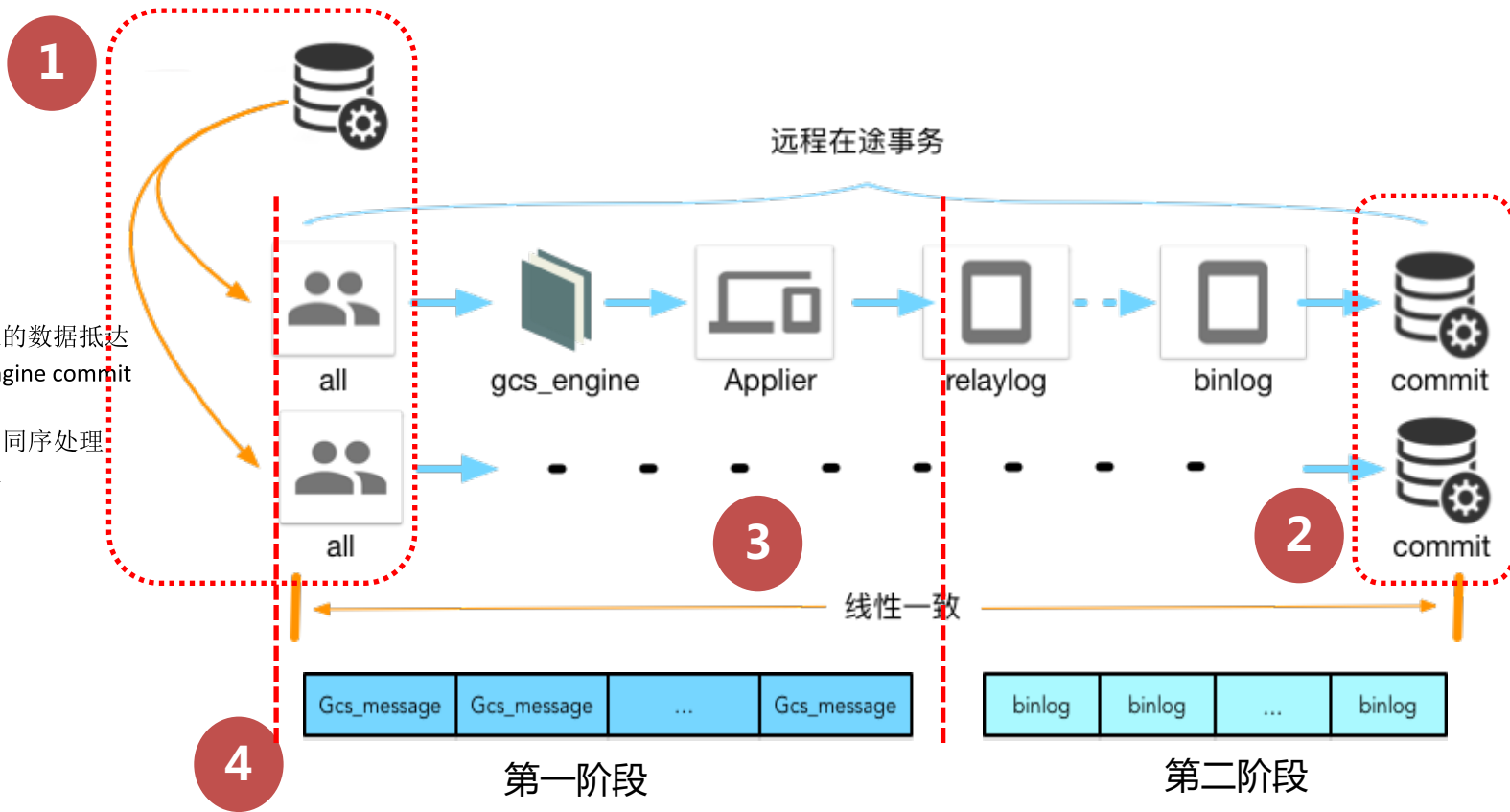
1. 读一致性的定义
2. 现有的通用一致性方式是什么
3. 我们可以用哪些方式来做到读一致性
4. 我们的最佳选择

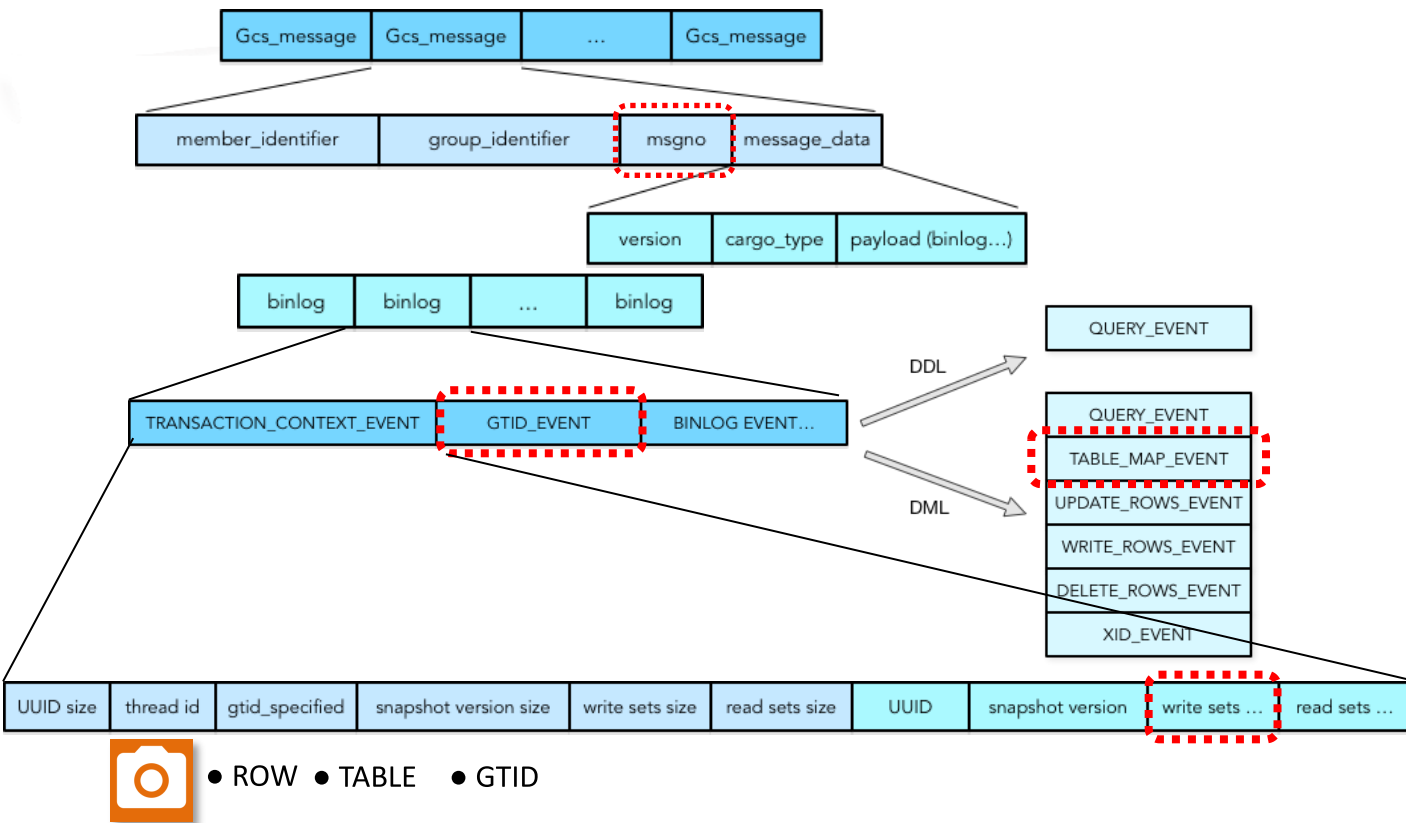


02

设计演进

1. 全同步，保证写入的数据抵达后查询
2. 读达到时，等待engine commit
3. 各个节点保证全局同序处理
4. 远程在途事务拆解

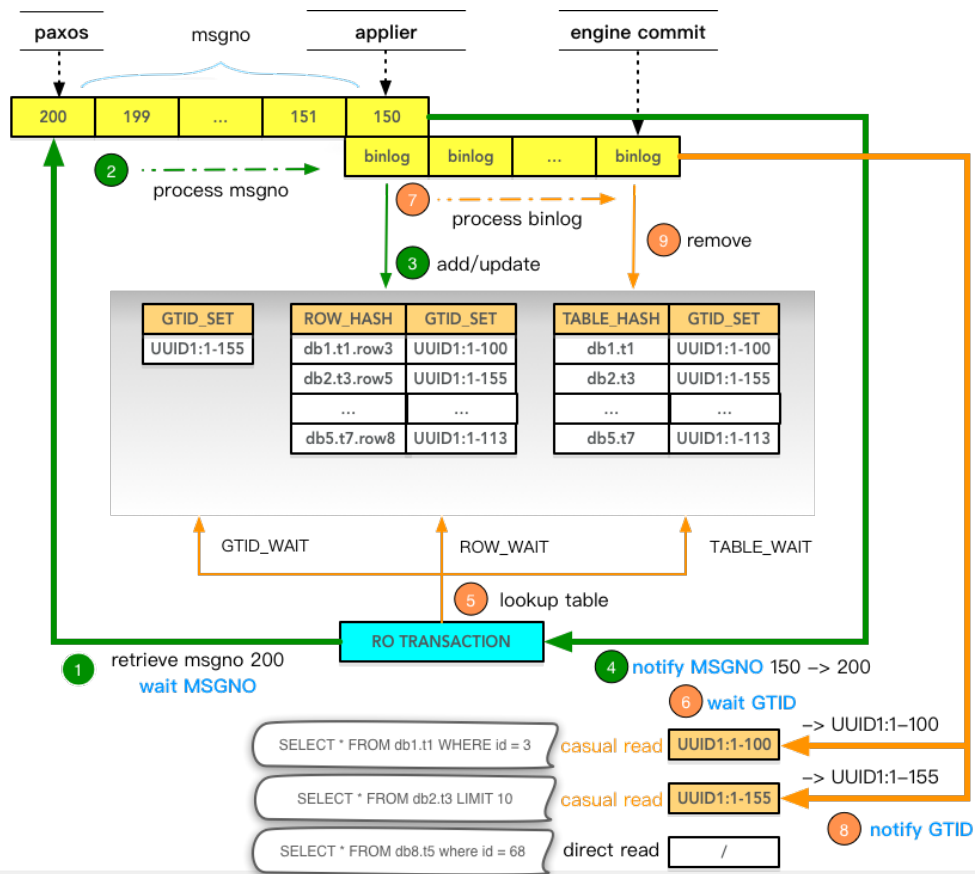




03

两阶段&Casual Read

1. 分析gcs_message的组成
2. 分析binlog的组成
3. 因果读减少等待成本



04

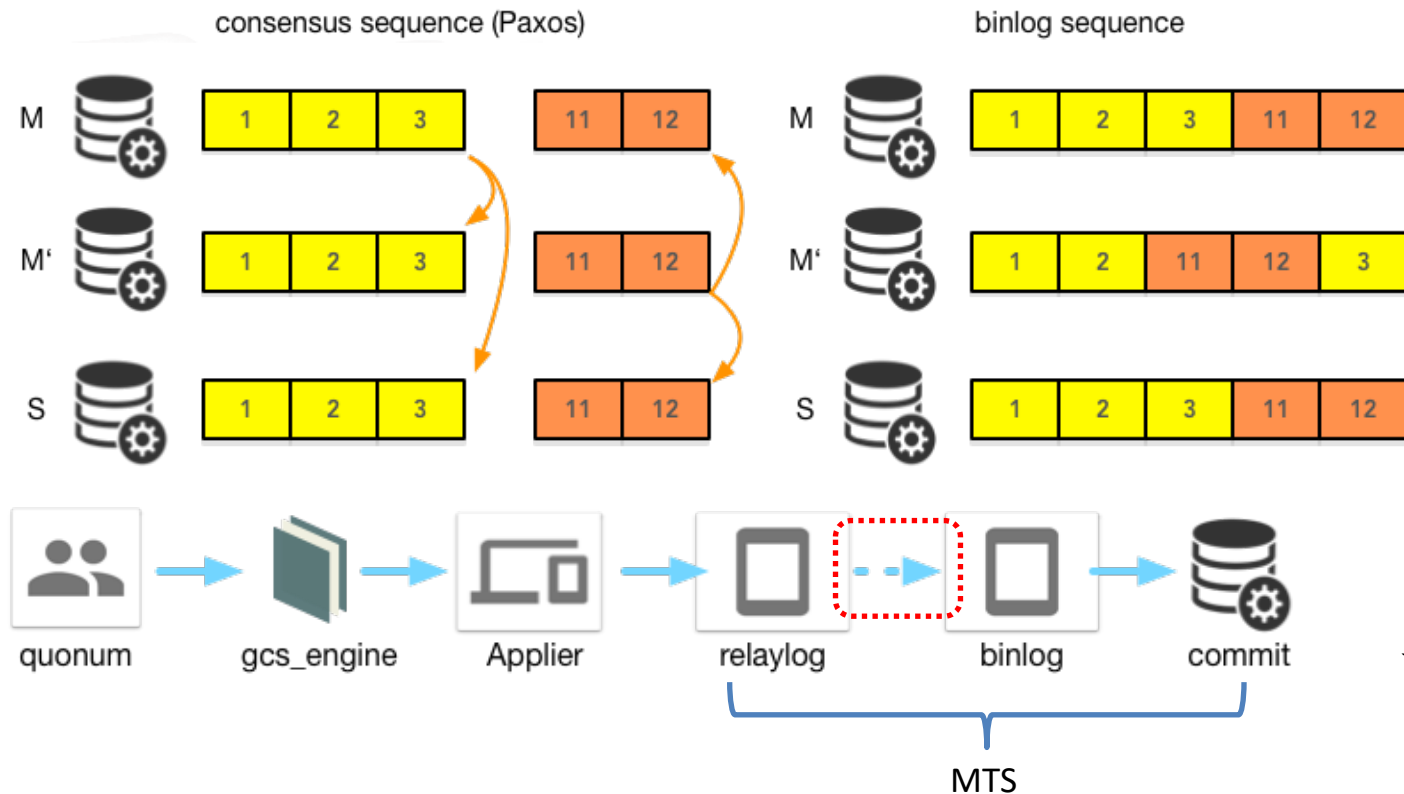
详细流程

1. 第一阶段 停-等 msgno
2. 第二阶段 停-等 gtid

Conclusion

Make simple

appear as only one copy of data

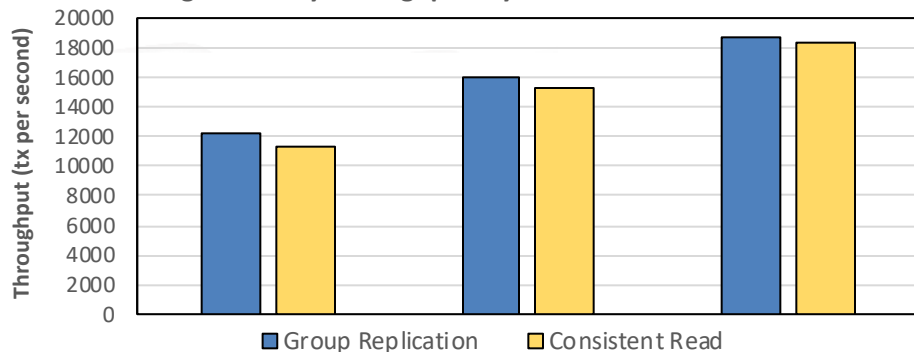


05

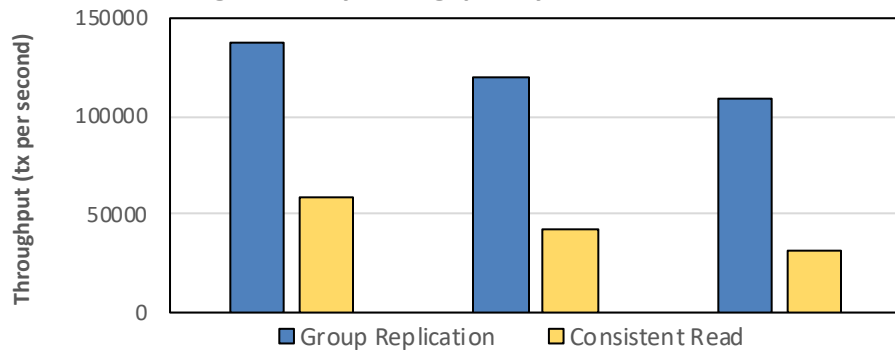
主从切换场景

1. 一致性协议保证传播线性一致
2. 应用层的日志binlog有可能不一致

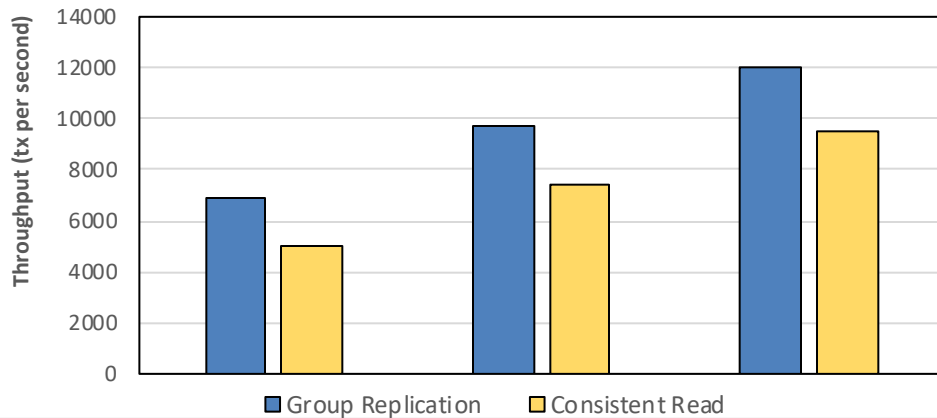
Single-Primary Throughput: Sysbench OLTP RW 1 IDC



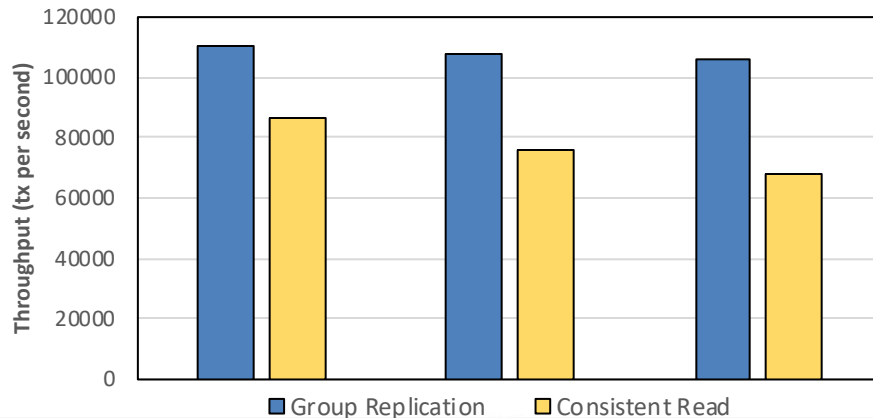
Single-Primary Throughput: Sysbench OLTP RO 1 IDC



Single-Primary Throughput: Sysbench OLTP RW 2 IDC



Single-Primary Throughput: Sysbench OLTP RO 2 IDC





内存使用

冲突检测库，在途事务，恢复时对增量数据的缓存
Paxos消息缓存，writesets开销



节点管理

失败节点重新加入
3+1容灾



读一致性

错误节点摘除读流量
开启流量控制、大事务限制



4

PART FOUR

Blueprint



Optimization

1

冲突检测库的gc优化
内存容量管理
writesets版本

Multiple Primary

2

多写试点
整合一致性读
跨机房优化

Enhance Consistency Read

3

优化性能
提前回放
版本管理

Port with MySQL 8.0.14

4

一致性增强
自动rejoin
指定主库



THANKS