

### 网易MGR使用和优化实践

温正湖-网易杭研-数据库内核

hzwenzhh@corp.netease.com

# 目录

#### MGR技术实现

- MGR简介
- Paxos全局排序
- 事务认证和回放

#### MGR局限与优化

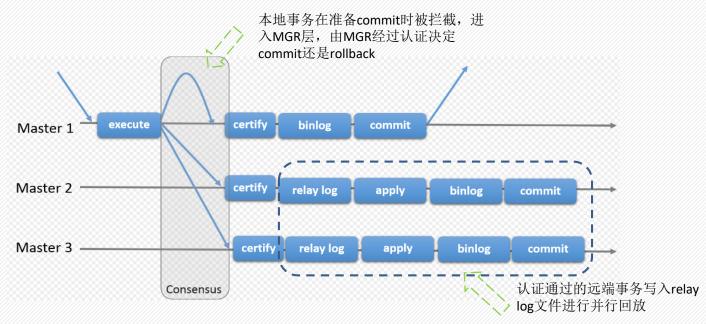
- 事务认证优化
- Paxos实现优化
- 其他优化

#### MGR使用方案

- MGR on RDS
- DDB + MGR
- 考拉跨机房部署

# MGR简介

#### ■ 事务执行流程



# MGR简介

- 事务认证的必要性
- > 多主模式有冲突:
- 1. 多节点同时更新了同一条记录
- 2. 节点间数据延迟,基于旧数据版 本进行更新
- ▶ MGR解决方案:
- 1. 基于记录版本的事务认证机制
  - 在内存中维护冲突检测数据库
  - 各节点独立进行冲突检测,但规则相同,所以结果一样

#### 待认证的事务

- 1. 事务节点server uuid
- 2. 事务执行时的节点gtid\_executed信息snapshot\_version
- 3. 执行事务的线程thread\_id
- 4. 事务gtid是否已存在gtid\_specified
- 5. 事务所修改的记录集的主键列表write set

	PK HASH	GTID_SET	sequence_number
	db1:tb1:1	group_name:1-50	120
冲突检查数据库	db1:tb1:2	group_name:1-20	33
certification_info	db3:tb2:ab	group_name:1-36	90
	db2:tb4:10	group_name:1-100	386

▶ 周期性清理冲突检测数据库中的垃圾数据 (社区版60s一次,硬编码)

# 事务全局排序

■ 事务封装和广播

PAX\_MSG

op : client\_msg, msg\_type : normal

Gcs\_internal\_message\_header::CT\_USER\_DATA

Plugin\_gcs\_message::CT\_TRANSACTION\_MESSAGE

进入Paxos的事务消息封装结构

#### Transaction msg

#### Transaction\_context\_log\_event

- 1. 事务节点server uuid
- 2. 事务执行时的节点gtid\_executed信息snapshot\_version
- 3. 执行事务的线程thread id
- 4. 事务gtid是否已存在gtid\_specified
- 5. 事务所修改的记录集的主键列表write\_set

#### Gtid log event

- 1. 事务节点server\_id
- 2. 事务是否为dml
- 3. 事务组提交信息last\_committed
- 4. 事务组提交信息sequence\_number
- 5、事务是否包含非row格式日志may\_have\_sbr\_stmts

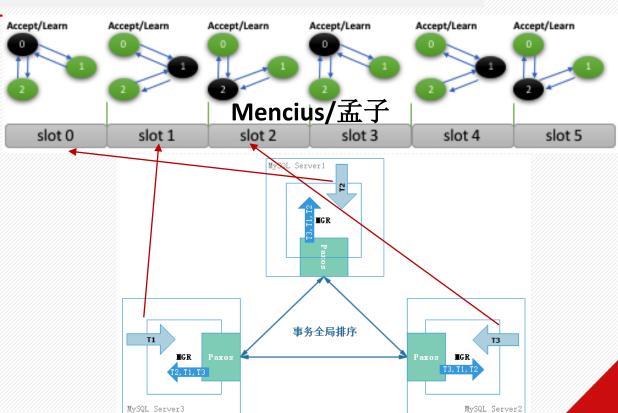
#### 事务数据(log event group)

- 1. Query\_log\_event: "Begin"
- Row\_log\_event: Write/Update/Delete
- 3. .....
- N. Xid log event

# 事务全局排序

- 各节点事务全局排序
  - 不同节点同时提交的事务在 Paxos中进行全局排序
    - round-robin方式

■ 各节点Paxos以相同的次序返回 给GCS模块进行认证



# MGR事务认证

- 事务认证执行框架
  - pipeline



Event cataloger

#### Transaction context log event

- 1.1、标记事务开始;
- 1.2、设transaction\_discarded 为false

#### pipeline



Certification handler

#### Transaction context log event

1.3、缓存携带的事务认证信息



Applier\_handler



#### Gtid log event

#### 2.1、获取缓存的事务信息进行事务冲突检查

2.2a、认证通过:确定gtid,远程事务确定 Nast commited和sequence number

2.2b、认证失败:置transaction\_discarded 为true

#### Gtid\_log\_event

2.3、认证通过的远程事务: 写入relav-log文件





3.1a、transaction\_discarded为true:表示未通过认证,返回

3.1b、transaction\_discarded为

false:继续pipeline处理

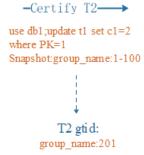
#### 事务数据(log event group)

3.2、认证通过的远程事务: 写入relay-log文件

### MGR事务认证

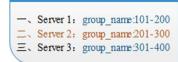
■ 认证并分配GTID (T2属于Server2)

PK HASH	GTID_SET	sequence_number
db1:t1:1	group_name:1-50	120
db1:t12	group_name:1-20	33
db3:t2ab	group_name:1-36	90
db2:t4:10	group_name:1-100	386



PK HASH	GTID_SET	sequence_number	
db1:t1:1	group_name:1-100,201	120	
đb1:t1:2	group_name:1-20	33	
db3:t2:ab	group_name:1-36	90	
đb2:t4:10	group_name:1-100	386	





Avaiable Gtid Set



group\_gtid\_executed: group\_name:1-100,201

#### Avaiable Gtid\_Set

- -. Server 1: group\_name:101-200
- □. Server 2: group\_name:202-300
- ≡. Server 3: group\_name:301-400

# MGR事务认证

■ 事务回放 - 分配 (远端) 事务组提交信息

parallel\_applier\_last\_committed\_global: 60 parallel\_applier\_sequence\_number: 387

PK HASH	GTID_SET	sequence_number
db1:t1:1	group_name:1-101	120
db1:t1:2	group_name:1-20	33
db3:t2:ab	group_name:1-36	90
db2:t4:10	group_name:1-100	386

T2 group\_name:1-101

T2 group commit info:
last\_commited: 120
sequence\_number: 387

parallel\_applier\_last\_committed\_global: 60 parallel\_applier\_sequence\_number: 388

	PK HASH	GTID_SET	sequence_number
	db1:t1:1	group_name:1-101	387
	db1:t1:2	group_name:1-20	33
	db3:t2:ab	group_name:1-36	90
•	db2:t4:10	group_name:1-100	386

# 目录

#### MGR技术实现

- MGR简介
- Paxos全局排序
- 事务认证和回放

#### MGR局限与优化

- 事务认证优化
- Paxos模块优化
- 其他优化

#### MGR使用方案

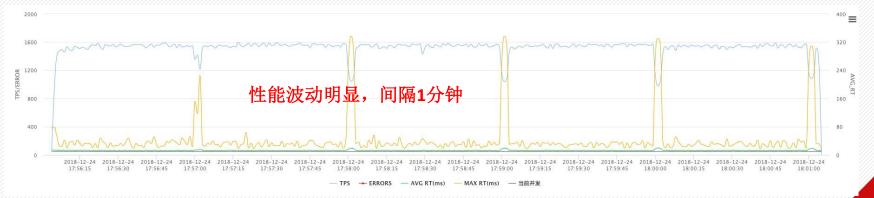
- MGR on RDS
- DDB + MGR
- 考拉跨机房部署

# 事务认证模块优化

- 冲突检测数据库大小不可控
  - Writeset过多导致性能波动
  - 所占内存过大导致OOM
  - 垃圾清理周期不可调



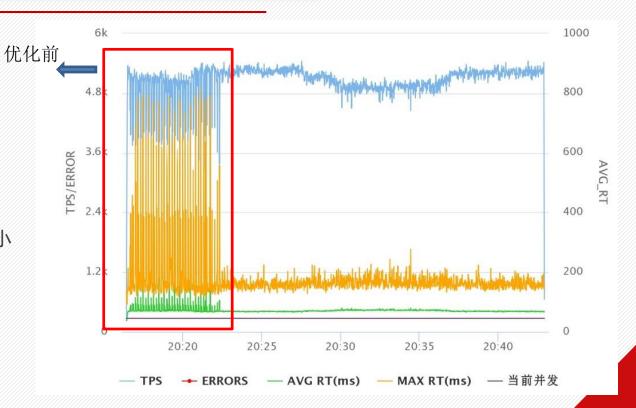
- 将gtid\_executed广播和垃圾清理解耦
- 动态调整广播周期和清理周期
- Writeset数目减半 (5.7版本)
- Writeset个数纳入流控管理
- · 垃圾清理不依赖gtid\_executed (单主模式)
- Snapshot\_version置空 (单主模式)



# 事务认证模块优化 - 效果



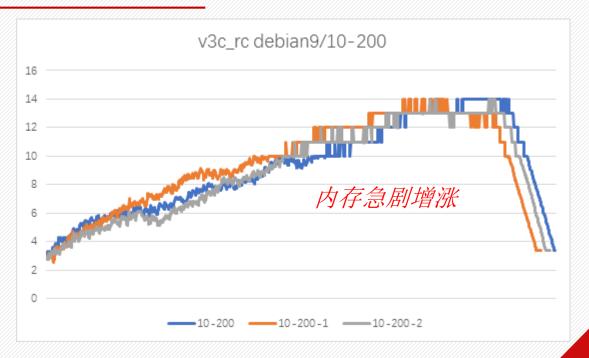
- 性能有所提升
- 一般事务场景下内存波动小



### Paxos模块优化

#### ■ 现状

- Paxos cache大小不透明,不可调
- · 网络瓶颈下, 大事务导致OOM



NDC 10并发,200记录batch,每条个batch约30MB

### Paxos模块优化

#### ■ Paxos cache优化

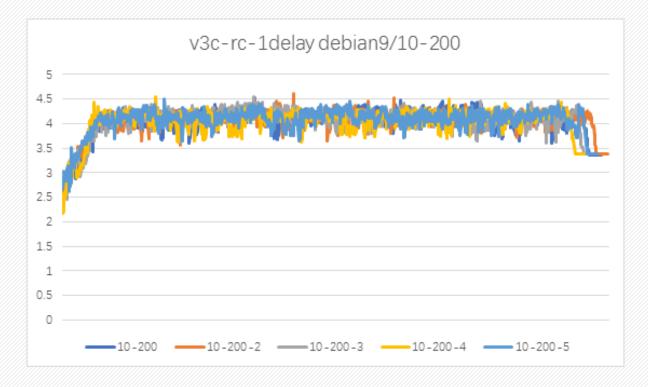
- · 动态调整paxos cache上限阈值,显示当前大小
- 达到阈值时从LRU而不是空闲队列分配lru\_machine

#### ■ 大事务场景优化

- 限制单个消息大小:减小MAX\_BATCH\_SIZE
- Proposer: 增大paxos消息propose retry周期
- Executor: 增大paxos消息获取周期,减小获取个数

# Paxos模块优化 – 效果





# 其他限制和优化

#### 1、MGR细粒度流控机制:

新增7个节点流控参数,支持在Group各节点分配负载配额、设置预留配额、设置故障恢复的配额,控制流控的作用周期等。详见官

#### 2、MGR状态及统计信息显示增强:

支持在一个MGR节点查看其他节点的状态和统计信息,replication\_group\_members表增加2个字段分别表示节点角色和版本信息, 待验证和已经执行的异地事务数量,本地事务总数和被回滚的数量。详见<u>https://mysqlhighavailability.com/group-replication-extending</u>

#### Group Replication模块

- 基于GTID\_EXECUTED的选主策略:同等条件选择gtid\_executed最大的节点为Primary
- 故障恢复时数据源选择策略优化:优先选择非Primary节点作为复制源
- 支持无被选举权的节点角色: group\_replication\_voteless\_member
- 保持Group重配置后Primary节点只读状态: group\_replication\_readonly\_after\_recofig
- 新增GCS/XCOM层配置参数: 视图变更等待group\_replication\_components\_wait\_time和Paxos客户端连接数group\_replication\_max\_handlers
- 增加更多状态信息:当前冲突检测状态group\_replication\_conflict\_detection\_status,本周期配额group\_replication\_flow\_control\_quota\_size及已使用量 group\_replication\_flow\_control\_quota\_used
- 网络分区超过默认30s后,节点退出MGR变为只读group\_replication\_unreachable\_majority\_timeout

# 提交的buglist

ID#	<u>Date</u>	<u>Updated</u>	<u> Type</u>	<u>Summary</u>
<u>89582</u> ☑	2018-02-08 5:04	2018-12-19 12:16	MySQL Server: Group Replication	Node may not switch to ONLINE under consistent load
<u>89938</u> ☑	2018-03-07 3:13	2018-07-11 9:47	MySQL Server: Group Replication	Rejoin old primary node may duplicate key when recovery
90213	2018-03-26 8:02	2018-08-02 10:53	MySQL Server: Group Replication	remote nodes COUNT_TRANSACTIONS_ROWS_VALIDATING are incorrect
91042	2018-05-28 6:19	2018-08-08 15:45	MySQL Server: Group Replication	stop group_replication hang with select replication_connection_status
91397	2018-06-25 4:38	2018-11-05 3:36	MySQL Server: Group Replication	certification db size should limited or add to flow control
<u>91646</u> ☑	2018-07-16 4:40	2019-03-19 7:48	MySQL Server: XA transactions	xa command still operation when super_read_only is true
91671 M	2018-07-17 2:08	2018-11-27 12:05	MySQL Server: Group Replication	stop slave sql_thread for channel 'group_replication_applier' could not return
91864 M	2018-08-02 8:17	2018-09-06 10:04	MySQL Server: Group Replication	Group_member_info::conflict_detection_enable is always OFF in multi master mode
93004	2018-10-30 7:24	2018-11-05 14:56	MySQL Server: Group Replication	MGR failed to boot in VPC(ECS) when using public ip
94814	2019-03-28 11:12	2019-04-15 10:05	MySQL Server: XA transactions	slave replication lock wait timeout because of wrong trx order in binlog file
<u>94881</u>	2019-04-03 4:37	2019-04-23 12:34	MySQL Server: InnoDB storage engine	slave replication lock wait timeout because of supremum record

# 目录

#### MGR技术实现

- MGR简介
- Paxos全局排序
- 事务认证和回放

#### MGR局限与优化

- 事务认证优化
- Paxos实现优化
- 其他优化

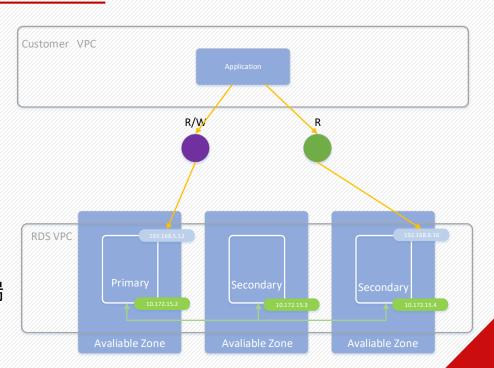
#### MGR使用方案

- MGR on RDS
- DDB + MGR
- 考拉跨机房部署

#### MGR on RDS

#### 故障切换:

- Agent 节点每秒心跳信息汇总成员状态 (replication group members)
- 管控服务识别view发生变化 (Priamry 节点变化)
- 等待所有远端同步的日志回放
   Count\_transactions\_in\_queue +
   Count\_transactions\_remote\_in\_applier\_
   queue
- 管控服务通过调用弹性网关服务切换客户端的状态



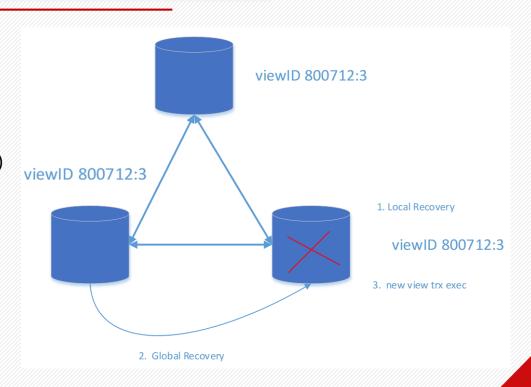
#### MGR on RDS

#### 故障切换的修复:

- 尝试重启宕机节点
- Local Recovery(本地relay log)
- Global Recovery(Remote Binary log)
- 缓存事务执行(当前view执行)

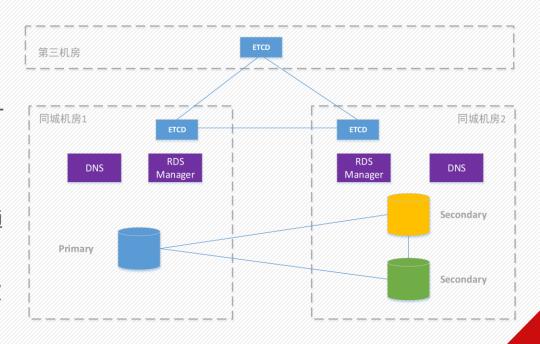
#### 无法重启或者binlog已删除:

- 选择一个Secondary 全量备份
- 重新加入Group Replicaiton



# 考拉MGR跨机房

- 基于Majority Quorum,机房级别故障,可以实现数据一致
- DNS基于ETCD实现单机房故障,另外一个机房可以修改DNS地址,跨机房高可用
- 机房级别故障场景下,单Primary节点通过Force members 强制服务
- 网络5秒以内的抖动会造成业务写入间歇 性抖动

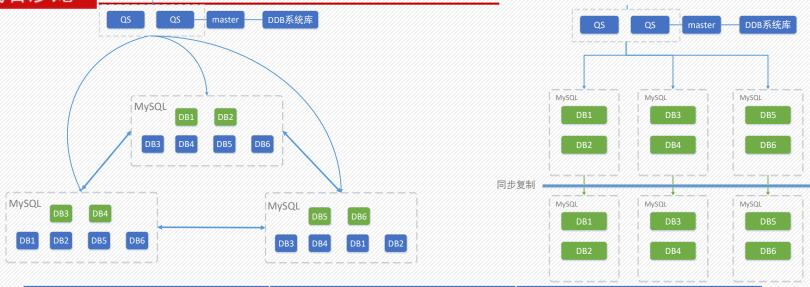


# 考拉MGR实例

名称	可用区	实例类型	THE PERSON NAMED IN
haitao-jd-price-offstation	可用区 A	金融版	Ì
haitao-jd-cachemanage	可用区 A	金融版	I
haitao-jd-mini	可用区 A	金融版	1

名称	可用区	实例类型	数据库引擎	规格	状态
haitao-jdb-ddb-stable4	可用区 A	金融版	MySQL 5.7.20	8核32GB	❷ 运行中
haitao-jdb-ddb-stable3	可用区 A	金融版	MySQL 5.7.20	8核32GB	❷ 运行中
haitao-jdb-ddb-stable02	可用区 A	金融版	MySQL 5.7.20	8核32GB	❷ 运行中
haitao-jdb-ddb-stable01	可用区 A	金融版	MySQL 5.7.20	8核32GB	❷ 运行中
haitao-jdb-ddb-perf04	可用区 A	金融版	MySQL 5.7.20	16核32GB	❷ 运行中
haitao-jdb-ddb-perf03	可用区 A	金融版	MySQL 5.7.20	16核32GB	❷ 运行中
haitao-jdb-ddb-perf02	可用区 A	金融版	MySQL 5.7.20	16核32GB	❷ 运行中
haitao-jdb-ddb-perf01	可用区 A	金融版	MySQL 5.7.20	16核32GB	❷ 运行中

### **DDB + Multi-Master MGR**



	DDB+MGR	DDB+VSR
MySQL节点数	3	6
数据副本数	3	2
QPS	13526	11616



# Q&A

#### 参考资料:

- MySQL官方公开的文档和PPT
- 知乎专栏:数据库内核 https://zhuanlan.zhihu.com/c\_2060713 40/settings/posts