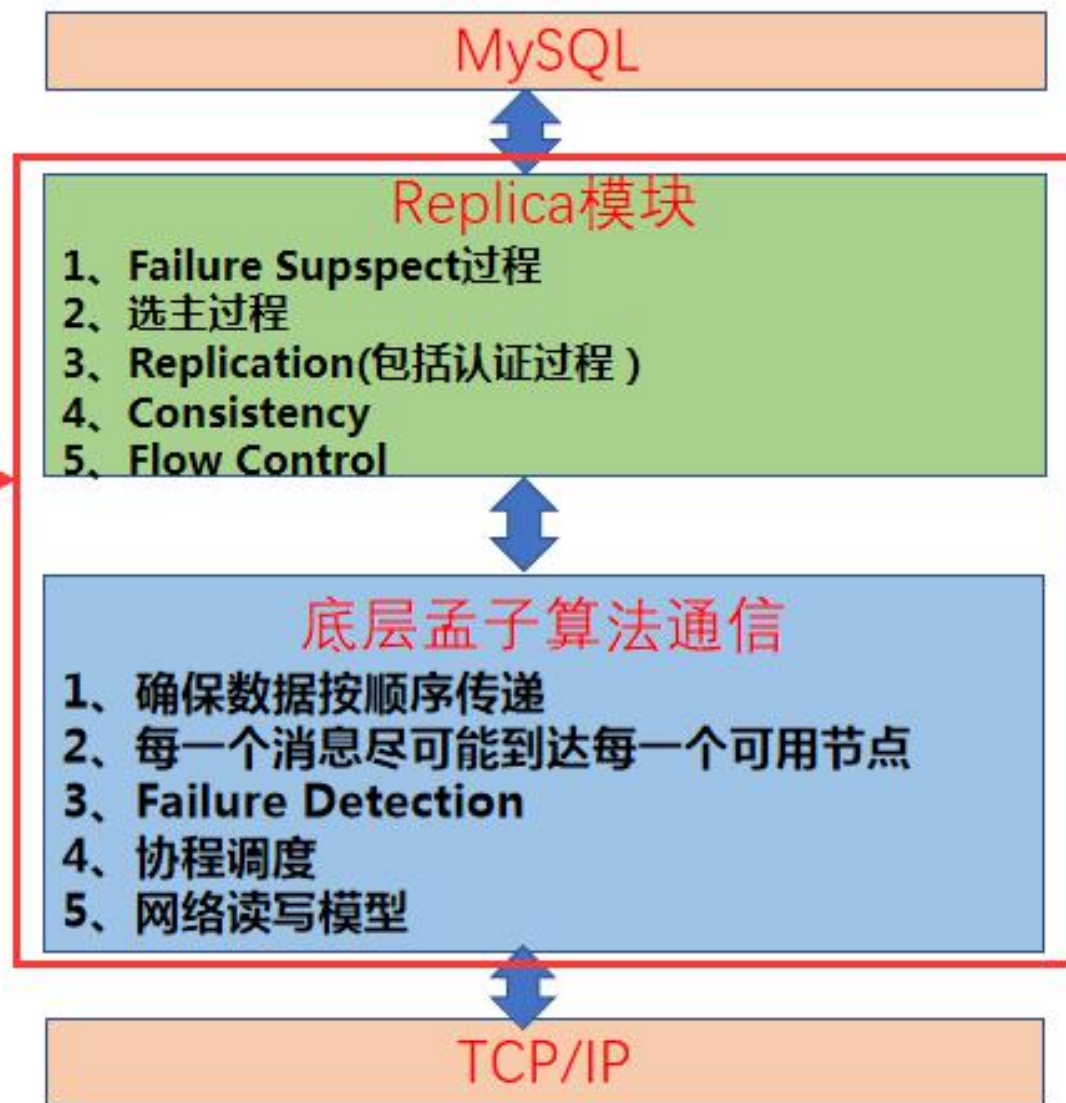# MySQL高可用组件MGR之深度分析

万里数据库

王斌
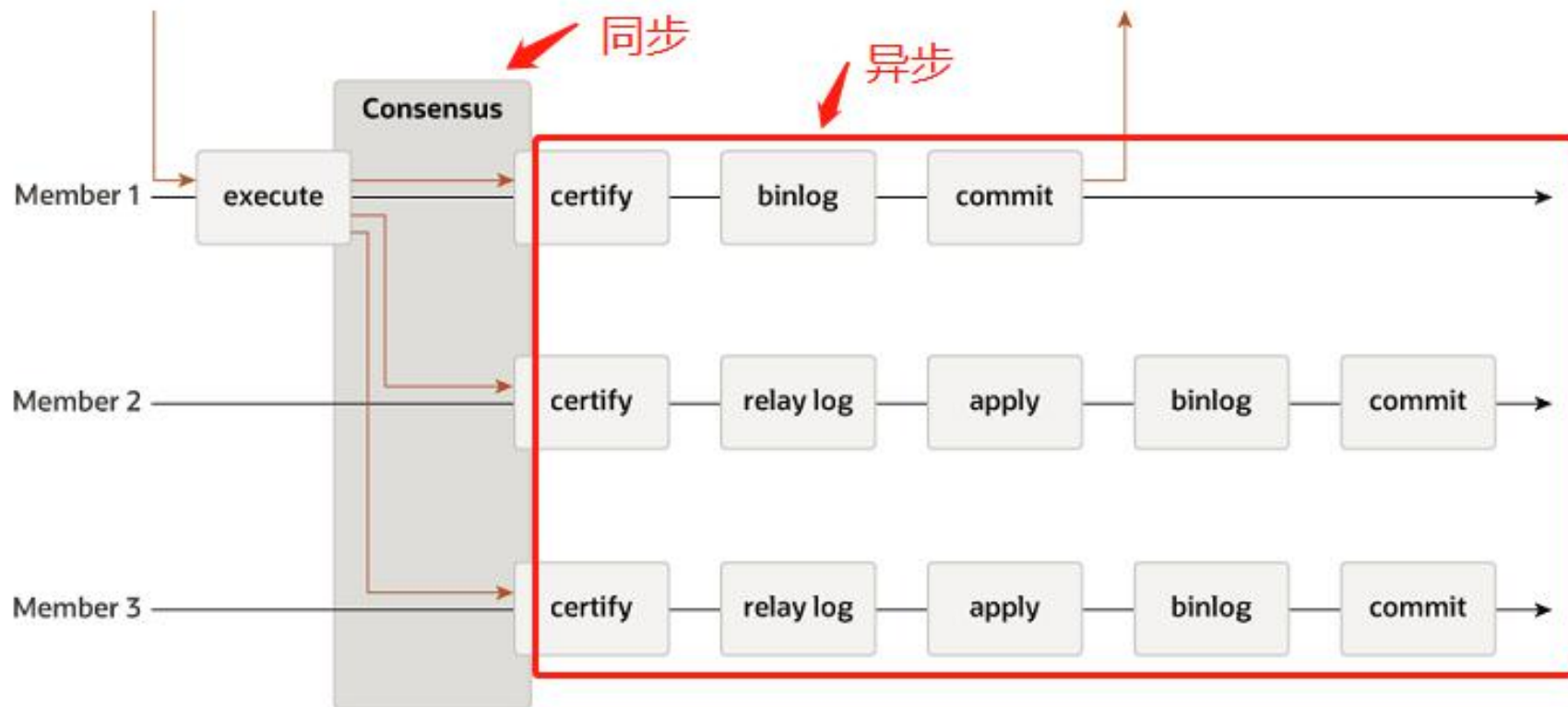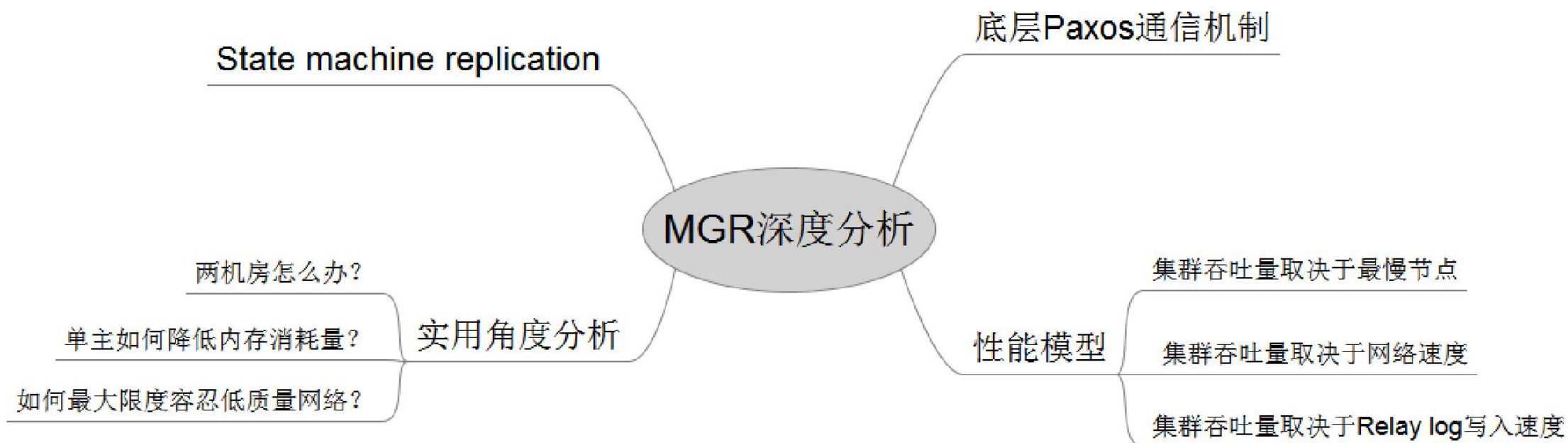
数/造/未/来

# MySQL高可用组件MGR

深 度 分 析 内 容

底层Paxos通信机制

参考 **孟子算法**

# A rotating leader protocol for multi-site systems

为**Wide Area Networks**而生的**Paxos**算法

两个节点之间RTT=**10ms**

```
Threads started!

[ 1s ] thds: 1 tps: 84.51 qps: 512.04 (r/w/o: 0.00/342.02/170.02) lat (ms,95%): 11.87 err/s: 0.00 reconn/s: 0.00
[ 2s ] thds: 1 tps: 86.06 qps: 516.38 (r/w/o: 0.00/344.25/172.13) lat (ms,95%): 12.30 err/s: 0.00 reconn/s: 0.00
[ 3s ] thds: 1 tps: 85.84 qps: 514.03 (r/w/o: 0.00/342.35/171.68) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 4s ] thds: 1 tps: 86.00 qps: 517.01 (r/w/o: 0.00/345.01/172.00) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 5s ] thds: 1 tps: 85.17 qps: 511.01 (r/w/o: 0.00/340.67/170.34) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 6s ] thds: 1 tps: 82.84 qps: 497.07 (r/w/o: 0.00/331.38/165.69) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 7s ] thds: 1 tps: 86.16 qps: 516.95 (r/w/o: 0.00/344.64/172.32) lat (ms,95%): 11.87 err/s: 0.00 reconn/s: 0.00
[ 8s ] thds: 1 tps: 84.84 qps: 509.05 (r/w/o: 0.00/339.37/169.68) lat (ms,95%): 13.95 err/s: 0.00 reconn/s: 0.00
[ 9s ] thds: 1 tps: 86.16 qps: 516.99 (r/w/o: 0.00/344.66/172.33) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 10s ] thds: 1 tps: 85.83 qps: 514.98 (r/w/o: 0.00/343.32/171.66) lat (ms,95%): 11.87 err/s: 0.00 reconn/s: 0.00
SQL statistics:
    queries performed:
```

多写场景下均匀写入节点1

```
Threads started!

[ 1s ] thds: 1 tps: 83.74 qps: 507.45 (r/w/o: 0.00/338.96/168.49) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 2s ] thds: 1 tps: 85.99 qps: 515.96 (r/w/o: 0.00/343.97/171.99) lat (ms,95%): 11.87 err/s: 0.00 reconn/s: 0.00
[ 3s ] thds: 1 tps: 85.05 qps: 510.31 (r/w/o: 0.00/340.21/170.10) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 4s ] thds: 1 tps: 86.93 qps: 521.58 (r/w/o: 0.00/347.72/173.86) lat (ms,95%): 11.87 err/s: 0.00 reconn/s: 0.00
[ 5s ] thds: 1 tps: 85.11 qps: 510.66 (r/w/o: 0.00/340.44/170.22) lat (ms,95%): 11.87 err/s: 0.00 reconn/s: 0.00
[ 6s ] thds: 1 tps: 82.93 qps: 497.60 (r/w/o: 0.00/331.73/165.87) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
[ 7s ] thds: 1 tps: 86.09 qps: 516.52 (r/w/o: 0.00/344.34/172.17) lat (ms,95%): 11.87 err/s: 0.00 reconn/s: 0.00
[ 8s ] thds: 1 tps: 84.98 qps: 509.88 (r/w/o: 0.00/339.92/169.96) lat (ms,95%): 14.21 err/s: 0.00 reconn/s: 0.00
[ 9s ] thds: 1 tps: 85.92 qps: 515.51 (r/w/o: 0.00/343.67/171.84) lat (ms,95%): 12.08 err/s: 0.00 reconn/s: 0.00
SQL statistics:
```

多写场景下，均匀写入节点2

```
Threads started!

[ 1s ] thds: 1 tps: 41.73 qps: 252.34 (r/w/o: 0.00/167.90/84.44) lat (ms,95%): 33.12 err/s: 0.00 reconn/s: 0.00
[ 2s ] thds: 1 tps: 42.04 qps: 255.26 (r/w/o: 0.00/171.18/84.09) lat (ms,95%): 33.12 err/s: 0.00 reconn/s: 0.00
[ 3s ] thds: 1 tps: 43.94 qps: 263.62 (r/w/o: 0.00/175.75/87.87) lat (ms,95%): 33.72 err/s: 0.00 reconn/s: 0.00
[ 4s ] thds: 1 tps: 43.00 qps: 258.00 (r/w/o: 0.00/172.00/86.00) lat (ms,95%): 33.12 err/s: 0.00 reconn/s: 0.00
[ 5s ] thds: 1 tps: 44.07 qps: 264.41 (r/w/o: 0.00/176.27/88.14) lat (ms,95%): 33.12 err/s: 0.00 reconn/s: 0.00
[ 6s ] thds: 1 tps: 45.00 qps: 269.97 (r/w/o: 0.00/179.98/89.99) lat (ms,95%): 31.94 err/s: 0.00 reconn/s: 0.00
[ 7s ] thds: 1 tps: 45.01 qps: 270.03 (r/w/o: 0.00/180.02/90.01) lat (ms,95%): 22.28 err/s: 0.00 reconn/s: 0.00
[ 8s ] thds: 1 tps: 45.93 qps: 275.57 (r/w/o: 0.00/183.71/91.86) lat (ms,95%): 24.83 err/s: 0.00 reconn/s: 0.00
[ 9s ] thds: 1 tps: 45.06 qps: 270.35 (r/w/o: 0.00/180.23/90.12) lat (ms,95%): 22.69 err/s: 0.00 reconn/s: 0.00
[ 10s ] thds: 1 tps: 45.01 qps: 270.04 (r/w/o: 0.00/180.03/90.01) lat (ms,95%): 31.94 err/s: 0.00 reconn/s: 0.00
SQL statistics:
    queries performed:
```

多写场景下只写入一个节点

MGR实现的孟子算法不适合单主场景

To summarize, Mencius temporarily stalls when any of the servers fails while Paxos temporarily stalls only when the leader fails. Also, the throughput of Mencius drops after a failure because of a reduction on available bandwidth, while the throughput of Paxos does not change since it does not use all available bandwidth.
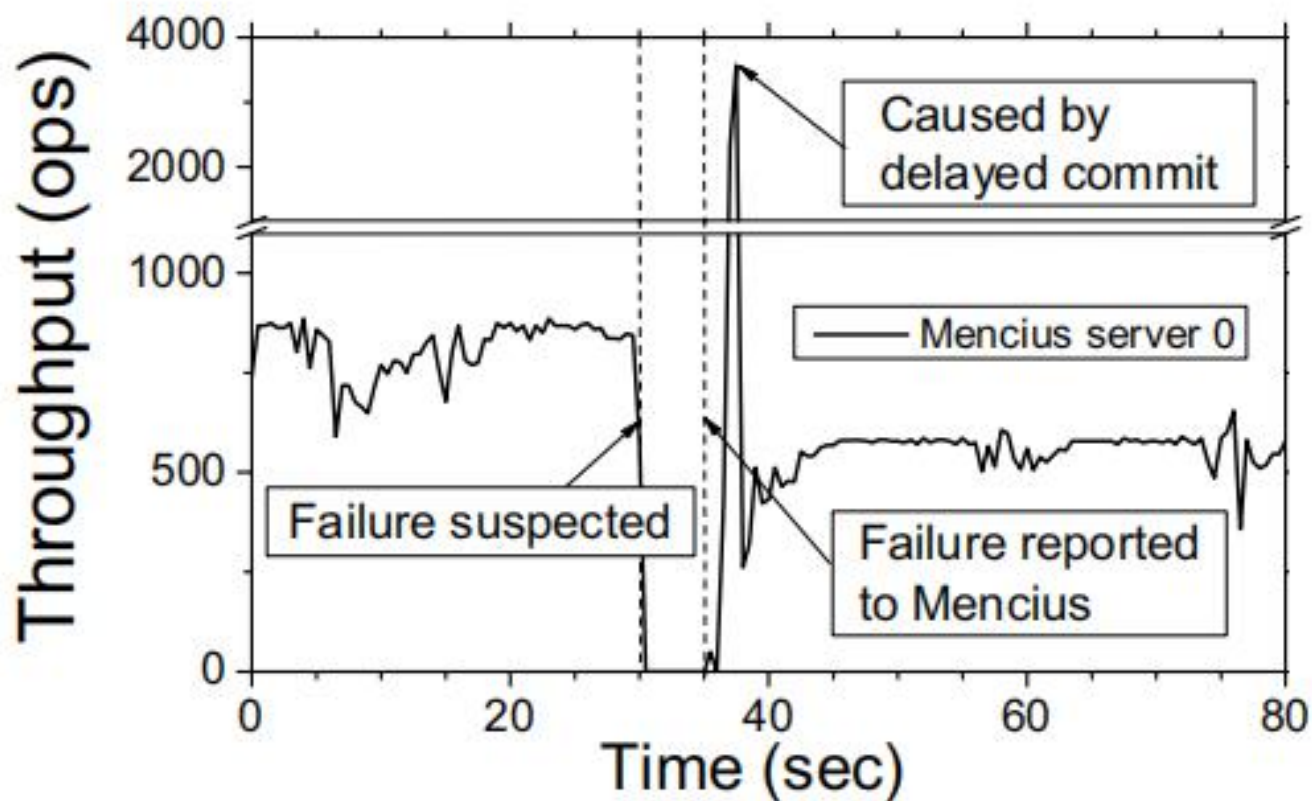
Figure 5.9: The throughput of Mencius when a server crashes

性　能　模　型

**Coordinator allocation:** Mencius's commit latency is limited by the slowest server. A solution to this problem is to have coordinators at only the fastest $f+1$ servers and have the slower $f$ servers forward their requests to the other sites.

# MGR采用了最慢节点性能模型

在POC测试场景下，

最慢节点性能模型是**没有前途**的

**MGR fast mode = 1**

| 并发 | 数据 | Cpu(usr+sys) | Disk | Net (r/s) |
|---|---|---|---|---|
| 600 | 489267 | 31% + 5% | 91% 470M | 270M/320M |
| 1200 | 737304 | 47% + 8% | 90% 500M | 410M/470M |
| 800*3 | 82.5w | 60% + 11% | 89% 490M | 470M/520M |
| 1200*3 | 75.5w | 50% + 10% | 90% 520M | 410M/440M |

**MGR fast mode = 2**

| 并发 | 数据 | Cpu(usr+sys) | Disk | Net (r/s) |
|---|---|---|---|---|
| 600 | 480731 | 27% + 5% | 93% 420M | 270M/300M |
| 1200 | 751829 | 45% + 7% | 90% 380-460M | 420M/465M |
| 800*3 | 89w | 60% + 11% | 90% 450M | 480M/520M |
| 1200*3 | 82w | 49% + 8% | 90% 420M | 430M/470M |

**异步复制**

| 并发 | 数据 | Cpu(usr+sys) | Disk | Net (r/s) |
|---|---|---|---|---|
| 600 | 525099 | 31% + 5% | 93% 470M | 260M/280M |
| 1200 | 768483 | 47% + 8% | 93% 480M | 370M/400M |
| 800*3 | 98.5w | 68% + 11% | 85% 430M | 415M/465M |
| 1200*3 | 96.5w | 71% + 13% | 80% 410M | 380M/440M |

4.3

performance_schema=OFF

innodb_thread_concurrency=0

| 场景 | 600 | 1200 | 2400 | 3600 |
|---|---|---|---|---|
| 4.0 +replication + （xa=1） | 534215.99 | 571137.68 | 546499.63 | 552304.06 |

关闭 performance_schema=OFF 约有 10%左右的性能提升

半同步复制

4.5

Semi-sync = off

innodb_thread_concurrency=0

performance_schema=ON

异步复制

| 场景 | 600 | 1200 | 2400 | 3600 |
|---|---|---|---|---|
| 4.0 +replication + （xa=1） | 758979.23 | 939809.83 | 770283.11 | 653396.29 |

关闭半同步后，减少了半同步的 rpl_semi_sync_master_wait_point = AFTER_SYNC 设置中的 ack
通信与在 slave 上的 relay log 刷盘，tpmC 提升明显

# State Machine Replication

**并不是**每一个MySQL操作都是状态机操作

```
| CHANNEL_NAME             | MEMBER_ID                            | MEMBER_HOST | MEMBER_PORT | MEMBER_STATE | MEMBER_ROLE | MEMBER_VERSION |
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
| group_replication_applier | e07288b3-d88b-11eb-8912-e454e8995a1e | 127.0.0.1   |       63306 | ONLINE       | PRIMARY     | 8.0.22         |
| group_replication_applier | e5fd0466-d88b-11eb-b8c3-e454e8995a1e | 127.0.0.1   |       53306 | ONLINE       | SECONDARY   | 8.0.22         |
| group_replication_applier | eb8b713b-d88b-11eb-acaa-e454e8995a1e | 127.0.0.1   |       43306 | ONLINE       | SECONDARY   | 8.0.22         |
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
3 rows in set (0.01 sec)

mysql> create user testuser identified with mysql_native_password by 'TT123t$';
ERROR 1396 (HY000): Operation CREATE USER failed for 'testuser'@'%'
```

==某些操作失败了，会更新缓存，但没有同步到从库，导致从库和主库状态不一样，破坏了主从一致性==

```
mysql> select * from performance_schema.replication_group_members;
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
| CHANNEL_NAME             | MEMBER_ID                            | MEMBER_HOST | MEMBER_PORT | MEMBER_STATE | MEMBER_ROLE | MEMBER_VERSION |
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
| group_replication_applier | e07288b3-d88b-11eb-8912-e454e8995a1e | 127.0.0.1   |       63306 | ONLINE       | PRIMARY     | 8.0.22         |
| group_replication_applier | e5fd0466-d88b-11eb-b8c3-e454e8995a1e | 127.0.0.1   |       53306 | ONLINE       | SECONDARY   | 8.0.22         |
| group_replication_applier | eb8b713b-d88b-11eb-acaa-e454e8995a1e | 127.0.0.1   |       43306 | ONLINE       | SECONDARY   | 8.0.22         |
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
3 rows in set (0.00 sec)

mysql> create user testuser identified with mysql_native_password by 'TT123t$';
Query OK, 0 rows affected (0.06 sec)

mysql> select * from performance_schema.replication_group_members;
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
| CHANNEL_NAME             | MEMBER_ID                            | MEMBER_HOST | MEMBER_PORT | MEMBER_STATE | MEMBER_ROLE | MEMBER_VERSION |
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
| group_replication_applier | e07288b3-d88b-11eb-8912-e454e8995a1e | 127.0.0.1   |       63306 | ONLINE       | PRIMARY     | 8.0.22         |
+--------------------------+--------------------------------------+-------------+-------------+--------------+-------------+----------------+
1 row in set (0.00 sec)
```
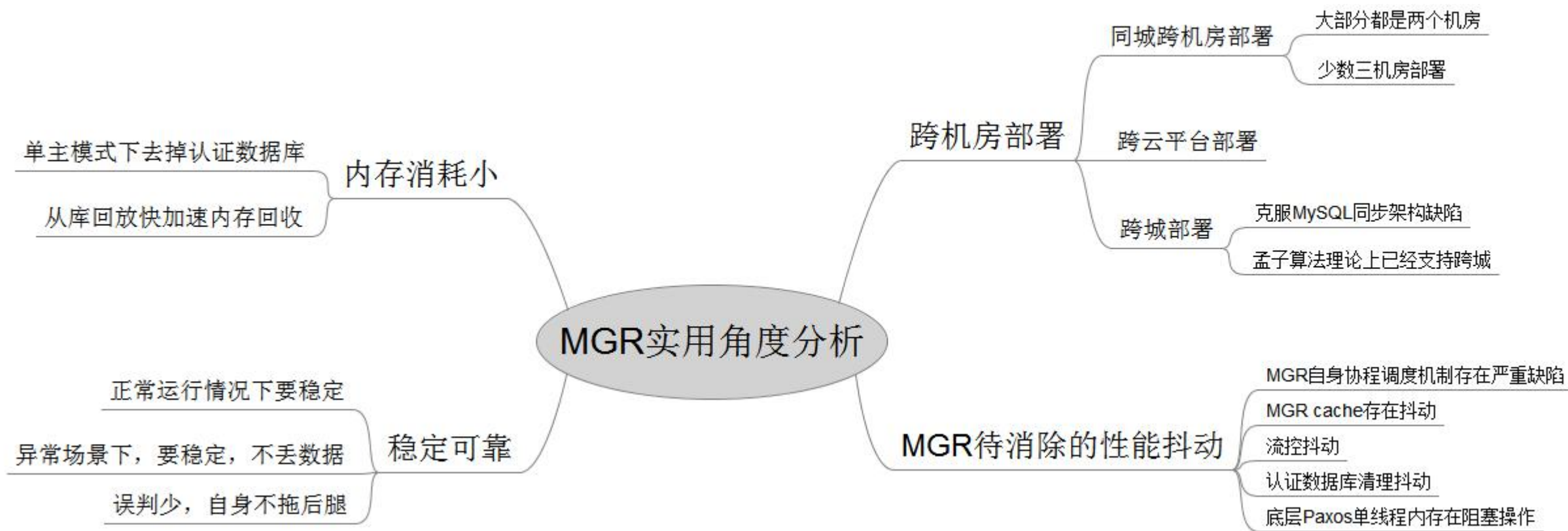
# 实 用 角 度 分 析

# GreatSQL已经实现

## 大部分上述用户诉求

# GreatSQL实现

## 基于地理标签的paxos通信机制

跨 机 房 部 署

对MGR中的底层孟子算法进行了改进
正常情况下一个RTT完成事务

# 消除多处性能抖动

更加公平的协程调度算法
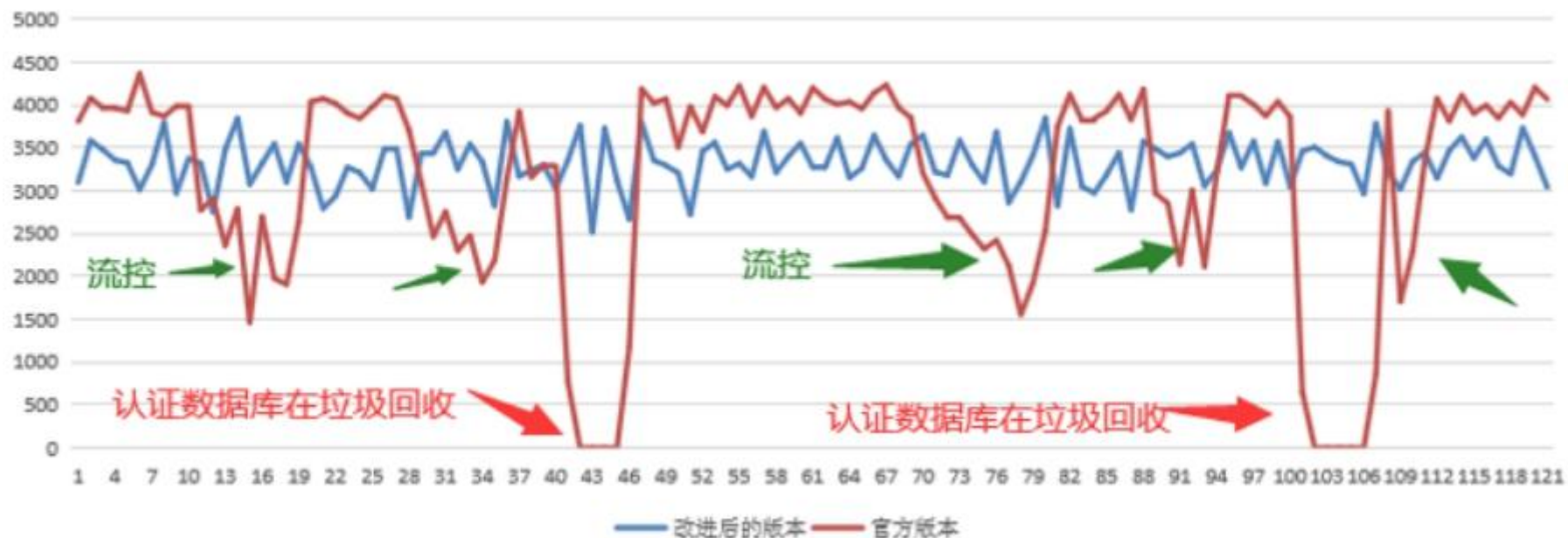
# GreatSQL实现

## 快速单主模式+多主重构
## （新算法+新数据结构）

每秒订单数随时间关系图

**数据库产品供应商**

**操作系统产品供应商**

**北京万里开源软件有限公司**

**专注国产自主可控基础软件产品研发与服务**

GreatDB
万里数据库

**基本情况**

公司成立于2000年10月，创意信息旗下企业，全资控股拓林思软件，20多年基础软件技术沉淀，100%内资背景。

**核心竞争力**

核心研发人员来自MySQL研发中心，熟练掌握数据库源代码；分布式数据库成功应用在金融、电信和能源等行业核心应用系统。

**服务团队**

经验丰富的售后实施团队，提供开发支持、系统优化和驻场支持等多种服务内容，给客户提供原厂7*24小时技术服务。