

基于MySQL 8 做分布式数据库， 有哪些坑？

演讲人：万里数据库 王斌

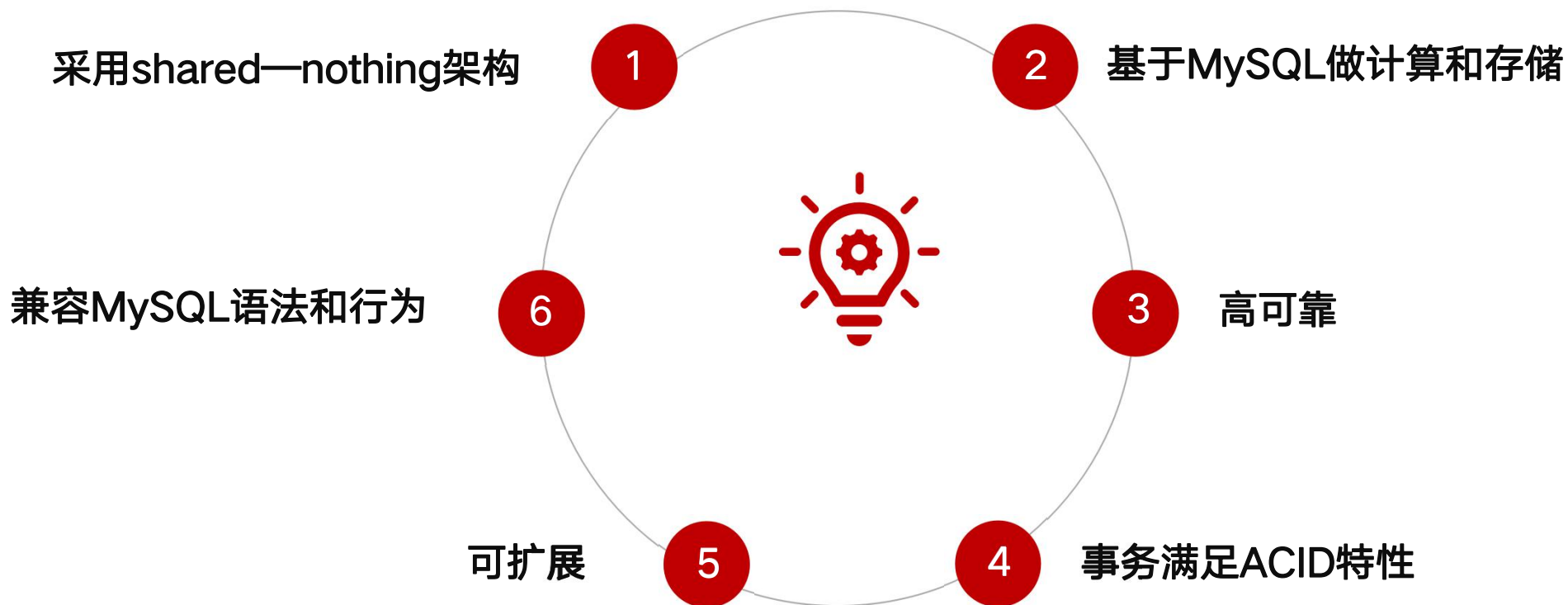


0
1



什么是**基于MySQL**的分布式数据库？

基于MySQL的分布式数据库的限定条件






0
2

基于MySQL做分布式数据库，有什么**优点**？

基于MySQL做分布式数据库的优点

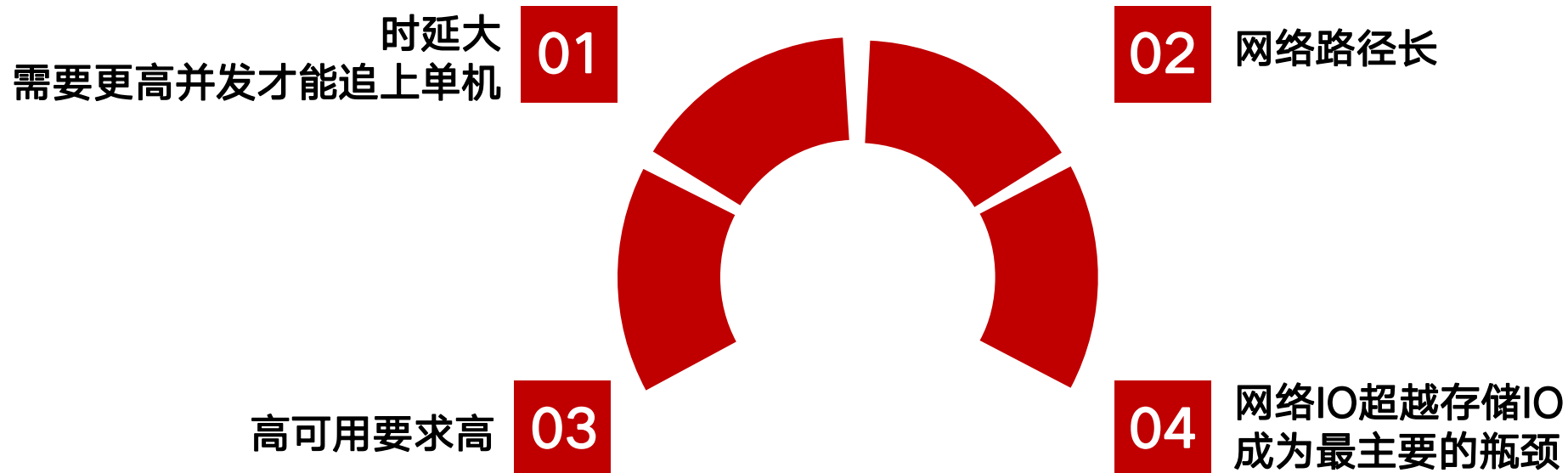




0
3

基于MySQL做分布式数据库，有哪些**特点**？

基于MySQL做分布式数据库，相比单机，有哪些特点？





0
4



基于MySQL做分布式数据库，有哪些坑呢？

MySQL 8 有哪些坑？



Coarse-grained latch 锁调度机制不合理



默认字符集不同导致的性能大坑



支持分布式事务的XA
代码存在大量bug



高可用机制不完善



执行计划只考虑
单机

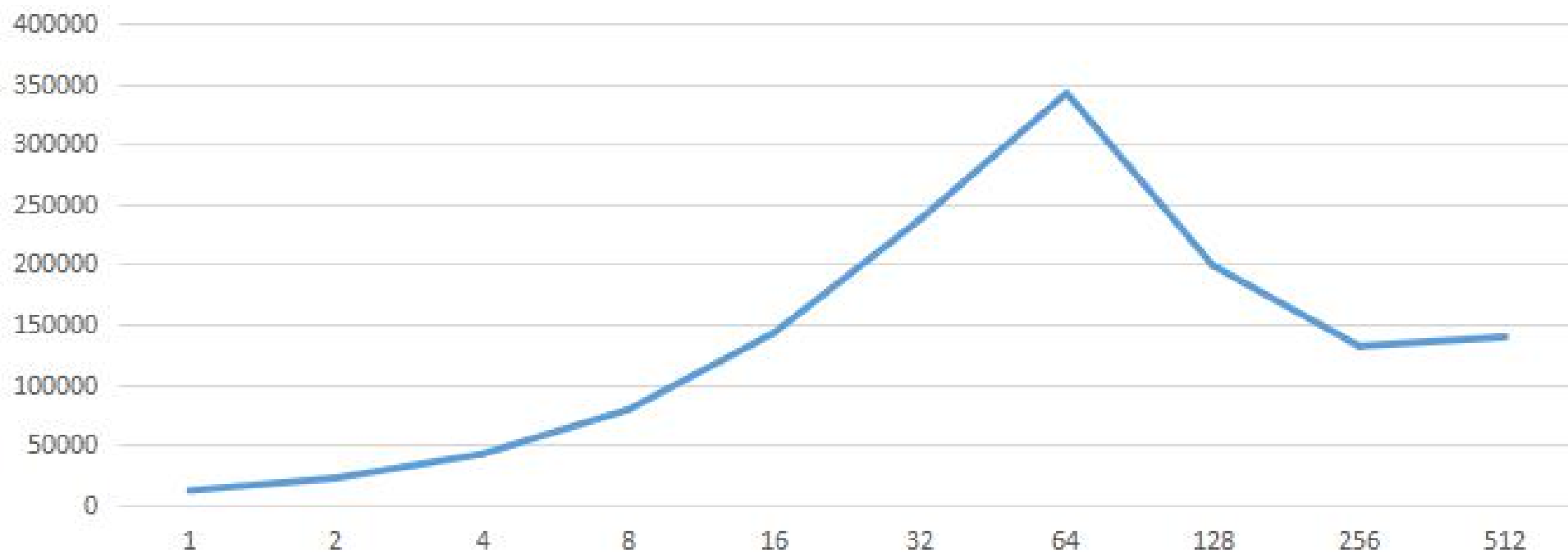


部分场景下
比5.7性能还差

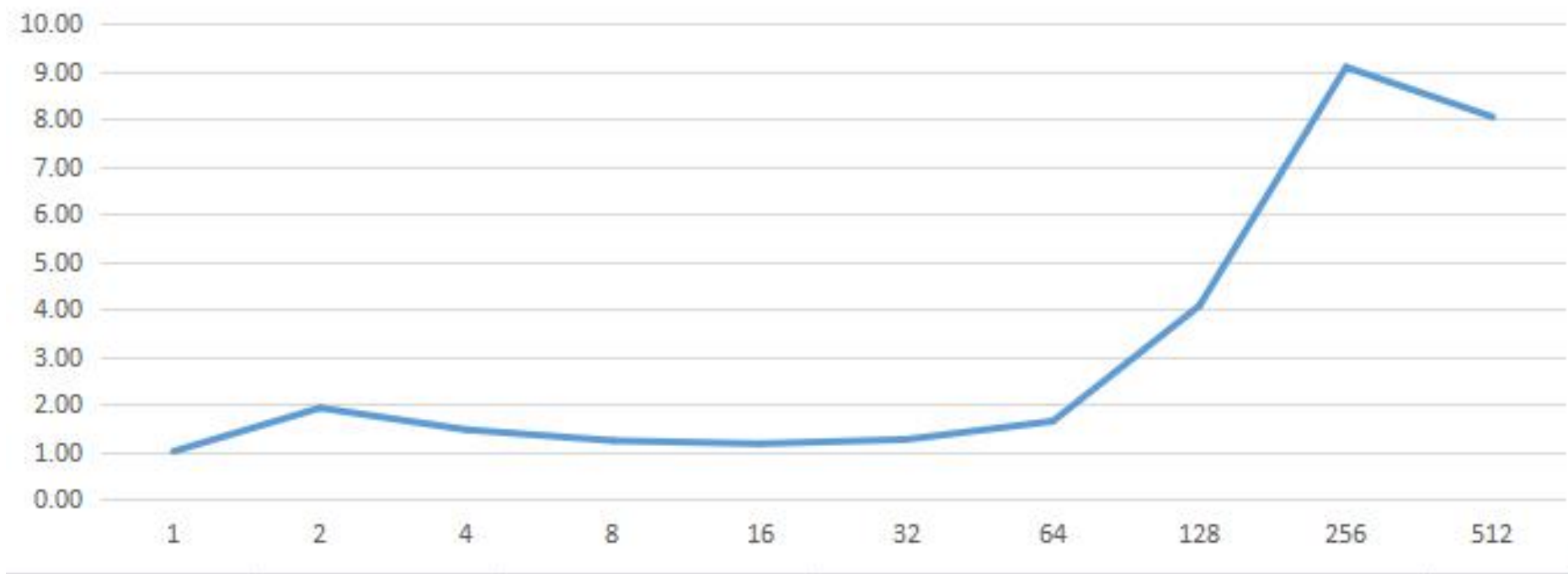
Latch粒度和锁调度机制决定了MySQL并发能力

MySQL并发能力

8.0.25版本，读提交隔离级别下，吞吐量随并发数量关系图



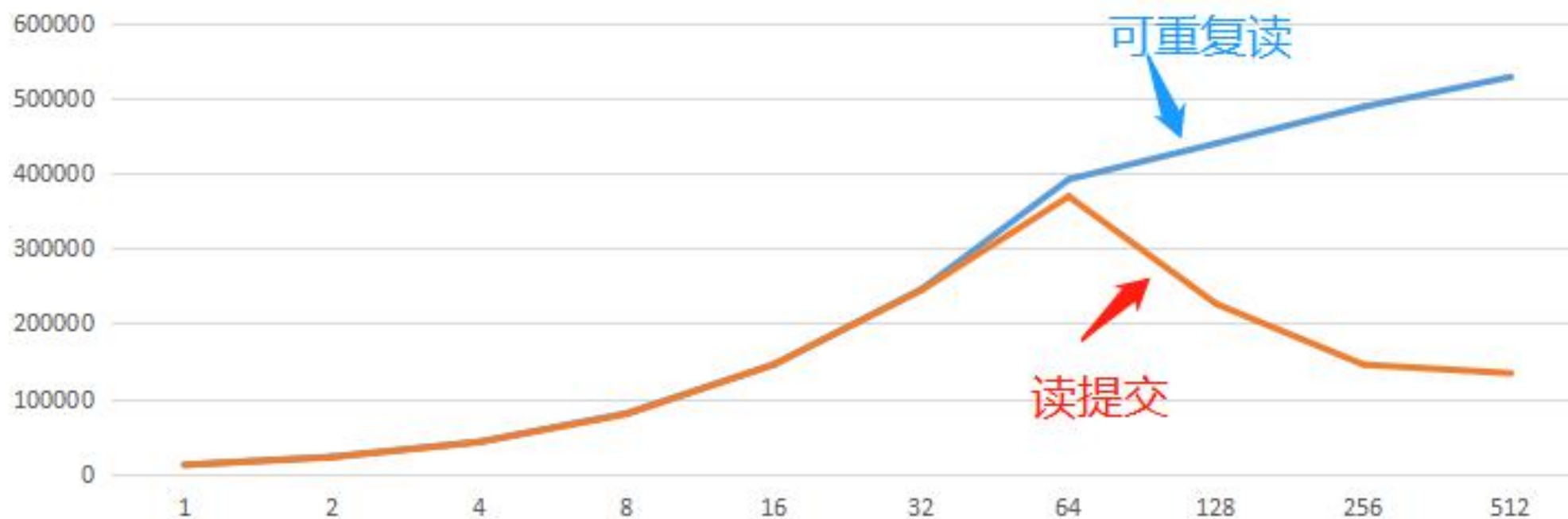
8.0.25版本，读提交隔离级别下，单个请求cpu消耗随并发关系图



主要原因是MySQL trx_sys所用的latch粒度太粗，
影响了读提交的并发扩展性

官方MySQL读提交 vs 官方MySQL可重复读

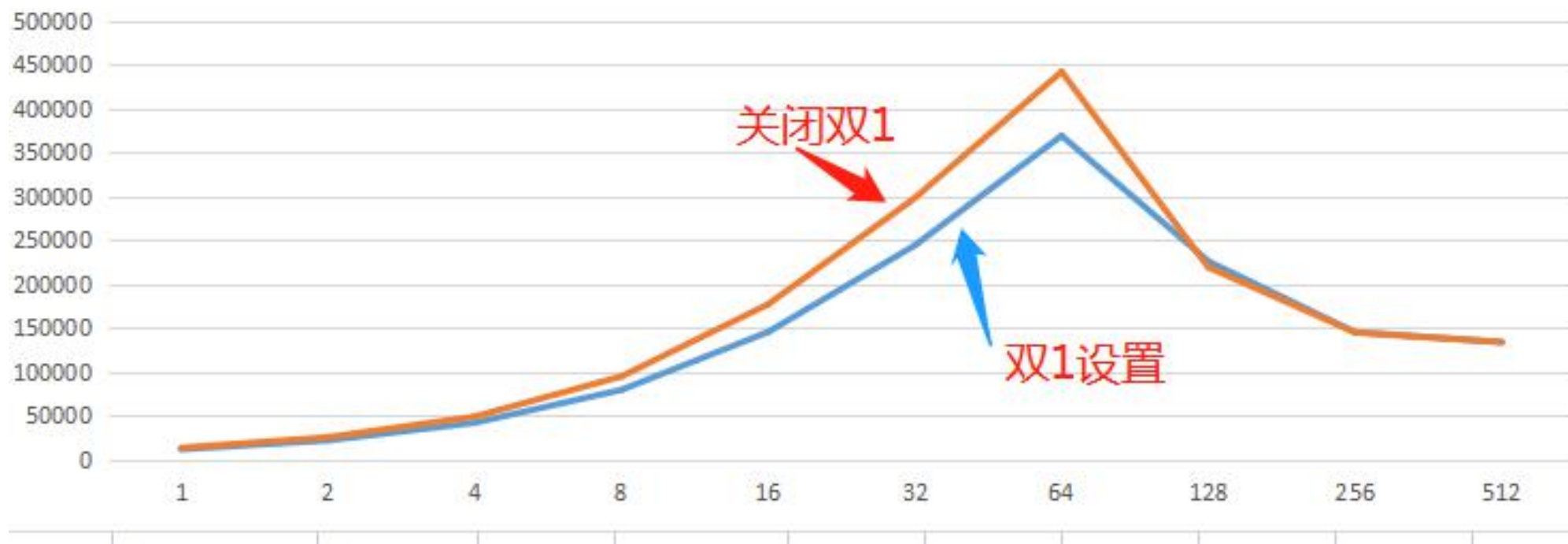
8.0.25版本，不同隔离级别下，吞吐量随并发关系图



通过trx_sys的latch拆分+ Read view优化来解决RC扩展性问题

MySQL灵异现象

8.0.25版本，读提交隔离级别下，吞吐量随并发关系图



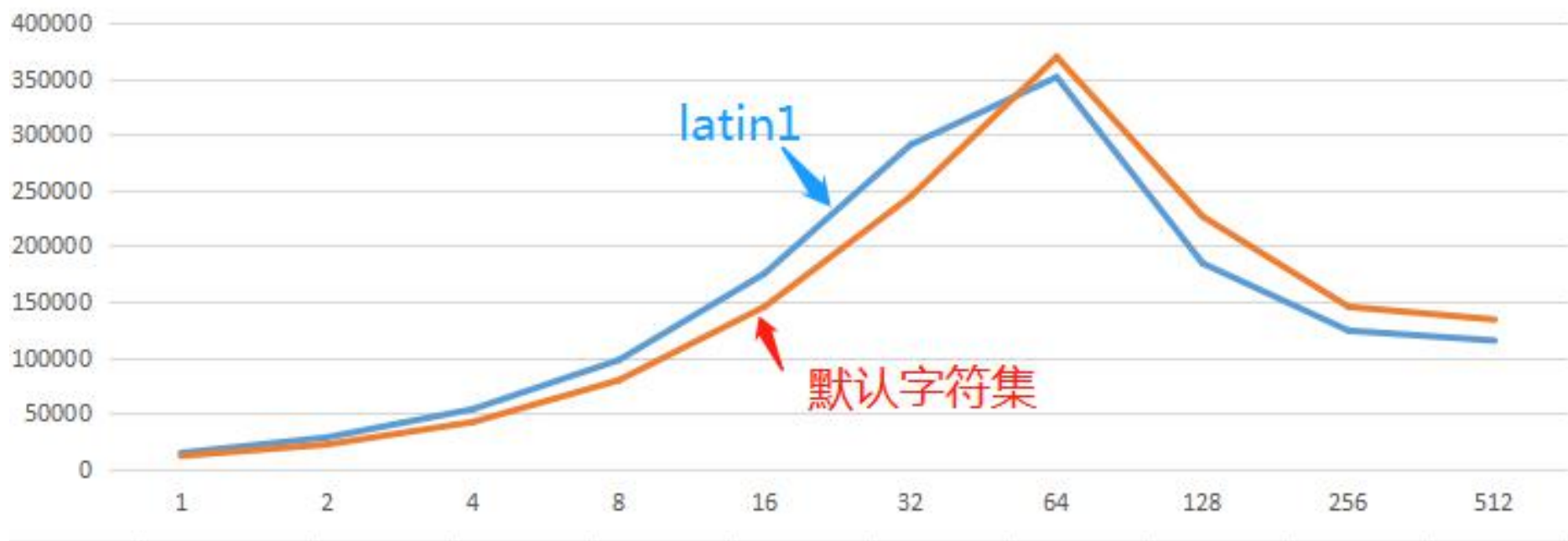
针对读提交，高并发下MySQL设置**双1读写综合性能**往往更好

锁调度机制不合理 + 锁竞争激烈 + Group commit高并发效率高

性能优化在复杂软件面前，往往是**追求一种平衡**

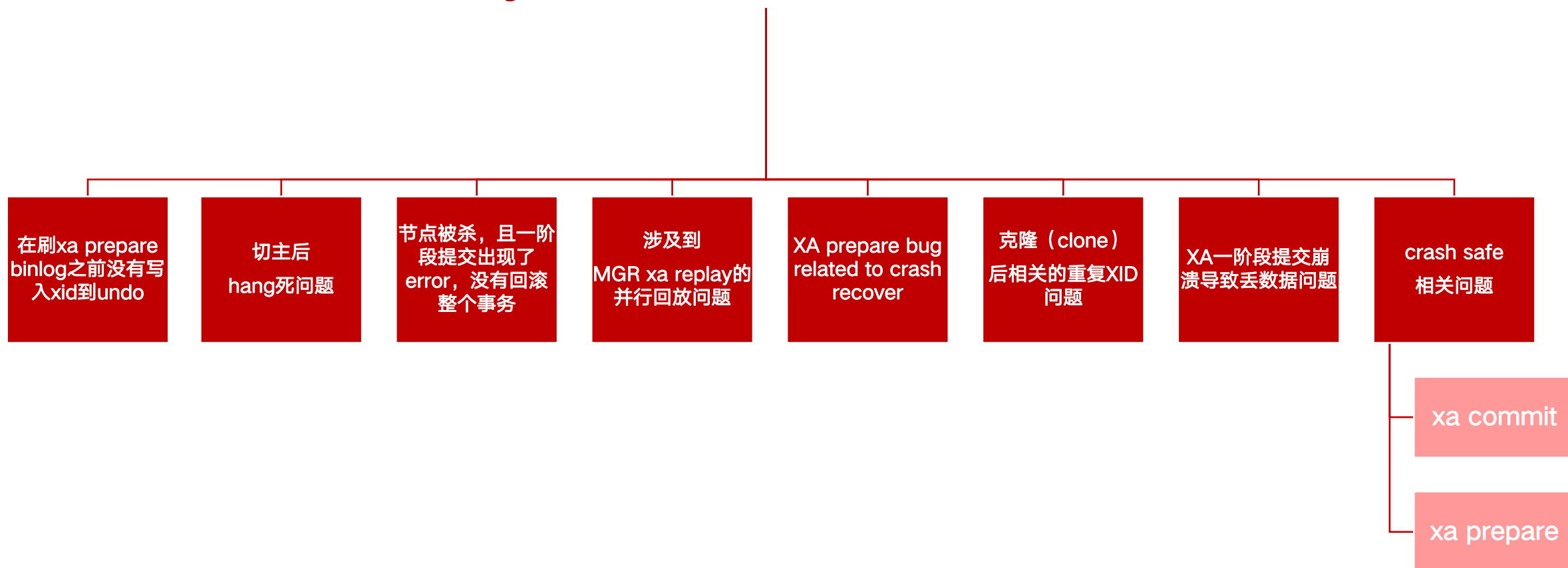
不同字符集带来的性能差异

8.0.25版本, latin1字符集 pk 默认字符集, 读提交, 吞吐量随并发关系图



MySQL XA机制是实现分布式事务的基础

MySQL XA有哪些坑？



MySQL Group Replication高可用机制不完善

MGR完善

需要重新设计

采用没有
认证数据
库的fast
工作模式

根据孟子
算法博士
论文改变
mgr里面
的paxos
算法

设计新的
流控算法

为多主设
计新的认
证数据库
数据结构

增加地理
标签机制

需要大量改造

relay log
batching
入盘机制

xcom
cache内
存分配静
态化机制

改进协程
调度机制

为multi-
master实
现完备的
OCC并发
控制算法

需要修复大量bug

一致性读
写各种同
步问题

纠正多主
模式下的
工作流程

recoveri
ng状态
下异常导
致的数据
紊乱问题

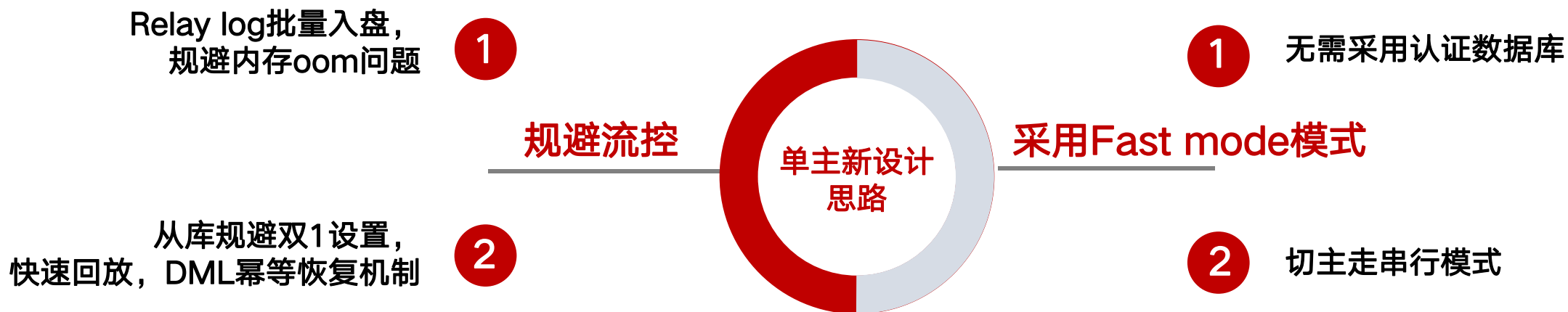
各种视图
问题

MGR改造样例

每秒订单数随时间关系图



单主新设计思路



MySQL执行计划只考虑**单机**

分布式数据库，往往**网络是主要瓶颈**

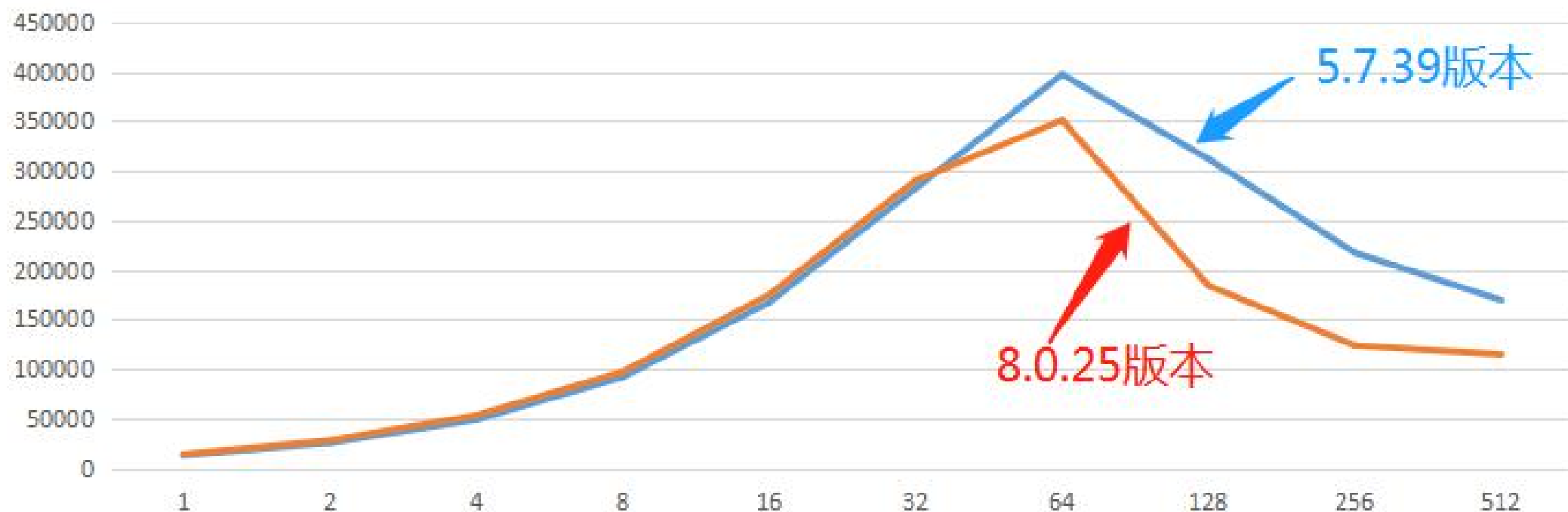
MySQL NDB都无法利用网络来优化执行计划

MySQL 8对5.7的性能改造目前**并不总是很成功**

同一字符集下，不同版本可重复读隔离级别下，吞吐量随并发关系图



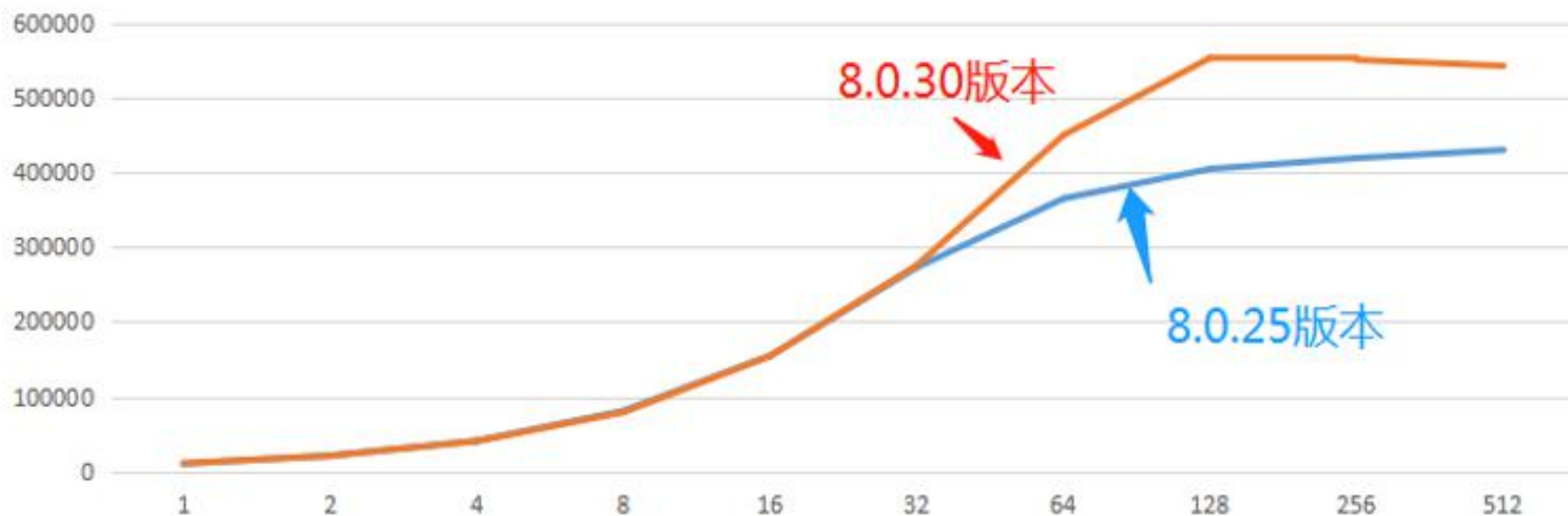
同一字符集下，不同版本读提交隔离级别下，吞吐量随并发关系图



根源在于MySQL 8.0.25的trx_sys所用的latch粒度太粗

最新版本在**trx_sys latch**进行了的一些简易拆分

同一字符集下，不同版本可重复读隔离级别下，吞吐量随并发关系图



/sql/8.0/en/news-8-0-26.html

🔗 ☆ ⚙️ □ 👤 更新

correctness issues. One of the issues addressed caused an XA transaction to be described incorrectly as “recovered”, which occurred when a client session disconnected from an XA transaction after XA PREPARE. (Bug #31870582)

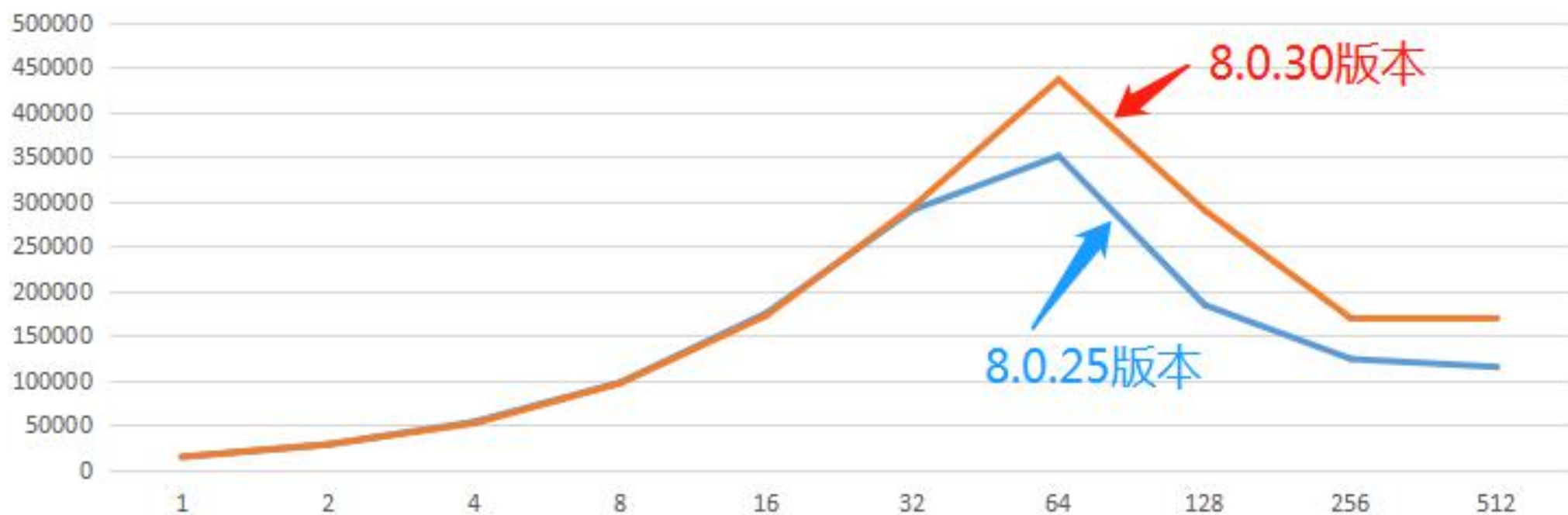
- **InnoDB:** TempTable debug assertion code for an `Indexed_cells` member function (`cell_from_mysql_buf_index_read()`) did not account for non-nullable columns with zero length. (Bug #31091089)
- **InnoDB:** Using the InnoDB memcached plugin, attempting to retrieve multiple values in a single get command returned an incorrect value. (Bug #29675958, Bug #95032)
- **InnoDB:** The `trx_sys_t::serialisation_mutex` was introduced to reduce contention on the on the `trx_sys_t::mutex`. The new mutex protects the `trx_sys_t::serialisation_list` when a transaction number is assigned, which was previously protected by the `trx_sys_t::mutex`.

Thanks to Zhai Weixiang for the contribution. (Bug #27933068, Bug #90643)

明显改善了可重复读的高并发能力

- **Partitioning:** When a table was partitioned by TIMESTAMP and a timestamp literal with a time zone offset was used in the WHERE clause of a SELECT statement, it was possible for a partition to be omitted from the result set.

同一字符集下，不同版本读提交隔离级别下，吞吐量随并发关系图



sql/8.0/en/news-8-0-26.html



- InnoDB: The `lock_sys` sharded `rw_lock` index used random index values generated by the `ut_Idx_IdxVal1()` function, which was not optimal for low-concurrency workloads. (Bug #32880577)
- **InnoDB:** A string value setting for the `innodb_redo_log_encrypt` variable was not handled properly. (Bug #32851525)
 - **InnoDB:** Read-write transaction set (`trx_sys->rw_trx_set`) shards, each with a dedicated mutex, were introduced to alleviate transaction system mutex (`trx_sys->mutex`) contention caused by transaction set insertions and removals. Related enhancements include moving transaction set modifiers to less critical locations, eliminating heap allocation inside of the `TrxUndoRsegs` constructor, converting transaction state (`trx->state`) and transaction start time (`trx->start_time`) fields to `std::atomic` fields, and new assertion code to validate threads that operate on transactions. (Bug #32832196)
 - **InnoDB:** Record buffer logic for the InnoDB memcached GET command was revised. (Bug #32828352)
 - **InnoDB:** The `ut_list` base member in the InnoDB sources now locates list nodes using the element portion of the list type rather than storing a member pointer in the base node of a list at runtime, which wasted resources. The patch also includes other `ut_list` related code improvements. (Bug #32820458)

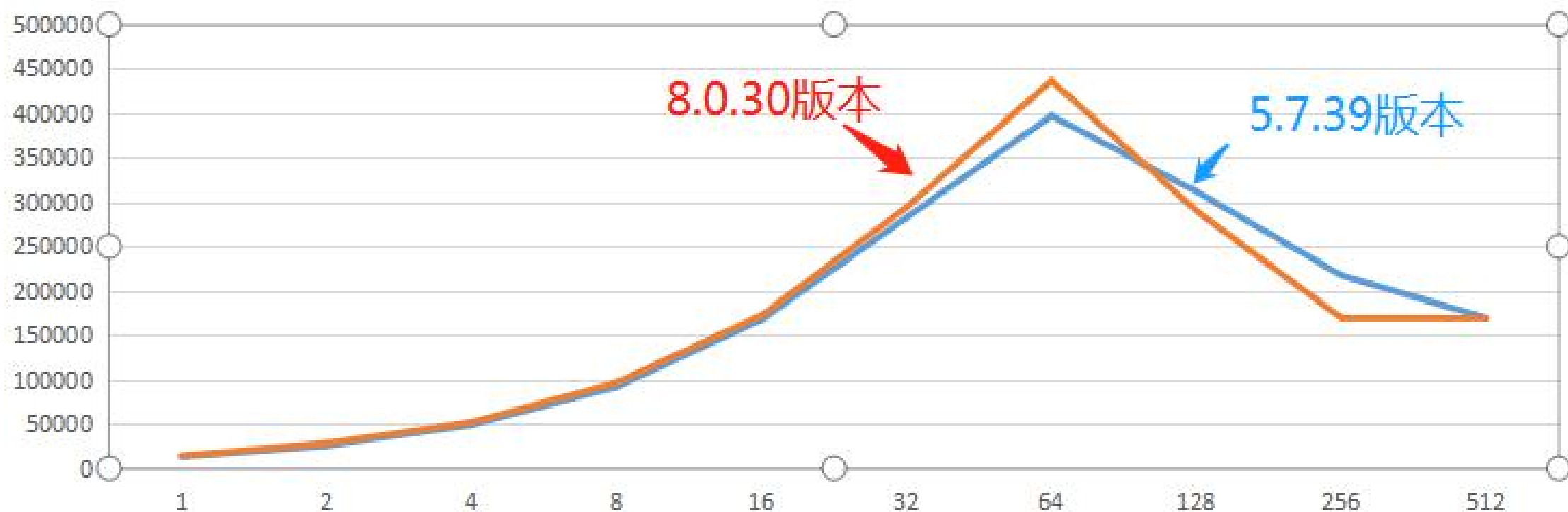
大幅提升了读提交的高并发能力

MySQL 5.7.39 vs MySQL 8.0.30

同一字符集下，不同版本可重复读隔离级别下，吞吐量随并发关系图



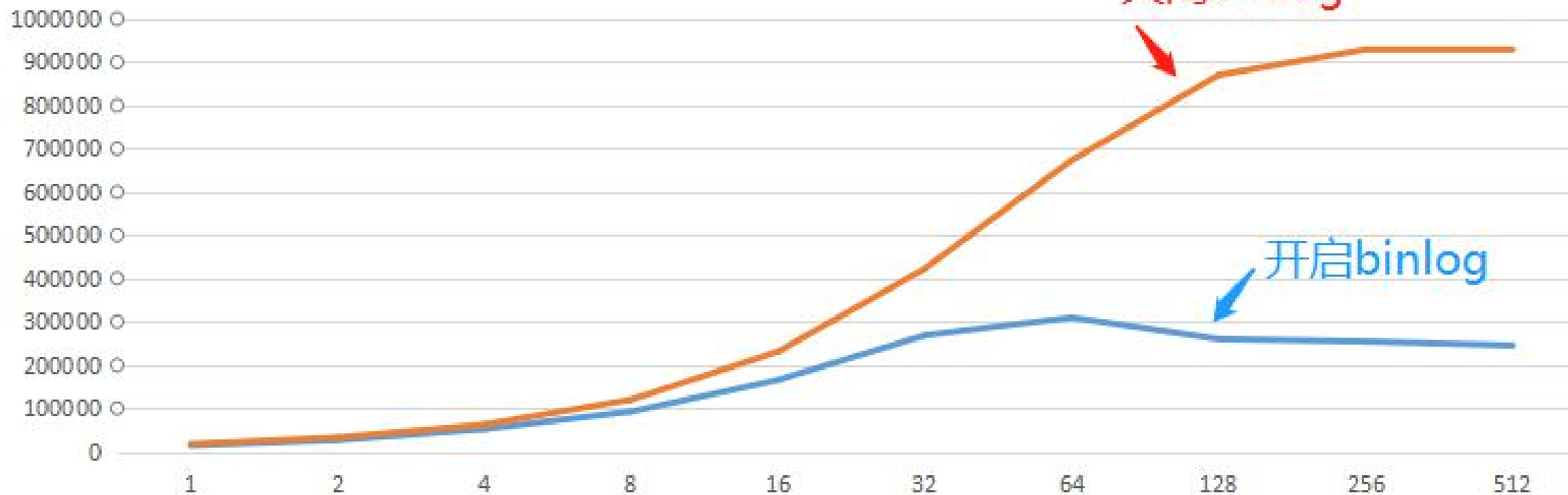
同一字符集下，不同版本读提交隔离级别下，吞吐量随并发关系图



trx_sys所利用的latch仍然可以进一步拆分，但**难度越来越大**

调度机制在MySQL 8多个地方仍然有**巨大的优化空间**

8.0.30版本，读提交，只写事务，吞吐量随并发关系图



困难是暂时的，前途是光明的

开源项目 > 数据库相关 > 数据库服务

GVP 万里数据库 / GreatSQL

Watch

49

Star

410

Fork

64

代码

Issues 6

Pull Requests 1

Wiki

统计

流水线

服务

全部

搜索 Issue



看板

里程碑

+ 新建 Issue

全部

开启的 5

进行中 1

已完成 6

已拒绝 0

创建者

负责人

标签

项目

里程碑

优先级

排序