

## Table of Contents

Introduction	1.1
0.概述	1.2
0.1 dble 简介与整体架构	1.2.1
0.2 dble对MyCat做的增强	1.2.2
0.3 快速开始	1.2.3
0.4 docker镜像快速开始	1.2.4
0.5 docker-compose快速开始	1.2.5
0.6 数据拆分简介	1.2.6
1.配置文件	1.3
1.1 rule.xml	1.3.1
1.2 schema.xml	1.3.2
1.3 server.xml	1.3.3
1.4 wrapper.conf	1.3.4
1.5 log4j2.xml	1.3.5
1.6 cache配置	1.3.6
1.6.1 cache配置	1.3.6.1
1.6.2 ehcache配置	1.3.6.2
1.7 全局序列配置	1.3.7
1.7.1 MySQL offset-step方式	1.3.7.1
1.7.2 时间戳方式(类Snowflake)	1.3.7.2
1.7.3 分布式时间戳方式(类Snowflake)	1.3.7.3
1.7.4 分布式offset-step方式	1.3.7.4
1.8 myid.properties	1.3.8
1.9 自定义拆分算法	1.3.9
1.10 配置文件变更记录	1.3.10
1.11 自定义告警	1.3.11
1.11 自定义全局表一致性检查	1.3.12
2.功能描述	1.4
2.1 管理端命令	1.4.1
2.1.1 select命令	1.4.1.1
2.1.2 set命令	1.4.1.2
2.1.3 show命令	1.4.1.3
2.1.4 switch命令	1.4.1.4
2.1.5 kill命令	1.4.1.5
2.1.6 stop命令	1.4.1.6
2.1.7 reload命令	1.4.1.7
2.1.8 rollback命令	1.4.1.8
2.1.9 offline命令	1.4.1.9
2.1.10 online命令	1.4.1.10
2.1.11 file命令	1.4.1.11
2.1.12 log命令	1.4.1.12
2.1.13 配置检查命令	1.4.1.13
2.1.14 pause & resume 命令	1.4.1.14
2.1.15 慢查询日志相关命令	1.4.1.15
2.1.16 创建/删除物理库命令	1.4.1.16
2.1.17 check @@@metadata命令	1.4.1.17
2.1.18 release @@@metadata命令	1.4.1.18
2.1.19 split命令	1.4.1.19
2.1.20 flow_control 命令	1.4.1.20
2.2 全局序列	1.4.2
2.2.1 MySQL offset-step方式	1.4.2.1
2.2.2 时间戳方式	1.4.2.2
2.2.3 分布式时间戳方式	1.4.2.3
2.2.4 分布式offset-step方式	1.4.2.4
2.3 读写分离	1.4.3
2.4 注解	1.4.4
2.5 分布式事务	1.4.5
2.5.1 XA事务概述	1.4.5.1
2.5.2 XA事务的提交以及回滚	1.4.5.2
2.5.3 XA事务的后续补偿以及日志清理	1.4.5.3
2.5.4 XA事务的记录	1.4.5.4
2.5.5 一般分布式事务概述	1.4.5.5
2.6 连接池管理	1.4.6
2.7 内存管理	1.4.7
2.8 集群同步协调&状态管理	1.4.8
2.9 grpc 告警	1.4.9
2.10 表meta数据管理	1.4.10
2.10.1 Meta信息初始化	1.4.10.1
2.10.2 Meta信息维护	1.4.10.2
2.10.3 一致性检测	1.4.10.3
2.10.4 View Meta	1.4.10.4
2.11 统计管理	1.4.11
2.11.1 查询条件统计	1.4.11.1
2.11.2 表状态统计	1.4.11.2
2.11.3 用户状态统计	1.4.11.3
2.11.4 命令统计	1.4.11.4
2.11.5 heartbeat统计	1.4.11.5
2.11.6 网络读写统计	1.4.11.6
2.12 故障切换	1.4.12
2.13 前后端连接检查	1.4.13
2.14 ER表	1.4.14
2.15 global表	1.4.15

2.16 缓存的使用	1.4.16
2.17 执行计划	1.4.17
2.18 性能观测和调整	1.4.18
2.19 智能计算reload	1.4.19
2.20 慢查询日志	1.4.20
2.21 单条SQL性能trace	1.4.21
2.22 KILL @@DDL_LOCK	1.4.22
2.23 外部高可用联动	1.4.23
2.23.1 外部后端MYSQL-HA连接	1.4.23.1
2.23.2 命令的使用说明	1.4.23.2
2.23.3 命令的实现细节	1.4.23.3
2.23.4 简单的HA交互使用案例	1.4.23.4
2.24 超时控制	1.4.24
2.25 流量控制	1.4.25
3.语法兼容	1.5
3.1 DDL	1.5.1
3.1.1 DDL&Table Syntax	1.5.1.1
3.1.2 DDL&View Syntax	1.5.1.2
3.1.3 DDL&Index Syntax	1.5.1.3
3.1.4 DDL透传	1.5.1.4
3.1.5 DDL&Database_Syntax	1.5.1.5
3.2 DML	1.5.2
3.2.1 INSERT	1.5.2.1
3.2.2 REPLACE	1.5.2.2
3.2.3 DELETE	1.5.2.3
3.2.4 UPDATE	1.5.2.4
3.2.5 SELECT	1.5.2.5
3.2.6 SELECT JOIN syntax	1.5.2.6
3.2.7 SELECT UNION Syntax	1.5.2.7
3.2.8 SELECT Subquery Syntax	1.5.2.8
3.2.9 LOAD DATA	1.5.2.9
3.2.10 不支持的语句	1.5.2.10
3.3 Prepared SQL Syntax	1.5.3
3.4 Transactional, Savepoint and Locking Statements	1.5.4
3.4.1 Lock&unlock	1.5.4.1
3.4.2 XA 事务语法	1.5.4.2
3.4.3 一般事务语法	1.5.4.3
3.4.4 SET TRANSACTION Syntax	1.5.4.4
3.4.5 SAVEPOINT, ROLLBACK TO SAVEPOINT, and RELEASE SAVEPOINT Syntax	1.5.4.5
3.5 DAL	1.5.5
3.5.1 SET	1.5.5.1
3.5.2 SHOW	1.5.5.2
3.5.3 KILL	1.5.5.3
3.6 存储过程支持方式	1.5.6
3.7 Utility Statements	1.5.7
3.8 Hint	1.5.8
3.9 其他不支持语句	1.5.9
3.10 函数与操作符支持列表(alpha版本)	1.5.10
3.11 导入导出方式	1.5.11
4.协议兼容	1.6
4.1 基本包	1.6.1
4.2 连接建立	1.6.2
4.3 文本协议	1.6.3
4.4 二进制协议 (Prepared Statements)	1.6.4
4.5 服务响应包	1.6.5
5.已知限制	1.7
5.1 druid引发的限制	1.7.1
5.2 其他已知限制	1.7.2
6.与MySQL Server的差异化描述	1.8
6.1 事务中遇到主键冲突需要显式回滚	1.8.1
6.2 INSERT不能显式指定自增序列	1.8.2
6.3 增加"show all tables"	1.8.3
6.4 去除了增删改的message信息	1.8.4
6.5 information_schema等库的支持	1.8.5
7.开发者须知	1.9
7.1 SQL开发编写原则	1.9.1
7.2 dble连接Demo	1.9.2
7.3 其他注意事项	1.9.3
8.配置示例	1.10
8.1 时间戳方式全局序列的配置	1.10.1
8.2 MySQL-offset-step 方式全局序列的配置	1.10.2
9.sysbench压测dble示例	1.11
9.1 测试环境及架构	1.11.1
9.2 修改dble配置	1.11.2
9.3 使用sysbench进行压测	1.11.3
A.Faq	1.12
A.1 ErrorCode	1.12.1
max Connections	1.12.1.1
Out Of Memory Error	1.12.1.2
The Problem Of Hint	1.12.1.3
NestLoop Parameters Lead To Temptable Exception	1.12.1.4
Can't Get Variables From DataNode	1.12.1.5
Port already in use:1984	1.12.1.6
Sharding Column Cannot Be Null	1.12.1.7
A.2 原理解释	1.12.2

How To Use Explain To Resolve The Distribution Rules Of Group Gy	1.12.2.1
Hash And ConsistentHashing And Jumpstringhash	1.12.2.2
A.3 使用说明	1.12.3
ToBeContinued2	1.12.3.1

## dble 中文技术参考手册

### 目录

参考 [gitbook](#) 左侧目录区 或 [SUMMARY.md](#)

### PDF下载

[PDF下载](#)

### 中文公开课

[dble中文公开课](#)

### 官方技术支持:

- 代码库 [github: github.com/actiontech/dble](#)
- 自动化测试库 [github: github.com/actiontech/dble-test-suite](#)
- 文档库 [github: github.com/actiontech/dble-docs-cn](#)
- 文档库 [github pages: actiontech.github.io/dble-docs-cn](#)
- 网站: [DBLE官方网站](#)
- QQ group: 669663113
- 开源社区微信公众号



### 注意

本分支上的手册适用于2.20.04.x版dble，其他版本的文档请参考对应tag分支或者release版文档。

### 提示

[如果您使用了dble，请告诉我们。](#)

### 联系我们

如果想获得dble 的商业支持，您可以联系我们：

- 全国支持: 400-820-6580
- 华北地区: 86-13718877200, 王先生
- 华南地区: 86-18503063188, 曹先生
- 华东地区: 86-18930110869, 梁先生
- 西南地区: 86-13540040119, 洪先生

## 0 概览

- [0.1 dble 简介与整体架构](#)
- [0.2 dble对MyCat做的增强](#)
- [0.3 快速开始](#)
- [0.4 docker镜像快速开始](#)
- [0.5 docker-compose快速开始](#)
- [0.6 数据拆分简介](#)

## 0.1 dble 简介与整体架构

### 0.1.1 dble简介

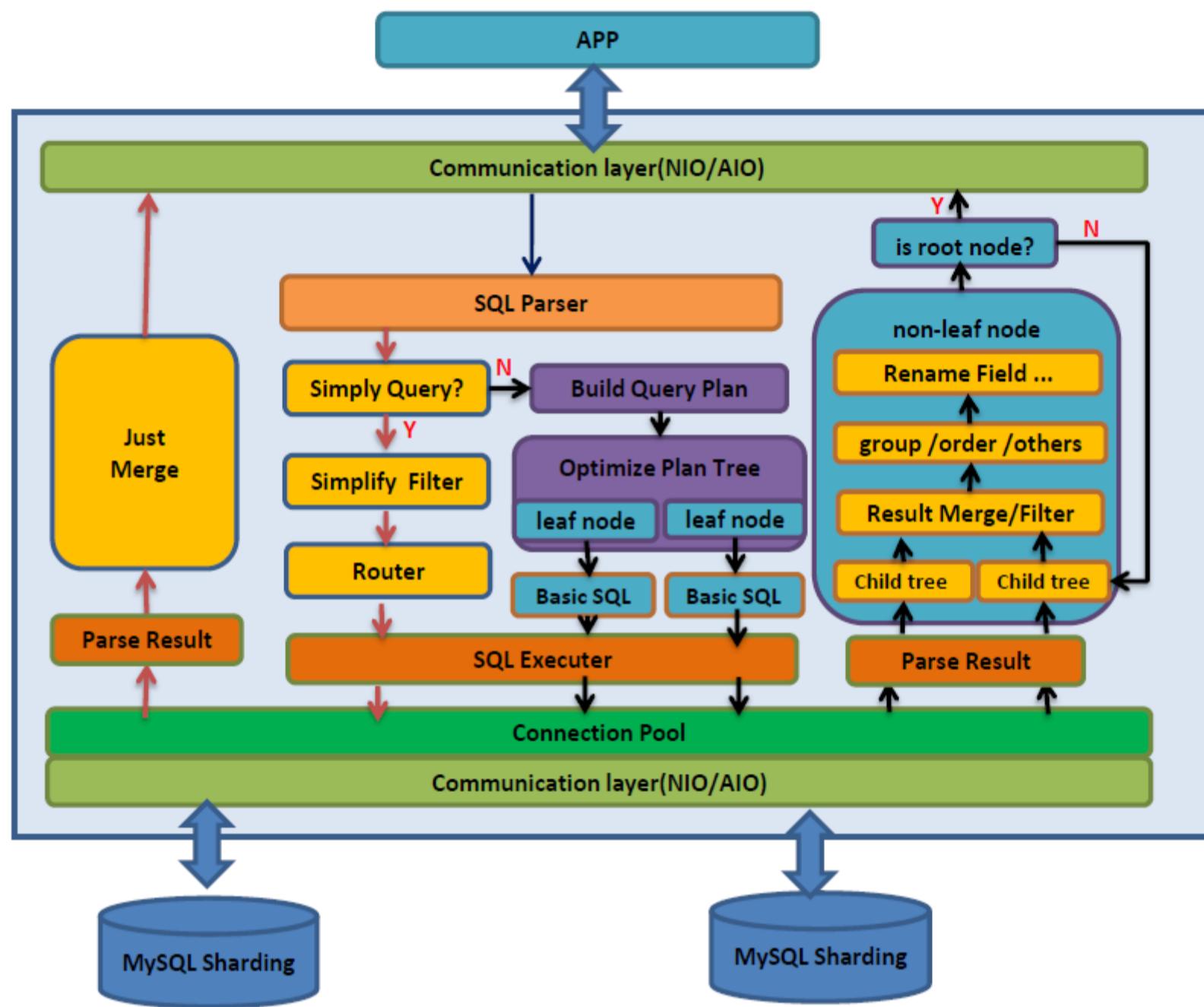
dble是[上海爱可生信息技术股份有限公司](#)基于mysql的高可扩展性的分布式中间件，存在以下几个优势特性：

- **数据水平拆分** 随着业务的发展，您可以使用dble来替换原始的单个MySQL实例。
- **兼容MySQL** 与MySQL协议兼容，在大多数情况下，您可以用它替换MySQL来为您的应用程序提供新的存储，而无需更改任何代码。
- **高可用性** dble服务器可以用作集群，业务不会受到单节点故障的影响。
- **SQL支持** 支持SQL 92标准和MySQL方言。我们支持复杂的SQL查询，如group by, order by, distinct, join, union, sub-query等等。
- **复杂查询优化** 优化复杂查询，包括但不限于全局表连接分片表，ER关系表，子查询，简化选择项等。
- **分布式事务支持** 使用两阶段提交的分布式事务。您可以为了性能选择普通模式或者为了数据安全采用XA模式。当然，XA模式依赖于MySQL-5.7的XA Transaction，MySQL节点的高可用性和数据的可靠性。

### 0.1.2 dble由来

- dble 是基于开源项目MyCat的，在此对于MyCat的贡献者们致以由衷的感谢。
- 对我们来说，专注于MySQL是一个更好的选择。所以我们取消了对其他数据库的支持，对兼容性，复杂查询和分布式事务的行为进行了深入的改进/优化。当然，还修复了一些bugs。详情可见[dble对MyCat做的改进](#)

### 0.1.3 dble内部架构



## 0.2 dble对MyCat做的增强

建议大家收看[免费课程](#)来初步认识dble的改进项

### 0.2.1 缺陷修复（目前mycat社区不活跃，不再追踪mycat的bug）

- 由于对堆外内存的使用不当，导致高并发操作时对同一片内存可能发生“double free”，从而造成JVM异常，服务崩溃。#4
- XA事务漏洞：包乱序导致客户端崩溃 #21
- where关键字写错时，会忽视后面的where条件，会得到错误的结果，比如select \* from customer wher id=1;#126
- 对于一些隐式分布式事务，例如insert into table values(节点1), (节点2)；原生mycat直接下发，这样当某个节点错误时，会造成该SQL执行了一部分
- 权限黑名单针对同一条sql只在第一次生效。#92
- 聚合/排序的支持度非常有限，而且在很多场景下还存在结果不正确、执行异常等问题 #43, #31, #44
- 针对between A and B语法，hash拆分算法计算出来的范围有误#23
- 开启全局表一致性检查时，对全局表的处理存在诸多问题，例如不能alter table、insert...on duplicate...时不更新时间戳、update...in ()报错等#24, #25, #26 ,#5
- 多值插入时，全局序列生成重复值 #1
- ER表在一个事务内被隔离，不能正确插入子表数据#13
- sharding-join结果集不正确#17

详情及其他修正请见[修复列表](#)

### 0.2.2 实现改进

- 对某些标准SQL语法支持不够好的方面作了改进，例如对create table if not exists...、alter table add/drop [primary] key...等语法的支持
- 对整体内部IO结构进行了大幅的改造和调优:参见[dble的IO结构](#) 或者见本节结尾附录2
- 禁止普通用户连接管理端口进行管理操作，增强安全性 #56
- 对全局序列做了如下改进
  - 删去无工程意义的本地文件方式
  - 改进数据库方式、ZK方式，使获取的序列号更加准确
  - 改进时间戳方式和ZK ID生成器方式，消除并发低时的数据分布倾斜问题
  - 修复了数据库方式全局序列中线程安全的问题#489
  - 移除自定义语法限制：全局序列值不能显式指定
    - 原来：insert into table1(id,name) values(next value for MYCATSEQ\_GLOBAL,'test');
    - 现在1：insert into table1(name) values('test');
    - 现在2：insert into table1 values('test');
    - 注意时间戳方式需要该字段是bigint
- 改进对ER表的支持，智能处理连接隔离，解决同一事务内不可以同时写入父子表的问题，并优化ER表的执行计划
- 系统通过智能判断，对于一些没有显式配置但实际符合ER条件的表视作ER表同样处理
- 在中间件内进行智能解析与判断，使用正确的schema，替换有缺陷的checkSQLschema 参数
- conf/index\_to\_charset.properties的内容固化到代码。
- 对于前端按照用户限制连接数，限制总连接数
- 改进原本的SQL统计，增加UPDATE/DELETE/INSERT

### 0.2.3 功能增强

- 提供了更强大的查询解析树，取代ShareJoin，使跨节点的语法支持度更广(join,union,subquery)，执行效率更高，同时聚合/排序也有了较大改进
- 提供科学的元数据管理机制，更好的支持show、desc等管理命令，支持不指定columns的insert语句 #7
- 元数据自动检查
  - 启动时对元数据进行一致性检查
  - 配置定时任务，对元数据进行一致性检查
- 提供更详实的执行计划，更准确的反映SQL语句的执行过程
  - 举例：mysql> explain select \* from sharding\_two\_node a inner join sharding\_four\_node b on a.id =b.id;

```
+-----+-----+-----+
| DATA_NODE | TYPE   | SQL/REF
+-----+-----+-----+
| dn1.0    | BASE SQL | select `a`.`id`, `a`.`c_char`, `a`.`ts`, `a`.`si` from `sharding_two_node` `a` ORDER BY `a`.`id` ASC |
| dn2.0    | BASE SQL | select `a`.`id`, `a`.`c_char`, `a`.`ts`, `a`.`si` from `sharding_two_node` `a` ORDER BY `a`.`id` ASC |
| dn1.1    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| dn2.1    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| dn3.0    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| dn4.0    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| merge.1  | MERGE   | dn1.0, dn2.0 |
| merge.2  | MERGE   | dn1.1, dn2.1, dn3.0, dn4.0 |
| join.1   | JOIN    | merge.1, merge.2 |
+-----+-----+-----+
9 rows in set (0.00 sec)
```

- set 系统变量语句的改进
- set charset/names 语句的支持
- [分布式事务](#):XA实现方式的异常处理的改进
- 大小写敏感支持
- 支持DUAL
- 支持单次请求[多语句](#)(部分客户端有使用,C和C++常见)
- 不断丰富的路由规则优化和条件优化
- 升级Druid,跟进最新的解析器
- 升级fastjson,避免安全问题
- 智能判断[reload](#)连接变更，热更新连接池
- 对MySQL协议及GUI工具/Driver更友好的支持
- 增加更多的[管理端命令](#),满足更多运维需要
- 缓存支持使用RocksDB
- 增加[慢查询日志功能](#),兼容mysqldumpslow 和 pt-query-digest
- Trace功能，用于分析单条查询的性能瓶颈
- 大小写敏感依赖于后端MySQL
- 支持[Prepared SQL Statement Syntax](#)
- 支持以下子查询
  - The Subquery as Scalar Operand
  - Comparisons Using Subqueries
  - Subqueries with ANY, IN, or SOME
  - Subqueries with ALL
  - Subqueries with EXISTS or NOT EXISTS
  - Derived Tables (Subqueries in the FROM Clause)
- 支持dble层面的 View
- 支持MySQL8.0 默认登陆验证插件
- 提供[自定义告警接口](#)
- 支持[自定义拆分算法](#)
- 支持自增列可以设置为非主键列
- 可以观察执行中的DDL
- 提供配置[预检查功能](#)
- 流量暂停和恢复功能

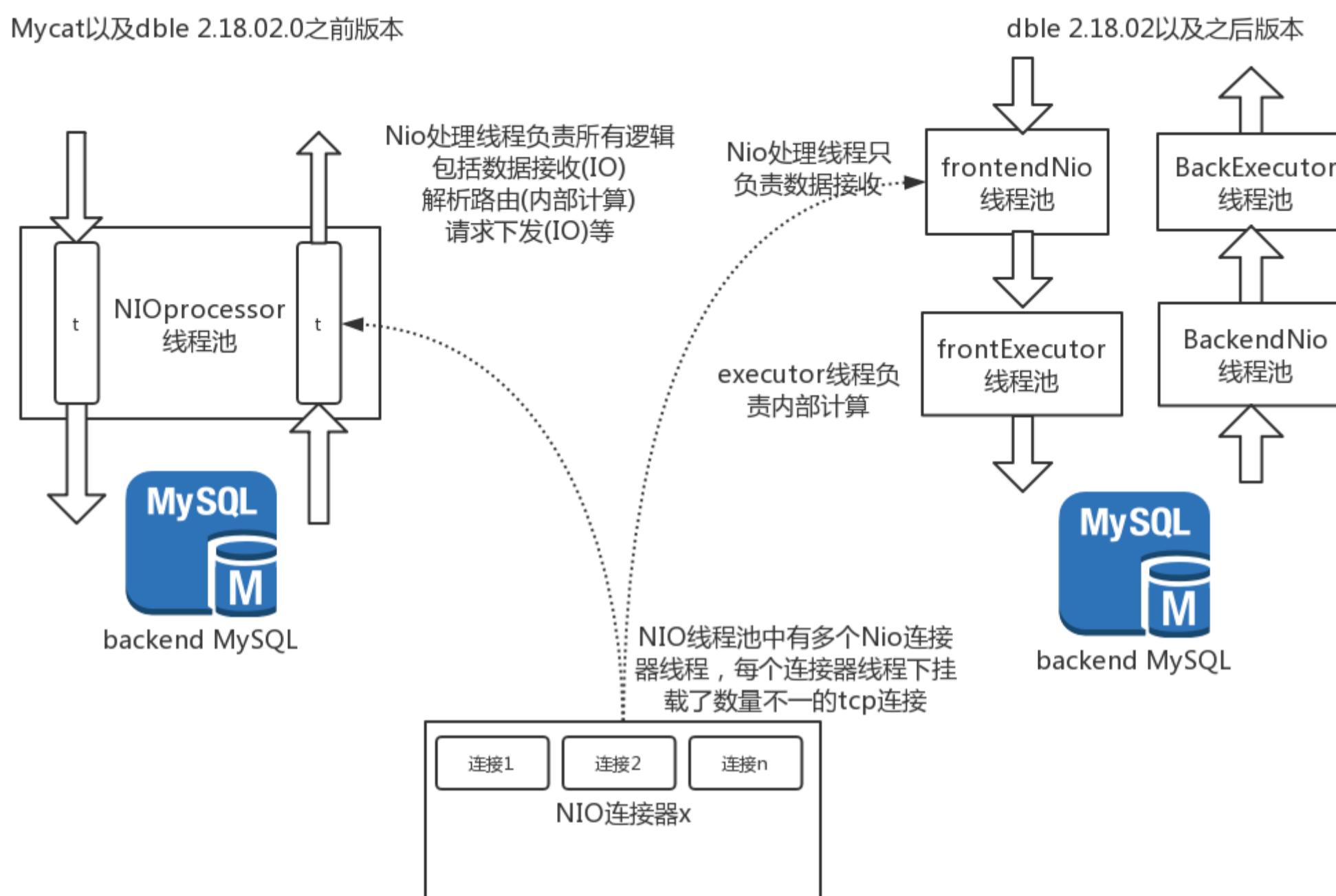
### 0.2.4 功能裁减

- 仅保留枚举、范围、HASH、日期等分片算法，对这几个算法进行了可用性的改进，使之更加贴合实际应用，项目需要时可以按需提供
- 移除异构数据库支持
- 禁止某些不支持的功能，这些功能客户端调用时不会报错，但结果并非用户想要的结果，例如无效的set语句
- 移除目前实现有问题的第一结点库内分表模式
- 移除writeType参数，等效于原来writeType = 1
- 移除handleDistributedTransactions 选项，直接支持分布式事务

### 0.2.5 附录

一些功能改进的详细描述，见开源数据库中间件DBLE对MyCAT的增强与改进

## 0.2.6 附录2



## 0.3 快速开始

### 0.3.1 关于本节

- 本节内容为您介绍如何使用dble安装包快速部署并启动一个dble服务，并简单了解dble的使用和管理

### 0.3.2 安装准备

以下部分将被需要作为dble启动的基础支撑

- 两个启动的MySQL实例  
dble是通过连接mysql数据库实例来进行数据的存储，所以请至少准备两个正在运行的mysql实例，假设您的机器上存在两个MySQL实例：  
A:\$url=ip1:3306,\$user=test,\$password=testPsw  
B:\$url=ip2:3306,\$user=test,\$password=testPsw  
请正确配置/etc/hosts，保证此MySQL实例可以正确访问，否则之后可能会报错 "NO ROUTE TO HOST"。  
并且在本例中需要在mysql实例1中新建几个数据库：

```
create database db_1;
create database db_3;
create database db_5;
```

并且在本例中需要在mysql实例2中新建几个数据库：

```
create database db_2;
create database db_4;
create database db_6;
```

- JVM环境  
dble是使用java开发的，所以需要启动dble您先需要在机器上安装java版本1.8或以上，并且确保JAVA\_HOME参数被正确的设置

### 0.3.3 下载并安装

- 通过此连接下载最新版本的安装包<https://github.com/actiontech/dble/releases>
- 解压并安装dble到指定文件夹中

```
mkdir -p $working_dir
cd $working_dir
tar -xvf actiontech-dble-$version.tar.gz
cd $working_dir/dble/conf
cp rule_template.xml rule.xml
cp schema_template.xml schema.xml
cp server_template.xml server.xml
```

### 0.3.4 dble的初始化配置

- 修改schema.xml，找到其中的writeHost将您对应的数据库信息进行替换

```
<writeHost host="hostM1" url="$url" user="$user" password="$password"/>
```

- 在本例中替换之后对应的部分因为

```
<writeHost host="hostM1" url="ip1:3306" user="test" password="testPsw"/>
```

另一个不赘述。

### 0.3.5 启动并连接

- 启动命令\$working\_dir/bin/dble start  
如果启动失败请使用此命令查看失败的详细原因 tail -f logs/wrapper.log
- 使用mysql客户端直接连接dble服务，默认密码123456 mysql -p -P8066 -h 127.0.0.1 -u root
- 您可以使用mysql一样的方式执行以下语句

```
use testdb;
drop table if exists tb_enum_sharding;
create table if not exists tb_enum_sharding (
    id int not null,
    code int not null,
    content varchar(250) not null,
    primary key(id)
)engine=innodb charset=utf8;
insert into tb_enum_sharding values(1,10000,'1'),(2,10010,'2'),(3,10000,'3'),(4,10010,'4');
```

- 您可以使用以下命令登录到dble的管理端口进行dble服务状态的查看和管理，默认密码654321  
mysql -p -P9066 -h 127.0.0.1 -u man1  
其中9066是管理端的端口，man1是server.xml中默认的管理用户，如果有变更，请自行修改  
详细的管理命令可以查询[2.1 管理命令文档](#)

## 0.4 快速开始(docker)

### 0.4.1 关于本节

- 本节内容为您介绍如何通过dockerhub上的dble镜像快速启动一个dble的demon

### 0.4.2 安装准备

- 安装docker
- 安装mysql连接工具，用于进行连接测试观察结果

### 0.4.3 安装过程

按照顺序依次执行以下docker命令：

```
docker network create -o "com.docker.network.bridge.name"="dble-net" --subnet 172.18.0.0/16 dble-net
docker run --name backend-mysql1 --ip 172.18.0.2 -e MYSQL_ROOT_PASSWORD=123456 -p 33061:3306 --network=dble-net -d mysql:5.7 --server-id=1
docker run --name backend-mysql2 --ip 172.18.0.3 -e MYSQL_ROOT_PASSWORD=123456 -p 33062:3306 --network=dble-net -d mysql:5.7 --server-id=2
sleep 30
docker run -d -i -t --name dble-server --ip 172.18.0.5 -p 8066:8066 -p 9066:9066 --network=dble-net actiontech/dble:latest
```

通过以上命令依次创建一个docker网络，两个分别映射到主机33061,33062的mysql服务，一个将服务端和管理端映射到主机8066和9066端口的服务。服务将在约一分钟之后被启动，这是由于为了进行快速的启动需要对于mysql和dble的配置进行一些初始化。

### 0.4.4 连接并使用

使用准备好的mysql连接工具连接主机的8066或者9066端口，在docker的默认配置中

8066 端口(服务端口能够执行SQL语句)的用户为 root/123456

9066 端口(管理端口能够执行管理语句)的用户为 man1/654321

此例子中准备了travelrecord、company、goods等表格并提前进行了表格创建，若需要连接更多的表格配置详情和使用方法

请在dble-server容器中查阅/opt/dble/conf/schema.xml文件

在虚拟机已经有安装mysql客户端的状态下，可以使用以下默认命令进行连接

```
#连接dble sql服务端口
mysql -P8066 -u root -p123456 -h 127.0.0.1
#连接dble 控制管理端口
mysql -P9066 -u man1 -p654321 -h 127.0.0.1
#连接后端mysql1
mysql -P33061 -u root -p123456 -h 127.0.0.1
#连接后端mysql2
mysql -P33062 -u root -p123456 -h 127.0.0.1
```

### 0.4.5 环境清理

使用完成或者进行环境重建的时候可以使用以下命令进行环境的清空

```
docker stop backend-mysql1
docker stop backend-mysql2
docker stop dble-server
docker rm backend-mysql1
docker rm backend-mysql2
docker rm dble-server
docker network rm dble-net
```

### 0.4.6 docker-compose 版本快速启动

docker-compose启动需要先从github项目下载对应的配置文件

```
wget https://raw.githubusercontent.com/actiontech/dble/master/docker-images/docker-compose.yml
```

通过使用docker-compose的配置脚本直接启动两个mysql镜像以及一个dble镜像，在配置文件存放目录执行

```
docker-compose up
```

同样的，默认状态下启动dble的状态和quick-start一致，dble-server容器会启动两个端口8066/9066对外提供服务。可以使用默认用户root/123456连接8066端口进行测试，同时两个后端mysql开端口于本机33061/33062，默认的用户也是root/123456

在使用或者测试完毕之后，在配置文件存放目录下使用以下指令方便回收对应资源

```
docker-compose stop
docker-compose rm
```

### 0.4.7 尝试使用本地配置启动docker-compose

#### 注意

本小结的内容需要用户在充分了解并掌握dble配置和结构的状态下进行，作为一种快速启动特定配置dble以供测试或调试使用。首次了解并使用dble的用户可以跳过此节。

本地配置启动dble的原理是通过volumes配置的映射将本地的配置目录映射到docker容器中，之后初始化的时候在初始化脚本中将本地配置目录中的文件复制到对应的dble/conf目录中，之后再进行初始化和dble启动。

首先需要在docker-compose.yml的最后一段dble-server容器配置中添加以下内容

```
volumes:
  - /opt/test/conf:/opt/self_conf
```

此命令将本地的/opt/test/conf目录映射到目标容器/opt/self\_conf目录中去，之后调整dble-server容器的启动命令

```
command: ["/opt/dble/bin/wait.sh", "backend-mysql1:3306", "--", "/opt/self_conf/docker_init_start.sh"]
```

将原本调用的/opt/dble/bin/docker\_init\_start.sh修改为本次初始化准备使用的self\_conf目录下的初始化脚本

这里对于本地启动的初始化脚本只给出如下一些建议

- 首先将需要修改的配置文件放置到/opt/dble/conf目录下
- 调用/opt/dble/bin/dble start启动dble服务
- 调用脚本/opt/dble/bin/wait-for-it.sh 脚本监听dble 8066 服务，等待dble服务启动完成
- 调用mysql命令对于dble的后端数据库以及初始数据进行初始化

## 0.5 快速开始(docker-compose)

### 0.5.1 关于本节

- 本节内容为你介绍如何快速使用dble的docker-compose文件来启动一个dble的quick start  
以及一个按照自定义的配置和sql脚本来启动dble quick start的用例

### 0.5.2 安装准备

- 安装docker
- 安装docker-compose
- 安装mysql连接工具，用于进行连接测试观察结果

### 0.5.3 安装过程

从dble项目中下载最新的docker-compose.yml文件  
<https://raw.githubusercontent.com/actiontech/dble/master/docker-images/docker-compose.yml>  
 使用docker-compose up命令直接启动dble-server, compose配置文件会从dockerhub拉取镜像并最终启动dble

### 0.5.4 连接并使用

使用准备好的mysql连接工具连接主机的8066或者9066端口，在docker的默认配置中  
 8066端口(服务端口能够执行SQL语句)的用户为root/123456  
 9066端口(管理端口能够执行管理语句)的用户为man1/654321  
 此例子中准备了travelrecord、company、goods等表格并提前进行了表格创建，若需要连接更多的表格配置详情和使用方法  
 请在dble-server容器中查阅/opt/dble/conf/schema.xml文件

### 0.5.5 使用自定义配置启动dble

在这里介绍下如何使用自定义的dble本地配置来启动docker-compose中的dble  
 首先，简单描述下默认docker-compose.yml执行的过程，dble-compose挨个启动容器，并且在最终启动dble-server的时候调用。存在于镜像actiontech/dble:latest目录下的/opt/dble/bin/wait-for-it.sh脚本等待指定的（默认为mysql1容器的mysql服务端口）TCP端口启动，待到指定TCP端口启动之后调用初始化脚本/opt/dble/bin/docker\_init\_start.sh对于dble-server这个容器进行初始化（替换配置文件，启动dble，执行sql文件）  
 下面是默认的docker\_init\_start.sh脚本

```
#!/bin/sh

echo "dble init&start in docker"

sh /opt/dble/bin/dble start
sh /opt/dble/bin/wait-for-it.sh 127.0.0.1:8066
mysql -P9066 -u man1 -h 127.0.0.1 -p654321 -e "create database @@dataNode ='dn1,dn2,dn3,dn4'"
mysql -P8066 -u root -h 127.0.0.1 -p123456 -e "source /opt/dble/conf/testdb.sql" testdb

echo "dble init finish"

/bin/bash
```

可以轻易的看出脚本的内容非常简单，启动dble->等待dble服务启动->执行dble管理命令create database在后端mysql数据库中创建database->执行初始化sql脚本在dble中创建表和插入初始化数据  
 所以当需要使用非默认的情况进行启动dble-server容器时需要以下几个基本步骤

- 需要按照以上的模块化步骤编写一个自己的初始化sh文件，然后照着需要的步骤对于dble进行启动和初始化
- 在dble-server的启动配置中添加volumes项，将linux本地的配置文件和初始化文件存放目录挂载至dble-server容器内部(注意dble本地放在/opt/dble下)
- 修改dble-server的启动命令，将初始化脚本修改为自定义的容器内部的初始化脚本

#### e.g.

docker-compose.yml修改如下部分

```
dble-server:
  image: actiontech/dble:latest
  container_name: dble-server
  hostname: dble-server
  privileged: true
  stdin_open: true
  tty: true
  volumes:
    - ./:/opt/init/
  command: ["/opt/dble/bin/wait-for-it.sh", "backend-mysql1:3306", "--", "/opt/init/customized_script.sh"]
  ports:
    - "8066:8066"
    - "9066:9066"
  depends_on:
    - "mysql1"
    - "mysql2"
  networks:
    net:
      ipv4_address: 172.18.0.5
```

对应本地目录./存在以下文件

```
schema.xml rule.xml server.xml init.sql customized_script.sh
```

脚本customized\_script.sh中的内容为

```
#!/bin/sh

echo "dble init&start in docker"

cp /opt/init/server.xml /opt/dble/conf/
cp /opt/init/schema.xml /opt/dble/conf/
cp /opt/init/rule.xml /opt/dble/conf/

sh /opt/dble/bin/dble start
sh /opt/dble/bin/wait-for-it.sh 127.0.0.1:8066
mysql -P9066 -u man1 -h 127.0.0.1 -p654321 -e "create database @@dataNode ='dn1,dn2,dn3,dn4'"
mysql -P8066 -u root -h 127.0.0.1 -p123456 -e "source /opt/init/init.sql" testdb

echo "dble init finish"

/bin/bash
```

### 0.5.6 环境清理

使用完成或者进行环境重建的时候使用docker-compose stop/down来进行环境的清理

## 0.6 数据拆分简介

### 0.6.1 数据拆分简介

- 在dble中将表格按照数据分片的大致方式将表格的归为三个种类
  - 全局表：设置为全局表的表格将会在每个mysql节点上存在一个实体，并且在每个表格实体里面存放的都是全量的数据，一般用作字典表之类的数据量小、和多表关联或者是作为数据字典的表格。
  - 拆分表：设置为拆分表的表格会根据具体选择的拆分算法类型将其中的数据按照一定的规则分别存放到不同的mysql节点中去，每个节点存放一部分的数据。一般用以存放超大数据量的业务类表格。通过对于业务类表格的水平数据拆分实现性能的优化。
  - 非拆分表：设置为非拆分表的表格会在指定的mysql节点上单独存在一个表格实体，在此表格中存放全量数据。一般用于存放数据压力不大的业务类表格，类似冷门功能的业务数据表。

### 0.6.2 规划拆分方案

- 系统开发之前，应该对业务特点进行深度分析，对数据规模进行较准确的评估，根据表的数据量、数据特点和表与表之间的关系，决定哪些表需要分片。数据节点的数量应该根据数据量和访问性能要求合理配置，过多的节点数量不仅浪费资源，而且增加运维复杂度，有时不仅不能提升性能，还会降低性能。
- 对于小表（数据量不大，如千万以下），尽量不要配置成拆分表。如果表比较独立，与其他的表基本上不进行join运算，可以作为非拆分表处理，性能不能满足要求的时候可以通过配置读写分离机制来提高性能。如果需要与拆分表进行join运算，可以配置成全局表。
- 关于每个实例的规格，选择的根据有两个因素：
  - 需要存储的最大磁盘空间，需要通过拆分算法计算，根据存储数据量最大的节点来计算，并预估未来2到3年的数据增长；
  - 估算一个实例需要的最大 QPS 和 TPS，要根据最慢的节点估算。
- QPS·TPS与实际的SQL语句、数据量、数据结构和数据节点的规格有关，根据经验来估计的话，很可能偏差较大。系统建设之前，应该配置同等或者接近的环境，进行针对性的性能测试，从而做出准确的判断。

### 0.6.3 数据拆分表格的配置方法

- 数据拆分配置包括节点配置和拆分规则配置。节点配置决定了数据的物理存储方式，由主机和节点组成，主机代表具体的数据库实例，节点代表实例中的数据库。拆分规则决定了数据在不同的分区上读写的规则，由拆分算法和应用拆分算法的逻辑库、表、字段组成。
- 数据写入时，系统对指定的字段值应用拆分算法得到目标节点，然后将数据写入目标节点。数据读取时，系统对查询条件应用拆分算法得到数据源节点，从源节点获取数据，在中间层进行结果合并之后返回请求方，对于不同的查询条件，源节点可能有一个或多个。

## 1.dble配置文件基础

- 配置文档列表以及相关对应功能
  - [rule.xml](#): 分片规则定义
  - [schema.xml](#): 具体分片定义，规定table和schema以及dataNode之间的关系，指定每个表格使用哪种类型的分片方法，定义每个dataNode的连接信息等
  - [server.xml](#): dble server相关参数定义，包括dble性能，定时任务，端口，用户配置等
  - [wrapper.conf](#): wrapper配置文件，配置有JVM虚拟机等参数
  - [log4j.xml](#): log4j2.xml，配置日志参数
  - [cache配置](#): 配置缓存参数
  - 全局序列：全局序列相关配置
  - [myid.properties](#): 集群相关配置
  - [自定义拆分算法](#)
  - [自定义告警](#)
  - [自定义全局表一致性检查](#)
- 重要日志及文件：
  - [/logs/wrapper.log](#): 启动日志，如果dble启动失败，将会有日志出现在这个文件中
  - [/logs/dble.log](#): dble日志，日志记录并反馈dble执行过程中的重要信息
- 配置文件变更记录：
  - [配置文件变更记录](#)

## 1.1 rule.xml

### 1.1.1 相关文件

这部分与rule.dtd和rule.xml相关。

在2.18.12.0之后的版本中rule.xml在文件头中添加了version属性，以供运维人员区分配置创建或修改的版本([xml配置version对照表](#)) version字段不匹配时，启动和dryrun会给出NOTICE的提示，但不会影响功能

rule.dtd定义解析规则，仅与开发相关。如有疑问，请参看[xml\\_dtd\\_intro](#)。

rule.xml定义实际用到的分区算法的配置，它包括如下两类信息：

#### A.分区规则定义

分区规则定义有如下形式：

- tableRule

配置名称	配置内容	可设多值	说明
name	规则名称	否	tableRule属性，该规则名将被schema.xml中表配置引用，必须在整个文件中唯一
rule	规则详情	否	rule子结构

- rule

配置名称	配置内容	可设多值	可选项/默认值
columns	拆分列名	否	
algorithm	执行的function名字	否	只能是function声明过的元素

举例

```
<tableRule name="auto-sharding-long">
    <rule>
        <columns>id</columns>
        <algorithm>rang-long</algorithm>
    </rule>
</tableRule>
```

#### B.分区算法定义

分区算法定义有如下形式：

name: 定义分区算法名，在分区规则定义中被引用。 class: 指定分区算法实现类。每一种分区算法要求的参数个数，类型各不相同，property name部分用于指定相应分区算法的参数。这部分请参考各分区算法描述。

- function

配置名称	配置内容	可设多值	说明
name	函数的名称	否	在分区规则定义中被引用
class	拆分算法	否	只能是Enum,NumberRange,Hash,StringHash,Date,PatternRange,jumpStringHash之一
property	根据具体的function代码示例的属性进行配置参数	是	

举例：

```
<function name="rang-long" class="com.actiontech.dble.route.function.AutoPartitionByLong">
    <property name="mapFile">auto-sharding-long.txt</property>
    ...
</function>
```

### 1.1.2 支持的分区算法

目前，已支持的分区算法有：hash, stringhash, enum, numberrange, patternrange, date, jumpstringhash。

#### 1.hash分区算法

function的 class属性设置为“hash”或者“com.actiontech.dble.route.function.PartitionByLong”的分区规则应用该算法。具体配置如下：

```
<function name="hashLong" class="hash">
    <property name="partitionCount">C1[, C2, ...Cn]</property>
    <property name="partitionLength">L1[, L2, ...Ln]</property>
</function>
```

partitionCount:指定分区的区间数，具体为 C1 [+C2 + ... + Cn]。

partitionLength:指定各区间长度，具体区间划分为 [0, L1], [L1, 2L1], ..., [(C1-1)L1, C1L1], [C1L1, C1L1+L2], [C1L1+L2, C1L1+2L2], ... 其中，每一个区间对应一个数据节点。

例如，配置F1：

```
<property name="partitionCount">2,3</property>
<property name="partitionLength">100,50</property>
```

将划分如下的分区：

[0 , 100) [100, 200) [200, 250) [250, 300) [300, 350)

再如，配置F2：

```
<property name="partitionCount">2</property>
<property name="partitionLength">1000</property>
```

将划分如下的分区： [0 , 1000) [1000, 2000)

根据具体配置，模的基数M有如下计算公式：C1L1 + ... + CnLn. 上面的例子中F1 的M值为350, F2的M值为2000。在进行分片查找时，将分区字段key和M值进行求模运算： value = key mod M 得到的value值再从区间分布中找到自己数据节点的序号。例如，当配置为F1, key = 805 时,value = 105,那么从5个区间内发现对应的数据节点的序号为1(从0开始)。

结点的个数N 记为 C1 [+C2 + ... + Cn]。上面的例子中F1 的N值为5, F2的N值为2。

注意事项：

1. M不能大于2880。2880的原因是这样的：2, 3, 4, 5, 6, 8, 9, 10, 12, 15, 16, 18, 20, 24, 30, 32, 36, 40, 45, 48, 60, 64, 72, 80, 90, 96, 120, 144, 160, 180, 192, 240, 288, 320, 360, 480, 576, 720, 960, 1440是2880的约数,这样预分片扩容方便。
2. N必须等于 schema.xml 中使用该分区算法的逻辑表的dataNode属性指定的DataNode数量之和，如dataNode="dn1,dn2,dn3,dn4"中，N必须等于4。
3. Cn和Ln的个数必须相等。
4. 分区字段必须为整型字段，如果是其他类型，要求值可转化为数字。
5. 当partitionLength为1时，hash分区算法退化为求模算法，M及N均为partitionCount的值。
6. NULL作为分片列的值的时候数据的结果恒落在0号节点（第一个节点上）

#### 2.stringhash分区算法

class属性设置为“stringhash”或者“com.actiontech.dble.route.function.PartitionByString”的分区规则应用该算法。具体配置如下：

```
<function name="hashString" class="stringhash">
    <property name="partitionCount">c1[,c2, ...cn]</property>
    <property name="partitionLength">l1[,l2, ...ln]</property>
    <property name="hashSlice">l:r</property>
</function>
```

**partitionCount, partitionLength**的具体意义参看hash分区算法。 **hashSlice**: 指定参与hash值计算的key的子串。字符串从0开始索引计数。

子串截取有如下计算步骤和格式:

步骤1. 子串索引区间确定

格式	条件	区间
n	$n \geq 0$	(0, n)
n	$n < 0$	(n, 0)
:r		(0, r)
:l:		(l, 0)
:		(0:0)
l:r		(l, r)

步骤2. 子串索引区间修订

a. 左边界l修订

第一次修订			第二次修订	
修订前	条件	修订后	条件	修订后
	$l \geq 0$			
	$l < 0$	$l = l + length$	$l < 0$	$l = 0$

b. 右边界r修订

第一次修订			第二次修订	
修订前	条件	修订后	条件	修订后
	$r > 0$		$r > length$	$r = length$
	$r \leq 0$	$r = r + length$		

注: **length**是分区字段实际串的长度。

步骤3. 结果确定

条件	结果子串	hash值
$l < r$	索引在[l, r)范围的子串	通过子串计算的hash值
$l \geq r$	空字符串	0

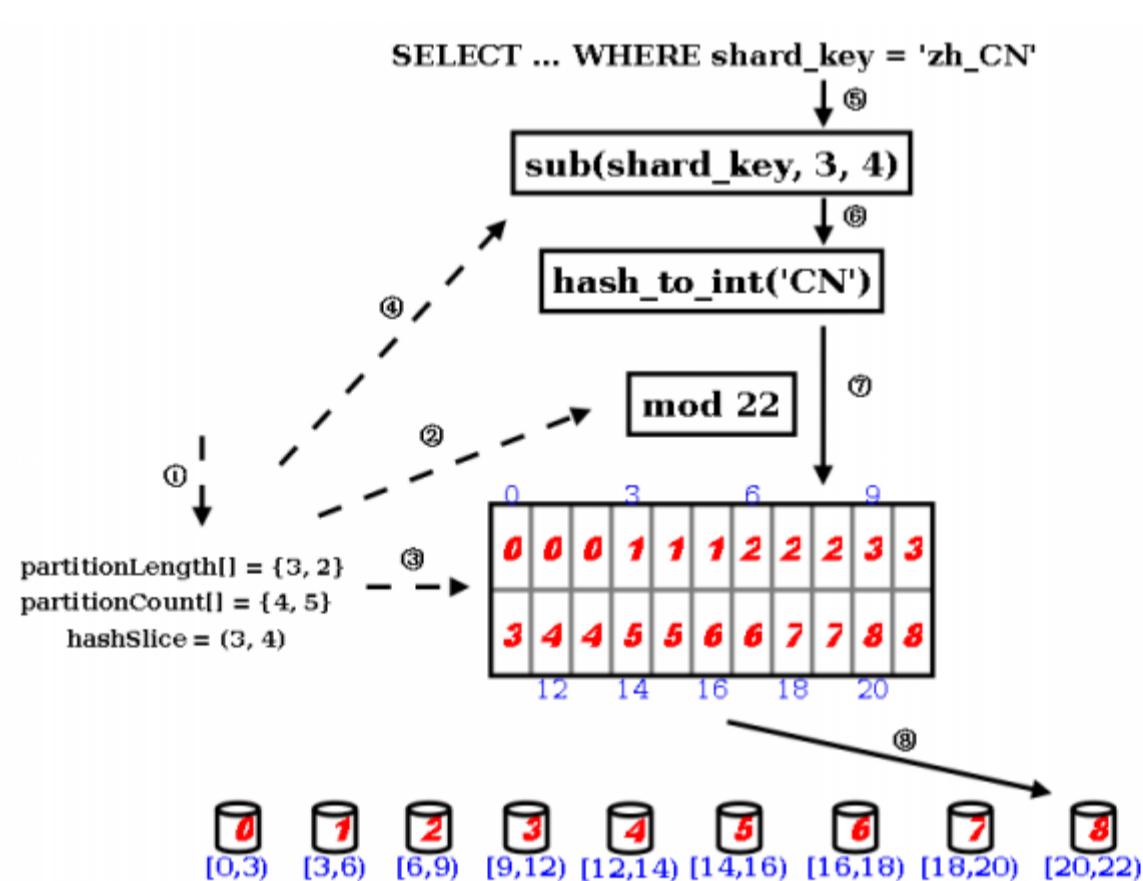
这个子串截取算法只是看起来比较复杂。用简洁但不太准确的语言可描述为: 区间边界为负值的, 从分区字段串尾部开始计数; 区间边界为正值的, 从分区字段串首部开始计数; 然后做一些纠错处理。

该算法与hash分区算法类似, 但针对的分区字段为字符串类型。在进行分片查找时, 要先对分区字段求hash值, 然后将得到的hash值做hash分区算法的分区查找运算。

注意事项:

1. 该分区算法和hash分区算法有同样的限制(注意事项3除外)
2. 分区字段为字符串类型。

算法解析:  
stringhash 按照用户定义的起点和终点去截取分片索引字段中的部分字符, 根据当中每个字符的二进制 unicode 值换算出一个长整型数值, 然后就直接调用内置 hash 算法求解分片路由: 先求模得到逻辑分片号, 再根据逻辑分片号直接映射到物理分片。



- 用户需要在 rule.xml 中定义 partitionLength[] 和 partitionCount[] 两个数组和 hashSlice 二元组。
- 在 DBLE 的启动阶段, 点乘两个数组得到模数, 也是逻辑分片的数量。
- 并且根据两个数组的叉乘, 得到各个逻辑分片到物理分片的映射表 (物理分片数量由 partitionCount[] 数组的元素值之和)。
- 此外根据 hashSlice 二元组, 约定把分片索引值中的第 4 字符到第 5 字符 (字符串以 0 开始编号, 编号 3 到编号 4 等于第 4 字符到第 5 字符) 字符串用于“字符串->整型”的转换。
- 在 DBLE 的运行过程中, 用户访问使用这个算法的表时, WHERE 子句中的分片索引值会被提取出来, 取当中的第 4 个字符到第 5 字符, 送入下一步。
- 设置一个初始值为 0 的累计值, 逐个取字符, 把累计值乘以 31, 再把这个字符的 unicode 值当成长整型加入到累计值中, 如此类推直至处理完截取出来的所有字符, 此时的累计值就能够代表用户的分片索引值, 完成了“字符串->整型”的转换。
- 对上一步的累计值进行求模, 得到逻辑分片号。
- 再根据逻辑分片号, 查映射表, 直接得到物理分片号。

### 3.enum分区算法

class属性设置为“enum”或者“com.actiontech.dble.route.function.PartitionByFileMap”的分区规则应用该算法。具体配置如下:

```
<function name="enum" class="enum">
    <property name="mapFile">partition.txt</property>
    <property name="defaultNode">-1</property>
    <property name="type">0</property>
</function>
```

**mapFile**:指定配置文件名。其格式将在下面做详细说明。 **defaultNode**: 指定默认节点号。默认值为-1, 不指定默认节点。 **type**: 指定配置文件中key的类型。0: 整型; 其它: 字符串。

配置文件格式如下: a. type值为0时, #comment //comment this line will be skiped int1=node0 int2=node1 ...

b. type值为非0时, #comment //comment this line will be skiped string1=node0 string2=node1 ...

在进行分片查找时，分片字段的枚举值对应的数据节点既是目的节点。如果不能在配置映射中找到枚举值对应的数据节点：如果配置了defaultNode，则defaultNode既是目的数据节点，否则出错。

注意事项：

1. 不包含“=”的行将被跳过。
2. node索引为数据节点索引号。
3. 重复的枚举值的分区数据节点以最后一个配置为准。
4. 分片字段为该枚举类型。
5. 分片字段为NULL时，数据落在defaultNode节点上，若此时defaultNode没有配置，则会报错；当真实存在于mysql的字段值为not null的时候，报错 "Sharding column can't be null when the table in MySQL column is not null"

#### 4. numberrange分区算法

class属性设置为“numberrange”或者“com.actiontech.dble.route.function.AutoPartitionByLong”的分区规则应用该算法。具体配置如下：

```
<function name="rangeLong" class="numberrange">
    <property name="mapFile">partition.txt</property>
    <property name="defaultNode">0</property>
</function>
```

**mapFile**:指定配置文件名。其格式将在下面做详细说明。**defaultNode**: 指定默认节点号。默认值为-1，不指定默认节点。

配置文件格式如下：#comment //comment this line will be skiped start1-end1=node1 start2-end2=node2 ...

该分区算法相当于定义了区间序列 [start1, end1], [start2, end2], ... 每一个区间对应一个数据节点。在进行分片查找时，分片字段的值落在的区间对应的数据节点即是目的节点。如果不能在配置映射中找到分片字段的值所落的区间时，分两种情况：1.配置了defaultNode，则defaultNode是目的数据节点；2.没配defaultNode，出错。注意事项：

1. 不包含“=”的行将被跳过。
2. node索引为数据节点索引号。
3. 如果区间存在重合，在对重合部分的分片字段值进行分片查找时在配置文件中最先定义的区间对应的数据节点为目的节点。
4. 分片字段为整型。
5. 分片字段为NULL时，数据落在defaultNode节点上，若此时defaultNode没有配置，则会报错；当真实存在于mysql的字段值为not null的时候，报错 "Sharding column can't be null when the table in MySQL column is not null"

#### 5. patternrange分区算法

class属性设置为“patternrange”或者“com.actiontech.dble.route.function.PartitionByPattern”的分区规则应用该算法。具体配置如下：

```
<function name="pattern" class="patternrange">
    <property name="mapFile">partition.txt</property>
    <property name="patternValue">1024</property>
    <property name="defaultNode">0</property>
</function>
```

**mapFile**:指定配置文件名。其格式将在下面做详细说明。**patternValue**:指定模值，默认为1024。**defaultNode**: 指定默认节点号。默认值为-1，不指定默认节点。

配置文件格式如下：#comment //comment this line will be skiped start1-end1=node1 start1-end2=node2 ...

该分区算法类似于numberrange分区算法。但在进行分片查找时先对分片字段值进行模patternValue运算，然后将得到的模数进行等同于numberrange分区算法的分区查找。注意事项：

1. 不包含“=”的行将被跳过。
2. node索引为数据节点索引号。
3. 如果区间存在重合，在对重合部分的分片字段值进行分片查找时在配置文件中最先定义的区间对应的数据节点为目的节点。
4. 分片字段的内容必须可以转化为整数。如果不能转化为整数：如果配置了defaultNode，目的数据节点为defaultNode；否则，出错。
5. 分片字段为NULL时，数据落在defaultNode节点上，若此时defaultNode没有配置，则会报错；当真实存在于mysql的字段值为not null的时候，报错 "Sharding column can't be null when the table in MySQL column is not null"

#### 6.date分区算法

class属性设置为“date”或者“com.actiontech.dble.route.function.PartitionByDate”的分区规则应用该算法。具体配置如下：

```
<function name="partbydate" class="date">
    <property name="dateFormat">yyyy-MM-dd</property>
    <property name="sBeginDate">2015-01-01</property>
    [<property name="sEndDate">2015-01-31</property>]
    <property name="sPartitionDay">10</property>
    <property name="defaultNode">0</property>
</function>
```

**dateFormat**:指定日期的格式。**sBeginDate**: 指定日期的开始时间。**sEndDate**: 指定日期的结束时间。该性质可以不配置或配置为空("")。**sPartitionDay**: 指定分区的间隔，单位是天。**defaultNode**: 指定默认节点号。默认值为-1，不指定默认节点。

该算法有两种工作模式：**模式1**：不配置sEndDate或者将sEndDate配置为""在这种模式下，该算法将时间以sBeginDate为开始，以sPartitionDay为间隔，进行区间划分，每个区间对应一个数据节点。在进行区间查找时，分区字段落在的区间对应的数据节点即是目的数据节点。如果分区字段值小于sBeginDate：如果配置了defaultNode，则数据会被路由至defaultNode；如果没有配置defaultNode，则会报错。**模式2**：配置sEndDate且不为""在这种模式下，该算法将时间以sBeginDate为开始，以sPartitionDay为间隔，以sEndDate为终点，进行区间划分，划分为N个区间，每个区间对应一个数据节点。在进行区间查找时：如果分区字段值小于等于sEndDate，则计算过程和结果等同于模式1。如果分区字段值大于sEndDate，则分片查找公式为：index=((key - sBeginDate)/sPartitionDay)%N，其中key为分片字段值，index为映射到的数据节点索引。这等同于一个环形映射。如果分区字段值大于sBeginDate，则会检查是否设置了defaultNode。如果设置了，数据会路由至defaultNode；否则报错。

注意事项：

1. 分片字段必须是符合dateFormat的日期字符串。
2. 区间划分不以日历时间为基准，无法对应到日历时间。内部进行区间划分时，会将sPartitionDay转化为以86400000毫秒为一天进行运算。
3. 在模式2的情况下，如果(sEndDate - sBeginDate)不是sPartitionDay的整数倍，则索引为0的数据节点承载更多的数据。
4. 分片字段为NULL时，数据落在defaultNode节点上，若此时defaultNode没有配置，则会报错；当真实存在于mysql的字段值为not null的时候，报错 "Sharding column can't be null when the table in MySQL column is not null"

#### 7. 跳增字符串算法

class属性设置为“jumpstringhash”或者“com.actiontech.dble.route.function.PartitionByJumpConsistentHash”的分区规则应用该算法，具体配置如下

```
<function name="jumphash"
    class="jumpStringHash">
    <property name="partitionCount">2</property>
    <property name="hashSlice">0:2</property>
</function>
```

**partitionCount**:分片数量 **hashSlice**:分片截取长度 该算法来自于Google的一篇文章A Fast, Minimal Memory, Consistent Hash Algorithm其核心思想是通过概率分布的方法将一个hash值在每个节点分布的概率变成 $1/n$ ，并且可以通过更简便的方法可以计算得出，而且分布也更加均匀

注意事项：

1. 分片字段值为NULL时，数据恒落在0号节点之上；当真实存在于mysql的字段值为not null的时候，报错 "Sharding column can't be null when the table in MySQL column is not null"

### 1.1.3 完整配置举例

```

<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE dble:rule SYSTEM "rule.dtd">
<db:rule xmlns:db="http://dbe.cloud/" version="9.9.9.9">
  <tableRule name="sharding-by-enum">
    <rule>
      <columns>id</columns>
      <algorithm>enum</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-range">
    <rule>
      <columns>id</columns>
      <algorithm>rangeLong</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-hash">
    <rule>
      <columns>id</columns>
      <algorithm>hashLong</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-hash2">
    <rule>
      <columns>id</columns>
      <algorithm>hashLong2</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-hash3">
    <rule>
      <columns>id</columns>
      <algorithm>hashLong3</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-mod">
    <rule>
      <columns>id</columns>
      <algorithm>hashmod</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-hash-str">
    <rule>
      <columns>id</columns>
      <algorithm>hashString</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-date">
    <rule>
      <columns>caldate</columns>
      <algorithm>partbydate</algorithm>
    </rule>
  </tableRule>

  <tableRule name="sharding-by-pattern">
    <rule>
      <columns>id</columns>
      <algorithm>pattern</algorithm>
    </rule>
  </tableRule>

  <!-- enum partition -->
  <function name="enum">
    <class>Enum</class>
    <property name="mapFile">partition-hash-int.txt</property>
    <property name="defaultNode"><!-- the default is -1,means unexpected value will report error-->
    <property name="type"><!--0 means key is a number, 1 means key is a string-->
  </function>

  <!-- number range partition -->
  <function name="rangeLong" class="NumberRange">
    <property name="mapfile">autopartition-long.txt</property>
    <property name="defaultNode"><!--he default is -1,means unexpected value will report error-->
  </function>

  <!-- Hash partition,when partitionLength=1, it is a mod partition-->
  <!--MAX(sum(count*length[i]) must not more than 2880-->
  <function name="hashLong" class="Hash">
    <property name="partitionCount">&lt;></property>
    <property name="partitionLength">128</property>
    <!-- <property name="partitionCount">2,3</property>
    <property name="partitionLength">4,5</property-->
  </function>

  <!-- Hash partition,when partitionLength=1, it is a mod partition-->
  <!--MAX(sum(count*length[i]) must not more than 2880-->
  <function name="hashLong2" class="Hash">
    <property name="partitionCount">&lt;></property>
    <property name="partitionLength">512</property>
    <!-- <property name="partitionCount">2,3</property>
    <property name="partitionLength">4,5</property-->
  </function>

  <!-- Hash partition,when partitionLength=1, it is a mod partition-->
  <!--MAX(sum(count*length[i]) must not more than 2880-->
  <function name="hashLong3" class="Hash">
    <property name="partitionCount">2,1</property>
    <property name="partitionLength">256,512</property>
    <!-- <property name="partitionCount">2,3</property>
    <property name="partitionLength">4,5</property-->
  </function>

  <!-- eg: mod 4 -->
  <function name="hashmod" class="Hash">
    <property name="partitionCount">4</property>
    <property name="partitionLength">1</property>
  </function>

  <!-- Hash partition for string-->
  <function name="hashString" class="StringHash">
    <property name="partitionCount">8</property>
    <property name="partitionLength">128</property>
    <property name="hashSlice">0:2</property>
    <!--<property name="hashSlice">4:0</property> -->
  </function>

  <!-- date partition -->
  <!-- 4 case:
  1.set sEndDate and defaultNode: input <sBeginDate ,router to defaultNode; input>sEndDate ,mod the period
  2.set sEndDate, but no defaultNode:input <sBeginDate report error; input>sEndDate ,mod the period
  3.set defaultNode without sEndDate: input <sBeginDate router to defaultNode;input>sBeginDate + (node size)*sPartitionDay-1 will report error(expected is defaultNode,but can't control now)
  4.sEndDate and defaultNode are all not set: input <sBeginDate report error;input>sBeginDate + (node size)*sPartitionDay-1 will report error
  -->
  <function name="partbydate" class="Date">
    <property name="dateFormat">yyyy-MM-dd</property>
    <property name="sBeginDate">2015-01-01</property>
    <property name="sEndDate">2015-01-31</property>
    <!--<property name="sEndDate">then in fact ,the sEndDate = sBeginDate+ (node size)*sPartitionDay-1 -->
    <property name="sPartitionDay">10</property>
    <property name="defaultNode"><!-- the default is -1-->
  </function>

  <!-- pattern partition -->
  <!--mapfile must contains all value of 0-patternValue-1,key and value must be Continuous increase-->
  <function name="pattern">
    <class>PatternRange</class>
    <property name="mapfile">partition-pattern.txt</property>
    <property name="patternValue">1024</property>
    <property name="defaultNode"><!--contains string which is not number,router to default node-->
  </function>

```

### 1.6.1 cache配置

```
</function>  
</dble:rule>
```

## 1.2 schema.xml

### 1.2.1 相关文件

这部分与schema.dtd和schema.xml相关。

在2.18.12.0之后的版本中server.xml在文件头中添加了version属性，以供运维人员区分配置创建或修改的版本([xml配置version对照表](#))

version字段不匹配时，启动和dryrun会给出NOTICE的提示，但不会影响功能

schema.dtd定义解析规则，仅与开发相关。如有疑问，请参看[xml\\_dtd\\_intro](#)。

schema.xml包含具体的数据主机配置，数据节点配置，数据库配置。注意：在2.20.04.0之后的版本中，我们对切换功能做了比较大的调整，不再支持自动switch，也不再支持一个datahost内部有多个writehost。

### 1.2.2 dataHost配置

- dataHost

配置名称	配置内容&示例	多节点	可选项/默认值	详细描述
name	节点名称	否	必需项，无默认值	dataHost的唯一标识，不允许重复
minCon	空闲时保有最小连接数	否	必需项，无默认值	1.在服务的启动后会根据dataNodeIdleCheckPeriod的时间周期注册的定时任务中会根据这个配置的值来创建/销毁连接，如非有大量请求，是缓慢式增长到最小值的。当空闲连接大于这个值，那么就会删除，如果空闲连接小于这个值并小于maxCon，那么就会创建连接到补齐这个数。注意：当datahost对应的schma的个数大于等于minCon时，建立的默认个数为schema个数+1(用于空schema)
maxCon	最大连接数,实际作用于每个子host	否	必需项，无默认值；当maxCon的数值低于修正后的minCon的值时，将maxCon修改成修正后的minCon，并且仅打印日志进行提示	host的连接总容量阈值 2.弹性伸缩时的最大边界
balance	读操作的负载均衡模式	否	必需项，无默认值，候选值0/1/2	在进行读负载均衡的时候会根据这个配置进行 0：不做均衡，直接分发到writeHost, readHost将被忽略，不会尝试建立连接池，但会有心跳连接 1：读操作在所有readHost中均衡 2：读操作在所有readHost和writeHost中均衡。具体拓扑结构见负载均衡相关章节
slaveThreshold	指定主从延迟阀值	否	默认-1，表示无延迟	1:在进行读取负载均衡的时候会根据最近一次的心跳状态以及读库和主库的延迟进行判断，如果延迟超过slaveThreshold配置，则认为此节点不适合进行读取，依赖于心跳为show slave status 2:此配置会影响到进行读负载均衡的时候延迟检测的开启，如果slaveThreshold=-1那么读负载均衡选取的时候不会进行延迟检测
tempReadHostAvailable	写库宕机之后读库是否可以提供服务	否	默认0,否	对于读请求来说，如果在进入读负载均衡分配的时候发现写库挂掉，就会通过配置进行判断，如果配置为1那么就会在依然存在读库中实现请求下发，否则服务不可用
heartbeat	子元素,心跳语句	否	必选项	该配置会在服务启动时设置的心跳任务里面被使用到，用于进行mysql实例状态的判断。 该配置有以下几种建议值： 1.普通心跳只是用于探活，建议使用select 1 2.使用 select @@read_only 探测结点可用性以及可写性 3.使用show slave status,可以探活，检查复制是否正常，以及延迟检测。如果Seconds_Behind_Master返回的状态有延迟，那么会被记入mysql实例的主从延迟中，影响读请求的路由分发，延迟超过指定限制读写分离会变为只读主库。
writeHost	子元素，表示写节点，每个datahost只能配置一个，配置见下面writeHost的描述	是	空	具体的物理节点配置

- heartbeat

配置名称	配置内容&示例	多节点	可选项/默认值	详细描述
timeout	heartbeat子元素,心跳超时阈值，单位：秒	否	默认0	心跳超时阈值。前置知识：dble会按照dataNodeHeartbeatPeriod的间隔向datasource发送心跳。心跳发起时候检查上次心跳是否为不正常的心跳，如果上次心跳尚未返回，并且距离最近的正常心跳的时间大于timeout，则标记该结点不可达。例如：心跳周期(dataNodeHeartbeatPeriod)2秒，第一次心跳正常，2s后的第二次心跳未返回，4s后第三次心跳发起时候发现上次不正常，不会真的再次下发，而是会根据4s和timeout的大小来确定该节点是否真正超时(该节点使用时候才会真正用到)。如果未超时，则什么也不做，继续下一个周期。如果超时了，则尝试杀掉超时的连接，无论是否杀成功，都会在下一个周期换一个连接继续做心跳。极端情况下会消耗很多连接。 2.心跳连续返回失败后，dble使用该结点时，会检查距离第一次失败的时间差，如果大于timeout，则报该结点不可达。
errorRetryCount	heartbeat子元素,心跳失败后的尝试次数	否	默认0，表示不重试	心跳失败后，开始重试errorRetryCount次。 1.重试期间成功，则标记回OK。 (目的：防止网络抖动或者连接异常断开场景) 2.重试期间超时，按照超时逻辑处理。

- writeHost

配置名称	配置内容&示例	多节点	可选项/默认值	详细描述
host	写节点名称	否	空	节点名称作为标识
url	写节点地址ip:port	否	空	被分成IP和PORT用于连接数据库
user	写节点用户	否	空	用于连接数据库
password	写节点用户密码	否	空	用于连接数据库
usingDecrypt	是否启用加密password	否	候选值0/1,默认值0	如果设置为1，password属性值应该为用工具encrypt.sh加密串 1:{host}:{user}:{password} 得到的串
weight	节点权重(负载均衡时候使用)	否	默认0	负载均衡过程中会查看所有节点的权重是否相等，如果不相等，那么就会根据权重来配置压力
readHost	子元素,对应读写分离读节点配置信息，详见readHost	是	空	用来关联读写节点关系

- readHost

配置名称	配置内容&示例	多节点	可选项/默认值	详细描述
host	写节点名称	否	空	节点名称作为标识
url	写节点地址ip:port	否	空	被分成IP和PORT用于连接数据库
user	写节点用户	否	空	用于连接数据库
password	写节点用户密码	否	空	用于连接数据库
usingDecrypt	是否启用加密password	否	候选值0/1,默认值0	如果设置为1，password属性值应该为用工具encrypt.sh加密串 1:{host}:{user}:{password} 得到的串
weight	节点权重(负载均衡时候使用)	否	默认0	负载均衡过程中会查看所有节点的权重是否相等，如果不相等，那么就会根据权重来配置压力

举例如下：

```
<dataHost name="localhost1" maxCon="1000" minCon="10" balance="2" slaveThreshold="100" [tempReadHostAvailable=""]>
    <heartbeat errorRetryCount="3" timeout ="5" >select user()</heartbeat>
    <writeHost host="hostM1" url="192.168.2.177:3307" user="root" password="root" usingDecrypt="" />
        <readHost host="hosts1" url="192.168.2.177:3309" user="root" password="root" weight="" usingDecrypt="" />
    ...
</writeHost>
...
</dataHost>
```

### 1.2.3 dataNode配置

- dataNode

配置名称	配置内容&示例	多节点	详细描述
name	数据节点名称, 唯一, 例如"dn,dn\$0-5"	否	作为数据节点的标识以及键, 节点个数的计算方法为: 从值出发, 以'.'(逗号) 分隔字符串, 如果其中有连续几项拥有相同的字符串前缀X(不能为空)并且后续其他几位为连续的数字时 (比如0到5), 可以以"X\$0-5"来省略表示, 个数为: 以逗号分隔的字符串个数加上包含\$的连续个数。name的个数必须等于database与dataHost的个数之积。
database	dataNode对应的存在于mysql物理实例中的schema, 可以配置单个或多个使用, 例如"db,db\$0-5"	否	所使用的详细数据库节点, 节点个数的计算方法为: 从值出发, 以'.'(逗号) 分隔字符串, 如果其中有连续几项拥有相同的字符串前缀X(不能为空)并且后续其他几位为连续的数字时 (比如0到5), 可以以"X\$0-5"来省略表示, 个数为: 以逗号分隔的字符串个数加上包含\$的连续个数。
dataHost	dataNode对应的数据库实例, 可以配置单个或多个使用, 例如"dh,dh\$0-5"	否	用于关联对应的Host节点, 节点个数的计算方法为: 从值出发, 以'.'(逗号) 分隔字符串, 如果其中有连续几项拥有相同的字符串前缀X(不能为空)并且后续其他几位为连续的数字时 (比如0到5), 可以以"X\$0-5"来省略表示, 个数为: 以逗号分隔的字符串个数加上包含\$的连续个数。

例如:

&lt;dataNode name="dn1" dataHost="localhost1" database="db1" /&gt;

name, datahost, database均可用如下格式在单个配置中配置多个节点: xxx\$n0-n1, xxx, 这种格式的意义为: xxxn0, ..., xxnm, ..., xxnn1, xxx, 其中

n0 &lt; nm &lt; n1.

例如: 配置

&lt;dataNode name="dn1\$0-19" dataHost="localhost1\$0-9" database="db1\$0-1" /&gt;

等同于:

&lt;dataNode name="dn10" dataHost="localhost10" database="db10" /&gt;
&lt;dataNode name="dn11" dataHost="localhost10" database="db11" /&gt;
&lt;dataNode name="dn12" dataHost="localhost11" database="db10" /&gt;
&lt;dataNode name="dn13" dataHost="localhost11" database="db11" /&gt;
...
&lt;dataNode name="dn119" dataHost="localhost19" database="db11" /&gt;

注意: 如果是使用通配符的配置, 那么dataNode(name)的通配符展开个数必须等于dataHost通配符展开个数与database通配符展开个数之积, 上例中, name的个数为20, dataHost的个数为10, database的个数为2; 又例如

&lt;dataNode name="dn,dn1\$0-19,dnx" dataHost="localhost,localhost1\$0-9" database="db1\$0-1" /&gt;

中, name的个数为22, dataHost的个数为11, database的个数为2。

注意: 自版本2.19.07.0开始, 若出现两个不同的dataNode拥有同样的database以及dataHost, 在配置检查的时候会报错(包括从通配符批量生成的dataNode)

## 1.2.4 schema配置

- schema

配置名称	配置内容&示例	多节点	可选项/默认值	详细描述
name	schema名称	否		schema的唯一标识, 不允许重复
dataNode	涉及的数据点	否	缺省无, 最多一个	未配置则只加载xml中的table子节点; 若配置后: 1.table子节点的配置会覆盖schema中的配置, 2.物理schema下存在的, 并且不在配置内的table被视为single node table
sqlMaxLimit	最大返回结果集限制	否	-1	当且仅当查询的SQL符合下列条件时, 产生效果 1 整个查询属于单表查询, 包括简单查询以及(order/group/聚合函数) 2 表格对应的schema需要有对应的sqlMaxLimit配置 3 表格的needAddLimit配置选项不能为false 4 表格的查询过程中不含有缓存键(配置cacheKey)过滤条件
table	详见table配置选项	是		每个表格的详细配置信息

- table

配置名称	配置内容&示例	多节点	可选项/默认值	详细描述
name	表格名称	否	必须项	表名, 可以配置多个使用';'分割
cacheKey	表格缓存唯一键	否	默认空	主键缓存时使用, 请确保业务上的key值唯一性
incrementColumn	表格自增列	否	默认空	指定表格自增列, 指定自增列的表格会有自增
needAddLimit	是否需要加返回结果集限制	否	true	与schemas 中的对应配置效果相同, 但是覆盖schema中配置
type	表格类型	否	默认default global	标记表格是全局表还是拆分表
dataNode	表格涉及的数据节点	否	空	两种格式: 1.xxx\$n0-n1 此种格式指定xxxn0, ..., xxnm, ..., xxnn1作为该表的数据节点。 2.distribute(xxx\$n0-n1)此种格式同样指定xxxn0, ..., xxnm, ..., xxnn1作为该表的数据节点, 但根据主机进行重排。 例如,,dn1关联host1, dn2关联host2, dn3关联host3, dn4关联host4, dn5关联host5, dn6关联host6, 如果用格式1,dataNode="dn1,dn2, dn3, dn4, dn5, dn6",其结果为: dn1, dn2, dn3, dn4, dn5, dn6, 如果用格式2, dataNode=distribute("dn1,dn2, dn3, dn4, dn5, dn6"),其结果为: dn1, dn3, dn5, dn2, dn4, dn6
rule	表格使用的分片规则	否	空	引用rule.xml中的拆分规则
globalCheck	是否启用全局表检查	否	false	true启用检查, false不检查
globalCheckClass	全局表检查类	否	CHECKSUM	全局表检查自定义类名或者是缩写 db自带CHECKSUM和COUNT两种默认实现
cron	全局表一致性检查周期	否	0 0 0 * * ?	quartz定时任务时间设置 详见: <a href="http://www.quartz-scheduler.org/api/2.4.0-SNAPSHOT/org/quartz/CronScheduleBuilder.html">http://www.quartz-scheduler.org/api/2.4.0-SNAPSHOT/org/quartz/CronScheduleBuilder.html</a>
ruleRequired	是否有绑定具体的rule	否	默认false	如果是false不需要进行存在规则检查, global表可为false
childTable	关联子表信息, 详见childTable选项	是	空	路由是通过父子关系进行ER关联

- childTable

配置名称	配置内容&示例	多节点	可选项/默认值	详细描述
name	表格名称	否	必需项	表名, 可以配置多个使用';'分割
cacheKey	表格缓存唯一键	否	默认空	主键缓存时使用, 请确保业务上的key值唯一性
incrementColumn	表格自增列	否	默认空	显式指定表格自增列
needAddLimit	是否需要加返回结果集限制	否	true	与schemas 中的对应配置效果相同, 但是覆盖schema中配置
joinKey	指定同父表进行join操作时的join键	否	必需项	子表和父表关联的字段
parentKey	指定进行join操作时父表中的join键	否	必需项	如果父表为非子表, 在父表中该字段必须与其拆分规则/拆分键有对等关系。
childTable	关联子表信息, 详见childTable选项	是	空	路由是通过父子关系进行ER关联

举例:

```
<schema name="TESTDB" [sqlMaxLimit="100"] [dataNode="dn1"]>
  <table name="payed" [cacheKey="id"] [autoIncrement="true"] [needAddLimit="true"] [type="global"] [rule="auto-sharding-long"] dataNode="dn1,dn2"/>
  ...
  <table name="customer" [cacheKey="id"] [incrementColumn="id"] [autoIncrement="true"] [needAddLimit="true"] [type="global"] [rule="auto-sharding-long"] dataNode="dn1,dn2">
    <childTable name="orders" [cacheKey="id"] [autoIncrement="true"] [needAddLimit="true"] joinKey="customer_id" parentKey="id">
      <childTable ...>
      ...
    </childTable>
  </childTable>
  ...
</table>
...
</schema>
```

### 1.2.5 注意事项

标签dataHost的属性dbDriver, dbTyp, writeType已废弃。  
标签schema的属性checkSQLschema已废弃。  
一个物理表不能配置给多个逻辑表使用, 否则使用后结果不可知, 比如#862。

### 1.2.6 举例

下面是一个schema.xml的完整例子:

```
<?xml version="1.0"?>
<!DOCTYPE dble:schema SYSTEM "schema.dtd">
<db: schema xmlns:db="http://db.ble.cloud/" version="9.9.9.9">

<schema name="TESTDB">
    <!-- auto sharding by id (long) -->
    <table name="travelrecord" dataNode="dn1,dn2" rule="sharding-by-hash2"/>

    <!-- global table is auto cloned to all defined data nodes , so can join
        with any table whose sharding node is in the same data node -->
    <table name="company" cacheKey="ID" type="global" dataNode="dn1,dn2,dn3,dn4"/>
    <table name="goods" cacheKey="ID" type="global" dataNode="dn1,dn2"/>
    <!-- random sharding using mod sharding rule -->
    <table name="hotnews" cacheKey="ID" autoIncrement="true" dataNode="dn1,dn2,dn3,dn4" rule="sharding-by-mod"/>
    <table name="customer" cacheKey="ID" dataNode="dn1,dn2" rule="sharding-by-mod">
        <childTable name="orders" cacheKey="ID" joinKey="customer_id" parentKey="id">
            <childTable name="order_items" joinKey="order_id" parentKey="id"/>
        </childTable>
        <childTable name="customer_addr" cacheKey="ID" joinKey="customer_id" parentKey="id"/>
    </table>
</schema>
<!-- <dataNode name="dn1$0-743" dataHost="localhost1" database="db$0-743"/> -->
<dataNode name="dn1" dataHost="localhost1" database="db1"/>
<dataNode name="dn2" dataHost="localhost1" database="db2"/>
<dataNode name="dn3" dataHost="localhost1" database="db3"/>
<dataNode name="dn4" dataHost="localhost1" database="db4"/>
<dataHost name="localhost1" maxCon="1000" minCon="10" balance="0" slaveThreshold="100">
    <heartbeat>select user()</heartbeat>
    <!-- can have multi write hosts -->
    <writeHost host="hostM1" url="localhost:3306" user="root" password="123456">
        <!-- can have multi read hosts -->
        <readHost host="hostS2" url="192.168.1.200:3306" user="root" password="xxx"/>
    </writeHost>
    <writeHost host="hostS1" url="localhost:3316" user="root" password="123456"/>
    <!-- <writeHost host="hostM2" url="localhost:3316" user="root" password="123456"/> -->
</dataHost>
</db: schema>
```

### 1.2.7 MySQL用户权限说明

权限项目	作用
SELECT	数据查询权限
INSERT	新增数据权限
UPDATE	更新数据权限
DELETE	删除数据权限
FILE	load data等数据导出导入的权限
CREATE	建库(管理端), 建表, 建索引权限
DROP	删除表的权限
ALTER	变更表结构的权限
LOCK TABLES	lock tables的权限
ALTER ROUTINE	(hint)变更/删除存储过程的权限
CREATE ROUTINE	(hint)创建存储过程的权限
EXECUTE	(hint)执行存储过程的权限
INDEX	建立/删除索引权限
SUPER	用于KILL功能
SHOW DATABASES	部分GUI工具会使用INFORMATION_SCHEMA SCHEMATA表
PROCESS	用于管理端show processlist 功能
REPLICATION CLIENT	可能场景: 1.用于配置了主从关系的数据结点, 使用了读写分离 2.使用show slave status作为心跳 3.需要使用show @@binlog_status功能
REFERENCES	外键约束(纯语法支持, 无使用意义)

## 1.3 server.xml 配置

### 整体XML结构

- system (目前reload不会做变更)
  - properties
- user (可多值,reload可变更)
- firewall (reload可变更)
  - whitehost
  - blacklist
  - host

### 配置选项

注意：“实例/全局属性”指的是如果集群配置多个中间件，是否要求配置在全局级别保持一致，或者实例间可以不同

#### version

在2.18.12.0之后的版本中server.xml在文件头中添加了version属性，以供运维人员区分配置创建或修改的版本([xml配置version对照表](#))  
version字段不匹配时，启动和dryrun会给出NOTICE的提示，但不影响功能

#### property

模块	配置名称	配置内容	默认值/单位	详细作用原理或应用	配置范围	范围	版本变更
后端连接socket配置	backSocketSoRcvbuf	后端套接字接收缓冲区大小	1024×1024×4,单位字节	在创建后端管道的时候作为buffer大小使用	正整数	实例	-
	backSocketSoSndbuf	后端套接字发送缓冲区大小	1024×1024,单位字节	在创建后端管道的时候作为buffer大小使用	正整数	实例	-
	backSocketNoDelay	后端Nagle算法是否禁用	默认1/单位无	在创建后端管道的时候禁用延迟加载，会影响网络包的情况 详见 <a href="#">相关资料</a>	1-是,0-否	实例	-
前端连接socket配置	frontSocketSoRcvbuf	前端套接字接受缓冲区大小	10241×1024,单位字节	在读取网络传输信息的时候作为每次缓冲的大小使用	正整数	实例	-
	frontSocketSoSndbuf	前端套接字发送缓冲区大小	1024×1024×4,单位字节	在创建前端管道的时候作为buffer大小使用	正整数	实例	-
	frontSocketNoDelay	前端Nagle算法是否禁用	默认1	在创建前端管道的时候禁用延迟加载 <a href="#">相关资料</a>	1-是,0-否	实例	-
Session预留内存配置	orderMemSize	session中的复杂查询order预留内存	默认4,单位M	在session初始化的时候创建内存分配对象，在复杂查询order by的时候使用到	正整数	实例	2.18.02.0引入
	otherMemSize	session中的复杂查询其他预留内存	默认4,单位M	在session初始化的时候创建内存分配对象，在复杂查询subQuery以及distinctd的时候使用	正整数	实例	2.18.02.0引入
	joinMemSize	session中的复杂查询join预留内存	默认4, 单位M	在session初始化的时候创建内存分配对象，在复杂查询join使用到	正整数	实例	2.18.02.0引入
堆外内存管理	bufferPoolChunkSize	内存池中分配的最小粒度	默认4096,单位字节	内存池中分配的最小粒度，需要的大小除以此粒度，向上取整		实例	-
	bufferPoolPageNumber	预分配内存池页数量	默认 0.8 × MaxDirectMemorySize / bufferPoolPageSize(default 2M), 向下取整	在初始化的时候通过和bufferPoolPageSize的相乘确定缓冲池最后的大小，内存配置建议见 <a href="#">1.4_wrapper.conf</a>		实例	-
	bufferPoolPageSize	预分配内存池页大小	默认1024 * 1024 * 2,单位字节	在初始化的时候通过和bufferPoolPageNumbe的相乘确定缓冲池最后的大小,  注意：虚拟机参数MaxDirectMemorySize( <a href="#">见1.4_wrapper.conf</a> )需要大于bufferPoolPageNumber * bufferPoolPageSize，否则会触发OOM		实例	-
	mappedFileSize	文件映射区单个文件最大体积	默认1024×1024×64,单位字节	在初始化的时候此参数确定文件映射区最大容量,参见内存管理章节		实例	2.17.04.0引入
	useOffHeapForMerge	处理跨分配结果合并是不是采用Direct Memory 暂时废弃，待重审	默认1	在初始化服务的时候会根据此配置初始化一个内存管理对象，并检查剩余内存的大小  并在执行中影响内存请求时请求到的内存的管理方式，在合并多节点返回数据使用		实例	2.17.11.0废弃
	memoryPageSize	每页内存的大小 暂时废弃，待重审	默认1m(写法带单位)	影响在程序中申请内存的最小单位，在多节点结果合并使用		实例	2.17.11.0废弃
	spillsFileBufferSize	写磁盘缓存大小 暂时废弃，待重审	默认 2K(写法带单位)	在合并结果集内存不够用的时候通过写磁盘来进行缓冲，  这个时候spillsFileBufferSize就是写磁盘流buffer的大小		实例	2.17.11.0废弃
	dataNodeSortedTempDir	写磁盘目录 暂时废弃，待重审	默认相对路径/sortDirs	在合并结果集内存不够用的时候写磁盘		实例	2.17.11.0废弃
统计管理	bufferUsagePercent	是否清理大结果集阈值	默认80,单位百分号	定时任务resultSetMapClear使用，周期clearBigSQLResultSetMapMs，定时清理统计的结果集，当定时任务执行时发现结果集统计超过阈值，触发清理结果集的行为	0-100	实例	-
	useSqlStat	是否启用SQL统计	默认1/单位无	启用之后会对于下发的查询进行SQL的统计，分别按照用户、表格、查询条件进行存放在内存中 并且开启之后会随之开启recycleSqlStat定时任务以固定5秒一次的周期回收SQL统计的结果	1-是0-否	实例	-
	clearBigSQLResultSetMapMs	定期大结果清理时间	默认600×1000,单位毫秒	定时任务resultSetMapClear的执行周期，定时清理记录的查询结果集,兼容往期错误拼写clearBigSQLResultSetMapMs	正整数	实例	-
	sqlRecordCount	慢查询记录阈值	默认10,单位条	在定时任务recycleSqlStat中会进行sql记录的清理,当发现记录的慢查询SQL数量超过阈值时，会仅保留阈值数个元素	正整数	实例	-

系统服务 基本参数	maxResultSet	大结果集阈值	默认512×1024,单位字节	当查询的结果集超过这个阈值时,查询的SQL和查询结果集的大小才会被记录到结果集统计里面	正整数	实例	-
	useCostTimeStat	是否启用查询耗时统计	默认0/单位无	开启之后以一定的比例统计查询过程中的各个步骤的耗时情况,可以使用BTraceCostTime.java进行观测,也可在管理端使用show @@cost_time观察	1-是0-否	实例	2.18.02.0 引入
	maxCostStatSize		默认100	show @@cost_time结果最近保留的行数		实例	2.18.02.0 引入
	costSamplePercent	查询采样百分比	默认1/单位%	在耗时采样统计中实际采样百分比为costSamplePercent		实例	2.18.02.0 引入
	useThreadUsageStat	开启线程使用率统计	默认0/单位无	开启之后能在管理端通过管理命令show @@thread_used查看各个部分的线程使用情况	1-是0-否	实例	2.18.02.0 引入
	bindIp	服务IP	默认 "0.0.0.0"	在服务初始化的时候作为侦听的IP	有效IP地址,推荐默认	实例	-
	serverPort	服务端口	默认8066	在服务初始化的时候作为服务侦听的端口	机器空闲端口	实例	-
	managerPort	控制端口	默认9066	在服务初始化的时候作为控制侦听的	机器空闲端口	实例	-
	maxCon	控制最大连接数	默认0	默认不做限制。若maxCon大于0,建立的连接数大于maxCon之后,建立连接会失败,注意当各个用户的maxcon总和值大于此值时,以当前值为准。全局maxCon不作用于manager用户	正整数	实例	2.18.09.0 引入
	processors	NIO前端处理器的数量	默认java虚拟机核数,单位个	进行前端网络IO吞吐的线程数	正整数	实例	-
功能性配置	backendProcessors	NIO后端处理器的数量	默认java虚拟机核数,单位个	进行后端网络IO吞吐的线程数	正整数	实例	2.18.02.0 引入
	fakeMySQLVersion	dble模拟mysql版本号	默认NULL	模拟成正常的MySQL版本在进行前端协议交互的时候能够使用到	MySQL版本号	最好全局,实例也可	-
	processorExecutor	前端业务处理线程池数量	默认(单核为2,否则等于宿主机核数)	进行前端具体业务处理的线程池大小,负责解析路由下发	正整数	实例	-
	backendProcessorExecutor	后端业务处理线程池数量	默认(单核为2,否则等于宿主机核数)	进行后端具体业务处理的线程池大小,负责回收结果集并合并	正整数	实例	2.18.02.0 引入
	complexExecutor	复杂查询后端业务线程池数量	默认(单核为2,否则等于宿主机核数,宿主机核数大于8时,数量为8)	负责复杂查询或者子命令结果集的回收	正整数	实例	2.18.02.0 引入
	sequenceHandlerType	全局序列处理器的方式	默认2,单位无	在初始化的时候根据这个配置选择不同的序列生成器进行加载  1. MySQL offset-step序列方式, sequence信息存储在数据库中  2. 时间戳方式(类Snowflake), 依赖sequence_time_conf.properties  3. 分布式time序列(类Snowflake)  4. 分布式offset-step序列  由于历史原因,错误的拼写sequenceHandlerType也可兼容	1, 2, 3, 4	全局	-
	serverNodeId	服务编号	默认1	在分布式事务的时候构造 XATXID, 同组中必须不同其形式为\$ServerName\$.serverNodeId.xid	正整数	实例	-
功能性配置	serverBacklog	前端tcp连接 backlog	默认2048	前端tcp连接 backlog	正整数	实例	2.17.04.0 引入
	showBinlogStatusTimeout	拉取一致性binlog线的超时时间	默认60000,单位毫秒	拉取一致性binlog线的超时时间	正整数	全局	2.17.08.0 引入
	usePerformanceMode	是否启用性能模式	默认0/单位无	开启之后Dble会大量占用CPU资源,并提供更高的性能体现,慎用	1-是0-否	实例	2.18.02.0 引入
	useCompression	是否启用数据压缩	默认0否	使用mysql压缩协议	1 - 是 0 - 否	全局	-
	usingAIO	是否启用AIO	默认0	在初始化服务的时候将会作为判断启用AIO或是NIO的依据	1 - 是 0 - 否	实例	-

	charset	字符集	utf8mb4	服务启动后的默认字符集于所有字符集相关的部分，包括前端连接和后端连接	有效字符集	全局	-
连接缺省值	maxPacketSize	包大小限制	默认 4×1024×1024	限制请求的包大小，启动时候dble会拉取并尝试同步（此值+1024）到每个datahost，如果同步失败，就取配置值与各个datahost中最小的那个值-1024.留出1024的冗余用于对SQL改写或者上下文同步的支持	正整数	全局	-
	txIsolation	隔离级别	默认 3	执行具体SQL的时候会比较前后端连接。 启动时候dble会拉取并尝试同步此值到每个datahost，如果同步失败或者session级别重新设置该值，session在SQL下发之前，会执行session级别的隔离级别set	1- READ_UNCOMMITTED 2-READ_COMMITTED 3-REPEATABLE_READ 4-SERIALIZABLE	全局	-
	autocommit	是否自动提交	默认 1, 自动提交	启动时候dble会拉取并尝试同步此值到每个datahost 如果同步失败或者session级别重新设置该值，执行具体SQL的时候会比较，如果不一致将会执行session级别的set	0/1	全局	2.19.11.0 引入
一致性检查	checkTableConsistency	表格一致性检查	默认0	如果值为1，那么在服务初始化的时候会启动一个定时任务，在定时任务会检查DB是不是存在，表格是不是存在，表结构是否一致	1-是,0-否	实例	-
	checkTableConsistencyPeriod	表格一致性检查周期	默认30×60×1000 ,单位毫秒	表格一致性检查周期	正整数	实例	-
心跳任务周期	dataNodeIdleCheckPeriod	后端空闲连接心跳周期	默认5×60×1000 ,单位毫秒	后端空闲连接心跳检查，超时关闭，调整容量	正整数	实例	-
	dataNodeHeartbeatPeriod	数据节点心跳任务周期	10×1000 ,单位毫秒	根据这个周期在服务初始化的时候注册心跳任务	正整数	实例	-
processor 内部前端连接检查	sqlExecuteTimeout	后端连接执行超时时间	默认 300 ,单位秒	如果超过这个时间没有完毕，就直接关闭连接	正整数	实例	-
	idleTimeout	(前端) 连接无响应超时时间	默认 30×60 × 1000 ,单位毫秒	在processor定时连接检查时，发现前端连接上一次的读写距今超过阀值，会直接关闭连接	正整数	实例	-
	processorCheckPeriod	processor定时任务检查周期	1000 ,单位毫秒	根据此配置定时的检查在processor中的前后端连接的状态	正整数	实例	-
普通事务日志相关	recordTxn	事务log记录	默认0	在初始化服务的时候会注册一个类，其作用就是将事务的log写到一个指定的文件中	1-是, 0-否	实例	2.17.04.0 引入
	transactionLogBaseDir	事务log目录	默认当前路径/tlogs	当开启日志log记录时，记录文件会被存放在对应目录下	绝对路径	实例	-
	transactionLogBaseName	事务log文件名称	默认server-tx	事务记录存储文件的文件名	符合运行系统文件的命名规范	实例	-
	transactionRotateSize	事务日志单个文件大小。	默认16 ,单位M	由于历史原因，兼容拼写transactionRateSize	正整数	实例	-
视图相关参数	viewPersistenceConfBaseDir	视图记录本地文件路径	dble目录/viewConf	用于存放视图本地记录文件的文件路径	绝对路径	集群配置时无意义	2.17.11.0 引入
	viewPersistenceConfBaseName	视图记录本地文件名	viewJson	视图记录的文件文件名	符合运行系统文件的命名规范	集群配置时无意义	2.17.11.0 引入
XA 事务	xaRecoveryLogBaseDir	xa的tm日志路径	dble目录/tmlogs/	此日志涉及到XA事务状态的记录，并且在Dble意外重启之后需要从里面获取重启之前的xa事务状态，切勿自行修改	绝对路径	实例	-
	xaRecoveryLogBaseName	xa的tm日志名称	tmlog		符合运行系统文件的命名规范	实例	-
	xaSessionCheckPeriod	XA定时任务执行周期	默认1000, 单位ms	在server开始的时候会注册一个定时任务以此参数为执行周期 (注：定时任务必定会被注册) 如果有尝试多次没有成功提交的session在之后的定时任务会被重复提交	正整数	实例	2.17.04.0 引入
	xaLogCleanPeriod	定时XAlog清除周期	默认1000, 单位ms	在server开始的时候会根据这个周期注册一个定时任务 (注：定时任务必定会被注册) 定时清XA log，主要是将已经回滚和提交成功的部分从记录中删除	正整数	实例	2.17.04.0 引入

	xaRetryCount	后台重试XA次数	默认0	后台定时任务重试XA次数，0为无限重试，达到设定次数后，停止重试	正整数	实例	2.19.03.0 引入
队列大小参数	joinQueueSize	join时,左右结点的暂存数据行数的队列大小	1024	当行数大于此值而又没有及时被消费者消费掉，将会阻塞，目的是防止接收数据量太大，堆积在内存中	正整数	实例	2.17.04.0 引入
	mergeQueueSize	merge时,左右结点的暂存数据行数的队列大小	1024	当行数大于此值而又没有及时被消费者消费掉，将会阻塞，目的是防止接收数据量太大，堆积在内存中	正整数	实例	2.17.04.0 引入
	orderByQueueSize	排序时,时,左右结点的暂存数据行数的队列大小	1024	当行数大于此值而又没有及时被消费者消费掉，将会阻塞，目的是防止接收数据量太大，堆积在内存中	正整数	实例	2.17.04.0 引入
使用Nest Loop优化	useJoinStrategy	是否使用nest loop优化	默认不使用	开启之后会尝试判断join两边的where来重新调整查询SQL下发的顺序	true 开启 false 不开启	实例	2.17.04.0 引入
	nestLoopConnSize	临时表阈值	默认4	若临时表行数大于这两个值乘积，则报告错误	正整数	实例	2.17.04.0 引入
	nestLoopRowsSize	临时表阈值	默认2000				
慢查询日志相关配置	enableSlowLog	慢查询日志开关	默认为0, 关闭	慢查询日志开关	0或者1	实例	2.18.09.0 引入
	slowLogBaseDir	慢查询日志存储文件夹	dble根目录/slowlogs	慢查询日志存储文件夹	文件夹路径	实例	2.18.09.0 引入
	slowLogBaseName	慢查询日志存储文件名前缀	slow-query	慢查询日志存储文件名前缀(后缀名是.log)	合法文件名	实例	2.18.09.0 引入
	flushSlowLogPeriod	日志刷盘周期, 单位秒	1	日志刷盘周期, 每隔这个周期, 会强制将内存数据刷入磁盘	正整数	实例	2.18.09.0 引入
	flushSlowLogSize	日志刷盘条数阈值	1000	日志刷盘条数阈值, 内存中每次写出这么多条日志, 会强制刷盘1次	正整数	实例	2.18.09.0 引入
	sqlSlowTime	慢日志时间阈值, 单位毫秒	100	慢日志时间阈值, 大于此时间的查询会记录下来	正整数	实例	2.18.09.0 引入
load data相关配置	maxCharsPerColumn	每列所允许最大字符数	默认为65535	每行所允许最大字符数	正整数	实例	2.19.03.0 引入
	maxRowSizeToFile	需要持久化的最大行数	默认为10000	当load data的数据行数超过阈值后, 会将数据保存在文件中以防OOM	正整数	实例	2.19.03.0 引入
高可用联动相关配置	useOuterHa	是否启用外部高可用联动, 如关闭此功能并且dble部署方式为单机, 将使用默认的切换方式, 详情请见切换相关章节	默认为true, 若此时不设置外部高可用, 将不做切换		true/false	实例	2.19.09.0 引入
流量控制相关参数	enableFlowControl	是否启用流量控制, true/false, 具体流量控制请参见相关功能描述章节	默认为false		true/false	实例	2.20.04.0 引入
	flowControlStartThreshold	流量控制触发队列长度阈值	默认为4096	当部分连接的写出队列超出阈值时触发流量控制	正整数	实例	2.20.04.0 引入
	flowControlStopThreshold	流量控制取消队列长度阈值	默认为256	当流量控制的连接写出队列长度小于阈值取消流量控制	正整数	实例	2.20.04.0 引入
废弃配置	processorBufferPoolType	缓冲池类型	默认0	在服务初始化的时候根据这个值确定缓冲池的类型  0, 直接缓冲块 1, 6个队列缓冲  (实际1未完成, 删除)			2.17.04.0 废弃
	useStreamOutput	启用结果集流输出, 不经过合并模块	默认 0	有bug, 无应用场景, 废弃			2.17.04.0 废弃
	systemReserveMemorySize	系统预留内存	384M	在启用useOffHeapForMerge配置的时候如果系统剩余的内存小于这个配置, 服务将不会启动 (仅在配置useOffHeapForMerge的时候有效)			2.17.04.0 废弃
	sqlInterceptor	SQL拦截器	默认 DefaultSqlInterceptor	拦截SQL, 进行统计等处理, 默认实现极影响性能, 废弃			2.17.04.0 废弃
	catletClassCheckSeconds	类加载重载检查时间	60	Catlet 类加载器 (废弃)			2.17.04.0 废弃

<code>sqlInterceptorFile</code>	SQL解析器sqllogs	默认 ./logs/sql.txt	在每个SQL解析的之前都会通过写文件进行记录（废弃）			2.17.04.0 废弃
<code>defaultSqlParser</code>	默认SQL解析器	默认druidparser	解析器不支持其他，废弃			2.17.04.0 废弃
<code>handleDistributedTransactions</code>	分布式事务开关	-	0为不过滤分布式事务，1为过滤分布式事务（如果分布式事务内只涉及全局表，则不过滤），2为不过滤分布式事务，但是记录分布式事务日志			2.17.04.0 废弃

## 用户配置

### user(管理用户配置)

配置名称	配置内容	是否能配多个元素	可选项/默认值	详细作用原理或应用	全局/实例属性
<code>name</code>	用户名	否	符合mysql用户名规范的字符串	用户唯一标识，用于登录校验	实例
<code>manager</code>	是否为manager用户	否	true-是 false-否 默认否	是否是manager用户，若一个用户是manager用户，则只能登录管理端口，并且只能执行管理端命令，依旧受到白名单制约	实例
<code>password</code>	密码	否		用户密码校验	实例
<code>maxCon</code>	负载限制，默认不做限制	否		用户的连接数限制，会在用户验证登录的时候进行校验，默认0，表示不做限制。特别的，当系统级别的maxCon已经到达上限之后，本用户的maxCon会失效，不能新建连接	实例
<code>usingDecrypt</code>	是否启用加密	否	默认 0 否	启用加密password项配置通过执行脚本encrypt.sh 0:{user}:{password}的结果进行配置 举例： encrypt.sh 0:xxx:123456 fP/nl3XPXrSfWjpQzit5llOrRU1QRXuLTYtATUG0fGW2k5kdXUhKL5zf02hE6nGjdnSWrufVkJPUZpbQ2qX9uQ== 配置项： password fP/nl3XPXrSfWjpQzit5llOrRU1QRXuLTYtATUG0fGW2k5kdXUhKL5zf02hE6nGjdnSWrufVkJPUZpbQ2qX9uQ== user xxx 登录项： -u root -p123456	实例

### user(业务用户配置)

配置名称	配置内容	是否能配多个元素	可选项/默认值	详细作用原理或应用	全局/实例属性
<code>name</code>	用户名	否	符合mysql用户名规范的字符串	用户唯一标识，用于登录校验	实例
<code>password</code>	密码	否		用户密码校验	实例
<code>maxCon</code>	负载限制，默认不做限制	否		用户的连接数限制，会在用户验证登录的时候进行校验，默认0，表示不做限制。特别的，当系统级别的maxCon已经到达上限之后，本用户的maxCon会失效，不能新建连接	实例
<code>readOnly</code>	用户只读标志	否		在执行SQL的时候进行判断，如果是readonly不执行查询以外的SQL	实例
<code>schemas</code>	用户对应使用的schema.xml中的方案	否		用户有权限的schemas列表，通过,进行分割(db1,db2,db3)	实例
<code>usingDecrypt</code>	是否启用加密	否	默认 0 否	启用加密password项配置通过执行脚本encrypt.sh 0:{user}:{password}的结果进行配置 举例： encrypt.sh 0:xxx:123456 fP/nl3XPXrSfWjpQzit5llOrRU1QRXuLTYtATUG0fGW2k5kdXUhKL5zf02hE6nGjdnSWrufVkJPUZpbQ2qX9uQ== 配置项： password fP/nl3XPXrSfWjpQzit5llOrRU1QRXuLTYtATUG0fGW2k5kdXUhKL5zf02hE6nGjdnSWrufVkJPUZpbQ2qX9uQ== user xxx 登录项： -u root -p123456	实例
<code>privileges</code>	权限	否		详情见user.privileges	实例

### user.privileges (用户权限)

配置名称	配置内容	多节点	可选项/默认值	详细作用原理或应用	全局/实例
<code>check</code>	是否需要检查	否	默认否 true-是 false-否	如果在进行DML解析并进行用户权限检查的时候，配置为是会严格检查权限情况，配置为否则不会检查	实例
<code>schema</code>	用户的schema节点权限情况	是	空	配置为空的情况下默认拥有所有有效schema的权限	实例

### user.privileges.schema

配置名称	配置内容	多节点	可选项/默认值	详细作用原理或应用	全局/实例
<code>name</code>	schema名称	否		用以标识对应schema	实例
<code>dml</code>	dml权限	否	0000	权限判断，每一位分别表示INSERT UPDATE SELECT DELETE四种权限 1- 拥有权限 0-没有权限 例如拥有所有权限为1111	实例
<code>table</code>		是		如果没有配置，则table继承schema的权限	实例

### user.privileges.schema.table

配置名称	配置内容	多节点	可选项/默认值	详细作用原理或应用	全局/实例
name	表格名称	否		在权限判断的时候作为key值	实例
dml	dml权限	否	0000	权限判断，每一位分别表示INSERT UPDATE SELECT DELETE四种权限 1- 拥有权限 0-没有权限 例如拥有所有权限为1111	实例

**firewall配置**

配置名称	是否可以配置多个	配置内容	默认值	配置范围/可选项	详细作用原理或应用	范围
whitehost	否	白名单	无		白名单的总结点	实例
blacklist	否	黑名单	无		黑名单的总结点	实例

**firewall.whitehost(host白名单配置)**

配置名称	节点是否可配多个	属性元素	元素是否可配多个	配置内容	默认值	配置范围/可选项	详细作用原理或应用	实例/全局属性
Host	是	host	否	白名单某个具体的IP host	无	有效IP地址 如需配置localhost可使用IP host = "0:0:0:0:0:0:1"	ip值是多个平行配置的唯一标识,ip重复配置会导致部分配置丢失。如单个IP想配置多个用户请在user属性里用逗号隔开。 黑名单功能类似mysql的权限控制,指定某些IP只能有某些特定的用户登录进行登录。 如启用白名单,除白名单标注的IP用户之外,其他dbie用户均无法登录  例如user="user1,user2,user3"	实例
		user	否	对应IP host可以允许连接的用户	无	有效dbie用户名, 使用';'进行分割		

**firewall.blacklist(黑名单配置)**

配置类型	配置名称	属性元素	是否可以配置多个	配置内容	默认值	配置范围/可选项	详细作用原理或应用	实例/全局属性
blacklist元素属性	blacklist	check	否	是否进行黑名单校验	false	true/false	如果开启黑名单校验,黑名单中规定的行为SQL都被禁止 影响范围全局,所有用户都无法豁免	实例
blacklist下挂property	property	name	是	详细的黑名单校验规则	无		如果开启黑名单校验具体的校验规则将有所有property来确定	实例

**firewall.blacklist.property(黑名单配置)**

配置名称	配置内容	默认值	可选项	详细作用原理或应用	备注
selectHavingAlwayTrueCheck	是否允许复杂select having 条件结果恒为真	true	true - 允许 false - 禁止	1、having部分的结果为真 2、SQL语句需要以注释结尾(不能有换行之类的) 3、条件部分不是简单SQL（单个条件、含有简单数值对等或大小比较、直接是真假值的表达式等） 符合以上三个条件的查询在校验的时候会被阻止， 举例：select * from test having id = id and hujh = hujh /*xxxxx*/;	暂不具有实际意义
selectWhereAlwayTrueCheck	是否允许复杂select where 条件结果恒为真	true	true - 允许 false - 禁止	1、where部分的结果为真 2、SQL语句需要以注释结尾(不能有换行之类的) 3、条件部分不是简单SQL（单个条件、含有简单数值对等或大小比较、直接是真假值的表达式等） 符合以上三个条件的查询在校验的时候会被阻止， 举例：select * from test where id = id and hujh = hujh /*xxxxx*/;	暂不具有实际意义
doPrivilegedAllow	druid内部权限控制使用	false	true - 允许 false - 禁止	druid内部函数调用flag，在dble中没有作用	暂不具有实际意义
wrapAllow	是否允许调用 isWrapFor和unwrap方法	true	true - 允许 false - 禁止	druid内部函数调用flag，在dble中没有作用	暂不具有实际意义
metadataAllow	是否允许调用getmetadata方 法	true	true - 允许 false - 禁止	druid内部函数调用flag，在dble中没有作用	暂不具有实际意义
completeInsertValuesCheck	在dble依赖的1.0.31版本中没 有效果	false	true - 允许 false - 禁止	druid内部函数调用flag，在dble中没有作用	暂不具有实际意义
mergeAllow	是否允许merge语句 (在mysql中不支持)	true	true - 允许 false - 禁止	会校验是否是merge语句	在Dble中没有效果
conditionLikeTrueAllow	是否允许like之后包含永真条件	true	true - 允许 false - 禁止	会根据SQL里面的内容进行判断，如果发现有like '%就会抛出异常	可用
conditionDoubleConstAllow	是否允许连续两个常量判断	false	true - 允许 false - 禁止	会根据SQL里面的内容进行判断，如果发现有两个常量判断抛出异常 select * from suntest asdf where 1 = 1 and 2 = 1;	可用
conditionAndAlwayFalseAllow	是否允许and连接的语句存在 恒为false的条件	false	true - 允许 false - 禁止	会根据where之后and跟随的条件进行判断，如果发现恒为假的情况会抛出异常 举例：select * from suntest where id = 567 and 1 != 1;	可用
conditionAndAlwayTrueAllow	是否允许and连接的语句存在 恒为true的条件	false	true - 允许 false - 禁止	会根据where之后and跟随的条件进行判断，如果发现恒为真的情况会抛出异常 举例：select * from suntest where id = 567 and 1 = 1;	可用
selectAllColumnAllow	是否允许查询所有列	true	true - 允许 false - 禁止	会根据查询sql进行判断，如果发现直接查询*会有这个异常， 但是如果带有别名(x.* )则不在此列	存在问题
multiStatementAllow	是否允许一次提交多条sql	false	true - 允许 false - 禁止	会检查sql的数量，如果超过一则抛出异常	Dble不支持
constArithmeticAllow	是否允许常量运算	true	true - 允许 false - 禁止	在SQL中如果发现包含有常量运算会抛出异常 select * from suntest asdf where id = 2 - 1;	可用
alterTableAllow	是否允许alter table 语句	true	true - 允许 false - 禁止	在执行alter table之前会纯粹根据SQL进行语句的校验， 发现SQL是alter table语句会抛出异常返回错误信息	可用
commitAllow	是否允许commit语句	true	true - 允许 false - 禁止	在执行commit之前会纯粹根据SQL进行语句的校验， 发现commit语句会抛出异常返回错误信息	可用
createTableAllow	是否允许create table 语句	true	true - 允许 false - 禁止	在执行create table之前会纯粹根据SQL进行语句的校验， 发现SQL是create table语句会抛出异常返回错误信息	可用
deleteAllow	是否允许delete语句	true	true - 允许 false - 禁止	在执行delete之前会纯粹根据SQL进行语句的校验，发现delete语句会抛出异常 返回错误信息	可用
dropTableAllow	是否允许drop table 语句	true	true - 允许 false - 禁止	在执行drop table之前会纯粹根据SQL进行语句的校验， 发现SQL是drop table语句会抛出异常返回错误信息	可用
insertAllow	是否允许insert	true	true - 允许 false - 禁止	在执行insert之前会纯粹根据SQL进行语句的校验，发现insert语句会抛出异常 返回错误信息	可用
intersectAllow	是否支持intersect	true	true - 允许 false - 禁止	在执行intersect之前会纯粹根据SQL进行语句的校验， 发现intersect语句会抛出异常返回错误信息	可用
lockTableAllow	是否允许lock tables语句	true	true - 允许 false - 禁止	在执行lock tables之前会纯粹根据SQL进行语句的校验， 发现lock语句会抛出异常返回错误信息	可用
minusAllow	是否支持minus语句	true	true - 允许 false - 禁止	在执行minus之前会纯粹根据SQL进行语句的校验， 发现minus语句会抛出异常返回错误信息	可用
callAllow	是否允许call语句	true	true - 允许 false - 禁止	在执行query之前会纯粹根据SQL进行语句的校验， 发现SQL是CALL语句会抛出异常返回错误信息	可用
selectIntoOutfileAllow	是否允许SELECT ... INTO OUTFILE	false	true - 允许 false - 禁止	在执行query之前会纯粹根据SQL进行语句的校验，发现SQL是 SELECT ... INTO OUTFILE语句会抛出异常返回错误信息	Dble本身不支 持
selectIntoAllow	是否允许select into 语句	true	true - 允许 false - 禁止	在执行query之前会纯粹根据SQL进行语句的校验， 发现SQL是select into语句会抛出异常返回错误信息	Dble本身不支 持
selectAllow	是否允许select语句	true	true - 允许 false - 禁止	在执行query之前会纯粹根据SQL进行语句的校验， 发现SQL是select语句会抛出异常返回错误信息	可用
renameTableAllow	是否允许rename table 语句	true	true - 允许 false - 禁止	在执行rename table之前会纯粹根据SQL进行语句的校验， 发现SQL是rename table语句会抛出异常返回错误信息	可用
replaceAllow	是否允许replace语句	true	true - 允许 false - 禁止	在执行replace之前会纯粹根据SQL进行语句的校验， 发现replace语句会抛出异常返回错误信息	可用

<code>rollbackAllow</code>	是否允许 <code>rollback</code>	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行 <code>rollback</code> 之前会纯粹根据SQL进行语句的校验，发现 <code>rollback</code> 语句会抛出异常返回错误信息	可用 <code>Dble慎用</code>
<code>setAllow</code>	是否允许 <code>set</code> 语句	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行 <code>set</code> 之前会纯粹根据SQL进行语句的校验，发现 <code>set</code> 语句会抛出异常返回错误信息	可用
<code>describeAllow</code>	是否支持 <code>describe</code> 语句	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现 <code>describe</code> 语句会抛出异常返回错误信息	可用
<code>limitZeroAllow</code>	是否允许出现 <code>limit 0</code> 的情况	<code>false</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现 <code>limit 0</code> 语句会抛出异常返回错误信息	可用
<code>showAllow</code>	是否允许 <code>show</code> 语句	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现 <code>show</code> 语句会抛出异常返回错误信息	可用
<code>hintAllow</code>	是否允许SQL包含 <code>hint</code>	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现SQL是包含 <code>hint</code> 语句会抛出异常返回错误信息	可用
<code>commentAllow</code>	是否允许在SQL中存在注释	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现注释语句会抛出异常返回错误信息	可用
<code>mustParameterized</code>	是否必须参数化	<code>false</code>	<code>true - 是</code> <code>false - 否</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现类似 <code>name = 'sdfasdf', id = 1</code> 语句会抛出异常返回错误信息	可用
<code>conditionOpXorAllow</code>	是否允许SQL中使用关系符XOR	<code>false</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现类似运算符语句会抛出异常返回错误信息	可用
<code>conditionOpBitwiseAllow</code>	查询条件中是否允许有“&”、“~”、“ ”、“^”运算符。	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行SQL之前会纯粹根据SQL进行语句的校验，发现类似运算符语句会抛出异常返回错误信息	可用
<code>startTransactionAllow</code>	是否允许 <code>START TRANSACTION</code>	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行 <code>START TRANSACTION</code> 之前会纯粹根据SQL进行语句的校验，发现 <code>START TRANSACTION</code> 语句会抛出异常返回错误信息 注：现阶段如果开启黑名单检查 <code>begin</code> 无法通过校验，这个是由于在druid中不支持的缘故	可用
<code>truncateAllow</code>	是否允许 <code>truncate</code> 语句	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行 <code>truncate</code> 之前会纯粹根据SQL进行语句的校验，发现 <code>truncate</code> 语句会抛出异常返回错误信息	可用
<code>updateAllow</code>	是否允许 <code>update</code> 语句	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行 <code>update</code> 之前会纯粹根据SQL进行语句的校验，发现 <code>update</code> 语句会抛出异常返回错误信息	可用
<code>useAllow</code>	是否允许 <code>use</code> 语句	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在执行 <code>use</code> 之前会纯粹根据SQL进行语句的校验，发现 <code>use</code> 语句会抛出异常返回错误信息	可用
<code>blockAllow</code>	是否允许语句块	<code>true</code>	<code>true - 允许</code> <code>false - 禁止</code>	在解析SQL阶段会判断SQL是否属于SQL语句块，如果是的话就会抛出错误 举例： <code>BEGIN select * from suntest;END;/</code>	可用
<code>deleteWhereNoneCheck</code>	是否允许 <code>delete</code> 语句没有 <code>where</code> 条件	<code>false</code>	<code>true - 启用</code> <code>false - 不启用</code>	如果发现 <code>delete</code> 语句没有限定条件会抛出异常 举例： <code>delete from suntest;</code>	可用
<code>updateWhereNoneCheck</code>	是否允许 <code>update</code> 语句没有 <code>where</code> 条件	<code>false</code>	<code>true - 启用</code> <code>false - 不启用</code>	如果发现 <code>update</code> 语句没有限定条件会抛出异常 举例： <code>update suntest set name = '33';</code>	可用
<code>deleteWhereAlwayTrueCheck</code>	是否允许 <code>delete</code> 语句存在恒真条件	<code>true</code>	<code>true - 启用</code> <code>false - 不启用</code>	如果解析发现 <code>delete</code> 语句存在恒真条件，并且满足sql以注释结尾，并且 <code>where</code> 条件不是简单条件的，会抛出异常 举例： <code>delete from suntest where id = id and name = name /*sdfaasdf*/;</code>	暂不具有实际意义
<code>updateWhereAlayTrueCheck</code>	是否允许 <code>update</code> 语句存在恒真条件	<code>true</code>	<code>true - 启用</code> <code>false - 不启用</code>	如果解析发现 <code>update</code> 语句存在恒真条件，并且满足sql以注释结尾，并且 <code>where</code> 条件不是简单条件的，会抛出异常 举例： <code>update suntest set name = '33' where id = id, name = name /*sdfaasdf*/;</code>	暂不具有实际意义
<code>selectIntersectCheck</code>	是否进行 <code>intersect</code> check	<code>true</code>	<code>true - 进行</code> <code>false - 不进行</code>	如果进行校验，则不允许 <code>except</code> 语句，当且仅当 <code>left sql</code> 的 <code>from</code> 不是空并且 <code>right from</code> 为空的时候不能通过校验 举例： <code>select * from sbtest1 where name = 'ff' INTERSECT select * from dual;</code>	暂不具有实际意义
<code>selectExceptCheck</code>	是否进行 <code>except</code> check	<code>true</code>	<code>true - 进行</code> <code>false - 不进行</code>	如果进行校验，则不允许 <code>except</code> 语句，当且仅当 <code>left sql</code> 的 <code>from</code> 不是空并且 <code>right from</code> 为空的时候不能通过校验 举例： <code>select * from sbtest1 where name = 'ff' except select * from dual;</code>	暂不具有实际意义
<code>selectMinusCheck</code>	是否进行 <code>MINUS</code> check	<code>true</code>	<code>true - 进行</code> <code>false - 不进行</code>	如果进行校验，则不允许 <code>MINUS</code> 语句，当且仅当 <code>left sql</code> 的 <code>from</code> 不是空并且 <code>right from</code> 为空的时候不能通过校验 举例： <code>select * from sbtest1 where namec = 'fff' minus select * from dual;</code>	暂不具有实际意义
<code>selectUnionCheck</code>	是否进行 <code>union</code> check	<code>true</code>	<code>true - 进行</code> <code>false - 不进行</code>	如果进行校验，则不允许 <code>union</code> 语句 举例： <code>select * from sbtest1 union select * from suntest;</code>	可用
<code>caseConditionConstAllow</code>	是否允许复杂查询中外部是一个常量	<code>false</code>	<code>true - 允许</code> <code>false - 禁止</code>	是否允许复杂查询中外部是一个常量，如果子查询外部对应的是常量那么就在SQL检查的时候抛出异常 举例： <code>delete from suntest where id = 123 and 'name' = (select case 'fname' when 'dsome' else 'good' end from xtest) /*sdfaasdf*/;</code>	暂不具有实际意义
<code>strictSyntaxCheck</code>	是否启用严格语法检查	<code>true</code>	<code>true - 是</code> <code>false - 否</code>	是否进行严格的语法检测，Druid SQL Parser在某些场景不能覆盖所有的SQL语法，属于调试级别的参数，在正常的使用中不建议更改	暂不具有实际意义
<code>schemaCheck</code>	检测是否使用了禁用的 schema	<code>true</code>	<code>true - 启用</code> <code>false - 禁止</code>	这个需要配合drui的配置模式使用，在dble此功能无法被使用	暂不具有实际意义
<code>tableCheck</code>	检测是否使用了禁用的 table	<code>true</code>	<code>true - 启用</code> <code>false - 禁止</code>	这个需要配合drui的配置模式使用，在dble此功能无法被使用	暂不具有实际意义
<code>functionCheck</code>	检测是否使用了禁用的 function	<code>true</code>	<code>true - 启用</code> <code>false - 禁止</code>	这个需要配合drui的配置模式使用，在dble此功能无法被使用	暂不具有实际意义
<code>objectCheck</code>	检测是否使用了禁用的 object	<code>true</code>	<code>true - 启用</code> <code>false - 禁止</code>	这个需要配合drui的配置模式使用，在dble此功能无法被使用	暂不具有实际意义

<b>variantCheck</b>	检测是否使用了禁用的变量	true false - 禁止	true - 启用 false - 禁止	这个需要配合drui的配置模式使用，在dble此功能无法被使用	暂不具有实际意义
---------------------	--------------	--------------------	-------------------------	---------------------------------	----------

**alarm(ucore告警grpc接口)**

此项不再需要单独配置，所需要的配置项从myid.properties 中读取，下面表格并非配置内容，而是内存里存储的值

名称	内容	默认值	详细作用原理或应用	实例/全局属性
url	grpc告警的url	myid.properties 里的ipAddress	在发送grpc的时候作为IP地址使用	实例
port	告警端口	myid.properties 里的port	grpc发送的目的端口	实例
level	告警等级	warn	如果配置error只会发送error等级的告警，如果配置warn会发送warn以及error的告警信息	实例
serverId	服务器ID	\$ushard-id(ip1,ip2) ,其中\$ushard-id 是myid.properties 里的myid	接口参数	实例
componentId	组件ID	\$ushard-id 即myid.properties 里myid	接口参数	实例
componentType	组件类型	ushard	接口参数	实例

**配置实例**

```

<?xml version="1.0" encoding="UTF-8"?>
<!- - Licensed under the Apache License, Version 2.0 (the "License");
- you may not use this file except in compliance with the License. - You
may obtain a copy of the License at - - http://www.apache.org/licenses/LICENSE-2.0
- - Unless required by applicable law or agreed to in writing, software -
distributed under the License is distributed on an "AS IS" BASIS, - WITHOUT
WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. - See the
License for the specific language governing permissions and - limitations
under the License. -->
<!DOCTYPE dble:server SYSTEM "server.dtd">
<db:server xmlns:db="http://dbie.cloud/" version="9.9.9.9">

<system>
    <!-- base config -->
    <!--<property name="bindIp">0.0.0.0</property>-->
    <!--<property name="serverPort">8066</property> -->
    <!--<property name="managerPort">9066</property> -->
    <!--<property name="processors">1</property>-->
    <!--<property name="processorExecutor">32</property> -->
    <!--<property name="fakeMySQLVersion">5.6.20</property>-->
    <property name="sequenceHandlerType">2</property>
    <!-- serverBacklog size,default 2048-->
    <property name="serverBackLog">2048</property>
    <!--<property name="serverNodeId"></property>-->
    <!--<property name="showBinlogStatusTimeout">60000</property>-->

    <!--option-->
    <!--<property name="useCompression">1</property>-->
    <!--<property name="usingAIO">0</property>-->

    <!--connection -->
    <!--<property name="charset">utf-8</property>-->
    <!--<property name="maxPacketSize">4194304</property>-->
    <!--<property name="txIsolation">3</property>-->

    <!--consistency-->
    <!-- check the consistency of table structure between nodes,default not -->
    <property name="checkTableConsistency">0</property>
    <!-- check periodt, he default period is 60000 milliseconds -->
    <property name="checkTableConsistencyPeriod">60000</property>

    <!-- heartbeat check period -->
    <property name="dataNodeIdleCheckPeriod">300000</property>
    <property name="dataNodeHeartbeatPeriod">10000</property>

    <!-- processor check conn-->
    <property name="processorCheckPeriod">1000</property><!-- unit millisecond -->
    <property name="sqlExecuteTimeout">300</property><!-- unit second -->
    <property name="idleTimeout">1800000</property><!-- unit millisecond -->

    <!-- transaction log -->
    <!-- 1 enable record the transaction log, 0 disable -->
    <property name="recordTxn">0</property>
    <!--<property name="transactionLogBaseDir"/>txlogs</property>-->
    <!--<property name="transactionLogBaseName">server-tx</property>-->
    <!--<property name="transactionRotateSize">16</property>&lt;!-- unit M &ndash;&gt;-->

    <!-- XA transaction -->
    <!-- use XA transaction ,if the mysql service crash,the unfinished XA commit/rollback will retry for several times
it is the check period for ,default is 1000 milliseconds-->
    <property name="xaSessionCheckPeriod">1000</property>
    <!-- use XA transaction ,the finished XA log will removed. the default period is 1000 milliseconds-->
    <property name="xaLogCleanPeriod">1000</property>
    <!--<property name="XARecoveryLogBaseDir"/>tmlogs</property>-->
    <!-- XA Recovery Log name -->
    <!--<property name="XARecoveryLogBaseName">tmlog</property>-->

    <!-- true is use JoinStrategy, default false-->
    <property name="useJoinStrategy">true</property>
    <property name="nestLoopConnSize"></property>
    <property name="nestLoopRowsSize">2000</property>

    <!-- off Heap unit:bytes-->
    <property name="bufferPoolChunkSize ">4096</property>
    <property name="bufferPoolPageNumber ">512</property>
    <property name="bufferPoolPageSize ">2097152</property>

    <!-- sql statistics-->
    <!-- 1 means use SQL statistics, 0 means not -->
    <property name="useSqlStat">0</property>
    <!--<property name="bufferUsagePercent">80</property>-->
    <!--<property name="clearBigSQLResultSetMapMs">6000000</property>-->
    <!--<property name="sqlRecordCount">10</property>-->
    <!--<property name="maxResultSet">524288</property>-->

    <!-- backSocket unit:bytes-->
    <!--<property name="backSocketSoRcvbuf ">4194304</property>-->
    <!--<property name="backSocketSoSndbuf">1048576</property>-->
    <!--<property name="backSocketNoDelay ">1</property>-->

    <!-- frontSocket-->
    <!--<property name="frontSocketSoRcvbuf ">1048576</property>-->
    <!--<property name="frontSocketSoSndbuf">4194304</property>-->
    <!--<property name="frontSocketNoDelay ">1</property>-->

    </system>

    <!-- firewall config -->
    <!--
    <firewall>
        <whitehost>
            <host host="127.0.0.1" user="root"/>
            <host host="0:0:0:0:0:0:1" user="root"/>
        </whitehost>
        <blacklist check="true">
            <property name="selectAllow">false</property>
        </blacklist>
    </firewall>
    -->

    <user name="man1">
        <property name="password">654321</property>
        <property name="manager">true</property>
        <!-- manager user can't set schema-->
    </user>

    <user name="root">
        <property name="password">123456</property>
        <property name="schemas">TESTDB</property>

        <!-- table's DML privileges INSERT/UPDATE/SELECT/DELETE -->
        <!--
        <privileges check="false">
            <schema name="TESTDB" dml="0110" >
                <table name="tb01" dml="0000"></table>
                <table name="tb02" dml="1111"></table>
            </schema>
        </privileges>
        -->
    </user>

    <user name="user">
        <property name="password">user</property>
        <property name="schemas">TESTDB</property>
        <property name="readOnly">true</property>
    </user>
</db:server>
```

## 1.4 Wrapper.conf配置概述

在Dble中选择wrapper作为管理dble进程的工具，其大概配置都属于固定的部分，只有其中的JVM参数是用户可以和需要通过修改来达成配置目的的。JVM参数都统一配置在以下的部分配置模块中，additional之后的标号可以根据现实情况进行添

```
wrapper.java.additional.1=-DDBLE_HOME=.
wrapper.java.additional.2=-server
wrapper.java.additional.3=-XX:+AggressiveOpts
wrapper.java.additional.4=-Dfile.encoding=UTF-8
wrapper.java.additional.5=-Dcom.sun.management.jmxremote
wrapper.java.additional.6=-Dcom.sun.management.jmxremote.port=1984
wrapper.java.additional.7=-Dcom.sun.management.jmxremote.authenticate=false
wrapper.java.additional.8=-Dcom.sun.management.jmxremote.ssl=false
wrapper.java.additional.9=-Dcom.sun.management.jmxremote.host=127.0.0.1
wrapper.java.additional.10=-Xmx4G
wrapper.java.additional.11=-Xms1G
wrapper.java.additional.12=-XX:MaxDirectMemorySize=2G
wrapper.java.additional.13=-XX:+PrintHeapAtGC
wrapper.java.additional.14=-XX:+PrintGCDetails
wrapper.java.additional.15=-Xloggc:/logs/gc_%t_%p.log
wrapper.java.additional.16=-XX:+PrintGCTimeStamps
wrapper.java.additional.17=-XX:+PrintGCDetails
wrapper.java.additional.18=-XX:+HeapDumpOnOutOfMemoryError
wrapper.java.additional.19=-XX:HeapDumpPath=/logs/
```

大部分以上的配置都没有特殊的意义，仅仅是一般的JVM配置，关于JVM调优的部分需要以现实情况进行操作，在此仅介绍几个特殊情况

1. MaxDirectMemorySize需要根据机器的情况进行提前适配，不然会导致服务无法正常启动。具体的细节为需要大于bufferPoolPageNumber\*bufferPoolPageSize，这两个选项在server.xml中配置

bufferPoolPageNumber 的默认值= (MaxDirectMemorySize \* 0.8 /bufferPoolPageSize)，向下取整，上限值为32767

bufferPoolPageSize 的默认值= 2M

以下为建议值：

dble总内存=0.6 可用物理内存(刨除操作系统,驱动等的占用)

Xmx = 0.4 dble总内存

MaxDirectMemorySize = 0.6 \* dble总内存

另外，在1.3\_server.xml中的bufferPoolPageNumber 和bufferPoolPageSize 受MaxDirectMemorySize影响。

2. 为了调试方便在debug模式下存在三个可用的XA事务调试JVM参数

-DPREPARE\_DELAY=10

-DCOMMIT\_DELAY=10

-DROLLBACK\_DELAY = 10

单位分别是秒，当且仅当DEBUG模式下，在XA事务发生三段提交或者回滚之前会发生延迟，单位为秒

注：大部分情况下这个参数用以进行测试和本地调试

### 1.4.1 配置实例

```

*****#
# Wrapper Properties
#*****wrapper.java.mainclass
# Java Application
wrapper.java.command=java
wrapper.working.dir=.

# Java Main class. This class must implement the WrapperListener interface
# or guarantee that the WrapperManager class is initialized. Helper
# classes are provided to do this for you. See the Integration section
# of the documentation for details.
wrapper.java.mainclass=org.tanukisoftware.wrapper.WrapperSimpleApp
set.default.REPO_DIR=lib
set.APP_BASE=.

# Java Classpath (include wrapper.jar) Add class path elements as
# needed starting from 1
wrapper.java.classpath.1=lib/wrapper.jar
wrapper.java.classpath.2=conf
wrapper.java.classpath.3=%REPO_DIR%/*
wrapper.java.classpath.4=.

# Java Library Path (location of Wrapper.DLL or libwrapper.so)
wrapper.java.library.path.1=lib

# Java Additional Parameters
#wrapper.java.additional.1=
wrapper.java.additional.1=-DDBLE_HOME=.
wrapper.java.additional.2=-agentlib:jdwtp=transport=dt_socket,server=y,address=8088,suspend=n
wrapper.java.additional.3=-server
wrapper.java.additional.4=-XX:MaxPermSize=64M
wrapper.java.additional.5=-XX:+AggressiveOpts
wrapper.java.additional.6=-Dfile.encoding=UTF-8
wrapper.java.additional.7=-Dcom.sun.management.jmxremote
wrapper.java.additional.8=-Dcom.sun.management.jmxremote.port=1984
wrapper.java.additional.9=-Dcom.sun.management.jmxremote.authenticate=false
wrapper.java.additional.10=-Dcom.sun.management.jmxremote.ssl=false
wrapper.java.additional.11=-Dcom.sun.management.jmxremote.host=127.0.0.1
wrapper.java.additional.12=-Xmx4G
wrapper.java.additional.13=-Xms1G
wrapper.java.additional.14=-XX:MaxDirectMemorySize=2G

# Initial Java Heap Size (in MB)
#wrapper.java.initmemory=3

# Maximum Java Heap Size (in MB)
#wrapper.java.maxmemory=64

# Application parameters. Add parameters as needed starting from 1
wrapper.app.parameter.1=DbleStartup
wrapper.app.parameter.2=start

*****#
# Wrapper Logging Properties
#*****#
# Format of output for the console. (See docs for formats)
wrapper.console.format=PM

# Log Level for console output. (See docs for log levels)
wrapper.console.loglevel=INFO

# Log file to use for wrapper output logging.
wrapper.logfile=logs/wrapper.log

# Format of output for the log file. (See docs for formats)
wrapper.logfile.format=LPTM

# Log Level for log file output. (See docs for log levels)
wrapper.logfile.loglevel=INFO

# Maximum size that the log file will be allowed to grow to before
# the log is rolled. Size is specified in bytes. The default value
# of 0, disables log rolling. May abbreviate with the 'k' (kb) or
# 'm' (mb) suffix. For example: 10m = 10 megabytes.
wrapper.logfile.maxsize=0

# Maximum number of rolled log files which will be allowed before old
# files are deleted. The default value of 0 implies no limit.
wrapper.logfile.maxfiles=0

# Log Level for sys/event log output. (See docs for log levels)
wrapper.syslog.loglevel=NONE

*****#
# Wrapper Windows Properties
#*****#
# Title to use when running as a console
wrapper.console.title=Dble-server

*****#
# Wrapper Windows NT/2000/XP Service Properties
#*****#
# WARNING - Do not modify any of these properties when an application
# using this configuration file has been installed as a service.
# Please uninstall the service before modifying this section. The
# service can then be reinstalled.

# Name of the service
wrapper.ntservice.name=dble

# Display name of the service
wrapper.ntservice.displayname=Dble-server

# Description of the service
wrapper.ntservice.description=The project of Dble-server

# Service dependencies. Add dependencies as needed starting from 1
wrapper.ntservice.dependency.1=

# Mode in which the service is installed. AUTO_START or DEMAND_START
wrapper.ntservice.starttype=AUTO_START

# Allow the service to interact with the desktop.
wrapper.ntservice.interactive=false

wrapper.ping.timeout=120
configuration.directory.in.classpath.first=conf

```

## 1.5 log4j2.xml

### 1.5.1 配置详述

Dble中的整体配置和一般java项目的log4j2.xml没有什么区别

#### 1.5.1.1 日志滚动删除配置

日志滚动删除是通过DefaultRolloverStrategy的配置对于RollingRandomAccessFile的RolloverStrategy内容进行重载

```
<DefaultRolloverStrategy max="100">
    <Delete basePath="logs" maxDepth="2">
        <IfFileName glob="*/dble-*.*.log.gz">
            <IfLastModified age="2d">
                <IfAny>
                    <IfAccumulatedFileSize exceeds="1 GB" />
                    <IfAccumulatedFileCount exceeds="10" />
                </IfAny>
            </IfLastModified>
        </IfFileName>
    </Delete>
</DefaultRolloverStrategy>
```

上例中参数说明如下：

basePath：基准路径，日志的统计归类和删除会在此目录下进行

maxDepth：最大路径深度，在basePath路径下maxDepth深度的日志文件都会被扫描，譬如.../logs/2018-01-01 .../logs/2018-01-02此类路径也都会被扫描和统计

glob：日志格式，在所有扫描到的文件中符合此命名规范的文件都被认为是日志文件

age：举例最后一次修改文件的时间，只有修改时间超过限值的文件才会被考虑删除

IfAccumulatedFileSize：触发文件大小，符合上述条件的文件大小总量达到触发值则触发清理

IfAccumulatedFileCount：触发文件数量，符合上述条件的文件数量达到触发值则触发清理

例子详述：在此例中代表会监控logs目录下2层目录深度内所有命名符合“dble-.log.gz”并且最后修改时间已经超出2天的文件，当有符合上述条件的文件大小到达1 GB或者文件数量到达10的时候会触发清理

### 1.5.2 配置实例

```
<?xml version="1.0" encoding="UTF-8"?>
<Configuration status="WARN" packages="com.actiontech.dble.log">
    <Appenders>
        <Console name="Console" target="SYSTEM_OUT">
            <PatternLayout pattern="%d [%-5p][%t] %m %throwable{full} (%C:%F:%L) %n"/>
        </Console>

        <RollingRandomAccessFile name="RollingFile" fileName="logs/dble.log"
            filePattern="logs/${date:yyyy-MM}/dble-%d{MM-dd}-%i.log.gz">
            <PatternLayout>
                <Pattern>%d{yyyy-MM-dd HH:mm:ss.SSS} %5p [%t] (%l) - %m%n</Pattern>
            </PatternLayout>
            <Policies>
                <OnStartupTriggeringPolicy/>
                <SizeBasedTriggeringPolicy size="250 MB"/>
                <TimeBasedTriggeringPolicy/>
            </Policies>
            <DefaultRolloverStrategy max="100">
                <Delete basePath="logs" maxDepth="2">
                    <IfFileName glob="*/dble-*.*.log.gz">
                        <IfLastModified age="2d">
                            <IfAny>
                                <IfAccumulatedFileSize exceeds="1 GB" />
                                <IfAccumulatedFileCount exceeds="10" />
                            </IfAny>
                        </IfLastModified>
                    </IfFileName>
                </Delete>
            </DefaultRolloverStrategy>
        </RollingRandomAccessFile>

        </Appenders>
        <Loggers>
            <asyncRoot level="debug" includeLocation="true">
                <AppenderRef ref="Console"/>
                <AppenderRef ref="RollingFile"/>
            </asyncRoot>
        </Loggers>
    </Configuration>
```

## 1.6 cache配置

- [1.6.1 cache配置](#)
- [1.6.2 ehcache配置](#)

## 1.6.1 cache 配置

### 1.6.1.1 dble的cache使用

dble的cache使用有如下三类:

- SQLRouteCache: 从前端连接收到的SQL以及对应的路由结果 内容: schema\_user\_SQL -> 具体路由结果RouteResult
- ER\_SQL2PARENTID: 父子表辅助查询SQL以及对应的路由结果。在插入ER子表的时候需要根据子表joinKey(父表parentKey)计算它应该插入的结点, 所以需要辅助路由来查询, 然后将辅助路由及对应的结果缓存下来以备下次查询。内容为: schema:select \* from 父表 where parentKey = (value of joinKey) -> 对应数据dataNode
- TableID2DataNodeCache: 表格主键(schema.xml中声明的主键)和对应数据节点 内容 tableName -> primaryKeyValue -> 对应数据dataNode

### 1.6.1.2 dble的cache实现

dble的cache实现有如下几种:

- ehcache, 用ehcache缓存作为cache实现。
- leveldb, 用leveldb数据库作为cache实现。
- mapdb, 用MapDB数据库引擎作为cache实现。
- rocksdb, 用RocksDB数据库引擎作为cache实现。

### 1.6.1.3 dble的cache配置

dble的cache配置分为总配置和实现配置。总配置由文件cacheservice.properties进行设定。实现配置由各个实现具体指定, 具体详见各个实现的分章节说明。

总配置有如下格式:

设置缓存类型:

```
factory.cache_type=cache_type
```

设置分类缓存的具体值, key为缓存池名字, value是类型, 最大容量, 以及失效时间

#### A.SQL路由缓存

```
pool.SQLRouteCache=type,max_size,expire_seconds
```

#### B.ER表子表路由缓存

```
pool.ER_SQL2PARENTID=type,max_size,expire_seconds
```

#### C.主键缓存

主键缓存有两种设置方式:

方式1:

设置缺省主键缓存和指定表的主键缓存, 未指定的表将共用缺省主键缓存空间:

缺省主键缓存:

```
layedpool.TableID2DataNodeCache=type,max_size,expire_seconds
```

指定表的主键缓存的key的格式是用如下格式增加配置表名:

```
layedpool.TableID2DataNodeCache.`schema`_`table`
```

举例如下:

```
layedpool.TableID2DataNodeCache.`schema`_`table`=max_size,expire_seconds
```

方式2:

只设置指定表的主键缓存, 其他表不使用主键缓存 layedpool.TableID2DataNodeCacheType=type

指定表的主键缓存的key的格式是用如下格式增加配置表名:

```
layedpool.TableID2DataNodeCache.`schema`_`table`
```

举例如下:

```
layedpool.TableID2DataNodeCache.`schema`_`table`=max_size,expire_seconds
```

### 1.6.1.4 cache配置说明

总配置文件中各配置项说明:

a. 以#开头的行为注释, 被忽略。空行被忽略。

b. factory.cache\_type=cache\_type是cache的总开关。cache\_type指定cache类型, 具体可以为: ehcache, leveldb、mapdb 或者rocksdb。如果要用cache功能, 必须配置该配置项。这个配置项可以指定多个, 每行仅能指定一个。每一个指定一个cache实现。

例如:

配置,

```
factory.encache=ehcache
pool.SQLRouteCache=encache,10000,1800
pool.ER_SQL2PARENTID=encache,1000,1800
```

中的type就必须是ehcache。而配置:

```
factory.encache=ehcache
factory.leveldb=leveldb
pool.SQLRouteCache=encache,10000,1800
pool.ER_SQL2PARENTID=leveldb,1000,1800
```

中的type可以为encache或者leveldb。

c. pool.SQLRouteCache=type,max\_size,expire\_seconds和pool.ER\_SQL2PARENTID=type,max\_size,expire\_seconds分别配置SQLRouteCache和ER\_SQL2PARENTID的缓存功能。这两个配置项可以配置也可以不配置, 不配值则不使用相应的缓存功能。type指定缓存类型, 必须是已配置的缓存实现类型; max\_size指定缓存的最大大小, 单位是字节; expire\_seconds指定缓存项的生命周期, 单位是秒。

d. layedpool.TableID2DataNodeCache=type,max\_size,expire\_seconds和layedpool.TableID2DataNodeCacheType=type。

这是TableID2DataNodeCache缓存配置的两种模式, 且这两种配置模式不能共存, 仅能指定一种。layedpool.TableID2DataNodeCache=type,max\_size,expire\_seconds配置模式使用default缓存; layedpool.TableID2DataNodeCacheType=type配置模式不使用default缓存。type指定缓存类型, 必须是已配置的缓存实现类型; max\_size指定缓存的最大大小, 单位是字节; expire\_seconds指定缓存项的生命周期, 单位是秒。

e. default缓存用于缓存没有为其指定特定缓存的表的分区键值到数据所在节点的映射。

f. layedpool.TableID2DataNodeCache.`schema`\_`table`=max\_size,expire\_seconds指定特定表的分区键值到数据所在节点的映射。

schema用实际的库名代换; table用实际的表名代换。max\_size指定缓存的最大大小, 单位是字节; expire\_seconds指定缓存项的生命周期, 单位是秒。缓存的实现类型同 layedpool.TableID2DataNodeCache=type,max\_size,expire\_seconds或者layedpool.TableID2DataNodeCacheType=type的type。

g. TableID2DataNodeCache缓存也可以不配置, 不配值则不使用此缓存功能。

### 1.6.1.5 注意事项

- 使用 RocksDB作为 cache 实现时, 需要在dble目录下手工创建 rocksdb 目录, 否则dble启动失败。

## 1.6.2 ehcache缓存配置

要用ehcache实现缓存，必须在cacheservice.properties中配置使用ehcache，具体请参看上一节内容。

### 1.6.2.1 ehcache版本

目前，dble使用的是2.6.11.

### 1.6.2.2 ehcache配置

ehcache的配置通过文件ehcache.xml进行。

具体的缓存存储策略和配置请参看<http://www.ehcache.org/documentation/ehcache-2.6.x-documentation.pdf>。

例如：

```
<ehcache xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:noNamespaceSchemaLocation="ehcache.xsd" maxEntriesLocalHeap="100000000" maxBytesLocalDisk="50G" updateCheck="false">
    <defaultCache maxElementsInMemory="1000000" eternal="false" overflowToDisk="false" diskSpoolBufferSizeMB="30" maxElementsOnDisk="10000000" diskPersistent="false" diskExpiryThreadIntervalSeconds="120" memoryStor</defaultCache>
```

需要特殊说明是：

1.dble仅用ehcache配置的defaultCache级别创建cache。

2.maxEntriesLocalHeap

该属性指定允许存储元素的最大条数。如果设定了该值且不为0，则缓存大小限制为缓存的**条数限制**，具体限制由cacheservice.properties中**max\_size**指定，参见上一节内容。如果没有设定该参数，则缓存大小限制仍为大小限制，具体限制由cacheservice.properties中**max\_size**指定，参见上一节内容。

3.timeToIdleSeconds

该属性指定一个元素在不被请求的情况下允许在缓存中存在的最长时间。它将被cacheservice.properties中**expire\_seconds**取代。

4.dble将defaultCache配置应用到每一个创建的cache。

## 1.7 全局序列

在分库分表的情况下，数据库自增主键无法保证自增主键的全局唯一。为此，**dble**提供了全局序列，并且针对不同应用场景提供了多种实现方式。

要应用全局序列，首先需要在**server.xml**中配置(参见[1.3 server.xml](#))：

```
<system>
  <property name="sequenceHandlerType">1</property>
</system>
```

根据**sequenceHandlerType**的类型配置具体的实现：

- 1:[MySQL offset-step方式](#)
- 2:[时间戳方式\(类Snowflake\)](#)
- 3:[分布式时间戳方式\(类Snowflake\)](#)
- 4:[分布式offset-step方式](#)

其次，使用全局序列的表要在**schema.xml**中配置表格并指定其自增列(参见[1.2 schema.xml](#))。

**自增列的指定逻辑：**以显式声明的以**incrementColumn**中的指定为准

```
//有显式指定incrementColumn=pid，则pid作为自增列
<table name="table1" cacheKey="id" dataNode="dn1,dn2" rule="two_node_hash" incrementColumn="pid"/>
```

配置之后的使用示例：

```
insert into table1(name) values('test');
insert into table1 set name = 'test';
```

每一个全局序列的具体实现配置将在各个小节进行详细说明；其功能将在[2.2 全局序列](#)进行详细描述。

### 与MySQL差异

在MySQL中要求自增列必须含有唯一键，在**dble**中对于表格的自增列的唯一属性不做要求，但是在**dble**进行插入数据操作的时候不允许用户手动插入自增列的数值，自增列只能由**dble**自己生成的ID进行填充并且存在以下情形和mysql的行为不一致

```
table1拥有列aid,bid,cid,did 其中bid为在dble中的自增列
insert into table1 values(1,2,3)
和以下sql的效果相等
insert into table1 set aid = 1,cid = 2,did = 3
```

**特别提醒** 在**dble**中全局序列只在生成的时候确保其唯一性，在之后的过程中用户被允许使用**update**或者**replace**语句进行更新自增字段（特例：自增字段同时也是分片字段是除外）。这给予用户提供了充足的修改空间但同时要求用户在更新自增字段的时候有足够的谨慎和了解

### 1.7.1 MySQL-offset-step 方式

#### 1.7.1.1 MySQL-offset-step序列配置

mysql序列由文件sequence\_db\_conf.properties进行配置。具体配置有如下格式:

```
#this is comment
`schema1`.`table1`=node1
`schema1`.`table2`=node1
`schema2`.`table1`=node2
...
```

**schemaX**: 使用全局序列的db表所属的db库名。

**tableX**: 使用全局序列db的db表名。

**nodeX**: 实现序列功能的数据节点名。此节点具体配置见下节。

#### 1.7.1.2 数据节点配置

mysql序列的实现依赖一些存储函数。因此在序列可用之前必须在数据节点**nodeX**上创建这些存储函数，并针对使用此节点的表(`**schemaX**.`**tableX**)初始化序列初始值。

数据节点**nodeX**上创建必要的表和存储函数

具体步骤:

a. 用mysql客户端登录到**nodeX**上

```
mysql ...
```

b. 切换到**nodeX**上的后端存储数据库db(参见1.2 schema.xml)

```
use db;
```

c. 执行创建脚本dbsql.sql (dbseq.sql的内容详见1.7.2.3节)

```
source ../../dbseq.sql;
```

d. 初始化相应序列初始值

```
INSERT INTO DBLE_SEQUENCE VALUES (`schemaX.`tableX`, 1, 1);
```

```
...
```

#### 1.7.1.3 dbseq.sql内容

dbseq.sql的内容如下:

```

DROP TABLE IF EXISTS DBLE_SEQUENCE;
CREATE TABLE DBLE_SEQUENCE ( name VARCHAR(64) NOT NULL, current_value BIGINT(20) NOT NULL, increment INT NOT NULL DEFAULT 1, PRIMARY KEY (name) ) ENGINE=InnoDB;

-- -----
-- Function structure for `dble_seq_currval`-
-- -----
DROP FUNCTION IF EXISTS `dble_seq_currval`;
DELIMITER ;;
CREATE FUNCTION `dble_seq_currval`(seq_name VARCHAR(64)) RETURNS varchar(64) CHARSET latin1
DETERMINISTIC
BEGIN
DECLARE retval VARCHAR(64);
SET retval="-1,0";
SELECT concat(CAST(current_value AS CHAR),"",CAST(increment AS CHAR) ) INTO retval FROM DBLE_SEQUENCE WHERE name = seq_name;
RETURN retval ;
END
;;
DELIMITER ;

-- -----
-- Function structure for `dble_seq_nextval`-
-- -----
DROP FUNCTION IF EXISTS `dble_seq_nextval`;
DELIMITER ;;
CREATE FUNCTION `dble_seq_nextval`(seq_name VARCHAR(64)) RETURNS varchar(64) CHARSET latin1
DETERMINISTIC
BEGIN
DECLARE retval VARCHAR(64);
DECLARE val BIGINT;
DECLARE inc INT;
DECLARE seq_lock INT;
set val = -1;
set inc = 0;
SET seq_lock = -1;
SELECT GET_LOCK(seq_name, 15) into seq_lock;
if seq_lock = 1 then
SELECT current_value + increment, increment INTO val, inc FROM DBLE_SEQUENCE WHERE name = seq_name for update;
if val != -1 then
UPDATE DBLE_SEQUENCE SET current_value = val WHERE name = seq_name;
end if;
SELECT RELEASE_LOCK(seq_name) into seq_lock;
end if;
SELECT concat( CAST((val - inc + 1) as CHAR),"",CAST(inc as CHAR)) INTO retval;
RETURN retval;
END
;;
DELIMITER ;

-- -----
-- Function structure for `dble_seq_setvals`-
-- -----
DROP FUNCTION IF EXISTS `dble_seq_setvals`;
DELIMITER ;;
CREATE FUNCTION `dble_seq_setvals`(seq_name VARCHAR(64), count INT) RETURNS VARCHAR(64) CHARSET latin1
DETERMINISTIC
BEGIN
DECLARE retval VARCHAR(64);
DECLARE val BIGINT;
DECLARE seq_lock INT;
SET val = -1;
SET seq_lock = -1;
SELECT GET_LOCK(seq_name, 15) into seq_lock;
if seq_lock = 1 then
SELECT current_value + count INTO val FROM DBLE_SEQUENCE WHERE name = seq_name for update;
IF val != -1 THEN
UPDATE DBLE_SEQUENCE SET current_value = val WHERE name = seq_name;
END IF;
SELECT RELEASE_LOCK(seq_name) into seq_lock;
end if;
SELECT CONCAT( CAST((val - count + 1) as CHAR), "", CAST(count as CHAR)) INTO retval;
RETURN retval;
END
;;
DELIMITER ;

-- -----
-- Function structure for `dble_seq_setval`-
-- -----
DROP FUNCTION IF EXISTS `dble_seq_setval`;
DELIMITER ;;
CREATE FUNCTION `dble_seq_setval`(seq_name VARCHAR(64), value BIGINT) RETURNS varchar(64) CHARSET latin1
DETERMINISTIC
BEGIN
DECLARE retval VARCHAR(64);
DECLARE inc INT;
SET inc = 0;
SELECT increment INTO inc FROM DBLE_SEQUENCE WHERE name = seq_name;
UPDATE DBLE_SEQUENCE SET current_value = value WHERE name = seq_name;
SELECT concat( CAST(value as CHAR),"",CAST(inc as CHAR)) INTO retval;
RETURN retval;
END
;;
DELIMITER ;

```

## 1.7.2 时间戳方式

时间戳方式由sequence\_time\_conf.properties文件进行配置。具体配置如下：

```
#this is comment  
WORKID=01  
DATAACENTERID=01  
START_TIME=2010-10-04 09:42:54
```

每个配置项的含义为：

**WORKID**: 指定worker id值，必须为[0,31]之间的整数。

**DATAACENTERID**: 指定data center id值，必须为[0,31]之间的整数。

**START\_TIME**: 指定开始时间戳，格式必须为 YYYY-MM-dd HH:mm:ss，默认开始时间 2010-10-04 09:42:54。

注意事项： **WORKID**, **DATAACENTERID** 的配置必须使该dble实例在dble集群中唯一。

另外，使用这种方式需要对应字段为bigint来保证63位

### 配置示例

见[时间戳方式全局序列的配置](#)

### 1.7.3 分布式时间戳方式

分布式时间戳方式由文件sequence\_distributed\_conf.properties进行配置。具体配置格式如下:

```
#this is comment  
INSTANCEID=01  
CLUSTERID=01  
START_TIME=2010-11-04 09:42:54
```

**INSTANCEID**: 指定实例ID值, 可以为'ZK'或者n (n为区间[0, 31]中的一个整数)。

**CLUSTERID**: 指定组ID值, 可以为m(m为区间[0, 15]中的一个整数)。

**START\_TIME**: 指定开始时间, 时间格式固定, 必须为 2010-11-04 09:42:54。

注意事项:

1. 当**INSTANCEID**的值配置为'ZK'时, 必须配置zookeeper服务器(参见1.8 myid.properties)。
2. **INSTANCEID**, **CLUSTERID**相当于联合主键, 保证dble在集群中唯一。若在2.19.05.0之前(含)版本的dble中使用zk分配**INSTANCEID**, 在**CLUSTERID**相同的情况下, 会存在**INSTANCEID**重复的风险, 进而导致全局序列重复, 详情可以参考<https://github.com/actiontech/dble/issues/1184>。建议在此版本和更早的版本中使用不同**CLUSTERID**。
3. 使用这种方式需要对应字段为bigint来保证63位。

#### 1.7.4 分布式offset-step方式

分布式offset-step方式由文件sequence\_conf.properties进行配置。具体格式如下:

```
# this is comment
`schema1`.`table1`.MINID=1001
`schema1`.`table1`.MAXID=2000
`schema1`.`table1`.CURID=1000

`schema2`.`table2`.MINID=1001
`schema2`.`table2`.MAXID=20000
`schema2`.`table2`.CURID=1000
```

**schemaX**: 使用mysql序列的db表所属的db库名。

**tableX**: 使用mysql序列的db表名。

注意事项:

1. 每一个zk序列均需要指定当前区间内最小值MINID, 当前区间内最大值MAXID, 当前区间内当前值CURID。这些值仅在初始配置时有效。
2. 初始配置时MINID要比CURID大1, 否则序列值从MINID+1开始。
3. 值 (MAXID - MINID + 1) 为每次从zookeeper服务器获取的序列值数量。
4. 配置db使用zookeeper服务器, 具体见相关配置[1.8 myid.properties](#)。

## 1.8 myid.properties

- cluster  
三种选择使用ZK为zk，使用Ucore为ucore，都不使用就是false，其他信息会被忽略
- ipAddress,port  
zk或者是ucore的连接信息；如果使用的是ZK，请将ZK配置信息写到ipAddress，用逗号隔开，port配置将被忽略
- clusterId  
集群ID，处用同一集群的dble实例在这个配置上需要一致
- myid  
实例ID，唯一标识实例的ID，在所有dble配置都需不同
- clusterHa  
是否启用外部高可用的集群联动 true为是 false或其他为否，不配置默认为否

### 1.8.1 不使用的例子

```
#set false if not use cluster ucore/zk
cluster=false
```

### 1.8.2 使用ZK的例子

```
#set false if not use cluster ucore/zk
cluster=zk
#clinet info
ipAddress=10.186.19.aa:2281,10.186.60.bb:2281

#cluster namespace, please use the same one in one cluster
clusterId=cluster-1
#it must be different for every node in cluster
myid=server_02
```

注意：2.18.12.X版本此处有bug，需要增加一行

```
port=xxxx
```

### 1.8.3 使用ucore的例子

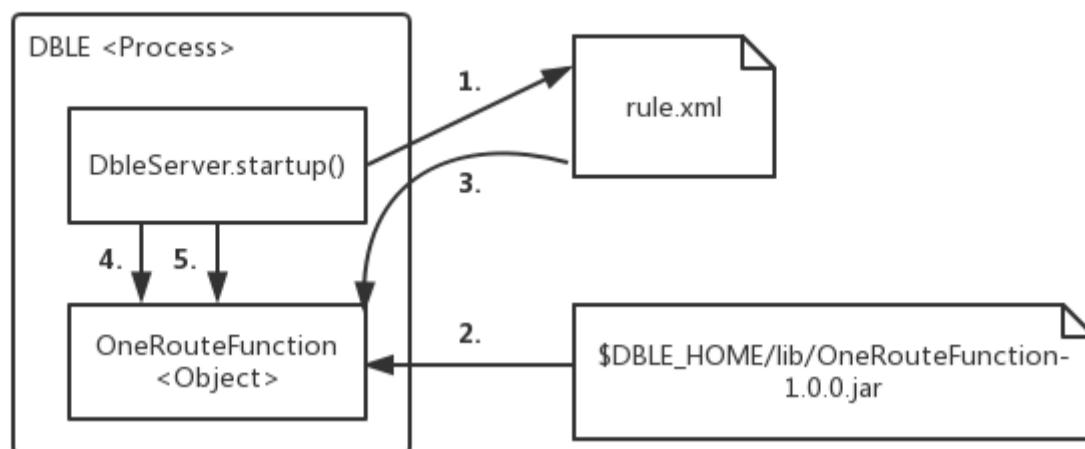
```
#set false if not use cluster ucore/zk
cluster=ucore
#clinet info
ipAddress=10.186.24.xx,10.186.24.yy
port=5700
#cluster namespace, please use the same one in one cluster
clusterId=cluster-1
#it must be different for every node in cluster
myid=server_02
# for ucore
serverID=server-udp1
```

## 1.9 自定义拆分算法

### 1.9.1 工作原理

#### 1.9.1.1 函数的加载

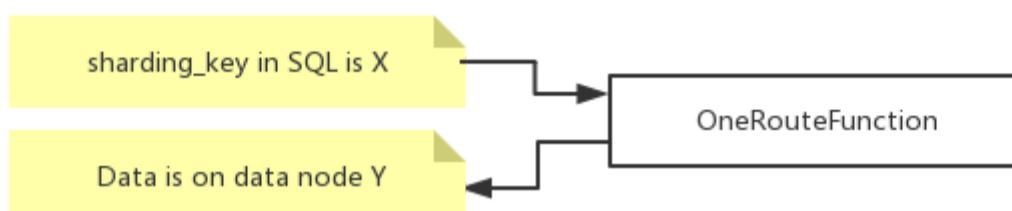
路由函数的加载发生在dble启动或重载时。



1. dble读取rule.xml时，根据用户配置的标签的class属性
2. dble通过Java的反射机制，从\$DBLE\_HOME/lib的jar包中，找到对应的jar（里的class文件），加载同名的类并创建对象
3. dble会逐个扫描中的标签，并根据name属性来调用路由函数的对应setter，以此完成赋值过程——例如，如果用户配置了2，那么dble就会尝试找路由函数中叫做setPartitionCount()的方法，并将字符串“2”传给它
4. dble调用路由函数的selfCheck()方法，执行函数编写者制定的检查动作，例如检查赋值得到的变量值是否有问题
5. dble调用路由函数的init()方法，执行函数编写者制定的准备动作，例如创建后面要用到的一些中间变量

#### 1.9.1.2 路由计算

路由函数接受用户SQL中的分片字段的值，计算出这个值对应的数据记录应该在哪个编号的数据分片（逻辑分片）上，DBLE从而知道把这个SQL准确发到这些分片上。



#### 1.9.1.3 参数查询

用户通过管理端口（默认9066），通过SHOW @@ALGORITHM WHERE SCHEMA=? AND TABLE=?来查询表上的路由算法时，dble调用路由算法的getAllProperties()方法，直接从内存中获取路由信息的配置。

```

mysql> show @@algorithm where schema=testdb and table=seqtest;
+-----+-----+
| KEY | VALUE |
+-----+-----+
| TYPE | SHARDING TABLE |
| COLUMN | ID |
| CLASS | com.actiontech.dble.route.function.PartitionByLong |
| partitionCount | 2 |
| partitionLength | 1 |
+-----+
5 rows in set (0.05 sec)

```

## 1.9.2 开发和部署

### 1.9.2.1 开发

开发时，理论上只需要引入AbstractPartitionAlgorithm抽象类和RuleAlgorithm接口及它们的依赖类就可以了。但实际上AbstractPartitionAlgorithm抽象类依赖了TableConfig类，由此开启了环游世界的依赖之旅。因此，现实的操作还是引用整个DBLE项目的源代码会比较直接方便。

开发一个新的路由函数时，必须给这个路由函数的开发新建项目，然后再引用DBLE项目（项目引用项目的方式）。而不应该直接打开DBLE的项目，然后在DBLE的项目里面直接新建源代码来直接开发（内嵌开发方式）。通过遵循这个做法，会有以下好处：

1. 路由函数可以独立打包，直接去看路由函数的jar包版本就能够确认函数版本；而把路由函数嵌到DBLE里的话，就很容易出现DBLE版本一样，但不清楚里面的函数是什么版本的窘况
2. 路由函数的递进可以更加自由，如果DBLE的AbstractPartitionAlgorithm抽象类和RuleAlgorithm接口没有变动，同一版本的路由函数可以延续使用好几个版本的DBLE，而不需要每次DBLE释放新版就得去重编译
3. 可以让路由函数中的受保护代码免受DBLE自身的开源协议影响

### 1.9.2.2 部署

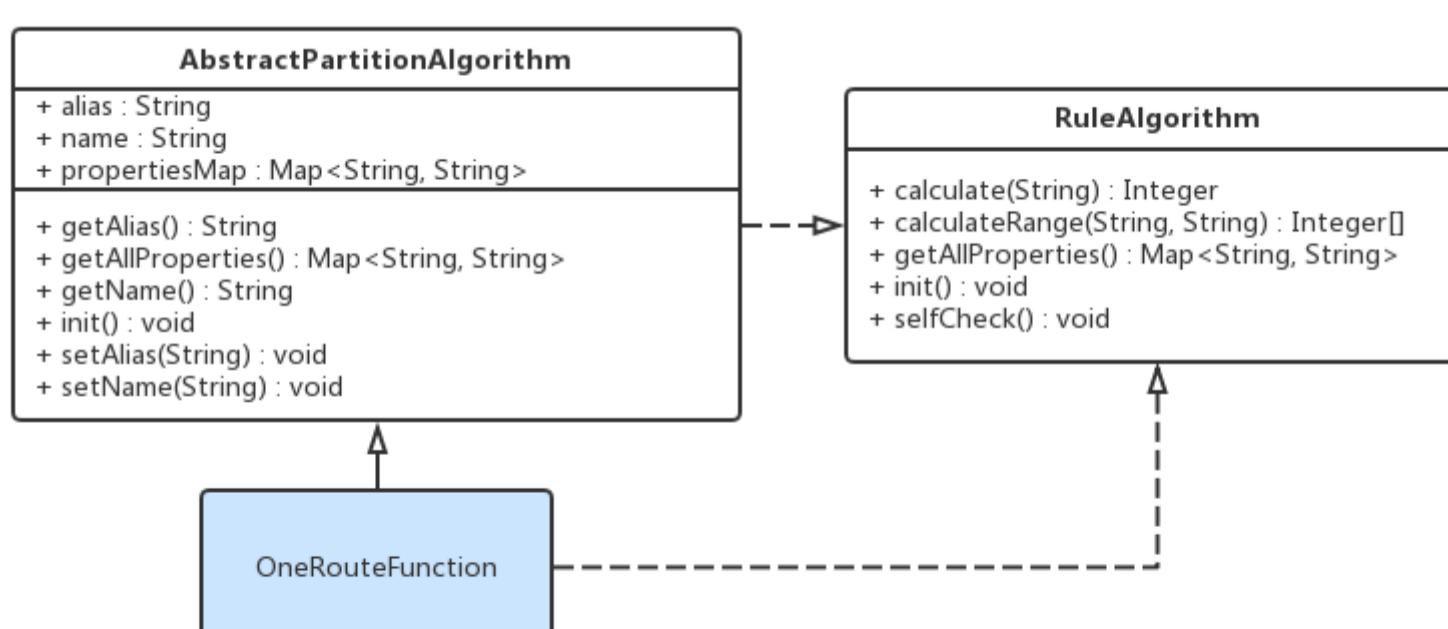
完成开发之后，成品打包成jar包进行发布，而不要直接发布class和依赖的library（其他项目的jar包或class文件）。

让DBLE使用上新的路由函数的过程：

1. 将成品jar包放入\$DBLE\_HOME/lib目录中
2. 调整jar包的所有者权限(chown)和文件权限(chmod)，使之与其他\$DBLE\_HOME/lib目录里的jar包一样
3. 按照原来的思路配置rule.xml，但需要注意标签的class属性必须要填写新的路由函数类的完全限定名（Fully Qualified Name），例如net.john.dble.route.functions.NewFunction
4. 配置逻辑表之类的必要信息，重启DBLE后，自动生效。

## 1.9.3 接口规范

每个路由函数本质上就是一个继承了AbstractPartitionAlgorithm抽象类，并且实现了RuleAlgorithm接口的一个类。下面以内置的com.actiontech.dble.route.function.PartitionByLong为例，介绍实现一个路由函数类所需要做的最小工作（必要工作）。



#### 1.9.3.1 配置项setters

在rule.xml中，我们需要配置partitionCount和partitionLength两个配置项。

```
<function name="hashmod" class="com">
<property name="partitionCount">4</property>
<property name="partitionLength">1</property>
</function>
```

为了让dble在函数加载过程中，能够认出这里的partitionCount（值为4）和partitionLength（值为1），因此PartitionByLong类中，就必须有属性设置方法（setter）setPartitionCount()和setPartitionLength()。而因为rule.xml是个文本型的XML文件，所以这些函数的传入参数就只能是一个String，数据类型转换和预处理的动作就由这些setter来处理了。

```
public void setPartitionCount(String partitionCount) {
    this.count = toIntArray(partitionCount);
    /* 参考本文的getAllProperties()的说明 */
    propertiesMap.put("partitionCount", partitionCount);
}

public void setPartitionLength(String partitionLength) {
    this.length = toIntArray(partitionLength);
    /* 参考本文的getAllProperties()的说明 */
    propertiesMap.put("partitionLength", partitionLength);
}
```

### 1.9.3.2 selfCheck()

在函数加载过程中，完成了配置项赋值之后，dble会调用这个路由函数对象的selfCheck()方法，让这个对象自我检查刚才读进来的配置项的值，放在一起是不是有问题。如果有问题的话，路由函数编写者在这时候，可以通过抛出RuntimeException来进行报错，并终止dble使用这个函数，当然，由于RuntimeException的霸道，dble自己也会因此而报错退出。

由于selfCheck()是RuleAlgorithm接口的要求，而且AbstractPartitionAlgorithm抽象类没实现它，对于想偷懒或者没有必要进行这个检查的人来说，还是需要自行定义一个空的同名方法来实现它。

```
@Override
public void selfCheck() {
```

### 1.9.3.3 init()

在函数加载过程的最后，dble调用这个路由函数对象的init()方法，让这个对象完成一些内部的初始化工作。

在我们的例子PartitionByLong里，通过init()方法准备了PartitionUtil对象，其中有一个哈希值的范围与逻辑分片号对应的数组，这样在后面的路由计算时就能通过查数组来加速得到结果。

```
@Override
public void init() {
    partitionUtil = new PartitionUtil(count, length);

    initHashCode();
}
```

### 1.9.3.4 calculate()和calculateRange()

dble执行用户SQL时，根据用户SQL的不同，调用calculate()或calculateRange()来确定用户的SQL应该发到哪个数据分片上去。

从IPO（Input-Process-Output）来分析，calculate()和calculateRange()的工作原理是一样的：

- **Input:** 用户SQL中的分片字段值
- **Output:** 用户SQL应该要发往的数据分片的编号
- **Process:** Input与Output转换的计算过程，由函数开发者编写

calculate()和calculateRange()的使用场景不同，导致它们存在着一些微小的差异。

函数名	调用场景	Input	Output
calculate()	用户SQL里分片字段的值是单值的情况，例如 ... WHERE sharding_key = 1	1个String	1个Integer
calculateRange()	用户SQL里分片字段的值是连续范围，例如 ... WHERE sharding_key BETWEEN 1 AND 5	2个String	Integer数组

```
@Override
public Integer calculate(String columnValue) {
    try {
        if (columnValue == null || columnValue.equalsIgnoreCase("NULL")) {
            return 0;
        }
        long key = Long.parseLong(columnValue);
        return calculate(key);
    } catch (NumberFormatException e) {
        throw new IllegalArgumentException("columnValue:" + columnValue + " Please eliminate any quote and non number within it.", e);
    }
}

@Override
public Integer[] calculateRange(String beginValue, String endValue) {
    long begin = 0;
    long end = 0;
    try {
        begin = Long.parseLong(beginValue);
        end = Long.parseLong(endValue);
    } catch (NumberFormatException e) {
        return new Integer[0];
    }
    int partitionLength = partitionUtil.getPartitionLength();
    if (end - begin >= partitionLength || begin > end) { //TODO: optimize begin > end
        return new Integer[0];
    }
    Integer beginNode = calculate(begin);
    Integer endNode = calculate(end);

    if (endNode > beginNode || (endNode.equals(beginNode) && partitionUtil.isSingleNode(begin, end))) {
        int len = endNode - beginNode + 1;
        Integer[] re = new Integer[len];
        for (int i = 0; i < len; i++) {
            re[i] = beginNode + i;
        }
        return re;
    } else {
        int split = partitionUtil.getSegmentLength() - beginNode;
        int len = split + endNode + 1;
        if (endNode.equals(beginNode)) {
            //remove duplicate
            len--;
        }
        Integer[] re = new Integer[len];
        for (int i = 0; i < split; i++) {
            re[i] = beginNode + i;
        }
        for (int i = split; i < len; i++) {
            re[i] = i - split;
        }
        return re;
    }
}
```

### 1.9.3.5 getAllProperties()

当用户找dble要路由函数的配置信息时，dble通过访问路由函数的getAllProperties()来获得一个<配置项, 配置值>的哈希表，然后将里面的内容逐项返回给用户。

getAllProperties()是RuleAlgorithm接口所规定要实现的，但为了简化编写新的路由函数的工作，在AbstractPartitionAlgorithm抽象类里，定义了propertiesMap这个私有变量，并且把“将propertiesMap交出去”作为实现了getAllProperties()方法的默认实现。一般来说，这个默认的实现能满足需求，而新路由函数编写者只需要在配置项setters处理用户配置时，将<配置项, 配置值>给put()进propertiesMap里就好了。

```
@Override
public Map<String, String> getAllProperties() {
    return propertiesMap;
}
```

#### 1.9.4 内置路由函数的缩写与类名对照表

DBLE内置的路由函数都位于com.actiontech.dble.route.function命名空间。但实际配置rule.xml的时候，却不用写那么长的完全限定名，这其实都是XMLRuleLoader类做了转换，因此实现了简写。下面就是7个内置函数的类名和它们的简写。

简写名	完整类名
date	PartitionByDate
enum	PartitionByFileMap
hash	PartitionByLong
jumpstringhash	PartitionByJumpConsistentHash
numberrange	AutoPartitionByLong
patternrange	PartitionByPattern
stringhash	PartitionByString

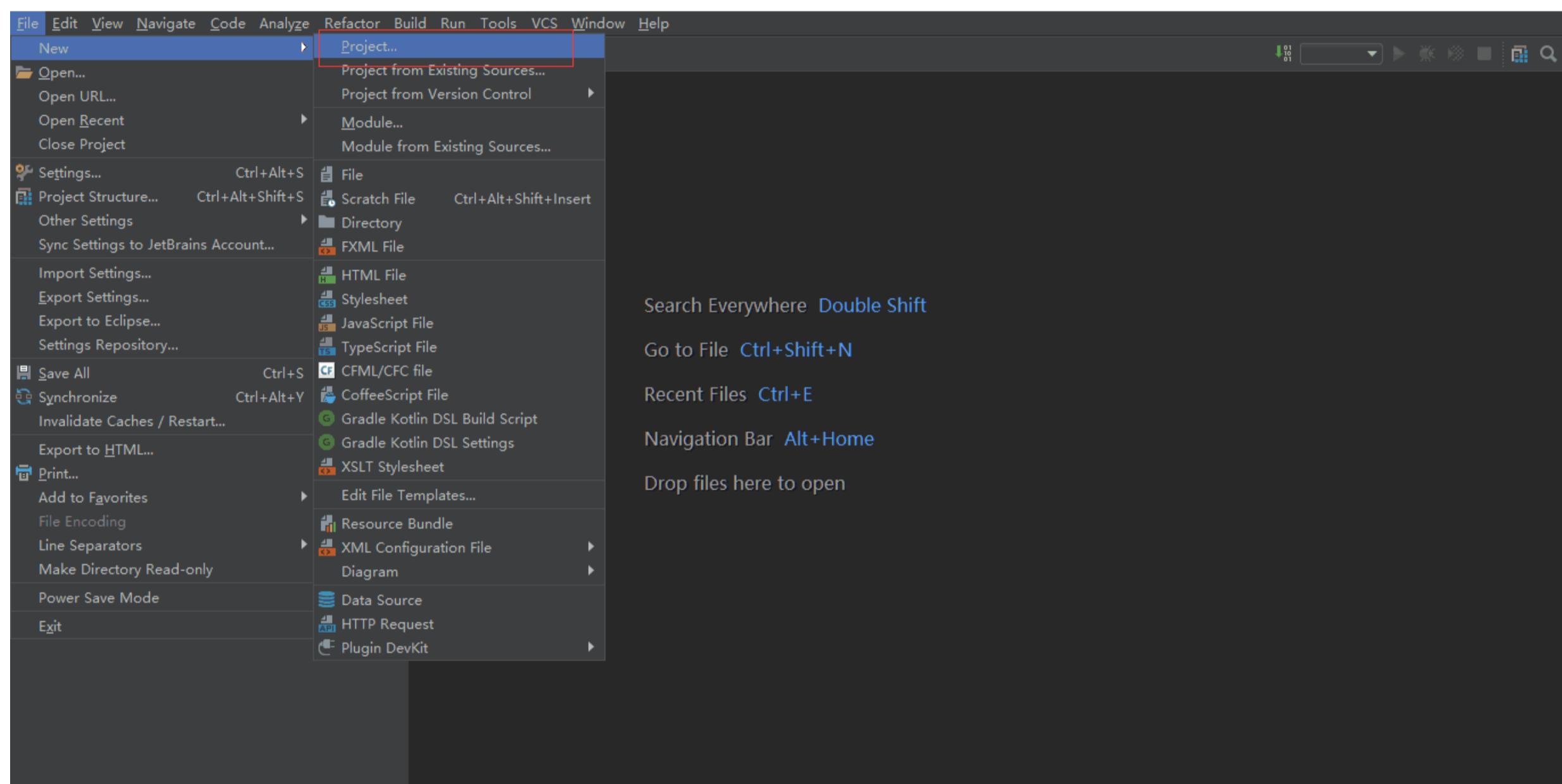
#### 1.9.5 IntelliJ IDEA中的实践

##### 1.9.5.0 前提

- 安装java开发环境
- 准备好dbe的最新的release版本jar包，以下是2.19.0.0 安装包下载链接：<https://github.com/actiontech/dble/releases/download/2.19.0.0%2Ftag/actiontech-dble-2.19.0.0.tar.gz>，解压后在lib目录中可以找到dbe的对应版本的jar包。

##### 1.9.5.1 创建 java 项目

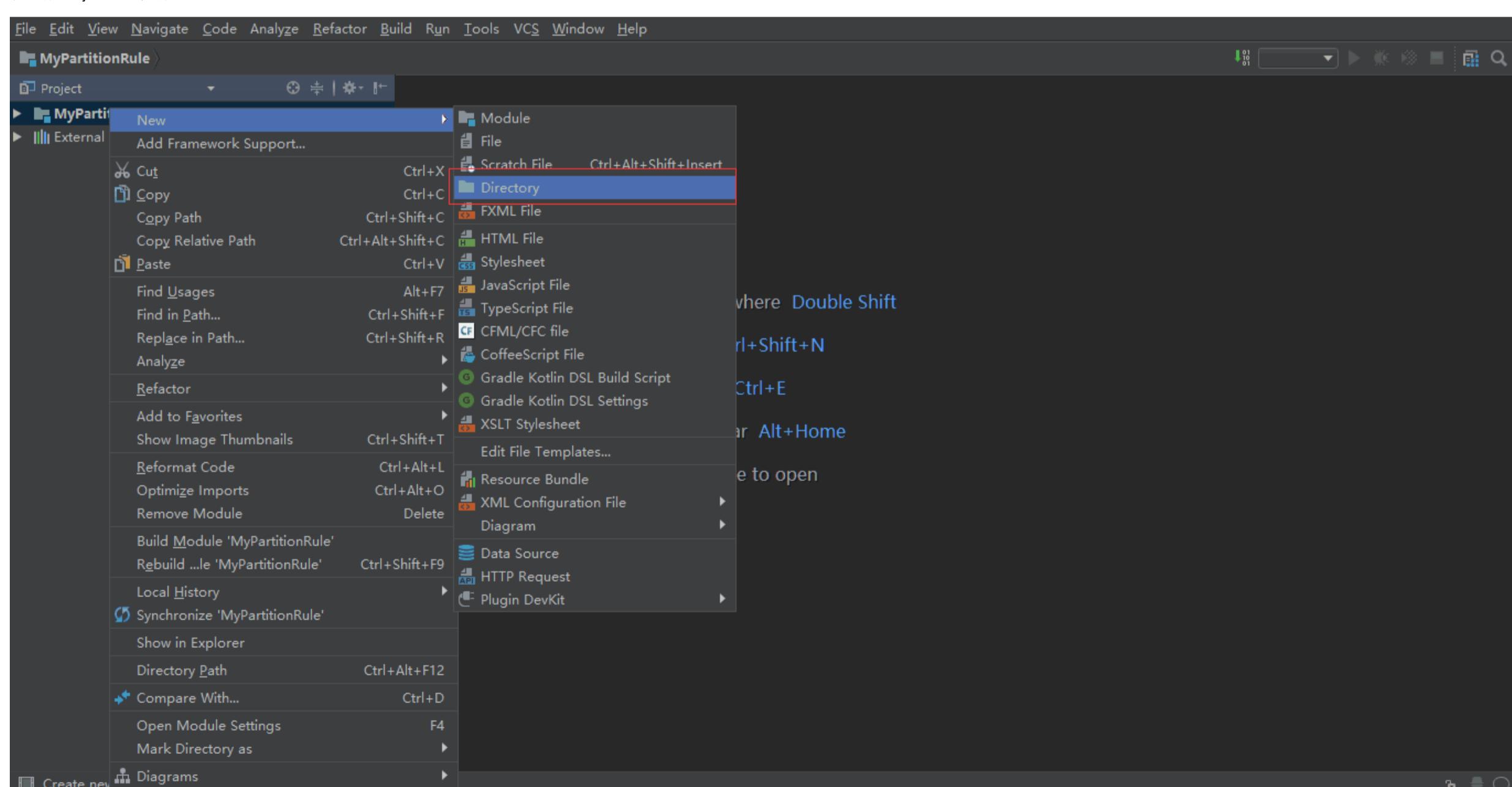
创建简单的java项目，点击红框中Project



点击Project后弹出框中默认选中java，若没有选中，手动选中后一路点击next，填写恰当的项目名即可。

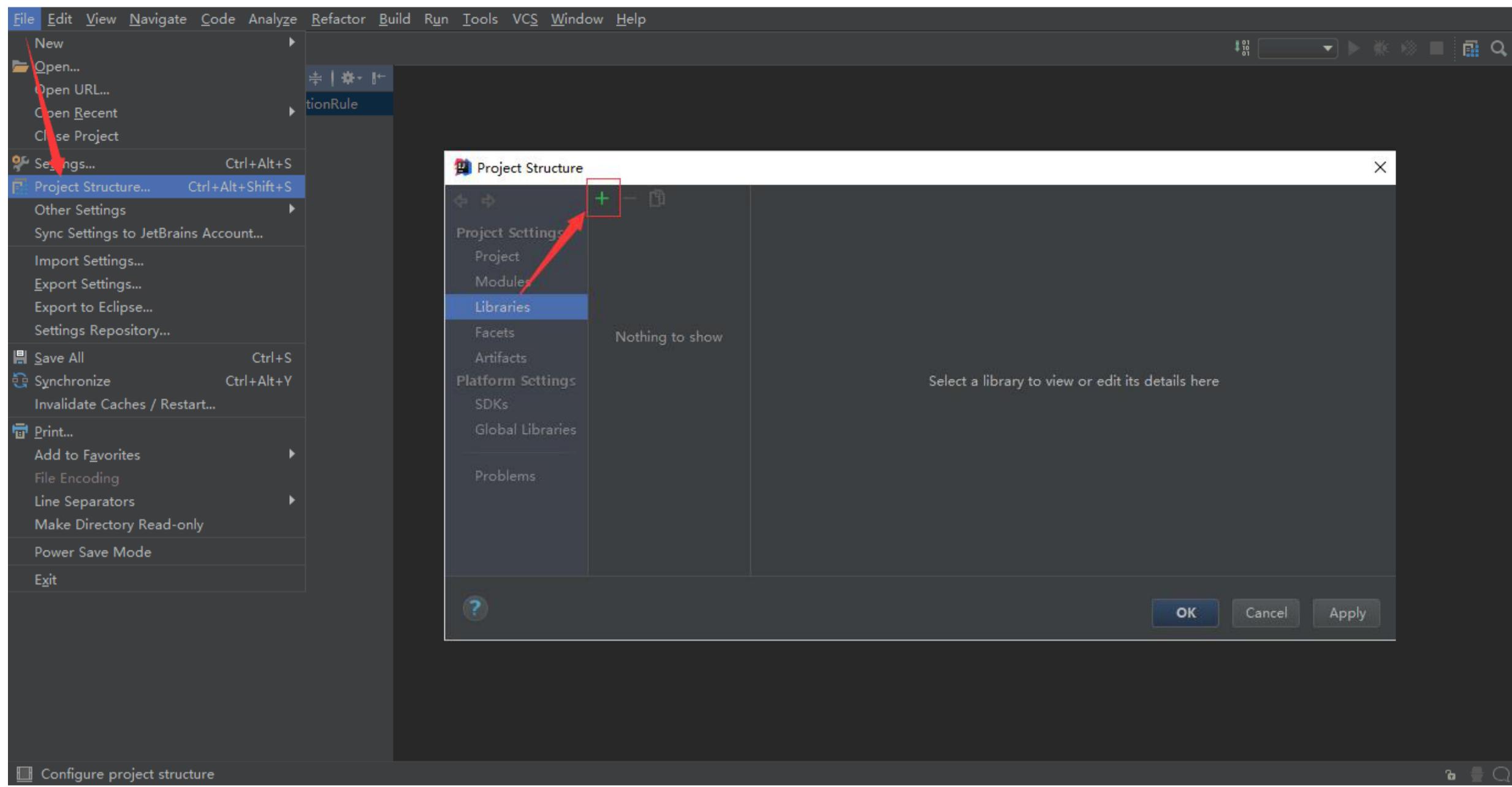
##### 1.9.5.2 引入 dble jar 包

右击新建的java项目，新建lib目录



将 dble 的jar包复制到lib目录下。

将lib中的dbe jar 包添加为项目的依赖。



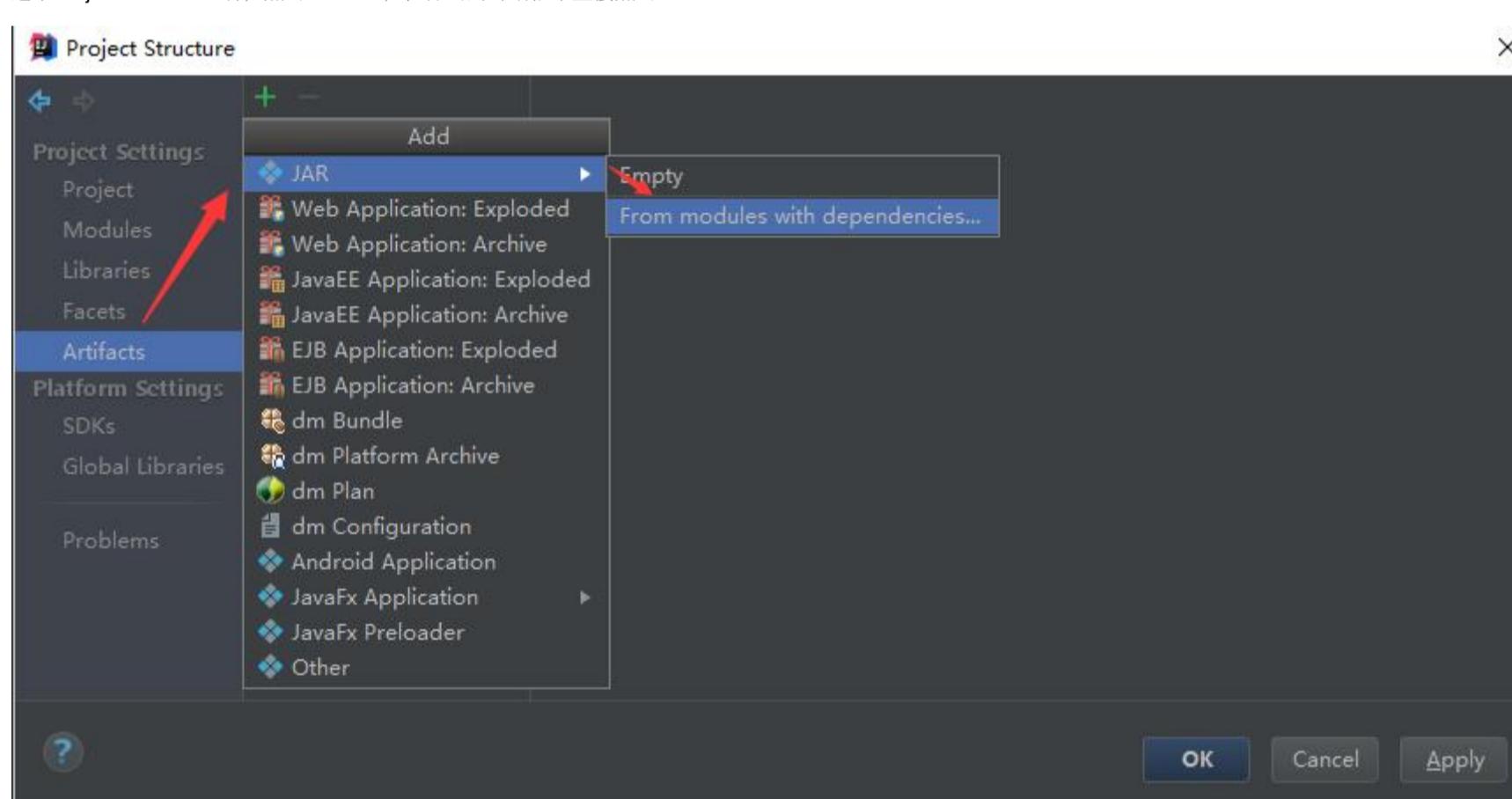
点击加号之后选中文件系统中的lib目录后加入即可

### 1.9.5.3 自定义分片算法类

在新建项目中自定义算法类，具体方式请参照上文。

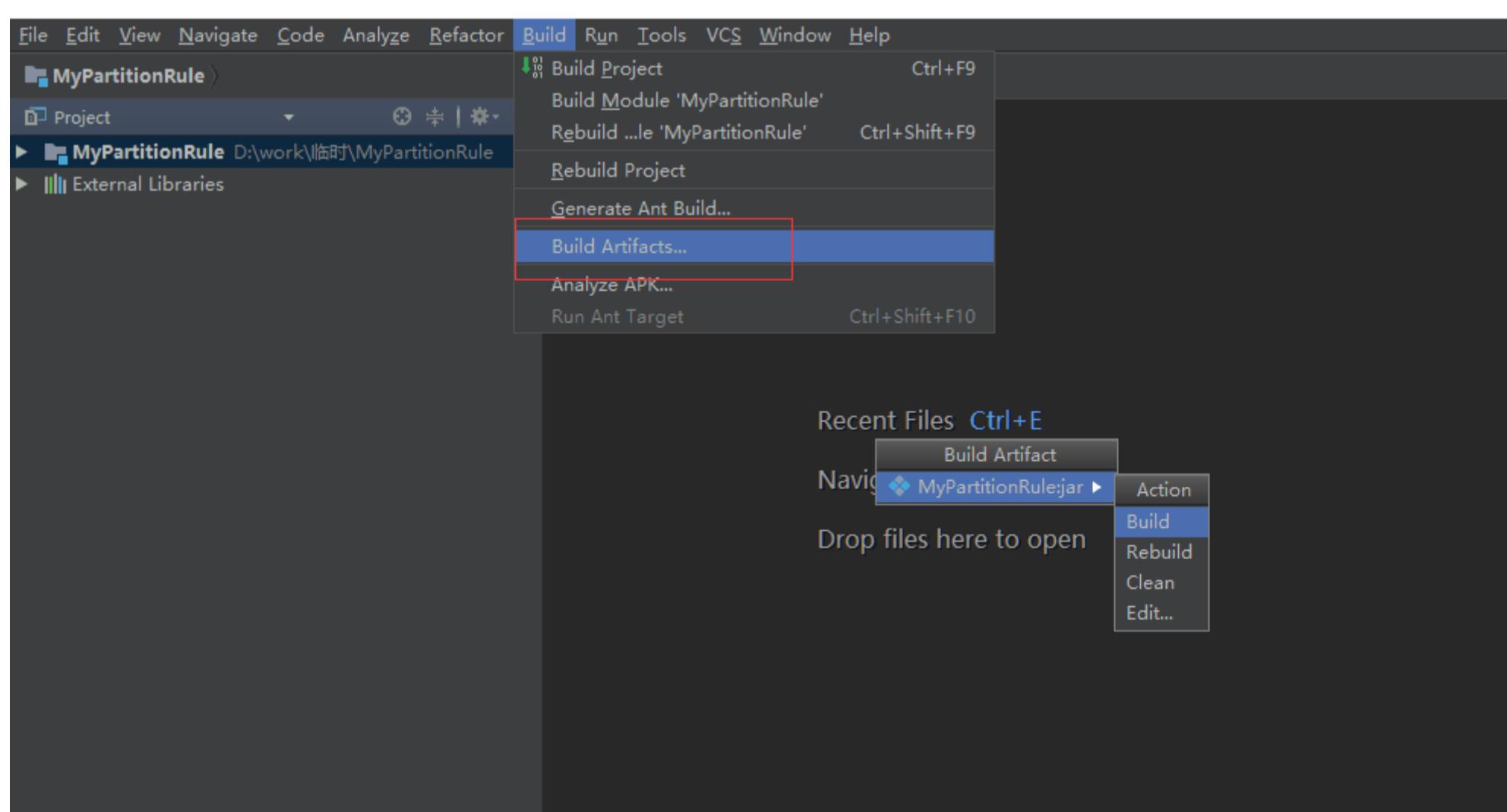
### 1.9.5.4 新建Artifacts

选中Project Strcuture后，点击Artifacts，在弹出的对话框中直接点击OK。



### 1.9.5.5 编译jar包

选择菜单栏中Build下的Build Artifacts，在弹出的对话框中选中自定义的Artifact后build。



编译完成之后，在项目路径下会生成一个out文件夹，在artifact子文件夹下可以找到生成的jar包。

## 1.10 配置文件版本变更

### 1.10.1 rule.xml

变更项	变更内容	变更版本	向后兼容性
新增	根version属性	2.19.01.0	兼容

### 1.10.2 schema.xml

变更项	变更内容	变更版本	向后兼容性	
新增	schema.table的incrementColumn属性，显式指定表格自增列	2.19.01.0	兼容	
新增	根version属性	2.19.01.0	兼容	
删除	schema.table的autoIncrement属性，自增与否只和incrementColumn有关	2.19.11.0	不兼容	
修改	schema.table的primaryKey属性修改为cacheKey，不再和自增属性有关	2.19.11.0	不兼容	
新增	schema.table的globalCheckClass属性，全局表检查类名	2.19.11.0	兼容	
新增	schema.table的cron属性，全局表检查周期	2.19.11.0	兼容	
删除	schema.datahost的switchType属性	2.20.04.0	不兼容	
变更	schema.datahost的balance属性	将原来的3取消，1等价于原来3的效果	兼容	新增 schema.table的globalCheck属性，全局表检查开关 2.19.11.0 兼容

### 1.10.3 server.xml

ID	变更版本	变更项	变更内容	向后兼容性
1	2.17.04.0	废弃	system的processorBufferPoolType	2.18.12.0以后不兼容
2	2.17.04.0	废弃	system的useStreamOutput	2.18.12.0以后不兼容
3	2.17.04.0	废弃	system的systemReserveMemorySize	2.18.12.0以后不兼容
4	2.17.04.0	废弃	system的sqlInterceptor	2.18.12.0以后不兼容
5	2.17.04.0	废弃	system的catletClassCheckSeconds	2.18.12.0以后不兼容
6	2.17.04.0	废弃	system的sqlInterceptorFile	2.18.12.0以后不兼容
7	2.17.04.0	废弃	system的defaultSqlParser	2.18.12.0以后不兼容
8	2.17.04.0	废弃	system的handleDistributedTransactions	2.18.12.0以后不兼容
9	2.17.04.0	新增	system的joinMemSize	兼容
10	2.17.04.0	新增	system的serverBacklog	兼容
11	2.17.04.0	新增	system的recordTxn	兼容
12	2.17.04.0	新增	system的xaSessionCheckPeriod	兼容
13	2.17.04.0	新增	system的xaLogCleanPeriod	兼容
14	2.17.04.0	新增	system的joinQueueSize	兼容
15	2.17.04.0	新增	system的mergeQueueSize	兼容
16	2.17.04.0	新增	system的orderByQueueSize	兼容
17	2.17.04.0	新增	system的useJoinStrategy	兼容
18	2.17.04.0	新增	system的nestLoopConnSize	兼容
19	2.17.04.0	新增	system的nestLoopRowsSize	兼容
20	2.17.04.0到2.17.11.0间存活	废弃	system的lowerCaseTableNames	2.18.12.0以后不兼容
21	2.17.08.0	新增	system的showBinlogStatusTimeout	兼容
22	2.17.09.0	废弃	system的useHandshakeV10	2.18.12.0以后不兼容
23	2.17.11.0	废弃	system的memoryPageSize	2.18.12.0以后不兼容
24	2.17.11.0	废弃	system的useOffHeapForMerge	2.18.12.0以后不兼容
25	2.17.11.0	废弃	system的数据节点排序临时目录	2.18.12.0以后不兼容
26	2.17.11.0	废弃	system的spillsFileBufferSize	2.18.12.0以后不兼容
27	2.17.11.0	新增	system的viewPersistenceConfBaseDir	兼容
28	2.17.11.0	新增	system的viewPersistenceConfBaseName	兼容
29	2.18.02.0	新增	system的orderMemSize	兼容
30	2.18.02.0	新增	system的otherMemSize	兼容
31	2.18.02.0	新增	system的useCostTimeStat	兼容
32	2.18.02.0	新增	system的最大成本统计大小	兼容
33	2.18.02.0	新增	system的成本采样百分比	兼容
34	2.18.02.0	新增	system的backendProcessors	兼容
35	2.18.02.0	新增	system的backendProcessorExecutor	兼容
36	2.18.02.0	新增	system的complexExecutor	兼容
37	2.18.02.0	新增	system的useThreadUsageStat	兼容
38	2.18.02.0	新增	system的usePerformanceMode	兼容
39	2.18.09.0	废弃	user的benchmark属性	不兼容
40	2.18.09.0	新增	system的最大连接数	兼容
41	2.18.09.0	新增	system的enableSlowLog	兼容
42	2.18.09.0	新增	system的慢查询日志基础目录	兼容
43	2.18.09.0	新增	system的慢查询日志基础名称	兼容
44	2.18.09.0	新增	system的flushSlowLogPeriod	兼容
45	2.18.09.0	新增	system的flushSlowLogSize	兼容
46	2.18.09.0	新增	system的sqlSlowTime	兼容
47	2.18.09.0	新增	user的最大连接数属性	兼容
48	2.18.10.0	新增	system的useOldMetalInit	兼容
49	2.19.01.0	新增	根version属性	兼容
50	2.19.03.0	新增	xaRetryCount	兼容
51	2.19.03.0	新增	maxCharsPerColumn	兼容
52	2.19.03.0	新增	maxRowSizeToFile	兼容
53	2.19.05.0	变更	system的最大连接数缺省值,由1024变为无限制	兼容
54	2.19.07.0	变更	system的charset的缺省值,由utf8变为utf8mb4	兼容
55	2.19.07.0	废弃	system的useOldMetalInit	不兼容
56	2.19.09.0	变更	system的最大包大小的缺省值,由16M变为4M	兼容
57	2.19.09.0	变更	system添加可选参数useOuterHa	兼容
58	2.19.11.0	变更	system添加参数autocommit	兼容
59	2.19.11.0	变更	system的bufferPoolPageNumber的缺省值,`由核数 × 20`变为`0.8 × MaxDirectMemorySize / bufferPoolPageSize(default 2M)`	兼容
60	2.20.04.0	删除	system的clusterHeartbeatPass属性	兼容
61	2.20.04.0	删除	system的clusterHeartbeatUser属性	兼容
62	2.20.04.0	删除	system的useZKSwitch属性	兼容
63	2.20.04.0	变更	system的sequnceHandlerType属性更名为sequenceHandlerType, 兼容旧值	兼容
64	2.20.04.0	变更	system的transactionRateSize属性更名为transactionRotateSize, 兼容旧值	兼容
65	2.20.04.0	变更	system的clearBigSqlResultSetMapMs属性更名为clearBigSQLResultSetMapMs, 兼容旧值	兼容
66	2.20.04.0	新增	system新增enableFlowControl属性	兼容
67	2.20.04.0	新增	system新增flowControlStartThreshold属性	兼容
68	2.20.04.0	新增	system新增flowControlStopThreshold属性	兼容
69	2.20.04.0	变更	system的useOuterHa默认属性值由false改为true	兼容

#### 1.10.4 wrapper.conf

变更项	变更内容	变更版本	向后兼容性
新增	gc相关配置-XX:+PrintHeapAtGC	2.19.01.0	兼容
新增	gc相关配置-XX:+PrintGCDates	2.19.01.0	兼容
新增	gc相关配置-Xloggc:./logs/gc.log	2.19.01.0	兼容
新增	gc相关配置-XX:+PrintGCTimeStamps	2.19.01.0	兼容
新增	gc相关配置-XX:+PrintGC	2.19.01.0	兼容
废弃	gc相关配置-XX:+PrintGC	2.19.03.0	兼容
新增	gc相关配置-XX:+PrintGCDetails	2.19.03.0	兼容
新增	OOM错误相关配置-XX:+HeapDumpOnOutOfMemoryError	2.20.04.0	兼容
新增	OOM错误相关配置-XX:HeapDumpPath=./logs/	2.20.04.0	兼容

### 1.10.5 cache配置

变更项	变更内容	变更版本	向后兼容性
新增	rocksdb	2.18.10.0	兼容

### 1.10.6 myid.properties

变更项	变更内容	变更版本	向后兼容性
删除项	zk模式的port属性	2.19.01.0	兼容
新增项	外部HA集群联动模式，可选参数clusterHa	2.19.09.0	兼容

### 1.10.7 schema/server/rule version 变更对照表

version版本号	dble变更版本	向后兼容性
2.19.01.0	2.19.01.0	兼容
2.19.03.0	2.19.03.0	兼容
2.19.05.0	2.19.05.0	兼容
1.0	2.19.07.0	兼容
2.0	2.19.11.0	不兼容
3.0	2.20.04.0	不兼容

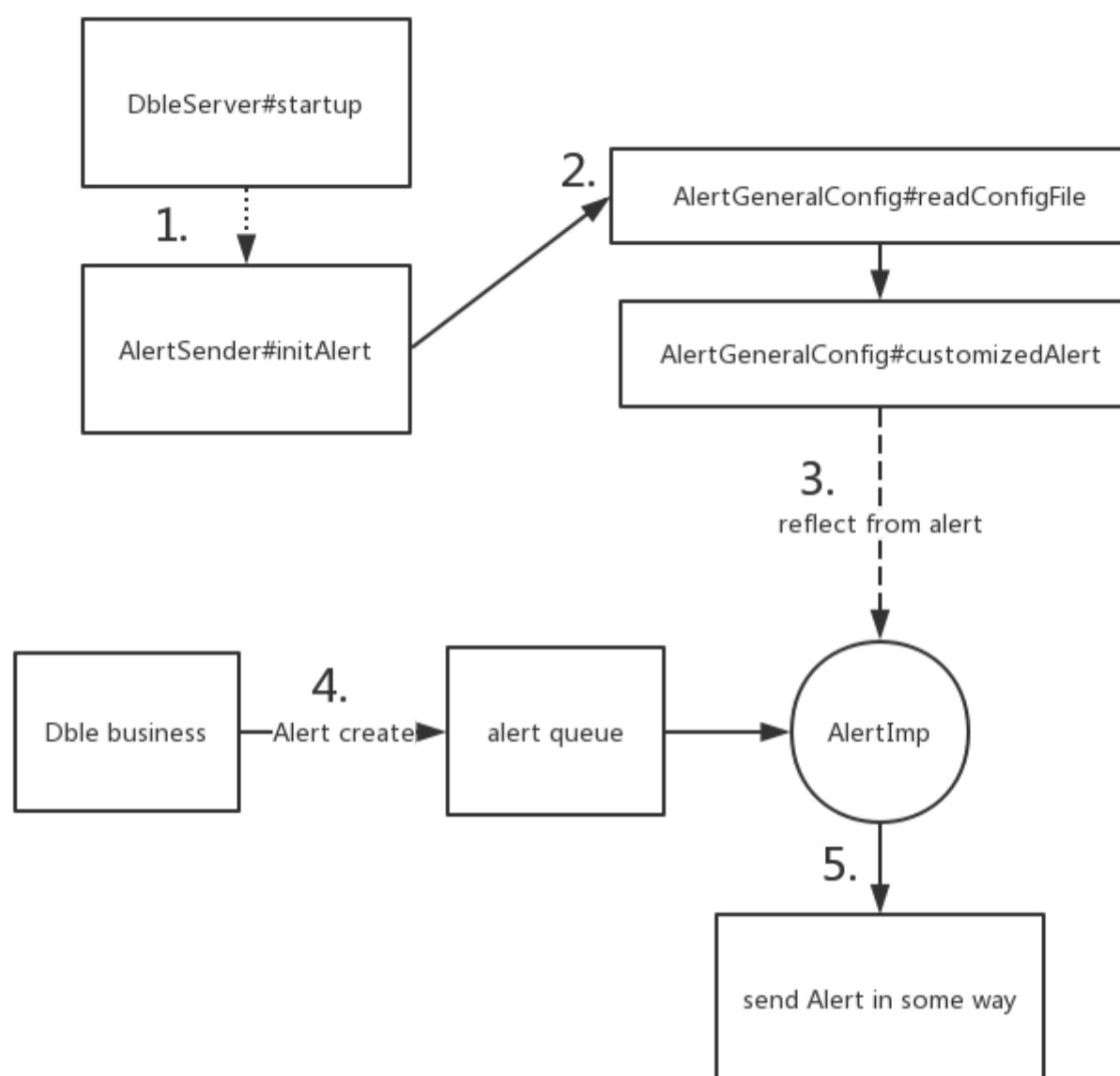
注意:在2.19.07.0版本中对于xml中的版本号进行了一次改造，新的版本号脱离dble版本号进行迭代，形式为X.Y。X代表配置的大版本号，Y代表配置的小版本号，当配置发生非兼容改变时X向上迭代，否则配置改变仅迭代Y

## 1.11 自定义告警模块

- 告警工作原理
  - 告警模块的加载
  - 告警发送和解除
  - 告警Interface Alert詳解
  - 告警发送对象詳解
- 自定义告警的开发和部署
  - 自定义告警的开发
  - 自定义告警的使用
  - 自定义告警示例以及示例文件

### 1.11.1 工作原理

#### 1.11.1.1 告警模块的加载



告警模块加载仅在dble启动的时候进行

1. dble启动，通过一系列加载之后，开始加载告警模块
2. 由对对象AlertGeneralConfig读取告警配置文件dble\_alert.properties
3. 根据读取到配置文件中的alert配置项进行对象反射，获取到用户自定义对象的实例
4. 当dble运行过程中发生告警或者告警解除事件，将告警任务放入任务队列中
5. 由一个dble线程循环消费队列中的内容，调用用户自定义对象对于告警任务进行发送

#### 1.11.1.2 告警发送和解除

在dble的告警模块中，将系统中的告警分成以下两种类型

- 可自动解决告警
- 不可自动解决告警

其中“可自动解决告警”为dble在运行过程中遭遇的偶发错误导致的告警，当错误被修正或者重试成功时告警被解除。例如：后端连接数到达最大值，XA事务提交单次失败等而“不可自动解决告警”则为dble不会重试也不能解决的错误，譬如：启动无法读取xa日志记录文件等

同时在dble中还将告警分成以下两种类型

- dble内部发生的错误告警
- dble和外部交互发生的错误告警

其中“dble内部发生的告警”指的是在dble对于内部任务在执行过程中发生的异常或者其他现象，比如写文件失败，kill后端连接失败等等而“dble和外部交互发生的错误告警”则指的是dble和外部具体节点交互时候的错误，譬如无法连接某后端节点，或者某后端节点心跳失败等等

根据以上两种分类，在dble的告警系统中存在着以下的四种告警信息

- 可自动解决-dble内部告警
- 不可自动解决-dble内部告警
- 可自动解决-dble与外部交互告警
- 不可自动解决-dble与外部交互告警

针对以上的四种告警的类别，我们可以发现，在dble的告警周期中，需要有以下四个逻辑上的方法：

- dble与外部交互告警发送方法
- dble内部告警发送方法
- dble与外部交互告警解除发送方法
- dble内部告警解除发送方法

#### 1.11.1.3 告警Interface Alert詳解

```

public interface Alert {
    void alertSelf(ClusterAlertBean bean);
    void alert(ClusterAlertBean bean);
    boolean alertResolve(ClusterAlertBean bean);
    boolean alertSelfResolve(ClusterAlertBean bean);
    void alertConfigCheck() throws Exception;
}
  
```

以上是整个告警接口的定义，可以看到除了一个比较特殊的方法alertConfigCheck（检查告警配置）其他的几个方法基本可以说是成对出现

- alert告警发送 --- dble与外部交互告警
- alertSelf告警发送 --- dble内部告警发送
- alertResolve告警解决发送 --- dble与外部告警解除发送
- alertSelfResolve内部告警解除发送 --- dble内部告警解除发送

四个信息发送的方法分别对应到上文对于几个dble告警系统中的方法需求

并在此解释下alert和alertSelf方法的区别：alert方法的调用输入中为ClusterAlertBean设定了alertComponentType和alertComponentId两个字段，而alertSelf方法则没有，设定适当的DBLEserver标识的动作交给Alert对象来完成当然在alertResolve和alertSelfResolve之间的区别也是同样的

#### 1.11.1.4 告警发送对象详解

```
public class ClusterAlertBean {
    String code; //告警的具体代码
    String level; //告警发生的具体级别
    String desc; //告警发送的详细描述
    String sourceComponentType; //告警发生的问题源头组件类型
    String sourceComponentId; //告警发送的问题源头组件ID
    String alertComponentType; //具体发送告警的组件类型
    String alertComponentId; //具体发送告警的组件ID
    String serverId; //告警发生的服务器ID
    long timestampUnix; //告警发送的时间戳
    long resolveTimestampUnix; //告警解除的时间，仅有解除的告警拥有此字段的值
    Map<String, String> labels; //告警的额外附加信息，补充信息
}
```

以上是在dble中对于告警对象的定义，每个字段的含义都在注释中进行了说明，在自定义告警发送的时候只需要对于告警对象中用户所关心的信息进行发送和处理即可

#### 1.11.2 自定义告警的开发和部署

##### 1.11.2.1 自定义告警的开发

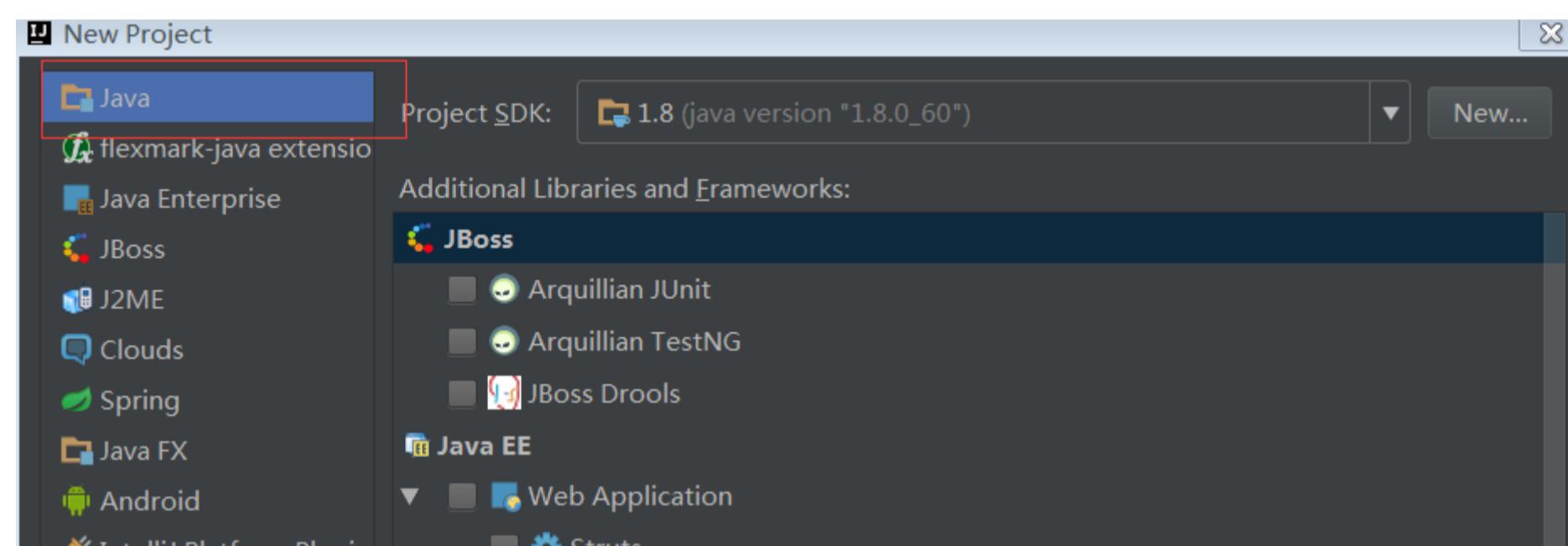
从上文中对于整个dble中告警系统的介绍，可以大致的将dble中的告警进行如下的总结：

- dble会在启动的时候根据配置文件中的信息加载一个Alert的对象
- dble的告警由dble内部的代码进行固定的触发，并将告警的任务发送到一个指定的告警队列中
- dble内部会有一个线程循环调用Alert的实现对象对于告警进行发送
- dble中的告警发送分为四种情况“内部告警”“外部告警”“内部告警解决”“外部告警解决”

按照上文中对于dble内部告警机制和原理的解释，可以整理出对于自定义告警的基本开发过程 ----- 创建一个Alert接口的实现，并且按照需要的方式将告警的信息发送到指定的地方

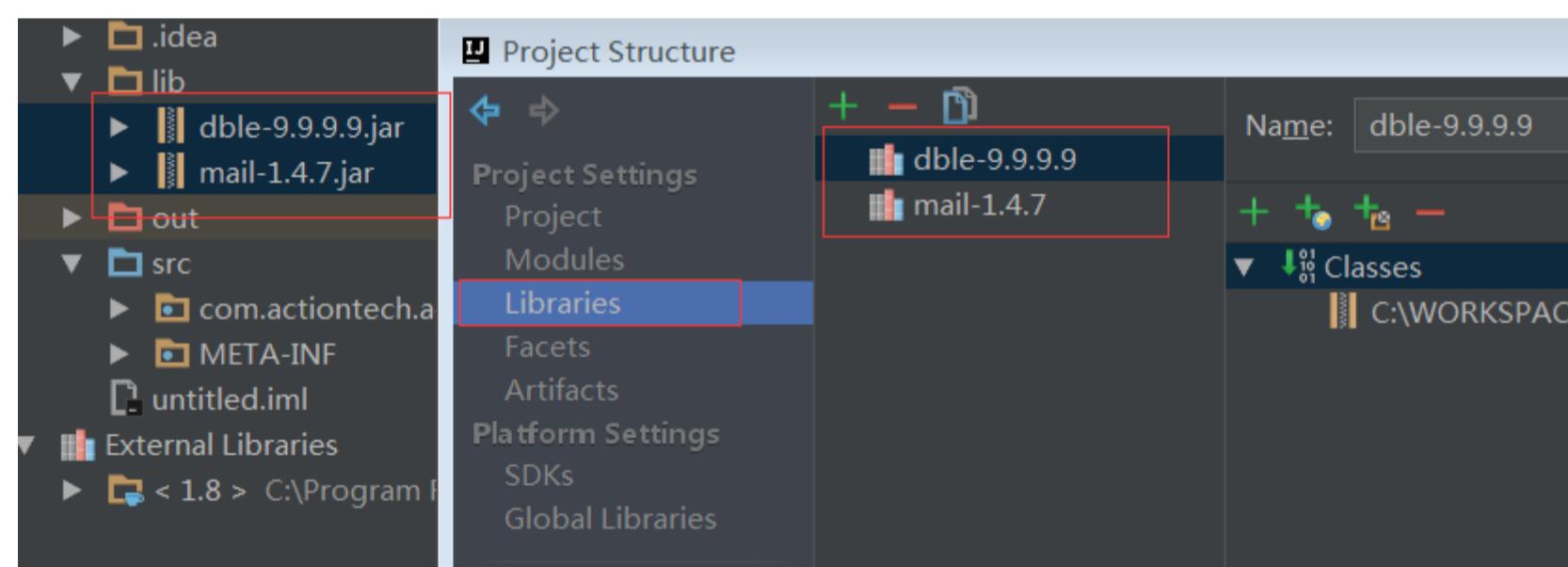
为了达成这个目的，可以按照以下的步骤进行逐一的开发，以下提供一个基于IDEA逐步操作的图文步骤，以一个通过邮件发送告警信息为例

##### 1 创建一个java项目



最简单的java项目就能完成这个告警自定义的开发

##### 2 将需要的依赖包copy到项目中，并且添加到lib



在此例中由于仅需额外依赖一个java的邮件发送包mail，所以这里仅导入了两个jar包作为依赖

##### 3.1 编写Alert实现类初始化方法

```
public MailAlert() {
    //init the mail data and read config file
    properties = AlertGeneralConfig.getInstance().getProperties();
}
```

在Alert实现类MailAlert中，可以通过调用AlertGeneralConfig.getInstance().getProperties();获取到dble启动时候加载到的告警配置文件(dble\_alert.properties)中的内容 所以，在这个地方如果自定义的告警需要使用到各种参数，都可以直接写在文件dble\_alert.properties的配置中，而在Alert实现类里面按照key值进行获取即可

##### 3.2 编写Alert的配置检查

```
@Override
public void alertConfigCheck() throws ConfigException {
    //check if the config is correct
    if (properties.getProperty(EMAIL_SENDER) == null
        || properties.getProperty(SENDER_PASSWORD) == null
        || properties.getProperty(EMAIL_SERVER) == null
        || properties.getProperty(EMAIL_RECEIVE) == null) {
        throw new ConfigException("alert check error, for some config is missing");
    }
}
```

在此例中由于是基于邮件进行告警发送，所以启动的时候需要在dble\_alert.properties文件中配置MAIL\_SENDER, SENDER\_PASSWORD, MAIL\_SERVER, MAIL\_RECEIVE四个选项 当配置不能满足自定义的这些选项的时候则抛出异常，dble启动加载的时候会打印对应的异常日志，并且切换为内置默认告警启动

上文中

```
properties.getProperty(EMAIL_RECEIVE)
```

这种方式则是从配置中获取自定义配置项的方法

##### 3.3 编写发送告警主体方法

由于此例中发送告警统一为邮件发送，包括告警和告警解决仅存在邮件内容中的差异，所以在这里仅用一个邮件发送方法send进行实现

```

private boolean sendMail(boolean isResolve, ClusterAlertBean clusterAlertBean) {
    try {
        Properties props = new Properties();
        props.setProperty("mail.debug", "true");
        props.setProperty("mail.smtp.auth", "true");
        props.setProperty("mail.host", properties.getProperty(MAIL_SERVER));
        props.setProperty("mail.transport.protocol", "smtp");

        MailSSLSocketFactory sf = new MailSSLSocketFactory();
        sf.setTrustAllHosts(true);
        props.put("mail.smtp.ssl.enable", "true");
        props.put("mail.smtp.socketFactory", sf);

        Session session = Session.getInstance(props);

        Message msg = new MimeMessage(session);
        msg.setSubject("DBLE告警 " + (isResolve ? "RESOLVE\n" : "ALERT\n"));
        String builder = new StringBuilder();
        builder.append(groupMailMsg(clusterAlertBean, isResolve));
        msg.setText(builder.toString());
        msg.setFrom(new InternetAddress(properties.getProperty(MAIL_SENDER)));

        Transport transport = session.getTransport();
        transport.connect(properties.getProperty(MAIL_SERVER), properties.getProperty(MAIL_SENDER), properties.getProperty(SENDER_PASSWORD));

        transport.sendMessage(msg, new Address[]{new InternetAddress(properties.getProperty(MAIL_RECEIVE))});
        transport.close();
        //send EMAIL SUCCESS return TRUE
        return true;
    } catch (Exception e) {
        e.printStackTrace();
    }
    //send fail return false
    return false;
}

private String groupMailMsg(ClusterAlertBean clusterAlertBean, boolean isResolve) {
    StringBuffer sb = new StringBuffer("Alert mail:\n");
    sb.append("    Alert type:" + clusterAlertBean.getCode() + " " + (isResolve ? "RESOLVE\n" : "ALERT\n"));
    sb.append("    Alert message:" + clusterAlertBean.getDesc() + "\n");
    sb.append("    Alert component:" + clusterAlertBean.getAlertComponentType() + "\n");
    sb.append("    Alert componentId:" + clusterAlertBean.getAlertComponentId() + "\n");
    sb.append("    Alert source:" + clusterAlertBean.getAlertComponentId() + "\n");
    sb.append("    Alert server:" + clusterAlertBean.getServerId() + "\n");
    sb.append("    Alert time:" + TimeStamp2Date(clusterAlertBean.getTimestampUnix()) + "\n");
    String detail = "|";
    if (clusterAlertBean.getLabels() != null) {
        for (Map.Entry<String, String> entry : clusterAlertBean.getLabels().entrySet()) {
            detail += entry.getKey() + ":" + entry.getValue();
        }
    }
    sb.append("    Other detail:" + detail + "|");
    return sb.toString();
}

```

具体如何使用java发送邮件我在此不进行赘述了，大概上面这段代码的原理就是从配置properties获取邮件配置的项，然后发送邮件到指定邮箱即可 而具体的groupMailMsg就是将告警发送过来的对象clusterAlertBean中的数据组织成邮件的文本内容，完全可以根据实际需要进行组织

### 3.4 实现具体的几个告警发送的方法

```

@Override
public void alertSelf(ClusterAlertBean clusterAlertBean) {
    alert(clusterAlertBean.setAlertComponentType(COMPARTMENT_TYPE).setAlertComponentId(properties.getProperty(COMPONENT_ID)));
}

@Override
public void alert(ClusterAlertBean clusterAlertBean) {
    clusterAlertBean.setSourceComponentType(COMPARTMENT_TYPE).
        setSourceComponentId(properties.getProperty(COMPONENT_ID)).
        setServerId(properties.getProperty(SERVER_ID)).
        setTimestampUnix(System.currentTimeMillis() * 1000000);
    sendMail(false, clusterAlertBean);
}

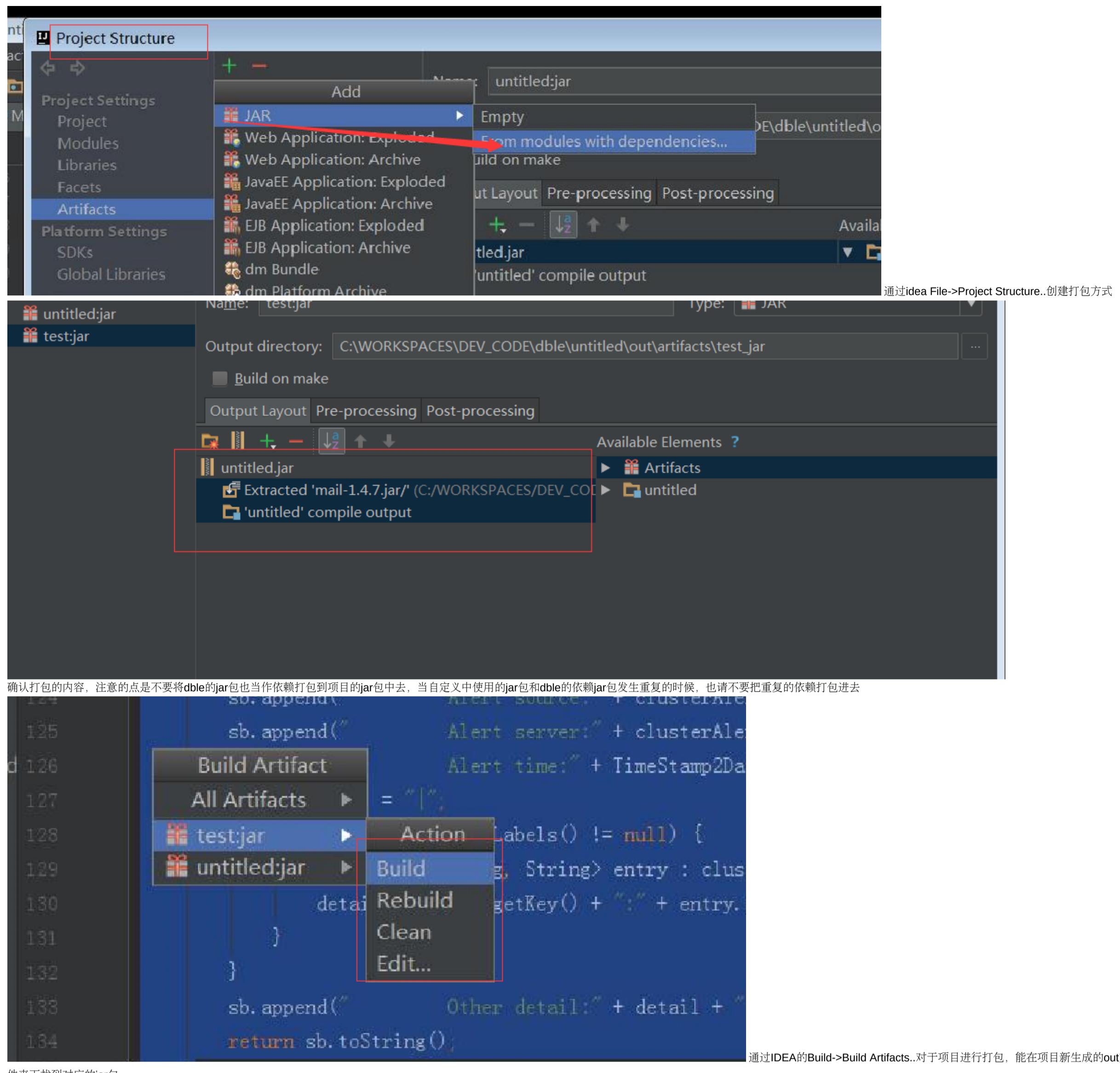
@Override
public boolean alertResolve(ClusterAlertBean clusterAlertBean) {
    clusterAlertBean.setSourceComponentType(COMPARTMENT_TYPE).
        setSourceComponentId(properties.getProperty(COMPONENT_ID)).
        setServerId(properties.getProperty(SERVER_ID)).
        setTimestampUnix(System.currentTimeMillis() * 1000000);
    return sendMail(true, clusterAlertBean);
}

@Override
public boolean alertSelfResolve(ClusterAlertBean clusterAlertBean) {
    return alertResolve(clusterAlertBean.setAlertComponentType(COMPARTMENT_TYPE).setAlertComponentId(properties.getProperty(COMPONENT_ID)));
}

```

以上方法实现就特别简单了，一般情况下都可以按照以上的示例进行实现，基本上就是补全部分缺失的信息，调用发送的方法即可

### 3.5 打包成jar进行使用



### 1.11.2.2 自定义告警的使用

有了上面打包出来的jar包之后就是该如何进行使用了，大致的步骤只有三个

- 将打包完成的jar包放入dble安装目录的lib目录下
- 配置文件dble\_alert.properties使得自定义的文件能被加载到，将类名配置到alert配置项中去
- 重启dble使得加载生效 在这里同样是衔接上文的案例，文件复制到lib下在此进行跳过 之后修改dble\_alert.properties配置文件，大致如下图所示

```
alert=com.actiontech.addtionAlert.MailAlert
mail_sender=123456789@qq.com
sender_pass=qwertryuiop
mail_server=smtp.qq.com
mail_receive=yyyyyyyyyyyy@actionsky.com
server_id=dble-server-001
component_id=DBLE-FOR-XXX-01
```

之后重启dble之后就会加载com.actiontech.addtionAlert.MailAlert这个类进行告警的发送了

### 1.11.2.3 自定义告警示例以及示例文件

在这里将上文中的实例以文件的形式给出，方便用户进行参考 [示例代码下载](#) [示例jar包下载](#)

### 1.11.3 附录：dble项目中各个告警CODE的含义

告警代码	解释	告警分类
DBLE_WRITE_TEMP_RESULT_FAIL	写出中间结果集到文件失败	内部/非自解决
DBLE_XA_RECOVER_FAIL	XA事务恢复失败	内部/非自解决
XA_READ_XA_STREAM_FAIL	读取XA事务记录文件失败	内部/非自解决
DBLE_XA_READ_DECODE_FAIL	解析XA事务记录文件失败	内部/非自解决
DBLE_XA_READ_IO_FAIL	读取XA事务记录文件失败	内部/非自解决
DBLE_XA_WRITE_IO_FAIL	写XA事务记录文件失败	内部/非自解决
DBLE_XA_WRITE_CHECK_POINT_FAIL	XA写出检查点信息失败	内部/自解决
DBLE_XA_BACKGROUND_RETRY_FAIL	XA后台重试提交失败	内部/自解决
DBLE_REACH_MAX_CON	节点连接数到达配置最大值，获取连接失败	内部/自解决
DBLE_TABLE_NOT_CONSISTENT_IN_DATAHOSTS	表结构在不同数据节点中不一致	内部/自解决
DBLE_TABLE_NOT_CONSISTENT_IN_MEMORY	表结构在数据节点中和dble内存中不一致	内部/自解决
DBLE_GLOBAL_TABLE_COLUMN_LOST	全局表检查列不存在	内部/自解决
DBLE_CREATE_CONN_FAIL	创建连接到后端mysql的连接失败	外部/自解决
DBLE_DATA_HOST_CAN_NOT_REACH	后端连接创建不可达	外部/非自解决
DBLE_KILL_BACKEND_CONN_FAIL	Kill后端连接执行失败	内部/非自解决
DBLE_NIOREACTOR_UNKNOWN_EXCEPTION	NIO意外报错	内部/非自解决
DBLE_NIOREACTOR_UNKNOWN_THROWABLE	NIO意外严重错误	内部/非自解决
DBLE_NIOPROXY_UNKNOWN_EXCEPTION	NIO后端连接创建器意外错误	内部/非自解决
DBLE_TABLE_LACK	配置中的表格未被创建	内部/自解决
DBLE_GET_TABLE_META_FAIL	获取表格的创建语句失败	内部/自解决
DBLE_TEST_CONN_FAIL	测试后端连接可达性失败	内部/非自解决
DBLE_HEARTBEAT_FAIL	心跳后端节点失败	外部/自解决
DBLE_DATA_NODE_LACK	缺少可用的数据节点	外部/自解决

## 自定义全局表检查

- 全局表一致性检查逻辑
- 全局表自定义方法详述
  - getCountSQL
  - getFetchCols
  - resultEquals
  - failResponse
  - resultResponse
- 自定义全局表检查操作步骤
- 自定义全局表检查配置

### 背景

全局表是dble中一种特殊类型的表格，一般来说认为在一个全局表table\_a所有分布的节点上，table\_a因同时满足以下两个条件：

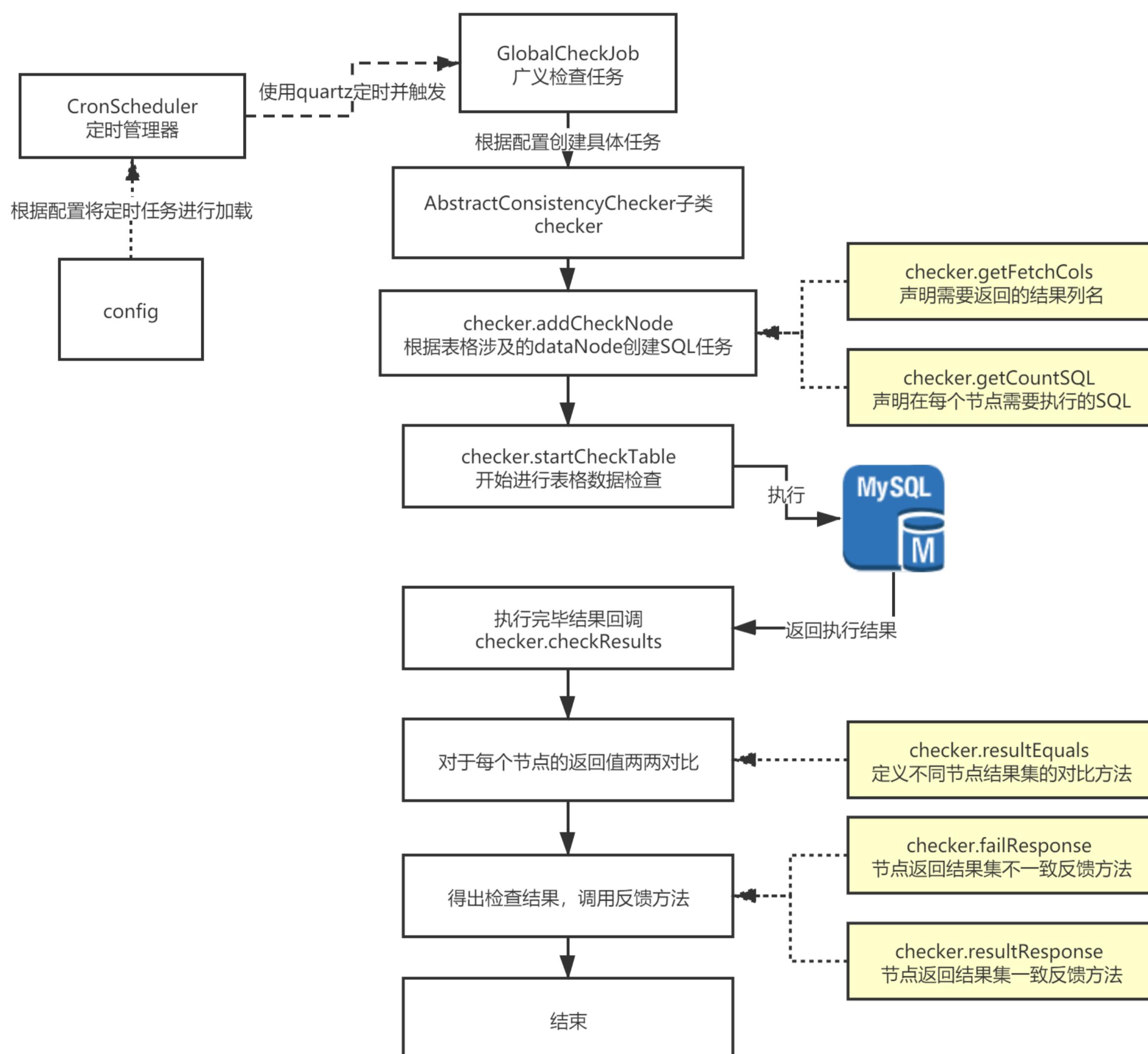
- 拥有相同的表格结构
- 拥有相同的表格数据

但事实上在系统和dble的运行过程中，可能由于一些不可避免分布式事务方面的误差，导致在长时间运行之后，不同节点上面的table\_a上面的数据不一致。为了及时的发现问题并不再造成更进一步的错误，dble中采用定时进行表格数据检查的方式对于table\_a中的数据一致性进行检查，并及时把检查的结果通知到运维人员。

### 工作原理大致逻辑

#### 全局表检查工作逻辑

全局表检查的大致逻辑如下图所示：



### 全局表检查方法详解

#### 1. 执行SQL定义String getCountSQL(String dbName, String tName)

**功能：**返回全局表检查需要对于表格执行的SQL内容

**输入：**SQL执行的MySQL中database的名称，所检查的表格的名称

**输出:** 检查具体需要执行的SQL

**举例:**

```
public String getCountSQL(String dbName, String tName) {
    //假如需要对于对应的table名字求checksum
    return "checksum table " + tName;
}
```

## 2.结果集定义 getFetchCols()

**功能:** 返回SQL执行完毕需要使用的结果集中的列名

**输入:** 无

**输出:** 需要收集的列名list

**举例:**

```
public String[] getFetchCols() {
    //checksum返回结果，我们只关心Checksum字段的返回值
    // mysql> checksum table suntest;
    //+-----+-----+
    //| Table | Checksum |
    //+-----+-----+
    //| db1.suntest | 1290812451|
    //+-----+-----+
    //所以return的内容只需要一个Checksum的列名即可
    return new String[]{"Checksum"};
}
```

## 3.SQL结果比较方法 boolean resultEquals(result1,result2)

**功能:** 用于判断两个不同节点的返回结果是否一致

**输入:** 不同节点的两个节点result, result1,result2

``` SQLQueryResult>> result

result | ----- row(List) | -----Key-Value(Field-Value) 例如checksum result | ----- row(List<1> checksum table suntest只有一行返回结果) | -----Key-Value(checksum - 1290812451 getFetchCols只取了一列)

\*\*输出: \*\* 需要收集的列名list

\*\*举例: \*\*

```
public boolean resultEquals(SQLQueryResult> or, SQLQueryResult> cr) { //因为checksum只有一行，并且即使表不存在也会有一行结果集 //所以直接取结果集的第一行即可 Map oresult = or.getResult().get(0); Map cresult = cr.getResult().get(0); //直接对比Map中checksum的值是不是一致即可 return (oresult.get("Checksum") == null && cresult.get("Checksum") == null) || (oresult.get("Checksum") != null && cresult.get("Checksum") != null && oresult.get("Checksum").equals(cresult.get("Checksum"))); }
```

##### \*\*4.失败行为接口 failResponse(resultList)\*\*

\*\*功能:\*\* 检查失败的通知/响应/其他自定义行为

\*\*输入:\*\* 检查结果列表

\*\*输出:\*\* 无

\*\*举例: \*\*

```
public void failResponse(List>> res) { //简单的情况下直接在日志中打印出对应信息 //如果有需要可以自行实现发邮件/发短信/发接口给告警系统等等 String errorMsg = "Global Consistency Check fail for table :" + schema + "." + tableName; System.out.println(errorMsg); for (SQLQueryResult>> r : res) { System.out.println("Checksum is :" + r.getResult().get(0).get("Checksum")); } }
```

##### \*\*5.其他结果通知 void resultResponse(errorList)\*\*

\*\*功能:\*\* 检查结果的通知/响应/其他自定义行为

\*\*输入:\*\* 检查(报错)结果列表

\*\*输出:\*\* 无

\*\*举例: \*\*

```
public void resultResponse(List>> elist) { //输入参数是检查过程中SQL执行报错的list，因为SQL自定义 //不同的检查SQL对于SQL报错的处理不同，具体报错应该别忽视 //或者应该视作不一致，由用户自己进行定义 String tableId = schema + "." + tableName;
```

```
if (elist.size() == 0) {
    System.out.println("Global Consistency Check success for table :" + schema + "-" + tableName);
} else {
    System.out.println("Global Consistency Check fail for table :" + schema + "-" + tableName);
    StringBuilder sb = new StringBuilder("Error when check Global Consistency, Table ");
    sb.append(tableName).append(" dataNode ");
    for (SQLQueryResult<List<Map<String, String>>> r : elist) {
        System.out.println("error node is : " + r.getTableName() + "-" + r.getDataNode());
        sb.append(r.getDataNode()).append(",");
    }
    sb.setLength(sb.length() - 1);
}
}
```

## 自定义全局表检查的开发和使用

#### 检查的自定义步骤

#### 1 创建一个java项目

![java项目创建](pic/create\_project\_global.png)

#### 2 将需要的依赖包copy到项目中，并且添加到lib

![添加依赖](pic/add\_library\_global.png)

#### 3 按照上一节的介绍逐个实现5个自定义方法

![方法实现](pic/function\_code\_global.png)

#### 4 打包成jar进行使用

![打包编译](pic/global\_package.png)

#### 5 示例文件及jar包

[示例代码下载](https://github.com/actiontech/dble-docs-cn/raw/master/1.config\_file/1.12\_customized\_global\_table\_check/CustomizeTest.java)

[示例jar包下载](https://github.com/actiontech/dble-docs-cn/raw/master/1.config\_file/1.12\_customized\_global\_table\_check/global\_customized.jar)

#### 自定义检查的配置

\*\*当前\*\*的全局表检查定义为schema.table级别，需要对于每个需要进行全局表一致性检查的表格进行配置，配置的下放带来一些繁琐的工作，但是却提供了一个重要的特性，用户可以根据不同全局表格的需要，或者是业务上面的特性，给与不同的全局表格不同的校验方式

\*\*注意:\*\* 检查方式的修改仅在reload或者重启之后生效\*\*\*

\*\*举例: \*\*

``` 自定义的jar包和其他dble内的自定义功能一样，将jar包放置于algorithm或者lib目录下就会在启动的时候被dble加载到，但是由于java中的类加载方式，如果由更新jar包内容和新增jar包的情况下，请先重启dble进程  
注意：当修改自定义jar包的时候请重启dble，此时reload可能无法得到预期的结果

## 2.功能描述

- 2.1 管理端命令
  - 2.1.1 select命令
  - 2.1.2 set命令
  - 2.1.3 show命令
  - 2.1.4 switch命令
  - 2.1.5 kill命令
  - 2.1.6 stop命令
  - 2.1.7 reload命令
  - 2.1.8 rollback命令
  - 2.1.9 offline命令
  - 2.1.10 online命令
  - 2.1.11 file命令
  - 2.1.12 log命令
  - 2.1.13 配置检查命令
  - 2.1.14 pause & resume 命令
  - 2.1.15 慢查询日志相关命令
  - 2.1.16 创建/删除物理库命令
  - 2.1.17 check @@metadata命令
  - 2.1.18 release @@reload\_metadata命令
- 2.2 全局序列
  - 2.2.1 MySQL offset-step方式
  - 2.2.2 时间戳方式
  - 2.2.3 分布式时间戳方式
  - 2.2.4 分布式offset-step方式
- 2.3 读写分离
- 2.4 注解
- 2.5 分布式事务
  - 2.5.1 XA事务概述
  - 2.5.2 XA事务的提交以及回滚
  - 2.5.3 XA事务的后续补偿以及日志清理
  - 2.5.4 XA事务的记录
  - 2.5.5 一般分布式事务概述
- 2.6 连接池管理
- 2.7 内存管理
- 2.8 集群同步协调&状态管理
- 2.9 grpc 告警
- 2.10 表meta数据管理
  - 2.10.1 Meta信息初始化
  - 2.10.2 Meta信息维护
  - 2.10.3 一致性检测
  - 2.10.4 View Meta
- 2.11 统计管理
  - 2.11.1 查询条件统计
  - 2.11.2 表状态统计
  - 2.11.3 用户状态统计
  - 2.11.4 命令统计
  - 2.11.5 heartbeat统计
  - 2.11.6 网络读写统计
- 2.12 故障切换
- 2.13 前后端连接检查
- 2.14 ER表
- 2.15 global表
- 2.16 缓存的使用
- 2.17 执行计划
- 2.18 性能观测和调整
- 2.19 智能计算reload
- 2.20 慢查询日志
- 2.21 单条SQL性能trace
- 2.22 KILL @@DDL\_LOCK
- 2.23 外部高可用联动
- 2.24 超时控制
- 2.25 流量控制

## 2.1 管理端命令集

- [2.1.1 select 命令](#)
- [2.1.2 set 命令](#)
- [2.1.3 show 命令](#)
- [2.1.4 switch 命令](#)
- [2.1.5 kill 命令](#)
- [2.1.6 stop 命令](#)
- [2.1.7 reload 命令](#)
- [2.1.8 rollback 命令](#)
- [2.1.9 offline 命令](#)
- [2.1.10 online 命令](#)
- [2.1.11 file 命令](#)
- [2.1.12 log 命令](#)
- [2.1.13 配置检查命令](#)
- [2.1.14 pause & resume 命令](#)
- [2.1.15 慢查询日志相关命令](#)
- [2.1.16 create database 命令](#)
- [2.1.17 check @@metadata 命令](#)
- [2.1.18 release @@reload\\_metadata 命令](#)
- [2.1.19 split 命令](#)
- [2.1.20 flow\\_control 命令](#)

### 2.1.1 select 命令

#### 2.1.1.1 select @@VERSION\_COMMENT

```
select @@VERSION_COMMENT;
```

描述: 查询db的版本信息;  
例:

```
MySQL [(none)]> select @@VERSION_COMMENT;
+-----+
| @VERSION_COMMENT      |
+-----+
| db Server (ActionTech) |
+-----+
1 row in set (0.02 sec)
```

列描述:  
略

#### 2.1.1.2 select @@SESSION.TX\_READ\_ONLY

```
select @@SESSION.TX_READ_ONLY;
```

描述: 无实际意义, 仅为了支持驱动连接管理端时下发的上下文  
结果: 定值, 永远返回0

#### 2.1.1.3 select @@max\_allowed\_packet

```
select @@max_allowed_packet;
```

描述: 无实际意义, 仅为了支持驱动连接管理端时下发的上下文  
结果: 定值, 永远返回1048576

## 2.1.2 set xxx

**set xxx;**

其中，xxx为要设置的变量。

描述：无实际意义，仅为了支持驱动连接管理端时下发的上下文

结果：永远返回OK

### 2.1.3 show命令

#### 2.1.3.1 show @@time.current

```
show @@time.current;
```

描述: 展示系统当前时间

结果: 略

#### 2.1.3.2 show @@time.startup

```
show @@time.startup;
```

描述: 展示系统启动时间

结果: 略

#### 2.1.3.3 show @@version

```
show @@version;
```

描述: 展示dble版本

结果: 略

#### 2.1.3.4 show @@server

```
show @@server;
```

描述: dble的当前信息

例:

```
mysql> show @@server;
+-----+-----+-----+-----+-----+-----+
| UPTIME | USED_MEMORY | TOTAL_MEMORY | MAX_MEMORY | RELOAD_TIME | ROLLBACK_TIME | CHARSET | STATUS |
+-----+-----+-----+-----+-----+-----+
| 1h 4m 47s | 17414592 | 87031808 | 1840250880 | 2017/10/17 16:42:09 | -1 | utf8 | ON |
+-----+-----+-----+-----+-----+-----+
1 row in set (0.05 sec)
```

列描述:

UPTIME: 服务已经启动时间  
USED\_MEMORY: 已使用堆内存  
TOTAL\_MEMORY: 总共的堆内存  
MAX\_MEMORY: 最大可用堆内存  
RELOAD\_TIME: 上次config加载时间  
ROLLBACK\_TIME: 上次config的rollback时间  
CHARSET: 字符集  
STATUS: 在线状态

#### 2.1.3.5 show @@threadpool

```
show @@threadpool;
```

描述: 展示当前线程池信息

例:

```
mysql> show @@threadpool;
+-----+-----+-----+-----+-----+-----+
| NAME | POOL_SIZE | ACTIVE_COUNT | TASK_QUEUE_SIZE | COMPLETED_TASK | TOTAL_TASK |
+-----+-----+-----+-----+-----+-----+
Timer	1	0	0	22596	22596
BusinessExecutor	8	1	0	216	217
complexQueryExecutor	0	0	0	0	0
+-----+-----+-----+-----+-----+
3 rows in set (0.03 sec)
```

列描述:

NAME: 线程池名称  
POOL\_SIZE: 线程池大小  
ACTIVE\_COUNT: 活动数量  
TASK\_QUEUE\_SIZE: 队列中的数量  
COMPLETED\_TASK: 已完成的任务数量  
TOTAL\_TASK: 总共任务数量

#### 2.1.3.6 show @@database

```
show @@database;
```

描述: 展示配置的schema名字

结果: 略

#### 2.1.3.7 show @@datanode

```
show @@datanode;
```

描述: 展示配置的所有datanode信息

例:

```
mysql> show @@datanode;
+-----+-----+-----+-----+-----+-----+
| NAME | DATHOST | SCHEMA_EXISTS | ACTIVE | IDLE | SIZE | EXECUTE | RECOVERY_TIME |
+-----+-----+-----+-----+-----+-----+
dn1	dh1/dble_test	true	0	0	1000	34	-1
dn2	dh2/dble_test	true	0	0	1000	34	-1
dn3	dh1/dble2_test	false	0	0	1000	26	-1
dn4	dh2/dble2_test	true	0	0	1000	26	-1
dn5	dh1/nosharding	true	0	0	1000	9	-1
+-----+-----+-----+-----+-----+-----+
5 rows in set (0.09 sec)
```

列描述:

NAME: 名称  
DATHOST: hostName/实际schema  
SCHEMA\_EXISTS: 对应后端物理库是否存在, true为存在, false为不存在。  
ACTIVE: 当前活动的后端连接数量  
IDLE: 当前空闲的后端连接数量(空闲容量维护疑似bug)  
SIZE: maxCon容量  
EXECUTE: 有过活动的后端连接数量统计  
RECOVERY\_TIME: 恢复心跳还需要秒数(stop @@heartbeat 中设置)

如果要查看某个schema相关的datanode信息, 执行:

```
show @@datanode where schema=xxx;
```

其中, xxx为要查看的schema的名字。

### 2.1.3.8 show @@datasource

```
show @@datasource;
```

描述: 展示配置的所有datasource信息

例:

```
mysql> show @@datasource;
+-----+-----+-----+-----+-----+-----+-----+-----+
| DATAHOST | NAME   | HOST    | PORT  | W/R   | ACTIVE | IDLE  | SIZE   |
+-----+-----+-----+-----+-----+-----+-----+-----+
| localhost2 | hostS1 | 10.18x.2x.63 | 3307 | W    | 1     | 9    | 100   |
| localhost1 | hostM1 | 10.18x.2x.64 | 3306 | W    | 1     | 9    | 100   |
+-----+-----+-----+-----+-----+-----+-----+-----+
2 rows in set (0.09 sec)
```

列描述

```
DATAHOST: dataSource所属datahost
NAME: dataSource名称
HOST: 主机名/ip
PORT: 端口
W/R: 读写结点标识
ACTIVE: 当前活动的后端连接数量, 按照dataSource统计
IDLE: 当前空闲的后端连接数量, 按照dataSource统计(空闲容量维护疑似bug)
SIZE: maxCon容量
EXECUTE: 有活动的后端连接数量统计, 按照dataSource统计
READ_LOAD: 读负载 (请求次数)
WRITE_LOAD: 写负载 (请求次数)
DISABLED: schema.xml中datasource中的配置 (2.19.09.0以前的版本没有此列, disabled为true的结点不显示)
```

如果要查看某个datanode对应的datasource信息, 执行:

```
show @@datasource where dataNode=xxx;
```

其中, xxx为要查看的datanode的名字。

### 2.1.3.9 show @@datasource.synstatus

```
show @@datasource.synstatus;
```

描述: 展示当前各datasource的同步信息

前提条件: heartbeat 配置了 show slave status (参见1.2 schema.xml)

例:

```
mysql> show @@datasource.synstatus;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| DATAHOST | NAME   | HOST    | PORT  | MASTER_HOST | MASTER_PORT | MASTER_USER | SECONDS_BEHIND_MASTER | SLAVE_IO_RUNNING | SLAVE_SQL_RUNNING | SLAVE_IO_STATE | CONNECT_RETRY | LAST_IO_ERROR |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| localhost1 | hostM1 | 10.18x.2x.63 | 3320 | 10.18x.2x.61 | 3320 | qrep | NULL | No | No | | 60 | |
| localhost2 | hostM2 | 10.18x.2x.64 | 3320 | 10.18x.2x.62 | 3320 | qrep | NULL | No | No | | 60 | |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
2 row in set (0.00 sec)
```

列描述:

```
DATAHOST: dataSource所属datahost
NAME: dataSource名称
HOST: 主机名/ip
PORT: 端口
```

其余列含义参见mysql中show slave status的命令。

### 2.1.3.10 show @@datasource.syndetail where name=?

```
show @@datasource.syndetail where name=xxx;
```

其中, xxx为要查看的datasource的名字。

描述: 展示24小时内各datasource的历次同步信息

例:

```
mysql> show @@datasource.syndetail WHERE name =hostM2;
+-----+-----+-----+-----+-----+-----+-----+-----+
| DATAHOST | NAME   | HOST    | PORT  | MASTER_HOST | MASTER_PORT | MASTER_USER | TIME           | SECONDS_BEHIND_MASTER |
+-----+-----+-----+-----+-----+-----+-----+-----+
localhost2	hostM2	10.18x.2x.64	3320	10.18x.2x.62	3320	qrep	2017-10-17 18:31:27	-1
localhost2	hostM2	10.18x.2x.64	3320	10.18x.2x.62	3320	qrep	2017-10-17 18:31:57	-1
localhost2	hostM2	10.18x.2x.64	3320	10.18x.2x.62	3320	qrep	2017-10-17 18:32:27	-1
localhost2	hostM2	10.18x.2x.64	3320	10.18x.2x.62	3320	qrep	2017-10-17 18:32:57	-1
+-----+-----+-----+-----+-----+-----+-----+-----+
4 row in set (0.05 sec)
```

列描述:

```
DATAHOST: dataSource所属datahost
NAME: dataSource名称
HOST: 主机名/ip
PORT: 端口
```

其余列含义参见mysql中show slave status的命令。

### 2.1.3.11 show @@datasource.cluster

```
show @@datasource.cluster;
```

描述: 此功能在2.20.04.0 版本已经废除。

### 2.1.3.12 show @@processor

```
show @@processor;
```

描述: 展示dble实例的processor信息

例:

```
mysql> show @@processor ;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| NAME   | NET_IN | NET_OUT | REACT_COUNT | R_QUEUE | W_QUEUE | FREE_BUFFER | TOTAL_BUFFER | BU_PERCENT | BU_WARNS | FC_COUNT | BC_COUNT |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
Processor0	6834	2749	0	0	0	687194767360	687194767360	0	0	0	0
Processor1	7221	2862	0	0	0	687194767360	687194767360	0	0	0	0
Processor2	6830	31141	0	0	0	687194767360	687194767360	0	0	1	0
Processor3	6375	2681	0	0	0	687194767360	687194767360	0	0	0	0
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
4 rows in set (0.06 sec)
```

列描述:

```

NAME: 名称
NET_IN: 接收流量
NET_OUT: 发送流量
REACT_COUNT: 固定值0
R_QUEUE: 固定值0
W_QUEUE: 写队列大小
FREE_BUFFER: BufferPool free大小
TOTAL_BUFFER: BufferPool 总大小
BU_PERCENT: BufferPool使用率百分比
BU_WARNS: 固定值0
FC_COUNT: 前端连接数量
BC_COUNT: 后端连接数量

```

### 2.1.3.13 show @@command

**show @@command;**

描述: processor对各个类型的数据包的分类统计信息

例:

```

mysql> show @@command;
+-----+-----+-----+-----+-----+-----+-----+-----+
| PROCESSOR | INIT_DB | QUERY | STMT_PREPARE | STMT_EXECUTE | STMT_CLOSE | PING | KILL | QUIT | OTHER |
+-----+-----+-----+-----+-----+-----+-----+-----+
Processor0	0	0	0	0	0	0	0	0	0
Processor1	0	0	0	0	0	0	0	0	0
Processor2	0	6	0	0	0	0	0	0	0
Processor3	0	4	0	0	0	0	0	0	0
+-----+-----+-----+-----+-----+-----+-----+-----+
4 rows in set (0.00 sec)

```

列描述:

```

PROCESSOR: processor名称
INIT_DB: COM_INIT_DB
QUERY: COM_QUERY
STMT_PREPARE: COM_STMT_PREEPARE
STMT_EXECUTE: COM_STMT_EXECUTE
STMT_CLOSE: COM_STMT_CLOSE
PING: COM_PING
KILL: COM_PROCESS_KILL
QUIT: COM_QUIT
OTHER: 其余

```

### 2.1.3.14 show @@connection where processor=? and front\_id=? and host=? and user=?

**show @@connection where processor=? and front\_id=? and host=? and user=?;**

描述: 查询前端连接信息, 可通过processor, front\_id, host和user进行过滤筛选, 条件可以任意组合搭配。

例:

```

mysql> show @@connection where processor='Processor2';
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| PROCESSOR | FRONT_ID | HOST | PORT | LOCAL_PORT | USER | SCHEMA | CHARACTER_SET_CLIENT | COLLATION_CONNECTION | CHARACTER_SET_RESULTS | NET_IN | NET_OUT | ALIVE_TIME(S) | RECV_BUFFER | SEND_QUEUE | TX_ISOLATION |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| Processor2 | 1 | 0:0:0:0:0:0:0:1 | 9066 | 54761 | man | NULL | utf8 | utf8_general_ci | utf8 | 237 | 23967 | 813 | 4096 | 0 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
1 row in set (0.11 sec)

```

结果列描述:

```

PROCESSOR: PROCESSOR名称
FRONT_ID: 前端连接ID
HOST: 客户端host
PORT: 本地端口(流量或者管理)
LOCAL_PORT: 客户端端口
USER: 用户
SCHEMA: 所在的schema
CHARACTER_SET_CLIENT: 字符集信息
COLLATION_CONNECTION: 字符集信息
CHARACTER_SET_RESULTS : 字符集信息
NET_IN: 接收流量
NET_OUT: 发送流量
ALIVE_TIME(S): 连接建立时长
RECV_BUFFER: 接收缓冲区大小(字节)
SEND_QUEUE: 发送缓冲区队列大小
TX_ISOLATION_LEVEL: 隔离级别
AUTOCOMMIT: 略
SYS_VARIABLES: 系统变量
USER_VARIABLES: 用户变量

```

### 2.1.3.15 show @@cache

**show @@cache;**

描述: 展示cache信息

例:

```

mysql> show @@cache;
+-----+-----+-----+-----+-----+-----+
| CACHE | MAX | CUR | ACCESS | HIT | PUT | LAST_ACCESS | LAST_PUT |
+-----+-----+-----+-----+-----+-----+
| ER_SQL2PARENTID | 1000 | 0 | 0 | 0 | 0 | | |
| SQLrouteCache | 10000 | 0 | 0 | 0 | 0 | | |
+-----+-----+-----+-----+-----+-----+
2 rows in set (0.09 sec)

```

列描述:

```

CACHE: cache名
MAX: 最大容量
CUR: 当前容量
ACCESS: 缓存查询次数
HIT: 命中次数
PUT: 加入缓存计数器
LAST_ACCESS: 上一次查询时间戳(格式为yyyy/mm/dd hh:mm:ss)
LAST_INPUT: 上一次加入缓存时间戳(格式为yyyy/mm/dd hh:mm:ss)

```

### 2.1.3.16 show @@backend where processor=? and backend\_id=? and mysqlid=? and host=? and port=?

**show @@backend where processor=? and backend\_id=? and mysqlid=? and host=? and port=?;**

描述: 查询活动的后端连接信息, 可与show @@session结合使用。该命令可通过processor, backend\_id, mysqlid, host和port进行过滤筛选, 条件可以任意组合搭配。

例:

```
mysql> show @@backend where host='10.18x.2x.63';
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| processor | BACKEND_ID | MYSQLID | HOST | PORT | LOCAL_TCP_PORT | NET_IN | NET_OUT | ACTIVE_TIME(S) | CLOSED | BORROWED | SEND_QUEUE | SCHEMA | CHARACTER_SET_CLIENT | COLLATION_CONNECTION | CHARACTER_SET_RESULT |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| Processor0 | 4 | 11094 | 10.18x.2x.63 | 3320 | 54734 | 600 | 171 | 166 | false | false | 0 | ares_test | utf8mb4 | utf8mb4_general_ci | utf8mb4
| Processor1 | 1 | 11077 | 10.18x.2x.63 | 3320 | 54701 | 865 | 236 | 166 | false | false | 0 | mycat_test | utf8mb4 | utf8mb4_general_ci | utf8mb4
| Processor2 | 2 | 11080 | 10.18x.2x.63 | 3320 | 54704 | 336 | 101 | 166 | false | false | 0 | nosharding | utf8mb4 | utf8mb4_general_ci | utf8mb4
| Processor3 | 3 | 11079 | 10.18x.2x.63 | 3320 | 54703 | 336 | 101 | 166 | false | false | 0 | NULL | utf8mb4 | utf8mb4_general_ci | utf8mb4
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
4 rows in set (0.00 sec)
```

列描述:

processor: processor名称  
 BACKEND\_ID: 后端连接ID  
 MYSQLID: mysql线程id(对应节点上的show processlist里的MYSQLID)  
 HOST: 主机名  
 PORT: 端口  
 LOCAL\_TCP\_PORT: tcp连接的本地端口  
 NET\_IN: 接收流量大小  
 NET\_OUT: 发送流量大小  
 ACTIVE\_TIME(S): 连接建立时间 (单位秒)  
 CLOSED: 是否被关闭  
 BORROWED: 是否正在使用  
 SEND\_QUEUE: 发送缓冲队列大小  
 SCHEMA: schema上下文  
 CHARACTER\_SET\_CLIENT: 字符集信息  
 COLLATION\_CONNECTION: 字符集信息  
 CHARACTER\_SET\_RESULTS: 字符集信息  
 TX\_ISOLATION\_LEVEL: 隔离级别 (新建未使用过的连接为-1, 表示未初始化)  
 AUTOCOMMIT: 是否自动提交  
 SYS\_VARIABLES: 系统变量  
 USER\_VARIABLES: 用户变量  
 XA\_STATUS: xa状态  
 DEAD\_TIME: 连接池被回收的时间, 连接在完成任务后也会关闭回收  
 USED\_FOR\_HEARTBEAT: 是否被用于心跳

### 2.1.3.17 show @@session

show @@session;

描述: 展示当前活动前端session的后端连接信息

例:

```
mysql> show @@session ;
+-----+-----+
| FRONT_ID | DN_COUNT | DN_LIST |
+-----+-----+
| 2 | 2 | MySQLConnection [backendId=59, lastTime=1508233042917, [...]] |
+-----+-----+
1 row in set (0.00 sec)
```

列描述:

FRONT\_ID: 前端连接ID  
 DN\_COUNT: 后端连接个数;  
 DN\_LIST: 后端连接的详情,

DN\_LIST例如:

```
MySQLConnection [backendId=59, lastTime=1508233042917, user=xxxx, schema=dble_test, old schema=dble_test, borrowed=true, fromSlaveDB=false, mysqlId=23201, character_set_client=utf8, character_set_results=utf8, collation_connection=latin1_swedish_ci, charsetName=latin1, txIsolation=REPEATABLE_READ, autocommit=false, attachment=, respHandler=, host=host ip, port=端口号, statusSync=同步上下文的类, writeQueue=写队列, modifiedSQLExecuted=是否是增删改]
```

含义:

backendId: 后端连接id  
 lastTime: 上次读写时间戳  
 user: 后端连接对应的用户  
 schema: 后端连接对应的schema  
 old\_schema: 上次连接被使用时的schema  
 borrowed: 被使用  
 fromSlaveDB: 是否是从数据库  
 mysqlId: 对应后端连接在数据库内的线程id (show processlist)  
 charsetName: 字符集  
 txIsolation: 隔离级别  
 autocommit: 自动提交  
 attachment: 路由结果集  
 respHandler: 收到回复时处理的handler类  
 host: host ip  
 port: 端口号  
 statusSync: 同步上下文的类  
 writeQueue: 写队列  
 modifiedSQLExecuted: 是否是增删改

### 2.1.3.18 show @@connection.sql

show @@connection.sql;

描述: 当前活动session的前端的SQL信息

例:

```
mysql> show @@connection.sql;
+-----+-----+-----+-----+-----+
| FRONT_ID | HOST | USER | SCHEMA | START_TIME | EXECUTE_TIME | SQL | STAGE |
+-----+-----+-----+-----+-----+
| 1 | 0:0:0:0:0:0:0:1 | man | NULL | 2017/10/17 17:00:58 | 139 | show @@connection.sql | Read SQL |
+-----+-----+-----+-----+-----+
1 row in set (0.13 sec)
```

列描述:

FRONT\_ID: 前端连接ID  
 HOST: 客户端host  
 USER: 用户  
 SCHEMA: 所在的schema  
 START\_TIME: 上次接收请求时间戳  
 EXECUTE\_TIME: 响应时间或者未完成SQL持续时间 (由于实现方式的原因, 可能出现正负20ms的误差)  
 SQL: 如果长度大于1024个字符, 将会被截断为1024  
 STAGE: 运行的当前阶段, 结束时会变成finished

### 2.1.3.19 show @@sql

show @@sql;

描述: 展示用户近期执行完的50条sql语句(多余的每5秒清理一次)

例:

```
mysql> show @@sql;
+-----+-----+-----+-----+
| ID | USER | START_TIME | EXECUTE_TIME | SQL |
+-----+-----+-----+-----+
| 1 | root | 2017/10/17 17:37:22 | 381 | select * from sharding_two_node LIMIT 100 |
+-----+-----+-----+-----+
1 row in set (0.02 sec)
```

列描述:

```
ID: 行号
USER: 用户
START_TIME: 上次接收请求时间戳
EXECUTE_TIME: 响应时间
SQL: 略
```

如果需要在查询后重置统计, 执行:

show @@sql true;

### 2.1.3.20 show @@sql.high

show @@sql.high;

描述: 展示各个用户的高频sql(容量1024, 超过会被定期清理, 清理周期5秒)

例:

```
mysql> show @@sql.high;
+-----+-----+-----+-----+-----+-----+
| ID   | USER  | FREQUENCY | AVG_TIME | MAX_TIME | MIN_TIME | EXECUTE_TIME | LAST_TIME      | SQL
+-----+-----+-----+-----+-----+-----+
| 1    | root   | 1          | 381     | 381     | 381     | 2017/10/17 17:37:23 | SELECT * FROM sharding_two_node LIMIT ? |
+-----+-----+-----+-----+-----+-----+
1 row in set (0.06 sec)
```

列描述:

```
ID: 行号
USER: 用户
FREQUENCY: sql曾被执行次数
AVG_TIME: 平均执行耗时
MAX_TIME: 最大执行耗时
MIN_TIME: 最小执行耗时
EXECUTE_TIME: 最近一次执行耗时
LAST_TIME: 最近一次执行时间戳
SQL: 略
```

如果需要在查询后重置统计, 执行:

show @@sql.high true;

### 2.1.3.21 show @@sql.slow

show @@sql.slow;

描述: 展示执行时间超过给定阈值(默认100毫秒, 可通过reload修改)的sql(默认10条, 可以通过设置系统参数sqlRecordCount修改, 多余的每5秒清理一次)

例:

```
mysql> show @@sql.slow;
+-----+-----+-----+
| USER | START_TIME      | EXECUTE_TIME | SQL
+-----+-----+-----+
| root | 2017/10/17 17:37:22 | 381          | select * from sharding_two_node LIMIT 100 |
+-----+-----+-----+
1 row in set (0.07 sec)
```

列描述:

```
USER: 用户
START_TIME:上次接收请求时间戳
EXECUTE_TIME:响应时间
SQL: 略
```

如果需要在查询后重置统计, 执行:

show @@sql.slow true;

### 2.1.3.22 show @@sql.resultset

show @@sql.resultset;

描述: 展示结果集大小超过某个阈值(默认512K, 可以通过maxResultSet配置)的sql, 结果集统计信息

例:

```
mysql> show @@sql.resultset;
+-----+-----+-----+-----+
| ID   | USER  | FREQUENCY | SQL           | RESULTSET_SIZE |
+-----+-----+-----+-----+
| 1    | root   | 1          | SELECT * FROM sharding_two_node | 1048576        |
+-----+-----+-----+-----+
1 row in set (0.05 sec)
```

列描述:

```
ID:行号
USER: 用户
FREQUENCY:sql曾被执行次数
SQL: 略
RESULTSET_SIZE:结果集的大小
```

### 2.1.3.23 show @@sql.sum

show @@sql.sum;

描述: 展示用户的sql执行情况, 是否带.user结果是一样的.带参数true, 表示查询结束后清空已经缓存的结果

例:

```
mysql> show @@sql.sum;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| ID   | USER  | R    | W    | R%   | MAX  | NET_IN | NET_OUT | TIME_COUNT | TTL_COUNT | LAST_TIME      |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1    | root   | 1    | 0    | 1.00 | 1     | 41    | 840    | [0, 0, 1, 0] | [0, 0, 1, 0] | 2017/10/17 17:37:23 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
1 row in set (0.26 sec)
```

列描述:

```
ID:行号
USER:用户
R:读的次数
W:写的次数, 恒为零(因为未实现统计)
R%:因为W为0, 此值恒为100%
MAX:最大并发数
NET_IN:网络流入量
NET_OUT:网络流出量
TIME_COUNT:query在四个时间区间的个数分布, 四个区间分别是前一天22-06 夜间, 06-13 上午, 13-18下午, 18-22 晚间
TTL_COUNT:query耗时在四个时间级别的个数分布, 四个区间分别是10毫秒内, 10 - 200毫秒内, 1秒内, 超过 1秒
LAST_TIME:上次SQL执行时间戳
```

如果需要在查询后重置统计, 执行:

```
show @@sql.sum true;
```

### 2.1.3.24 show @@sql.sum.user

等同于:

```
show @@sql.sum;
```

如果需要在查询后重置统计, 执行:

```
show @@sql.sum.user true;
```

### 2.1.3.25 show @@sql.sum.table

```
show @@sql.sum.table;
```

描述: 展示各个表的读写情况

例:

```
mysql> show @@sql.sum.table;
+----+-----+-----+-----+-----+-----+
| ID | TABLE          | R   | W   | Rk  | RELATABLE | RELACOUNT | LAST_TIME      |
+----+-----+-----+-----+-----+-----+
| 1  | sharding_two_node | 1  | 0  | 1.00 | NULL     | NULL       | 2017/10/17 17:37:23 |
+----+-----+-----+-----+-----+-----+
1 row in set (0.00 sec)
```

列描述:

ID:行号  
TABLE:表名(注:解析器实现很简单,可能有bug)  
R:读的次数  
W:写的次数,恒为零(因为未实现统计)  
Rk:因为W为0,此值恒为1000  
RELATABLE:关联表的名称(目前拆分表关联查询都使用查询计划树,此值为NULL)  
RELACOUNT:关联表的个数(目前拆分表关联查询都使用查询计划树,此值为NULL)  
LAST\_TIME:上次SQL执行时间戳

如果需要在查询后重置统计, 执行:

```
show @@sql.sum.table true;
```

### 2.1.3.26 show @@heartbeat

```
show @@heartbeat;
```

描述: 展示datasource的heartbeat信息

例:

```
mysql> show @@heartbeat;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| NAME  | HOST        | PORT | RS_CODE | RETRY | STATUS | TIMEOUT | EXECUTE_TIME | LAST_ACTIVE_TIME | STOP  | RS_MESSAGE |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| hostM1 | 10.18x.2x.63 | 3320 | OK    | 0 | idle   | 0 | 8,8,8 | NULL           | false | NULL      |
| hostM2 | 10.18x.2x.64 | 3320 | OK    | 0 | idle   | 0 | 9,9,9 | NULL           | false | NULL      |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
2 rows in set (0.07 sec)
```

列描述:

NAME: dataSource名称  
HOST:主机名/IP  
PORT:端口  
RS\_CODE:状态码, 有以下四种状态: INIT, OK, ERROR, TIMEOUT  
RETRY:重试错误次数  
STATUS:checking/idle  
TIMEOUT:心跳超时阈值(始终为0, bug?)  
EXECUTE\_TIME:最近3个时段的平均响应时间, 默认1, 10, 30分钟  
LAST\_ACTIVE\_TIME:上次收到心跳回复时间戳  
STOP:是否stop, 和stop命令相关  
RS\_MESSAGE:心跳失败信息, 当RS\_CODE为INIT, OK, TIMEOUT时, message为null, 只有当RS\_CODE为ERROR时, message才会显示最近一次心跳失败的信息

### 2.1.3.27 show @@heartbeat.detail where name=?

```
show @@heartbeat.detail where name=xxx;
```

其中, xxx为要查询的dataSource的名字。

描述: 展示指定dataSource的heartbeat的详细信息

前提条件:至少发生过一次心跳语句 (与dataNodeHeartbeatPeriod相关)

例:

```
mysql> show @@heartbeat.detail where name='hostM1';
+-----+-----+-----+-----+
| NAME  | HOST        | PORT | TIME          | EXECUTE_TIME |
+-----+-----+-----+-----+
| hostM1 | 10.18x.2x.63 | 3320 | 2017-10-17 17:31:58 | 7
| hostM1 | 10.18x.2x.63 | 3320 | 2017-10-17 17:32:59 | 9
+-----+-----+-----+-----+
2 row in set (0.00 sec)
```

列描述:

NAME: dataSource名称  
HOST:主机名/IP  
PORT:端口  
TIME:收到心跳时间戳  
EXECUTE\_TIME:心跳执行耗时(毫秒)

### 2.1.3.28 show @@sysparam

```
show @@sysparam;
```

描述: 展示sysconfig参数配置

结果: 略

### 2.1.3.29 show @@syslog limit=?

```
show @@syslog limit=N;
```

其中, N为整数。

描述: 按时间从新到旧展示db.log中的N条记录内容。

结果: 略

### 2.1.3.30 show @@white

```
show @@white;
```

描述: 展示配置的白名单信息

例:

```
mysql> show @@white;
+-----+-----+
| IP      | USER    |
+-----+-----+
| 0:0:0:0:0:0:0:1 | root,man,man2 |
| 127.0.0.1     | root      |
+-----+-----+
2 rows in set (0.02 sec)
```

列描述:  
略

### 2.1.3.31 show @@directmemory

**show @@directmemory;**

描述: 堆外内存使用总览  
结果集举例:

```
+-----+-----+-----+
| DIRECT_MEMORY_MAXED | DIRECT_MEMORY_POOL_SIZE | DIRECT_MEMORY_POOL_USED |
+-----+-----+-----+
| 3GB                | 1024MB             | 44KB                 |
+-----+-----+-----+
1 row in set (0.16 sec)
```

结果列描述:

DIRECT\_MEMORY\_MAXED: 通过-XX:MaxDirectMemorySize设置的值  
DIRECT\_MEMORY\_POOL\_SIZE: 内存池的大小, 等于bufferPoolPageSize和bufferPoolPageNumber的乘积  
DIRECT\_MEMORY\_POOL\_USED: 已经使用的内存池中的DirectMemory内存

### 2.1.3.32 show @@command.count

**show @@command.count;**

描述: 查询当前系统的查询数;

结果: 略

### 2.1.3.33 show @@connection.count

**show @@connection.count;**

描述: 查询当前的前端链接数;

结果: 略

### 2.1.3.34 show @@backend.statistics

**show @@backend.statistics;**

描述: 查询系统中进程的后端数据源信息;

例:

```
MySQL [(none)]> show @@backend.statistics;
+-----+-----+-----+
| HOST   | PORT   | ACTIVE  | TOTAL   |
+-----+-----+-----+
| 192.168.2.177 | 3307  | 0       | 10      |
| 192.168.2.177 | 3308  | 0       | 10      |
+-----+-----+-----+
2 rows in set (0.02 sec)
```

列描述:

HOST: 数据源的ip  
PORT: 数据源的端口  
ACTIVE: 数据源正在被使用的链接数  
TOTAL : 活的后端链接数。

### 2.1.3.35 show @@backend.old

**show @@backend.old;**

描述: reload @@config\_all之后待回收的活动的后端链接信息  
结果格式: 同show @@backend

### 2.1.3.36 show @@binlog.status

**show @@binlog.status;**

描述: 对server下所有节点拉一条一致性的binlog线。

例:

```
mysql> show @@binlog.status;
+-----+-----+-----+-----+-----+
| Url      | File      | Position | Binlog_Do_DB | Binlog_Ignore_DB | Executed_Gtid_Set |
+-----+-----+-----+-----+-----+
| 10.18.x.2x.63:3320 | mysql-bin.000024 | 14128    |          | 7ad71aab-de94-11e5-9488-3a935460da28:1-67646 |
| 10.18.x.2x.64:3320 | mysql-bin.000049 | 604440   |          | ba8f8b5c-debf-11e5-a87b-26b8a61f9012:1-91   |
+-----+-----+-----+-----+-----+
2 rows in set (0.11 sec)
```

列描述:

Url: 后端节点的连接Url值  
其余列: 等同于在对应节点上执行show master status的结果。

### 2.1.3.37 show @@help

**show @@help;**

描述: 展示帮助信息

结果: 略

### 2.1.3.38 show @@sql.large

**show @@sql.large;**

描述: 展示各个用户的结果集超过10000行的sql(容量为10,多的会被定时清理, 清理周期5秒)

例:

```
mysql> show @@sql.large;
+-----+-----+-----+
| USER | ROWS | START_TIME           | EXECUTE_TIME | SQL               |
+-----+-----+-----+
| root | 20000 | 2017/10/17 17:37:23 | 381          | SELECT * FROM sharding_two_node LIMIT ? |
+-----+-----+-----+
1 row in set (0.06 sec)
```

列描述:

```
USER: 用户
ROWS: 该查询的行数
START_TIME: 上次接收请求时间戳 EXECUTE_TIME: 响应时间
SQL: 略
```

如果需要在查询后重置统计, 执行:

show @@sql.large true;

**2.1.3.39 show @@sql.condition**

show @@sql.condition;

描述: 查询条件统计, 需要配合reload @@query\_cf 使用, 前者设置了table&column后, 运行此语句后展示sql查询条件统计信息。(最多100000条, 超出后不再统计)  
比如select from sharding\_two\_node where id =0; 和select from sharding\_two\_node where id =1;

例:

```
mysql> show @@sql.condition;
+-----+-----+-----+
| ID   | KEY      | VALUE | COUNT |
+-----+-----+-----+
2	sharding_two_node.id	0	1
3	sharding_two_node.id	1	2
2	sharding_two_node.id.valuekey	size	2
3	sharding_two_node.id.valuecount	total	3
+-----+-----+-----+
4 rows in set (0.05 sec)
```

列描述:

```
ID: 行号
KEY: schema.table 最后两行为schema.table.valuekey 和 schema.table.valuecount
VALUE: 对应key的value值
COUNT: 查询的次数
```

**2.1.3.40 show @@cost\_time;**

show @@cost\_time;

描述: 查询query耗时统计的结果, 需要在server.xml中开启useCostTimeStat选项之后才会有统计结果

例:

```
mysql> show @@cost_time;
+-----+-----+
| OVER_ALL(us) | FRONT_PREPARE          | BACKEND_EXECUTE          |
+-----+-----+
| 71496 | Id:9,Time:53135;Id:12,Time:54056 | Id:9,Time:16924;Id:12,Time:16006 |
| 15316 | Id:17,Time:2301;Id:11,Time:3196 | Id:17,Time:10691;Id:11,Time:11397 |
+-----+-----+
2 rows in set (0.05 sec)
```

列描述:

```
OVER_ALL: 总耗时
FRONT_PREPARE: 前端连接以及db表中的耗时
BACKEND_EXECUTE: 后端连接执行耗时
```

**2.1.3.41 show @@dataNodes where schema=? and table=?;**

show @@dataNodes

描述: 查询某具体表格的节点信息

例:

```
mysql> show @@dataNodes where schema=testdb and table=seqtest;
+-----+-----+-----+-----+-----+
| NAME | SEQUENCE | HOST     | PORT | PHYSICAL_SCHEMA | USER | PASSWORD |
+-----+-----+-----+-----+-----+
| dn1  | 0         | 10.186.24.113 | 3309 | db1           | root | 123456  |
| dn2  | 1         | 10.186.24.113 | 3309 | db2           | root | 123456  |
+-----+-----+-----+-----+-----+
2 rows in set (0.05 sec)
```

列描述:

```
NAME: 节点名称
SEQUENCE: 节点编号
HOST: 节点所在的IP
PORT: 节点对应的服务端口
PHYSICAL_SCHEMA: 节点所对应的物理库
USER: 节点连接的用户
PASSWORD: 节点连接的密码
```

**2.1.3.42 show @@algorithm where schema=? and table=?;**

show @@algorithm

描述: 查询某具体表格的分片算法信息, 由于不同算法会有不同的分片参数以及辅助文件及数据, 所以不同算法表格的输出分片事项都不相同

例:

```
mysql> show @@algorithm where schema=testdb and table=seqtest;
+-----+-----+
| KEY      | VALUE          |
+-----+-----+
TYPE	SHARDING TABLE
COLUMN	ID
CLASS	com.actiontech.dble.route.function.PartitionByLong
partitionCount	2
partitionLength	1
+-----+-----+
5 rows in set (0.05 sec)
```

行描述:

```
KEY: 分片事项
VALUE: 详细信息
```

**2.1.3.43 show @@thread\_used;**

show @@thread\_used;

描述: 查看各个主要业务处理线程的使用状况

例:

```
mysql> show @@thread_used;
+-----+-----+-----+-----+
| THREAD_NAME | LAST_QUARTER_MIN | LAST_MINUTE | LAST_FIVE_MINUTE |
+-----+-----+-----+-----+
BusinessExecutor3	0%	0%	0%
$_.NIO_REACTOR_BACKEND-2	0%	0%	0%
BusinessExecutor1	0%	0%	0%
$_.NIO_REACTOR_BACKEND-3	0%	0%	0%
$_.NIO_REACTOR_BACKEND-0	0%	0%	0%
$_.NIO_REACTOR_FRONT-0	0%	0%	0%
$_.NIO_REACTOR_BACKEND-1	0%	0%	0%
+-----+-----+-----+-----+
7 rows in set (0.00 sec)
```

行描述:

THREAD\_NAME:线程名称  
LAST\_QUARTER\_MIN: 最近15秒使用率  
LAST\_MINUTE: 最近一分钟使用率  
LAST\_FIVE\_MINUTE: 最近五分钟使用率

**2.1.3.44 show @@ddl;****show @@ddl;**

描述: 查看正在执行, 没有在dble内部释放锁的DDL

例:

```
mysql> show @@ddl;
+-----+-----+
| Schema | Table | Sql |
+-----+-----+
| testdb | sharding_two_node | alter table sharding_two_node add column id2 int |
| mytest | sharding_four_node | drop table sharding_four_node |
+-----+-----+
2 rows in set (0.00 sec)
```

行描述:

Schema:Schema名称  
Table: Table名称  
Sql: ddl sql语句

**2.1.3.45 show @@processlist;****show @@processlist;**

描述: 查看前端连接和后端连接对应关系, 若前端连接没有对应的后端连接, 显示NULL。

例:

```
mysql> show @@processlist;
+-----+-----+-----+-----+-----+-----+-----+
| Front_Id | Datanode | Mysqld | User | Front_Host | db | Command | Time | State | Info |
+-----+-----+-----+-----+-----+-----+-----+
1	dn2	2303	root	127.0.0.1:33222	db2	Sleep	17		NULL
2	NULL	man1	127.0.0.1:34882	NULL	NULL	0		NULL	
3	dn3	2259	root	127.0.0.1:33226	db1	Sleep	4		NULL
3	dn2	2308	root	127.0.0.1:33226	db2	Sleep	4		NULL
3	dn1	2304	root	127.0.0.1:33226	db1	Sleep	4		NULL
+-----+-----+-----+-----+-----+-----+-----+
5 rows in set (0.05 sec)
```

行描述:

Front\_Id: 前端连接ID  
Datanode: 前端连接下发操作的 datanode  
Mysqld: 后端连接对应的 mysql 线程ID  
User: 用户名  
Front\_Host: 客户端主机名  
db: 后端连接默认数据库, 来自于 mysql 'show processlist' 字段 db  
Command: mysql线程正在执行的指令类型, 来自于 mysql 'show processlist' 字段 Command  
Time: mysql线程处于当前state的时间, 来自于 mysql 'show processlist' 字段 Time  
State: mysql线程执行状态, 来自于 mysql 'show processlist' 字段 State  
Info: mysql线程执行语句, 来自于 mysql 'show processlist' 字段 Info

**2.1.3.46 show @@session.xa;****show @@session.xa;**

描述: 查看后台重试的xa事务信息。

例:

```
mysql> show @@session.xa;
+-----+-----+-----+
| FRONT_ID | XA_ID | XA_STATE | DATANODES |
+-----+-----+-----+
| 1 | 'Dble_Server.1.1' | TX_COMMIT_FAILED_STATE | dn1,dn3 |
+-----+-----+-----+
1 rows in set (0.00 sec)
```

行描述:

FRONT\_ID: 前端连接ID  
XA\_ID: xa事务id  
XA\_STATE: xa事务状态  
DATANODES: xa提交失败的datanode名称

**2.1.3.47 show @@reload\_status****show @@reload\_status**

描述: 查看dble中最近的reload信息

举例

```
+-----+-----+-----+-----+-----+-----+-----+
| INDEX | CLUSTER | RELOAD_TYPE | RELOAD_STATUS | LAST_RELOAD_START | LAST_RELOAD_END | TRIGGER_TYPE | END_TYPE |
+-----+-----+-----+-----+-----+-----+-----+
| 0 | false | RELOAD_ALL | NOT_RELOADING | 2019/08/19 14:28:04 | 2019/08/19 14:28:05 | LOCAL_COMMAND | RELOAD_END |
+-----+-----+-----+-----+-----+-----+-----+
```

行描述:

INDEX: reload对应的编号, 能与日志中的[RL]日志编号相对应  
CLUSTER: 当前dble使用的集群方式  
RELOAD\_TYPE: 最近的reload的类型 reload\_metadata/config/config\_all/rollback  
RELOAD\_STATUS: 最近一次reload的执行状态not\_reloading/self\_reload/meta\_reload/waiting\_others  
LAST\_RELOAD\_START: 起始时间  
LAST\_RELOAD\_END: 结束时间  
TRIGGER\_TYPE: 触发类型reload\_command/cluster\_notify  
END\_TYPE: 结束原因

此命令被用于配合命令 [release @@reload\\_metadata](#)

### 2.1.3.48 show @@user

**show @@user**

描述: 查看db1e 所有用户

举例

```
mysql> show @@user;
+-----+-----+-----+
| Username | Manager | Readonly | Max_con |
+-----+-----+-----+
man1	Y	N	no limit
root	N	N	no limit
user	N	N	no limit
+-----+-----+-----+
3 rows in set (0.03 sec)
```

行描述:

Username: 用户名  
Manager: 是否是管理用户  
 Readonly: 是否是只读用户  
Max\_con: 最大连接数

### 2.1.3.49 show @@user.privilege

**show @@user.privilege**

描述: 查看db1e 用户的权限信息, 不包含管理用户

举例

```
mysql> show @@user.privilege;
+-----+-----+-----+-----+-----+-----+
| Username | Schema | Table | INSERT | UPDATE | SELECT | DELETE |
+-----+-----+-----+-----+-----+-----+
root	testdb1	*	Y	Y	Y	Y
root	testdb	*	Y	Y	Y	Y
user	testdb	*	N	Y	Y	N
+-----+-----+-----+-----+-----+-----+
3 rows in set (0.01 sec)
```

行描述:

Username: 用户名  
Schema: 用户授权逻辑库  
Table: 用户显式指定dml权限的表名, 未指定的其他表使用\*表示  
INSERT: 插入权限位  
UPDATE: 更新权限位  
SELECT: 查询权限位  
DELETE: 删除权限位

### 2.1.3.50 show @@data\_distribution where table ='schema.table'

**show @@data\_distribution where table ='schema.table'**

描述: 查看某个表在各个节点上的数据分布情况

举例

```
+-----+-----+
| DATANODE | COUNT |
+-----+-----+
dn1	100
dn2	101
dn3	98
dn4	104
+-----+-----+
4 rows in set (0.09 sec)
```

行描述:

DATANODE: 数据节点名字  
COUNT: 数据量

### 2.1.3.51 show @@Questions

**show @@Questions**

描述: 查看自启动之后SQL服务端口执行的QUERY和Transaction数量

举例:

```
mysql> show @@Questions;
+-----+-----+
| Questions | Transactions |
+-----+-----+
| 0          | 0           |
+-----+-----+
```

行描述:

Questions: 收到的查询的数量  
Transactions: 执行事务的数量, 非事务查询算单语句事务

#### 2.1.4 switch命令

本小节部分自2.20.04.0 版本起废弃

## 2.1.5 kill命令

### 2.1.5.1 kill @@connection;

```
kill @@connection id1,id2,...;
```

其中，idx为前端连接id值，可以通过show @@connection 获取。

描述：关闭存在的前端链接，对不存在的前端链接，不会报错

结果：返回OK，关闭的前端链接数。

### 2.1.5.2 kill @@xa\_session;

```
kill @@xa_session id1,id2,...;
```

其中，idx为 session id值，可以通过show @@session.xa 获取。

描述：取消指定session后台重试xa事务，对不存在的session，不会报错

结果：返回OK，取消的session数量。

### 2.1.5.3 KILL @@DDL\_LOCK where schema=? and table=?

```
KILL @@DDL_LOCK where schema=? and table=?;
```

描述：释放指定schema下table的ddl锁，详细描述可参考 [2.22 KILL @@DDL\\_LOCK](#)

结果：返回OK。

## 2.1.6 stop命令

### 2.1.6.1 stop @@heartbeat

`stop @@heartbeat keys:datashost;`

其中: `keys:datashost`名字列表, 可以是多个, 用逗号隔开的key; `key`:可以直接是`datashost`, 也可以是`datashost$0-n` (如`datashost$0-2`实际会输出`datashost[0],datashost[1],datashost[2]`, 疑似BUG) ;

`value`:应当是个整数, 单位毫秒

描述: 设置`datashost`名为`key`的host停止`heartbeat` n秒

结果: 返回OK

注意: 未作异常处理, 慎用

## 2.1.7 reload 命令

### 2.1.7.1 reload @@config

`reload @@config;`

2.19.09.0(不含)版本以前:

描述: 重新加载数据资源配置并自检, 涉及rule.xml, schema.xml, server.xml 内容, 但不包含dataHosts 和dataNodes, 自检失败将返回ERROR。  
结果: OK 或者ERROR

相关影响: 当执行此命令时, 当有以下情况发生时, 涉及到的表的meta信息会被重载, 否则保持原有表meta信息。

- 有新增的表
- 有删除的表
- 表的datanode或者type属性发生变更
- 有新增的schema
- 有删除的schema
- schema 的默认datanode发生变更

2.19.09.0(含)版本之后,功能完全等同于`reload @@config_all`

### 2.1.7.2 reload @@config\_all

`reload @@config_all [-s] [-f] [-r];`

描述: 重新加载所有配置, 涉及rule.xml, schema.xml, server.xml 内容, 会重新加载配置dataHosts 和dataNodes,例外的是此命令对server.xml中 `<system>` 段中的配置不生效。

-s 跳过测试后端链接, 默认不加此参数会对配置中的后端链接进行测试, 测试失败将返回ERROR;  
-f 关闭所有变更的数据host (加-r参数所有datahost会被视为变更) 相关的处于事务中的前端链接,如果无此参数默认仅将相应后端链接放入旧链接池。  
-r 不做智能判断,将所有后端连接池全部重新加载一遍。不加此参数时, 将对新配置进行智能判断, 只会对增删改的连接池做变更, 不影响未作变更的连接池  
更多细节 参考 [2.19 智能计算reload\\_all](#)

结果: OK 或者ERROR

相关影响: 当执行此命令时, 当有以下情况发生时, 涉及到的表的meta信息会被重载, 否则保持原有表meta信息。

- 有新增的表
- 有删除的表
- 表的datanode或者type属性发生变更
- 表的datanode对应的物理节点或者对应的datasource发生变更
- 有新增的schema
- 有删除的schema
- schema 的默认datanode属性发生变更
- schema的datanode对应的物理节点或者对应的datasource发生变更。

另外, 如果包含-r参数则不做上述判断, 全部重新加载meta数据。

如果不包含-r但是包含-s参数,则对metadata是否需要重新加载的计算时, 忽略所有datasource的变更

注意,不能在配置变更中体现的某些变化是无法重新加载metadata的, 举例[#1002](#)

一个带有默认datanode的schema尝试通过删除配置将拆分表或者global表变成非拆分表是不符合规范的。应当避免这种操作。

注意: 如果使用默认的切换方式(即单实例部署并且system的outerHA属性为false), 需要做配置的重载时,需要人工保证流程是标准的,否则可能导致切换功能故障,具体请参看相关章节[2.12 故障切换](#)。

### 2.1.7.3 reload @@metadata

`reload @@metadata;`

描述: 重新加载所有元数据信息。

结果: 返回OK。

支持过滤表达式

`reload @@metadata where schema=? [ and table=? ]`

描述: 重新加载指定schema中所有表或指定表的元数据信息。

结果: 返回OK。

`reload @@metadata where table in ('schema1.table1','schema2.table2','schema1.table3',...)`

描述: 重新加载schema1中table1,table3和schema2中table2的元数据信息。

结果: 返回OK。

### 2.1.7.4 reload @@sqlslow=N;

`reload @@sqlslow=N;`

描述: 设定用户分析统计的slow sql时间阈值到N毫秒; 结果: OK

### 2.1.7.5 reload @@user\_stat

`reload @@user_stat;`

描述: 重置用户状态统计结果。

影响的命令:

```
show @@sql;
show @@sql.sum;
show @@sql.slow;
show @@sql.high;
show @@sql.large;
show @@sql.resultset;
```

结果: OK

### 2.1.7.6 reload @@query\_cf

`reload @@query_cf=table&column];`

其中, table为要统计的目标表的表名, column为目标表中目标列的列名。

描述: 重设show @@sql.conditiont要统计的条件。

结果: OK

如果要清除查询条件统计表列的设置执行命令:

`reload @@query_cf;`

或者

`reload @@query_cf=NULL;`

## 2.1.8 rollback命令

### 2.1.8.1 rollback @@config

`rollback @@config;`

描述: 针对reload @@config 或者reload @@config\_all 操作, 回滚为原有配置信息, 注意可能与配置文件不同, 是内存级别的修改。

结果: OK或者ERROR

### 2.1.9 offline命令

**offline;**

描述: 设置dble处于离线状态, 之后外部对dble的ping命令, select user, 心跳都会返回错误。

问题: 其他查询未处理, 合理性需要考量

结果: OK

### 2.1.10 online命令

**online;**

描述: 设置dbfile处于在线状态, 与offline相反  
结果: OK

## 2.1.11 file命令

### 2.1.11.1 file @@list

描述: 列出conf目录下文件

例:

```
MySQL [(none)]> file @@list;
+-----+
| DATA
+-----+
| 1 : wrapper.conf time:2019-10-17 11:13 |
| 2 : myid.properties time:2019-10-17 11:11 |
+-----+
2 rows in set (0.00 sec)
```

### 2.1.11.2 file @@show filename

其中, `filename`为目标文件的文件名

描述: 展示配置文件`filename`的内容 结果: 单列, 其中每行为`conf/filename`的每行内容

结果: 略

### 2.1.11.3 file @@upload filename content

描述: 加载`content`写到文件`fileName`中, `fileName`前后必须以空格隔开, `xml`文件会做校验。

结果: 单行单列, 返回加载是否成功

### 2.1.12 log命令

**log @@[file=logfile limit=numberOfRow key=keyword regex=regex]**

描述：展示logs下文件logfile的信息。格式错误导致参数列表为空，展示所有文件的信息。

- logfile: 不指定默认为dble.log
- rowLimit: 以逗号为分隔符，不指定默认为0:10000
- keyword: 筛选日志中包含key,如果需要筛选的key中包含空格，允许使用'key word'形式
- regexStr: 筛选符合正则表达式的结果，注意regexStr中不允许包含空格，如有需要使用\s进行替代

结果：单列，按行展示logs下文件logfile的信息。

### 2.1.13 配置检查命令dryrun

背景：当我们修改配置之后并且reload之前，可以通过dryrun来检查配置的正确性。

举例，当这样一个schema.xml准备reload时，先观察dryrun结果： schema.xml:

```
<?xml version="1.0"?>
<!--
  Copyright (C) 2016-2018 ActionTech.
  License: http://www.gnu.org/licenses/gpl.html GPL version 2 or higher.
-->

<!DOCTYPE dble:schema SYSTEM "schema.dtd">
<dbleschema xmlns:dbles="http://dbles.cloud/">
  <schema name="testDb" sqlMaxLimit="100" >
    <table name="sharding_two_node" dataNode="dn1,dn2" rule="two_node_hash"/>
    <table name="sharding_two_node2" dataNode="dn1,dn2" rule="two_node_hash"/>
    <table name="sharding_two_node3" dataNode="dn1,dn2" rule="two_node_hash"/>
    <table name="sharding_four_node" dataNode="dn1,dn2,dn3,dn4" rule="four_node_hash" />
    <table name="test_table" type="global" dataNode="dn$1-2"/>
    <table name="a_test" cacheKey="id" dataNode="dn1,dn2,dn3,dn4" rule="four_node_hash" />
    <table name="a_order" cacheKey="id" dataNode="dn1,dn2,dn3,dn4" rule="four_node_hash" />
    <table name="test_shard" dataNode="dn1,dn2,dn3,dn4" rule="four_node_hash"/>
    <table name="test_global" dataNode="dn1,dn2,dn3,dn4" type="global"/>

    <table name="sbtest1" cacheKey="id" dataNode="dn1,dn2,dn3,dn4" rule="four_node_hash" />

  </schema>

  <schema name="nosharding_test" sqlMaxLimit="100" dataNode="dn5">
    </schema>
    <dataNode name="dn1" dataHost="dh1" database="ares_test" />
    <dataNode name="dn2" dataHost="dh2" database="dbles_test" />
    <dataNode name="dn3" dataHost="dh1" database="mycat_test" />
    <dataNode name="dn4" dataHost="dh2" database="mycat_test" />
    <dataNode name="dn5" dataHost="dh1" database="nosharding" />
    <dataNode name="dn8" dataHost="dh1" database="xxxxxx" />
    <dataNode name="dn9" dataHost="dh1" database="xxxxxx2" />
    <dataHost name="dh1" maxCon="10" minCon="1" balance="2" slaveThreshold="100" tempReadHostAvailable="0">
      <heartbeat>show slave status</heartbeat>
      <writeHost host="hostM1" url="10.186.xx.xxx:3306" user="action"
                 password="action" />
    </writeHost>
  </dataHost>

  <dataHost name="dh2" maxCon="10" minCon="1" balance="1" slaveThreshold="100" >
    <heartbeat>show slave status</heartbeat>
    <writeHost host="hostM2" url="10.186.xx.xxx:3306" user="qrep"
               password="qrep" />
  </dataHost>
</dbles:schema>
```

dryrun结果如下：

```
mysql> dryrun;
+-----+-----+
| TYPE | LEVEL   | DETAIL           |
+-----+-----+
| Xml  | WARNING | dataNode dn9 is useless |
| Xml  | WARNING | dataNode dn8 is useless |
+-----+-----+
2 rows in set (0.58 sec)
```

列名含义：

TYPE: 错误类型，比如XML表示xml配置错误，BACKEND表示后端连接错误

LEVEL: 错误级别：分为WARNING 和ERROR表,一般来说WARNING错误不影响启动和使用，但需要注意。

DETAIL :错误详情

## 2.1.14 dataNode级别的流量暂停和恢复功能

背景：当我们做部分datanode的拆分或者是数据重分片时，有时希望不影响到其他无关的业务所涉及到的datahost和datanode，所以希望只停下部分datanode的流量。

### 2.1.14.0 功能描述

pause功能会在dble不停止的状态下停止对于指定后端节点的流量，在暂停期间的所有涉及到节点的路由结果会被挂起（为防止挂起的连接过多，queue和wait\_limit参数会控制挂起连接的数量和时间），直到恢复命令resume被执行之后，之前被挂起的查询才会继续进行。

pause的执行不一定成功，pause命令将会在指定的timeout时间内等待需要暂停的流量上的所有事务或者正在执行的sql结束，当目标节点上来自dble的sql或者事务一直无法在指定时间内结束的时候，本次暂停会返回失败。

pause（暂停）是一种dble server级别的全局状态，暂停节点不可追加或者逐步恢复，只能一起暂停指定的一个或几个dataNode并一起恢复流量，并且暂停期间执行reload命令不会恢复对应的流量状态。但是暂停操作不会被记录在文件之中，若发生dble重启，则暂停状态会被重置。

这里推荐谨慎的使用暂停并选择影响的范围，推荐逐个变动暂停，reload，恢复，再进行下一个变动。

#### 2.1.14.1 暂停流量：

```
pause @@DataNode = 'dn1,dn2' and timeout = 10 ,queue = 10,wait_limit = 10;
```

参数描述：

timeout:这个参数是等待涉及到的事务完成的时间，如果在到达timeout之后，仍然有事务未完成，本次pause失败，单位秒。

queue:这个参数表示暂停期间的阻塞前端连接的数量，超过此数量时，前端连接建立将会发生错误。

wait\_limit:是针对被阻塞的每个单个的前端的时间限制，如果被阻塞了超过wait\_limit的时间，将会返回错误，单位秒。

#### 2.1.14.2 恢复流量：

```
RESUME; 返回正常的OK 或者错误"No dataNode paused"
```

#### 2.1.14.3 查看当前暂停状态：

```
show @@pause;
```

```
mysql> show @@pause;
+-----+
| PAUSE_DATANODE |
+-----+
| dn1      |
| dn2      |
+-----+
2 rows in set (0.15 sec)
```

另外，商业版本支持集群操作，开源版本后续也将支持

## 2.1.15 慢日志相关命令

慢日志相关的命令:

### 2.1.15.1 查询慢查询日志的开启状态

```
mysql> show @@slow_query_log;
+-----+
| @@slow_query_log |
+-----+
| 0 |
+-----+
1 row in set (0.00 sec)
```

### 2.1.15.2 开启慢查询日志

```
mysql> enable @@slow_query_log;
Query OK, 1 row affected (0.09 sec)
enable slow_query_log success
```

### 2.1.15.3 关闭慢查询日志

```
mysql> disable @@slow_query_log;
Query OK, 1 row affected (0.03 sec)
disable slow_query_log success
```

### 2.1.15.4 查看慢查询日志统计阈值

```
mysql> show @@slow_query.time;
+-----+
| @@slow_query.time |
+-----+
| 100 |
+-----+
1 row in set (0.00 sec)
```

### 2.1.15.5 修改慢查询日志统计阈值

```
mysql> reload @@slow_query.time=200;
Query OK, 1 row affected (0.10 sec)
reload @@slow_query.time success

mysql> show @@slow_query.time;
+-----+
| @@slow_query.time |
+-----+
| 200 |
+-----+
1 row in set (0.00 sec)
```

### 2.1.15.6 查看慢查询日志刷盘周期

```
mysql> show @@slow_query.flushperiod;
+-----+
| @@slow_query.flushperiod |
+-----+
| 1 |
+-----+
1 row in set (0.00 sec)
```

### 2.1.15.7 修改慢查询日志刷盘周期

```
mysql> reload @@slow_query.flushperiod=2;
Query OK, 1 row affected (0.05 sec)
reload @@slow_query.flushPeriod success

mysql> show @@slow_query.flushperiod;
+-----+
| @@slow_query.flushperiod |
+-----+
| 2 |
+-----+
1 row in set (0.00 sec)
```

### 2.1.15.8 查看慢查询日志刷盘条数阈值

```
mysql> show @@slow_query.flushsize;
+-----+
| @@slow_query.flushsize |
+-----+
| 1000 |
+-----+
1 row in set (0.01 sec)
```

### 2.1.15.9 修改慢查询日志刷盘条数阈值

```
mysql> reload @@slow_query.flushsize=1100;
Query OK, 1 row affected (0.03 sec)
reload @@slow_query.flushsize success

mysql> show @@slow_query.flushsize;
+-----+
| @@slow_query.flushsize |
+-----+
| 1100 |
+-----+
1 row in set (0.00 sec)
```

### 2.1.15.10 查看某个连接的当前执行状态

show @@connection.sql.status where FRONT\_ID=?; 此功能需要开启慢日志才有效，当对应的连接当前query已经执行完毕时，执行此命令的结果与 trace 功能相同。如果query正在执行，本结果将试图展示query执行到哪一个步骤了。例如，广播查询

```
mysql> show @@connection.sql.status where FRONT_ID= 1;
+-----+-----+-----+-----+-----+-----+
| OPERATION | START(ms) | END(ms) | DURATION(ms) | DATA_NODE | SQL/REF |
+-----+-----+-----+-----+-----+-----+
| Read_SQL | 0.0 | 0.082598 | 0.082598 | - | - |
| Parse_SQL | 0.082598 | 0.676424 | 0.593826 | - | - |
| Route_Calculation | 0.676424 | 0.895382 | 0.218958 | - | - |
| Prepare_to_Push/Optimize | 0.895382 | 6743.838628 | 6742.943246 | - | - |
| Execute_SQL | 6743.838628 | 6753.488422 | 9.649794 | dn1 | select * from sharding_4_t1 |
| Execute_SQL | 6743.838628 | 6751.472835 | 7.634207 | dn3 | select * from sharding_4_t1 |
| Execute_SQL | 6743.838628 | 6750.981646 | 7.143018 | dn4 | select * from sharding_4_t1 |
| Execute_SQL | 6743.838628 | 6753.31394 | 9.475312 | dn2 | select * from sharding_4_t1 |
| Fetch_result | 6753.488422 | 6754.383316 | 0.894894 | dn1 | select * from sharding_4_t1 |
| Fetch_result | 6751.472835 | 6751.656604 | 0.183769 | dn3 | select * from sharding_4_t1 |
| Fetch_result | 6750.981646 | 6751.188385 | 0.206739 | dn4 | select * from sharding_4_t1 |
| Fetch_result | 6753.31394 | 6754.286055 | 0.972115 | dn2 | select * from sharding_4_t1 |
| Write_to_Client | 6750.981646 | unfinished | unknown | - | - |
+-----+-----+-----+-----+-----+-----+
13 rows in set (0.04 sec)
```

再比如join

```
mysql> show @@connection.sql.status where FRONT_ID= 1;
+-----+-----+-----+-----+-----+-----+
| OPERATION | START(ms) | END(ms) | DURATION(ms) | DATA_NODE | SQL/REF |
+-----+-----+-----+-----+-----+-----+
| Read_SQL | 0.0 | 0.039588 | 0.039588 | - | - |
| Parse_SQL | 0.039588 | 0.756578 | 0.71699 | - | - |
| Route_Calculation | 0.756578 | 1.5547 | 0.798122 | - | - |
| Prepare_to_Push/Optimize | 1.5547 | 3.428551 | 1.873851 | - | - |
| Execute_SQL | 3.428551 | 2362.10579 | 2358.677239 | dn1_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| Fetch_result | 2362.10579 | unfinished | unknown | dn1_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| Execute_SQL | 3.428551 | 2362.122407 | 2358.693856 | dn2_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| Fetch_result | 2362.122407 | unfinished | unknown | dn2_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| Execute_SQL | 3.428551 | 2362.307153 | 2358.878602 | dn3_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| Fetch_result | 2362.307153 | unfinished | unknown | dn3_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| Execute_SQL | 3.428551 | 2364.523615 | 2361.095064 | dn4_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| Fetch_result | 2364.523615 | unfinished | unknown | dn4_0 | select `a`.`age`, `a`.`id`, `a`.`name` from `sharding_4_t1` `a` ORDER BY `a`.`id` ASC |
| MERGE_AND_ORDER | 2362.639012 | unfinished | unknown | merge_and_order_1 | dn1_0; dn2_0; dn3_0; dn4_0 |
| SHUFFLE_FIELD | 4178.383366 | unfinished | unknown | shuffle_field_1 | merge_and_order_1 |
| Execute_SQL | 3.428551 | 2365.71371 | 2362.285159 | dn1_1 | select `b`.`id` from `sharding_2_t1` `b` ORDER BY `b`.`id` ASC |
| Fetch_result | 2365.71371 | unfinished | unknown | dn1_1 | select `b`.`id` from `sharding_2_t1` `b` ORDER BY `b`.`id` ASC |
| Execute_SQL | 3.428551 | 2365.952707 | 2362.524156 | dn2_1 | select `b`.`id` from `sharding_2_t1` `b` ORDER BY `b`.`id` ASC |
| Fetch_result | 2365.952707 | unfinished | unknown | dn2_1 | select `b`.`id` from `sharding_2_t1` `b` ORDER BY `b`.`id` ASC |
| MERGE_AND_ORDER | 2366.164823 | unfinished | unknown | merge_and_order_2 | dn1_1; dn2_1 |
| SHUFFLE_FIELD | not started | unfinished | unknown | - | - |
| JOIN | not started | unfinished | unknown | - | - |
| SHUFFLE_FIELD | not started | unfinished | unknown | - | - |
| Write_to_Client | not started | unfinished | unknown | - | - |
+-----+-----+-----+-----+-----+-----+
23 rows in set (0.04 sec)
```

## 2.1.16 创建/删除物理数据库命令

### 2.1.16.1 创建物理数据库

用于dble启动后发现有些datanode对应的物理库还未建立，可以使用后端命令一次性建立。

命令格式：

```
create database @@dataNode ='dn.....'
```

dataNode值支持 dn\$1-4 这种形式。

当所包含的datanode至少有一个不在配置文件当中时，返回错误： DataNode \$Name does not exists.

否则会对所有 datanode 执行 create database if not exists \$databaseName，执行完成之后返回OK。

### 2.1.16.2 删除物理数据库

用于删除某些datanode对应的物理库，可以使用后端命令一次性删除。

命令格式：

```
drop database @@dataNode ='dn.....'
```

dataNode值支持 dn\$1-4 这种形式。

当所包含的datanode至少有一个不在配置文件当中时，返回错误： DataNode \$Name does not exists.

否则会对所有 datanode 执行 drop database if exists \$databaseName，执行完成之后返回OK。若在执行过程中发生错误，会将show @@datanode里面的SCHEMA\_EXISTS置为false，需要用户人工确认下是否物理库已删除成功。

## check系列命令

### 2.1.17.0 check @@metadata 命令

用于检查meta信息是否存在以及加载的时间。

命令格式:

- 第一种形式 `check @@metadata`
  - 返回上一次 `reload @@metadata` 的datatime (或者上一次 `reload @@config_all` 的datatime 或者启动时加载meta的datatime)
- 第二种形式 `check full @@metadata`, 并且支持以下过滤条件:
  - `where schema=? and table=?`
  - `where schema=?`
  - `where reload_time='yyyy-MM-dd HH:mm:ss', where reload_time>='yyyy-MM-dd HH:mm:ss', where reload_time<='yyyy-MM-dd HH:mm:ss'`
  - `where reload_time is null`
  - `where consistent_in_data_nodes=0`
  - `where consistent_in_data_nodes = 1`
  - `where consistent_in_memory=0`
  - `where consistent_in_memory = 1`
  - `If no where, retrun all results.`
- `check full @@metadata` 结果集如下:

schema	table	reload_time	table_structure	consistent_in_data_nodes	consistent_in_memory
schema	table	2018-09-18 11:01:04	CREATE TABLE table`(...	1	1

column `table_structure` 和 `show create table` 命令结果的形式一样

column `consistent_in_data_nodes` 表示不同分片之间的一致性, 0为不一致, 1为一致

column `consistent_in_memory` 表示内存中meta与实际后端结点的一致性, 0为不一致, 1为一致

当`table_structure`列为null时, `consistent_in_data_nodes`列和`consistent_in_memory`列没有意义。

当`consistent_in_data_nodes`为0时, `consistent_in_memory`没有意义。

### 2.1.17.1 check @@global schema = " [and table = "]

用于进行手动全局表检查的命令, 当即触发一次特定范围的全局表检查, 并且将检查结果作为返回值进行展示。

结果如下所示:

```
mysql> check @@global schema = 'testdb';
+-----+-----+-----+
| SCHEMA | TABLE | DISTINCT_CONSISTENCY_NUMBER | ERROR_NODE_NUMBER |
+-----+-----+-----+
| testdb | tb_global1 | 0 | 0 |
+-----+-----+-----+
```

**SCHEMA:** 所检查的SCHEMA名字

**TABLE:** 所检查的TABLE名字

**DISTINCT\_CONSISTENCY\_NUMBER:** 返回有几个不同的检查结果的版本

**ERROR\_NODE\_NUMBER:** 在检查过程中有几个节点执行SQL报错

具体最终结果是否符合用户的预期, 还是要根据表格的检查SQL和返回信息来进行判断, 但一般情况下认为当`DISTINCT_CONSISTENCY_NUMBER`的值大于1的情况下, 表格中的内容一定不一致

## 2.1.18 release 命令

### 2.1.18.1 release @@reload\_metadata

描述：打断正在进行中的reload\_metadata步骤。当reload/reload config\_all/reload metadata/rollback config进行到最后一个步骤时，都会进行reload metadata  
就经验而言这个步骤存在hang死的风险，所以当出现reload metadata步骤迟迟未返回的情况，可以使用这个命令进行终止，使得dble能够继续提供服务 结果：OK 或者ERROR

注意：

- 此命令会导致config可能和table meta不匹配的情况，若执行此命令，请务必在后面追加命令reload @@metadata更新最新的meta信息
- 当一个reload过程被打断时，执行reload的进程会返回编号为5999的SQL错误，请慎重对待此错误CODE，这个错误代表着“新的配置被应用，meta信息在更新的过程中终止”
- 在执行运维操作的过程中可能会需要查看dble当前的reload状态，可通过具体命令show @@reload\_status进行查看

## 2.1.19 split 命令

### 背景

在进行POC时，现场人员进行数据导入时经常遇到各种问题，比较典型的是dble在导入文件时，对部分sql语句的不支持。另外，未分片的历史数据通过dble导入，旧数据会路由分片，在数据量较大时，耗时会比较长，在此过程中出现错误的话也很难排查。基于以上原因，2.19.0.0版本提供工具对mysqldump导出的源数据文件按照后端分片节点进行分割。分割后的数据文件可以在每个后端分片执行导入，适配数据按分片导入的需求。

### dump文件语句处理

1. create database: 会将逻辑数据库转换为物理库。
2. ddl语句: 根据表的分片节点写入到对应后端节点的dump文件中，对于自增列，会将自增列的数据类型修改为bigint。
3. insert: 全局序列列值会被dble替换为全局序列，再按照拆分列根据拆分算法路由到对应后端节点的dump文件中。
4. 一些属性设置的语句会根据最近一次解析的ddl来决定下发到具体的后端节点的dump文件中。
5. 会跳过对视图的处理
6. 会跳过对子表的处理

### 使用方法

#### 命令

在管理端口执行以下命令：

```
mysql > split src dest [-sschema] [-r500] [-w500] [-l10000]
```

- src: 表示原始dump文件名
- dest: 表示生成的dump文件存放的目录
- -s: 表示默认逻辑数据库名，当dump文件中不包含schema的相关语句时，会默认导出到该schema
- -r: 表示设置读文件队列大小，默认500
- -w: 表示设置写文件队列大小，默认500
- -l: 表示split后一条insert中最多包含的values，只针对分片表，默认4000

生成的分片文件以格式：源文件名-datanode名-时间戳.dump，最新的文件时间戳最大。

例如：我的原始dump文件是 /tmp/mysql\_dump.sql，我想切分以后输出到/tmp/dump/目录下：命令就是：

```
split /tmp/mysql_dump.sql /tmp/dump/
```

### 日志

默认情况下，split过程中生成的日志打印到在dble.log中，提供配置让split命令产生的日志单独存放，若需要开启，则需修改log4j.xml文件。

```
<Configuration status="WARN">
<Appenders>
    <!-- 将下面的此段配置追加至已安装dble的log4j.xml中的Appenders下 -->
    <RollingFile name="DumpFileLog" fileName="logs/dump.log"
        filePattern="logs/${date:yyyy-MM}/dump-%d{MM-dd}-%i.log.gz">
        <PatternLayout>
            <Pattern>%d{yyyy-MM-dd HH:mm:ss.SSS} %5p [%t] (%l) - %m%n</Pattern>
        </PatternLayout>
    <Policies>
        <onStartupTriggeringPolicy/>
        <SizeBasedTriggeringPolicy size="250 MB"/>
        <timeBasedTriggeringPolicy/>
    </Policies>
    <DefaultRolloverStrategy max="10"/>
    </RollingFile>
</Appenders>
<Loggers>
    <!-- 将下面的此段配置追加至已安装dble的log4j.xml中的Loggers下，可通过调整level为debug来调整性能 -->
    <Logger name="dumpFileLog" level="info" additivity="false" includeLocation="false" >
        <AppenderRef ref="DumpFileLog" />
        <AppenderRef ref="RollingFile"/>
    </Logger>
</Loggers>
</Configuration>
```

可通过日志中的“dump file has been read d%”关键字来查看解析进度。

可开启日志的debug级别来调整性能

### 任务停止

执行dump file任务的管理连接不受 idletimeout 参数的限制。用户可以通过kill @@connection id 方式杀掉管理连接从而停止dump file的任务的执行。

### 使用限制

1. 数据导入之后需要运维检查下数据完整性。
2. 对于使用全局序列的表，表原先全局序列中的值会被擦除，替换成全局序列，需要注意。
3. 暂时不支持子表的dump操作。

## 2.1.20 flow\_control 命令

### 2.1.20.1 查询流量控制当前配置状态

```
mysql> flow_control @@show;
+-----+-----+-----+
| FLOW_CONTROL_ENABLE | FLOW_CONTROL_START | FLOW_CONTROL_END |
+-----+-----+-----+
| false | 4096 | 256 |
+-----+-----+-----+
1 row in set (0.00 sec)
```

列描述:

- FLOW\_CONTROL\_ENABLE 是否有开启流量控制功能
- FLOW\_CONTROL\_START 流量控制功能触发值
- FLOW\_CONTROL\_END 流量控制功能取消值

### 2.1.20.2 修改流量控制当前配置状态

```
flow_control @@set [enableFlowControl = true/false] [flowControlStart = ?] [flowControlEnd = ?]
```

```
MySQL [(none)]> flow_control @@set enableFlowControl = true flowControlStart= 100 flowControlEnd = 30;
Query OK, 0 rows affected (0.02 sec)
```

通过此命令可以修改对应内存中生效的值，不会影响配置文件，对应关系如下：

- enableFlowControl server.xml中enableFlowControl参数
- flowControlStart server.xml中flowControlStartThreshold参数
- flowControlEnd server.xml中flowControlStopThreshold参数

注：以上命令三个参数顺序相关

### 2.1.20.3 查看当前流量控制中的连接

```
MySQL [(none)]> flow_control @@list;
+-----+-----+-----+
| CONNECTION_TYPE | CONNECTION_ID | CONNECTION_INFO | WRITE_QUEUE_SIZE |
+-----+-----+-----+
| ServerConnection | 2 | 127.0.0.1:60772/testdb user = root | 2920 |
+-----+-----+-----+
```

列描述:

- CONNECTION\_TYPE 连接的类型MySQLConnection/ServerConnection
- CONNECTION\_ID 连接在db中ID信息，可以通过ID查找日志
- CONNECTION\_INFO 连接详细信息，使用端口，IP地址，用户，MySQL中的连接ID等
- WRITE\_QUEUE\_SIZE 当前连接的写队列长度

## 2.2 全局序列

在分库分表的情况下，数据库自增主键无法保证自增主键的全局唯一。

为此，dble提供了全局序列，并且针对不同应用场景提供了多种实现方式。要应用全局序列，需要在server.xml中配置：

```
<system>
  <property name="sequenceHandlerType">1</property>
</system>
```

其中sequenceHandlerType包括如下种类：

- 1: MySQL offset-step方式
- 2: 时间戳方式
- 3: 分布式时间戳方式
- 4: 分布式offset-step方式

全局序列用于给自增列赋值，写入数据时字段值由系统自动生成，不可以再指定值，使用示例：

```
/*id为主键，配置为自增长*/
insert into table1(name) values('test');
```

### 2.2.1 MySQL offset-step方式

MySQL offset-step方式利用了数据库的并发和并发更新互斥的特性。

在数据库中建立一张表，保存序列的当前值和步长，使用该序列时，系统从表中读取序列的当前值，加上步长，返回给客户端的同时更新表中的当前值。线程获取序列的执行步骤如下：

1. 检查是否有请求序列未用的缓存的序列值。如果没有，执行步骤2；否则，执行步骤4。
2. 到相应的数据节点上执行 `SELECT dble_seq_nextval('seqName');` 其中，`seqName`为请求序列的序列名。此步骤获取的序列值数依赖于该序列的步长`increment`，为区间`[current_value, current_value+increment]`中的值。
3. 缓存步骤2得到的序列值。
4. 返回缓存的序列值中第一个未用的序列值。

以上提到的部分配置内容参见[1.7.1 offset-step方式](#)

## 2.2.2 时间戳方式

这种方式下，全局序列在dble服务实例本地产生，只能生成全局唯一的ID，不能实现连续自增。

使用这种方式需要对应字段为bigint来保证63位(63位的原因是Java没有无符号整数类型，所以最高位恒为0，保证全局序列是个正数)

序列值是63bits的整数。整数的位模式如下：

a.29bits	b.5bits	c.5bits	d.12bits	e.12bits
----------	---------	---------	----------	----------

其中，

- a - e为从高位到低位。
- a为系统当前时间截的低41位中的高29位。
- b为5位data center id。
- c为5位 worker id。
- d为12位自增长值
- e为系统当前时间截的低41位中的低12位。

注意事项：

1. data center id和worker id的最大值均为31。
2. 每毫秒时间内允许的最多序列值为4095。为了保证序列值的唯一性，在毫秒时间内请求超过4095个序列值时系统会进行等待到下一毫秒开始。
3. 因为java没有无符号整数，实际使用41位时间截相对于1288834974657L(2010年左右)的偏移量。
4. 相对于偏移量的处理过够后的41位时间截可供使用69年。

### 2.2.3 分布式时间戳方式

本方式提供一个基于Zookeeper(也可以本地配置)的分布式ID生成器，可以生成全局唯一的63位二进制ID。  
PS:63位的原因是Java没有无符号整数类型，所以最高位恒为0，保证全局序列是个正数

#### 2.2.3.1 位模式

序列值是63bits的整数。整数的位模式如下:

a.9bits	b.5bits	c.4bits	d.6bits	e.39bits
---------	---------	---------	---------	----------

其中:

- a - e为从高位到低位。
- a为线程id的低9位值。
- b为5位实例 id值(根据配置，此值为配置文件中的**INSTANCEID**值或者从zookeeper服务器获取的值, 参见[1.7.3 分布式time序列](#))。
- c为4位数据中心id值(即配置文件中的**CLUSTERID**的值, 参见[1.7.3 分布式time序列](#))。
- d为6位自增长值
- e为系统当前时间戳的低39位值(可以使用17年)。

#### 2.2.3.2 退化的分布式时间序列

如果配置文件中的**INSTANCEID**值不为'ZK', 序列的维护仅依赖于单实例 (主要是**INSTANCEID**值的维护)，此时序列类似于时间戳方式(参见[2.2.2 时间戳方式](#))。

#### 2.2.3.3 分布式时间序列

如果配置文件中的**INSTANCEID**值为'ZK', 序列的维护(主要是**INSTANCEID**值的维护)用zookeeper的临时自增节点来维持。每次生成全局序列时,向zk申请一个临时自增节点,通过计算自增节点数 % 32 获取**INSTANCEID**.

#### 2.2.4 分布式offset-step方式

分布式offset-step方式利用zookeeper的分布式锁功能进行实现。

线程在获取序列值时执行如下步骤:

1. 根据序列名(`schemaX`.`tableX`)获取相应的分布式锁。
2. 从zookeeper服务器上获取未用序列值的最小值min。
3. 根据配置(参见[1.7.4 分布式offset-step序列](#)),计算能获取的最大未用序列值max。
4. 用max+1更新zookeeper服务器上的未用序列值的最小值min。
5. 释放相应的分布式锁。

## 2.3 读写分离

### 2.3.1 读写分离条件

要实现读写分离必须满足如下条件:

1. 必须在schema.xml中配置readHost(参见1.2 schema.xml)而且balance配置不为0。
  2. SQL语句为select 或者show。
  3. 在非事务中。
- 当然, 也可以通过注释`#!dbe:db_type=slave, ... /`或者`!dbe:db_type=slave, ... /`强制发从(参见2.4 Hint).

### 2.3.2 负载均衡

dble通过在writeHost下配置多个readHost为读操作提供负载均衡(参见1.2 schema.xml)。负载均衡通过如下的链接获取来实现。链接获取有如下两个步骤:

- 1.参与读写分离的资格检查, 确定符合读写分离的候选主机集合
- 2.参与者承担读任务的负载均衡算法

#### 2.3.2.1 参与读写分离的资格检查

该算法在每次连接获取时提供可用的mysql实例集.

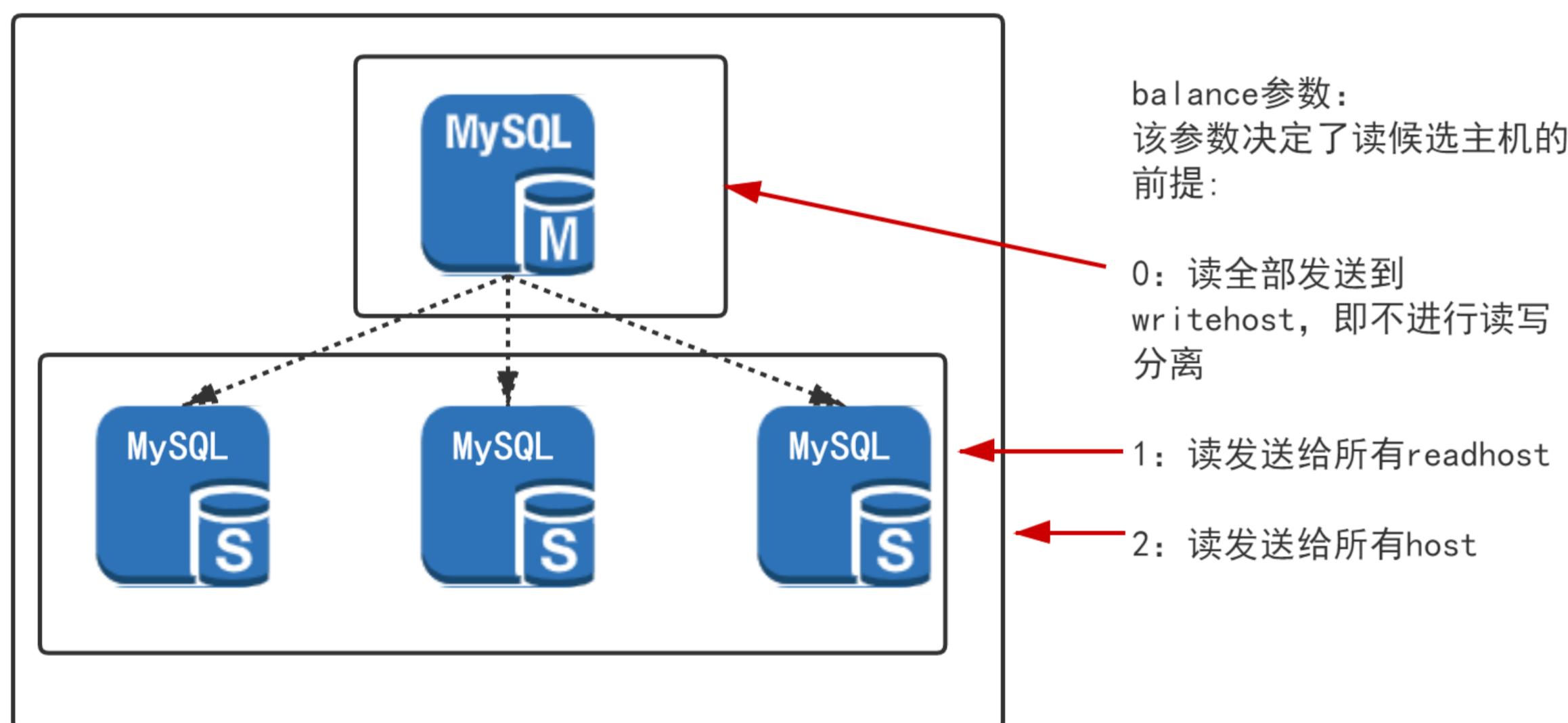
- 写节点正常(heartbeat状态正常)
  - 写节点参与均衡, 则有资格参与读分离中的读, 加入候选主机集合
  - 读节点
    - 节点正常且需要同步状态检测, 检查同步状态确定资格, 决定是否加入候选主机集合
    - 节点正常且不需要同步状态检测, 直接符合资格, 加入候选主机集合
- 写节点异常
  - 如果配置读节点临时可用(参数tempreadhostavailable), 是否符合资格, 参考写节点正常时读结点的资格检查方案.
  - 如果配置读节点不可用(参数tempreadhostavailable关闭), 不能加入加入候选主机集合.

#### 2.3.2.2 负载均衡算法

该算法在候选主机集中选择一个mysql实例以便获取连接.

- 候选主机集为空, 选择当前写主机.
- 候选主机集为非空
  - 有权重设置(weight参数), 但不是所有权重等值, 依权重随机选择.
  - 无权重设置或所有权重等值, 等权随机选择. 此种情况指示上面情况的特例.

#### 2.3.2.3 写节点是否参与均衡与datahost的balance属性有关, 具体见下见下图



## 2.4 注解/Hint

**Hint**, 即注解。我们定义为:在要执行的SQL语句前添加额外的一段由注解组织的代码, 这样SQL就能按照编写者的意图执行, 这段代码称之为“注解”。注解的作用如下:

- 指定路由, 比如强制读写分离。
- 帮助dble支持一些不能实现的语句, 如单节点内存储过程的创建和调用, 如insert...select...;

原理解释: 分布式环境下执行SQL语句的流程是先进行SQL解析处理, 计算出路由信息后, 然后到对应的物理库上去执行; 若传入的SQL语句无法解析, 则不会去执行。注解则可以告诉解析器按照注解内的SQL (称之为注解SQL) 去进行解析处理, 解析出路由信息后, 将真正要执行的SQL语句 (称之为原始SQL) 发送到对应的物理库上去执行。

### 2.4.1 Hint语法

Hint语法有两种形式:

1. /\*!dbe:type=...\*/
  2. /\*#dbe:type=...\*/
- ```
/*#dbe: */ for mybatis and /*!dbe: */ for mysql
```

其中, type有3种值可选: datanode, db\_type, sql。每一种值的功能和形式详见各个部分的具体说明。

### 2.4.2 类型datanode

1. 形式  
datanode=node  
其中, node为单个数据节点名, 不能为多值(node定义参见配置1.2 schema.xml)。
2. 功能  
为不方便路由或者不能路由的语句指定具体的目的数据节点。

### 2.4.3 类型db\_type

1. 形式  
db\_type=master或者db\_type=slave
2. 功能  
帮助实现正确的业务逻辑, 强制读写分离。
3. 注意事项  
delete, insert, replace, update, ddl语句不能使用db\_type=slave进行注解。

### 2.4.4 类型sql

1. 形式  
sql=sql\_statement
2. 功能  
用sql\_statement的路由结果集作为实际sql语句的执行数据节点。支持存储过程。

### 2.4.5 注意事项

写注解需要注意如下事项:

- dble的注解和MySQL原生注解含义不同, 想通过MySQL原生注解来设置变量或者指定索引是无法得到预期结果的。如[#1169](#)
- 使用select语句作为注解SQL, 不要使用delete/update/insert等语句。delete/update/insert等语句虽然也能用在注解中, 但这些语句在SQL处理中有一些额外的逻辑判断, 会降低性能, 不建议使用;
- 注解SQL本身禁用表关联语句, 注解目的是路由计算, 如果本身写得过于复杂, 会影响路由计算;
- 使用hint做DDL需要额外执行reload @@metadata
- 使用hint做session级别的系统变量和环境变量可能不会生效, 请慎用
- 使用注解并不额外增加的执行时间; 从解析复杂度以及性能考虑, 注解SQL应尽量用最简单的SQL语句, 如select id from tab\_a where id='10000';
- 能不用注解也能够解决的场景, 尽量不用注解。

## 2.5 分布式事务

- 2.5.1 XA事务概述
- 2.5.2 XA事务的提交以及回滚
- 2.5.3 XA事务的后续补偿以及日志清理
- 2.5.4 XA事务的记录
- 2.5.5 一般分布式事务概述

## 2.5.1 XA事务概述

### 2.5.1.1 XA事务概述

普通事务只能在单节点内部保证事务的完整性，如果事务要在不同的节点上执行，无法保证数据一致性，具有跨节点的数据一致性要求的场景需要使用2阶段协议实现分布式事务。由于以上普通事务的缺陷所以Dble引入了Mysql的XA事务来解决这个问题，MySQL5.7之前版本的XA事务存在一些问题，启用需要MySQL 5.7版本，否则无法保证数据不丢失。dble提供的分布式事务采用两阶段提交协议，目前还是弱XA事务，不能绝对保证跨节点数据的强一致性。

1. 事务开始前需要设置手动提交: `set autocommit=0;`
2. 使用命令开启XA 事务: `set xa=on;`
3. 执行相应的SQL 语句;
4. 对事务提交或者回滚，事务结束: `commit/rollback;`

分布式事务的性能开销比较大，所以只推荐在全局表的事务以及其他一些对一致性要求比较高的场景中使用。一般情况下尽量不要在同一个事务中运行跨节点的SQL语句。dble支持大事务，但一方面大事务会造成事务执行时间上升，事务信息规模扩大，导致系统性能下降；另一方面，大事务也容易造成潜在的事务数据不一致问题。因此，大事务一定要谨慎使用，一般建议单个事务的 SQL 语句不要超过100条。

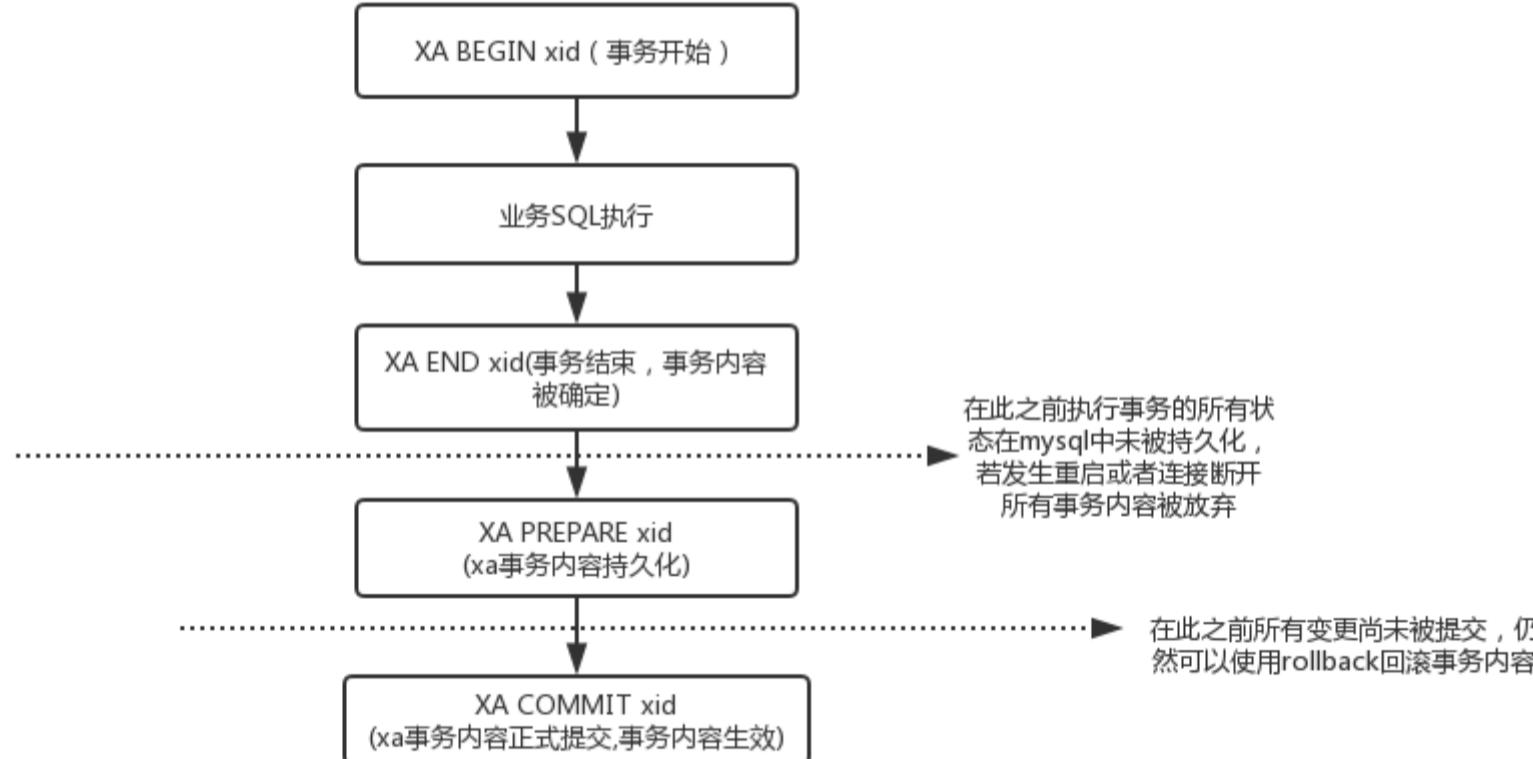
以下提供一个JDBC使用XA事务的实例

```
public class XaDemo {
    public static final String URL = "jdbc:mysql://localhost:8066/testdb";
    //在这里也可以使用jdbc:mysql://127.0.0.1:8066?sessionVariables=xa=1
    //进行替代，在这种情况下不需要执行set xa = 1
    public static final String USER = "root";
    public static final String PASSWORD = "123456";

    public static void main(String[] args){
        try {
            //1.加载驱动程序
            Class.forName("com.mysql.jdbc.Driver");
            //2. 获得数据库连接
            Connection conn = DriverManager.getConnection(URL, USER, PASSWORD);
            //3.操作数据库，实现增删改查
            Statement stmt = conn.createStatement();
            stmt.execute("set xa = 1");
            //开始一个xa事务
            stmt.execute("begin");
            try {
                //执行相关数据操作的时候需要对于可能出现的错误进行catch
                //并在错误出现的时候对于整个事务进行 rollback
                stmt.execute("insert into xa_test set id = 11,name = '3333'");
                stmt.execute("insert into xa_test set id = 22,name = '333'");
                stmt.execute("insert into xa_test set id = 3,name = '33'");
                //数据执行完成提交
                stmt.execute("commit");
            }catch (Exception e){
                System.out.println(" error "+e);
                //如果在数据操作的时候出现错误，将整个事务的操作回滚
                stmt.execute("rollback");
            }finally {
                stmt.close();
                conn.close();
            }
        }catch(Exception e){
        }
    }
}
```

### 2.5.1.2 XA事务的基础

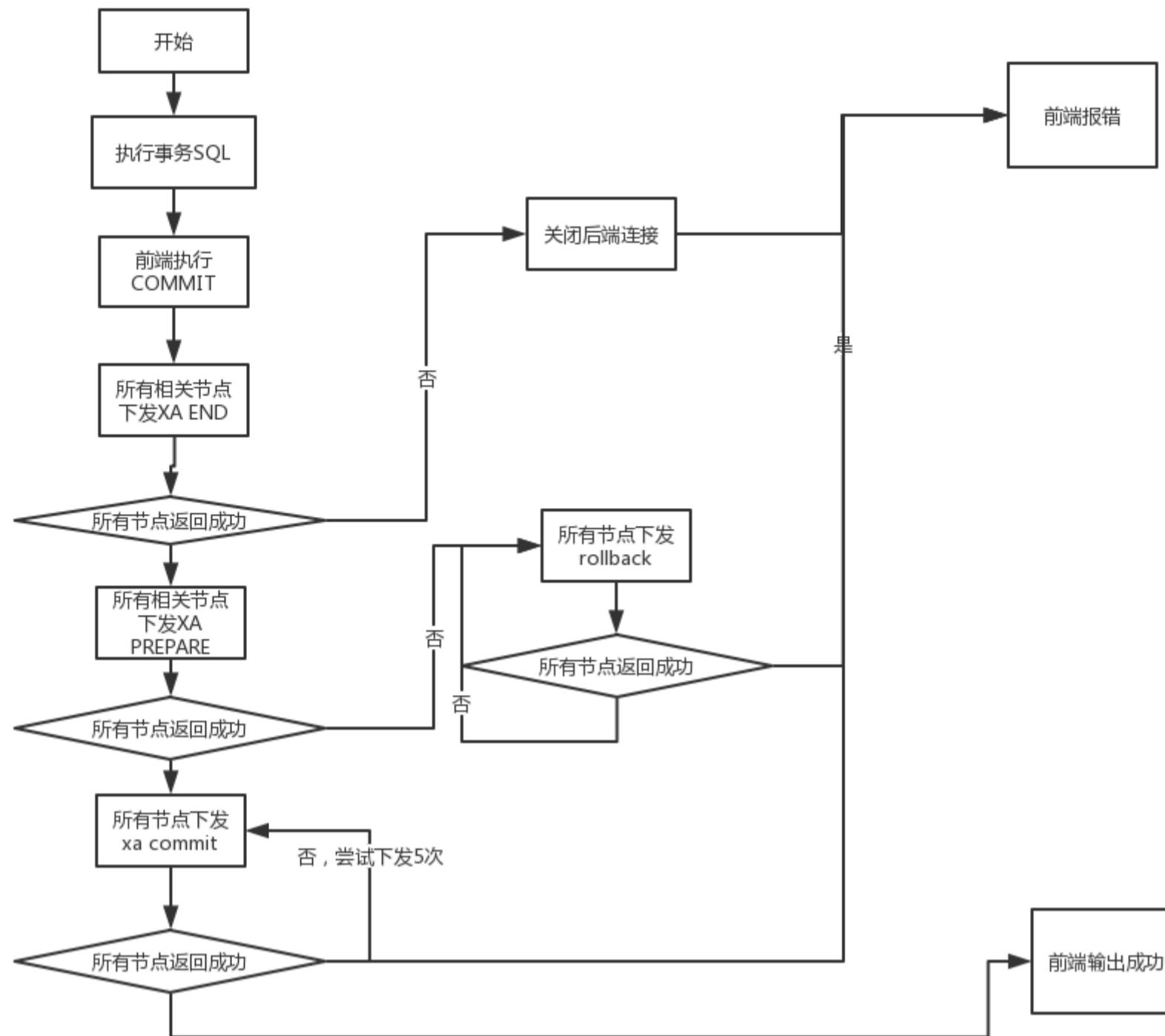
xa事务基于mysql5.7的xa事务的特性，其流程和特性由以下图所示:



## 2.5.2 XA事务的提交以及回滚

### 2.5.2.1 XA事务的提交

在Dble中采用二段提交的方式对于XA事务进行提交 具体的逻辑可见下图

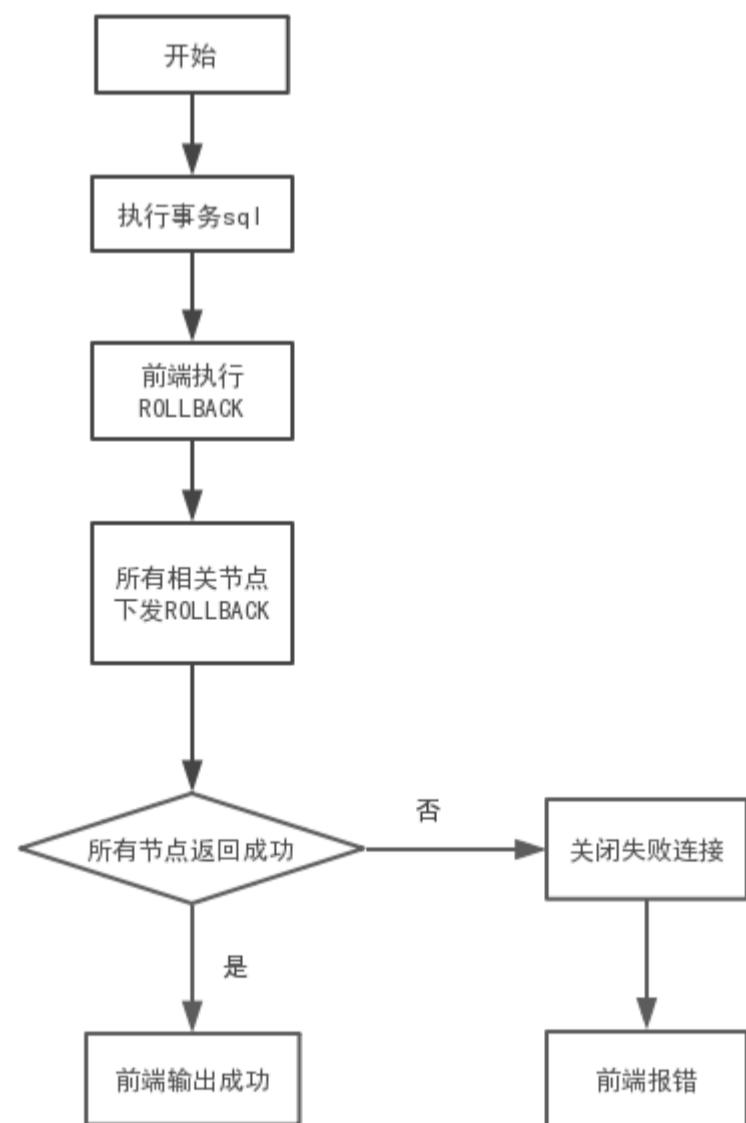


整体来说的处理原则如下：

1. 将XA事务的提交分为END PREPARE COMMIT三个部分
2. 如果在PREPARE下发之前有节点失败或报错，关闭所有后端连接放弃事务数据
3. 如果在PREPARE下发过程中发生失败，则回滚事务，所有节点下发ROLLBACK
4. 如果在COMMIT节点发生失败，则尝试重新下发，几次尝试未果将事务交给定时任务来继续重试

### 2.5.2.2 XA事务的回滚

相对来说回滚的逻辑就容易的多，直接在所有节点下发rollback。如果失败，直接关闭失败的连接。



### 2.5.2.3 XA事务重试机制

2.19.03.0之前的dble版本对于失败的事务，策略是将事务交给后台进行无限重试。2.19.03.0版本对这一过程进行了可配置化，并提供管理命令控制这一过程。

#### 2.5.2.3.1 配置

2.19.03.0版本可以通过server.xml文件中的xaRetryCount属性配置xa后台重试策略：

1. 当 xaRetryCount 等于0时，后台无限重试
2. 当 xaRetryCount 大于0时，后台尝试次数达到xaRetryCount后，重试停止

若重试失败，会发出一条告警，重试成功后自动解决相应告警。

#### 2.5.2.3.2 相关命令

2.19.03.0版本新增两个管理命令

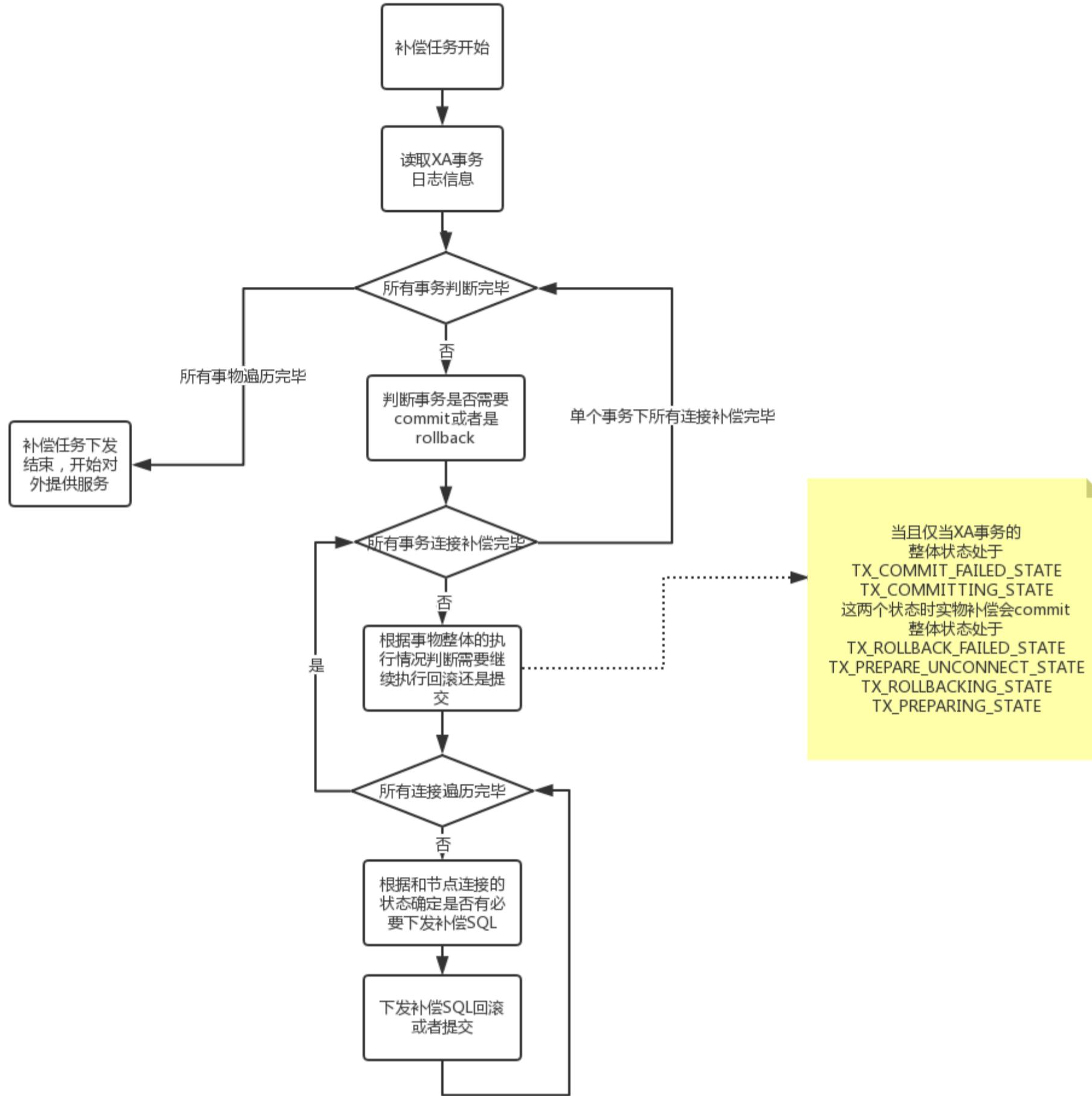
#### 1.6.1 cache配置

1. 通过 `show @@session.xa` 查看后台重试xa事务信息
2. 通过 `kill @@xa_session id1,id2...` 取消指定session后台尝试xa事务

### 2.5.3 XA事务的后续补偿以及日志清理

#### 2.5.3.1 XA事务的补偿逻辑

由于XA事务是通过多次提交来达成最终的提交和回滚的，所以就会出现由于服务下线或者其他原因导致的提交或是回滚任务被执行了一半。在这种情况下发生的时候dble通过XA事务记录的日志进行日志的处理以及补偿。详细通过在dble重启的时候在重启之前对于XA日志进行读取，并根据里面的内容对于XA事务进行回滚或者补偿。具体的补偿流程如下图所示：



注：XA事务的补偿不保证所有事务都能在补偿期间正确提交，若补偿期间依然因为SQL执行造成失败，只会在dble普通日志中体现或者告警（如果有配置告警功能）

#### 2.5.3.2 XA事务的日志清理逻辑 (xaLogClean)

xa事务的日志只保存最近一段时间的，所以需要有定期的清理已经正确提交或者是正确回滚的日志数据。会根据server.xml中的配置信息xaLogCleanPeriod定时将没有后续作用的数据做清除。

#### 2.5.3.3 XA事务定期连接检查 (xaSessionCheck)

在dble提供服务期间如果发生commit(狭义commit,特指XA事务中prepare之后的commit动作)失败或者是 rollback失败，那么失败的事务将被存储到内存队列中，并且进行定时的提交或回滚，直到事务被正确的提交或回滚。这是由于在xa事务的逻辑中prepare全部完成之后的事务已经都被成功持久化，仅需要提交即可。

## 2.5.4 XA事务的记录

### 2.5.4.1 XA 事务过程中记录的内容

由于在Dble中采用两段提交的分布式事务，所以使用XA事务的时候对于DBLE本身就拥有了状态。状态就需要有文件或者其他方式的记录，其中关于XA事务细节的记录主要是记录以下几个部分

1. 事务ID
2. 事务状态
3. 事务中每个节点的连接host
4. 事务中每个节点的连接端口
5. 事务中每个节点连接最后的事务状态
6. 事务中每个节点连接的过期状态(没有实际作用)
7. 事务中每个节点连接对应的后端数据库

这里举例一个记录的实例

```
{
  "id": "'Dble_Server.1.15'",
  "state": "8",
  "participants": [
    {
      "host": "10.186.24.37",
      "port": "3308",
      "p_state": "8",
      "expires": 0,
      "schema": "db3"
    },
    {
      "host": "10.186.24.37",
      "port": "3306",
      "p_state": "8",
      "expires": 0,
      "schema": "db2"
    },
    {
      "host": "10.186.24.37",
      "port": "3308",
      "p_state": "8",
      "expires": 0,
      "schema": "db2"
    },
    {
      "host": "10.186.24.37",
      "port": "3306",
      "p_state": "8",
      "expires": 0,
      "schema": "db1"
    }
  ]
}
```

### 2.5.4.2 XA事务中status的标识字典

| status | 状态                         | 解释                                            |
|--------|----------------------------|-----------------------------------------------|
| 0      | TX_INITIALIZE_STATE        | XA事务处于初始化状态                                   |
| 1      | TX_STARTED_STATE           | XA事务处于开始状态，在事务开始直到提交或者回滚之前<br>XA事务的状态一直会保持此状态 |
| 2      | TX_ENDED_STATE             | XA END下发成功状态                                  |
| 3      | TX_PREPARED_STATE          | XA PREPARED成功状态                               |
| 4      | TX_PREPARE_UNCONNECT_STATE | XA PREPARED下发过程中连接被断开                         |
| 5      | TX_COMMIT_FAILED_STATE     | XA COMMIT 下发失败                                |
| 6      | TX_ROLLBACK_FAILED_STATE   | XA ROLLBACK 失败                                |
| 7      | TX_CONN_QUIT               | 后端mysql连接失败                                   |
| 8      | TX_COMMITED_STATE          | XA 事务提交成功                                     |
| 9      | TX_ROLBACKED_STATE         | XA 事务回滚成功                                     |
| 10     | TX_COMMITTING_STATE        | XA 事务正在提交                                     |
| 11     | TX_ROLLBACKING_STATE       | XA 事务正在回滚                                     |
| 12     | TX_PREPARING_STATE         | XA 事务正在下发prepare                              |

### 2.5.4.3 XA事务记录的存储方式

#### 一、本地文件方式

顾名思义在这种方式下，xa事务的状态将以本地记录文件的方式被存放到对应的文件中，具体的路径和文件名配置是

配置与server.xml中的{xaRecoveryLogBaseDir}/{XaRecoveryLogBaseName}.log默认条件下文件会被储存在./tmlogs/tmlogs-1.log

一般只在Dble单机状态下使用本地文件方式，使用集群时本地文件的方式将在集群状态下造成不可预知的错误

#### 二、ZK存储方式

ZK存储方式不需要额外的配置，当Dble使用ZK配置时，自动默认XA事务记录的存储方式也会是ZK存储

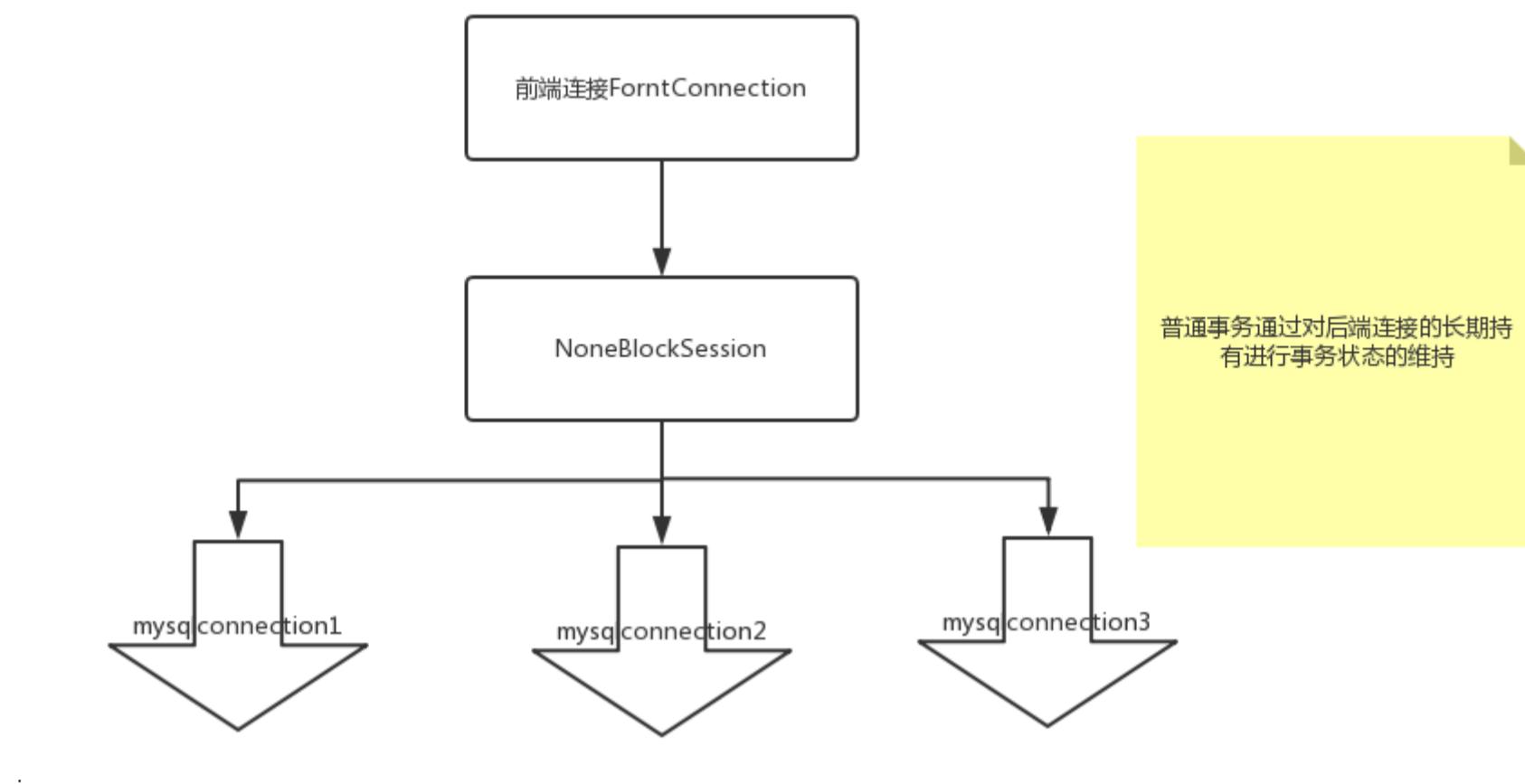
具体的XA事务记录的内容保持不变，记录在dble/{clusterId}/XALOG/{myid} 的Key值中

## 2.5.5 一般分布式事务概述

在分布式事务中整体的逻辑和mysql的事务逻辑类似，通过长期持有的连接来进行，每个前端连接frontconnection对应一个session，在dble的每个session中有对应的事务状态以及session所持有的后端连接集合target，在非事务状态下或者是autocommit状态下每次后端连接被使用完毕之后就会被移除target并释放回空闲连接池，但是在事务开启的状态下，在SQL执行完毕的时候connection会在target中长期储存，直到session发起commit或者是rollback。

综上我们可以看到，Dble中的普通分布式事务其实就是后端mysql事务的集合，并且这个事务是没有文件记录的，由于mysql的事务特性在事务发生的过程中若断开连接等同于放弃事务，所以可能出现在commit的过程中由于各种意外情况导致事务的部分提交，例如连接后端四个节点dn1,dn2,dn3,dn4在提交commit进行依次下发的时候dn1,dn2,dn3都提交成功，但是dn4由于网络意外提交失败，导致了预期执行的部分内容丢失，并且由于dn1,dn2,dn3已经提交成功无法进行数据回滚，只能进行人工的数据补偿。

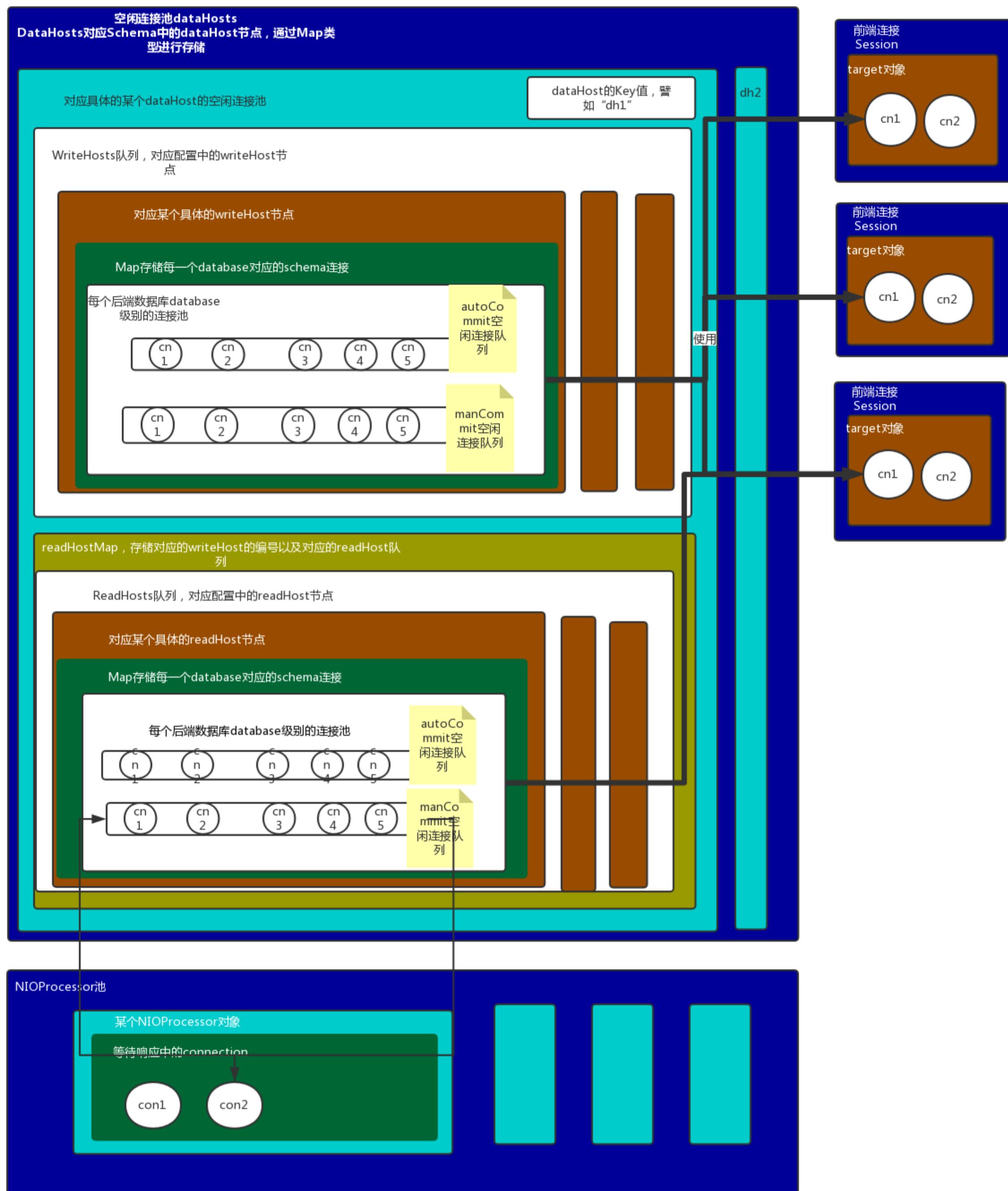
其逻辑图如下：



## 2.6 连接池管理

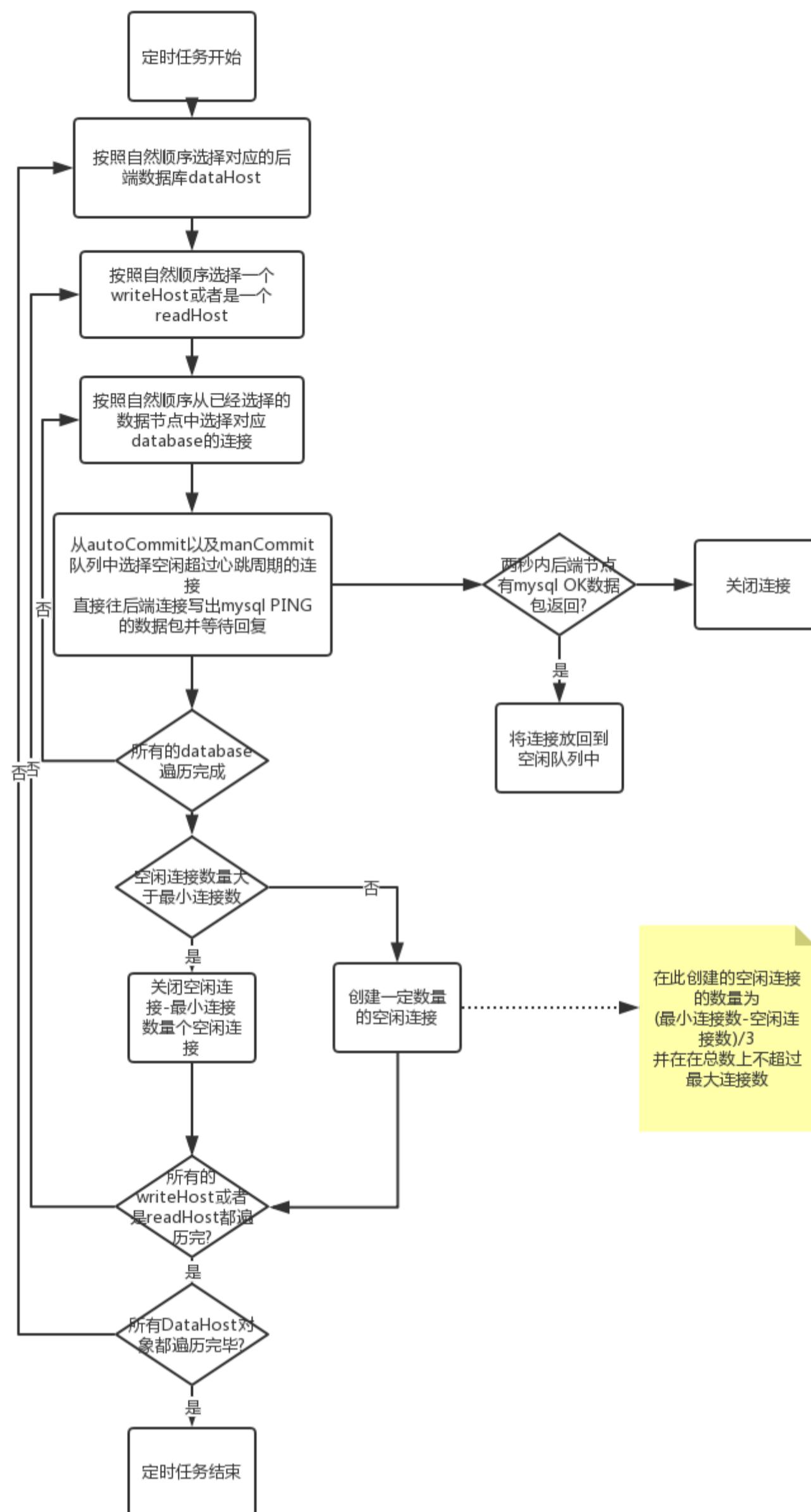
### 2.6.1 Dble后端连接池的存储结构

Dble的后端连接池基本分为两个部分，一个是在空闲的部分，一个是在等待响应的部分 在空闲的部分基本结构符合在schema.xml配置文件里面的配置结构，在接近底层融合了配置信息内的dataNode的结构，最终形成面向具体实例具体database的空闲连接队列：



### 2.6.2 Dble后端连接池的心跳管理

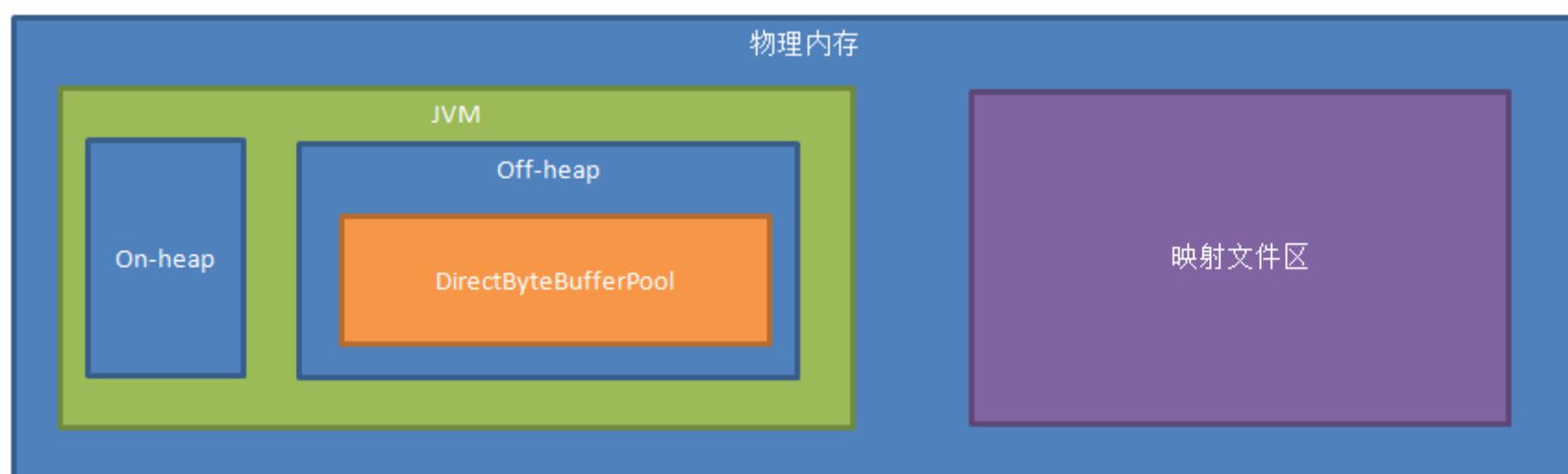
Dble后端节点的管理和心跳是通过定时任务dataNodeConHeartBeatCheck来完成的，此定时任务以server.xml中的dataNodeIdleCheckPeriod配置为周期，定时查看空闲的定时任务 其流程基本逻辑如下：



1. 遍历所有物理节点对应所有database的空闲连接队列
2. 选取其中的空闲超过一个心跳周期的连接进行心跳
3. 心跳在2S内能得到响应则放回，否则关闭连接
4. 通过最小连接数和当前空闲连接数之前的差值确定需要提前准备空闲连接或者是关闭多余的空闲连接

## 2.7 内存管理

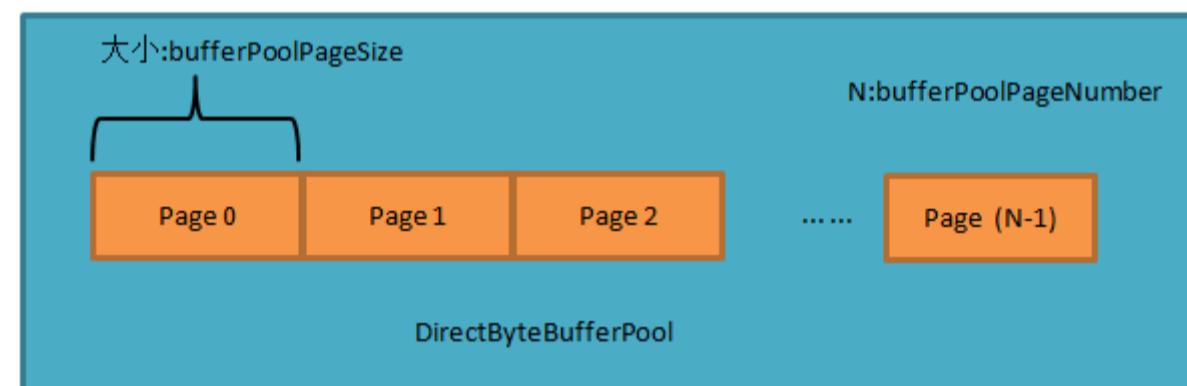
### 2.7.1 内存结构概览:



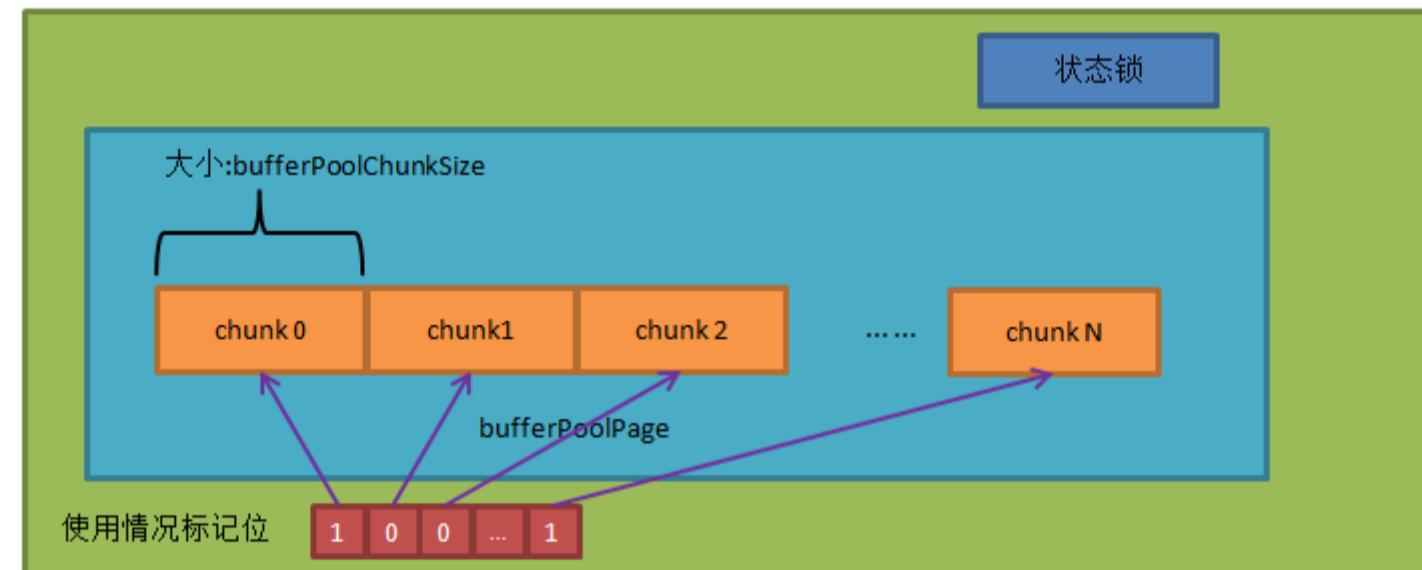
- On-Heap 大小由JVM 参数Xms ,Xmx 决定，就是正常服务需要的内存，由jvm自动分配和回收。
  - Off-Heap大小由JVM 参数 XX:MaxDirectMemorySize 大小确定。
  - DirectByteBufferPool 大小 = bufferPoolPageNumber\*bufferPageSize
- bufferPoolPageNumber和bufferPageSize可在Server.xml 配置， bufferPageSize默认为2M, bufferPoolPageNumber默认为： MaxDirectMemorySize \* 0.8 /bufferPageSize
- 映射文件区不在JVM之内  
大小计算方法求得值tmpMin= Min(物理内存的一半, free内存)  
可在映射区保存的文件的个数=(向下取整(tmpMin/mappedFileSize))  
真实使用值= 可在映射区保存的文件的个数\*mappedFileSize (mappedFileSize默认64M)  
其中mappedFileSize可由Server.xml 指定  
注：这里是通过单例模式，一次性确定映射文件区的大小，如果运行过程中，内存被其他进程占用，这里可能存在风险

### 2.7.2 DirectByteBufferPool

BufferPool结构概览:



单个bufferPoolPage的内部结构:



#### 1.用途:

- 1.1 网络读写时使用  
1.2 中间结果集暂存时用于写数据前的缓存buffer

#### 2.初始化:

按照参数，初始化为bufferPoolPageNumber个页，每个页大小为bufferPageSize

#### 3.内存分配

##### 3.1 分配大小

如果不指定分配大小，则默认分配一个最小单元(最小单元由bufferPoolChunkSize决定，默认大小为 4k,请最好设为bufferPageSize 的约数，否则最后一个会造成浪费)。

如果指定分配大小，则分配放得下分配大小的最小单元的整数倍(向上取整)。

总之，大小为 M\*bufferPoolChunkSize

##### 3.2 分配方式

遍历 缓冲池从N+1页到bufferPoolPageNumber-1页(上次分配过的记为第N页)

对单页加锁在每个页中从头寻找未被使用的连续M个最小单元  
  
如果没找到，再从第0页找到第N页  
  
以上成功后更新上次分配页，标记分配的单元  
  
如果找不到可存放的单页(比如大于bufferPageSize )，直接分配On-Heap内存

#### 4.内存回收

##### 4.1 如果是On-Heap内存

直接clear，等GC回收

##### 4.2 如果是Off-Heap内存

遍历所有页，找到对应页

对单页加锁，到对应页的对应块的位置，标记为未使用

### 2.7.3 处理中间结果集过大时内存使用

目前dble 对每个session的内存管理如下：

- Join内存管理，最大上限为4M，目前写死，用于管理join操作暂存的数据
- Order内存管理，最大上限为4M，目前写死，用于管理排序操作暂存的数据
- Other内存管理，最大上限为4M，目前写死，用于管理distinct group，nestloop操作暂存的数据

每一个复杂查询中，当子查询单元需要在中间件当中暂存数据的时候，数据会存在Heap内存当中，但如果当前存储使用的内存大于4M，则需要写内存映射文件。

如果内存映射文件个数达到上限（参见概览中的可在映射区保存的文件的个数），则会去写硬盘。

写文件的时候，当单个文件大于mappedFileSize时，会将文件拆分。

注1：内存映射文件总大小在第一次使用时就确定，有风险

注2：写硬盘时候的缓冲区是从DirectByteBufferPool申请的chunk大小的

## 2.8 集群同步协调&状态管理

### 2.8.1 概述

大多数时候，**dble**结点是无状态的，所以可以用常用的高可用/负载均衡软件来接入各个结点。  
这里不讨论各个负载均衡软件的使用。  
主要讨论一下某些情况下需要同步状态的操作和细节。  
注：本部分内容需要额外部署zookeeper用于管理集群的状态和同步。

### 2.8.2 配置

配置文件 myid.properties:

```
#不使用/单节点部署请设置为false
loadZk=true
#zk的clinet端入口
zkURL=127.0.0.1:2181
#zk上的命名空间,同一集群使用同一个
clusterId=server-cluster-1
#每个结点必须不同
myid=server_fz_01
```

### 2.8.3 初始化状态

#### 方式1.

通过执行脚本init\_zk\_data.sh方式将某个结点的配置文件等数据写入ZK,所有结点启动时都从ZK拉取配置数据。

#### 方式2.

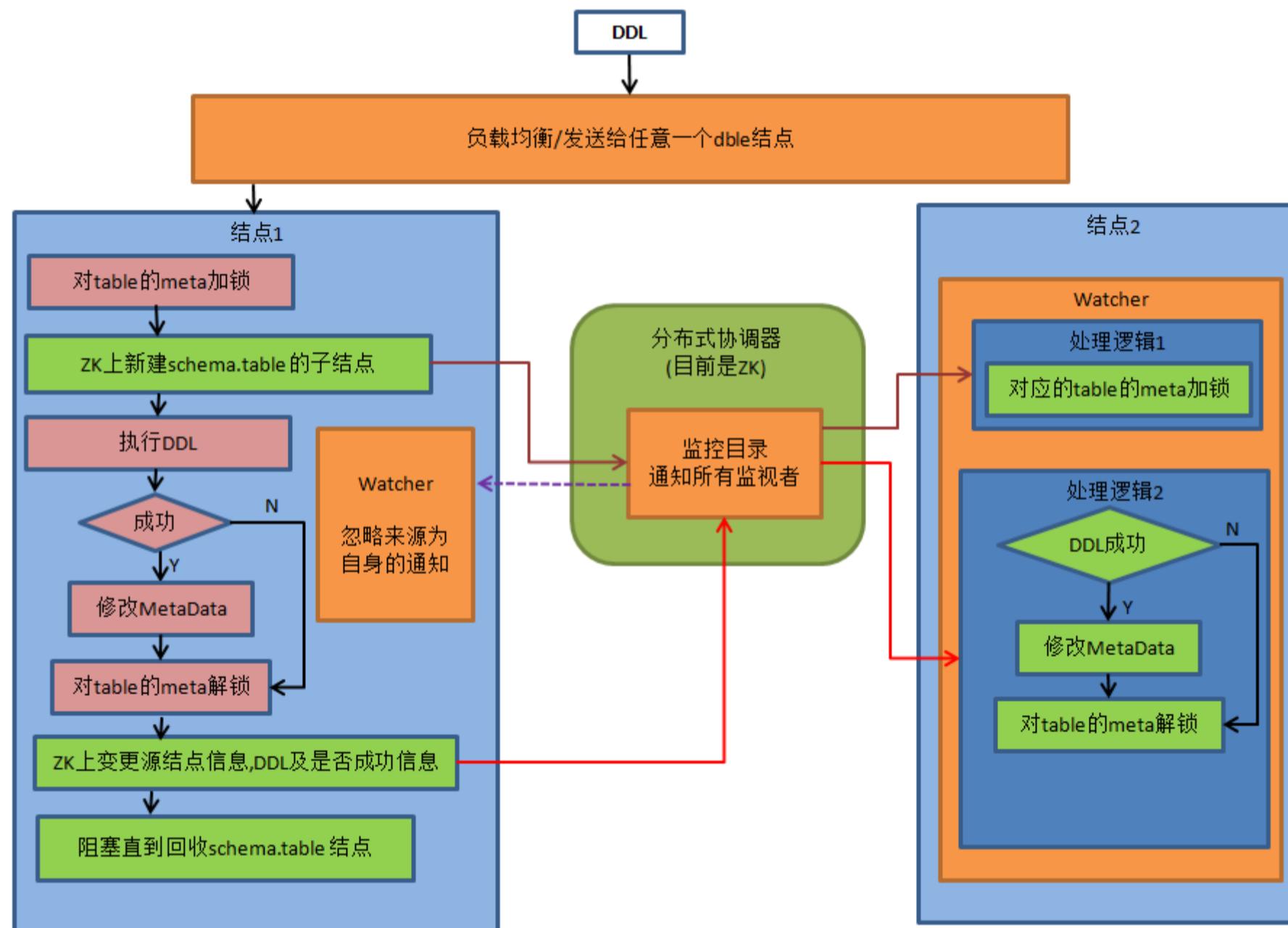
在第一个结点第一次启动时自动将自己的配置文件写入ZK,其他结点启动时从ZK拉取。

第一个结点的判定：用分布式锁抢占的方式，未抢占到结点会阻塞等待直到获取到分布式锁，如果此时初始化标记已经被设置，则从ZK拉取配置，否则将自己本地配置写入。

### 2.8.4 状态同步

#### A.DDL

做DDL时候会在执行的某个节点成功后，将消息推给ZK，ZK负责通知其他节点做变更。流程如下图：

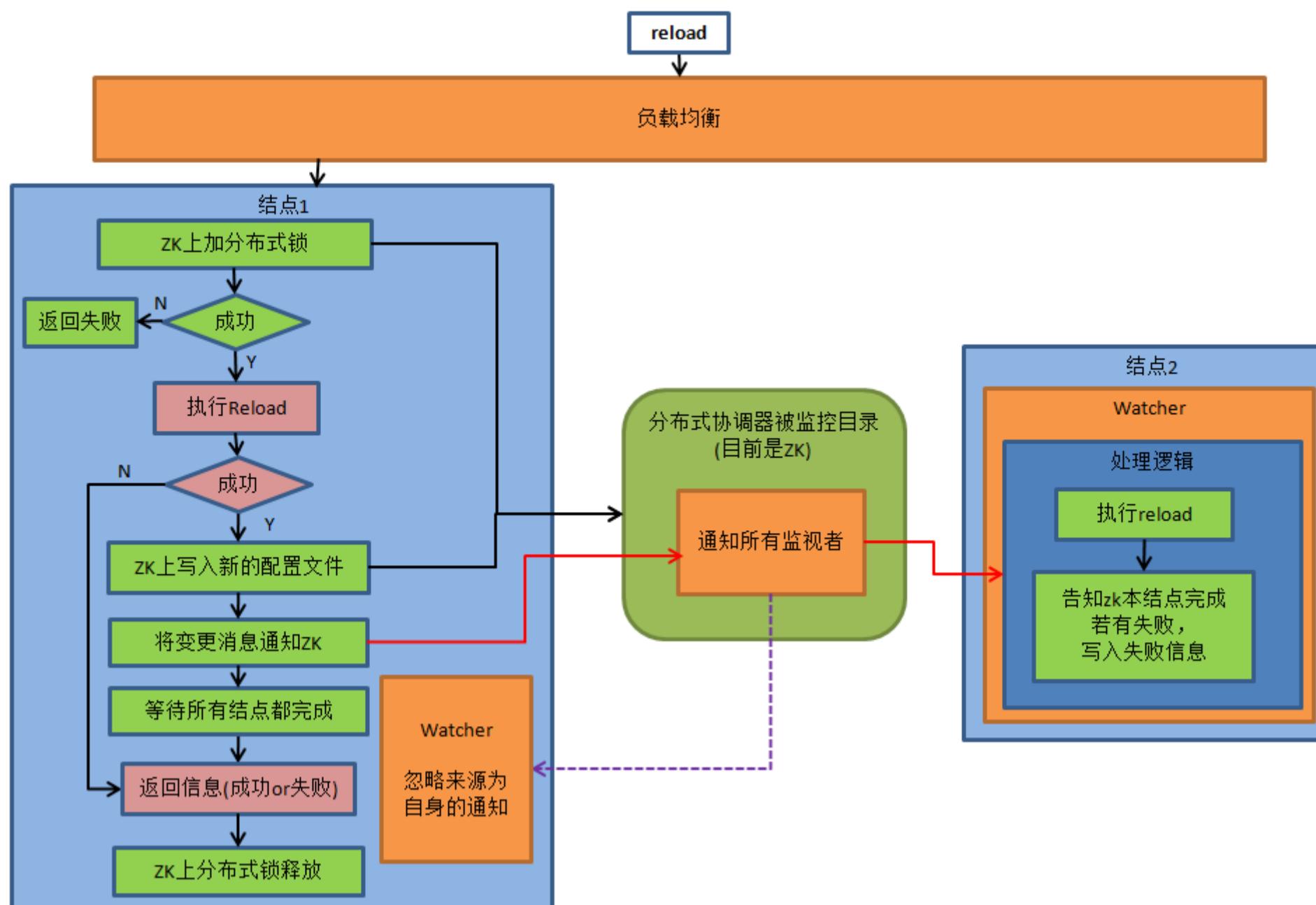


注：

- 新结点启动时，加载元数据前会检测是否有其他结点在做DDL变更，如果有，则等待。  
否则加分布式锁防止加载元数据期间其他结点做DDL变更，直到元数据加载完释放分布式锁。
- 回收指的是：每个节点ddl完成后会告知分布式协调器自己已经完成，并检查如果所有结点均完成，则删除schema.table 结点。此操作为原子操作。
- 发起结点故障：如果执行DDL的结点故障下线，其他结点会侦听到此消息，保证解开对应结点的tablemeta锁，并记录故障告警（如果配置了告警通道），需要运维人工介入修改ZK对应ddl结点的状态，检查各个结点meta数据状态，可能需要reload metadata。
- 逻辑上不应该有某个监听节点上加载meta失败的情况，如果发生了，告警处理  
(人工介入对应结点的meta是不正确的,需要reload meta)
- 注:view目前是异步模式,可能存在某个间隙view修改成功，查询仍旧拿到旧版的view结构。

#### B.reload @@config/ reload @@config\_all/ rollback @@config

执行流程如下图：



注：如果在部分结点失败，则会返回错误及错误原因以及结点名。

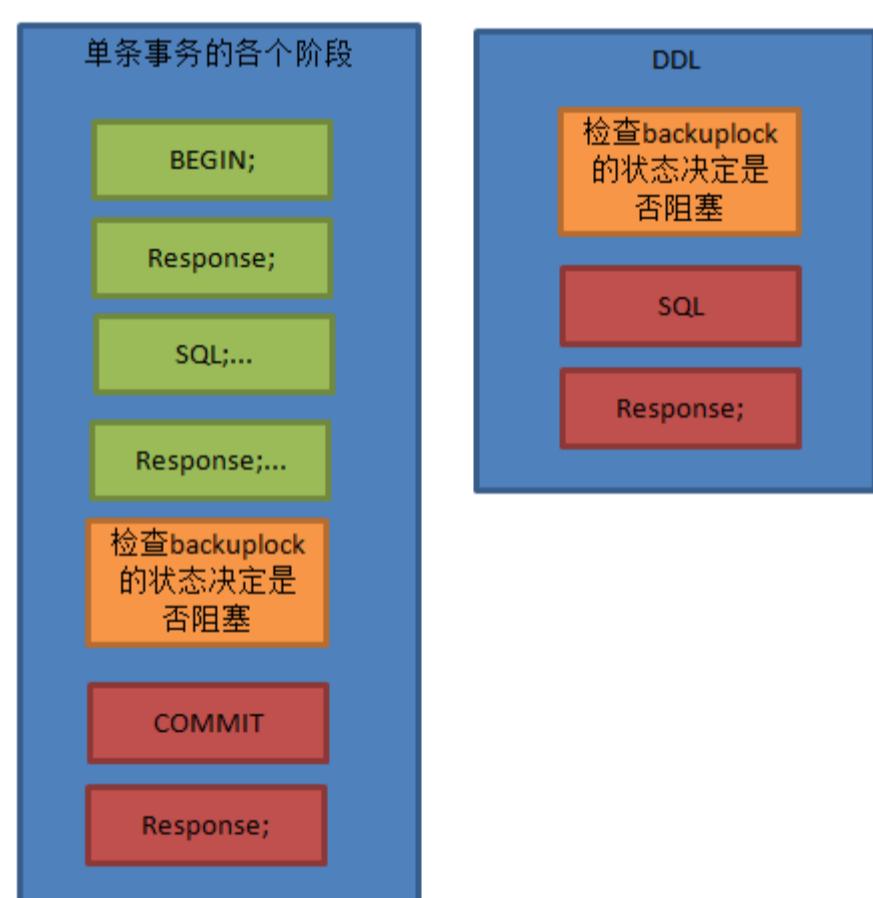
### C. 拉取一致性的binlog线

目的：

获得后端数据库实例的一致性binlog位置。由于两阶段提交的第二阶段执行在各结点无法保证时序性和同步性，所以直接下发表show master status获取binlog可能会造成不一致。

实现方式：

如下图，当前端收到show @@binlog.status语句时，遍历当前所有活动session查看状态。  
若session处于绿色区域，则在进入红色区域前等待知道show @@binlog.status结果返回。  
若存在session处于红色区域，则需要等待所有红色区域的session返回结果走出红色区域后下发表show master status。



超时处理：

此处有可能有死锁发生。

场景：session1 正在更新tableA，处于绿色区域，session2 下发有关于tableA 的DDL，等待 metaLock解锁，处于红色区域.session3 下发show @@binlog.status。  
此时session1 等待session3, session2等待session1, session3等待session2。

因此引入超时机制。如果session3 等待超过showBinlogStatusTimeout(默认60s, 可配置)，自动放弃等待，环状锁解除。

集群协调：

1. 收到请求后同步通知ZK，先等待本身结点准备工作结束，之后ZK通知其他结点处理。
2. 所有结点遍历各自的活动的session，进入红色区域的等待处理完成，绿色区域的暂停进入红色区域。
3. 结点将准备好/超时将状态上报给ZK
4. 主节点等待所有结点状态上报完成之后，判断是否可以执行任务，若是，则执行show @@binlog.status并返回结果，否则报告本次执行失败。
5. 主节点通过ZK通知各结点继续之前的任务

集群超时处理：

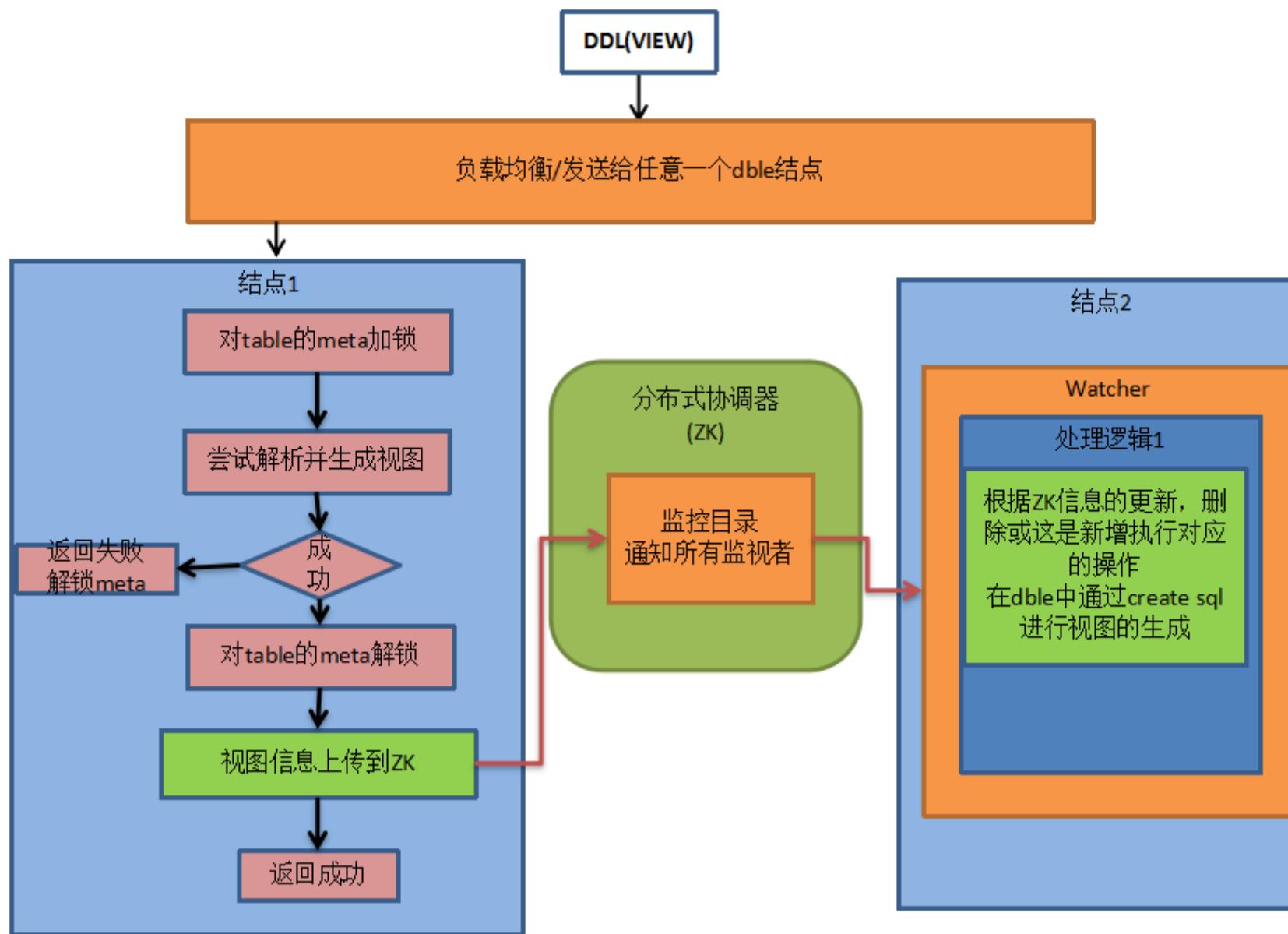
若有结点超时未准备好，主节点会报超时错误，并通过ZK通知各结点继续之前的任务。

故障处理：

主节点执行过程中故障下线，其他结点会感知，保证自己结点一定时间后自动解锁继续原有任务。  
ZK状态需要人工干预。人工等待所有结点超时之后，手动删除/修改ZK上的状态信息以便下次执行时不出问题。

### D. View管理

在使用集群模式时，使用ZK进行view视图信息的管理，使得整个db表集群能够进行视图信息的同步。由于视图只用于数据查询且不会造成数据异常的属性，在视图同步时采用异步同步的方法。



注 在zk上的view数据信息如下形式 key schema.table value {"serverId":"create\_Server\_id", "createSql":"view\_create\_sql"}

#### E.其他(暂未实现)

cache配置变化 (暂未实现专门入口) :  
reload 时会将变化写入zk并通知其他结点同步

新增自增序列表 (暂未实现专门入口) :  
reload未比对并添加新表的配置。  
包括sequence\_conf.properties 和sequence\_db\_conf.properties 方式

#### 2.8.5 XA日志管理

在使用集群模式时, 未完成的XA日志会存放在zookeeper上, 更加安全, 防止某台机器硬盘物理损坏导致日志丢失 (此处可能会有并发高吞吐引发的性能及其他问题, 待测试)。

#### 2.8.6 ZK整体目录结构

```

dble
  server-cluster-1
    conf
      init
      status
      schema
        schema
        dataNode
        dataHost
      server
        default
        user
        firewall
      rules
        tableRule
        function
    cache
      ehcache.xml
      cacheservice.properties
    sequences
      leader
      instance
      incr_sequence//批量步长方式
        table_name
      common
        sequence_conf.properties
        sequence_db_conf.properties
        sequence_distributed_conf.properties(INSTANCEID配置为非zk的, 不同步)
    binlog_pause
      instance
      status
    lock
      syncMeta.lock(tmp,启动时)
      confInit.lock
      confChange.lock
      binlogStatus.lock
    online
      myid(tmp)
    ddl
      `schema`.`table` (ddl运行时)
        instance
        myid(tmp)
      `schema`.`table2`
    xalog
      node1
      node2
    view
      schema.table
  
```

#### 2.8.7 全局序列

类twitter snowflake 方式, ZK完成的工作是生成每个节点的instanceID。  
类offset-step 方式, ZK完成的工作是存储当前的Step值。

#### 2.8.8 附录

|           |
|-----------|
| 单节点部署执行步骤 |
| 多结点部署额外步骤 |

图例

## 2.9 Grpc告警功能

### 2.9.1 告警功能概述

Dble拥有和商业项目ucore进行告警对接的功能，当dble触发某些重要的报错信息时，会通过ucore提供的grpc接口将对应的告警信息发送到ucore告警中，免去了运维人员在日志文件中的大量搜索，能够直观展示在页面上。

### 2.9.2 告警配置依赖

依赖myid.properties 的告警的基础信息

| 名称            | 内容         | 默认值                                                         | 详细作用原理或应用          | 实例/全局属性 |
|---------------|------------|-------------------------------------------------------------|--------------------|---------|
| url           | grpc告警的url | myid.properties 里的ipAddress                                 | 在发送grpc的时候作为IP地址使用 | 实例      |
| port          | 告警端口       | myid.properties 里的port                                      | grpc发送的目的端口        | 实例      |
| serverId      | 服务器ID      | \$ushard-id(ip1,ip2) ,其中\$ushard-id 是myid.properties 里的myid | 接口参数               | 实例      |
| componentId   | 组件ID       | \$ushard-id 即myid.properties 里myid                          | 接口参数               | 实例      |
| componentType | 组件类型       | ushard                                                      | 接口参数               | 实例      |

## 2.10 meta数据管理

Meta数据的管理包含以下部分：

- [2.10.1 Meta信息初始化](#)
- [2.10.2 Meta信息维护](#)
- [2.10.3 一致性检测](#)
- [2.10.4 View Meta](#)

## 2.10.1 Meta信息初始化

Dble在启动时会逐个同步抓取相关表、视图的信息，并对分区表做一致性检测。

### 2.10.1.1 表信息的初始化

表信息是从后端数据节点抓取的。Dble在启动时会根据schema、table、数据节点的配置（参见1.2 schema.xml）对逐个数据节点执行：

```
show tables;  
show create table ...
```

获取后端表的实际创建信息并做解析以获取需要的信息。

### 2.10.1.2 视图信息的初始化

视图的信息是属于Dble中间件的状态信息，是为数不多的Dble独有的状态之一，存在两种形式的初始化

1 非集群模式状态下通过本地文件进行view的存储和初始化

2 集群状态下通过ZK进行view的存储和初始化

关于其他view meta的实现以及使用细节详见2.10.4 view meta

### 2.10.1.3 分区表的一致性规则

分区表分布在多个实例上的后端表的表结构不必完全相同，但最好完全形同。多个后端表被认为一致的规则为：

- 有相同的列；
- 有相同的主键；
- 有相同的唯一键；
- 有同样的索引；
- 有相同的键。

只要这五个方面相同就认为后端的表具有一致性，一致性检测通过。

## 2.10.2 Meta信息维护

dble在每次执行某些类型ddl语句之后都会更新相关表，视图的元信息。目前，更新元信息的ddl语句类型如下：

- create table语句
- drop table语句
- alter table语句
- truncate table语句
- create index语句
- drop index语句

根据是否配置zookeeper服务器服务(参见1.8 myid.properties)，Meta信息的维护逻辑分如下两种情况：

### 1. 不用zookeeper服务

在此种情况下，由于仅有一个dble运行实例，本地更新已是全部信息，不必做进一步的维护逻辑。

### 2. 利用zookeeper服务

此种情况下的更新逻辑为：

- a. 在启动时每一个dble实例都向zookeeper服务器注册监听事件，监听系统中表元信息的改变。
- b. 当某一个dble实例更新了某些表的元信息之后，它负责向zookeeper服务器广播更新事件。
- c. 其他dble实例在监听到更新事件后更新自己维护的元信息。

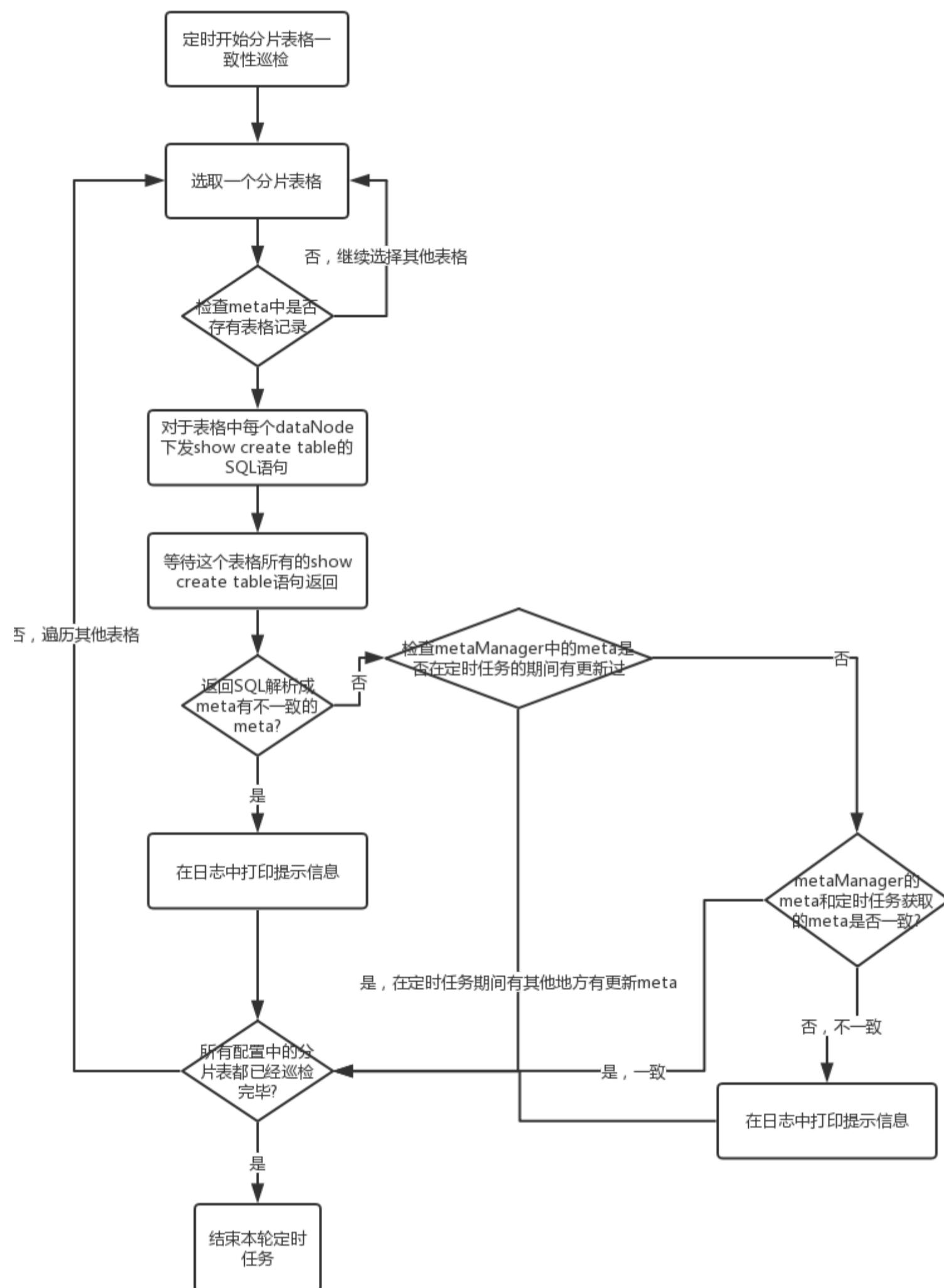
### 2.10.3 一致性检测

由于Dble对于分片的应用和处理可以看到，对于某张具体的分片表，`dble`默认在具体数据库节点上的所有表格的实体拥有同样的结构，由于可能在Dble运行的过程中的种种问题导致了数据库节点的表格结构不一致的情况，在Dble中有对于各个分片表的表格结构进行定期检查确认的机制。

分片表结构一致性检查的定时任务为`tableStructureCheckTask`，其周期为`checkTableConsistencyPeriod`默认值 $30 \times 60 \times 1000$ （30分钟），任务开关为`checkTableConsistency`默认为0（关闭），以上配置均配置在`server.xml`中。大致的表格一致性检查任务逻辑如下：

- 循环配置信息中的`schema`以及`table`
- 检查对应数据库和表格的`meta`信息是否存在
- 对于每个有效的节点下发送SQL “`show create table`”
- 根据收到的建表语句创建新的`meta`，并通过`Set`进行去重
- 判断有几个不同的`meta`，若`meta`数量超过一个则进行在日志中打印告警信息
- 将新老的`meta`进行对比，如果发现新老`meta`之间的`meta`也不一致则在日志中打印告警信息

整体的流程图如下：



#### 全局表的一致性检查

全局表在Dble的使用中被认为每个节点都存在全量的数据，所以在全局表的检查中不光检查了全局表的表格结构，还在需要另外检查表格的数据一致性，功能开启由`server.xml`参数`useGlobbleTableCheck`控制。

## 2.10.4 view meta

### 2.10.4.1 view 种类

- mysql中的view
- dble中的view

### 2.10.4.2 view meta概述

在Dble 2.18.11.0 版本中新增了对view的支持，作为一个中间件支持view采用的方法是和复杂查询类似的逻辑，将VIEW创建的语句中的select部分进行解析，并将解析所得到的解析树进行存储，当有query调用到视图的时候使用解析树进行局部的替代还原成完整的查询SQL，从而达到view实现的效果。

在MySQL中的View meta信息是一种持久化的存储。Dble中的view实现方式也类似，这样Dble中view的实现会给Dble本身带来状态，这个状态需要Dble自身进行保存和维护，所以在view中采取了和XA事务日志相似的方式，并且以视图创建SQL的形式进行以下两种方式的存储：

- 本地文件存储
- ZK存储

在每次发生重启的时候Dble会在初始化meta的时候，通过读取文件或者是ZK中的信息将创建视图的SQL进行重新的解析，最终以解析树的形式保存到meta中去。

在Dble 2.19.10.0版本中新增了对mysql view创建的支持，mysql view会被dble直接下发到后端mysql中执行。不过，创建mysql view有着各种限制，首先view 所属的schema必须采取以下配置，即垂直分片的方式：

```
<schema name="schema2" sqlMaxLimit="100" dataNode="dn5">
</schema>
```

view中涉及的表必须是dataNode dn5 中的表。对于mysql view，dble则没有持久化，但是依然会有相对应的view meta。mysql view 支持通过dble修改，也支持在后端节点手动修改并通过reload方式加载到dble。

### 2.10.4.2 view meta保存

#### 本地文件存储

本地文件的存储中，对应的信息以JSON的格式存放在本地文件中，文件的存储路径以及文件名称通过server.xml中的viewPersistenceConfBaseDir和viewPersistenceConfBaseName两个参数进行配置（默认存储于./viewConf/viewJson中），具体的文件内容举例如下

```
[{
  "schema": "testdb",
  "list": [
    {
      "name": "view_test",
      "sql": "create view view_test as select * from a_test"
    },
    {
      "name": "vt2",
      "sql": "create or replace view vt2 as select * from suntest"
    },
    {
      "name": "sunttest",
      "sql": "create view sunttest as select * from sbtest"
    }
  ]
}]
```

#### ZK KV存储

在集群状态下Dble的各个内部状态需要同步，包括view的创建删除和修改，具体的key值会存储在ZK BASE\_PATH/view下，并使用key值schema\_name:view\_name，在view中使用json的格式进行create sql和修改的serverID的存储，以下给出一个举例：

```
{
  "serverId": "10010",
  "createSql": "create view view_test as select * from a_test"
}
```

此JSON字符串当作value存储在 /.../view/testdb:view\_test

## 2.11 统计管理

统计管理包含以下几个部分

- [2.11.1 查询条件统计](#)
- [2.11.2 表状态统计](#)
- [2.11.3 用户状态统计](#)
- [2.11.4 命令统计](#)
- [2.11.5 heartbeat统计](#)
- [2.11.6 网络读写统计](#)

## 2.11.1 查询条件统计

查询条件统计是在处理查询结果时进行的。统计的是关于某个表中关于某列的条件。

dble系统每次只能进行关于一个表中的一个列的条件统计。

要进行此类统计首先必须开启useSqlStat(参见1.3 server.xml)，其次要设置要统计的表和列，详见下节。

### 2.11.1.1 查询条件统计表列设置

查询条件统计表列的设置命令：

```
reload @@query_cf=table&column;
```

其中**table**为要统计的目标表的表名，**column**为目标表中目标列的列名。

如果要清除查询条件统计表列的设置执行命令：

```
reload @@query_cf;
```

### 2.11.1.2 查看查询条件统计结果

要查看查询条件统计结果执行如下命令：

```
show @@sql.condition;
```

## 2.11.2 表状态统计

表状态统计是在处理查询结果时进行的。要进行此类统计必须开启useSqlStat(参见[1.3 server.xml](#))。

### 2.11.2.1 统计内容

表状态主要统计如下内容:

- 读操作的次数
- 写操作的次数
- 表关系，即在同一个语句中出现多个表，以第一个表作为操作的主表，其他表是和主表有关系的表。
- 对主表最后一次操作的时刻。

### 2.11.2.2 统计结果查看

表状态的统计结果可以通过如下命令查看:

```
show @@sql.sum.table;
```

如果要重置表状态统计，执行如下命令:

```
show @@sql.sum.table true;
```

### 2.11.3 用户状态统计

用户状态统计是在处理查询结果时进行的。要进行此类统计必须开启useSqlStat(参见[1.3 server.xml](#))。

#### 2.11.3.1 统计内容

用户状态统计如下内容:

1. 每个用户的最大并发
2. 每个用户的前N条最慢SQL的sql语句, 语句执行的开始时间。N通过sqlRecordCount进行设置(参见[1.3 server.xml](#)), 默认值为10。慢SQL的标准为执行时间大于等于T的sql语句。  
通过命令:  
`reload @@sqlslow=t;`  
进行设置。t为整数, 单位为毫秒。
3. 网络读写字节数
4. 每用户最后执行的50条语句
5. 每用户结果集大于10000行的前10条select语句
6. 每用户频度最大的前1024条sql
7. 每用户结果集大小 (单位: 字节) 大于M的sql语句。M通过maxResultSet进行设置(参见[1.3 server.xml](#))。

#### 2.11.3.2 用户状态统计服务的命令

- `show @@sql;`
- `show @@sql.high;`
- `show @@sql.large;`
- `show @@sql.resultset;`
- `show @@sql.slow;`
- `show @@sql.sum.user;`

以上命令的使用请参看[2.1 管理端命令集](#)。

#### 2.11.3.3 用户状态统计重置

要清空用户状态统计结果而重新进行统计, 执行如下命令:

```
reload @@user_stat;
```

#### 2.11.4 命令统计

命令统计是分类统计dble执行的命令计数，在执行各个命令类的命令时进行统计。统计的命令类如下：

1. initDB
2. query
3. stmtPrepare
4. stmtSendLongData
5. stmtReset
6. stmtExecute
7. stmtClose
8. ping
9. kill
10. quit
11. heartbeat
12. other，除以上命令类命令之外的所有命令。

##### 2.11.4.1 命令统计结果查看

命令统计结果查看执行如下命令：

- show @@command;
- show @@command.count; 具体命令的使用请参看[2.1 管理端命令集](#)

## 2.11.5 heartbeat统计

heartbeat统计后端mysql实例的heartbeat状态信息，在对后端mysql实例进行heartbeat检测时进行统计。

### 2.11.5.1 统计内容

heartbeat统计每一次从heartbeat查询发送到查询结果接受之间的时间差，同步状态。

### 2.11.5.2 统计结果查看

heartbeat统计的查看执行如下命令：

- show @@heartbeat;
- show @@heartbeat.detail;
- show @@datasource.synstatus;
- show @@datasource.syndetail where name=xxx;其中，xxx为datasource名字。
- show @@datasource.cluster;

以上命令的使用请参看参看[2.1 管理端命令集](#)

## 2.11.6 网络读写统计

dble的每一个前端在进行网络通信时对读写的数据量进行统计。

### 2.11.6.1 统计内容

- 网络读字节数
- 网络写字节数
- 最后一次进行读/写的时刻

### 2.11.6.2 服务的命令

- show @@connection;
- show @@backend;
- show @@connection.sql;

以上命令的使用请参看[2.1 管理端命令集](#)

## 2.12 故障切换

版本2.20.04.0起对切换功能进行了重定义,废除原有的切换功能,请知悉。

新的MySQL高可用切换分为两类:

一个是单实例部署的dble会内置一个自带的高可用切换的python3脚本,跟随dble启动和停止,需要设置server.xml中system的useOuterHa 参数为false。

另一个是支持第三方的高可用切换接口功能,支持集群部署。当然,单实例dble也可以采用这种方式,需要设置system的useOuterHa 参数为true。

这种方式情况下,如果没有真的配置第三方高可用切换组件,则什么也不会发生。具体请参考[高可用切换接口](#)

### 2.12.1 前提条件

- 单实例部署的dble, 设置server.xml中system的useOuterHa 参数为false。
- 安装好相应的python3环境, 详情参考本章“自定义python脚本”部分

### 2.12.2 工作方式

dble启动时跟随dble启动python3脚本进程,脚本会读取schema.xml内部的配置,检查后端数据库结点的状态,发生异常后会执行切换命令来通知dble

dble提供以下运维命令:

```
show @@custom_mysql_ha
```

用于查看或者监控当前进程是否存活

当python脚本意外终止时, 可以手工运行命令来启动:

```
enable @@custom_mysql_ha
```

当然, 也可以手动关闭该脚本

```
disable @@custom_mysql_ha
```

### 2.12.2 注意事项

#### 2.12.2.1 此功能目前只支持在linux环境下使用。

#### 2.12.2.2 reload注意事项

做reload @@config 之前需要修改配置文件, 这时候如果不终止python脚本, 可能会读到中间状态的文件, 所以我们建议的标准流程如下:

1. disable @@custom\_mysql\_ha 关闭切换功能
2. 修改配置文件
3. reload @@config 重新加载配置
4. enable @@custom\_mysql\_ha 开启切换功能

#### 2.12.2.3 自定义python脚本

custom\_mysql\_ha.py 脚本在dble安装目录的bin目录下, 使用python3编写, 您可以根据实际情况自行修改。

为了使它能正常工作, 需要做以下的准备工作:

1、安装Python3, 确认python3是否已安装可通过如下命令

```
/usr/local/bin/python3 --version  
/usr/local/bin/pip3 --version
```

2、安装mysqlclient及依赖

CentOS 依赖

```
yum install mysql-devel
```

or Ubuntu 依赖

```
apt-get install libmysqlclient-dev
```

然后

```
pip3 install mysqlclient
```

3、six

```
pip3 install six
```

or

```
pip3 install six -i http://pypi.douban.com/simple --trusted-host pypi.douban.com
```

4、coloredlogs

```
pip3 install coloredlogs
```

or

```
pip3 install coloredlogs -i http://pypi.douban.com/simple --trusted-host pypi.douban.com
```

5、rsa

```
pip3 install rsa
```

or

```
pip3 install rsa -i http://pypi.douban.com/simple --trusted-host pypi.douban.com
```

## 2.13 前后端连接检查

### 2.13.1 前端连接空闲检查

- 根据在server.xml中配置的processorCheckPeriod进行定时的前端连接检查
- 根据配置在server.xml中的idleTimeout判断前端连接是否存在空闲时间超限的现象
- 如果发现空闲时间超过限度则关闭连接

### 2.13.2 后端连接SQL超时检查

- 根据在server.xml中配置的processorCheckPeriod进行定时的前后端连接检查
- 根据配置sqlExecuteTimeout检查所有正在执行的后端连接，是否有执行时间超限的情况
- 关闭所有执行时间超过限度的非DDL后端连接

## 2.14 ER 拆分

### 2.14.1 普通ER拆分

当我们需要两个表要进行join的时候,由于数据被分配到不同的节点上,普通的nest loop方式可能效率会很低。如果我们能把需要join的表按照统一规则划分到相同的区上,就能大概率的解决这一问题。举个简单的例子如下:

销售单			销售详情单	
流水号	顾客ID	日期	流水号	商品ID
201712010001	1	20171201	201712010001	1
201712010002	4	20171201	201712010001	2
201712020001	15	20171202	201712010002	7
201712020002	4	20171202	201712010002	31

这样两张表具有外键关系,我们可以仿照聚簇的方式按照对应列来进行拆分,这样就可以不跨库实现join了。

如下图:

分片node1			销售详情单	
销售单(按日期拆分表1)			流水号	商品ID
流水号	顾客ID	日期	201712010001	1
201712010001	1	20171201	201712010001	2
201712010002	4	20171201	201712010002	7
分片node2			销售详情单	
销售单(按日期拆分表2)			流水号	商品ID
流水号	顾客ID	日期	201712020001	103
201712020001	15	20171202	201712020001	2304
201712020002	4	20171202	201712020002	27
			201712020002	91

要实现这样的ER功能,需要如下配置。

```
<table name="sales" dataNode="dn1,dn2" rule="sharding">
<childTable name="sales_detail" joinKey="sales_detail_pos_num" parentKey="sales_pos_num"/>
</table>
```

### 2.14.2 智能的ER关系

当我们有多个不同的表时,上面的配置方式有点难以使用了。

这种情况下,如果2张或者多张表在mycat上的分片规则相同并且具体分片也相同,即使没有配置ER关系,也会当作ER关系来处理。举例如下配置(片段):

```
<!--schema片段-->
<table name="tableA" dataNode="dn1,dn2" rule="ruleA" />
<table name="tableB" dataNode="dn1,dn2" rule="ruleB" />
<table name="tableC" dataNode="dn2,dn1" rule="ruleC" />
<table name="tableD" dataNode="dn3,dn4" rule="ruleA" />
<table name="tableE" dataNode="dn1,dn2" rule="ruleA" />
<table name="tableF" dataNode="dn1,dn2" rule="ruleF" />

<!--rule片段-->
<tableRule name="ruleA">
<rule>
  <columns>id_a</columns>
  <algorithm>hash_function</algorithm>
</rule>
</tableRule>
<tableRule name="ruleB">
<rule>
  <columns>id_b</columns>
  <algorithm>hash_function</algorithm>
</rule>
</tableRule>
<tableRule name="ruleC">
<rule>
  <columns>id_c</columns>
  <algorithm>hash_function</algorithm>
</rule>
</tableRule>
<tableRule name="ruleF">
<rule>
  <columns>id_a</columns>
  <algorithm>enum_par</algorithm>
</rule>
</tableRule>
<function name="enum_par">
  class="io.mycat.route.function.PartitionByFileMap"
  <property name="mapFile">partition-hash-int.txt</property>
</function>
<function name="hash_function" class="io.mycat.route.function.PartitionByLong">
  <property name="partitionCount">2</property>
  <property name="partitionLength">512</property>
</function>
```

最终会得出这样一个映射关系,识别分组会根据数据分布结点和function的唯一性将表分为几组,同一组的才会有ER关系。

table名	拆分列	数据分布结点	function	识别分组
tableA	id_a	dn1,dn2	hash_function	1
tableB	id_b	dn1,dn2	hash_function	1
tableC	id_c	dn2,dn1	hash_function	2
tableD	id_a	dn3,dn4	hash_function	3
tableE	id_a	dn1,dn2	hash_function	1
tableF	id_a	dn1,dn2	enum_par	4

即,ER关系集合为:

PS: 此处略去了schema,实际实现需要标识schema防止重复,ER关系是到列的,如果关联关系不是上述表对应的列,也不会视为ER.

## 2.15 global 表

在一些业务系统中，存在着类似字典表的表格。它们与业务表之间可能有关系，这种关系，可以理解为“标签”，而不应理解为通常的“主从关系”。这些表具有以下几个特性：

- 变动不频繁
- 数据量总体变化不大
- 数据规模不大，很少有超过数十万条记录。

鉴于此，`dble` 定义了一种特殊的表，称之为“全局表”，全局表具有以下特性：

- 全局表的插入、更新操作会实时在所有节点上执行，保持各个分片的数据一致性
- 全局表的查询操作，只从一个节点获取
- 全局表可以跟任何一个表进行JOIN 操作

将字典表或者符合字典表特性的一些表定义为全局表，某种程度上部分解决了数据JOIN的难题。

举例如下：



对于数据量不大的字典表（例：超市商品），在多个分片上都有一份同样的副本

相关JOIN语句可以直接下发给各个结点，直接合并结果集就行。

JOIN 例子(伪SQL):

```
SELECT 日期,商品名,COUNT(*) AS 订单量
FROM 商品表
JOIN 销售详单 USING(商品ID)
WHERE 日期范围(跨结点)
GROUP BY 日期,商品名.
```

## 2.16 cache 的使用

cache 的配置请参见1.6节的内容。

### 2.16.1 主键缓存

主键缓存里key, value键值对是[schema.table.主键值, 路由节点]

适合主键路由缓存的场景:

- 1.条件都包含主键
- 2.拆分里和主键列不同, 先通过拆分列查出数据 (主键结果进缓存), 再通过主键列查询。

另外, 为了使缓存生效, 可能影响查询延迟。

**方式一:开启默认缓存和具体表的缓存。**

KEY是layerpool.TableID2DataNodeCache,

eg. layerpool.TableID2DataNodeCache=encache,10000,18000

如果不再添加新的配置, 那么schema.xml 所有配置了主键的table将会共用此配置

当然, 也可以指定对某个表格进行缓存的配置 格式如下: layerpool.TableID2DataNodeCache.`testdb`\_`testtable`= 容量, 超时时间

**方式二:KEY是layerpool.TableID2DataNodeCacheType, 仅开启具体表的缓存。**

eg. layerpool.TableID2DataNodeCacheType=encache

然后再对具体某个表格进行缓存配置:

layerpool.TableID2DataNodeCache.`testdb`\_`testtable`= 容量, 超时时间

注意: 方式一和方式二不能同时开启。

### 2.16.2 路由缓存

路由缓存里键值对是[sql, 路由节点]

配置的值里KEY是pool.SQLRouteCache

VALUE 是逗号隔开的三个值, 分别为cachefactory的name,容量, 超时时间。

### 2.16.3 ER子表计算缓存

路由缓存里键值对是[子表对应的joinkey, 路由节点]

KEY是pool.ER\_SQL2PARENTID

VALUE是逗号隔开的三个值, 分别为cachefactory的name,容量, 超时时间

## 2.17 执行计划

### 2.17.1 执行计划的意义

对执行计划进行分析，可以了解中间件和节点是否对SQL语句生成了最优的执行计划，是否有优化的空间，从而为SQL优化提供重要的参考信息。

### 2.17.2 执行计划的分类

dble的执行计划分为两个层次：dble层的执行计划与节点层的执行计划。

dble层的执行计划：在SQL语句执行前，dble会根据SQL语句的基本信息，判断该SQL语句应该在哪些节点上执行，将SQL改写成在节点上执行的具体形式，并决定采用何种策略进行数据合并与计算等。

节点层的执行计划：就是原生的MySQL执行计划。

### 2.17.3 dble层的执行计划

dble用EXPLAIN指令来查看dble层的执行计划。如：`explain select id,accountno from account where userid=2;`

EXPLAIN指令的执行结果包括语句下发的节点，实际下发的SQL语句和数据的合并操作的信息。这些信息是系统静态分析产生的，并没有真正的执行语句。

另外，复杂查询的查询计划也会有所反映，可以通过计划来优化查询语句

如

```
mysql> explain select * from sharding_two_node a inner join sharding_four_node b on a.id =b.id;
+-----+-----+
| DATA_NODE | TYPE      | SQL/REF
+-----+-----+
| dn1.0    | BASE SQL | select `a`.`id`, `a`.`c_char`, `a`.`ts`, `a`.`si` from `sharding_two_node` `a` ORDER BY `a`.`id` ASC |
| dn2.0    | BASE SQL | select `a`.`id`, `a`.`c_char`, `a`.`ts`, `a`.`si` from `sharding_two_node` `a` ORDER BY `a`.`id` ASC |
| dn1.1    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| dn2.1    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| dn3.0    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| dn4.0    | BASE SQL | select `b`.`id`, `b`.`c_flag`, `b`.`c_decimal` from `sharding_four_node` `b` ORDER BY `b`.`id` ASC |
| merge.1  | MERGE     | dn1.0, dn2.0
| merge.2  | MERGE     | dn1.1, dn2.1, dn3.0, dn4.0
| join.1   | JOIN      | merge.1, merge.2
+-----+-----+
9 rows in set (0.00 sec)
```

再举例如：

```
mysql> explain select id from single union all select b.si from sharding_four_node a inner join sharding_two_node b on a.id =b.id;
+-----+-----+
| DATA_NODE | TYPE      | SQL/REF
+-----+-----+
| dn1.0    | BASE SQL | select 'single'.`id` from 'single'
| dn1.1    | BASE SQL | select `a`.`id` from `sharding_four_node` `a` ORDER BY `a`.`id` ASC
| dn2.0    | BASE SQL | select `a`.`id` from `sharding_four_node` `a` ORDER BY `a`.`id` ASC
| dn3.0    | BASE SQL | select `a`.`id` from `sharding_four_node` `a` ORDER BY `a`.`id` ASC
| dn4.0    | BASE SQL | select `a`.`id` from `sharding_four_node` `a` ORDER BY `a`.`id` ASC
| dn1.2    | BASE SQL | select `b`.`si`, `b`.`id` from `sharding_two_node` `b` ORDER BY `b`.`id` ASC
| dn2.1    | BASE SQL | select `b`.`si`, `b`.`id` from `sharding_two_node` `b` ORDER BY `b`.`id` ASC
| merge.2  | MERGE     | dn1.1, dn2.0, dn3.0, dn4.0
| merge.3  | MERGE     | dn1.2, dn2.1
| join.1   | JOIN      | merge.2, merge.3
| merge.1  | MERGE     | dn1.0
| union_all.1 | UNION_ALL | join.1, merge.1
+-----+-----+
12 rows in set (0.01 sec)
```

### 2.17.4 节点层的执行计划

通过EXPLAIN2命令可查看指定节点上的执行计划。如：

explain2 datanode=dn1 sql=select id,accountno from account where userid=2;

explain2会将sql语句加上explain下发到指定的datanode执行，并把节点上explain的结果返回调用者。

## 2.18 性能观测以及调试概览

- Btrace脚本性能观察(观察查询过程中每个不同阶段的耗时)
- manager命令show @@thread\_used观察(观察不同线程的负载情况)

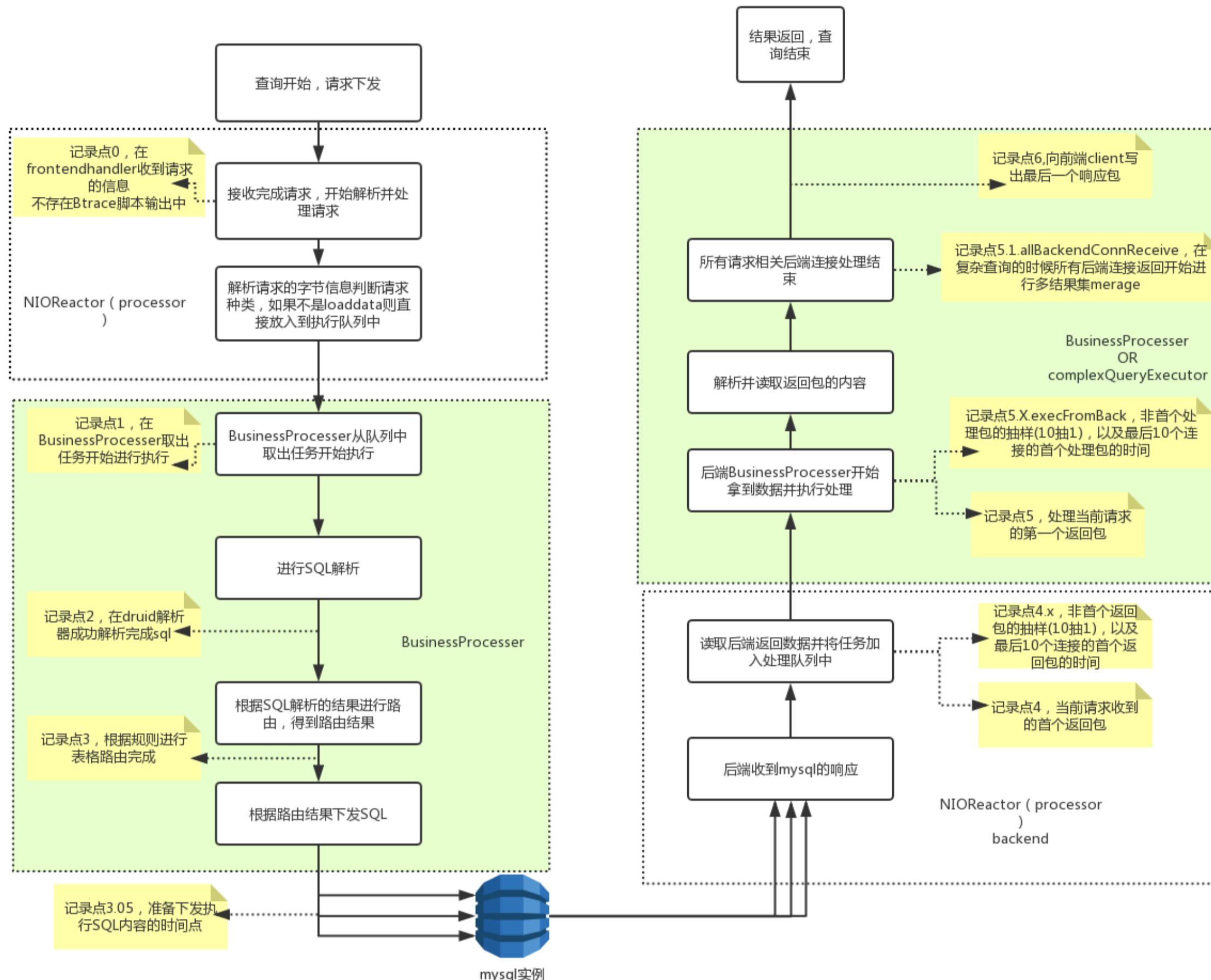
### 2.18.1 Btrace脚本观察

#### 2.18.1.1 观察

Dble源码中自带观测脚本,BTraceCostTime.java文件是专用的Btrace观测脚本,当前脚本适用Btrace v1.3

Btrace 相关资料<https://github.com/btracelio/btrace>

开启性能观察的统计需要配置server.xml中的两个参数useCostTimeStat和costSamplePercent, 分别是启用耗时监控的标志位以及耗时监控的采样百分比, useCostTimeStat = 1的状态下耗时监控被开启, 并且以costSamplePercent的概率进行统计, 默认的采样率是1%, 采样率过低有可能导致btrace脚本输出没有统计结果的情况, 采样率过高则会影响性能本身。脚本中包含几个统计点如下图所示:



通过执行btrace监控正在运行的dble, 可以获得类似以下的结果图:

profiling:	Block	Invocations	SelfTime.Total	SelfTime.Avg	SelfTime.Min	SelfTime.Max	WallTime.Total	WallTime.Avg	WallTime.Min	WallTime.Max
	request->1.startProcess	9073	-638142734	-70334	-1051058	493260	202952071	22368	10565	493260
	request->2.endParse	9073	234134936	25805	13206	523393	437087007	48174	23771	1016653
	request->3.endRoute	9073	404389553	44570	20123	121474	841476560	92745	43894	1075553
	request->4.resFromBack	9073	592398	65	-649019	1602901	4805691843	529669	261612	1602901
	request->5.startExecuteBackend	9073	-56808823	-6261	-1749483	2020297	5047581273	556329	350530	2020297
	request->6.response	9073	59150286	6519	3045	366620	5107315454	562913	353575	2386917

以上输出值有效信息为Block,Invocations,WallTime., 而SelfTime,由于btrace的原因 不准确, 可以忽略。

其中通过记录一个查询在每个记录点的时间点来观察是否存在某个中间步骤耗时过长, 并进行针对性的优化和调整

### 2.18.1.2 节点描述

- 0 : 逻辑时间0, 来自客户端的请求首次被dble接收到的时间点
- 1.startProcess : 开始处理前端请求的时间节点
- 2.endParse : SQL被解析完成的时间节点
- 3.endRoute : SQL被路由完成的时间节点
- 3.05 readyToDeliver : SQL准备开始下发的时间点
- 4.resFromBack : 前端请求首次由任意个后端连接返回信息
- 4.X.resFromBack : 前端请求后端连接的首个返回包的采样时间点, 例: 4.3.resFromBack指的是第三个后端连接首次返回给dble的时间点, 具体的后端连接的编号为首次返回的顺序, 即首个有网络包返回的后端连接编号为1
- 5.startExecuteBackend : 前端请求后端连接首个包进入处理阶段的时间点
- 5.X.startExecuteBackend : 前端请求后端连接的首个被处理的采样时间点, 例: 5.3.startExecuteBackend指的是第三个后端连接首次进入到后端数据处理的阶段, 具体的编号为首次进入处理阶段的顺序
- 5.1.allBackendConnReceive : 在复杂查询的情况下开始merge的时间点
- 6.response : 给前端开始最终返回包的时间点

注意:Btrace关于性能跟踪的脚本设计专为单一SQL模式下使用, 在多个SQL混合查询的情况下, 由于每个SQL涉及的采样点可能不同, 会出现数据上的异常, 甚至是节点顺序和时间点倒挂的现象

### 2.18.1.3 调整策略

- 1-0的耗时增长: 需增大server.xml中processorExecutor的值 (其中0表示dble开始读到某条语句时的时间, 即初始值: 0, 所以1-0的时间就是request->1.startProcess的输出值, 不需要进行计算)
- 4-3的耗时增长: 需增大server.xml中backendProcessors的值
- 5-4的耗时增长: 需增大server.xml中backendBusinessExecutor的值

### 2.18.2 Manager命令观察

#### 2.18.2.1 观察

Dble 在18.02.0版本中新添加了manager端口的性能观测命令，可以通过命令查看各个线程的负载情况，需要配合在配置文件server.xml中的新增参数useThreadUsageStat进行使用  
使用命令show @@thread\_used会返回各个dble中关键线程最近时间的负载情况，如下示例

THREAD_NAME	LAST_QUARTER_MIN	LAST_MINUTE	LAST_FIVE_MINUTE
backendBusinessExecutor2	0%	0%	0%
backendBusinessExecutor1	0%	0%	0%
backendBusinessExecutor0	0%	0%	0%
BusinessExecutor3	0%	0%	0%
\$_NIO_REACTOR_BACKEND-2	0%	0%	0%
BusinessExecutor1	0%	0%	0%
\$_NIO_REACTOR_BACKEND-3	0%	0%	0%
BusinessExecutor2	12%	3%	3%
\$_NIO_REACTOR_BACKEND-0	0%	0%	0%
\$_NIO_REACTOR_FRONT-0	0%	0%	0%
\$_NIO_REACTOR_BACKEND-1	0%	0%	0%
BusinessExecutor0	0%	0%	0%

12 rows in set (0.00 sec)

- BusinessExecutorX  
面向前端的业务处理线程，主要处理前端请求的解析，sql解析路由，下发查询到mysql实例等
- backendBusinessExecutorX  
面向后端的业务处理线程，主要处理后端mysql查询结果的返回解析，结果聚合，并发回结果到client
- \$\_NIO\_REACTOR\_FRONT\_X  
前端请求接受线程，负责请求的接收和读取，之后把数据处理交给BusinessExecutor进行
- \$\_NIO\_REACTOR\_BACKEND\_X  
后端请求接收线程，负责mysql返回信息的接收和读取，之后把数据交给backendBusinessExecutor处理

### 2.18.2.2 调整策略

当输出的统计结果中，有一个或者多个类型的线程使用率过高（经验值为超过80%），可以适当调整对应的处理线程的数量

- processors 参数控制前端\_NIO\_REACTOR\_FRONT\_X的数量
- backendProcessors 参数控制后端\_NIO\_REACTOR\_BACKEND\_X的数量
- backendProcessorExecutor 参数控制backendBusinessExecutorX数量
- processorExecutor 参数控制BusinessExecutorX数量

## 2.19 智能计算reload

我们对reload @@config\_all 做了重构，增加了对datahost的变化计算，使得reload行为对整个系统的影响变到最小。

变更的datahost，等效于删除旧的datahost，新增新的datahost

### 2.19.1 默认reload @@config\_all

连接池行为描述：

#### 2.19.1.1 不变的datahost

此种情况下,datahost连接不发生变化，如果关联的schmea发生变更,在需要使用时候进行连接新建或者偷取并同步上下文的方式进行更新。

#### 2.19.1.2 新增的datahost

建立新的连接池供使用

#### 2.19.1.3 删除的datahost

遍历当前连接池，如果没有事务正在使用连接，则回收，否则放回后端待回收连接池（在show backend中任可以看到放入回收池的时间），等事务结束时候连接被关闭

### 2.19.2 reload @@config\_all -f

强制回收所有正在使用的链接

连接池行为描述：

#### 2.19.2.1 不变的datahost

此种情况下,datahost连接不发生变化，正在使用的连接会被回收

#### 2.19.2.2 新增的datahost

建立新的连接池供使用

#### 2.19.2.3 删除的datahost

遍历如果当前连接池，如果没有事务正在使用连接，则回收，否则关闭对应的前端连接以及相关的后端连接

### 2.19.3 reload @@config\_all -r

不计算datahost的变化，相当于原本所有的连接池会被删除，然后新建所有的连接池

遍历旧连接池，如果没有事务正在使用连接，则回收，否则放回后端待回收连接池（在show backend中任可以看到放入回收池的时间），等事务结束时候连接被关闭

### 2.19.4 reload @@config\_all -s

跳过datahost连接检查（与节点建立连接，连接成功，检测通过）

## 2.20 慢查询日志

类似于MySQL的慢查询日志，可以全局开启并设置，记录db1e server运行过程当中慢查询日志，日志格式兼容MySQL慢查询分析工具(已测试过MySQL官方工具mysqldumpslow和Percona的pt-query-digest)。此外，开启了慢查询日志工具之后，也可以查询某个连接的当前SQL的执行状态，相关命令为：show @@connection.sql.status where FRONT\_ID=?;

### 2.20.1 在server.xml里增加了6个参数，用于启动时候控制慢查询日志的行为

```
<!-- 是否开启慢查询日志 -->
<property name="enableSlowLog">0</property>
<!-- 慢查询日志保存文件目录 -->
<!--<property name="slowLogBaseDir">./slowlogs</property>-->
<!-- 慢查询日志保存文件前缀名称 -->
<!--<property name="slowLogBaseName">slow_query</property>-->
<!-- 日志两次刷盘之间的最大周期，单位是秒 -->
<property name="flushSlowLogPeriod">1</property>
<!-- 日志两次刷盘之间内存中的最大条数阈值 -->
<property name="flushSlowLogSize">1000</property>
<!-- 慢查询统计阈值，大于此值会被认为是慢查询，单位是毫秒 -->
<property name="sqlSlowTime">100</property>
```

### 2.20.2 管理端口增加命令，用于运行过程中动态修改慢查询日志统计行为

```
enable @@slow_query_log; -- 开启慢查询日志
show @@slow_query_log; -- 查询慢查询日志的开启状态
disable @@slow_query_log; -- 关闭慢查询日志
show @@slow_query_log; -- 再次查询慢查询日志的开启状态

show @@slow_query.time; -- 查看慢查询日志统计阈值
reload @@slow_query.time=200; -- 修改慢查询日志统计阈值

show @@slow_query.flushperiod; -- 查看慢查询日志刷盘周期
reload @@slow_query.flushperiod=2; -- 修改慢查询日志刷盘周期

show @@slow_query.flushsize; -- 查看慢查询日志刷盘条数阈值
reload @@slow_query.flushsize=1100; -- 修改慢查询日志刷盘条数阈值
```

### 2.20.3 支持慢查询日志分析工具：MySQL的 mysqldumpslow 工具和Percona 的 pt-query-digest 工具

慢查询日志大概是这样的：

```

/FAKE_PATH/mysql, Version: FAKE_VERSION. started with:
Tcp port: 3320 Unix socket: FAKE_SOCK
Time           Id Command  Argument
# Time: 2018-08-23T17:40:10.149000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.132709  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000350  Prepare_Push: 0.116678  dn1_First_Result_Fetch: 0.013686  dn1_Last_Result_Fetch: 0.001422  Write_Client: 0.001995
SET timestamp=1535017210149;
select * from sharding_two_node where id =1;
# Time: 2018-08-23T17:40:10.200000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.035600  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000062  Prepare_Push: 0.006733  dn2_First_Result_Fetch: 0.012524  dn1_First_Result_Fetch: 0.010971  dn2_Last_Result_Fetch: 0.015368  dn1_I
SET timestamp=1535017210200;
select * from sharding_two_node;
# Time: 2018-08-23T17:40:10.282000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.045337  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000166  Prepare_Push: 0.003941  dn1_First_Result_Fetch: 0.039652  dn1_Last_Result_Fetch: 0.000300  Write_Client: 0.001578
SET timestamp=1535017210282;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:10.315000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.031232  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.005467  Prepare_Push: 0.001989  dn2_First_Result_Fetch: 0.020240  dn2_Last_Result_Fetch: 0.001900  Write_Client: 0.003536
SET timestamp=1535017210315;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:10.432000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.116672  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.013625  Prepare_Push: 0.024767  dn2_First_Result_Fetch: 0.056395  dn1_First_Result_Fetch: 0.026420  dn2_Last_Result_Fetch: 0.000743  dn1_I
SET timestamp=1535017210432;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:10.772000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.338569  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000082  Prepare_Push: 0.258365  dn1_0_First_Result_Fetch: 0.047494  dn1_0_Last_Result_Fetch: 0.029018  dn2_0_First_Result_Fetch: 0.042964
SET timestamp=1535017210772;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:10.821000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.046745  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000059  Prepare_Push: 0.025401  dn1_0_First_Result_Fetch: 0.011755  dn1_0_Last_Result_Fetch: 0.001180  Generate_New_Query: 0.001706  dn1_I
SET timestamp=1535017210821;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:12.061000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.036952  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.001111  Prepare_Push: 0.001132  dn1_First_Result_Fetch: 0.034266  dn1_Last_Result_Fetch: 0.000084  Write_Client: 0.000443
SET timestamp=1535017212061;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:12.091000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.028213  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000666  Prepare_Push: 0.001206  dn2_First_Result_Fetch: 0.025991  dn2_Last_Result_Fetch: 0.000101  Write_Client: 0.000349
SET timestamp=1535017212091;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:12.132000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.040363  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000223  Prepare_Push: 0.001172  dn2_First_Result_Fetch: 0.019852  dn1_First_Result_Fetch: 0.019810  dn2_Last_Result_Fetch: 0.000901  dn1_I
SET timestamp=1535017212132;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:12.145000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.012196  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000115  Prepare_Push: 0.001403  dn1_0_First_Result_Fetch: 0.006714  dn1_0_Last_Result_Fetch: 0.002561  dn2_0_First_Result_Fetch: 0.006787
SET timestamp=1535017212145;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:12.164000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.016979  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000224  Prepare_Push: 0.002236  dn1_0_First_Result_Fetch: 0.006678  dn1_0_Last_Result_Fetch: 0.000703  Generate_New_Query: 0.000866  dn1_I
SET timestamp=1535017212164;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:13.134000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.010213  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000145  Prepare_Push: 0.001520  dn1_First_Result_Fetch: 0.007996  dn1_Last_Result_Fetch: 0.000201  Write_Client: 0.000551
SET timestamp=1535017213134;
select * from sharding_two_node where id =1;
# Time: 2018-08-23T17:40:13.153000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.014257  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000080  Prepare_Push: 0.002394  dn2_First_Result_Fetch: 0.008839  dn1_First_Result_Fetch: 0.008837  dn2_Last_Result_Fetch: 0.001424  dn1_I
SET timestamp=1535017213153;
select * from sharding_two_node;
# Time: 2018-08-23T17:40:13.212000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.029822  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000063  Prepare_Push: 0.001128  dn1_First_Result_Fetch: 0.028277  dn1_Last_Result_Fetch: 0.000109  Write_Client: 0.000355
SET timestamp=1535017213212;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:13.321000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.076093  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000231  Prepare_Push: 0.001334  dn2_First_Result_Fetch: 0.035072  dn1_First_Result_Fetch: 0.035074  dn2_Last_Result_Fetch: 0.018756  dn1_I
SET timestamp=1535017213321;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:13.348000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.026278  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000335  Prepare_Push: 0.001249  dn1_0_First_Result_Fetch: 0.011028  dn1_0_Last_Result_Fetch: 0.009279  dn2_0_First_Result_Fetch: 0.019200
SET timestamp=1535017213348;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:13.381000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.029912  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000662  Prepare_Push: 0.003189  dn1_0_First_Result_Fetch: 0.014453  dn1_0_Last_Result_Fetch: 0.001013  Generate_New_Query: 0.000911  dn1_I
SET timestamp=1535017213381;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:14.163000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.012540  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000197  Prepare_Push: 0.001303  dn2_First_Result_Fetch: 0.006452  dn1_First_Result_Fetch: 0.007858  dn2_Last_Result_Fetch: 0.004065  dn1_I
SET timestamp=1535017214163;
select * from sharding_two_node;
# Time: 2018-08-23T17:40:14.220000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.027587  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000051  Prepare_Push: 0.000744  dn1_First_Result_Fetch: 0.026441  dn1_Last_Result_Fetch: 0.000104  Write_Client: 0.000350
SET timestamp=1535017214220;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:14.253000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.031984  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000070  Prepare_Push: 0.001144  dn2_First_Result_Fetch: 0.030202  dn2_Last_Result_Fetch: 0.000182  Write_Client: 0.000568
SET timestamp=1535017214253;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:14.292000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.037327  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000272  Prepare_Push: 0.001316  dn2_First_Result_Fetch: 0.014299  dn1_First_Result_Fetch: 0.014331  dn2_Last_Result_Fetch: 0.001148  dn1_I
SET timestamp=1535017214292;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:14.303000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.010244  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000050  Prepare_Push: 0.001101  dn1_0_First_Result_Fetch: 0.004540  dn1_0_Last_Result_Fetch: 0.002781  dn2_0_First_Result_Fetch: 0.004708
SET timestamp=1535017214303;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:14.327000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.021078  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000109  Prepare_Push: 0.002098  dn1_0_First_Result_Fetch: 0.006720  dn1_0_Last_Result_Fetch: 0.000748  Generate_New_Query: 0.001158  dn1_I
SET timestamp=1535017214327;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:15.254000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.010569  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000076  Prepare_Push: 0.001050  dn1_First_Result_Fetch: 0.008330  dn1_Last_Result_Fetch: 0.000146  Write_Client: 0.001113
SET timestamp=1535017215254;
select * from sharding_two_node where id =1;
# Time: 2018-08-23T17:40:15.321000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.024216  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000081  Prepare_Push: 0.001295  dn1_First_Result_Fetch: 0.021938  dn1_Last_Result_Fetch: 0.000422  Write_Client: 0.000902
SET timestamp=1535017215321;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:15.351000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2
# Query_time: 0.027790  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000666  Prepare_Push: 0.000760  dn2_First_Result_Fetch: 0.025984  dn2_Last_Result_Fetch: 0.000094  Write_Client: 0.000386
SET timestamp=1535017215351;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:15.392000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1]  Id:  2

```

```

# Query_time: 0.039805 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000217 Prepare_Push: 0.000804 dn2_First_Result_Fetch: 0.017410 dn1_First_Result_Fetch: 0.017468 dn2_Last_Result_Fetch: 0.001490 dn1_I
SET timestamp=1535017215392;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:15.410000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.017384 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000152 Prepare_Push: 0.001183 dn1_0_First_Result_Fetch: 0.005037 dn1_0_Last_Result_Fetch: 0.007164 dn2_0_First_Result_Fetch: 0.008156
SET timestamp=1535017215410;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:15.434000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.021341 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000318 Prepare_Push: 0.002764 dn1_0_First_Result_Fetch: 0.010897 dn1_0_Last_Result_Fetch: 0.000544 Generate_New_Query: 0.000798 dn1_I
SET timestamp=1535017215434;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:16.322000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.030106 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000217 Prepare_Push: 0.001253 dn1_First_Result_Fetch: 0.028330 dn1_Last_Result_Fetch: 0.000086 Write_Client: 0.000306
SET timestamp=1535017216322;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:16.353000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.030005 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.001306 Prepare_Push: 0.001004 dn2_First_Result_Fetch: 0.027242 dn2_Last_Result_Fetch: 0.000140 Write_Client: 0.000453
SET timestamp=1535017216403;
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.049618 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.001459 Prepare_Push: 0.000830 dn2_First_Result_Fetch: 0.024286 dn1_First_Result_Fetch: 0.025469 dn2_Last_Result_Fetch: 0.001726 dn1_I
SET timestamp=1535017216403;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:16.526000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.121702 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000337 Prepare_Push: 0.000889 dn1_0_First_Result_Fetch: 0.009370 dn1_0_Last_Result_Fetch: 0.002010 dn2_0_First_Result_Fetch: 0.009160
SET timestamp=1535017216526;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:16.560000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.030306 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.001534 Prepare_Push: 0.001759 dn1_0_First_Result_Fetch: 0.011846 dn1_0_Last_Result_Fetch: 0.001663 Generate_New_Query: 0.003223 dn1_I
SET timestamp=1535017216560;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:17.325000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.017545 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.006231 Prepare_Push: 0.002335 dn1_First_Result_Fetch: 0.008121 dn1_Last_Result_Fetch: 0.000277 Write_Client: 0.000857
SET timestamp=1535017217325;
select * from sharding_two_node where id =1;
# Time: 2018-08-23T17:40:17.390000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.026216 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000854 Prepare_Push: 0.000904 dn1_First_Result_Fetch: 0.024157 dn1_Last_Result_Fetch: 0.000081 Write_Client: 0.000301
SET timestamp=1535017217390;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:17.410000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.020098 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000699 Prepare_Push: 0.000711 dn2_First_Result_Fetch: 0.017634 dn2_Last_Result_Fetch: 0.000132 Write_Client: 0.001051
SET timestamp=1535017217411;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:17.491000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.078505 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.001702 Prepare_Push: 0.000763 dn2_First_Result_Fetch: 0.018547 dn1_First_Result_Fetch: 0.018482 dn2_Last_Result_Fetch: 0.036637 dn1_I
SET timestamp=1535017217491;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:17.518000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.026112 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000686 Prepare_Push: 0.000872 dn1_0_First_Result_Fetch: 0.007054 dn1_0_Last_Result_Fetch: 0.001072 dn2_0_First_Result_Fetch: 0.005839
SET timestamp=1535017217518;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:17.558000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.038199 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000966 Prepare_Push: 0.005189 dn1_0_First_Result_Fetch: 0.013137 dn1_0_Last_Result_Fetch: 0.001134 Generate_New_Query: 0.003973 dn1_I
SET timestamp=1535017217558;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:18.353000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.019048 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.003008 Prepare_Push: 0.000844 dn2_First_Result_Fetch: 0.006415 dn1_First_Result_Fetch: 0.009082 dn2_Last_Result_Fetch: 0.000323 dn1_I
SET timestamp=1535017218353;
select * from sharding_two_node;
# Time: 2018-08-23T17:40:18.410000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.025498 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000060 Prepare_Push: 0.000696 dn1_First_Result_Fetch: 0.024394 dn1_Last_Result_Fetch: 0.000084 Write_Client: 0.000348
SET timestamp=1535017218410;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:18.430000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.018794 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000047 Prepare_Push: 0.001301 dn2_First_Result_Fetch: 0.017073 dn2_Last_Result_Fetch: 0.000099 Write_Client: 0.000373
SET timestamp=1535017218430;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:18.471000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.039810 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000052 Prepare_Push: 0.000661 dn2_First_Result_Fetch: 0.019799 dn1_First_Result_Fetch: 0.019923 dn2_Last_Result_Fetch: 0.000698 dn1_I
SET timestamp=1535017218471;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:18.484000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.012214 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000047 Prepare_Push: 0.001782 dn1_0_First_Result_Fetch: 0.007109 dn1_0_Last_Result_Fetch: 0.001544 dn2_0_First_Result_Fetch: 0.005518
SET timestamp=1535017218484;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:18.507000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.019695 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000046 Prepare_Push: 0.001448 dn1_0_First_Result_Fetch: 0.006244 dn1_0_Last_Result_Fetch: 0.000988 Generate_New_Query: 0.001564 dn1_I
SET timestamp=1535017218507;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:19.351000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.020937 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000059 Prepare_Push: 0.000800 dn1_First_Result_Fetch: 0.019607 dn1_Last_Result_Fetch: 0.000169 Write_Client: 0.000472
SET timestamp=1535017219351;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:19.370000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.018011 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.001184 Prepare_Push: 0.000583 dn2_First_Result_Fetch: 0.015894 dn2_Last_Result_Fetch: 0.000129 Write_Client: 0.000351
SET timestamp=1535017219370;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:19.412000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.041319 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000689 Prepare_Push: 0.000573 dn2_First_Result_Fetch: 0.017735 dn1_First_Result_Fetch: 0.017876 dn2_Last_Result_Fetch: 0.000601 dn1_I
SET timestamp=1535017219412;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:19.423000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.010063 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000200 Prepare_Push: 0.001136 dn1_0_First_Result_Fetch: 0.006820 dn1_0_Last_Result_Fetch: 0.000694 dn2_0_First_Result_Fetch: 0.003519
SET timestamp=1535017219423;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:19.454000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.027592 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000182 Prepare_Push: 0.012798 dn1_0_First_Result_Fetch: 0.005960 dn1_0_Last_Result_Fetch: 0.000530 Generate_New_Query: 0.000811 dn1_I
SET timestamp=1535017219454;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:20.312000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.025903 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.001470 Prepare_Push: 0.000887 dn1_First_Result_Fetch: 0.022114 dn1_Last_Result_Fetch: 0.000197 Write_Client: 0.001433
SET timestamp=1535017220312;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:20.342000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.028503 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.005643 Prepare_Push: 0.001172 dn2_First_Result_Fetch: 0.021342 dn2_Last_Result_Fetch: 0.000074 Write_Client: 0.000346
SET timestamp=1535017220342;
delete from sharding_two_node where id =-519;
# Time: 2018-08-23T17:40:20.381000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.037424 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000641 Prepare_Push: 0.000959 dn2_First_Result_Fetch: 0.015139 dn1_First_Result_Fetch: 0.015238 dn2_Last_Result_Fetch: 0.000795 dn1_I
SET timestamp=1535017220408;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:21.214000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.023376 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000073 Prepare_Push: 0.001306 dn1_First_Result_Fetch: 0.021694 dn1_Last_Result_Fetch: 0.000081 Write_Client: 0.000302

```

```

SET timestamp=1535017221214;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:21.241000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.025408  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000083  Prepare_Push: 0.001029  dn2_First_Result_Fetch: 0.023856  dn2_Last_Result_Fetch: 0.000122  Write_Client: 0.000440
SET timestamp=1535017221241;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:21.281000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.038482  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000087  Prepare_Push: 0.000871  dn2_First_Result_Fetch: 0.016690  dn1_First_Result_Fetch: 0.016708  dn2_Last_Result_Fetch: 0.000579  dn1_Last_Result_Fetch: 0.000122
SET timestamp=1535017221281;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:21.293000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.011657  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000615  Prepare_Push: 0.001320  dn1_0_First_Result_Fetch: 0.006906  dn1_0_Last_Result_Fetch: 0.001589  dn2_0_First_Result_Fetch: 0.005105
SET timestamp=1535017221293;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:21.312000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.017169  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000635  Prepare_Push: 0.001609  dn1_0_First_Result_Fetch: 0.006997  dn1_0_Last_Result_Fetch: 0.000728  Generate_New_Query: 0.001037  dn1_0_Generate_New_Query: 0.0001037
SET timestamp=1535017221312;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:22.150000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.026153  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000180  Prepare_Push: 0.000771  dn1_First_Result_Fetch: 0.024940  dn1_Last_Result_Fetch: 0.000061  Write_Client: 0.000261
SET timestamp=1535017221250;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:22.170000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.019181  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000180  Prepare_Push: 0.000642  dn2_First_Result_Fetch: 0.018060  dn2_Last_Result_Fetch: 0.000088  Write_Client: 0.000299
SET timestamp=1535017222170;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:22.220000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.049834  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000261  Prepare_Push: 0.000735  dn2_First_Result_Fetch: 0.019862  dn1_First_Result_Fetch: 0.019807  dn2_Last_Result_Fetch: 0.000418  dn1_Last_Result_Fetch: 0.000148
SET timestamp=1535017222220;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:22.240000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.019128  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000934  Prepare_Push: 0.002731  dn1_0_First_Result_Fetch: 0.013430  dn1_0_Last_Result_Fetch: 0.000521  dn2_0_First_Result_Fetch: 0.003296
SET timestamp=1535017222240;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:22.270000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.028479  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.003233  Prepare_Push: 0.004986  dn1_0_First_Result_Fetch: 0.009870  dn1_0_Last_Result_Fetch: 0.001172  Generate_New_Query: 0.001590  dn1_0_Generate_New_Query: 0.0001590
SET timestamp=1535017222270;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:23.097000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.053956  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000056  Prepare_Push: 0.001034  dn1_First_Result_Fetch: 0.052343  dn1_Last_Result_Fetch: 0.000174  Write_Client: 0.000523
SET timestamp=1535017223097;
select * from sharding_two_node where id =1;
# Time: 2018-08-23T17:40:23.110000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.010839  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000499  Prepare_Push: 0.000680  dn2_First_Result_Fetch: 0.006349  dn1_First_Result_Fetch: 0.009082  dn2_Last_Result_Fetch: 0.000270  dn1_Last_Result_Fetch: 0.000182
SET timestamp=1535017223110;
select * from sharding_two_node;
# Time: 2018-08-23T17:40:23.181000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.027573  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000109  Prepare_Push: 0.000980  dn2_First_Result_Fetch: 0.026156  dn2_Last_Result_Fetch: 0.000086  Write_Client: 0.000328
SET timestamp=1535017223181;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:23.231000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.049380  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.002435  Prepare_Push: 0.000670  dn2_First_Result_Fetch: 0.025278  dn1_First_Result_Fetch: 0.025242  dn2_Last_Result_Fetch: 0.000392  dn1_Last_Result_Fetch: 0.000268
SET timestamp=1535017223231;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:23.268000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.025207  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000060  Prepare_Push: 0.001492  dn1_0_First_Result_Fetch: 0.007693  dn1_0_Last_Result_Fetch: 0.000752  Generate_New_Query: 0.001946  dn1_0_Generate_New_Query: 0.0001946
SET timestamp=1535017223268;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:24.121000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.027104  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.001558  Prepare_Push: 0.001107  dn1_First_Result_Fetch: 0.024084  dn1_Last_Result_Fetch: 0.000085  Write_Client: 0.000356
SET timestamp=1535017224121;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:24.141000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.019191  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000072  Prepare_Push: 0.000673  dn2_First_Result_Fetch: 0.017923  dn2_Last_Result_Fetch: 0.000092  Write_Client: 0.000522
SET timestamp=1535017224141;
select * from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:24.182000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.039883  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000483  Prepare_Push: 0.000584  dn2_First_Result_Fetch: 0.017241  dn1_First_Result_Fetch: 0.017320  dn2_Last_Result_Fetch: 0.000603  dn1_Last_Result_Fetch: 0.000482
SET timestamp=1535017224182;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:24.196000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.012406  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000073  Prepare_Push: 0.000958  dn1_0_First_Result_Fetch: 0.008102  dn1_0_Last_Result_Fetch: 0.001255  dn2_0_First_Result_Fetch: 0.007566
SET timestamp=1535017224196;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:24.218000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.021238  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000084  Prepare_Push: 0.001143  dn1_0_First_Result_Fetch: 0.010305  dn1_0_Last_Result_Fetch: 0.000532  Generate_New_Query: 0.001852  dn1_0_Generate_New_Query: 0.0001852
SET timestamp=1535017224218;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:24.25093000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.029579  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000106  Prepare_Push: 0.000882  dn1_First_Result_Fetch: 0.028241  dn1_Last_Result_Fetch: 0.000069  Write_Client: 0.000351
SET timestamp=1535017225093;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:25.121000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.027422  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.001867  Prepare_Push: 0.001330  dn2_First_Result_Fetch: 0.023887  dn2_Last_Result_Fetch: 0.000102  Write_Client: 0.000338
SET timestamp=1535017225121;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:25.161000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.038850  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000050  Prepare_Push: 0.000753  dn2_First_Result_Fetch: 0.019091  dn1_First_Result_Fetch: 0.019189  dn2_Last_Result_Fetch: 0.000582  dn1_Last_Result_Fetch: 0.000516
SET timestamp=1535017225161;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:25.191000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.016379  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000043  Prepare_Push: 0.001276  dn1_0_First_Result_Fetch: 0.007469  dn1_0_Last_Result_Fetch: 0.000678  Generate_New_Query: 0.001327  dn1_0_Generate_New_Query: 0.0001327
SET timestamp=1535017225191;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:26.026000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.024231  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.001105  Prepare_Push: 0.000469  dn2_First_Result_Fetch: 0.022188  dn2_Last_Result_Fetch: 0.000100  Write_Client: 0.000470
SET timestamp=1535017226051;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:26.051000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.024231  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.001105  Prepare_Push: 0.000469  dn2_First_Result_Fetch: 0.022188  dn2_Last_Result_Fetch: 0.000100  Write_Client: 0.000470
SET timestamp=1535017226051;
select count(*) from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:26.091000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.039762  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.001669  Prepare_Push: 0.001915  dn2_First_Result_Fetch: 0.018107  dn1_First_Result_Fetch: 0.018187  dn2_Last_Result_Fetch: 0.000633  dn1_Last_Result_Fetch: 0.000633
SET timestamp=1535017226091;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:26.105000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.012654  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000261  Prepare_Push: 0.000935  dn1_0_First_Result_Fetch: 0.007328  dn1_0_Last_Result_Fetch: 0.000733  dn2_0_First_Result_Fetch: 0.006229
SET timestamp=1535017226105;
select count(*) from sharding_two_node;
# Time: 2018-08-23T17:40:26.134000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.028335  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000046  Prepare_Push: 0.003442  dn1_0_First_Result_Fetch: 0.009563  dn1_0_Last_Result_Fetch: 0.001069  Generate_New_Query: 0.001856  dn1_0_Generate_New_Query: 0.0001856
SET timestamp=1535017226134;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
# Time: 2018-08-23T17:40:26.185000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:0:1] Id:  2
# Query_time: 0.014882  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0  Read_SQL: 0.000067  Prepare_Push: 0.001351  dn1_First_Result_Fetch: 0.013084  dn1_Last_Result_Fetch: 0.000137  Write_Client: 0.000381
SET timestamp=1535017226859;
```

```

select * from sharding_two_node where id =1;
# Time: 2018-08-23T17:40:26.874000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.010509 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000066 Prepare_Push: 0.001761 dn2_First_Result_Fetch: 0.006921 dn1_First_Result_Fetch: 0.008256 dn2_Last_Result_Fetch: 0.000279 dn1_L
SET timestamp=1535017226874;
select * from sharding_two_node;
# Time: 2018-08-23T17:40:26.931000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.028690 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000237 Prepare_Push: 0.001126 dn1_First_Result_Fetch: 0.026194 dn1_Last_Result_Fetch: 0.000640 Write_Client: 0.001133
SET timestamp=1535017226931;
delete from sharding_two_node where id =15;
# Time: 2018-08-23T17:40:26.951000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.018818 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000342 Prepare_Push: 0.001671 dn2_First_Result_Fetch: 0.016482 dn2_Last_Result_Fetch: 0.000063 Write_Client: 0.000323
SET timestamp=1535017226951;
delete from sharding_two_node where id =519;
# Time: 2018-08-23T17:40:26.991000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.039399 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000082 Prepare_Push: 0.000706 dn2_First_Result_Fetch: 0.019233 dn1_First_Result_Fetch: 0.019167 dn2_Last_Result_Fetch: 0.000426 dn1_L
SET timestamp=1535017226991;
insert into sharding_two_node values(15,'15',15),(519,'519',519);
# Time: 2018-08-23T17:40:27.032000Z
# User@Host: root[root] @ [0:0:0:0:0:0:0:1] Id: 2
# Query_time: 0.029495 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.000064 Prepare_Push: 0.001349 dn1_0_First_Result_Fetch: 0.005745 dn1_0_Last_Result_Fetch: 0.000632 Generate_New_Query: 0.001056 dn1_L
SET timestamp=1535017227032;
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);

```

### 2.20.3.1 mysqldumpslow 结果:

```

Reading mysql slow query log from /tmp/slow3.log
Count: 17 Time=0.05s (0s) Lock=0.00s (0s) Rows=0.0 (0), root[root]@[0:0:0:0:0:0:0:1]
insert into sharding_two_node values(N,'S',N),(N,'S',N)

Count: 13 Time=0.05s (0s) Lock=0.00s (0s) Rows=0.0 (0), root[root]@[0:0:0:0:0:0:0:1]
select count(*) from sharding_two_node

Count: 6 Time=0.04s (0s) Lock=0.00s (0s) Rows=0.0 (0), root[root]@[0:0:0:0:0:0:0:1]
select * from sharding_two_node where id =N

Count: 33 Time=0.03s (0s) Lock=0.00s (0s) Rows=0.0 (0), root[root]@[0:0:0:0:0:0:0:1]
delete from sharding_two_node where id =N

Count: 17 Time=0.03s (0s) Lock=0.00s (0s) Rows=0.0 (0), root[root]@[0:0:0:0:0:0:0:1]
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=N)

Count: 6 Time=0.02s (0s) Lock=0.00s (0s) Rows=0.0 (0), root[root]@[0:0:0:0:0:0:0:1]
select * from sharding_two_node

```

### 2.20.3.2 pt-query-digest 结果:

```

# 710ms user time, 70ms system time, 23.35M rss, 68.36M vsz
# Current date: Thu Aug 23 17:48:25 2018
# Hostname: 10-186-24-63
# File: /tmp/slow_query4.log
# Overall: 92 total, 6 unique, 5.41 QPS, 0.18x concurrency -----
# Time range: 2018-08-23T17:40:10 to 2018-08-23T17:40:27
# Attribute          total      min      max      avg     95%    stddev   median
# ======          ======      ==      ==      ==      ==     ==      ==      ==
# Exec time          3s       10ms    339ms    34ms    75ms    37ms    27ms
# Lock time          0        0       0       0       0       0       0       0
# Rows sent          0        0       0       0       0       0       0       0
# Rows examine       0        0       0       0       0       0       0       0
# Query size        4.91k    31      94      54.64   92.72   20.87   42.48
# Generate New      0.03    0.00    0.00    0.00    0.00    0.00    0.00
# Prepare Push       0.56    0.00    0.26    0.01    0.01    0.03    0.00
# Read SQL          0.07    0.00    0.01    0.00    0.00    0.00    0.00
# Write Client       0.70    0.00    0.11    0.01    0.02    0.01    0.00
# dn1 0 First       0.29    0.00    0.05    0.01    0.01    0.01    0.01
# dn1 0 Last R      0.07    0.00    0.03    0.00    0.01    0.01    0.00
# dn1 1 First       0.12    0.00    0.02    0.01    0.01    0.00    0.01
# dn1 1 Last R      0.03    0.00    0.01    0.00    0.00    0.00    0.00
# dn1 First Re      0.93    0.01    0.05    0.02    0.03    0.01    0.02
# dn1 Last Res      0.04    0.00    0.01    0.00    0.00    0.00    0.00
# dn2 0 First       0.13    0.00    0.04    0.01    0.02    0.01    0.01
# dn2 0 Last R      0.08    0.00    0.03    0.01    0.02    0.01    0.00
# dn2 First Re      0.80    0.01    0.06    0.02    0.03    0.01    0.02
# dn2 Last Res      0.09    0.00    0.04    0.00    0.01    0.01    0.00

# Profile
# Rank Query ID          Response time Calls R/Call V/M   I
# ====== ======          ====== ====== ====== ====== ====== ======
# 1 0x13233F8ADA41C6E20889AEE0C2B...  0.8815 28.1%  33  0.0267  0.00 DELETE sharding_two_node
# 2 0xF68D16B46E487184E8FD3BB3912...  0.8525 27.1%  17  0.0501  0.01 INSERT sharding_two_node
# 3 0x846D813C53609C853F7CBA6D2DB...  0.6306 20.1%  13  0.0485  0.15 SELECT sharding_two_node
# 4 0x3FB41587E746A475282C1ED2606...  0.4335 13.8%  17  0.0255  0.00 SELECT sharding_two_node
# 5 0x04CD91DDFC4E1DD7A22E312C72...  0.2399  7.6%   6  0.0400  0.05 SELECT sharding_two_node
# MISC 0xMISC           0.1028  3.3%   0.0  0.0171  0.0 <1 ITEMS>

# Query 1: 2.06 QPS, 0.06x concurrency, ID 0x13233F8ADA41C6E2D889AEE0C2BC6CB5 at byte 943
# Scores: V/M = 0.00
# Time range: 2018-08-23T17:40:10 to 2018-08-23T17:40:26
# Attribute          pct total      min      max      avg     95%    stddev   median
# ======          ====== ======      ==      ==      ==      ==     ==      ==
# Count            35      33
# Exec time        28     882ms    18ms    45ms    27ms    31ms    5ms    27ms
# Lock time         0        0       0       0       0       0       0       0
# Rows sent         0        0       0       0       0       0       0       0
# Rows examine      0        0       0       0       0       0       0       0
# Query size       27     1.37k    42      43     42.52   42.48   0.59   42.48
# Prepare Push      6     0.04    0.00    0.00    0.00    0.00    0.00    0.00
# Read SQL         35     0.03    0.00    0.01    0.00    0.00    0.00    0.00
# Write Client      2     0.02    0.00    0.00    0.00    0.00    0.00    0.00
# dn1 First Re     45     0.42    0.02    0.04    0.03    0.03    0.01    0.03
# dn1 Last Res     7     0.00    0.00    0.00    0.00    0.00    0.00    0.00
# dn2 First Re     47     0.38    0.02    0.03    0.02    0.03    0.00    0.02
# dn2 Last Res     3     0.00    0.00    0.00    0.00    0.00    0.00    0.00
# String:
# Hosts          0:0:0:0:0:0:1
# Users          root
# Query_time distribution
#   1us
#   10us
#   100us
#   1ms
#   10ms #####
#   100ms #####
#   1s
#   10s+
# Tables
#   SHOW TABLE STATUS LIKE 'sharding_two_node'\G
#   SHOW CREATE TABLE `sharding_two_node`\G
delete from sharding_two_node where id =15\G
# Converted for EXPLAIN
# EXPLAIN /*!50100 PARTITIONS*/
select * from sharding_two_node where id =15\G

# Query 2: 1.66 QPS, 0.05x concurrency, ID 0xF68D16B46E487184E8FD3BB3912A3658 at byte 1690
# Scores: V/M = 0.01
# Time range: 2018-08-23T17:40:10 to 2018-08-23T17:40:26
# Attribute          pct total      min      max      avg     95%    stddev   median
# ======          ====== ======      ==      ==      ==      ==     ==      ==
# Count            18      17
# Exec time        27     853ms    37ms   117ms    50ms    78ms    21ms   40ms
# Lock time         0        0       0       0       0       0       0       0
# Rows sent         0        0       0       0       0       0       0       0
# Rows examine      0        0       0       0       0       0       0       0
# Query size       21     1.06k    64      64     64.00   64.00   0       64
# Prepare Push      7     0.04    0.00    0.02    0.00    0.00    0.01    0.00
# Read SQL         32     0.02    0.00    0.01    0.00    0.00    0.00    0.00
# Write Client     63     0.45    0.02    0.06    0.03    0.05    0.01    0.02
# dn1 First Re     37     0.35    0.01    0.04    0.02    0.03    0.00    0.02
# dn1 Last Res     40     0.01    0.00    0.00    0.00    0.00    0.00    0.00
# dn2 First Re     46     0.37    0.01    0.06    0.02    0.03    0.01    0.02
# dn2 Last Res     72     0.07    0.00    0.04    0.00    0.02    0.01    0.00
# String:
# Hosts          0:0:0:0:0:0:1
# Users          root
# Query_time distribution
#   1us
#   10us
#   100us
#   1ms
#   10ms #####
#   100ms #####
#   1s
#   10s+
# Tables
#   SHOW TABLE STATUS LIKE 'sharding_two_node'\G
#   SHOW CREATE TABLE `sharding_two_node`\G
insert into sharding_two_node values(15,'15',15) /*... omitted ...*/\G

# Query 3: 0.81 QPS, 0.04x concurrency, ID 0xB46D813C53609C853F7CBA6D2DB4047C at byte 2152
# Scores: V/M = 0.15
# Time range: 2018-08-23T17:40:10 to 2018-08-23T17:40:26
# Attribute          pct total      min      max      avg     95%    stddev   median
# ======          ====== ======      ==      ==      ==      ==     ==      ==
# Count            14      13
# Exec time        20     631ms    10ms    339ms    49ms   116ms   84ms   12ms
# Lock time         0        0       0       0       0       0       0       0
# Rows sent         0        0       0       0       0       0       0       0
# Rows examine      0        0       0       0       0       0       0       0
# Query size       9     494     38      38     38.00   38.00   0       38
# Prepare Push     48     0.27    0.00    0.26    0.02    0.00    0.07    0.00
# Read SQL          5     0.00    0.00    0.00    0.00    0.00    0.00    0.00
# Write Client     24     0.17    0.00    0.11    0.01    0.02    0.03    0.00
# dn1 0 First      48     0.14    0.00    0.05    0.01    0.01    0.01    0.01
# dn1 0 Last R     80     0.06    0.00    0.03    0.00    0.01    0.01    0.00
# dn2 0 First      100    0.13    0.00    0.04    0.01    0.02    0.01    0.01
# dn2 0 Last R     100    0.08    0.00    0.03    0.01    0.02    0.01    0.00
# String:
# Hosts          0:0:0:0:0:0:1
# Users          root
# Query_time distribution
#   1us
#   10us
#   100us
#   1ms
#   10ms #####
#   100ms #####
#   1s
#   10s+
# Tables
#   SHOW TABLE STATUS LIKE 'sharding_two_node'\G
#   SHOW CREATE TABLE `sharding_two_node`\G

```

```

# EXPLAIN /*!50100 PARTITIONS*/
select count(*) from sharding_two_node\G

# Query 4: 1 QPS, 0.03x concurrency, ID 0x3FB41587E746A475282C1ED2606795FB at byte 2596
# Scores: V/M = 0.00
# Time range: 2018-08-23T17:40:10 to 2018-08-23T17:40:27
# Attribute   pct  total    min     max   avg   95% stddev median
# ====== ====== ====== ====== ====== ====== ====== ====== ======
# Count      18    17
# Exec time   13   434ms   16ms   47ms   26ms   38ms   8ms   26ms
# Lock time    0     0     0     0     0     0     0     0     0
# Rows sent    0     0     0     0     0     0     0     0     0
# Rows examine  0     0     0     0     0     0     0     0     0
# Query size   31   1.56k   94     94   94     94     0     94
# Generate New 100   0.03   0.00   0.00   0.00   0.00   0.00   0.00
# Prepare Push 13   0.07   0.00   0.03   0.00   0.01   0.01   0.00
# Read SQL    12   0.01   0.00   0.00   0.00   0.00   0.00   0.00
# Write Client  2   0.02   0.00   0.00   0.00   0.00   0.00   0.00
# dn1 0 First  51   0.15   0.00   0.01   0.01   0.01   0.00   0.01
# dn1 0 Last R 19   0.01   0.00   0.00   0.00   0.00   0.00   0.00
# dn1 1 First 100   0.12   0.00   0.02   0.01   0.01   0.00   0.01
# dn1 1 Last R 100   0.03   0.00   0.01   0.00   0.00   0.00   0.00
# dn1 1 Last R 100   0.03   0.00   0.01   0.00   0.00   0.00   0.00
# String:
# Hosts          0:0:0:0:0:0:1
# Users         root
# Query_time distribution
#   1us
# 10us
# 100us
# 1ms
# 10ms #####
# 100ms #####
# 1s
# 10s+
# Tables
#   SHOW TABLE STATUS LIKE 'sharding_two_node'\G
#   SHOW CREATE TABLE `sharding_two_node`\G
# EXPLAIN /*!50100 PARTITIONS*/
select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1)\G

# Query 5: 0.38 QPS, 0.01x concurrency, ID 0x04CDF91DDFC4E1DD7A22E312C72C268D at byte 0
# Scores: V/M = 0.05
# Time range: 2018-08-23T17:40:10 to 2018-08-23T17:40:26
# Attribute   pct  total    min     max   avg   95% stddev median
# ====== ====== ====== ====== ====== ====== ====== ====== ======
# Count      6    6
# Exec time   7   240ms   10ms   133ms   40ms   128ms   43ms   35ms
# Lock time    0     0     0     0     0     0     0     0     0
# Rows sent    0     0     0     0     0     0     0     0     0
# Rows examine  0     0     0     0     0     0     0     0     0
# Query size   5   258     43     43   43     43     0     43
# Generate New 22   0.12   0.00   0.12   0.02   0.12   0.04   0.00
# Prepare Push  9   0.01   0.00   0.01   0.00   0.01   0.00   0.00
# Read SQL    11   0.10   0.01   0.05   0.02   0.05   0.02   0.01
# Write Client  6   0.00   0.00   0.00   0.00   0.00   0.00   0.00
# dn1 Last Res 6   0.00   0.00   0.00   0.00   0.00   0.00   0.00
# String:
# Hosts          0:0:0:0:0:0:1
# Users         root
# Query_time distribution
#   1us
# 10us
# 100us
# 1ms
# 10ms #####
# 100ms #####
# 1s
# 10s+
# Tables
#   SHOW TABLE STATUS LIKE 'sharding_two_node'\G
#   SHOW CREATE TABLE `sharding_two_node`\G
# EXPLAIN /*!50100 PARTITIONS*/
select * from sharding_two_node where id =1\G

```

### 2.20.3.3 限制

当mysqldumpslow版本为5.6或更早的版本时，在解析慢日志的过程中，由于mysqldumpslow只能识别'190428 10:28:16'格式的Time字段，而dbe却使用了'2019-04-28T10:28:16.515000Z'字段格式，因此解析生成的慢日志会出错。详情参见issue:<https://github.com/actiontech/dble/issues/908>

### 2.20.4 附：慢查询日志格式解析

#### 2.20.4.1 MySQL 慢查询日志格式

先放一段正常记录的MySQL慢日志

```
/usr/local/mysql5.7.11/bin/mysqld-debug, Version: 5.7.11-debug-log (MySQL Community Server - Debug (GPL)). started with:
Tcp port: 3320 Unix socket: /tmp/mysql_3320.sock
Time           Id Command    Argument
# Time: 2018-05-15T10:53:23.798040Z
# User@Host: action[action] @ [192.168.2.206]  Id: 436
# Query_time: 296.145816  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0
use test;
SET timestamp=1526381603;
drop table sharding_two_node;
# Time: 2018-05-15T11:32:25.549290Z
# User@Host: action[action] @ [192.168.2.206]  Id: 451
# Query_time: 129.555883  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0
use nosharding;
SET timestamp=1526383945;
drop table test4;
# Time: 2018-05-15T11:32:25.550190Z
# User@Host: action[action] @ [192.168.2.206]  Id: 454
# Query_time: 84.316518  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0
SET timestamp=1526383945;
insert into test4 values(1,'1');
# Time: 2018-05-15T11:37:01.079214Z
# User@Host: action[action] @ [192.168.2.206]  Id: 483
# Query_time: 49.571983  Lock_time: 0.000000 Rows_sent: 0  Rows_examined: 0
SET timestamp=1526384221;
drop table test3;
# Time: 2018-07-11T05:28:34.161405Z
# User@Host: action[action] @ [192.168.2.206]  Id: 16421
# Query_time: 10.035706  Lock_time: 0.000000 Rows_sent: 1  Rows_examined: 0
use test;
SET timestamp=1531286914;
insert into test4 values(1,'1');
```

我们从MySQL源代码上来分析慢日志的格式:

### 1.1 文件头: (包含版本信息等)

```
/usr/local/mysql5.7.11/bin/mysqld-debug, Version: 5.7.11-debug-log (MySQL Community Server - Debug (GPL)). started with:
Tcp port: 3320 Unix socket: /tmp/mysql_3320.sock
```

Time Id Command Argument

如果MySQL实例没有专门设置参数log-short-format, 则会有time行和session信息行

### 1.2 time行

例如:

```
# Time: 2018-05-15T10:53:23.798040Z
```

### 1.3 session信息行

例如:

```
# User@Host: action[action] @ [192.168.2.206] Id: 436
```

### 1.4 执行时间行

注意用两个空格分隔不同键值对, 用单个空格区分key和value

```
# Query_time: 296.145816  Lock_time: 0.000000  Rows_sent: 0  Rows_examined: 0
```

### 1.5 如果有Database changed, 则有转换schema行:

```
use nosharding;
```

### 1.6 set行(注: mysqldumpslow官方工具只处理了SET timestamp=)

举例

```
SET timestamp=1526383945;
```

#### 1.6.1 last\_insert\_id

如果设置了

```
stmt_depends_on_first_successful_insert_id_in_prev_stmt (含义待调查)
```

会有last\_insert\_id=某个值

#### 1.6.2 insert\_id

如果没有设置log-short-format, 并且

```
auto_inc_intervals_in_cur_stmt_for_binlog.nb_elements() (含义待调查)
```

会有last\_id=某个值

#### 1.6.3 timestamp=

设置时间戳

### 1.7 命令行

例如:

```
# administrator command
```

如果是is\_command (含义待调查), 则有此行

### 1.8 SQL语句行

例如

```
insert into test4 values(1,'1');
```

## 2.20.4.2 dble 慢日志格式

为了兼容mysqldumpslow和pt-query-digest 工具, dble的慢日志格式如下:

## 2.1 文件头

```
/FAKE_PATH/mysql, Version: FAKE_VERSION. started with:  
Tcp port: 3320 Unix socket: FAKE_SOCK  
Time Id Command Argument
```

## 2.2 time行

因为java8目前只能获取到毫秒级别的绝对时间戳，所以时间戳后三位为0，与mysql不同，例如：

```
# Time: 2018-08-23T17:40:10.149000Z
```

## 2.3 session信息行

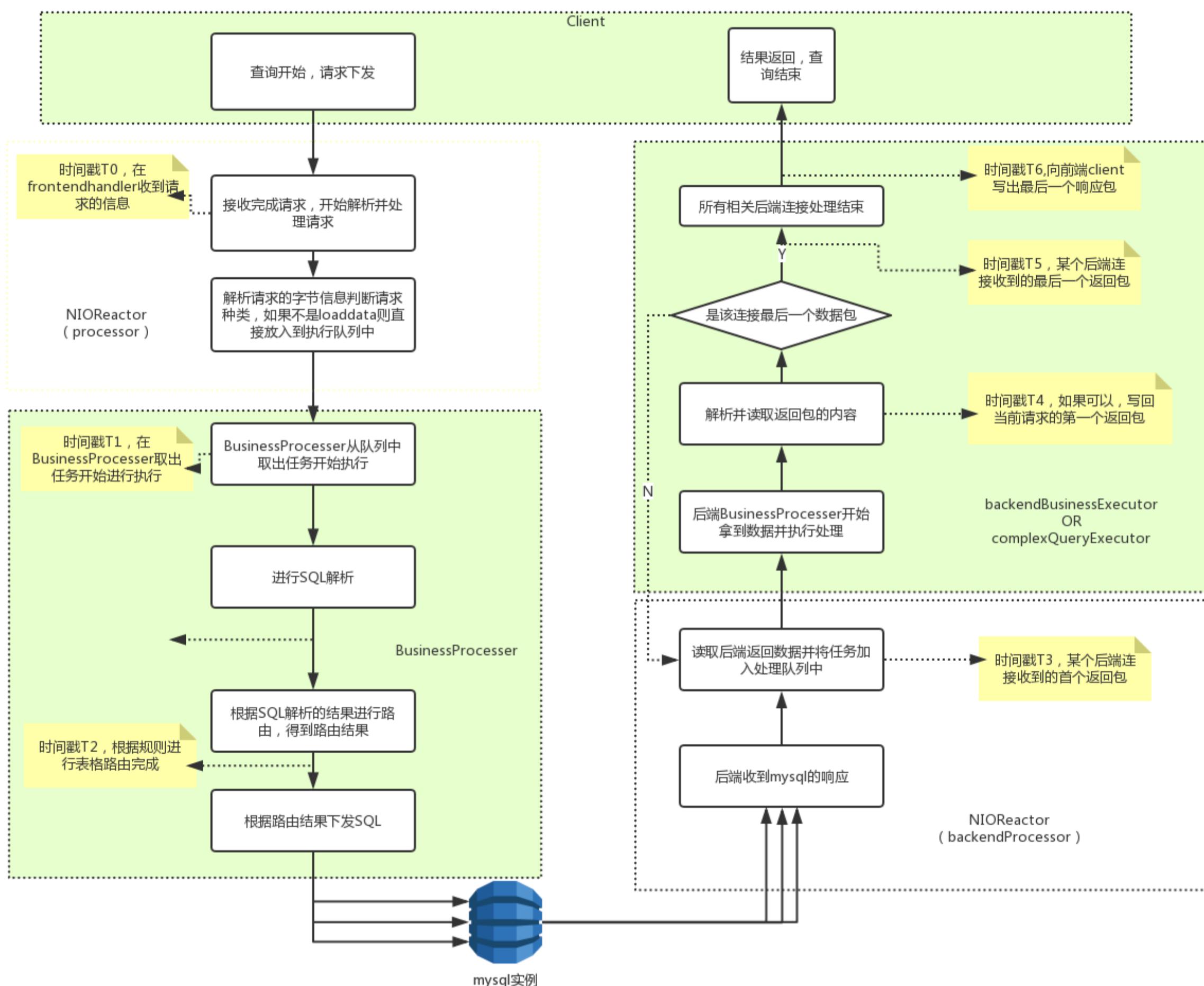
例如：

```
# User@Host: root[root] @ [0:0:0:0:0:0:1] Id: 2
```

## 2.4 执行时间行

注意用两个空格分隔不同键值对，用单个空格区分key和value。

如图，这是dble的内部流程，根据此图：



我们在兼容MySQL的情况下增加了如下字段：

**Read\_SQL:** 接受到SQL字节数据开始到将字节数据转为可供解析的SQL字符串的耗时，即T1-T0

**Prepare\_Push:** 解析/路由/尝试优化及其他准备工作耗时，与Read\_SQL是线性顺序的，即T2-T1

**\${datanodeName}\_First\_Result\_Fetch:** 某个datanode收到的第一个数据包，相对于Prepare\_Push结束时的耗时，即T3-T2。另外：不同的datanode之间是并行的

**\${datanodeName}\_Last\_Result\_Fetch:** 某个datanode收到的最后一个数据包，相对于自己的\${datanodeName}\_First\_Result\_Fetch的耗时，即T5-T3

**Write\_Client:** 可以开始返回客户端（收到一个后端返回数据包）到最后一个数据包准备写出的时间，如果是复杂查询，可能包含合并等操作的结果，即T6-T4

另外，MySQL原有的Query\_time为dble开始读取SQL到准备写回最后一个数据包的时间间隔，即T6-T0

**Lock\_time, Rows\_sent, Rows\_examined**用0填充

例如：

```
# Query_time: 0.116672 Lock_time: 0.000000 Rows_sent: 0 Rows_examined: 0 Read_SQL: 0.013625 Prepare_Push: 0.024767 dn2_First_Result_Fetch: 0.056395 dn1_First_Result_Fetch: 0.026420 dn2_Last_Result_Fetch: 0.000743 dn1_Last_Result_Fetch: 0.001700 Write_Client: 0.051861
```

## 2.5 set行(只处理了SET timestamp=)

例如：

```
SET timestamp=1535017210432;
```

## 2.6 SQL语句行

例如

```
insert into sharding_two_node values(15,'15',15),(519,'519',519);
```

## 2.21 单条SQL性能trace

这个功能与MySQL的profile功能类似，用于统计一条查询在db引擎内部的各个阶段的耗时情况，用来发现性能瓶颈，改进SQL或者db引擎配置，达到提高性能的目的。此功能为session级别。

### 1. 可以通过此命令检查trace功能是否开启

```
mysql> select @@trace;
+-----+
| @@trace |
+-----+
| 0      |
+-----+
1 row in set (0.02 sec)
```

### 2. 开启 trace 功能

```
mysql> set trace =1;
Query OK, 0 rows affected (0.09 sec)

mysql> select @@trace;
+-----+
| @@trace |
+-----+
| 1      |
+-----+
1 row in set (0.00 sec)
```

### 3. 单节点查询的trace结果举例

```
mysql> select * from sharding_two_node where id =1;
+----+-----+-----+
| id | c_flag | c_decimal |
+----+-----+-----+
| 1  | 1     | 1.0000 |
+----+-----+-----+
1 row in set (0.02 sec)

mysql> show trace;
+-----+-----+-----+-----+-----+
| OPERATION | START(ms) | END(ms) | DURATION(ms) | DATA_NODE | SQL/REF |
+-----+-----+-----+-----+-----+
| Read SQL  | 0.0       | 0.1085  | 0.1085    | -         |          |
| Parse SQL | 0.1085   | 0.49607 | 0.38757   | -         |          |
| Route Calculation | 0.49607 | 1.274142 | 0.778072 | -         |          |
| Prepare to Push | 1.274142 | 1.560543 | 0.286401 | -         |          |
| Execute SQL | 1.560543 | 18.711851 | 17.151308 | dn1       | select * from sharding_two_node where id =1 |
| Fetch result | 18.711851 | 18.978213 | 0.266362 | dn1       | select * from sharding_two_node where id =1 |
| Write to Client | 18.711851 | 19.276344 | 0.564493 | -         |          |
| Over All    | 0.0       | 19.276344 | 19.276344 | -         |          |
+-----+-----+-----+-----+-----+
8 rows in set (0.00 sec)
```

### 4. 多节点查询的trace结果举例

```
mysql> select * from sharding_two_node ;
+----+-----+-----+
| id | c_flag | c_decimal |
+----+-----+-----+
| 513 | 513   | 513.0000 |
| 514 | 514   | 514.0000 |
| 515 | 515   | 515.0000 |
| 516 | 516   | 516.0000 |
| 1  | 1     | 1.0000 |
| 2  | 2     | 2.0000 |
| 3  | 3     | 3.0000 |
| 4  | 4     | 4.0000 |
| 5  | 5     | 5.0000 |
| 7  | 7     | 7.0000 |
| 8  | 8     | 8.0000 |
| 9  | 9     | 9.0000 |
| 10 | 10   | 10.0000 |
| 11 | 11   | 11.0000 |
| 12 | 12   | 12.0000 |
+----+-----+-----+
15 rows in set (0.01 sec)

mysql> show trace;
+-----+-----+-----+-----+-----+
| OPERATION | START(ms) | END(ms) | DURATION(ms) | DATA_NODE | SQL/REF |
+-----+-----+-----+-----+-----+
| Read SQL  | 0.0       | 0.079175 | 0.079175   | -         |          |
| Parse SQL | 0.079175 | 0.637315 | 0.55814    | -         |          |
| Route Calculation | 0.637315 | 1.046389 | 0.409974   | -         |          |
| Prepare to Push | 1.046389 | 1.465238 | 0.418849   | -         |          |
| Execute SQL | 1.465238 | 8.141409 | 6.676171   | dn1       | SELECT * FROM sharding_two_node LIMIT 100 |
| Execute SQL | 8.141409 | 7.59109  | 6.125852   | dn2       | SELECT * FROM sharding_two_node LIMIT 100 |
| Fetch result | 7.59109  | 8.366718 | 0.775628   | dn1       | SELECT * FROM sharding_two_node LIMIT 100 |
| Fetch result | 7.59109  | 9.324157 | 1.733067   | dn2       | SELECT * FROM sharding_two_node LIMIT 100 |
| Write to Client | 7.59109 | 9.324157 | 9.324157   | -         |          |
| Over All    | 0.0       | 9.324157 | 9.324157   | -         |          |
+-----+-----+-----+-----+-----+
10 rows in set (0.00 sec)
```

### 5. 多节点写入的SQL 的trace结果举例,事实上,这是一个隐式分布式事务

```
mysql> insert into sharding_two_node values(15,'15',15),(519,'519',519);
Query OK, 2 rows affected (0.06 sec)

mysql> show trace;
+-----+-----+-----+-----+-----+
| OPERATION | START(ms) | END(ms) | DURATION(ms) | DATA_NODE | SQL/REF |
+-----+-----+-----+-----+-----+
| Read SQL  | 0.0       | 0.131959 | 0.131959   | -         |          |
| Parse SQL | 0.131959 | 0.601637 | 0.469678   | -         |          |
| Route Calculation | 0.601637 | 0.825479 | 0.223842   | -         |          |
| Prepare to Push | 0.825479 | 1.025374 | 0.199895   | -         |          |
| Execute SQL | 1.025374 | 27.095675 | 26.070301  | dn1       | INSERT INTO sharding_two_node VALUES (15, '15', 15) |
| Execute SQL | 27.095675 | 25.023911 | 23.998537  | dn2       | INSERT INTO sharding_two_node VALUES (519, '519', 519) |
| Fetch result | 25.023911 | 27.405046 | 0.309371   | dn1       | INSERT INTO sharding_two_node VALUES (15, '15', 15) |
| Fetch result | 27.405046 | 26.478398 | 1.454487   | dn2       | INSERT INTO sharding_two_node VALUES (519, '519', 519) |
| Distributed Transaction Prepare | 27.405046 | 27.736411 | 0.331365   | -         |          |
| Distributed Transaction Commit | 27.736411 | 57.426311 | 29.6899   | -         |          |
| Write to Client | 57.426311 | 57.428266 | 32.404355  | -         |          |
| Over All    | 0.0       | 57.428266 | 57.428266  | -         |          |
+-----+-----+-----+-----+-----+
12 rows in set (0.00 sec)
```

### 6. 复杂查询的trace结果举例

```
mysql> select count(*) from sharding_two_node;
+-----+
| COUNT(*) |
+-----+
| 20 |
+-----+
1 row in set (0.01 sec)

mysql> show trace;
+-----+-----+-----+-----+-----+
| OPERATION | START(ms) | END(ms) | DURATION(ms) | DATA_NODE | SQL/REF |
+-----+-----+-----+-----+-----+
| Read SQL | 0.0 | 0.08553 | 0.08553 | - | - |
| Parse SQL | 0.08553 | 0.56987 | 0.48434 | - | - |
| Try Route Calculation | 0.56987 | 0.71698 | 0.14711 | - | - |
| Try to Optimize | 0.71698 | 1.237487 | 0.529507 | - | - |
| Execute SQL | 1.237487 | 0.991029 | 7.853542 | dn1.0 | select COUNT(*) as `$_COUNT$_.rpda_0` from `sharding_two_node` LIMIT 100 |
| Fetch result | 0.991029 | 10.186782 | 1.095753 | dn1.0 | select COUNT(*) as `$_COUNT$_.rpda_0` from `sharding_two_node` LIMIT 100 |
| Execute SQL | 1.237487 | 8.348635 | 7.111148 | dn2.0 | select COUNT(*) as `$_COUNT$_.rpda_0` from `sharding_two_node` LIMIT 100 |
| Fetch result | 8.348635 | 9.342241 | 0.993606 | dn2.0 | select COUNT(*) as `$_COUNT$_.rpda_0` from `sharding_two_node` LIMIT 100 |
| MERGE | 8.721543 | 10.289905 | 1.568362 | merge.1 | dn1.0; dn2.0 |
| ORDERED_GROUP | 8.726918 | 10.424309 | 1.69739 | ordered_group.1 | merge.1 |
| LIMIT | 9.020162 | 10.499574 | 1.479412 | limit.1 | ordered_group.1 |
| SHUFFLE_FIELD | 9.023584 | 10.501529 | 1.477945 | shuffle_field.1 | limit.1 |
| Write to Client | 9.072457 | 11.52055 | 2.448093 | - | - |
| Over All | 0.0 | 11.52055 | 11.52055 | - | - |
+-----+-----+-----+-----+-----+
14 rows in set (0.03 sec)
```

## 7. 子查询的trace结果举例

```
mysql> select count(*) from sharding_two_node where id =(select id from sharding_two_node where id=1);
+-----+
| COUNT(*) |
+-----+
| 1 |
+-----+
1 row in set (0.03 sec)

mysql> show trace;
+-----+-----+-----+-----+-----+
| OPERATION | START(ms) | END(ms) | DURATION(ms) | DATA_NODE | SQL/REF |
+-----+-----+-----+-----+-----+
| Read SQL | 0.0 | 0.063047 | 0.063047 | - | - |
| Parse SQL | 0.063047 | 0.491182 | 0.428135 | - | - |
| Try Route Calculation | 0.491182 | 0.799576 | 0.308394 | - | - |
| Try to Optimize | 0.799576 | 2.347412 | 1.547833 | - | - |
| Execute SQL | 2.347412 | 11.183808 | 8.836396 | dn1.0 | select `sharding_two_node`.`id` as `autoalias_scalar` from `sharding_two_node` where id = 1 LIMIT 2 |
| Fetch result | 11.183808 | 12.360691 | 1.176883 | dn1.0 | select `sharding_two_node`.`id` as `autoalias_scalar` from `sharding_two_node` where id = 1 LIMIT 2 |
| MERGE | 11.889546 | 12.436445 | 0.546899 | merge.1 | dn1.0 |
| LIMIT | 11.894923 | 12.483364 | 0.588441 | limit.1 | merge.1 |
| SHUFFLE_FIELD | 11.896389 | 12.484843 | 0.588441 | shuffle_field.1 | limit.1 |
| SCALAR_SUB_QUERY | 12.038123 | 12.485808 | 0.447688 | scalar_sub_query.1 | shuffle_field.1 |
| Generate New Query | 12.485808 | 13.824463 | 1.338655 | - | - |
| Execute SQL | 13.824463 | 26.749647 | 12.925184 | dn1.1 | scalar_sub_query.1; select COUNT(*) as `$_COUNT$_.rpda_0` from `sharding_two_node` where sharding_two_node.id = 1 |
| Fetch result | 26.685134 | 28.753476 | 2.068342 | dn1.1 | scalar_sub_query.1; select COUNT(*) as `$_COUNT$_.rpda_0` from `sharding_two_node` where sharding_two_node.id = 1 |
| MERGE | 26.954918 | 29.091683 | 2.136765 | merge.2 | dn1.1 |
| ORDERED_GROUP | 26.977889 | 29.563316 | 2.585427 | ordered_group.1 | merge.2 |
| SHUFFLE_FIELD | 27.568285 | 29.567226 | 1.998941 | shuffle_field.2 | ordered_group.1 |
| Write to Client | 27.72517 | 30.014911 | 2.289741 | - | - |
| Over All | 0.0 | 30.014911 | 30.014911 | - | - |
+-----+-----+-----+-----+-----+
18 rows in set (0.01 sec)
```

## 2.22 KILL @@DDL\_LOCK

### 背景

在运维dble集群的过程中，有时会遇到 "There is other session is doing DDL" 或者 "xxx is doing DDL" 等问题，这些问题会导致某些表无法被操作，无法进行 reload 等。此时，需要人工干预保证集群运行正常。

### 逻辑判断

在 dble 集群中，可以通过 ZK 中某些键值对推断正在执行 DDL 的节点状态，比如通过键 `universe/dble/{cluster-id}/ddl/{schema.table}` 的值内容推断。该值理论上只会被发起者节点所修改，其他节点会订阅该键值，该值的内容为json 格式，其中有一个 `status` 字段：

1. `status` 为 INIT 时，发起动作的节点会锁住本节点对应 `table` 的 `table meta` 并获取当前 `key` 的分布式锁。其他的 `dble` 节点订阅到该值后会锁住本节点上对应 `table` 的 `table meta`；
2. 待发起节点执行 DDL 之后，会根据 DDL 执行结果更新当前键的值，此时根据 DDL 执行结果 `status` 分为两种情形：SUCCESS 和 FAILED。其他节点订阅到 `status` 为 SUCCESS 时，才会真正执行对应的 DDL 操作，否则会取消执行 DDL 并释放本节点对应 `table` 的 `table meta` 锁；
3. 每个节点都会插入一条 `universe/dble/{cluster-id}/ddl/{schema.table}/{dble-id}:SUCCESS` 数据。

到此，除发起者的节点执行已经结束，但是发起者节点会等待其他节点响应，只有全部节点都汇报执行完成之后，当前 `ddl` 才算完成。

发起者节点通过什么来判断其他节点都响应完成？根据以下两个键的内容：

- `universe/dble/{cluster-id}/ddl/{schema.table}/{dble-id}`
- `universe/dble/{cluster-id}/online/下属结点`

所有执行当前 DDL 操作的节点都会在 `universe/dble/{cluster-id}/ddl/{schema.table}/` 下面插入记录，发起节点通过判断 `online` 节点是否在 `universe/dble/{cluster-id}/ddl/{schema.table}/` 下面留下记录判断 DDL 是否被该节点所应用。只有两者相符合时，`ddl` 才算真正意义上的完成。

此时，发起者会释放 `table meta` 锁和分布式锁，并且删除 `universe/dble/{cluster-id}/ddl/{schema.table}` kv 树。

### kill ddl\_lock

当前指令只会释放对应 `ddl` 在当前节点中所持有的元数据锁，但是不会影响执行该 `ddl` 的线程的状态。若出现上述问题，可以参考 上面如何判断其他节点都响应完成的方式找出哪些没有响应，在这些节点和发起节点上执行当前命令，并手动删除 `universe/dble/{cluster-id}/ddl/{schema.table}` kv 树。

### 注意事项

1. 需要注意 `kill` 指令的执行顺序：先在各个从节点上执行该指令，最后再在主节点上执行。
2. 如果不是特殊因素，此命令请不要随意使用。

## 2.23 外部高可用联动

- 2.23.1 外部后端MySQL-HA连接
- 2.23.2 命令的使用说明
- 2.23.3 命令的实现细节
- 2.23.4 简单的HA交互使用案例

## 2.23.1 外部后端MySQL-HA连接

dble在版本2.19.0开始提供对应的外部接口，基于以下几个基础要点进行设计：

- 外部HA组件能够通过一定的方式和dble进行信息的主动沟通（包括脚本，网络通信，命令等方式）
- 外部HA组件能够做出对与MySQL可用性的明确判断（包括哪些节点停流量，哪些节点主从切换，哪些节点启动流量）

在基于以上的两个条件的情况下，基于提供外部接口通用性以及易用性的考虑，决定通过管理端用户命令的方式添加对应的交互接口，具体的提供以下的三个交互接口：

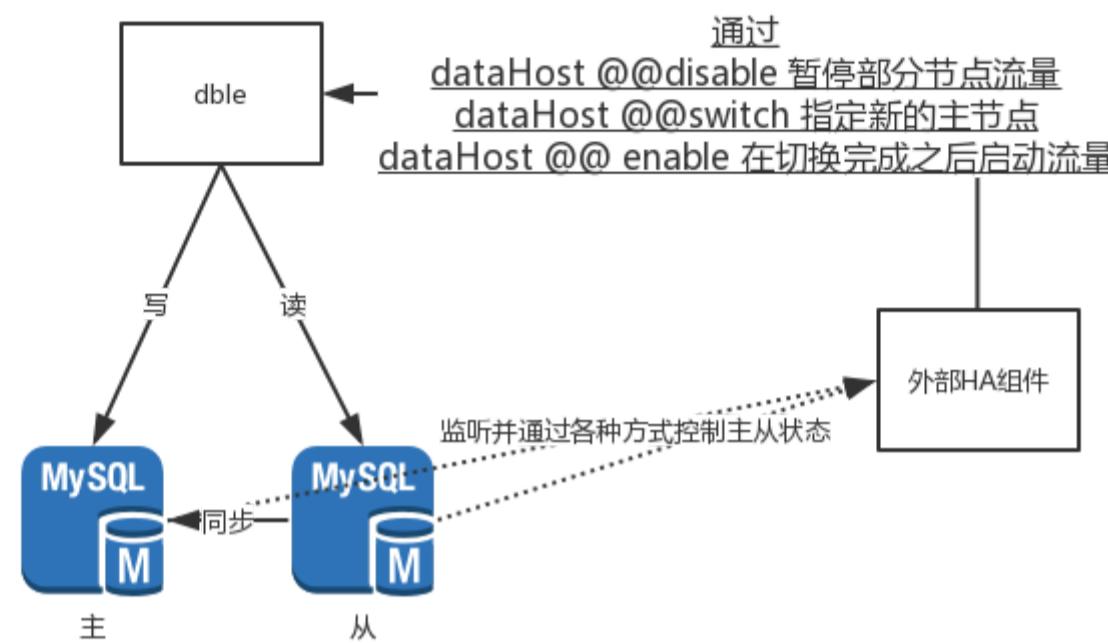
- 停止流量命令(dataHost @@disable)
- 启动流量命令(dataHost @@enable)
- 主从切换命令(dataHost @@switch)

其次为了方便进行管理和观测，提供了一个当前切换状态的查询命令：

- 切换状态查询命令(dataHost @@events)

外部HA组件可以通过 mysql命令、脚本调用mysql命令、或者程序的DB连接执行mysql命令的方式进行交互。

同时此功能设计考虑到zk集群用户的需求，当用户使用zk集群状态下的dble进行高可用切换时，会有集群的联动机制。外部HA组件只需要通知其中一个节点即可。



以下提供一个外部HA和dble交互设计图：

更加详细的情况请查询以下文档：  
[命令的使用说明](#)  
[命令的实现细节](#)  
[简单的HA交互使用案例](#)

## 2.23.2 dataHost命令的具体使用方法和解释

### 外部ha启用参数

当项目准备使用dble的外部ha联动方式来完成高可用切换或者同步的时候，需要提前进行如下配置

#### server.xml中system

```
<property name="useOuterHa">true</property>
```

当此参数配置时，对于dble的配置出现以下的限制条件：

- 单个dataHost下面只允许配置一个writeHost
- 基于dble内部心跳的高可用切换不可用

注意：本参数的调整需要重启dble服务

#### myid.properties

配置参数clusterHa = true

- 当此参数启用时，集群状态的dble将会在集群中同步自身的dataSource的状态
- 此配置在server.xml中useOuterHa参数为false时不生效
- 当useOuterHa参数为true但clusterHa不会true时，dble可以执行高可用切换的所有指令，但是其行为退化为单机dble，需要人工进行集群中多个dble的状态同步

注意：此参数的调整需要重启dble服务

### 关于dble中对于后端MySQL状态的解释

dataSource的状态“disabled/enable”仅表示dble层面对于具体每个MySQL后端节点是否允许有流量的标识，和具体的MySQL存活状态无关

### dataHost @@disable

命令细节：

```
dataHost @@disable name = 'dataHost_name' [node = 'dataSource_name']
```

- 其中的dataHost\_name指的是在schema.xml中配置的数据Host的名称，而dataSource\_name指的是在一个dataHost下面具体的writeHost/readHost的hostName
- 当此命令不指定node = '..'的内容时，默认将此dataHost下所有writeHost/readHost的状态置为disabled
- 被标记为disabled的节点无法提供正常的查询，即使对应的mysql真实的存活着
- 具体的dataHost的disabled状态可以通过管理端命令show @@datasource进行查询
- 当dble服务不在useOuterHa = true的模式的情况下，此命令无法执行，返回报错
- 若当前dble服务尚存连接被disabled的连接，在命令执行过程中会被全部关闭，包括正在新建过程中的连接，可能会导致少量的查询报错

### dataHost @@enable

命令细节：

```
dataHost @@enable name = 'dataHost_name' [node = 'dataSource_name']
```

- 其中的dataHost\_name指的是在schema.xml中配置的数据Host的名称，而dataSource\_name指的是在一个dataHost下面具体的writeHost/readHost的hostName
- 当此命令不指定node = '..'的内容时，默认使得整个dataHost下的所有writeHost/readHost的状态enable
- 当dble服务不在useOuterHa = true的模式的情况下，此命令无法执行，返回报错

### dataHost @@switch

命令细节： dataHost @@switch name = 'dataHost\_name' master = 'dataSource\_name'

- 其中的dataHost\_name指的是在schema.xml中配置的数据Host的名称，而dataSource\_name指的是在一个dataHost下面具体的writeHost/readHost的hostName
- 此命令name和master内容都为必填，在缺少任意元素的状态下会报错
- 当dble服务不在useOuterHa = true的模式的情况下，此命令无法执行，返回报错
- 此命令的作用会导致writeHost和readHost之间发生位置互换，被选中成为新master的节点将变成writeHost，并且其他所有dataSource变成readHost
- 此命令不会导致所有dataSource节点的disabled状态变化，但如果命令使得一个dataSource从writeHost退化成为readHost，此dataSource上的所有既有连接都会被关闭，以确保新的写请求不会被写入到错误的dataSource上面去，这可能导致一些前端连接的查询报错和事务失败

### 2.23.3 高可用联动命令的逻辑细节

#### 简述

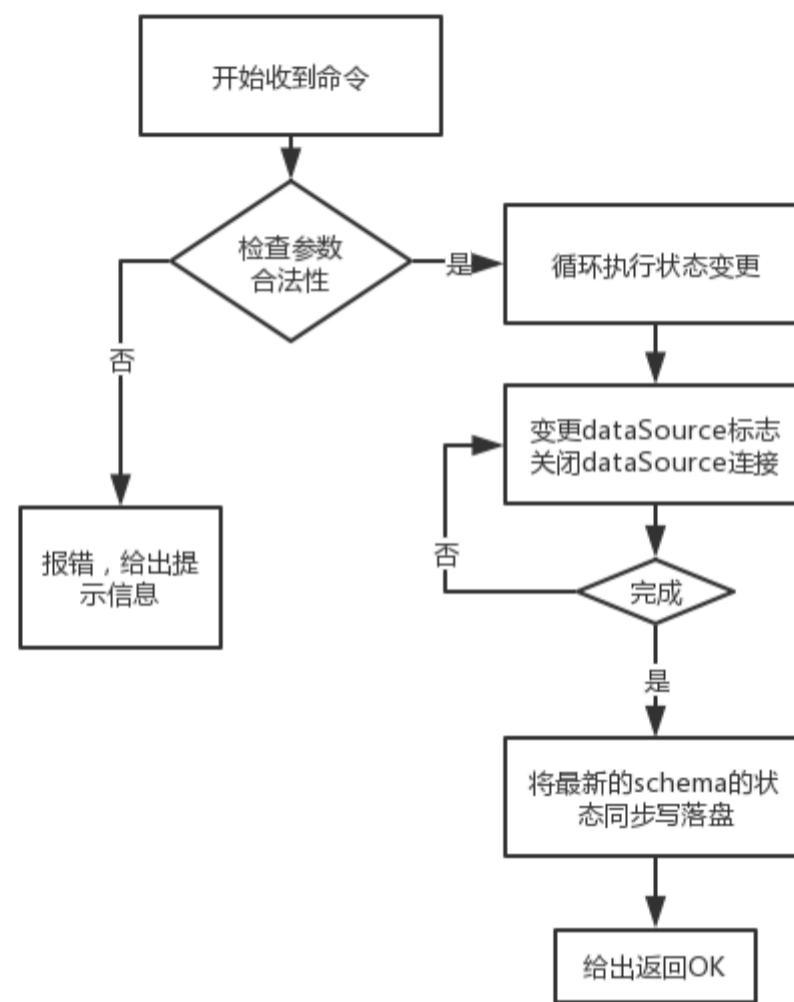
在dble高可用协同的几个接口中基本逻辑都是相似的，只是去更新dataHost的属性，只不过在细节上每个命令都有一些自身的特殊行为

- disable命令在执行过程中会断开当前所有连接，并且在zk集群状态下会要求其他节点同步响应
- enable命令没有任何附加行为
- switch命令会在切换过程中断开旧writeHost的所有连接

#### dataHost @@disable

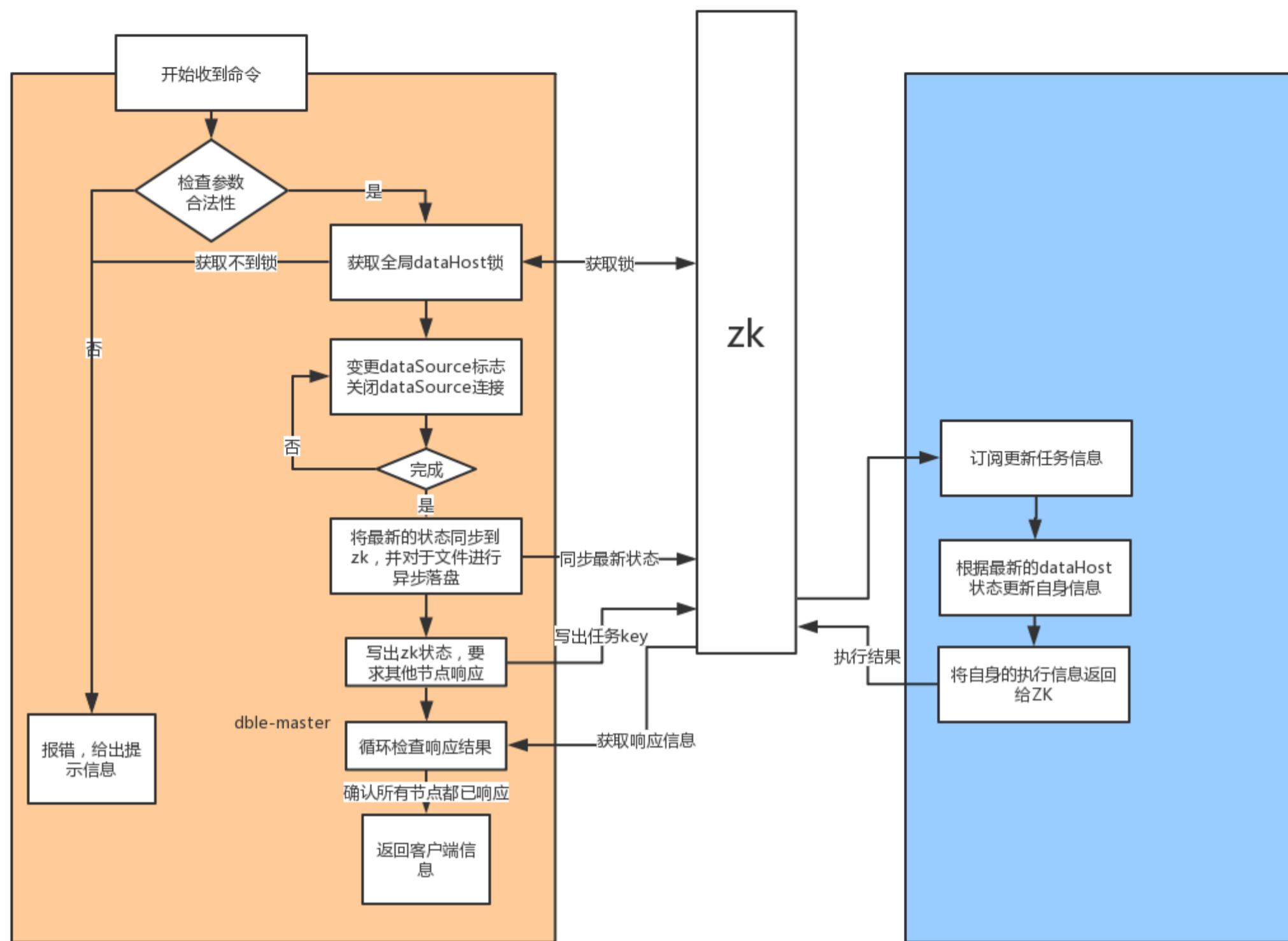
dble在单机情况下大致逻辑如下：

- 检查命令参数是否符合条件
- 更新dataSource的状态
- 关闭所有已经存在的连接
- 同步更新配置文件，将dataHost的最新状态落盘
- 返回OK信息



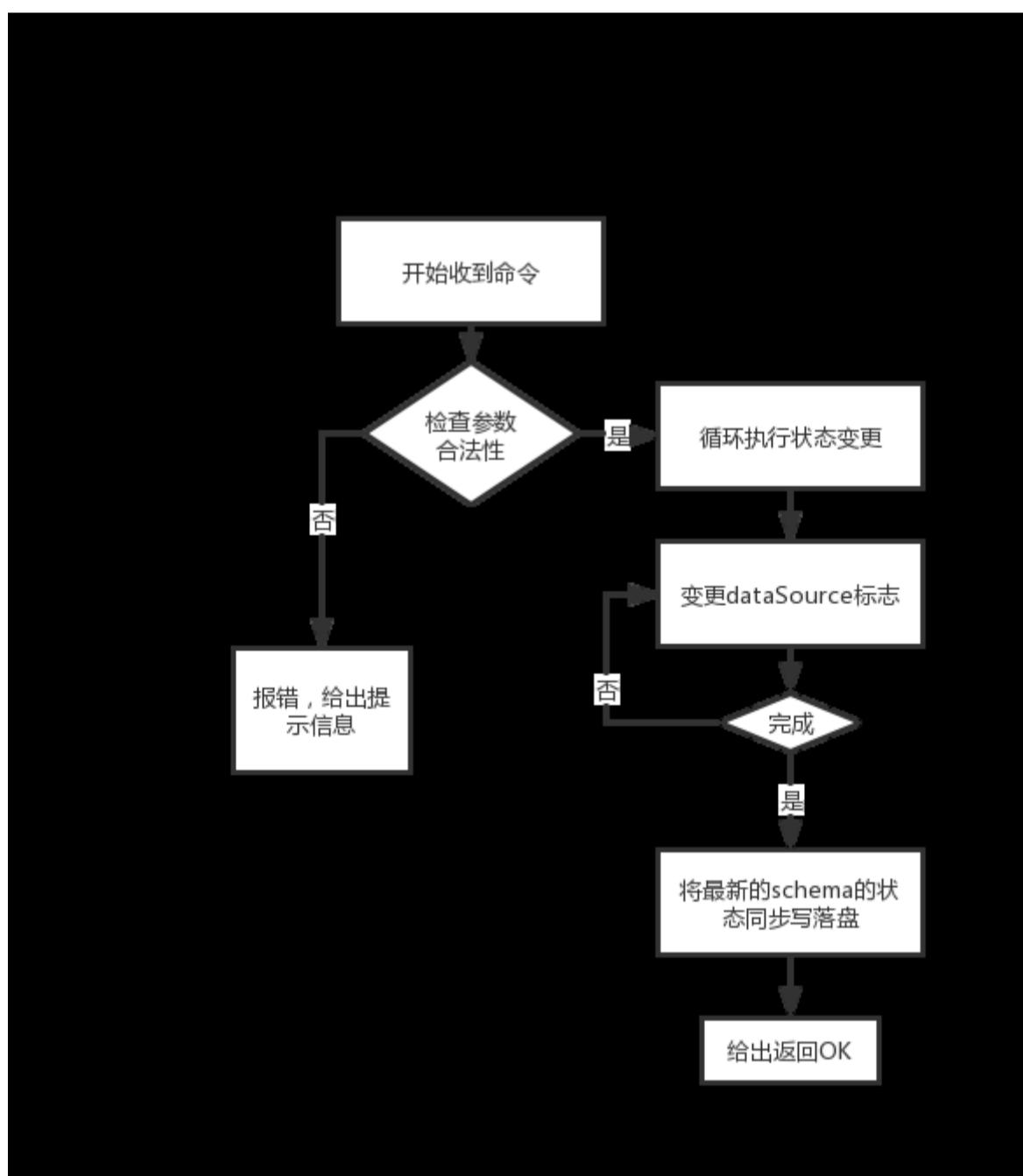
在集群状态下的逻辑有所不同：

- 检查命令参数是否符合条件
- 更新dataSource的状态
- 关闭所有已经存在的连接
- 异步更新配置文件，将dataHost的最新状态落盘
- 将最新的dataHost的状态信息同步到zk上
- 写出一个zk上面的key任务，等待其他节点响应
- 其他节点响应任务，更新自身写出响应结果
- 发起dble检查所有响应结果，等待所有节点响应完成
- 返回最终的执行结果

**dataHost @@enable**

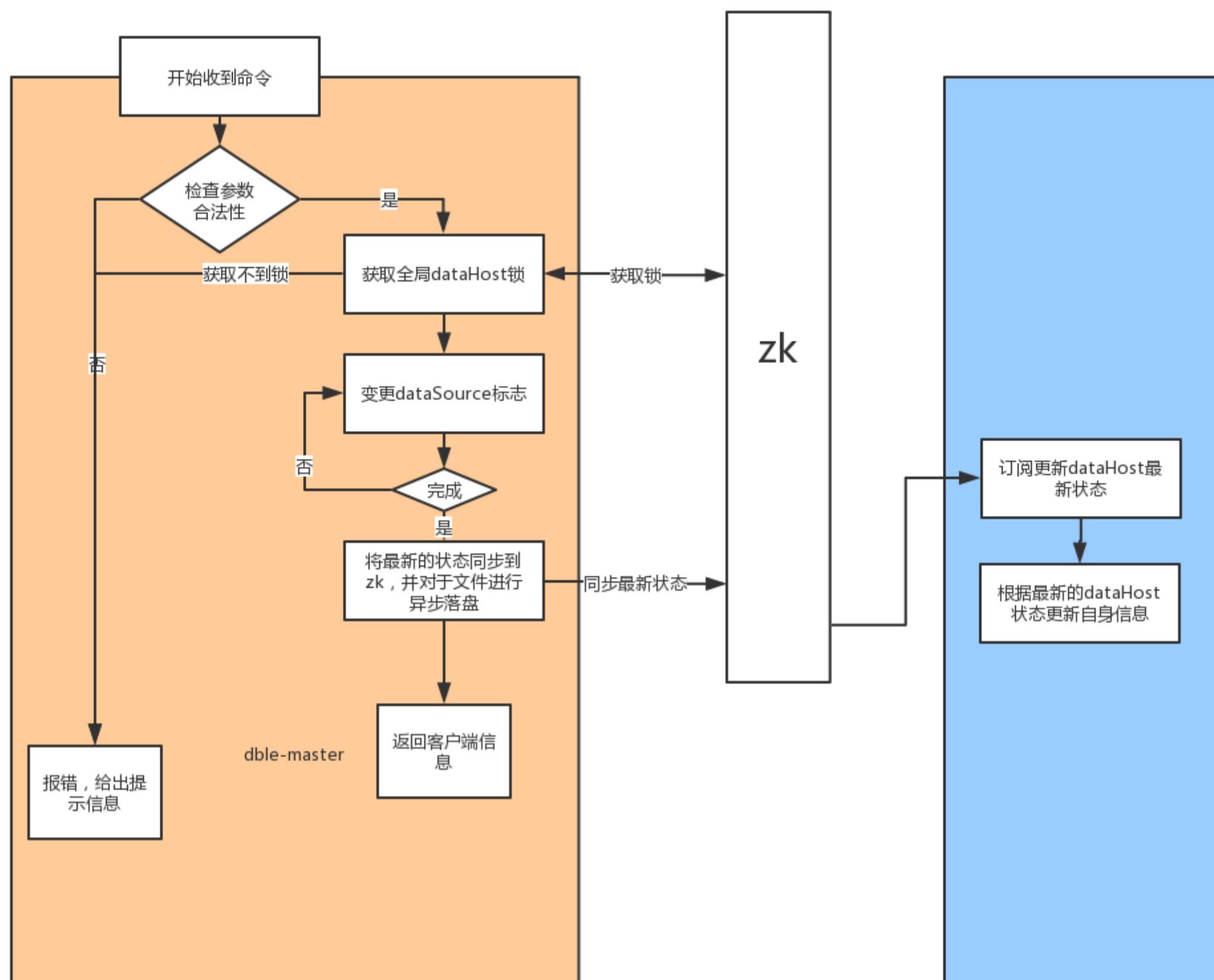
dble在单机情况下的逻辑大致如下:

- 检查命令参数是否符合条件
- 更新dataSource的状态
- 同步更新配置文件，将dataHost的最新状态落盘
- 返回OK信息



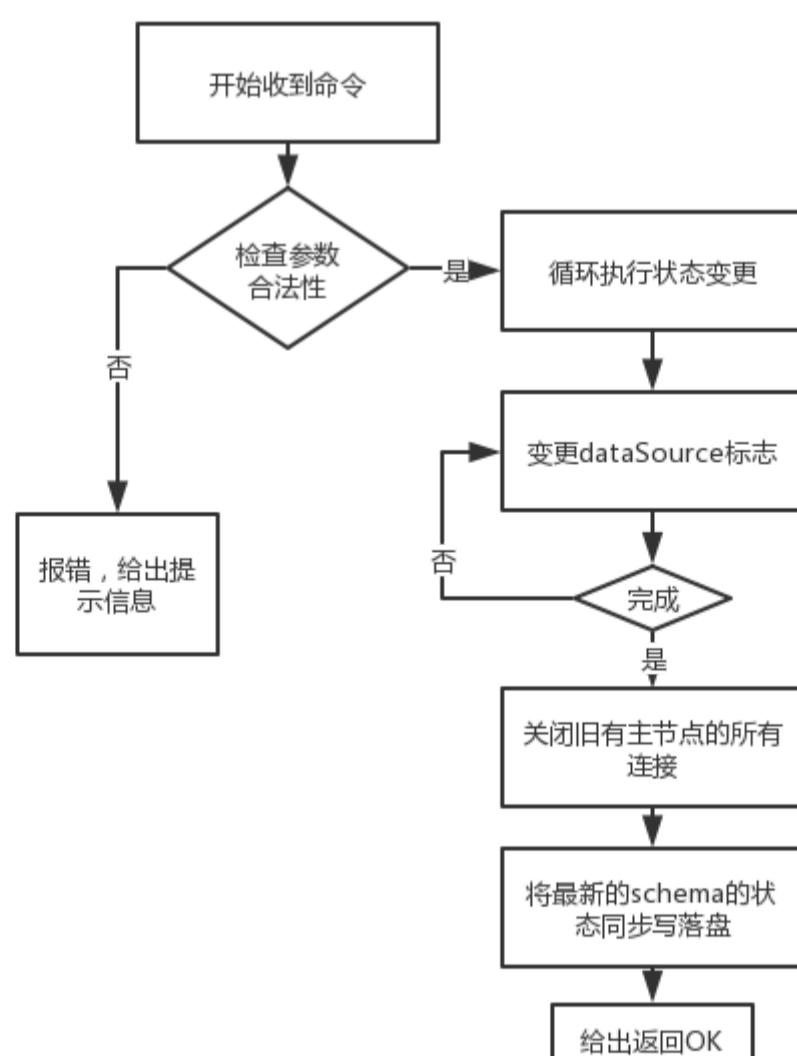
在ZK集群的情况下逻辑大致如下:

- 检查命令参数是否符合条件
- 更新dataSource的状态
- 关闭所有已经存在的连接
- 异步更新配置文件，将dataHost的最新状态落盘
- 将最新的dataHost的状态信息同步到ZK上
- 写出一个ZK上面的key任务，等待其他节点响应
- 其他节点响应任务，更新自身写出响应结果
- 发起dble检查所有响应结果，等待所有节点响应完成
- 返回最终的执行结果

**dataHost @@switch**

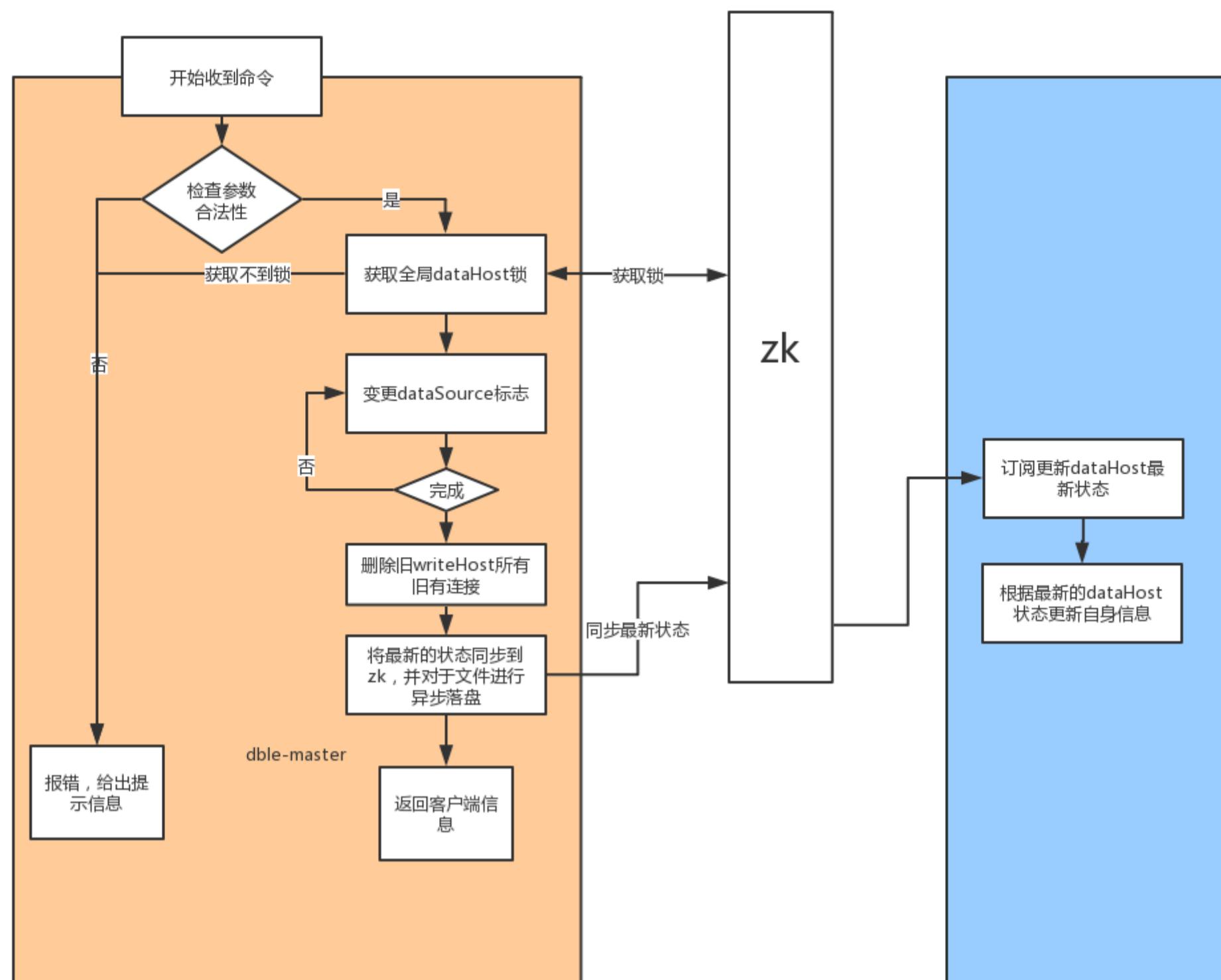
dble在单机情况下的逻辑大致如下:

- 检查命令参数是否符合条件
- 更新dataSource的状态
- 关闭旧writeHost的所有连接
- 同步更新配置文件，将dataHost的最新状态落盘
- 返回OK信息



在ZK集群的状况下逻辑大致如下:

- 检查命令参数是否符合条件
- 更新dataSource的状态
- 关闭所有已经存在的连接
- 异步更新配置文件，将dataHost的最新状态落盘
- 将最新的dataHost的状态信息同步到zk上
- 写出一个zk上面的key任务，等待其他节点响应
- 其他节点响应任务，更新自身写出响应结果
- 发起dble检查所有响应结果，等待所有节点响应完成
- 返回最终的执行结果

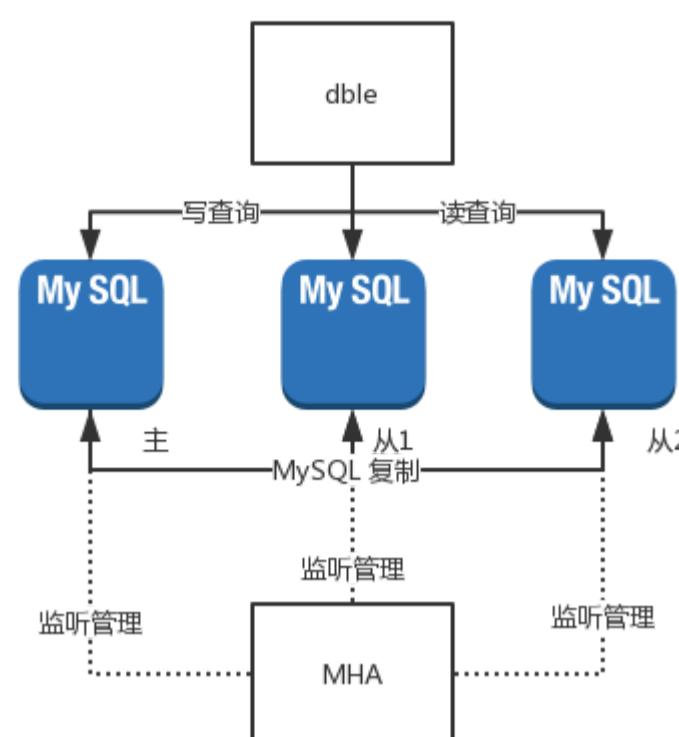


## 2.23.4 mha-dble高可用联动实例

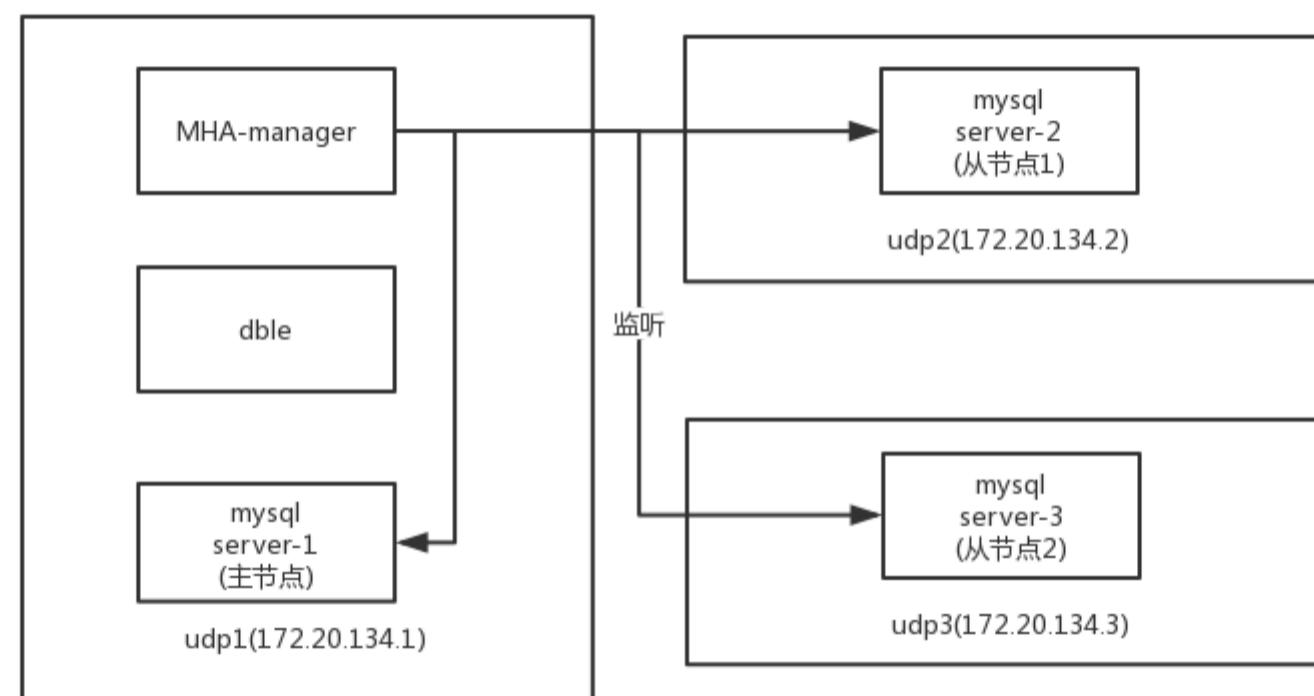
### 基础介绍

#### 服务结构图

整体的服务结构如下图所示，dble将一个一主两从的MySQL视作一个读写分离的dataHost，然后通过MHA对此MySQL进行复制的管理：



下图是对于各个服务在每个docker容器上面的分布图：



### 实验过程

- 搭建一个MySQL主从(一主两从)的基础环境
- 搭建MHA监管MySQL复制组的环境
- 搭建dble环境使用这MySQL组进行启动
- 手动kill MySQL主实例让MHA托管的高可用组发生切换
- 验证dble和MySQL的高可用切换过程

#### 实验目的

- 本实验原则上存在两个目的：
- 验证dble在2.19.09.0版本提供的高可用接口功能
  - 搭建简单的MHA和dble高可用的环境范例

### 前期环境准备

前期准备环境：  
三个docker容器，安装并部署mysql实例，并让三个mysql之间形成一主两从的复制关系(操作略)

### 操作过程

#### mha环境搭建

- 在每个docker容器中开启ssh服务

```
/usr/sbin/sshd -D
```

- 给每个容器之间创建免密登录,去到所有容器中执行,包括需要给自身容器创建

```
ssh-keygen -t rsa
ssh-copy-id -i ~/.ssh/id_rsa.pub root@other1
ssh-copy-id -i ~/.ssh/id_rsa.pub root@other2
ssh-copy-id -i ~/.ssh/id_rsa.pub root@self
```

- 确保/usr/bin/mysqlbinlog和/usr/bin/mysql文件存在,若不存在,请使用yum安装mysql即可
- 创建所有节点上的MHA工作目录

```
mkdir /etc/masterha/app1 -p
mkdir /var/log/masterha/app2 -p
mkdir /var/log/masterha/app1 -p
```

- 给所有mysql节点统一复制权限

```
grant all on *.* to root@'%' identified by '123456' with grant option;
grant replication slave on *.* to repl@'%' identified by 'repl';
```

- 下载并安装MHA的rpm包

```
从项目的文档页面能够直接找到对应的下载地址
https://github.com/yosinorim/mha4mysql-manager/wiki/Downloads
下载一个MHA Manager 0.56 rpm RHEL6
以及一个MHA Node 0.56 rpm RHEL6
并将下载完成的rpm包上传到对应的容器中
注意使用yum localinstall 命令对于rpm包进行安装，这样在安装过程中可以由yum来进行依赖的处理
Node的包需要先于Manager的包安装，不然会有依赖方面的报错
在所有容器上安装完成这两个rpm
当yum无法正确安装rpm依赖，请更新yum源到阿里云的yum
```

## 切换脚本准备

mha这个项目是使用perl写的，同时mha这个项目给用户提供了大量的自定义插入接口，可以允许用户在一些关键步骤进行自定义的操作响应以及通知，在自动高可用切换过程中，mha提供了一个切换脚本的入口

master\_ip\_failover\_script脚本配置

详细的参数介绍可以参考官网的解释[https://github.com/yosinorim/mha4mysql-manager/wiki/Parameters#master\\_ip\\_failover\\_script](https://github.com/yosinorim/mha4mysql-manager/wiki/Parameters#master_ip_failover_script)

大体上来说这个脚本的三种命令会在以下几个时间节点得到调用

- 启动检查HA的状态，调用脚本master\_ip\_failover的status命令
- 发现MySQL master节点失效，调用master\_ip\_failover的stop命令，并输入失效master的信息
- 在新选择的master节点补偿数据完成，并进行写恢复设置read\_only=0，调用master\_ip\_failover脚本的start命令

在本次实验中，通过自定义的脚本行为来进行mha和dble之间的事件交互，所以我们对于原有的样例脚本进行以下的修改

- 原有MySQL master节点失效时，将dataHost上面对应写节点的状态修改成disable，这要求在stop命令中对于dble发送dataHost @@disable命令
- 当新的MySQL master被选出来并且上线时，在脚本start阶段通过调用dataHost @@switch命令将新的master节点切换成为dataHost中的主节点

按照以上的逻辑对于修改完成的master\_ip\_failover脚本如下(注意给创建的脚步提供执行权限)

```

#!/usr/bin/env perl

# Copyright (C) 2011 DeNA Co.,Ltd.
#
# This program is free software; you can redistribute it and/or modify
# it under the terms of the GNU General Public License as published by
# the Free Software Foundation; either version 2 of the License, or
# (at your option) any later version.
#
# This program is distributed in the hope that it will be useful,
# but WITHOUT ANY WARRANTY; without even the implied warranty of
# MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
# GNU General Public License for more details.
#
# You should have received a copy of the GNU General Public License
# along with this program; if not, write to the Free Software
# Foundation, Inc.,
# 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301 USA

## Note: This is a sample script and is not complete. Modify the script based on your environment.

use strict;
use warnings FATAL => 'all';

use Getopt::Long;
use MHA::DBHelper;

my (
    $command,           $ssh_user,           $orig_master_host,
    $orig_master_ip,   $orig_master_port, $new_master_host,
    $new_master_ip,    $new_master_port,  $new_master_user,
    $new_master_password
);
GetOptions(
    'command=s'         => \$command,
    'ssh_user=s'        => \$ssh_user,
    'orig_master_host=s'=> \$orig_master_host,
    'orig_master_ip=s'  => \$orig_master_ip,
    'orig_master_port=i'=> \$orig_master_port,
    'new_master_host=s'  => \$new_master_host,
    'new_master_ip=s'   => \$new_master_ip,
    'new_master_port=i' => \$new_master_port,
    'new_master_user=s' => \$new_master_user,
    'new_master_password=s'=> \$new_master_password,
);
exit &main();

sub main {
    if ( $command eq "stop" || $command eq "stopssh" ) {

        # $orig_master_host, $orig_master_ip, $orig_master_port are passed.
        # If you manage master ip address at global catalog database,
        # invalidate orig_master_ip here.
        my $exit_code = 1;
        eval {

            # dataHost @@disable name = "dataHost1" node='$orig_master_host'
            # 调用对应的disable命令，使得部分节点不可写
            $orig_master_host =>tr './';
            system "mysql -P9066 -u man1 -p654321 -h 172.20.134.1 -e \"dataHost @@disable name = 'dataHost1' node='".$orig_master_host."'\"";
            $exit_code = 0;
        };
        if ($@) {
            warn "Got Error: $@\n";
            exit $exit_code;
        }
        exit $exit_code;
    }
    elsif ( $command eq "start" ) {

        # all arguments are passed.
        # If you manage master ip address at global catalog database,
        # activate new_master_ip here.
        # You can also grant write access (create user, set read_only=0, etc) here.
        my $exit_code = 10;
        eval {
            my $new_master_handler = new MHA::DBHelper();

            # args: hostname, port, user, password, raise_error_or_not
            $new_master_handler->connect( $new_master_ip, $new_master_port,
                $new_master_user, $new_master_password, 1 );

            ## Set read_only=0 on the new master
            $new_master_handler->disable_log_bin_local();
            print "Set read_only=0 on the new master.\n";
            $new_master_handler->disable_read_only();

            ## Creating an app user on the new master
            print "Creating app user on the new master..\n";
            $new_master_handler->enable_log_bin_local();
            $new_master_handler->disconnect();

            ## try to switch the dataHost master into new master
            ## 调用dataHost switch的命令，将新的new_master_host节点提升
            $new_master_host =>tr './';
            system "mysql -P9066 -u man1 -p654321 -h 172.20.134.1 -e \"dataHost @@switch name = 'dataHost1' master='".$new_master_host."'\"";
            $exit_code = 0;
        };
        if ($@) {
            warn $@;

            # If you want to continue failover, exit 10.
            exit $exit_code;
        }
        exit $exit_code;
    }
    elsif ( $command eq "status" ) {
        # test for start command
        exit 0;
    }
    else {
        &usage();
        exit 1;
    }
}

sub usage {
    print
"Usage: master_ip_failover --command=start|stop|stopssh|status --orig_master_host=host --orig_master_ip=ip --orig_master_port=port --new_master_host=host --new_master_ip=ip --new_master_port=port\n";
}

```

将脚本存放到指定目录/etc/masterha/app1目录下

并且创建MHA最终使用的app1.conf配置文件在/etc/masterha/app1目录，具体创建的配置文件内容如下所示：

```
#mha manager工作目录
manager_workdir = /var/log/masterha/app1
manager_log = /var/log/masterha/app1/app1.log
remote_workdir = /var/log/masterha/app2
master_ip_failover_script=/etc/masterha/app1/master_ip_failover
# master_ip_online_change_script=/etc/masterha/app1/master_ip_online_change
# MySQL管理帐号和密码
user=root
password=123456
# 系统ssh用户
ssh_user=root

# 复制帐号和密码
repl_user=repl
repl_password= repl

# 监控间隔(秒)
ping_interval=1
manager_log=/var/log/masterha/app1/manager.log

[server1]
hostname=172.20.134.1
master_binlog_dir = /opt/3306/
port=3306

[server2]
# 每个机器上面的mysql实例的信息
hostname=172.20.134.2
master_binlog_dir = /opt/3306/
candidate_master=1
check_repl_delay=0
port=3306

[server3]
# 每个机器上面的mysql实例的信息
hostname=172.20.134.3
master_binlog_dir = /opt/3306/
candidate_master=1
check_repl_delay=0
port=3306
```

最终通过以下命令对于MHA的监听线程进行启动

```
nohup masterha_manager --conf=/etc/masterha/app1/app1.conf >> /var/log/masterha/app1/manager.log 2>&1 &
```

## dble配置

配置dble在这里是比较简单的内容，我们仅使用默认的配置即可，从项目的release页面下载最新的2.19.09.0安装包，解压并将conf文件下的rule\_temp.xml, server\_temp.xml, schema\_temp.xml三个配置文件分别重命名成为rule.xml, server.xml, schema.xml  
并按照当前的配置需求对于其中的server.xml以及schema.xml进行如下的配置调整

server.xml

```
<property name="useOuterHa">true</property>

schema.xml

<dataNode name="dn1" dataHost="dataHost1" database="db_1"/>
<dataNode name="dn2" dataHost="dataHost1" database="db_2"/>
<dataNode name="dn3" dataHost="dataHost1" database="db_3"/>
<dataNode name="dn4" dataHost="dataHost1" database="db_4"/>
<dataNode name="dn5" dataHost="dataHost1" database="db_5"/>
<dataNode name="dn6" dataHost="dataHost1" database="db_6"/>
<dataHost balance="0" maxCon="1000" minCon="10" name="dataHost1" slaveThreshold="100">
    <heartbeat>show slave status</heartbeat>
    <writeHost host="172_20_134_1" url="172.20.134.1:3306" password="123456" user="root" disabled="false" id="udp-1" weight="0"/>
        <readHost host="172_20_134_3" url="172.20.134.3:3306" password="123456" user="root" disabled="false" id="udp-3" weight="0"/>
        <readHost host="172_20_134_2" url="172.20.134.2:3306" password="123456" user="root" disabled="false" id="udp-2" weight="0"/>
    </writeHost>
</dataHost>
```

启动dble，并通过管理用户man1查看基础状态下的dble后端节点状态

```
MySQL [(none)]> show @@datasource;
+-----+-----+-----+-----+-----+-----+-----+-----+
| DATAHOST | NAME   | HOST   | PORT  | W/R   | ACTIVE | IDLE  | SIZE   | EXECUTE | READ_LOAD | WRITE_LOAD | DISABLED |
+-----+-----+-----+-----+-----+-----+-----+-----+
| dataHost1 | 172_20_134_1 | 172.20.134.1 | 3306 | W    | 1 | 0 | 1000 | 1 | 0 | 0 | false |
| dataHost1 | 172_20_134_3 | 172.20.134.3 | 3306 | R    | 1 | 0 | 1000 | 0 | 0 | 0 | false |
| dataHost1 | 172_20_134_2 | 172.20.134.2 | 3306 | R    | 1 | 0 | 1000 | 0 | 0 | 0 | false |
+-----+-----+-----+-----+-----+-----+-----+-----+
```

可见在初始状态下，节点的使用属性和配置文件一致，172.20.134.1节点作为写节点，并且所有节点的使用状态都是“可用”

## 最终效果

- 通过ps命令在172.20.134.1容器中找到对应的mysqld进程
- 通过kill -9 命令对于172.20.134.1容器中的mysqld进程进行关闭
- 重新通过命令检查dble中此时的后端节点状态，发现如下的执行结果
 

```
MySQL [(none)]> show @@datasource;
+-----+-----+-----+-----+-----+-----+-----+-----+
| DATAHOST | NAME   | HOST   | PORT  | W/R   | ACTIVE | IDLE  | SIZE   | EXECUTE | READ_LOAD | WRITE_LOAD | DISABLED |
+-----+-----+-----+-----+-----+-----+-----+-----+
+ dataHost1 | 172_20_134_2 | 172.20.134.2 | 3306 | W | 1 | 0 | 1000 | 0 | 0 | 0 | false || dataHost1 | 172_20_134_3 | 172.20.134.3 | 3306 | R | 1 | 0 | 1000 | 0 | 0 | 0 | false || dataHost1 | 172_20_134_1 | 172.20.134.1 | 3306 | R | 0 | 0 | 1000 | 0 | 0 | 0 | true | +-----+-----+-----+-----+-----+-----+-----+
```

+ 检查当前dble的配置文件，发现配置文件中的后端节点位置发生扭转，172\_20\_134\_2变成新的写节点

show slave status

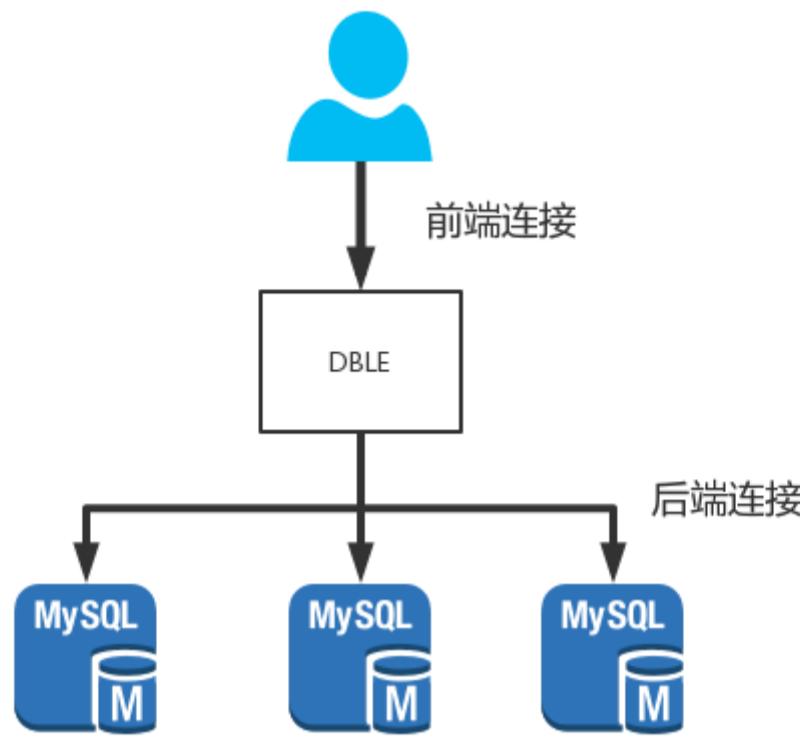
```
+ 检查存在于172.20.134.3上mysql从的状态，发现其复制指向已经发生切换，切换为新主172.20.134.2，与dble中最新的配置文件相符
```

```
MySQL [(none)]> show slave status\G * Slave_IO_State: Waiting for master to send event Master_Host: 172.20.134.2 Master_User: repl Master_Port: 3306 Connect_Retry: 60 Master_Log_File: mysql-bin.000001 Read_Master_Log_Pos: 2113 Relay_Log_File: udp3-relay-bin.000002 Relay_Log_Pos: 320 Relay_Master_Log_File: mysql-bin.000001 Slave_IO_Running: Yes Slave_SQL_Running: Yes .... ``
```

## 2.24 超时（连接/执行）控制

### 概述

在dble中连接分为两种，前端连接，分别对应客户端连接dble的网络连接和dble中连接mysql的网络连接。大致的情况如下图所示：



由于网络TCP连接的特性，dble中需要对于这两种连接都有超时控制，具体在dble中将连接的超时控制分为**空闲超时**和**执行超时**两种情况。当超时发生的时候，dble会通过切断超时连接的方法进行控制。

### 超时的检查和实现

在配置文件server.xml中存在三个配置项

- sqlExecuteTimeout(后端连接执行超时)
- idleTimeout (前端连接空闲超时)
- processorCheckPeriod (超时检查周期)

以下描述具体超时检查的实现逻辑

- dble按照配置的processorCheckPeriod作为周期去周期性检查所有的连接
- 检查所有后端连接，若满足以下条件就关闭后端连接
  - 此后端连接已经被某个前端连接借走（正在执行或持有）
  - 此后端连接距离上一次收发包超过sqlExecuteTimeout
  - 此后端连接没有在执行DDL或者执行XA事务
- 检查所有前端连接，若满足以下条件就关闭前端连接
  - 此前端连接正在执行XA事务，并不处于commit失败补偿/rollback失败补偿
  - 此前端连接距离上一次收发包超过idleTimeout

综上总结为超时关闭连接的逻辑可以描述为：

- 后端连接非DDL,XA事务的情况下被前端连接借走，并且超过sqlExecuteTimeout没有收发包（包括普通事务执行间隔，后端连接被长期持有的时间超时）
- 前端连接非XA事务特殊阶段，超过idleTimeout没有收发包（包括在load data大文件过程中发生的长时间空档）

### SQL执行超时注意事项

- 部分语言或者框架的连接池会长期持有连接，并且没有设定或者定时发送网络包，可能会导致空闲的连接自动断开，然后应用在取用到对应断开连接的时候可能会报错
- 当开启一个事务但是长期不执行SQL会导致后端连接被判断定为执行超时，当应用再次在连接上执行内容的时候，会给出后端连接已关闭的报错信息
- 当前端连接执行一个大文件内容load data的时候，由于后端可能执行较慢，前端连接存在超过idleTimeout导致连接关闭的情况，当有类似数据导入计划的时候，考虑暂时放宽idleTimeout的限制
- 连接的超时检查和关闭是通过一个循环来实现的，所以并不是一个实时检查的值，最坏的情况对于超时的检查会迟到一个processorCheckPeriod周，当应用对于SQL超时需要一个严格的时间设定时，请谨慎使用

## 2.25 dble流量控制

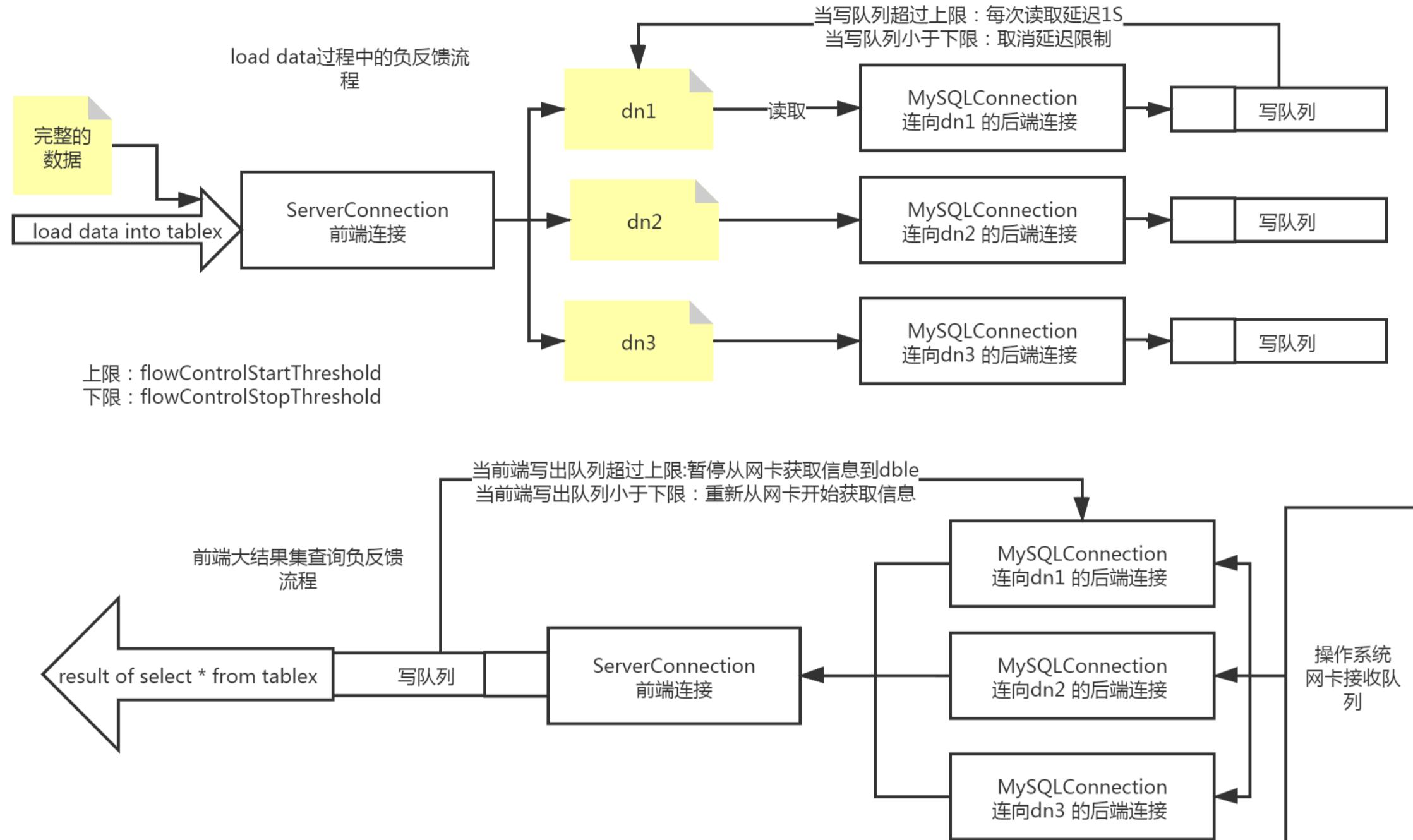
### 背景

在之前的dble版本中，当进行大文件load以及大结果集查询的过程中，都有可能由于数据的发送不及时造成数据在dble内存中堆积，当条件足够的时候甚至有可能造成dble服务的OOM，进而影响服务的稳定运行。在本版2.20.04中由社区开发者@ssxlulu提供了对于这部分流量控制的实现，通过连接级别的写队列长度，进行数据加载/获取的负反馈调节，从而实现在数据load和大结果集查询过程中的内存使用情况稳定。具体的issue详情请见[issue](#)

### 原理

dble中的流量控制通过连接级别的写出队列进行负反馈调节，具体在生效的时候分成两种具体的形式：

- load data过程中后端连接写出队列过长负反馈
- Query过程中前端连接写出队列过长负反馈 下图中分别描述了两种负反馈具体实现的逻辑：



注：流量控制功能的生效级别为连接级别，不同连接之间的队列长度不会互相影响

### 使用举例

本功能为默认关闭功能，需要在server.xml中进行相关配置进行显式开启，或者通过管理端的辅助命令进行实时的调整。在server.xml中使用下列参数使得功能开启并生效：

```
enableFlowControl(是否开启流量控制开关)
flowControlStartThreshold(写队列上限阈值, 写队列超限时开启流量控制)
flowControlStopThreshold(写队列下限阈值, 写队列低于阈值时取消流量控制)
```

附带三个管理端命令可以修改当前的流量控制参数：

- flow\_control @@show (展示当前流控配置参数信息)
- flow\_control @@set [enableFlowControl = true/false] [flowControlStart = ?] [flowControlEnd = ?] (修改配置参数)
- flow\_control @@list (展示当前正在被流量控制的连接)

### 3.语法兼容

- [3.1 DDL](#)
  - [3.1.1 DDL&Table Syntax](#)
  - [3.1.2 DDL&View Syntax](#)
  - [3.1.3 DDL&Index Syntax](#)
  - [3.1.4 DDL透传](#)
- [3.2 DML](#)
  - [3.2.1 INSERT](#)
  - [3.2.2 REPLACE](#)
  - [3.2.3 DELETE](#)
  - [3.2.4 UPDATE](#)
  - [3.2.5 SELECT](#)
  - [3.2.6 SELECT JOIN syntax](#)
  - [3.2.7 SELECT UNION Syntax](#)
  - [3.2.8 SELECT Subquery Syntax](#)
  - [3.2.9 LOAD DATA](#)
  - [3.2.10 不支持的语句](#)
- [3.3 Prepared SQL Syntax](#)
- [3.4 Transactional and Locking Statements](#)
  - [3.4.1 Lock&unlock](#)
  - [3.4.2 XA 事务语法](#)
  - [3.4.3 一般事务语法](#)
  - [3.4.4 SET TRANSACTION Syntax](#)
- [3.5 DAL](#)
  - [3.5.1 SET](#)
  - [3.5.2 SHOW](#)
  - [3.5.3 KILL](#)
- [3.6 存储过程支持方式](#)
- [3.7 Utility Statements](#)
- [3.8 Hint](#)
- [3.9 其他不支持语句](#)
- [3.10 函数与操作符支持列表\(alpha版本\)](#)
- [3.11 导入导出方式](#)

### 3.1 DDL

DDL 包含以下几部分内容

注意，当DDL执行时，涉及到的相同的表上的DML和DDL操作将会报错。

- [3.1.1 DDL&Table Syntax](#)
- [3.1.2 DDL&View Syntax](#)
- [3.1.3 DDL&Index Syntax](#)
- [3.1.4 DDL透传](#)

### 3.1.1 TABLE DDL

#### 3.1.1.1 CREATE TABLE Syntax

```

CREATE TABLE [IF NOT EXISTS] tbl_name
  (create_definition,...)
  [table_options]
  [partition_options]

create_definition:
  col_name column_definition

column_definition:
  data_type [NOT NULL | NULL] [DEFAULT default_value]
  [AUTO_INCREMENT] [UNIQUE [KEY] | [PRIMARY] KEY]
  [COMMENT 'string']

data_type:
  BIT[(length)]
  | TINYINT[(length)] [UNSIGNED] [ZEROFILL]
  | SMALLINT[(length)] [UNSIGNED] [ZEROFILL]
  | MEDIUMINT[(length)] [UNSIGNED] [ZEROFILL]
  | INT[(length)] [UNSIGNED] [ZEROFILL]
  | INTEGER[(length)] [UNSIGNED] [ZEROFILL]
  | BIGINT[(length)] [UNSIGNED] [ZEROFILL]
  | REAL[(length,decimals)] [UNSIGNED] [ZEROFILL]
  | DOUBLE[(length,decimals)] [UNSIGNED] [ZEROFILL]
  | FLOAT[(length,decimals)] [UNSIGNED] [ZEROFILL]
  | DECIMAL[(length[,decimals])] [UNSIGNED] [ZEROFILL]
  | NUMERIC[(length[,decimals])] [UNSIGNED] [ZEROFILL]
  | DATE
  | TIME[(fsp)]
  | TIMESTAMP[(fsp)]
  | DATETIME[(fsp)]
  | YEAR
  | CHAR[(length)]
  | VARCHAR(length)
  | BINARY[(length)]
  | VARBINARY(length)
  | TINYBLOB
  | BLOB
  | MEDIUMBLOB
  | LONGBLOB
  | TINYTEXT
  | TEXT
  | MEDIUMTEXT
  | LONGTEXT
  | ENUM(value1,value2,value3,...)

table_options:
  table_option [,] table_option] ...

table_option:
  ENGINE [=] engine_name
  | [DEFAULT] CHARACTER SET [=] charset_name
  | CHECKSUM [=] {0 | 1}
  | [DEFAULT] COLLATE [=] collation_name
  | COMMENT [=] 'string'
  | CONNECTION [=] 'connect_string'
  | KEY_BLOCK_SIZE [=] value
  | MAX_ROWS [=] value
  | MIN_ROWS [=] value
  | PASSWORD [=] 'string'
  | ROW_FORMAT [=] {DEFAULT|DYNAMIC|FIXED|COMPRESSED|REDUNDANT|COMPACT}
  | STATS_AUTO_RECALC [=] {DEFAULT|0|1}
  | STATS_PERSISTENT [=] {DEFAULT|0|1}

partition_options:
  {{LINEAR} HASH(expr)
  | PARTITION BY [linear] KEY (column_list)
  | RANGE(expr) | COLUMNS(column_list)
  | LIST(expr) | COLUMNS(column_list)
  }
  [(partition_definition [, partition_definition] ...)]

```

注意:

- engine\_name仅能为忽略大小写的“InnoDB”
- 不建议含有枚举类型的表作为分片表。比如表结构: `CREATE TABLE `test` (id enum('1','2','3') DEFAULT '1')`。因为此种表在插入id列时,既可以使用枚举值插入,也可以使用枚举值的下标,‘1’的下标是1,以此类推。若用户以枚举值进行分片,但是插入时确使用枚举值下标,因为db表不会将下标转换为枚举值,所以分片会出现问题,详细可参考[issue](https://github.com/actiontech/db表/issues/816) : <https://github.com/actiontech/db表/issues/816>。

例:

```

create table if not exists test(
  id bigint primary key AUTO_INCREMENT,
  col1 int not null default 5,
  col2 int null COMMENT 'info for col1',
  col3 varchar(20) not null,
  col4 varchar(20) unique key
);

create table test(
  id int primary key,
  col_bit    BIT(1),
  col_tinyint TINYINT(2) UNSIGNED ZEROFILL,
  col_smallint SMALLINT(3) UNSIGNED ZEROFILL,
  col_mediumint MEDIUMINT(4) UNSIGNED ZEROFILL,
  col_int INT(5) UNSIGNED ZEROFILL,
  col_integer INTEGER(6) UNSIGNED ZEROFILL,
  col_bigint BIGINT(7) UNSIGNED ZEROFILL,
  col_real REAL(8,1) UNSIGNED ZEROFILL,
  col_double DOUBLE(9,2) UNSIGNED ZEROFILL,
  col_float FLOAT(10,3) UNSIGNED ZEROFILL,
  col_decimal DECIMAL(11,4) UNSIGNED ZEROFILL,
  col_numeric NUMERIC(12,5) UNSIGNED ZEROFILL,
  col_date DATE,
  col_time TIME(3),
  col_timestamp TIMESTAMP(4),
  col_datetime DATETIME(5),
  col_year YEAR,
  col_char CHAR(10) ,
  col_varchar VARCHAR(20) ,
  col_binary BINARY(30),
  col_varbinary VARBINARY(40),
  col_tinyblob TINYBLOB,
  col_blob BLOB,
  col_mediumblob MEDIUMBLOB,
  col_longblob LONGBLOB,
  col_tinytext TINYTEXT ,
  col_text TEXT ,
  col_mediumtext MEDIUMTEXT ,
  col_longtext LONGTEXT ,
  col_enum ENUM('a','b','c')
);

```

或者

```

create table test(
    id int primary key,
    col1 varchar(20)
)ENGINE = innodb
AVG_ROW_LENGTH = 20
DEFAULT CHARACTER SET = utf8
CHECKSUM = 1
DEFAULT COLLATE = utf8_general_ci
COMMENT = 'info of table test'
CONNECTION = '111111'
DELAY_KEY_WRITE = 1
INSERT_METHOD = LAST
KEY_BLOCK_SIZE = 65536
MAX_ROWS = 3
MIN_ROWS = 2
PACK_KEYS = 1
ROW_FORMAT = DEFAULT;

```

### 3.1.1.2 ALTER TABLE Syntax

```

ALTER [IGNORE] TABLE tbl_name
    [alter_specification [, alter_specification] ...]

alter_specification:
    | ADD [COLUMN] col_name column_definition
        [FIRST | AFTER col_name]
    | ADD [COLUMN] (col_name column_definition,...)
    | CHANGE [COLUMN] old_col_name new_col_name column_definition
        [FIRST|AFTER col_name]
    | MODIFY [COLUMN] col_name column_definition
        [FIRST | AFTER col_name]
    | DROP [COLUMN] col_name
    | ADD [INDEX|KEY] [index_name] (index_col_name,...)
    | DROP {INDEX|KEY} index_name
    | ADD PRIMARY KEY (index_col_name,...)
    | DROP PRIMARY KEY

```

例:

```

alter table test add column col5 int not null default 1 first,add column col6 int after col4;
alter table test change column col1 col1_new int after col3;
alter table test modify column col1_new varchar(20) after id;
alter table test drop column col6;
alter table test add key idx_col4(col4);
alter table test drop key idx_col4;
alter table test drop primary key;
alter table test add primary key (id);

```

### 3.1.1.3 DROP TABLE Syntax

```

DROP TABLE [IF EXISTS]
    tbl_name [, tbl_name] ...
    [RESTRICT | CASCADE]

```

例:

```

drop table if exists test cascade;
drop table test restrict;

```

### 3.1.1.4 TRUNCATE TABLE Syntax

```

TRUNCATE [TABLE] tbl_name

```

例:

```

truncate table test;

```

### 3.1.2 VIEW DDL

Syntax

**create view :**

```
CREATE [OR REPLACE] VIEW
    view_name [(column_list)]
    AS select_statement
```

**alter view :**

```
ALTER VIEW
    view_name [(column_list)]
    AS select_statement
```

**drop view:**

```
DROP VIEW [IF EXISTS] view_name [, view_name]
```

**show create view :**

```
SHOW CREATE VIEW view_name;
```

### 3.1.3 INDEX DDL

#### 3.1.3.1 CREATE INDEX Syntax

```
CREATE [UNIQUE|FULLTEXT] INDEX index_name  
[index_type]  
ON tbl_name (index_col_name,...)  
  
index_col_name:  
    col_name [(length)] [ASC | DESC]  
  
index_type:  
    USING {BTREE | HASH}
```

例:

```
create unique index idx1 using btree on test(col1);  
create index idx2 using hash on test(col2);  
create fulltext index idx3 on test(col4);  
create fulltext index idx4 on test(col4(10));
```

#### 3.1.3.2 DROP INDEX Syntax

```
DROP INDEX index_name ON tbl_name
```

例:

```
drop index idx1 on test;
```

### 3.1.4 DDL透传

除了前面章节中的DDL语句，**dble**在非注解方式方式下不支持mysql中的其它DDL语句(eg. ALTER EVENT)。但通过如下的注解方式，**dble**支持所有的mysql中的DDL语句：

```
/*!dble:sql=select ... from tbx where id=M*/ ddl statement;
```

其中，**tbx**为分区表，**id**为分区列，**M**为分区列的某个值。

例如：

```
MySQL [TESTDB]> /*!dble:sql=select * from a_test where id=2*/CREATE PROCEDURE account_count()
BEGIN    SELECT 'Number of accounts:', COUNT(*) FROM mysql.user;
END//
```

### 3.1.5 DDL&Database\_Syntax

## 3.2 DML

DML 包含以下内容

- [3.2.1 INSERT](#)
- [3.2.2 REPLACE](#)
- [3.2.3 DELETE](#)
- [3.2.4 UPDATE](#)
- [3.2.5 SELECT](#)
- [3.2.6 SELECT JOIN syntax](#)
- [3.2.7 SELECT UNION Syntax](#)
- [3.2.8 SELECT Subquery Syntax](#)
- [3.2.9 LOAD DATA](#)
- [3.2.10 不支持的语句](#)

### 3.2.1 INSERT

#### 3.2.1.1 Syntax

```
INSERT [INTO] tbl_name
[(col_name,...)]
{VALUES | VALUE} ({expr },...),(...),...
[ ON DUPLICATE KEY UPDATE
  col_name=expr
  [, col_name=expr] ... ]
OR
INSERT [INTO] tbl_name
SET col_name=(expr | DEFAULT), ...
[ ON DUPLICATE KEY UPDATE
  col_name=expr [, col_name=expr] ... ]
```

#### 3.2.1.2 举例

```
insert into test (col1,col3) values(1,'cust1'),(2,'cust2');
insert into test (col1,col3) values(default,'cust3');
insert into test set col1=4,col3='cust4';
insert into test set col1=default,col3='cust5';
insert into test (col1,col3) values(default,cast(now() as char));
```

#### 3.2.1.3 限制

- 在插入ER关系的子表时，每个语句之允许插入一个ROW
- 全局序列在插入时不允许指定值，全部由dble序列生成
- 对于含有枚举类型的分片表，比如表结构：CREATE TABLE `test` (`id` enum('1','2','3') DEFAULT '1')，在插入id列时，既可以使用枚举值插入，也可以使用枚举值的下标，'1'的下标是1，以此类推。若用户以枚举值进行分片，但是插入时确使用枚举值下标，因为dble不会将下标转换为枚举值，所以分片会出现问题，详细可参考issue：<https://github.com/actiontech/dble/issues/816>。

### 3.2.2 REPLACE

#### 3.2.2.1 Syntax

REPLACE

```
[INTO] tbl_name [(col_name [, col_name] ...)]  
{VALUES | VALUE} {value_list} [, {value_list}] ...
```

OR

REPLACE

```
[INTO] tbl_name SET assignment_list
```

#### 3.2.2.2 举例

```
REPLACE INTO test VALUES (1, 'Old', '2014-08-20 18:47:00');  
REPLACE INTO test set id = 1, type= 'Old',create_date = '2014-08-20 18:47:00';
```

#### 3.2.2.3 限制

- 由于replace的语义为如果存在则替换，如果不存在则新增，所以在使用表格自增主键的时候
- 如果对于自增表格使用replace且ID不存在，那么就会插入一条指定ID的数据，并不会自动生成ID

### 3.2.3 DELETE

#### 3.2.3.1 Syntax

```
DELETE [IGNORE]  
FROM tbl_name [WHERE where_condition]
```

#### 3.2.3.2 举例

```
delete from test where id>5;
```

#### 3.2.3.3 限制

- Delete语句中的where\_condition部分只允许出现简单的条件，不能支持计算表达式以及子查询
- 不支持多表Join 的DELETE

### 3.2.4 UPDATE

#### 3.2.4.1 Syntax

```
UPDATE table_reference  
SET col_name1={expr1} [, col_name2={expr2}] ...  
[WHERE where_condition]
```

#### 3.2.4.2 举例

```
UPDATE test SET VALUE =1 where id=5;
```

#### 3.2.4.3 限制

- UPDATE语句中的where\_condition部分只允许出现简单的条件，不能支持计算表达式以及子查询
- 不支持多表Join 的UPDATE

### 3.2.5 SELECT

#### 3.2.5.1 Syntax

```
SELECT  
[ALL | DISTINCT | DISTINCTROW ]  
select_expr  
[, select_expr ...]  
[FROM table_references [WHERE where_condition]  
[GROUP BY {col_name | expr | position} [ASC | DESC], ...]  
[HAVING where_condition] [ORDER BY {col_name | expr | position} [ASC | DESC], ...]  
[LIMIT {[offset,] row_count | row_count OFFSET offset}]
```

#### 3.2.5.2 举例

```
select id,col1,col3 from test where id=3;  
select distinct col1,col3 from test where id>=3;  
select count(*),max(id),col1 from test group by col1 desc having(count(*)>1) order by col1 desc;  
select id,col1,col3 from test order by id limit 2 offset 2;  
select id,col1,col3 from test order by id limit 2,2;  
select 1+1,'test',id,col1*1.1,now() from test limit 3;  
select current_date,current_timestamp;
```

### 3.2.6 JOIN Syntax:

table\_references:

```
table_reference [, table_reference] ...
```

table\_reference:

```
table_factor | join_table
```

table\_factor:

```
tbl_name [[AS] alias]
| table_subquery [AS] alias
| ( table_references )
```

join\_table:

```
table_reference [INNER | CROSS] JOIN table_factor [join_condition]
| table_reference STRAIGHT_JOIN table_factor
| table_reference STRAIGHT_JOIN table_factor ON conditional_expr
| table_reference {LEFT|RIGHT} [OUTER] JOIN table_reference join_condition
| table_reference NATURAL [{LEFT|RIGHT} [OUTER]] JOIN table_factor
```

join\_condition:

```
ON conditional_expr
| USING (column_list)
```

### 3.2.7 UNION Syntax:

```
SELECT ...
UNION [ALL | DISTINCT] SELECT ...
[UNION [ALL | DISTINCT] SELECT ...]
```

### 3.2.8 Subquery

#### 3.2.8.1 The Subquery as Scalar Operand

For example :

```
SELECT (SELECT s2 FROM t1);
SELECT (SELECT s1 FROM t2) FROM t1;
SELECT UPPER((SELECT s1 FROM t1)) FROM t2;
```

#### 3.2.8.2 Comparisons Using Subqueries

The most common use of a subquery is in the form:

```
non_subquery_operand comparison_operator (subquery)
```

Where comparison\_operator is one of these operators:

```
= > < >= <= <> != <=>
```

MySQL also permits this construct:

```
non_subquery_operand LIKE (subquery)
```

#### 3.2.8.3 Subqueries with ANY, IN, or SOME

Syntax:

```
operand comparison_operator ANY (subquery)
operand IN (subquery)
operand comparison_operator SOME (subquery)
```

Where comparison\_operator is one of these operators:

```
= > < >= <= <> !=
```

#### 3.2.8.4 Subqueries with ALL

Syntax:

```
operand comparison_operator ALL (subquery)
```

#### 3.2.8.5 Subqueries with EXISTS or NOT EXISTS

For example:

```
SELECT column1 FROM t1 WHERE EXISTS (SELECT * FROM t2);
```

Not support Correlated Subqueries for now.

#### 3.2.8.6 Derived Tables (Subqueries in the FROM Clause)

```
SELECT ... FROM (subquery) [AS] tbl_name ...
```

### 3.2.9 LOAD DATA

#### 3.2.9.1 Syntax

```
LOAD DATA
[LOCAL]
INFILE 'file_name' INTO TABLE tbl_name
[CHARACTER SET charset_name]
[{{FIELDS | COLUMNS}
[TERMINATED BY 'string']
[[OPTIONALLY] ENCLOSED BY 'char']
[ESCAPED BY 'char']] ]
[LINES [STARTING BY 'string']
[TERMINATED BY 'string']] ]
```

#### 3.2.9.2 举例

```
load data infile 'data.txt' into table test_table CHARACTER SET 'utf8' FIELDS TERMINATED by ',';
```

#### 3.2.9.3 原理

dble解析MySQL协议之后，会根据数据路由拆分文件，每满足maxRowSizeToFile(可通过server.xml配置)就写到文件中，再通过 load data local infile 的方式导入到后端结点。所以，这里local\_infile这个参数会影响到load data的正确性。

可参考[#1085](#)

#### 3.2.9.4 限制

- 存在BUG导致在dble中CHARACTER SET charset\_name必填
- 在mysql中如果插入的数据不符合规范会插入部分数据，相对来说dble的load data对正确性有更高的要求，一个错误的发生会导致整体操作的回滚
- 在ENCLOSED BY的时候存在BUG，使用的时候会导致转义的数据无法被正确转义存储到数据库中
- 由于当前解析设定的限制，loaddata默认每列最大字节数65535，可通过server.xml 中的maxCharsPerColumn配置修改
- 使用load data导入数据时，若导入表是分片表，应保证导入文件中分片键数据符合分片规则的要求，否则dble将会报错，详细参考issue：<https://github.com/actiontech/dble/issues/770>
- load data语句需要严格按照语法书写，由于druid解析器的缘故，若语法或关键字错误，druid会解析出错误的语句，从而导致结果错误，详细请参照issue：<https://github.com/actiontech/dble/issues/1248>
- load data语句读字段的时候，如果遇到行结束符，会认为是本行结束了，需要注意：<https://github.com/actiontech/dble/issues/1507>
- load data中使用用户变量后，再次查询用户变量的值是不正确的，可参考issue：<https://github.com/actiontech/dble/issues/1761>

### 3.2.10 不支持的DML语句

**DO Syntax**

**HANDLER Syntax**

**LOAD XML Syntax**

### 3.3 PREPARE SQL Syntax

#### 3.3.1 PREPARE Syntax

```
PREPARE stmt_name FROM preparable_stmt
```

例:

```
prepare stmt1 from "select * from a_test where id=?";
```

#### 3.3.2 EXECUTE Syntax

```
EXECUTE stmt_name  
[USING @var_name [, @var_name] ...]
```

例:

```
SET @a = 1;  
EXECUTE stmt1 USING @a;
```

#### 3.3.3 DEALLOCATE PREPARE Syntax

```
{DEALLOCATE | DROP} PREPARE stmt_name
```

例:

```
DROP PREPARE stmt1;
```

## 3.4 Transactional and Locking Statements

Transactional and Locking Statements 包含以下内容

- [3.4.1 Lock&unlock](#)
- [3.4.2 XA 事务语法](#)
- [3.4.3 一般事务语法](#)
- [3.4.4 SET TRANSACTION Syntax](#)

### 3.4.1 Lock&unlock

#### 3.4.1.1 Syntax

```
LOCK TABLES tbl_name [[AS] alias] lock_type
```

lock\_type: READ | WRITE

UNLOCK TABLES

#### 3.4.1.2 举例

```
lock tables test_table read;  
unlock tables;
```

#### 3.4.1.3 限制

1. 当前session加锁后访问其他表可能不会被阻止或者报错。
2. 加写锁后，复杂查询可能不会返回正确结果。
3. 如果在xa事务中使用lock table语句，db1e会将隐式提交并将xa事务关闭。

### 3.4.2 XA 事务语法

#### 3.4.2.1 Syntax

开启XA

```
set xa = {0|1}
```

开启事务

```
START TRANSACTION;
```

```
BEGIN
```

```
SET autocommit = {0 | 1}
```

提交事务

```
COMMIT
```

回滚事务

```
ROLLBACK
```

#### 3.4.2.2 限制

1. 在Dble中一旦开启了SQL黑名单检查，不能使用BEGIN开启事务(druid解析器不支持)。
2. 使用 lock table 语句后会将 xa 事务关闭。

### 3.4.3 一般事务语法

#### 3.4.3.1 Syntax

开启事务

```
START TRANSACTION;
```

BEGIN

```
SET autocommit = {0 | 1}
```

提交事务

```
COMMIT
```

回滚事务

```
ROLLBACK
```

#### 3.4.3.2 限制

- 在Dble中一旦开启了SQL黑名单检查，不能使用BEGIN开启事务(druid解析器不支持)
- 2PC实现的分布式事务(非xa方式)可能会出现commit时部分提交的情况,如需保障最终一致性，需要开启XA

#### 3.4.4 SET TRANSACTION Syntax

```
SET SESSION TRANSACTION ISOLATION LEVEL level
```

**level:**

```
REPEATABLE READ  
| READ COMMITTED  
| READ UNCOMMITTED  
| SERIALIZABLE  
SET @@SESSION.TX_ISOLATION = 'level_str'
```

```
level_str: REPEATABLE-READ  
| READ-COMMITTED  
| READ-UNCOMMITTED  
| SERIALIZABLE
```

注: 因为隔离级别不加**session**关键字语义不同, 暂不支持

### 3.4.5 SAVEPOINT, ROLLBACK TO SAVEPOINT, and RELEASE SAVEPOINT Syntax

#### 3.4.5.1 Syntax

```
SAVEPOINT identifier
ROLLBACK [WORK] TO [SAVEPOINT] identifier
RELEASE SAVEPOINT identifier
```

#### 3.4.5.1 举例

```
# start transaction
set autocommit = 0;

# savepoint
savepoint s0;
insert into test value(1);
savepoint s1;
insert into test value(2);
savepoint s2;
insert into test value(3);

# rollback to
rollback to s0

# release
release savepoint s0
```

#### 3.4.5.2 限制

1. 在mysql中,在事务外也可以定义savepoint,但是这些savepoint是没有意义的.因此在db中savepoint强制在事务中使用,否则报错.
2. ROLLBACK TO [SAVEPOINT] *identifier* 语句暂不支持 work 可选项.

### 3.5 DAL

DAL主要包含以下内容

- [3.5.1 SET](#)
- [3.5.2 SHOW](#)
- [3.5.3 KILL](#)

### 3.5.1 SET语句

#### 3.5.1.1 XA

```
set xa=value

value:
  0
  | off
  | false
  | 1
  | on
  | true
```

例:

```
set xa=1;
```

注意事项: XA设置不能在多变量设置语句中使用。

#### 3.5.1.2 AUTOCOMMIT

```
set autocommit=value

value:
  0
  | off
  | false
  | 1
  | on
  | true
```

例:

```
set autocommit=1;
```

注意事项: AUTOCOMMIT设置不能在多变量设置语句中使用。

#### 3.5.1.3 NAMES

```
SET NAMES {'charset_name' [COLLATE 'collation_name'] | DEFAULT}
```

例:

```
set names utf8;
set names utf8 collate utf8_general_ci;
set names default;
```

#### 3.5.1.4 CHARSET

```
SET {CHARACTER SET | CHARSET}
{'charset_name' | DEFAULT}
```

例:

```
set CHARACTER SET utf8;
```

#### 3.5.1.5 COLLATION\_CONNECTION/CHARSET\_SET\_X

```
SET COLLATION_CONNECTION='collation_name'
SET CHARSET_SET_CLIENT='charset_name'.
SET CHARSET_SET_RESULTS='charset_name' 其中, 'charset_name' 可以为NULL。
SET CHARSET_SET_CONNECTION='charset_name'
```

例:

```
set collation_connection=utf8_general_ci;
set CHARSET_SET_CLIENT=utf8;
set CHARSET_SET_RESULTS=utf8;
set CHARSET_SET_CONNECTION=utf8;
```

#### 3.5.1.6 TRANSACTION ACCESS MODE

```
SET SESSION { TX_READ_ONLY | TRANSACTION_READ_ONLY}=value

value:
  0
  | off
  | false
  | 1
  | on
  | true
```

例:

```
set session @@tx_read_only=1;
```

#### 3.5.1.7 TRANSACTION ISOLATION LEVEL

```
SET SESSION {TRANSACTION_ISOLATION | TX_ISOLATION}=level

level:
  READ-UNCOMMITTED | READ-COMMITTED | REPEATABLE-READ | SERIALIZABLE
```

例:

```
SET SESSION TX_ISOLATION=READ-COMMITTED;
```

#### 3.5.1.8 USER/SYSTEM VARIABLE

```
SET variable_assignment[, variable_assignment] ...  
  
variable_assignment:  
  
@user_var_name = expr  
| SESSION system_var_name = expr  
| system_var_name = expr  
| @@system_var_name = expr  
| @@session.system_var_name = expr
```

注意事项:

1. 不能设置全局系统变量。
2. 支持的系统变量为:

```

audit_log_current_session
audit_log_filter_id
auto_increment_increment
auto_increment_offset
autocommit
big_tables
binlog_direct_non.Transactional_updates
binlog_error_action
binlog_format
binlog_row_image
binlog_rows_query_log_events
binlogging_impossible_mode
block_encryption_mode
bulk_insert_buffer_size
character_set_client
character_set_connection
character_set_database
character_set_filesystem
character_set_results
character_set_server
collation_connection
collation_database
collation_server
completion_type
debug
debug_sync
default_storage_engine
default_tmp_storage_engine
default_week_format
disconnect_on_expired_password
div_precision_increment
end_markers_in_json
eq_range_index_dive_limit
error_count
explicit_defaults_for_timestamp
external_user
foreign_key_checks
group_concat_max_len
gtid_next
gtid_owned
identity
innodb_create_intrinsic
innodb_ft_user_stopword_table
innodb_lock_wait_timeout
innodb_optimize_point_storage
innodb_strict_mode
innodb_support_xa
innodb_table_locks
innodb_tmpdir
insert_id
interactive_timeout
join_buffer_size
keep_files_on_create
last_insert_id
lc_messages
lc_time_names
lock_wait_timeout
long_query_time
low_priority_updates
max_allowed_packet
max_delayed_threads
max_error_count
max_execution_time
max_heap_table_size
max_insert_delayed_threads
max_join_size
max_length_for_sort_data
max_seeks_for_key
max_sort_length
max_sp_recursion_depth
max_statement_time
max_user_connections
min_examined_row_limit
myisam_repair_threads
myisam_sort_buffer_size
myisam_stats_method
ndb-allow-copying-alter-table
ndb_autoincrement_prefetch_sz
ndb_blob-read-batch-bytes
ndb_blob-write-batch-bytes
ndb_deferred_constraints
ndb_force_send
ndb_fully_replicated
ndb_index_stat_enable
ndb_index_stat_option
ndb_join_pushdown
ndb_log_bin
ndb_log_bin
ndb_table_no_logging
ndb_table_temporary
ndb_use_copying_alter_table
ndb_use_exact_count
ndb_use_transactions
ndbinfo_max_bytes
ndbinfo_max_rows
ndbinfo_show_hidden
ndbinfo_table_prefix
net_buffer_length
net_read_timeout
net_retry_count
net_write_timeout
new
old_alter_table
old_passwords
optimizer_prune_level
optimizer_search_depth
optimizer_switch
optimizer_trace
optimizer_trace_features
optimizer_trace_limit
optimizer_trace_max_mem_size
optimizer_trace_offset
parser_max_mem_size
preload_buffer_size
profiling
profiling_history_size
proxy_user
pseudo_slave_mode
pseudo_thread_id
query_alloc_block_size
query_cache_type
query_cache_wlock_invalidate
query_prealloc_size
rand_seed1
rand_seed2
range_alloc_block_size
range_optimizer_max_mem_size
rbr_exec_mode
read_buffer_size
read_rnd_buffer_size
session_track_gtids
session_track_schema
session_track_state_change
session_track_system_variables
show_old_temporals
sort_buffer_size
sql_auto_is_null
sql_big_selects
sql_buffer_result
sql_log_bin
sql_log_off
sql_mode
sql_notes
sql_quote_show_create
sql_safe_updates

```

```
sql_select_limit
sql_warnings
storage_engine
thread_pool_high_priority_connection
thread_pool_prio_kickup_timer
time_zone
timestamp
tmp_table_size
transaction_alloc_block_size
transaction_allow_batching
transaction_prealloc_size
transaction_write_set_extraction
tx_isolation
tx_read_only
unique_checks
updateable_views_with_limit
version_tokens_session
version_tokens_session_number
wait_timeout
warning_count
```

例:

```
set @a=20;
SET SESSION sql_mode = 'TRADITIONAL';
SET sql_mode = 'TRADITIONAL';
```

1. insert\_id 在使用过程中可能会在另一个前端连接中被重置而导致主键冲突的问题。sql\_auto\_is\_null 和 insert\_id 是联合使用的，使用时也有此限制。详情参见issue: <https://github.com/actiontech/dble/issues/1252>.

### 3.5.1.9 TRACE

用于观察单条SQL的性能，打开此开关后，可以执行需要观察性能的查询语句，然后执行 `show trace` 来观察最后结果。

可以用 `select @@trace` 观察当前的开启状态。详情请见 [单条SQL性能trace](#)

```
set trace=value

value:
  0
  | off
  | false
  | 1
  | on
  | true
```

例:

```
set trace=1;
```

### 3.5.2 SHOW语句

#### 3.5.2.1 dble劫持的SHOW

- SHOW DATABASES  
将schema.xml 中的所有schema展示出来。
- SHOW CREATE DATABASE [IF NOT EXISTS] schema  
将schema.xml 中的指定schema的创建语句展示，创建语句为dble伪造，无实际意义。
- SHOW [FULL|ALL] TABLES [FROM db\_name] [LIKE 'pattern'] WHERE expr  
当schmea没有配置默认节点时，将schema下配置的tables直接展示出来。  
当schema有默认节点时，将语句转发至默认节点，然后将结果集与schema下配置的tables做一个去重合并，再返回给客户端。
- SHOW ALL TABLES [FROM db\_name] [LIKE 'pattern'] WHERE expr  
dble自有命令，与SHOW FULL TABLES 返回结果集类似，不同之处是Table\_type这列分为了 SHARDING TABLE, sharding table, GLOBAL TABLE。参见[6.Differernce\\_from\\_SQL\\_Server.md](#)。
- SHOW [FULL] {COLUMNS | FIELDS} FROM tbl\_name [{FROM|IN} db\_name] [LIKE 'pattern' | WHERE expr]  
将逻辑schema转为物理schema之后下发到表所在的任意节点。
- SHOW { INDEX | INDEXES | KEYS } {FROM | IN} tbl\_name [ {FROM | IN} db\_name ] [ WHERE expr ]  
将逻辑schema转为物理schema之后下发到表所在的任意节点。
- SHOW CREATE TABLE tbl\_name  
将逻辑schema转为物理schema之后下发到表所在的任意节点。
- SHOW [GLOBAL | SESSION] VARIABLES [LIKE 'pattern' | WHERE expr]  
随机转发到任意节点，收到结果集后，用本地变量进行覆盖（global不正确？）
- SHOW CREATE VIEW view\_name  
将dble层面的view展示出来
- SHOW CHARSET  
将之转为show character set 之后透传转发
- SHOW TABLE STATUS [{FROM | IN} db\_name] [LIKE 'pattern' | WHERE expr]  
只是为了支持SQLyog，其中name列逻辑和show tables一致，其他列均为伪造。
- SHOW TRACE  
观察trace结果，详情请见[单条SQL性能trace](#)

注意事项：

所有以上命令的explain结果可能不准确。

例：

```
show databases;
show full tables;
show columns from a_test;
show index from a_test;
show create table a_test;
show variables;
show charset;
```

#### 3.5.2.2 dble透传的SHOW

除了dble劫持的特定SHOW语句外，其它SHOW语句都透传，这些语句与mysql语法相同。

例：

```
SHOW CHARACTER SET;
SHOW CHARACTER SET like 'utf8';
SHOW CHARACTER SET where maxlen=2;
```

### 3.5.3 KILL

#### 3.5.3.1 KILL conn\_id

其中，`conn_id`为前端连接id值，可以通过运维命令`show @@connection` 获取。

##### 3.5.3.1.1 举例

```
kill 1;
```

##### 3.5.3.1.2 限制

- 在Kill自身连接的时候只会向自身写入OK包，不会有其他操作
- 如果Kill的连接在XA事务的提交或者回滚状态，不会直接关闭后端连接，会仅关闭前端连接
- 后端连接的关闭通过向MySQL节点发送 `KILL processlist_id` 来完成

#### 3.5.3.2 KILL query conn\_id

其中，`conn_id`为前端连接id值，可以通过运维命令`show @@connection` 获取。

##### 3.5.3.2.1 举例

```
kill query 1;
```

##### 3.5.3.2.3 说明

- dble 中`kill query`的实现是将正在执行语句的后端连接与前端连接相割离的方式来实现。
- 后端未执行完成的语句，取决与mysql自身的机制。

##### 3.5.3.2.2 限制

- 对于`ddl`语句，不保证一致性
- 对于未开启事务的`dml`操作不保证一致性

## 3.6 存储过程支持方式

### 3.6.1 Syntax

#### Create procedure

```
/Hint/ CREATE [DEFINER = { user | CURRENT_USER }]
PROCEDURE sp_name ([proc_parameter[...]])
[characteristic ...] routine_body
/Hint/ CREATE
[DEFINER = { user | CURRENT_USER }]
FUNCTION sp_name ([func_parameter[...]])
RETURNS type [characteristic ...] routine_body
```

#### drop procedure

```
/Hint/ DROP {PROCEDURE | FUNCTION} [IF EXISTS] sp_name
```

#### call procedure

```
[/Hint/] CALL sp_name([parameter[...]])
[/Hint/] CALL sp_name[]
```

### 3.6.2 举例

```
删除存储过程:
/*!dbe:sql=select 1 from account */drop procedure if exists proc_arc;
创建存储过程:
/*!dbe:sql=select 1 from account */create procedure proc_arc(userid1 int)
begin
    insert into account_arc select * from account where userid=userid1;
    update account set arc_flag=true,arc_time=now() where userid=userid1;
end;
调用存储过程:
/*!dbe:sql=select 1 from account */call proc_arc(1);
```

### 3.6.3 限制

- dble支持存储过程和自定义函数的透传，存储过程的开发完全使用MySQL的语法，开发、调试与部署的方法同单机MySQL相同。存储过程和自定义函数需要在所有节点上创建，节点扩容的时候也需要考虑存储过程和自定义函数的迁移。
- 存储过程和自定义函数是直接发送到节点上执行，中间件不参与运算，因此要慎重使用，需要保证过程的内部不出现跨节点运算。
- 存储过程调用时，要在调用语句之前增加注解，系统根据注解透传到节点运行，存储过程的执行路径以及执行结果的正确性由开发者保证。对于只是写入数据，不返回结果的存储过程，需要注意避免重复写入数据。对于返回结果的存储过程，需要特别注意返回结果的正确性。dble不会对存储过程的结果进行汇聚运算，只能由应用端自行完成。

## 3.7 Utility Statements

### 3.7.1 USE

```
USE db_name
```

例:

```
use TESTDB;
```

### 3.7.2 EXPLAIN

```
EXPLAIN explainable_stmt
```

```
SELECT statement
| DELETE statement
| INSERT statement
| REPLACE statement
| UPDATE statement
```

注意事项:

1. INSERT中表不能为自增序列表

例:

```
explain SELECT select * from a_test where id=1;
```

2. 在db中, EXPLAIN 不等价于DESC

### 3.7.3 EXPLAIN2

```
EXPLAIN2 DATANODE=node_name sql$sql_stmt
```

例:

```
explain2 datanode=dn2 sql=select * from a_test where id=1;
```

### 3.7.4 DESC

```
{DESCRIBE | DESC} tbl_name [col_name | wild]
```

例:

```
DESC a_test id;
```

注意: 在db中, EXPLAIN 不等价于DESC

## 3.8 Hint

### 3.8.1 Syntax

```
/* { ! | #}dble: {sql=SELECT select_expr FROM table_references WHERE where_condition
|datanode=datanode_name
|db_type={slave|master}}
*/ ordinary_sql
```

### 3.8.2 举例

```
/*!dble:sql=select 1 from sbtest */ call p_show_time();
/*!dble:datanode=dn1*/ update sbtest set name = 'test';
/*!dble:db_type=master*/ select count(*) from sbtest;
/*#dble:sql=select 1 from sbtest */ call p_show_time();
/*#dble:datanode=dn1*/ update sbtest set name = 'test';
/*#dble:db_type=master*/ select count(*) from sbtest;
```

### 3.9 其他不支持语句

- Compound-Statement Syntax
- Replication Statements
- DDL:
  - 不支持针对database的操作语句，包括alter database、drop database、create database 在管理端口支持，流量端口遇到会判断schema已经配置后返回ok，否则报错。
  - 不支持create table时的一些table option，如DATA DIRECTORY、ALGORITHM等，table option在alter table时也不能修改
  - 不支持ALTER TABLE ... LOCK ...
  - 不支持ALTER TABLE ... ORDER BY ...
  - 不支持create table ... like ... 和create table ... select ...
  - 库名、表名不可修改，拆分字段的名称和类型都不可以变更
  - 不支持外键关联
  - 不支持临时表
  - 不支持分布式级别的存储过程和自定义函数
  - 不支持触发器
- DML:
  - 对于INSERT... VALUES(expr)，不支持expr中含有子查询
  - 不支持INSERT DELAYED...
  - 不支持INSERT... SELECT...
  - 不支持不包含拆分字段的INSERT语句
  - 不支持HANDLER语句
  - 不支持UPDATE多张表
  - 不支持修改拆分字段的值
  - 不支持DELETE ... ORDER BY ... LIMIT ...
  - 不支持DELETE多张表
  - 不支持DELETE/UPDATE ...LIMIT路由到一个分片表的多个节点
  - 不支持DO语句
- 查询:
  - 不支持select ... use/ignore index ...
  - 不支持select ... group by ... with rollup
  - 不支持select ... for update | lock in share mode 正确语义
  - 不支持select ... into outfile ...
  - 不支持Row Subqueries
  - 不支持select ... union [all] select ... order by ..., 可写成(select ...) union [all] (select ...) order by ...
  - 不支持session变量赋值与查询，如set @rowid=0;select @rowid:="@rowid+1,id from user;
- 管理语句:
  - 不支持用户管理及权限管理语句
  - 不支持表维护语句，包括ANALYZE/CHECK/CHECKSUM/OPTIMIZE/REPAIR TABLE
  - 不支持INSTALL/UNINSTALL PLUGIN语句
  - 不支持BINLOG语句
  - 不支持CACHE INDEX/ LOAD INDEX INTO CACHE语句
  - 不支持除FLUSH TABLES [WITH READ LOCK]以外的其他FLUSH语句，FLUSH TABLE也仅语法支持无实际意义
  - 不支持RESET语句
  - 不支持大部分的运维SHOW语句，如SHOW PROFILES、SHOW ERRORS等

## 3.10 函数与操作符支持列表(alpha版本)

### 3.10.0 注意:

1. 如果能保证SQL会将函数整体下发给某个后端结点，函数支持度没有意义，支持度将托管于MySQL。
2. 在2.18.09.0版本之前（含），涉及到中文字符集的字符串相关函数有bug。
3. 除了聚合函数经过充分测试之外，其余函数只进行了简单的测试，并没有全覆盖测试，请使用前充分测试。
4. 未在列表中的函数，一概不支持。

### 3.10.1 Operators

Name	Description	Support
AND, &&	Logical AND	Y
=	Assign a value (as part of a SET statement, or as part of the SET clause in an UPDATE statement)	Y
:=	Assign a value	N
BETWEEN ... AND ...	Check whether a value is within a range of values	Y
BINARY	Cast a string to a binary string	N
&	Bitwise AND	Y
~	Bitwise inversion	Y
\		Bitwise OR Y
^	Bitwise XOR	Y
CASE	Case operator	Y
DIV	Integer division	Y
/	Division operator	Y
=	Equal operator	Y
<=>	NULL-safe equal to operator	Y
>	Greater than operator	Y
>=	Greater than or equal operator	Y
IS	Test a value against a boolean	Y
IS NOT	Test a value against a boolean	Y
IS NOT NULL	NOT NULL value test	Y
IS NULL	NULL value test	Y
->	Return value from JSON column after evaluating path; equivalent to JSON_EXTRACT().	N
->>	Return value from JSON column after evaluating path and unquoting the result; equivalent to JSON_UNQUOTE(JSON_EXTRACT()).	N
<<	Left shift	Y
<	Less than operator	Y
<=	Less than or equal operator	Y
LIKE	Simple pattern matching	Y
-	Minus operator	Y
%, MOD	Modulo operator	Y
NOT, !	Negates value	Y
NOT BETWEEN ... AND ...	Check whether a value is not within a range of values	Y
!=, <>	Not equal operator	Y
NOT LIKE	Negation of simple pattern matching	Y
NOT REGEXP	Negation of REGEXP	Y
\	\	, OR Logical OR Y
+	Addition operator	Y
REGEXP	Whether string matches regular expression	Y
>>	Right shift	Y
RLIKE	Whether string matches regular expression	N
SOUNDS LIKE	Compare sounds	N
*	Multiplication operator	N
-	Change the sign of the argument	Y
XOR	Logical XOR	Y
COALESCE()	Return the first non-NULL argument	Y
GREATEST()	Return the largest argument	Y
IN()	Check whether a value is within a set of values	Y
INTERVAL()	Return the index of the argument that is less than the first argument	Y
ISNULL()	Test whether the argument is NULL	Y
LEAST()	Return the smallest argument	Y
STRCMP()	Compare two strings	Y

### 3.10.2 Control Flow Functions

Name	Description	Support
CASE	Case operator	Y
IF()	If/else construct	Y
IFNULL()	Null if/else construct	Y
NULLIF()	Return NULL if expr1 = expr2	Y

### 3.10.3 String Functions

Name	Description	Support
ASCII()	Return numeric value of left-most character	Y
BIN()	Return a string containing binary representation of a number	N
BIT_LENGTH()	Return length of argument in bits	Y
CHAR()	Return the character for each integer passed	Y
CHAR_LENGTH()	Return number of characters in argument	Y
CHARACTER_LENGTH()	Synonym for CHAR_LENGTH()	Y
CONCAT()	Return concatenated string	Y
CONCAT_WS()	Return concatenate with separator	Y
ELT()	Return string at index number	Y
EXPORT_SET()	Return a string such that for every bit set in the value bits, you get an on string and for every unset bit, you get an off string	N
FIELD()	Return the index (position) of the first argument in the subsequent arguments	Y
FIND_IN_SET()	Return the index position of the first argument within the second argument	Y
FORMAT()	Return a number formatted to specified number of decimal places	Y
FROM_BASE64()	Decode base64 encoded string and return result	N
HEX()	Return a hexadecimal representation of a decimal or string value	Y
INSERT()	Insert a substring at the specified position up to the specified number of characters	Y
INSTR()	Return the index of the first occurrence of substring	Y
LCASE()	Synonym for LOWER()	Y
LEFT()	Return the leftmost number of characters as specified	Y
LENGTH()	Return the length of a string in bytes	Y
LIKE	Simple pattern matching	Y
LOAD_FILE()	Load the named file	N
LOCATE()	Return the position of the first occurrence of substring	Y
LOWER()	Return the argument in lowercase	Y
LPAD()	Return the string argument, left-padded with the specified string	Y
LTRIM()	Remove leading spaces	Y
MAKE_SET()	Return a set of comma-separated strings that have the corresponding bit in bits set	Y
MATCH	Perform full-text search	N
MID()	Return a substring starting from the specified position	N
NOT LIKE	Negation of simple pattern matching	Y
NOT REGEXP	Negation of REGEXP	Y
OCT()	Return a string containing octal representation of a number	N
OCTET_LENGTH()	Synonym for LENGTH()	N
ORD()	Return character code for leftmost character of the argument	Y
POSITION()	Synonym for LOCATE()	N
QUOTE()	Escape the argument for use in an SQL statement	Y
REGEXP	Whether string matches regular expression	Y
REPEAT()	Repeat a string the specified number of times	Y
REPLACE()	Replace occurrences of a specified string	Y
REVERSE()	Reverse the characters in a string	Y
RIGHT()	Return the specified rightmost number of characters	Y
RLIKE	Whether string matches regular expression	N
RPAD()	Append string the specified number of times	Y
RTRIM()	Remove trailing spaces	Y
SOUNDEX()	Return a soundex string	Y
SOUNDS LIKE	Compare sounds	Y
SPACE()	Return a string of the specified number of spaces	Y
STRCMP()	Compare two strings	Y
SUBSTR()	Return the substring as specified	Y
SUBSTRING()	Return the substring as specified	Y
SUBSTRING_INDEX()	Return a substring from a string before the specified number of occurrences of the delimiter	Y
TO_BASE64()	Return the argument converted to a base-64 string	N
TRIM()	Remove leading and trailing spaces	Y
UCASE()	Synonym for UPPER()	Y
UNHEX()	Return a string containing hex representation of a number	Y
UPPER()	Convert to uppercase	Y
WEIGHT_STRING()	Return the weight string for a string	N

### 3.10.4 Numeric Functions and Operators

Name	Description	Support
ABS()	Return the absolute value	Y
ACOS()	Return the arc cosine	Y
ASIN()	Return the arc sine	Y
ATAN()	Return the arc tangent	Y
ATAN2(), ATAN()	Return the arc tangent of the two arguments	Y
CEIL()	Return the smallest integer value not less than the argument	Y
CEILING()	Return the smallest integer value not less than the argument	Y
CONV()	Convert numbers between different number bases	Y
COS()	Return the cosine	Y
COT()	Return the cotangent	Y
CRC32()	Compute a cyclic redundancy check value	Y
DEGREES()	Convert radians to degrees	Y
DIV	Integer division	Y
/	Division operator	Y
EXP()	Raise to the power of	Y
FLOOR()	Return the largest integer value not greater than the argument	Y
LN()	Return the natural logarithm of the argument	Y
LOG()	Return the natural logarithm of the first argument	Y
LOG10()	Return the base-10 logarithm of the argument	Y
LOG2()	Return the base-2 logarithm of the argument	Y
-	Minus operator	Y
MOD()	Return the remainder	Y
%, MOD	Modulo operator	Y
PI()	Return the value of pi	Y
+	Addition operator	Y
POW()	Return the argument raised to the specified power	Y
POWER()	Return the argument raised to the specified power	Y
RADIANS()	Return argument converted to radians	Y
RAND()	Return a random floating-point value	Y
ROUND()	Round the argument	Y
SIGN()	Return the sign of the argument	Y
SIN()	Return the sine of the argument	Y
SQRT()	Return the square root of the argument	Y
TAN()	Return the tangent of the argument	Y
*	Multiplication operator	Y
TRUNCATE()	Truncate to specified number of decimal places	Y
-	Change the sign of the argument	Y

### 3.10.5 Date and Time Functions

Name	Description	Support
ADDDATE()	Add time values (intervals) to a date value	Y
ADDTIME()	Add time	Y
CONVERT_TZ()	Convert from one time zone to another	N
CURDATE()	Return the current date	Y
CURRENT_DATE()	Synonyms for CURDATE()	Y
CURRENT_TIME()	Synonyms for CURTIME()	Y
CURRENT_TIMESTAMP()	Synonyms for NOW()	Y
CURTIME()	Return the current time	Y
DATE()	Extract the date part of a date or datetime expression	Y
DATE_ADD()	Add time values (intervals) to a date value	Y
DATE_FORMAT()	Format date as specified	Y
DATE_SUB()	Subtract a time value (interval) from a date	Y
DATEDIFF()	Subtract two dates	Y
DAY()	Synonym for DAYOFMONTH()	N
DAYNAME()	Return the name of the weekday	Y
DAYOFMONTH()	Return the day of the month (0-31)	Y
DAYOFWEEK()	Return the weekday index of the argument	Y
DAYOFYEAR()	Return the day of the year (1-366)	Y
EXTRACT()	Extract part of a date	Y
FROM_DAYS()	Convert a day number to a date	Y
FROM_UNIXTIME()	Format Unix timestamp as a date	Y
GET_FORMAT()	Return a date format string	Y
HOUR()	Extract the hour	Y
LAST_DAY	Return the last day of the month for the argument	N
LOCALTIME()	Synonym for NOW()	Y
LOCALTIMESTAMP()	Synonym for NOW()	Y
MAKEDATE()	Create a date from the year and day of year	Y
MAKETIME()	Create time from hour, minute, second	Y
MICROSECOND()	Return the microseconds from argument	Y
MINUTE()	Return the minute from the argument	Y
MONTH()	Return the month from the date passed	Y
MONTHNAME()	Return the name of the month	Y
NOW()	Return the current date and time	Y
PERIOD_ADD()	Add a period to a year-month	Y
PERIOD_DIFF()	Return the number of months between periods	Y
QUARTER()	Return the quarter from a date argument	Y
SEC_TO_TIME()	Converts seconds to 'HH:MM:SS' format	Y
SECOND()	Return the second (0-59)	Y
STR_TO_DATE()	Convert a string to a date	Y
SUBDATE()	Synonym for DATE_SUB() when invoked with three arguments	Y
SUBTIME()	Subtract times	Y
SYSDATE()	Return the time at which the function executes	Y
TIME()	Extract the time portion of the expression passed	Y
TIME_FORMAT()	Format as time	Y
TIME_TO_SEC()	Return the argument converted to seconds	Y
TIMEDIFF()	Subtract time	Y
TIMESTAMP()	With a single argument, this function returns the date or datetime expression; with two arguments, the sum of the arguments	N
TIMESTAMPADD()	Add an interval to a datetime expression	Y
TIMESTAMPDIFF()	Subtract an interval from a datetime expression	Y
TO_DAYS()	Return the date argument converted to days	Y
TO_SECONDS()	Return the date or datetime argument converted to seconds since Year 0	Y
UNIX_TIMESTAMP()	Return a Unix timestamp	Y
UTC_DATE()	Return the current UTC date	Y
UTC_TIME()	Return the current UTC time	Y
UTC_TIMESTAMP()	Return the current UTC date and time	Y
WEEK()	Return the week number	Y
WEEKDAY()	Return the weekday index	Y
WEEKOFYEAR()	Return the calendar week of the date (1-53)	Y
YEAR()	Return the year	Y
YEARWEEK()	Return the year and week	Y
CURRENT_DATE	Synonyms for CURDATE()	N
CURRENT_TIME	Synonyms for CURTIME()	N
CURRENT_TIMESTAMP	Synonyms for NOW()	N
LOCALTIME	Synonym for NOW()	N
LOCALTIMESTAMP()	Synonym for NOW()	N

### 3.10.6 Cast Functions and Operators

Name	Description	Support
BINARY	Cast a string to a binary string	N
CAST()	Cast a value as a certain type	Y
CONVERT()	Cast a value as a certain type	Y

#### 3.10.6.1 CAST 不支持以下类型或语法

BINARY  
 CHAR[(N)] [charset\_info] 包含charset\_info时  
 JSON  
 SIGNED [INTEGER] 包含INTEGER时(druid解析问题)  
 UNSIGNED [INTEGER] 包含INTEGER时 (druid解析问题)

### 3.10.6.2 CONVERT 不支持以下类型或语法

BINARY  
 CHAR[(N)] [charset\_info] 包含charset\_info时  
 JSON  
 SIGNED [INTEGER] 包含INTEGER时(druid解析问题)  
 UNSIGNED [INTEGER] 包含INTEGER时 (druid解析问题)

### 3.10.7 Bit Functions and Operators

Name	Description	Support
BIT_COUNT()	Return the number of bits that are set	Y
&	Bitwise AND	Y
~	Bitwise inversion	Y
\		Bitwise OR
^	Bitwise XOR	Y
<<	Left shift	Y
>>	Right shift	Y

### 3.10.8 Aggregate (GROUP BY) Functions

Name	Description	Support
AVG()	Return the average value of the argument	Y
BIT_AND()	Return bitwise AND	Y
BIT_OR()	Return bitwise OR	Y
BIT_XOR()	Return bitwise XOR	Y
COUNT()	Return a count of the number of rows returned	Y
COUNT(DISTINCT)	Return the count of a number of different values	Y
GROUP_CONCAT()	Return a concatenated string	Y
JSON_ARRAYAGG()	Return result set as a single JSON array	N
JSON_OBJECTAGG()	Return result set as a single JSON object	N
MAX()	Return the maximum value	Y
MIN()	Return the minimum value	Y
STD()	Return the population standard deviation	Y
STDDEV()	Return the population standard deviation	Y
STDDEV_POP()	Return the population standard deviation	Y
STDDEV_SAMP()	Return the sample standard deviation	Y
SUM()	Return the sum	Y
VAR_POP()	Return the population standard variance	Y
VAR_SAMP()	Return the sample variance	Y
VARIANCE()	Return the population standard variance	Y

备注: STD 和VARIANCE相关函数, 因为分布式计算的局限性, 精度会有一些问题, 见 [STD\(\) / STDDEV\(\) / STDDEV\\_POP\(\) / STDDEV\\_SAMP\(\) / VAR\\_POP\(\) / VAR\\_SAMP\(\) / VARIANCE\(\)](#) result precision is not correct 在AVG/SUM等相关计算函数的时候, 由于dble的数据来源于MySQL的查询返回, 当MySQL默认返回的数据精度不够时, 可能出现最终查询的结果和MySQL的计算结果有少许差异的情况, 见 [for data type float, dble and mysql may get different results](#)

### 3.10.9 Full-Text Search Functions

not supported

### 3.10.10 XML Functions

not supported

### 3.10.11 Encryption and Compression Functions

not supported

### 3.10.12 Information Functions

not supported

### 3.10.13 Spatial Analysis Functions

not supported

### 3.10.14 JSON Functions

not supported

### 3.10.15 Functions Used with Global Transaction IDs

not supported

### 3.10.16 MySQL Enterprise Encryption Functions

not supported

### 3.10.17 Miscellaneous Functions

not supported

参考资料: [MySQL5.7 函数文档](#)

## 导入导出方式的支持

### 支持工具

1. workbench
- 2.dbeaver
3. mysqldump
4. navicat
5. 导入数据也可以使用mysql中的source和load data

### 注意点

1. 若使用mysqldump导出时, 请按照以下格式进行导出, 否则可能出现错误, 因为有些 mysqldump 参数dble不支持。

```
./mysqldump -h127.0.0.1 -utest -P3306 -p111111 --default-character-set=utf8mb4 --master-data=2 --single-transaction --set-gtid-purged=off --hex-blob --databases schema1 > export.sql
```

1. 导入时, 脚本中若存在非注释性的视图相关语句, 需要注释掉或删除。
2. 导出时, 因为dble对视图相关的一些语句不支持, 因此尽量确保导出的dble中不存在视图。

## 4 协议兼容

- [4.1 基本包](#)
- [4.2 连接建立](#)
- [4.3 文本协议](#)
- [4.4 二进制协议 \(Prepared Statements\)](#)
- [4.5 服务响应包](#)

## 4.1 基本包

- 标准包: 支持
- 大包（**16M以上**）：支持
- 压缩包: 支持, 需要全局配置, 参见 [1.3 server.xml](#)
- 压缩后的大包: 不支持

## 4.2 连接建立

### 4.2.1 Authentication Plugin

- a.空(即默认,等同于mysql\_native\_password 或者8.0以后的caching\_sha2\_password)
- b.mysql\_native\_password
- c.caching\_sha2\_password

### 4.2.2 Capabilities

因为db的连接模型是预先建立后端连接池，所以前端连接设置的一些权能标志位无法生效，这里列举一下权能标志位的使用情况。

当后端连接的标志位使用了之后，无法通过前端连接设置来重置。

如：select @@sql\_mode 会永远包含IGNORE\_SPACE

这样带来的限制是，某些和函数同名的词会当成保留字来处理。参见[MySQL文档](#) 以及相关issue-972

名称	标记值	描述	后端连接设置值	模拟服务端权能位
CLIENT_LONG_PASSWORD	1	Use the improved version of Old Password Authentication.Assumed to be set since 4.1.1.	Y	Y
CLIENT_FOUND_ROWS	2	Send found rows instead of affected rows in EOF_Packet	Y	Y
CLIENT_LONG_FLAG	4	Get all column flags.	Y	Y
CONNECT_WITH_DB	8	Database (schema) name can be specified on connect in Handshake Response Packet.	Y	Y
CLIENT_NO_SCHEMA	16	Don't allow database.table.column.	N	N
CLIENT_COMPRESS	32	Compression protocol supported	server.xml的useCompression选项控制	
CLIENT_ODBC	64	Special handling of ODBC behavior.No special behavior since 3.22.	Y	Y
CLIENT_LOCAL_FILES	128	Can use LOAD DATA LOCAL.	Y	Y
CLIENT_IGNORE_SPACE	256	Ignore spaces before '('.	Y	Y
CLIENT_PROTOCOL_41	512	New 4.1 protocol	Y	Y
CLIENT_INTERACTIVE	1024	This is an interactive client.	Y	Y
CLIENT_SSL	2048	Use SSL encryption for the session	N	N
CLIENT_IGNORE_SIGPIPE	4096	Client only flag.Not used.	Y	Y
CLIENT_TRANSACTIONS	8192	Client knows about transactions	Y	Y
CLIENT_RESERVED	16384	DEPRECATED:Old flag for 4.1 protocol.	N	N
CLIENT_RESERVED2	32768	DEPRECATED:Old flag for 4.1 authentication.	Y	Y
CLIENT_MULTI_STATEMENTS	65536	Enable/disable multi-stmt support	Y	Y
CLIENT_MULTI_RESULTS	131072	Enable/disable multi-results	Y	Y
CLIENT_PS_MULTI_RESULTS	262144	Multi-results and OUT parameters in PS-protocol	N	N
CLIENT_PLUGIN_AUTH	524288	Client supports plugin authentication.	N	Y
CLIENT_CONNECT_ATTRS	1048576	Client supports connection attributes.	N	N
CLIENT_PLUGIN_AUTH_LENENC_CLIENT_DATA	2097152	Enable authentication response packet to be larger than 255 bytes.	N	N
CLIENT_CAN_HANDLE_EXPIRED_PASSWORDS	4194304	Don't close the connection for a user account with expired password.	N	N
CLIENT_SESSION_TRACK	8388608	Capable of handling server state change information.	N	N
CLIENT_DEPRECATED_EOF	16777216	Client no longer needs EOF_Packet and will use OK_Packet instead.	N	N
CLIENT_SSL_VERIFY_SERVER_CERT	1UL << 30	Verify server certificate	N	N
CLIENT_REMEMBER_OPTIONS	1UL << 31	Don't reset the options after an unsuccessful connect.	N	N

权能标志位定义：

参考 [https://dev.mysql.com/doc/dev/mysql-server/8.0.13/group\\_group\\_cs\\_capabilities\\_flags.html](https://dev.mysql.com/doc/dev/mysql-server/8.0.13/group_group_cs_capabilities_flags.html)

## 4.3 文本协议

### 4.3.1 Supported

- COM\_INIT\_DB  
Specifies the default schema for the connection.
- COM\_PING  
Sends a packet containing one byte to check that the connection is active.
- COM\_QUERY  
Sends the server an SQL statement to be executed immediately. Support Multi-Statement.
- COM\_QUIT  
Client tells the server that the connection should be terminated.
- COM\_SET\_OPTION  
Enables or disables server option.
- COM\_CHANGE\_USER  
Resets the connection and re-authenticates with the given credentials.
- COM\_RESET\_CONNECTION  
Resets a connection without re-authentication.
  - 关闭后端连接(rollback & unlock)
  - 事务状态情况
  - 用户变量清空
  - 系统变量恢复成系统默认值
  - prepare清空
  - 上下文(字符集, 隔离级别)恢复成为默认值
  - LAST\_INSERT\_ID 置零

#### 4.3.1.1 Multi-Statement

- Supported
  - DML:select/insert/update/replace/delete
  - DDL
  - OTHER
    - BEGIN;
    - COMMIT;
    - LOCK TABLE
    - UNLOCK TABLES
    - START
    - KILL
    - USE
    - ROLLBACK
    - MYSQL\_CMD\_COMMENT
    - MYSQL\_COMMENT
    - SELECT VERSION\_COMMENT (SELECT @@VERSION\_COMMENT)
    - SELECT DATABASE (select database())
    - SELECT USER (select user())
    - SELECT VERSION (select version())
    - SELECT SESSION\_INCREMENT(select @@session.auto\_increment\_increment)
    - SELECT SESSION\_ISOLATION(select @@session.tx\_isolation)
    - SELECT LAST\_INSERT\_ID(select last\_insert\_id() as id)
    - SELECT IDENTITY(select @@identity)
    - SELECT SESSION\_TX\_READ\_ONLY (select @@session.tx\_read\_only)
- Not Supported
  - EXPLAIN
  - EXPLAIN2
  - DESCRIBE
  - SET
  - SHOW DATABASES/TABLES/TABLE\_STATUS/COLUMNS/INDEX/CREATE\_TABLE/VARIABLES/CREATE\_VIEW/CHARSET
  - HELP
  - LOAD\_DATA\_INFILE\_SQL
  - CREATE\_VIEW
  - REPLACE\_VIEW
  - ALTER\_VIEW
  - DROP\_VIEW

### 4.3.2 Not Supported

- COM\_DEBUG  
Forces the server to dump debug information to stdout
- COM\_STATISTICS  
Get internal server statistics.
- COM\_CREATE\_DB
- COM\_DROP\_DB

### 4.3.3 Internal

- COM\_SLEEP  
Used inside the server only.
- COM\_CONNECT an internal command in the server.
- COM\_TIME an internal command in the server.
- COM\_DAEMON an internal command in the server.
- COM\_DELAYED\_INSERT an internal command in the server.

### 4.3.4 Deprecated

- COM\_PROCESS\_INFO  
Deprecated from 5.7.11.
- COM\_PROCESS\_KILL  
Deprecated from 5.7.11.
- COM\_FIELD\_LIST  
Deprecated from 5.7.11.
- COM\_SHUTDOWN  
Deprecated from 5.7.9.
- COM\_REFRESH  
Deprecated from 5.7.11.

## 4.4 二进制协议 (Prepared Statements)

### 4.4.1 Supported

- COM\_STMT\_CLOSE  
Closes a previously prepared statement.
- COM\_STMT\_EXECUTE  
Executes a previously prepared statement.
- COM\_STMT\_RESET  
Resets a prepared statement on client and server to state after preparing.
- COM\_STMT\_SEND\_LONG\_DATA  
When data for a specific column is big, it can be sent separately.

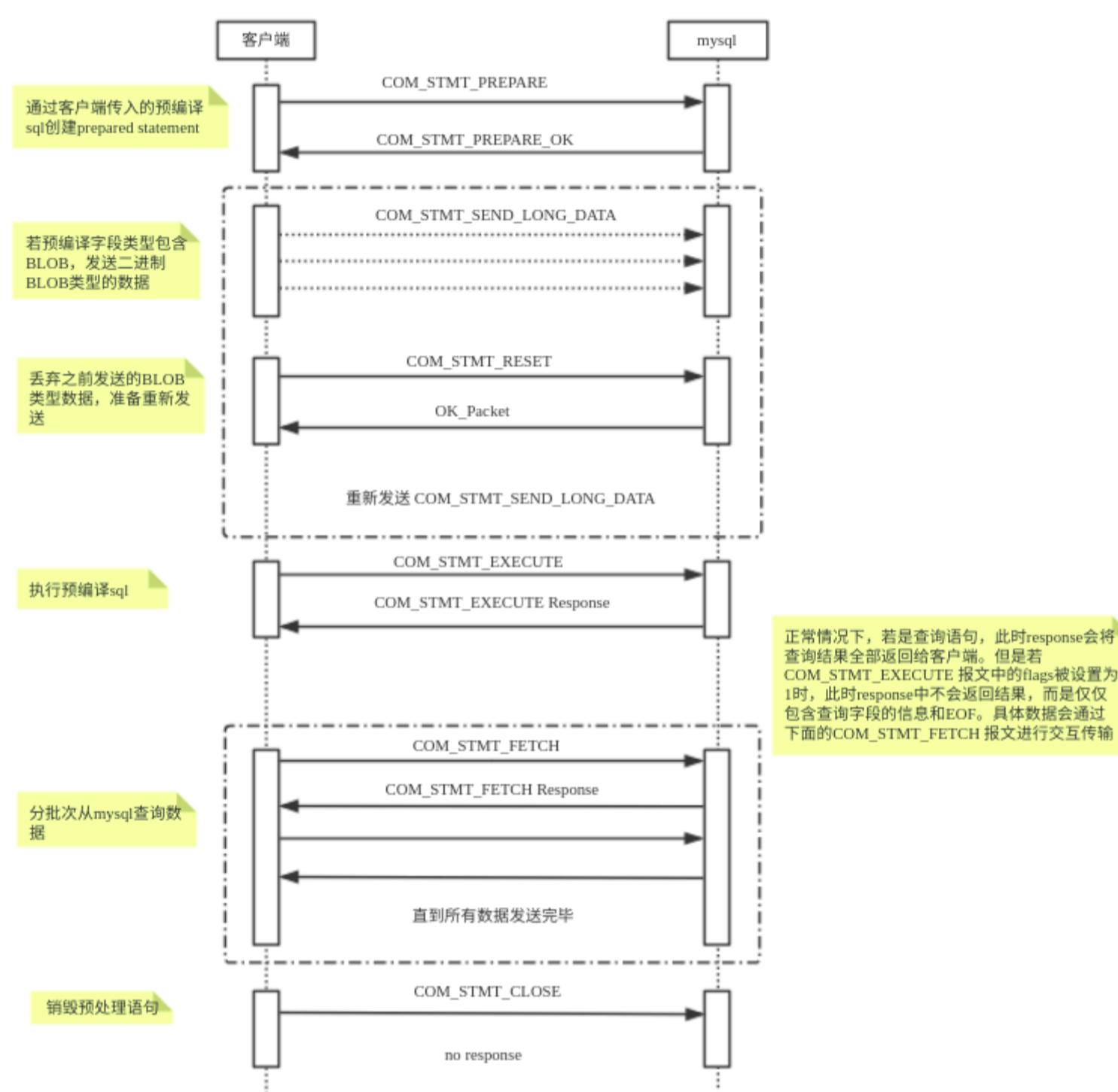
### 4.4.2 Faked

- COM\_STMT\_PREPARE  
Prepares a statement on the server  
COM\_STMT\_PREPARE response is faked. Parameters and numbers are faked.

### 4.4.3 Not Supported

- COM\_STMT\_FETCH  
Fetches rows from a prepared statement

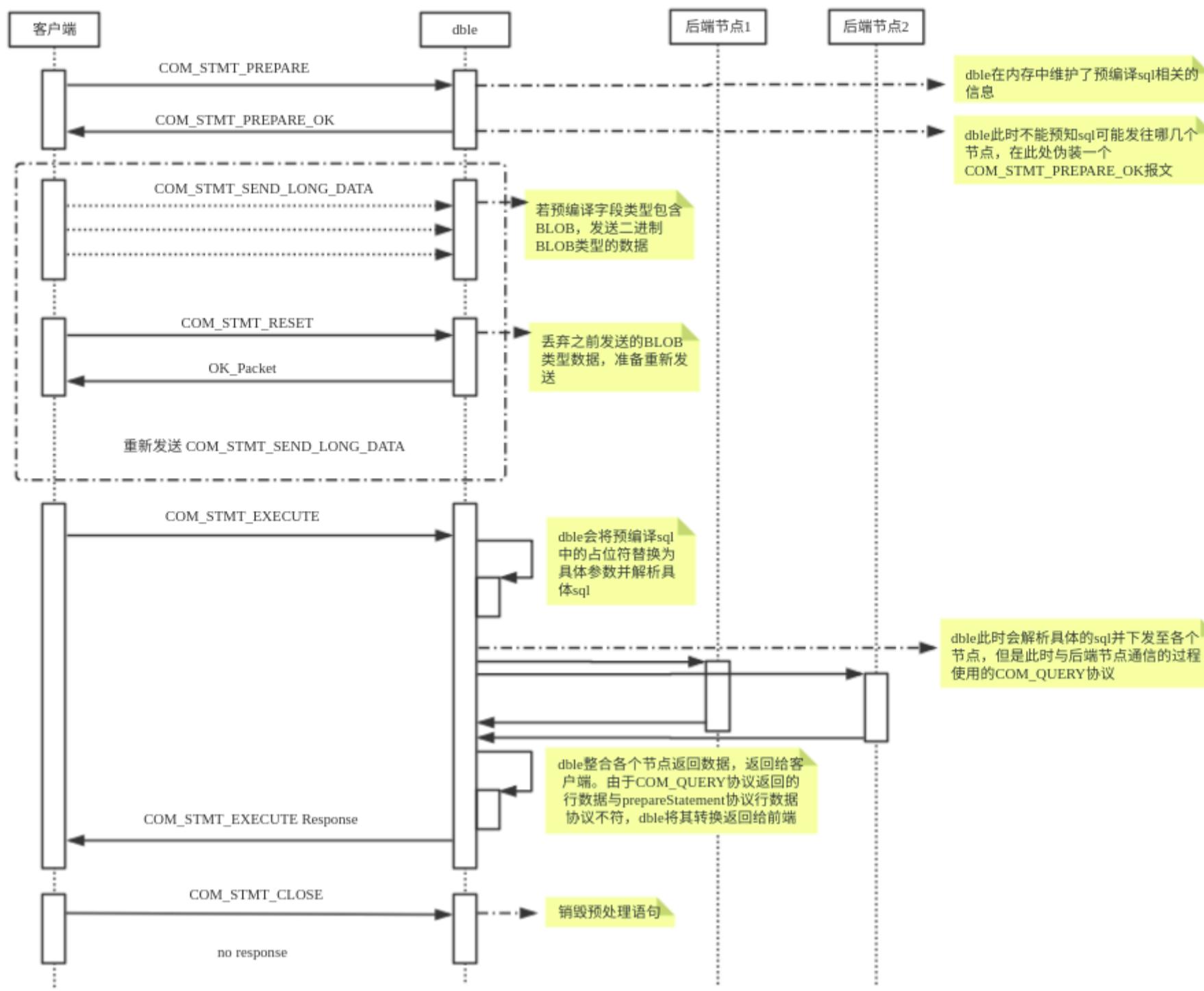
### 4.4.4 MySQL Flow



注意点:

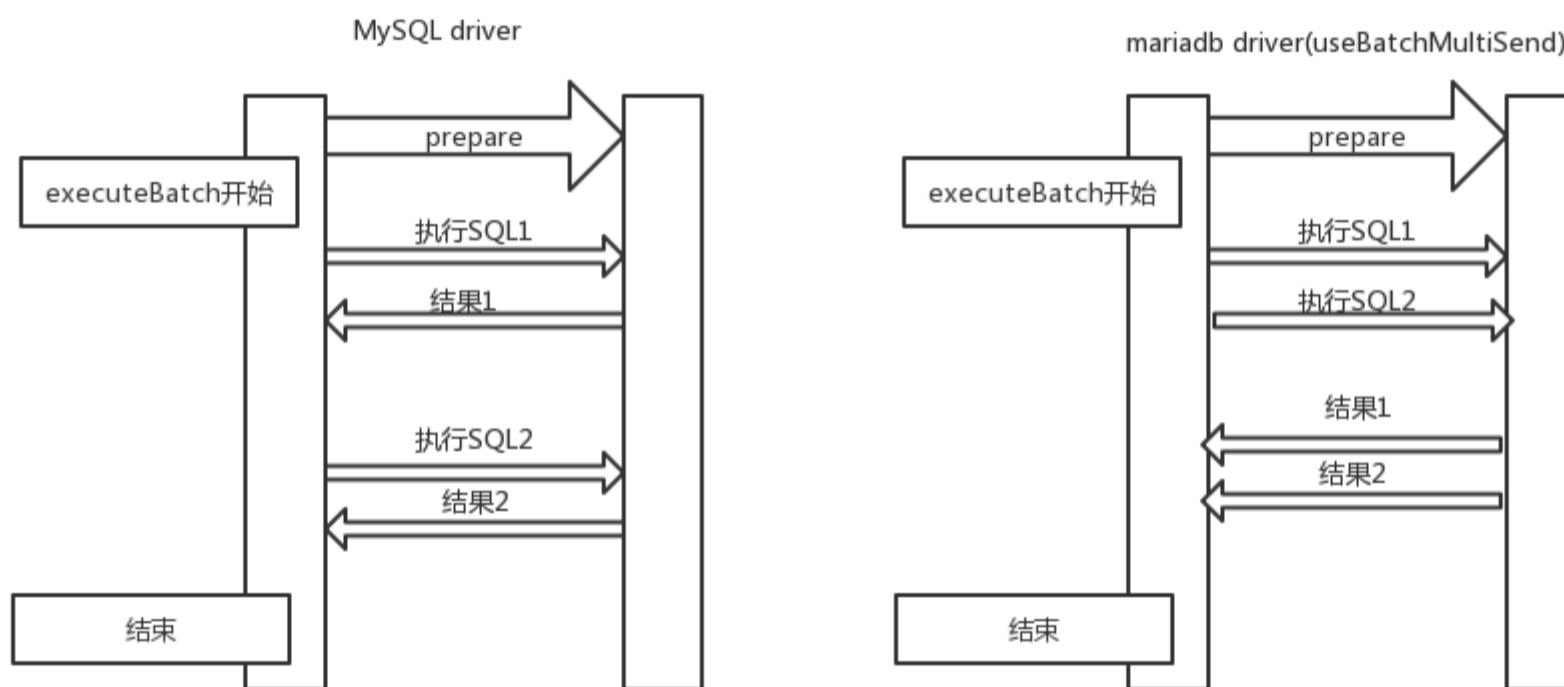
- 可通过url方式指定useCursorFetch=true, 开启分批从server查询数据。
- 但是jdbc中默认fetchSize为0, jdbc中必须fetchSize > 0才会发送fetch包分批查询数据, 否则和普通prepareStatement没有区别。

### 4.4.5 Dble Flow



#### 4.4.6 额外限制

在使用mariadb-java-client版本1.5.0及以上版本连接dble的时候，executeBatch使用的时候存在注意事项，需要禁用driver的useBatchMultiSend选项，否则会出现不固定的错误。原因是在mariadb-java-client版本1.5.0及以上版本中添加了一个新的优化功能useBatchMultiSend，具体的优化在executeBatch的时候driver在未收到上一条数据执行结果的情况下提前下发后续的执行语句，这个优化与当前版本的dble内部实现不兼容。大致情况如下图所示：



原文地址：<https://mariadb.com/kb/en/library/option-batchmultisend-description/>

## 4.5 服务响应包

兼容支持:

- EOF\_Packet
- ERR\_Packet
- OK\_Packet
- LOCAL\_INFILE Packet
- PACKET\_LOCAL\_INFILE
- PACKET\_RESULTSET

## 5. 已知限制

- [5.1 druid引发的限制](#)
- [5.2 其他已知限制](#)

## 5.1 druid引发的限制

1. INSERT ... VALUE ... 和 INSERT ... VALUES ...这个两个语句中, **VALUE[S]**关键字必须写正确, 否则会出现解析正常, 但结果错误。  
这个是由于druid的bug引起的。当VALUES[S]关键字错误时, druid会把这个错误的关键字当作别名处理。  
已向druid提[issue2218](#)。  
dble相关issue [dble\\_issue\\_379](#).
2. 其他  
参见 [dble\\_issue\\_788](#)

## 5.2 其他已知限制

1. 显式配置的父子表，在子表插入数据时，不能多值插入  
原因：子表插入数据时，如有parentkey不是父表的拆分列。需要去对应父表的拆分结点中反向查询路由规则，如果此时是多值插入，就会变成多值反查，查询功能较难完成，即使完成，性能也会很差。  
相关issue：<https://github.com/actiontech/dble/issues/12>
2. 不建议使用JDBC的rewriteBatchedStatements=true  
原因：  
insert：多条简单insert拼接成一条长的insert.. values(),一个com\_query包，对于dble来说，可能引发分布式事务，降低性能和数据一致性。
3. 使用JDBC的useServerPrepStmts=true会降低性能  
原因：dble是将前端的Binary Protocol 转为Text Protocol，收到结果集之后再反向转回协议，所以需要额外的工作，降低了性能。
4. lock/unlock 实现不完整  
相关issue：<https://github.com/actiontech/dble/issues/38>
5. 不支持在schema.xml里配置复合主键  
原因：schema.xml里的主键作用为主键路由和全局序列，这两者暂不支持多主键。  
相关issue：<https://github.com/actiontech/dble/issues/70>
6. 不保证主键唯一性 由于dble不要求配置在table中的主键在创建表格时一定为主键，并且由于dble后端分布式存储的方式，dble不对于主键做唯一性检查，并且允许用户对于非分片的主键字段进行任意更新
7. 并发更新多行数据/全局表数据可能导致死锁超时  
原因：并发情况下，分布式下发sql可能乱序。  
相关issue：<https://github.com/actiontech/dble/issues/85>
8. 防火墙导致无响应  
原因：防火墙可能drop包，java层面的tcp\_keepalive无法指定时间。  
相关issue：<https://github.com/actiontech/dble/issues/87>
9. 方差/标准差精度问题  
原因：方差计算方式导致  
相关issue：<https://github.com/actiontech/dble/issues/100>
10. order by lock in share mode/for update , lock clause is ignored  
原因：无法支持。  
相关issue：<https://github.com/actiontech/dble/issues/127>
11. 不支持 \_charset\_name 'string' \_charset\_name+b'val'  
相关issue：<https://github.com/actiontech/dble/issues/262>  
相关issue：<https://github.com/actiontech/dble/issues/267>
12. 未能正确支持 set sql\_select\_limit  
相关issue：<https://github.com/actiontech/dble/issues/331>
13. 日期拆分算法中，sEndDate如果不配置，default node就无用  
原因：算法设计问题。  
相关issue：<https://github.com/actiontech/dble/issues/357>
14. selece @@sql\_mode 始终包含IGNORE\_SPACE  
原因：后端权限标志位设置导致，参见4.2节内容。  
相关issue：<https://github.com/actiontech/dble/issues/364>
15. replace ... into  
由于replace的语义为如果存在则替换，如果不存在则新增，所以在使用表格自增主键的时候 如果对于自增表格使用replace且ID不存在，那么就会插入一条指定ID的数据，并不会自动生成ID
16. kill 语句杀自身session，直接返回ok，不会有任何实质性操作
17. 由于2.19.01版本在rule/schema/server.xml中加入了version字段造成的BUG，导致在2.19.01使用zk集群进行同步化操作的时候要求version字段不能为空或者不填，此问题在后续版本会进行修复
18. 若mysql节点上设置了 set global local\_infile = 0 ,dble load data指令执行报错  
原因：dble会将load data指令转换为 load data local infile ... 指令下发至各个后端mysql，所以各个节点 local\_infile 参数需要开启。  
相关issue：<https://github.com/actiontech/dble/issues/1111>
19. 不能正确支持 set @@sql\_auto\_is\_null=on; 原因：set @@sql\_auto\_is\_null=on 的语义是用新生成的自增序列取代null,在dble层不做实现。  
相关issue：<https://github.com/actiontech/dble/issues/978>
20. 不能正确支持mariadb-java-client版本1.5.0及以上版本的useBatchMultiSend优化 原因： mariadb-java-client版本1.5.0及以上版本修改了prepare执行的时候的交互逻辑 相关issue：<https://github.com/actiontech/dble/issues/1244>
21. 复杂查询会透传难以理解的报错，建议结合explain语句分析 原因： 需要枚举所有错误来做替换.成本高，收益低 相关issue：<https://github.com/actiontech/dble/issues/1449>

## 6. 与MySQL Server的差异化描述

这里主要描述一些与MySQL Server不同的行为，这些行为不是bug，是分布式场景下的一些正常的行为表现，但与已知的MySQL行为不一致。

- [6.1 事务中遇到主键冲突需要显式回滚](#)
- [6.2 INSERT不能显式指定自增序列](#)
- [6.3 增加"show all tables"](#)
- [6.4 去除了增删改的message信息](#)
- [6.5 information\\_schema等库的支持](#)

## 6.1 事务中遇到主键冲突需要显式回滚，MySQL不需要

具体表现如下：

MySQL行为：

```
[test_yhq]>select * from char_columns_4;
+----+-----+
| id | c_char |
+----+-----+
| 1 | xx |
| 4 | z |
+----+-----+
2 rows in set (0.02 sec)

[test_yhq]>begin;
Query OK, 0 rows affected (0.01 sec)

[test_yhq]>insert into char_columns_4 values(1,'yy');
ERROR 1062 (23000): Duplicate entry '1' for key 'PRIMARY'
[test_yhq]>insert into char_columns_4 values(2,'yy');
Query OK, 1 row affected (0.00 sec)

[test_yhq]>commit;
Query OK, 0 rows affected (0.02 sec)
```

db表行为：

```
[testdb]>select * from sharding_four_node order by id;
+----+-----+-----+
| id | c_flag | c_decimal |
+----+-----+-----+
| 1 | 1_1   | 1.0000 |
| 2 | 2      | 2.0000 |
| 3 | 3      | 3.0000 |
+----+-----+-----+
3 rows in set (0.28 sec)

begin;
Query OK, 0 rows affected (0.01 sec)

[testdb]>insert into sharding_four_node values(1,'1',1.0);
ERROR 1062 (23000): Duplicate entry '1' for key 'PRIMARY'
[testdb]>insert into sharding_four_node values(13,'13',13.0);
ERROR 1003 (HY000): Transaction error, need to rollback.Reason:[ errNo:1062 Duplicate entry '1' for key 'PRIMARY']
[testdb]>commit;
ERROR 1003 (HY000): Transaction error, need to rollback.Reason:[ errNo:1062 Duplicate entry '1' for key 'PRIMARY']
```

## 6.2 INSERT自增序列表由db1e生成，不能显式指定自增列，MySQL可以。

具体表现如下：

MySQL行为：

```
desc mysql_autoinc;
+-----+-----+-----+-----+
| Field | Type   | Null | Key | Default | Extra
+-----+-----+-----+-----+
| c_char | char(255) | YES |    | NULL    |          |
| id     | bigint(20) | NO  | PRI | NULL    | auto_increment |
+-----+-----+-----+-----+
2 rows in set (0.02 sec)

[test_yhq]>insert into mysql_autoinc values('1',1);
Query OK, 1 row affected (0.01 sec)
```

db1e行为

```
desc sharding_four_node_autoinc;
+-----+-----+-----+-----+
| Field | Type   | Null | Key | Default | Extra
+-----+-----+-----+-----+
| c_char | char(255) | YES |    | NULL    |          |
| id     | bigint(20) | NO  | PRI | NULL    | auto_increment |
+-----+-----+-----+-----+
2 rows in set (0.01 sec)

[testdb]>insert into sharding_four_node_autoinc values('2',2);
ERROR 1064 (HY000): In insert Syntax, you can't set value for Autoincrement column!
```

### 6.3 ADD "show all tables"

The optional ALL modifier causes SHOW TABLES to display a second output column with values of BASE TABLE for a table ,VIEW for a view, SHARDING TABLE for a sharding table and GLOBAL TABLE for a global table.

具体表现如下:

```
[testdb]>show all tables;
+-----+-----+
| Tables in testdb | Table_type |
+-----+-----+
| global_four_node | GLOBAL TABLE |
| global_four_node_autoinc | GLOBAL TABLE |
| global_two_node | GLOBAL TABLE |
| sbtest1 | SHARDING TABLE |
| sharding_four_node | SHARDING TABLE |
| sharding_four_node2 | SHARDING TABLE |
| sharding_four_node_autoinc | SHARDING TABLE |
| sharding_two_node | SHARDING TABLE |
| single | SHARDING TABLE |
| customer | BASE TABLE |
| district | BASE TABLE |
+-----+-----+
11 rows in set (0.02 sec)
```

## 6.4 去除了增删改的message 信息

具体表现如下:

MySQL行为:

```
mysql> insert into sharding_two_node values(9,'9',9.0),(10,'10',10.0);
Query OK, 2 rows affected (0.24 sec)
Records: 2  Duplicates: 0  Warnings: 0
```

dble行为:

```
mysql> insert into sharding_two_node values(11,'11',11.0),(12,'12',12.0);
Query OK, 2 rows affected (0.49 sec)
```

## 6.5 information\_schema等库的支持

### 背景

用户使用 Navicat Premium 12 连接db时，Navicat Premium 12会查询 information\_schema, mysql 等数据库中系统表的数据。db 需要支持查询这些系统表的语句，保证这些 driver 的正常使用。

### Navicat Premium12

1. SELECT SCHEMA\_NAME, DEFAULT\_CHARACTER\_SET\_NAME, DEFAULT\_COLLATION\_NAME FROM information\_schema.SCHEMATA;

该语句查询当前 mysql 实例中的所有 schema 的名称，character set 以及 collation。该语句会影响 Navicat Premium 12 的使用。

db 中的 schema 是逻辑上的，可以通过 SchemeConfig 获取所有 schema 的名称，schema 的 character set 以及 collation 使用默认返回。

问题1：db 中的 schema 是逻辑上的，会对应多个 DataNode，会出现 DataNode 的 character set 以及 collation 与 默认 character set 和 collation 不同，需要运维安装 MySQL 时候保证。

问题2：当前 db 中处理 SCHEMATA 系统表时，只是对表名做了判断，并未对查询字段做检验。

### 以下语句不影响driver的使用

对此种语句的处理是根据请求字段，伪造一个行数为0的报文返回。

```

1. SELECT TABLE_SCHEMA, TABLE_NAME, TABLE_TYPE
   FROM information_schema.TABLES WHERE TABLE_SCHEMA = 'testdb'
   ORDER BY TABLE_SCHEMA, TABLE_TYPE
2. SELECT TABLE_SCHEMA, TABLE_NAME, COLUMN_NAME, COLUMN_TYPE
   FROM information_schema.COLUMNS
   WHERE TABLE_SCHEMA = 'testdb'
   ORDER BY TABLE_SCHEMA, TABLE_NAME
3. SELECT DISTINCT ROUTINE_SCHEMA, ROUTINE_NAME, PARAMS.PARAMETER
   FROM information_schema.ROUTINES LEFT JOIN
   ( SELECT SPECIFIC_SCHEMA, SPECIFIC_NAME,
      GROUP_CONCAT(CONCAT(DATA_TYPE, ' ', PARAMETER_NAME) ORDER BY ORDINAL_POSITION SEPARATOR ', ') PARAMETER, ROUTINE_TYPE
   FROM information_schema.PARAMETERS GROUP BY SPECIFIC_SCHEMA, SPECIFIC_NAME, ROUTINE_TYPE
   )PARAMS
   ON ROUTINES.ROUTINE_SCHEMA = PARAMS.SPECIFIC_SCHEMA AND
   ROUTINES.ROUTINE_NAME = PARAMS.SPECIFIC_NAME AND
   ROUTINES.ROUTINE_TYPE = PARAMS.ROUTINE_TYPE
   WHERE ROUTINE_SCHEMA = 'testdb' ORDER BY ROUTINE_SCHEMA
4. SELECT TABLE_NAME, CHECK_OPTION, IS_UPDATABLE, SECURITY_TYPE, DEFINER
   FROM information_schema.VIEWS
   WHERE TABLE_SCHEMA = 'testdb' ORDER BY TABLE_NAME ASC
5. SELECT * FROM information_schema.ROUTINES
   WHERE ROUTINE_SCHEMA = 'testdb' ORDER BY ROUTINE_NAME
6. SELECT EVENT_CATALOG, EVENT_SCHEMA, EVENT_NAME, DEFINER, TIME_ZONE,
   EVENT_DEFINITION, EVENT_BODY, EVENT_TYPE, SQL_MODE, STATUS, EXECUTE_AT,
   INTERVAL_VALUE, INTERVAL_FIELD, STARTS, ENDS, ON_COMPLETION, CREATED,
   LAST_ALTERED, LAST_EXECUTED, ORIGINATOR, CHARACTER_SET_CLIENT,
   COLLATION_CONNECTION, DATABASE_COLLATION, EVENT_COMMENT
   FROM information_schema.EVENTS WHERE EVENT_SCHEMA = 'testdb'
   ORDER BY EVENT_NAME ASC
7. SELECT COUNT(*) FROM information_schema.TABLES
   WHERE TABLE_SCHEMA = 'testdb' UNION
   SELECT COUNT(*)
   FROM information_schema.COLUMNS
   WHERE TABLE_SCHEMA = 'testdb' UNION
   SELECT COUNT(*) FROM information_schema.ROUTINES WHERE ROUTINE_SCHEMA = 'testdb'
```

## 7 开发者须知

- [7.1 SQL开发编写原则](#)
- [7.2 dble连接Demo](#)
- [7.3 其他注意事项](#)

## 7.1 SQL开发编写原则

- 分布式数据库处理的是分布式关系运算，其SQL的优化方法与单机关系数据库有所不同，侧重考虑的分布式环境中的网络开销。分布式执行计划中，应该尽量将SQL中的运算下推到底层各个节点执行，避免跨节点运算，从而减少网络开销、提升SQL执行效率。

### 1) 执行计划

访问数据时的一组有序的操作步骤集合，称为执行计划。**dble**的执行计划分为两个层次：**dble**层的执行计划与节点层的执行计划。对执行计划进行分析，可以了解中间件和节点是否对SQL语句生成了最优的执行计划，是否有优化的空间，从而为SQL优化提供重要的参考信息。在SQL语句执行前，**dble**会根据SQL语句的基本信息，判断该SQL语句应该在哪些节点上执行，将SQL改写成在节点上执行的具体形式，并决定采用何种策略进行数据合并与计算等，这就是**dble**层的执行计划。节点层的执行计划就是原生的MySQL执行计划。**dble**用EXPLAIN指令来查看**dble**层的执行计划。如：

```
explain select id,accountno from account where userid=2;
```

EXPLAIN指令的执行结果包括语句下发的节点，实际下发的SQL语句和数据的合并操作的信息。这些信息是系统静态分析产生的，并没有真正的执行语句。通过EXPLAIN2命令可查看指定节点上的执行计划。如：

```
explain2 datanode=dn1 sql=select id,accountno from account where userid=2;
```

explain2会将sql语句加上explain下发到指定的datanode执行，并把节点上explain的结果返回调用者。

### 2) SQL优化

对于SQL优化有以下原则：

- SQL语句中尽可能带有拆分字段，并且拆分字段过滤条件的取值范围越小，越有助于提高查询速度。在数据写入时，必须给拆分字段赋值。
- 拆分字段的查询条件尽可能为等值条件。
- 如果拆分字段的条件是IN子句，则IN后面的值的数目应尽可能少。特别注意，随着业务增长，某些IN子句的条件会随之增长。
- 如果SQL语句不带有拆分字段，那么DISTINCT、GROUP BY和ORDER BY在同一个SQL语句中尽量只出现一种。
- 数据查询时，应该尽量减少节点返回的结果数量，这样能够使消耗的网络带宽最小，使查询性能能够达到最优状态。

### 3) 跨节点查询

跨节点查询通常发生在下列语句中：

- 涉及多个分片的分布式Join；
- 涉及多个分片的聚合函数；
- 复杂子查询。

这些语句可以通过以下策略来优化：

- 合理调整查询语句，使查询条件中包含拆分字段的等值条件，这样就能够将语句下推到单一节点执行。
- 参与Join的表配置相同的拆分规则，查询时将拆分字段作为Join的关联条件，这样Join操作就可以在节点内完成。
- 将参与Join的表中数据量较小的表配置成全局表，通过数据冗余避免跨节点。
- 如果一定需要跨节点Join，尽量对驱动表添加更多的过滤条件，从而使参与跨节点Join的数据量尽可能的少。
- 如果查询涉及到了多个节点，则尽量不要使用limit a,b这样的子句。
- GROUP子句尽量包含拆分字段。
- 对语句进行改写，将复杂子查询分解为多条语句来执行。
- 子查询尽可能改写成Join的形式。

## 7.2 dble连接Demo

### 开发框架连接

#### ibatis

利用ibatis连接dble时，连接方式与MySQL相同。下面是一个简单示例。 JDBC配置：

```
jdbc.driverClass=com.mysql.jdbc.Driver
jdbc.jdbcUrl=jdbc:mysql://127.0.0.1:8066/TESTDB?useUnicode=true&characterEncoding=utf-8
jdbc.user=root
jdbc.password=123456
```

映射文件配置：

```
<?xml version="1.0" encoding="UTF-8" ?>
<!DOCTYPE mapper
PUBLIC "-//mybatis.org//DTD Mapper 3.0//EN"
"http://mybatis.org/dtd/mybatis-3-mapper.dtd">
<mapper namespace="com.mapper.UserMapper">
<insert id="saveUser" parameterType="com.bean.User">
    insert into user(id,name,phone,birthday)
    values (0,#{name},#{phone},#{birthday})
    <selectKey keyProperty="id" order="after" resultType="int">
        select last_insert_id() as id
    </selectKey>
</insert>
<delete id="deleteUserById" parameterType="java.lang.String">
    delete from user where id=#{id}
</delete>
<update id="updateUser" parameterType="com.bean.User">
    update user set name=#{name},phone=#{phone},birthday=#{birthday} where id=#{id}
</update>
<update id="updateUsers">
    /*!dble:sql=select * from user;*/update users set usercount=(select count(*) from user),ts=now()
</update>
<select id="getUserById" parameterType="java.lang.String" resultType="com.bean.User">
    select * from user where id=#{id}
</select>
<select id="getUsers" resultType="com.bean.User">
    select * from user
</select>
```

语句select last\_insert\_id() as id可用来获取新写入记录的ID。 updateUser方法用到了 dble的注解，由于ibatis中的符号#具有特殊含义，因此注解中不能含有#。

#### hibernate

利用hibernate连接dble时，连接方式与MySQL相同。下面是一个简单示例。 hibernate.cfg.xml:

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE hibernate-configuration PUBLIC
"-//Hibernate/Hibernate Configuration DTD 3.0//EN"
"http://www.hibernate.org/dtd/hibernate-configuration-3.0.dtd">
<hibernate-configuration>
<session-factory>
<property name="hibernate.connection.driver_class">com.mysql.jdbc.Driver</property>
<property name="hibernate.connection.url">jdbc:mysql://192.168.58.51:8066/testdb?useUnicode=true&characterEncoding=utf-8</property>
<property name="hibernate.connection.username">root</property>
<property name="hibernate.connection.password">123456</property>
<property name="hibernate.dialect">org.hibernate.dialect.MySQLInnoDBDialect</property>
<property name="hibernate.format_sql">true</property>
<property name="hibernate.hbm2ddl.auto">update</property>
<mapping resource="com/actiontech/test/News.hbm.xml"/>
</session-factory>
</hibernate-configuration>
```

News.hbm.xml:

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE hibernate-mapping PUBLIC
"-//Hibernate/Hibernate Mapping DTD 3.0//EN"
"http://www.hibernate.org/dtd/hibernate-mapping-3.0.dtd">
<hibernate-mapping>
<class name="com.actiontech.test.News" table="news_table">
<id name="id" type="java.lang.Integer">
    <column name="id" />
</id>
<property name="title" type="java.lang.String">
    <column name="title" />
</property>
<property name="content" type="java.lang.String">
    <column name="content" />
</property>
</class>
</hibernate-mapping>
```

News.java:

```

package com.actiontech.test;
public class News {
    private Integer id;
    private String title;
    private String content;
    public Integer getId() {
        return id;
    }
    public void setId(Integer id) {
        this.id = id;
    }
    public String getTitle() {
        return title;
    }
    public void setTitle(String title) {
        this.title = title;
    }
}

public String getContent() { return content; } public void setContent(String content) { this.content = content; } }</pre> NewsManager.java:

package com.actiontech.test;
import org.hibernate.Session;
import org.hibernate.SessionFactory;
import org.hibernate.Transaction;
import org.hibernate.cfg.Configuration;
public class NewsManager {
    public static void main(String[] args)
        throws Exception {
        Configuration config = new Configuration().configure();
        SessionFactory factory = config.buildSessionFactory();
        Session session = factory.openSession();
        Transaction transaction = session.beginTransaction();
        News news = new News();
        news.setId(10);
        news.setTitle("dble示例");
        news.setContent("Hibernate 连接dble的第一个例子");
        session.save(news);
        transaction.commit();
        session.close();
        factory.close();
    }
}

```

dble虽然支持Hibernate但不建议使用Hibernate，因为Hibernate无法控制SQL的生成，无法做到对查询SQL的优化，大量下可能会出现性能问题。

## JDBC

利用JDBC连接dble时，连接方式与MySQL相同。下面是一个简单示例：

```

package com.actiontech.test;
import java.sql.Connection;
import java.sql.DriverManager;
import java.sql.PreparedStatement;
import java.sql.ResultSet;
import java.sql.SQLException;
import java.sql.Statement;
import java.util.Properties;
import java.util.concurrent.CountDownLatch;
import java.util.concurrent.atomic.AtomicLong;
public class SingleMixEngine {
    public static void main(String[] args) throws Exception {
        Class.forName("com.mysql.jdbc.Driver");
        Properties props = new Properties();
        props.setProperty("user", "root");
        props.setProperty("password", "123456");
        SingleMixEngine engine = new SingleMixEngine();
        engine.execute(props, "jdbc:mysql://192.168.58.51:8066/testdb");
    }
    final AtomicLong tml = new AtomicLong();
    final String tableName="news_table";
    public void execute(Properties props,String url) {
        CountDownLatch cdl = new CountDownLatch(1);
        long start = System.currentTimeMillis();
        for (int i = 0; i < 1; i++) {
            TestThread insertThread = new TestThread(props,cdl, url);
            Thread t = new Thread(insertThread);
            t.start();
            System.out.println("Test start");
        }
        try {
            cdl.await();
            long end = System.currentTimeMillis();
            System.out.println("Test end,total cost:" + (end-start) + "ms");
        } catch (Exception e) {
        }
    }
    class TestThread implements Runnable {
        Properties props;
        private CountDownLatch countDownLatch;
        String url;
        public TestThread(Properties props,CountDownLatch cdl, String url) {
            this.props = props;
            this.countDownLatch = cdl;
            this.url = url;
        }
        public void run() {
            Connection connection = null;
            PreparedStatement ps = null;
            Statement st = null;
            try {
                connection = DriverManager.getConnection(url,props);
                connection.setAutoCommit(true);
                st = connection.createStatement();
                String dropSql = "drop table if exists " + tableName;
                System.out.println("Execute SQL:\n\t"+dropSql);
                st.execute(dropSql);

                String createSql = "create table " + tableName + "(id int,title varchar(20),content varchar(50))";
                System.out.println("Execute SQL:\n\t"+createSql);
                st.execute(createSql);

                String insertSql = "insert into " + tableName + " (id,title,content) values(?, ?, ?)";
                System.out.println("Prepared SQL:\n\t"+insertSql);
                ps = connection.prepareStatement(insertSql);
                for (int i = 1; i <= 3; i++) {
                    ps.setInt(i,i);
                    ps.setString(2, "测试"+i);
                    ps.setString(3, "这是第"+i+"条测试数据");
                    ps.execute();
                    System.out.println("Insert data:\t"+i+", "+ "测试"+i+", "+ "这是第"+i+"条测试数据");
                }

                String querySQL = "select * from " + tableName + " order by id";
                System.out.println("Execute SQL:\n\t"+querySQL);
                ResultSet rs = st.executeQuery(querySQL);
                int colcount = rs.getMetaData().getColumnCount();
                System.out.println("Current Data:");
                while(rs.next()){
                    for(int i=1;i<colcount;i++){
                        System.out.print("\t"+rs.getString(i));
                    }
                    System.out.println();
                }

                String updateSql = "update " + tableName + " set title='test1' where id=1";
                System.out.println("Execute SQL:\n\t"+updateSql);
                st.execute(updateSql);
                rs = st.executeQuery(querySQL);
                System.out.println("Current Data:");
                while(rs.next()){
                    for(int i=1;i<colcount;i++){
                        System.out.print("\t"+rs.getString(i));
                    }
                    System.out.println();
                }

                String deleteSql = "delete from " + tableName + " where id=2";
                System.out.println("Execute SQL:\n\t"+deleteSql);
                st.execute(deleteSql);
                rs = st.executeQuery(querySQL);
                System.out.println("Current Data:");
                while(rs.next()){
                    for(int i=1;i<colcount;i++){
                        System.out.print("\t"+rs.getString(i));
                    }
                    System.out.println();
                }

                String createIndexSql = "create index idx_1 on " + tableName + "(title)";
                System.out.println("Execute SQL:\n\t"+createIndexSql);
                st.execute(createIndexSql);

                String dropIndexSql = "drop index idx_1 on " + tableName;
                System.out.println("Execute SQL:\n\t"+dropIndexSql);
                st.execute(dropIndexSql);
            } catch (Exception e) {
                System.out.println(new java.util.Date().toString());
                e.printStackTrace();
            } finally {
                if (ps != null)
                    try {
                        ps.close();
                    } catch (SQLException e1) {
                        e1.printStackTrace();
                    }
                if (connection != null)
                    try {
                        connection.close();
                    } catch (SQLException e1) {
                        e1.printStackTrace();
                    }
                this.countDownLatch.countDown();
            }
        }
    }
}

```

### 7.3 其他注意事项

- 在应用分布式数据库时，首先要考虑清楚，是否真的需要分片，在所有环境下分片都是不得已的选择，必然会增加研发、运维、管理的复杂度，现在很多技术例如分区、复制、缓存等都能提高系统的处理速度和吞吐量，也许我们应用这些技术能够实现技术目标而无需分片
- 数据分片之前，请确保已经对单机数据库做了优化，包括系统架构、硬件、MySQL版本、数据读写模式等，数据分片之后这些优化依然非常有意义
- 配置拆分规则时，还应该同时考虑运维的复杂性，以及未来系统的扩容。
- 数据量较大时，DDL语句会耗时较多，分布式环境下这个问题会更加突出，因此应尽量在业务空闲时进行DDL操作。

## 8 配置示例

- [8.1 时间戳方式全局序列的配置](#)
- [8.2 MySQL-offset-step 方式全局序列的配置](#)

## 8.1 时间戳方式全局序列的配置

配置表tb的id列为时间戳方式全局序列，并按id列分片

### 1) schema.xml

```
<?xml version='1.0' encoding='utf-8'?>
<!DOCTYPE dble:schema SYSTEM "schema.dtd">
<schema name="myschema" dataNode="dn1">
    <table name="tb" dataNode="dn1,dn2" rule="rule-tb" cacheKey="id" incrementColumn="id" />
</schema>
<dataNode name="dn1" dataHost="host_1" database="db1" />
<dataNode name="dn2" dataHost="host_2" database="db1" />
<dataHost name="host_1" maxCon="1000" minCon="1000" balance="0" slaveThreshold="100">
    <heartbeat>select user()</heartbeat>
    <writeHost host="hostM1" url="172.100.10.101:3306" user="test1" password="test1"></writeHost>
</dataHost>
<dataHost name="host_2" maxCon="1000" minCon="1000" balance="0" slaveThreshold="100">
    <heartbeat>select user()</heartbeat>
    <writeHost host="hostM2" url="172.100.10.102:3306" user="test1" password="test1"></writeHost>
</dataHost>
<dble:schema xmlns:dble="http://dble.cloud/">
```

### 2) rule.xml

```
<?xml version='1.0' encoding='utf-8'?>
<!DOCTYPE dble:rule SYSTEM "rule.dtd">
<tableRule name="rule-tb">
    <rule>
        <columns>id</columns>
        <algorithm>mod</algorithm>
    </rule>
</tableRule>
<function name="mod" class="Hash">
    <property name="partitionCount">2</property>
    <property name="partitionLength">1</property>
</function>
<dble:rule xmlns:dble="http://dble.cloud/">
```

### 3) server.xml

```
<?xml version='1.0' encoding='utf-8'?>
<!DOCTYPE dble:server SYSTEM "server.dtd">
<system>
    <property name="sequenceHandlerType">2</property>
    <property name="processors">4</property>
    <property name="backendProcessors">4</property>
    <property name="processorExecutor">8</property>
    <property name="backendProcessorExecutor">4</property>
    <property name="sqlExecuteTimeout">3000000</property>
</system>
<user name="abc">
    <property name="password">abc</property>
    <property name="schemas">myschema</property>
    <property name="benchmark">1000000</property>
</user>
<user name="test">
    <property name="password">test</property>
    <property name="manager">true</property>
</user>
<dble:server xmlns:dble="http://dble.cloud/">
```

### 4) sequence\_time\_conf.properties

```
WORKID=01
DATAACENTERID=01
START_TIME=2010-10-01 09:42:54
```

配置要点:

- schema.xml:

```
<table name="tb" dataNode="dn1,dn2" rule="rule-tb" cacheKey="id" incrementColumn="id" />
```

- server.xml:

```
<property name="sequenceHandlerType">2</property>
```

- sequence\_time\_conf.properties(默认配置同上): 如果使用集群, WORKID, DATAACENTERID的配置必须使该dble实例在dble集群中唯一, 且必须为[0,31]之间的整数

## 8.2 MySQL-offset-step 方式全局序列的配置

配置表sbtest1的id列为MySQL-offset-step方式全局序列，并按id列分片

### 1) schema.xml

```
<?xml version='1.0' encoding='utf-8'?>
<!DOCTYPE dble:schema SYSTEM "schema.dtd"><db: schema xmlns:db="http://dble.cloud/">
<schema dataNode="dn5" name="mytest" sqlMaxLimit="100">
    <table dataNode="dn1,dn2,dn3,dn4" name="sbtest1" cacheKey="id" incrementColumn="id" rule="hash-four" />
</schema>

<dataNode dataHost="172.100.9.5" database="db1" name="dn1" />
<dataNode dataHost="172.100.9.6" database="db1" name="dn2" />
<dataNode dataHost="172.100.9.5" database="db2" name="dn3" />
<dataNode dataHost="172.100.9.6" database="db2" name="dn4" />
<dataNode dataHost="172.100.9.5" database="db3" name="dn5" />

<dataHost balance="0" maxCon="1000" minCon="10" name="172.100.9.5">
    <heartbeat>select user()</heartbeat>
    <writeHost host="hostM1" password="111111" url="172.100.9.5:3306" user="test">
    </writeHost>
</dataHost>

<dataHost balance="0" maxCon="1000" minCon="10" name="172.100.9.6">
    <heartbeat>select user()</heartbeat>
    <writeHost host="hostM2" password="111111" url="172.100.9.6:3306" user="test">
    </writeHost>
</dataHost>
</db: schema>
```

### 2) rule.xml

```
<?xml version='1.0' encoding='utf-8'?>
<!DOCTYPE dble:rule SYSTEM "rule.dtd"><db:rule xmlns:db="http://dble.cloud/">
<tableRule name="hash-four">
    <rule>
        <columns>id</columns>
        <algorithm>four-long</algorithm>
    </rule>
</tableRule>
<function class="Hash" name="four-long">
    <property name="partitionCount">4</property>
    <property name="partitionLength">1</property>
</function>
</db:rule>
```

### 3) server.xml

```
<?xml version='1.0' encoding='utf-8'?>
<!DOCTYPE dble:server SYSTEM "server.dtd"><db:server xmlns:db="http://dble.cloud/">
<system>
    <property name="sequenceHandlerType">1</property>
    <property name="useGlobaleTableCheck">1</property>
</system>

<user name="test">
    <property name="password">111111</property>
    <property name="schemas">mytest</property>
</user>
<user name="root">
    <property name="password">111111</property>
    <property name="manager">true</property>
</user>
</db:server>
```

### 4) sequence\_db\_conf.properties

```
#sequence stored in datanode
`mytest`.`sbtest1`=dn1
```

mytest, sbtest1, dn1均为在schema.xml配置的值

在dn1分片对应的后端dataHost/db1上执行dble安装目录下的conf/dbseq.sql(路径根据情况自行修改)。

```
mysql -h172.100.9.5 -utest -p111111 -Ddb1
mysql>source conf/dbseq.sql
```

在上述sql文件执行成功后向创建的表DBLE\_SEQUENCE插入自增相关的配置数据:

```
mysql -h172.100.9.5 -utest -p111111 -Ddb1
mysql>INSERT INTO DBLE_SEQUENCE VALUES ('`mytest`.`sbtest1`', 16, 1);
```

DBLE\_SEQUENCE列说明:

- name: 在sequence\_db\_conf.properties中配置的逻辑数据库和表名
- current\_value: 全局序列的当前值
- increment: 每次取出多少值用于全局序列, 注意全局序列递增的步长固定是1

登录dble业务端口创建设置了全局序列并以其分片的表:

```
mysql -utest -p111111 -h127.0.0.1 -P8066 -Dmytest
mysql> drop table if exists sbtest1;
Query OK, 0 rows affected (0.05 sec)
mysql> create table sbtest1(id int, k int unsigned not null default '0', primary key(id));
Query OK, 0 rows affected (0.05 sec)

mysql> insert into sbtest1 values(2);
Query OK, 1 row affected (0.11 sec)

mysql> select * from sbtest1;
+---+---+
| id | k |
+---+---+
| 17 | 2 |
+---+---+
1 row in set (0.01 sec)
```

从上面的sql可以看到, 在设置DBLE\_SEQUENCE表时, current\_value设置的是16, 在insert后变为了17。

配置要点:

- schema.xml:

```
<table dataNode="dn1,dn2,dn3,dn4" name="sbtest1" cacheKey="id" incrementColumn="id" rule="hash-four" />
```

- server.xml:

```
<property name="sequenceHandlerType">1</property>
```

- sequence\_db\_conf.properties:

```
'mytest'.'sbtest1'=dn1
```

- 在sequence\_db\_conf.properties配置的后端分片dn1对应的后端数据库上执行dbseq.sql，并插入全局序列表对应的记录

## 9 sysbench压测dble示例

- [9.1 测试环境及架构](#)
- [9.2 修改dble配置](#)
- [9.3 使用sysbench进行压测](#)

## 9.1 测试环境及架构

### 版本信息

- Sysbench version: 1.0
- Dble version: 5.6.23-dble-2.19.11.0-2d7c4911b7a4fecaa9eb0299f49c32ec11e97c42-20200228124218
- MySQL version: 5.7.25

### 测试架构

- sysbench独立服务器 (172.20.134.1)
- dble独立服务器 (172.20.134.2)
- 3台mysql独立服务器 (172.20.134.3, 172.20.134.4, 172.20.134.5)

## 9.2 修改dble配置

说明：此配置仅为示例配置，并非调优配置，请根据运行环境自行调优，调优步骤参考：[2.18 性能观测以及调试概览](#)

### 1) schema.xml

```
<?xml version="1.0"?>
<!DOCTYPE dble:schema SYSTEM "schema.dtd">
<dbles:schema xmlns:dbles="http://dble.cloud/">
<!--benchmarks-->
<schema name="sbtest">
<table name="sbtest1" primaryKey="id" dataNode="dn$1-9" rule="hash-sysbench" />
</schema>
<dataNode name="dn$1-3" dataHost="host_1" database="dbledb$1-3" />
<dataNode name="dn$4-6" dataHost="host_2" database="dbledb$4-6" />
<dataNode name="dn$7-9" dataHost="host_3" database="dbledb$7-9" />
<dataHost name="host_1" maxCon="1000" minCon="100" balance="0" slaveThreshold="-1">
<heartbeat>select user()</heartbeat>
<writeHost host="hostM1" url="172.20.134.3:3306" user="test1" password="test1"></writeHost>
</dataHost>

<dataHost name="host_2" maxCon="1000" minCon="100" balance="0" slaveThreshold="-1">
<heartbeat>select user()</heartbeat>
<writeHost host="hostM2" url="172.20.134.4:3306" user="test1" password="test1"></writeHost>
</dataHost>

<dataHost name="host_3" maxCon="1000" minCon="100" balance="0" slaveThreshold="-1">
<heartbeat>select user()</heartbeat>
<writeHost host="hostM3" url="172.20.134.5:3306" user="test1" password="test1"></writeHost>
</dataHost>
</dbles:schema>
```

### 2) rule.xml

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE dble:rule SYSTEM "rule.dtd">
<dbles:rule xmlns:dbles="http://dble.cloud/">
<tableRule name="hash-sysbench">
<rule>
<columns>id</columns>
<algorithm>hash-sysbench</algorithm>
</rule>
</tableRule>
<function name="hash-sysbench" class="Hash">
<property name="partitionCount">9</property>
<property name="partitionLength">1</property>
</function>
</dbles:rule>
```

### 3) server.xml

```
<?xml version='1.0' encoding='utf-8'?>
<!DOCTYPE dble:server SYSTEM "server.dtd">
<dbles:server xmlns:dbles="http://dble.cloud/">
<system>
<property name="processors">10</property>
<property name="backendProcessors">10</property>
<property name="processorExecutor">8</property>
<property name="backendProcessorExecutor">6</property>
<property name="sqlExecuteTimeout">3000000</property>
</system>
<user name="test">
<property name="password">111111</property>
<property name="schemas">sbtest</property>
<property name="benchmark">1000000</property>
</user>
<user name="root">
<property name="password">111111</property>
<property name="manager">true</property>
</user>
</dbles:server>
```

在后端mysql节点创建相应物理库

- 172.20.134.3:3306 创建库: dbledb1, dbledb2, dbledb3
- 172.20.134.4:3306 创建库: dbledb4, dbledb5, dbledb6
- 172.20.134.5:33066 创建库: dbledb7, dbledb8, dbledb9

## 9.3 使用sysbench进行压测

说明:

如果在压测中遇到问题请参考: [how about a sysbench-testing-quick-start](#); 如果在此参考中未包括您遇到的问题请在[issue](#)下添加评论, 将问题及配置做出说明。

注意:

使用sysbench压测语句, 在高并发下会出现主键冲突的报错。

```
FATAL: mysqlDrv_query() returned error 1062 (Duplicate entry '49823' for key 'PRIMARY') for query 'INSERT INTO sbtest1 (id, k, c, pad) VALUES (49823, 58210, '27111667985-11552069038-79242882109-05602914209-02374993639-3224261'
```

解决方式见如上链接。

数据清理:

```
/usr/share/sysbench/oltp_read_write.lua --mysql-db=sbtest --mysql-host=172.20.134.2 --mysql-port=8066 --mysql-user=test1 --mysql-password=test1 --auto_inc=off --tables=1 --table-size=100000 --threads=4 --time=30 --report-interval=1 --max-requests=0 --percentile=95 --db-ps-mode=disable --skip-trx=on cleanup
```

数据准备:

```
/usr/share/sysbench/oltp_read_write.lua --mysql-db=sbtest --mysql-host=172.20.134.2 --mysql-port=8066 --mysql-user=test1 --mysql-password=test1 --auto_inc=off --tables=1 --table-size=100000 --threads=4 --time=30 --report-interval=1 --max-requests=0 --percentile=95 --db-ps-mode=disable --skip-trx=on prepare
```

执行压测:

```
/usr/share/sysbench/oltp_read_write.lua --mysql-db=sbtest --mysql-host=172.20.134.2 --mysql-port=8066 --mysql-user=test1 --mysql-password=test1 --auto_inc=off --tables=1 --table-size=100000 --threads=4 --time=30 --report-interval=1 --max-requests=0 --percentile=95 --db-ps-mode=disable --skip-trx=on run
```

- [A.1 ErrorCode](#)
- [A.2 原理解释](#)
- [A.3 使用说明](#)

- Max Connections
- Out Of Memory Error
- The Problem Of Hint
- NestLoop Parameters Lead To Temptable Exception
- Can't Get Variables From DataNode
- Port Already In Use 1984
- Sharding Column Cannot Be Null

## db1e-MaxConnections

### Issue

- [Err] 3009 - java.io.IOException: the max activeConnnections size can not be max than maxconnections.

### Resolution

1. 检查是否使用了过多的短链，导致频繁建立连接，链接池不够用

**注意：**不要过多的使用短链，容易消耗db1e资源

2. maxCon配置过小，调大maxCon

配置名称	配置内容	默认值/单位	作用原理或应用
maxCon	控制最大连接数	默认1024	大于此连接数之后,建立连接会失败.注意当各个用户的maxcon总和值大于此值时, 以当前值为准 全局maxCon不作用于manager用户

### Root Cause

配置的链接池数量不够用，导致报错。

### Relevant Content

#### 数据库连接池

##### 1. maxWait

从连接池获取连接的超时等待时间，单位毫秒。

**注意：**maxActive不要配置过大，虽然业务量飙升后还能处理更多的请求，但其实连接数的增多在很多场景下反而会减低吞吐量

##### 2. connectionProperties

可以配置 connectTimeout 和 socketTimeout，单位都是毫秒。connectTimeout 配置建立 TCP 连接的超时时间；socketTimeout 配置发送请求后等待响应的超时时间。

**注意：**不设置这两项超时时间，服务会有高的风险。

例如网络异常下 socket 没办法检测到网络错误，如果没有设置 socket 网络超时，连接就会一直等待 DB 返回结果，造成新的请求都无法获取到连接。

##### 3. maxActive

最大连接池数量，允许的最大同时使用中的连接数。

**注意：**当业务出现大流量涌入时，连接池耗尽，maxWait未配置或者配置为 0 时将无限等待，导致等待队列越来越长，表现为业务接口大量超时，实际吞吐越低。

## dble-OutOfMemoryError

### Setting

load data语句，一共有17G的数据，每条语句4K  
load data相关参数值均为默认

### Issue

- INFO | jvm 1 | 2019/06/28 14:55:37 | Exception in thread "backendBusinessExecutor17" java.lang.OutOfMemoryError: GC overhead limit exceeded.

### Resolution

- 调小maxRowSizeToFile参数值，减少占用内存量

配置名称	配置内容	默认值	详细作用原理或应用
maxCharsPerColumn	每列所允许最大字符数	默认为65535	每行所允许最大字符数
maxRowSizeToFile	需要持久化的最大行数	默认为10000	当load data的数据行数超过阈值后，会将数据保存在文件中以防OOM

- 调大wrapper.conf中的Xmx值

注意：发生OOM，可以关注一下这个值的大小是否配置合适

### Root Cause

- maxRowSizeToFile参数设置过大，对于机器来说超出了承载的数据量，导致内存溢出
- 配置文件中Xmx设置过小，内存不够

注意：On-Heap 大小由JVM 参数Xms ,Xmx 决定，就是正常服务需要的内存，由jvm自动分配和回收。

### Relevant Content

#### load data相关配置

设置load data相关参数时，不要过度的调大，防止OOM

#### JVM配置

以下为建议值：

- dble总内存=0.6 可用物理内存(刨除操作系统,驱动等的占用)
- Xmx = 0.4 dble总内存
- MaxDirectMemorySize = 0.6 \* dble总内存

#### 堆内存分配

- JVM初始分配的内存由-Xms指定，默认是物理内存的1/64;
- JVM最大分配的内存由-Xmx指定，默认是物理内存的1/4。
- 默认空余堆内存小于40%时，JVM就会增大堆直到-Xmx的最大限制；空余堆内存大于70%时，JVM会减少堆直到-Xms的最小限制。

## dble-TheProblemOfHint

### Setting

- mysql> /\*dble:sql=select 1 from rp\_cre\_data\_mobile\_track\_cmcc / call update\_track();
- 预期：表为分片表，sql语句应该下发到所有节点
- 结果：只有一个节点在执行

### Issue

- mysql> show @@processlist;

Front_Id	Datanode	BconnID	user	Front_Host	db	Command	Time	State	Info
33	NULL	NULL	root	略	NULL	NULL	0	updating	NULL
34	NULL	NULL	root	略	NULL	NULL	0	updating	NULL
35	NULL	NULL	root	略	NULL	NULL	0	updating	NULL
41	dn9	9372	root	略	db9	Query	0	updating	NULL
42	NULL	NULL	root	略	NULL	NULL	0	updating	NULL
43	NULL	NULL	root	略	NULL	NULL	0	updating	NULL
30	NULL	NULL	admin	略	NULL	NULL	0	updating	NULL

### Resolution

- 将注解方式: /\*dble:type=....\*/
 改为: /\*dble:type=....\*/
  - 或者在mysql client端，加上 -c 选项
- 注意: mysql --help  
 -c, --comments  
 Preserve comments. Send comments to the server. The default is --skip-comments (discard comments), enable with --comments.

### Root Cause

- /\*dble:type=....\*/ 这种注释方式是mysql的标准注释
- 如果不加 -c 选项， 默认注释会被skip

### Relevant Content

#### hint作用

- 指定路由，比如强制读写分离。
- 帮助dble支持一些不能实现的语句，如单节点内存储过程的创建和调用，如insert...select...;

#### Hint语法

Hint语法有两种形式:

- /\*dble:type=....\*/
- /\*dble:type=...\*/

type有3种值可选: datanode, db\_type, sql。

type详情请见: [https://actiontech.github.io/dble-docs-cn/2.Function/2.04\\_hint.html](https://actiontech.github.io/dble-docs-cn/2.Function/2.04_hint.html)

#### Hint注意事项

- 使用select语句作为注解SQL，不要使用delete/update/insert 等语句。delete/update/insert 等语句虽然也能用在注解中，但这些语句在SQL处理中有一些额外的逻辑判断，会降低性能，不建议使用
- 注解SQL 禁用表关联语句
- 使用hint做DDL需要额外执行reload @@metadata
- 使用hint做session级别的系统变量和环境变量可能不会生效，请慎用
- dble的注解和MySQL原生注解含义不同，想通过MySQL原生注解来设置变量或者指定索引是无法得到预期结果的。如#1169
- 使用注解并不额外增加的执行时间；从解析复杂度以及性能考虑，注解SQL应尽量用最简单的SQL 语句，如select id from tab\_a where id='10000'；
- 能不用注解也能够解决的场景，尽量不用注解

## dble-NestLoop Parameters Lead To Temptable Exception

### Setting

- 启用了NestLoop优化，设置NestLoop值的范围为4
- 两个表join关联：student表和class表
- student表结构：id列、name列、class\_name列，主键id
- class表结构：id列、class\_name列、teacher\_name列
- SELECT class.teacher\_name FROM student LEFT JOIN class on student.class\_name=class.class\_name WHERE student.name='张三';

### Issue

- com.actiontech.dble.plan.common.exception.TempTableException: temptable too much rows,[rows size is 5].

### Resolution

- 调大NestLoop中的默认参数值，如下：

配置名称	配置内容	默认值/单位	详细作用原理或应用
useJoinStrategy	是否使用nestLoop优化	默认使用	开启之后会尝试判断join两边的where来重新调整查询SQL下发的顺序
nestLoopConnSize	临时表阈值	默认4	若临时表行数大于这两个值乘积，则报告错误
nestLoopRowsSize	临时表阈值	默认2000	若临时表行数大于这两个值乘积，则报告错误

- 或者关闭NestLoop，不使用其优化，如下：
  - <property name="useJoinStrategy">false</property>
- 为了更直观，本例提前调小参数值，引发场景复现

### Root Cause

- 使用了NestLoop优化，但是依据NestLoop规则选择出的表数据量太大，超出了NestLoop默认范围
- NestLoop优化规范

注意：

- 对于两表join时，NestLoop选择小表作为驱动表（同样适用于多join）
- 判断依据：是否有where条件，有where条件的作为小表
- 当两个表都有where条件或者都没有where条件，NestLoop无法判断，不起作用
- 在不确定哪个表为小表时，不建议开启NestLoop

- 针对本例SQL做详细说明：如下

- SELECT class.teacher\_name FROM student LEFT JOIN class on student.class\_name=class.class\_name WHERE student.name='张三';

根据dble中NestLoop优化规则可知：

- SQL中使用了student表和class表
- where条件限制指定于student表，所以NestLoop可以根据规则选择出student表作为驱动表（小表）
- 实际SQL中student表的数据量较大，超出了NestLoop值的范围，引发报错

### Relevant Content

#### MySQL的多表连接之NestLoop介绍

##### 1. NestLoop:

- 对于被连接的数据子集较小的情况，Nested Loop是个较好的选择。
- Nested Loop就是扫描一个表（外表），每读到一条记录，就根据Join字段上的索引去另一张表（内表）里面查找，若Join字段上没有索引查询优化器一般就不会选择 Nested Loop。
- 在Nested Loop中，内表（一般是带索引的大表）被外表（也叫“驱动表”，一般为小表——不紧相对其它表为小表，而且记录数的绝对值也较小，不要求有索引）驱动，外表返回的每一行都要在内表中检索找到与它匹配的行，因此整个查询返回的结果集不能太大。

##### 2. NestLoop优缺点

类别	使用条件	相关资源	特点	缺点
NestLoop	任何条件	CPU、磁盘I/O	当有高选择性索引或进行限制性搜索时效率比较高，能够快速返回第一次的搜索结果。	当索引丢失或者查询条件限制不够时，效率很低；当表记录数多时，效率低。

## dble-Can't get variables from data node

### Setting

- servre.xml

```
<property name="password">123456</property>
<user name="man1">
<property name="password">123456</property>
<property name="manager">true</property></user>
<user name="root">
<property name="password">123456</property>
<property name="schemas">TESTDB</property>
```

- schema.xml

```
<schema name="testdble">
<table name="tb1_user" primary Key="ID" dataNode="dn1,dn2" rule="rule_mod2" /> </schema>
<dataNode name="dn1" dataHost="localhost1" database="t_dble1"/>
<dataNode name="dn2" dataHost="localhost2" database="t_dble2"/>
<dataHost name="localhost1" maxCon="1000" minCon="10" balance="0" slaveThreshold="100">
<heartbeat>show slave status</heartbeat>
<writeHost host="hostM1" url="localhost:3306" user="root" password="nE7jA%5m"> </writeHost>
</dataHost>
<dataHost name="localhost2" maxCon="1000" minCon="10" balance="0" slaveThreshold="100">
<heartbeat>show slave status</heartbeat>
<writeHost host="hostM2" url="localhost:3306" user="root" password="nE7jA%5m"> </writeHost> </dataHost>
```

### Issue

- 查看dble启动日志:

```
Running dble-server...
wrapper | --> Wrapper Started as Console
wrapper | Launching a JVM...
jvm 1 | Wrapper (Version 3.2.3)
http://wrapper.tanukisoftware.org
jvm 1 | Copyright 1999-2006 Tanuki Software, Inc. All Rights Reserved.
jvm 1 |
jvm 1 | java.io.IOException:Can't get variables from data node ...
wrapper | <-- Wrapper Stopped
```

### Resolution

- 检查mysql是否正常启动, 如果启动正常见下一步;
- schema中的root用户能否通过配置文件的信息连接mysql, 连接成功见下一步;
- 检查root用户权限;
- 连接mysql, 并执行show variables命令, 未执行成功见下一步;
- 修改配置文件中datahost指定的后端数据库密码, 更新配置文件, dble正常启动

### Root Cause

- 通过schema配置信息成功连接mysql后, 并不能执行show variables
- 由于mysql 5.7 初始化之后, 首次使用随机密码登陆, 没有修改密码, 无法对数据库进行操作

### Relevant Content

- 安装好mysql5.7后, 第一次初始化数据库
- 随机密码登录mysql, 首次登录后, mysql要求必须修改默认密码, 否则不能执行任何其他数据库操作, 这样体现了不断增强的MySQL安全性。
- 第一次登陆后必须更改密码:
  - mysql> show databases;
  - ERROR 1820 (HY000): You must reset your password using ALTER USER statement before executing this statement.
  - mysql > set password = password('xxxxx');

## dble-Port already in use:1984

### Issue

- wrapper.log-Error1

```
STATUS | wrapper | 2019/07/23 16:37:06 | --> Wrapper Started as Daemon
STATUS | wrapper | 2019/07/23 16:37:06 | Launching a JVM...
INFO | jvm 1 | 2019/07/23 16:37:06 | OpenJDK 64-Bit Server VM warning: ignoring option MaxPermSize=64M; support was removed in 8.0
INFO | jvm 1 | 2019/07/23 16:37:07 | 错误: 代理抛出异常错误: java.rmi.server.ExportException: Port already in use: 1984; nested exception is:
INFO | jvm 1 | 2019/07/23 16:37:07 | java.net.BindException: Address already in use (Bind failed)
INFO | jvm 1 | 2019/07/23 16:37:07 | sun.management.AgentConfigurationError: java.rmi.server.ExportException: Port already in use: 1984;
```

- wrapper.log-Error2

```
STATUS | wrapper | 2019/07/26 16:12:48 | --> Wrapper Started as Daemon
STATUS | wrapper | 2019/07/26 16:12:49 | Launching a JVM...
INFO | jvm 1 | 2019/07/26 16:12:49 | Wrapper (Version 3.2.3) http://wrapper.tanukisoftware.org
INFO | jvm 1 | 2019/07/26 16:12:49 | Copyright 1999-2006 Tanuki Software, Inc. All Rights Reserved.
INFO | jvm 1 | 2019/07/26 16:12:49 |
INFO | jvm 1 | 2019/07/26 16:12:51 | java.net.BindException: Address already in use
INFO | jvm 1 | 2019/07/26 16:12:51 | at sun.nio.ch.Net.bind0(Native Method)
INFO | jvm 1 | 2019/07/26 16:12:51 | at sun.nio.ch.Net.bind(Net.java:433)
INFO | jvm 1 | 2019/07/26 16:12:51 | at sun.nio.ch.Net.bind(Net.java:425)
INFO | jvm 1 | 2019/07/26 16:12:51 | at sun.nio.ch.ServerSocketChannelImpl.bind(ServerSocketChannelImpl.java:223)
INFO | jvm 1 | 2019/07/26 16:12:51 | at com.actiontech.dble.net.NIOAcceptor.<init>(NIOAcceptor.java:46)
```

### Resolution

- 根据Error1:  
修改配置文件wrapper.conf: 修改被占用端口1984  
-Dcom.sun.management.jmxremote.port=1984
- 根据Error2:  
netstat -nap 查看程序运行的pid (8066和9066依然存在)  
kill -9 pid 杀掉进程  
启动dble成功

### Root Cause

- 已经启动过一个开启jmx服务的java程序后，再启动dble会报这个错；
- dble启动过程中会占用三个端口：业务端口，管理端口，jvm对外提供jmx服务端口；
- jmx可通过jconsole连接上jvm，观测jvm的运行状态。

### Relevant Content

#### 1. JVM

JVM是一种使用软件模拟出来的计算机，它用于执行Java程序，有一套非常严格的技术规范，是Java跨平台的依赖基础。  
Java虚拟机有自己想象中的硬件，如处理器，堆栈，寄存器等，还有相应的指令系统它允许Java程序就好像一台计算机允许c或c++程序一样。

#### 2. JMX

所谓JMX，是Java Management Extensions的缩写，是Java管理系统的一个标准、一个规范，也是一个接口，一个框架。  
它和JPA、JMS是一样的，就是通过将监控和管理涉及到的各个方面的问题和解决办法放到一起，统一设计，以便向外提供服务，以供使用者调用。

#### 3. Jconsole

Jconsole是JDK自带的监控工具，它用于连接正在运行的本地或者远程的JVM，对运行在java应用程序的资源消耗和性能进行监控，并画出大量的图表，提供强大的可视化界面。  
自身占用的服务器内存很小，甚至可以说几乎不消耗。

## dble-Sharding column can't be null

### Setting

- rule.xml部分配置如下:

```
<tablerule name="sharding-by-range">
<rule>
<columns>number</columns>
<algorithm>rangeLong</algorithm>
</rule>
</tablerule>
<function name="rangeLong" class="NumberRange">
<property name="mapFile">partition.txt</property>
<property name="defaultNode">0</property>
</function>
```

- create table account (id int(10),number int(10) not null,name varchar(20) not null);
- insert into account (id,number,name) values (1,NULL,'aaa');

### Issue

ERROR 1064 (HY000): Sharding column can't be null when the table in MySQL column is not null

### Resolution

- number列和name列的插入值不为NULL;
- 或者修改number列为允许插入NULL值;  
ALTER TABLE account MODIFY number VARCHAR (20);
- 注意:** 上一步的前提是:  
在server.xml中开启参数alterTableAllow;

```
</whitehost>
<blacklist check="true">
<property name="alterTableAllow">true</property>
</blacklist>
```

并修改sharding-by-range中的拆分列, dble不允许对分片键或ER键进行alter, 会造成无法分片;

```
<tablerule name="sharding-by-range">
<rule>
<columns>id</columns>
<algorithm>rangeLong</algorithm>
</rule>
```

- alter列值为允许插入空值后, 再将拆分列修改为原值。

### Root Cause

- 在MySQL中执行相同insert:  
报错: ERROR 1048 (23000): Column 'number' cannot be null
- desc查看表结构: number列和name列均定义为非空列, 不允许插入空值。

Field	Type	Null	Key	Default	Extra
id	int(10)	YES		NULL	
number	int(10)	NO		NULL	
name	varchar(20)	NO		NULL	

- [How To Use Explain To Resolve The Distribution Rules Of Group G](#)
- [Hash And ConsistentHashing And Jumpstringhash](#)

## dble-How To Use Explain To Resolve The Distribution Rules Of Group By

### Questions

一张数据表做了分表。如果查询里有group by分组统计，运行原理是按范围去各分表查询出数据后，再到中间件里进行分组统计的吗？

### Conclusions

- 会在中间件中做数据重聚合
  - 利用explain工具查看sql的执行过程
  - dble 在explain上做了大量改善，相比mycat能提供更详实的执行计划，更准确的反映SQL语句的执行过程

### For Example

- 配置好配置文件：

```
schema.xml: <table name="eee" dataNode="dn1,dn2" cacheKey="id" rule="sharding-by-hash"/>
rule.xml: <rule>
<columns>id</columns>
<algorithm>hashLong</algorithm>
</rule>

</tableRule>
<function name="hashLong" class="Hash"> <property name="partitionCount">2</property> <property name="partitionLength">128</property> </function>
```

1. 在dble client创建表 eee 并插入数据： mysql> select \* from eee;

| id | name || -- | -- | 1 | 上海 | 2 | 广州 | 3 | 杭州 | 4 | 北京 | 5 | 北京 | 130 | 北京 | 131 | 北京 | 132 | 上海 | 133 | 上海 | 134 | 上海 |

mysql> select name,count(name) from eee group by name;

name	COUNT(name)
上海	4
北京	4
广州	1
杭州	1

- 利用explain工具查看sql的执行过程

DATA_NOD	TYPE	SQL/REF
dn1_0	BASE SQL	select eee . name ,COUNT(name) as _\$COUNT\$_rpda_0 from eee GROUP BY eee . name ASC
dn2_0	BASE SQL	select eee . name ,COUNT(name) as _\$COUNT\$_rpda_0 from eee GROUP BY eee . name ASC
merge_1	MERGE	dn1_0; dn2_0
aggregate_1	AGGREGATE	merge_1
shuffle_field_1	SHUFFLE_FIELD	aggregate_1

### Instructions

由explain的结果可知：

- dble将sql语句下发到对应datanode执行
- 将对应datanode数据结果进行merge
- 对merge后的数据进行group by聚合
- SHUFFLE\_FIELD进行整理，达到用户预期的结果

注意：普通用户可以不关注SHUFFLE\_FIELD

### Relevant Content

#### dble的内部功能层

- 在dble内部，包括了三个部分：面向app的连接层，内部功能层，面向mysq的连接池。
- 内部功能层实现：前端请求接收，处理过后由后端协议层发出，将数据返回给用户。

内部功能涉及到了简单查询和复杂查询：

- 简单查询：直接下发表单个/多个节点
- 复杂查询：dble内部需要进行排序、聚合、join、group by，结果集计算
- 详细介绍：<https://opensource.actionsky.com/dble-lesson-one/>

#### dble内部各线程池简介

- 后端IO接收线程 处理来自MySQL连接的网络包（sql的执行结果、查询的结果集）。
- 后端业务处理线程 简单查询的后端MySQL返回处理，并将结果转发反馈给客户端。
- 复杂查询处理线程 处理复杂查询MySQL返回结果，包括结果集的聚合，对比，排序，去重，子查询语句下发等。

PS: dble中仅此线程池不限线程数量

## dble-Hash And ConsistentHashing And Jumpstringhash

### Questions

- dble中的hash和一致性hash是否相同
- 为什么没有一致性hash

### Conclusions

- dble中的hash并非一致性hash
- 从一致性hash的性能和均衡性来看，已被跳增一致性hash取代

### Instructions

#### dble-hash

- 配置如下:
- rule.xml

```
<function name="hashLong" class="hash">
<property name="partitionCount">1,2</property>
<property name="partitionLength">10,20</property>
</function>
```

- 或者

```
<function name="hashLong" class="hash">
<property name="partitionCount">4</property>
<property name="partitionLength">10</property>
</function>
```

注意: partitionCount: 指定分区的区间数

partitionLength: 指定各区间的长度

模值 (M) 为最后一个区间段的末尾值 (C1L1 + ... + CnLn)

对于dble-hash来说，该算法可以均匀地将数据分布到各个节点。

例如:

Count=2,Length=2 -> [0,2][2,4] -> 模值=4

key分别取值: 1,2,3,4,5,6,7,8

数据分布情况如下:

node1	node2
1,4,5,8	2,3,6,7

当增加一个节点，使Count=3,Length=2 -> [0,2][2,4][4,6] -> 模值=6 key分别取值: 1,2,3,4,5,6,7,8

数据分布情况如下:

node1	node2	node3
1,6,7	2,3,8	4,5

从这个例子中我们可以看出：当节点数增加时，大多数旧的数据都需要重新分布，而重新分布的成本就是需要在count数发生变化的时候，进行数据迁移，大多数的数据都需要重新移动。

#### summary

因此，对于这种算法，当node数发生变化（增加、移除）后，大多数旧的数据都需要重新分布，并进行数据迁移。

#### Consistent Hashing

一致性哈希，将整个哈希值空间组织成一个虚拟的圆环，整个空间按顺时针方向组织。例如我们有NodeA、Node B、Node C、Node D四个节点，有数据A、数据B、数据C、数据D四个数据对象，根据一致性哈希算法，数据A会被分布到Node A上，B被分布到Node B上，C被分布到Node C上，D被分布到Node D上：

```
数据A —> NodeA
数据B —> NodeB
数据C —> NodeC
数据D —> NodeD
```

现假设Node C不幸宕机，此时数据A、B、D不会受到影响，只有数据C被重分布到Node D。如果在系统中增加一台服务器Node X：

```
数据A —> NodeA
数据B —> NodeB
新增 —> NodeX
数据C —> NodeC
数据D —> NodeD
```

数据A、B、D不受影响，只有数据C需要重分布到新的Node X。也就是说旧数据之间不会发生数据的变动。

虽然一致性Hash算法解决了节点变化导致的数据迁移问题，但是数据项分布的均匀性不够好。一致性哈希算法分布不均匀的原因是因为：将node进行哈希后，这些值并没有均匀地落在环上，因此，这些节点所管辖的范围（每个节点实际占据环上的区间大小不）并不均匀，最终导致了数据分布的不均匀。

因此，为使得每个节点在环上所“管辖”更加均匀，一致性哈希算法引入了虚拟节点机制，即对每一个服务节点计算多个哈希，每个计算结果位置都放置一个服务节点，称为虚拟节点。

例如：

我们有NodeA、Node B；可以为每台服务器计算三个虚拟节点，于是可以分别计算 "Node A#1"、"Node A#2"、"Node A#3"、"Node B#1"、"Node B#2"、"Node B#3"的哈希值，形成六个虚拟节点。

同时数据定位算法不变，只是多了一步虚拟节点到实际节点的映射，例如定位到"Node A#1"、"Node A#2"、"Node A#3"三个虚拟节点的数据均定位到Node A上，"Node B#1"、"Node B#2"、"Node B#3"三个虚拟节点的数据均定位到Node B上，用来解决分布不均的问题。

```
Node A#1 —> NodeA
Node A#2 —> NodeA
Node A#3 —> NodeA
Node B#1 —> NodeB
Node B#2 —> NodeB
Node B#3 —> NodeB
```

因此，通过增加虚节点的方法，使得每个节点在环上所“管辖”更加均匀。这样就既保证了在节点变化时，尽可能小的影响数据分布的变化，而同时又保证了数据分布的均匀。也就是靠增加“节点数量”加强管辖区间的均匀。

#### summary

对于一致性hash算法来说，当node数发生变化（增加、移除）后，旧的数据之间没有发生变动，只是需要将部分旧数据重新分布到新的节点。

#### Jumpstringhash

- 配置如下:
- rule.xml

```
<function name="jumphash"
class="jumpStringHash">
<property name="partitionCount">2</property>
<property name="hashSlice">0:2</property>
</function>
```

注意: partitionCount: 分片数量 hashSlice: 分片截取长度

该算法该算法来自于Google的一篇文章A Fast, Minimal Memory, Consistent Hash Algorithm，核心思想是通过概率分布的方法将一个hash值在每个节点分布的概率变成 $1/n$ ，并且可以通过更简便的方法可以计算得出，并且分布也更加均匀。

设计目标是把对象均匀地分布在所有节点中（平衡性）；当节点数量变化时，只需要把一些对象从旧节点移动到新节点，不需要做其它移动（单调性）。

根据论文原理，可以这样说明：

- 假设我们有四个节点，比作为四个桶，分别是：0,1,2,3
- 数据落入第一个桶(0)的概率是1；落入第二个桶(1)的概率是 $1/(n+1)$ ，也就是 $1/2$ ；落入第三个桶(2)的概率是 $1/(n+1)$ ，也就是 $1/3$ ；落入第四个桶(3)的概率是 $1/(n+1)$ ，也就是 $1/4$ 。
- 由此，一般规律是：数据落入每个桶中的概率，有占比  $n/(n+1)$  的结果保持不变，而有  $1/(n+1)$  跳变为  $n+1$ 。
- 而每个数据都是落入至index (max) 中，比如数据a有概率分布至0,1,2,3四个桶中，那么最终会落入至3内。

当增加一个桶时：由0,1,2,3变为0,1,2,3,4，需要重新分布的数据仅是每个桶内有概率分布到新桶内的数据，旧数据之间不会发生变化。

## summary

由于Jumpstringhash采用的是概率分布的结果，因此计算相对于一致性hash较简单。

- Jumpstringhash相比于一致性hash算法来说，占用内存更小，计算更快，数据分布更均匀；
- Jumpstringhash和一致性hash算法，对于节点变动的情况下，都是将部分旧数据重新分布到新节点上，旧数据之间不会发生变动；
- Jumpstringhash和一致性hash算法相比于dble-hash来说，节点变动的情况下，旧数据之间不会发生变动；
- dble-hash相比于Jumpstringhash和一致性hash算法来说，最大的优势在于计算方便，可以人为的快速计算出结果。

- [ToBeContinued2](#)

待更新