

- 一、实验目的与环境
- 二、基础环境介绍
- 三、搭建主从集群
 - 1、理论基础
 - 2、同步的原理
 - 3、搭建主从集群
 - 3.1 配置master主服务器
 - 3.2 配置slave从服务
 - 3.3 主从集群测试
 - 3.4 集群搭建扩展：
 - 3.5、GTID同步集群
 - 4、集群扩容
 - 5、半同步复制
 - 1、理解半同步复制
 - 2、搭建半同步复制集群
 - 6、主从架构的数据延迟问题：
- 四、MySQL的高可用方案--了解
 - 1、MMM
 - 2、MHA
 - 3、MGR
- 五、分库分表
 - 1、分库分表有什么用
 - 2、分库分表的方式
 - 3、分库分表的缺点
 - 4、什么时候需要分库分表？
 - 5、常见的分库分表组件

一、实验目的与环境

实验目的：

MySQL是现在互联网最常用的开源数据库产品。但是我们平常开发使用，大都是用的单机服务。而在实际生产中，往往数据量会极为庞大，并且数据的安全性要求也更高，这样单机的MySQL，不管是性能还是安全都是打不到要求的。所以在生产环境中，MySQL必须是要搭建一套主从复制的架构，同时可以基于一些工具实现高

可用架构。然后，在此基础上，就可以基于一些中间件实现读写分离架构。最后如果数据量非常大，还必须可以实现分库分表的架构。

当然，MySQL的这些架构搭建还是比较复杂的，通常会由专门的运维人员来搭建。所以这次实验的目的，并不是要大家就学会去搭建MySQL集群，而是带大家在生产环境下的MySQL架构有一定的理解，能够在自己的生产项目中运用上MySQL的生产架构。同时，也是为后面学习ShardingSphere分库分表打下基础。

实验环境：

1、Linux服务器两台：centos7

2、mysql版本：mysql-8.0.20

说明：

关于MySQL的安装：Linux服务的基础安装，在之前学习MySQL优化课程时已经讲解过，这里就不再多做介绍了。MySQL服务器安装方式，大家可以参考官网的手册，自行进行。

MySQL安装手册地址：<https://dev.mysql.com/doc/refman/8.0/en/binary-installation.html>

如果有同学想要在Linux上进行实验而MySQL安装实在有问题的话，可以推荐大家使用宝塔面板，可以省掉非常多的MySQL安装的问题。

另外，对于熟悉Docker的同学，可以直接用Docker来搭建，更为简单高效，也是未来互联网企业进行服务搭建的主流方式。大家以后学Docker的时候也可以留意下。

二、基础环境介绍

两台服务器，均安装CentOS7。

1、192.168.232.128 作为mysql主节点部署

2、192.168.232.129 作为mysql从节点部署

mysql版本：mysql-8.0.20

为了便于使用，两个mysql服务需要打开远程登录权限，开启方式需要在本机登录mysql，执行以下语句。

```
1  #开启远程登录
2  use mysql;
3  update user set host='%' where user='root';
4  flush privileges;
```

三、搭建主从集群

1、理论基础

主从架构有什么用？通过搭建MySQL主从集群，可以缓解MySQL的数据存储以及访问的压力。

1、数据安全

给主服务增加一个数据备份。基于这个目的，可以搭建主从架构，或者也可以基于主从架构搭建互主的架构。

2、读写分离

对于大部分的JAVA业务系统来说，都是读多写少的，读请求远远高于写请求。这时，当主服务的访问压力过大时，可以将数据读请求转为由从服务来分担，主服务只负责数据写入的请求，这样大大缓解数据库的访问压力。

要理解，MySQL的主从架构只是实现读写分离的一个基础。实现读写分离还是需要一些中间件来支持，比如ShardingSphere。

3、故障转移-高可用

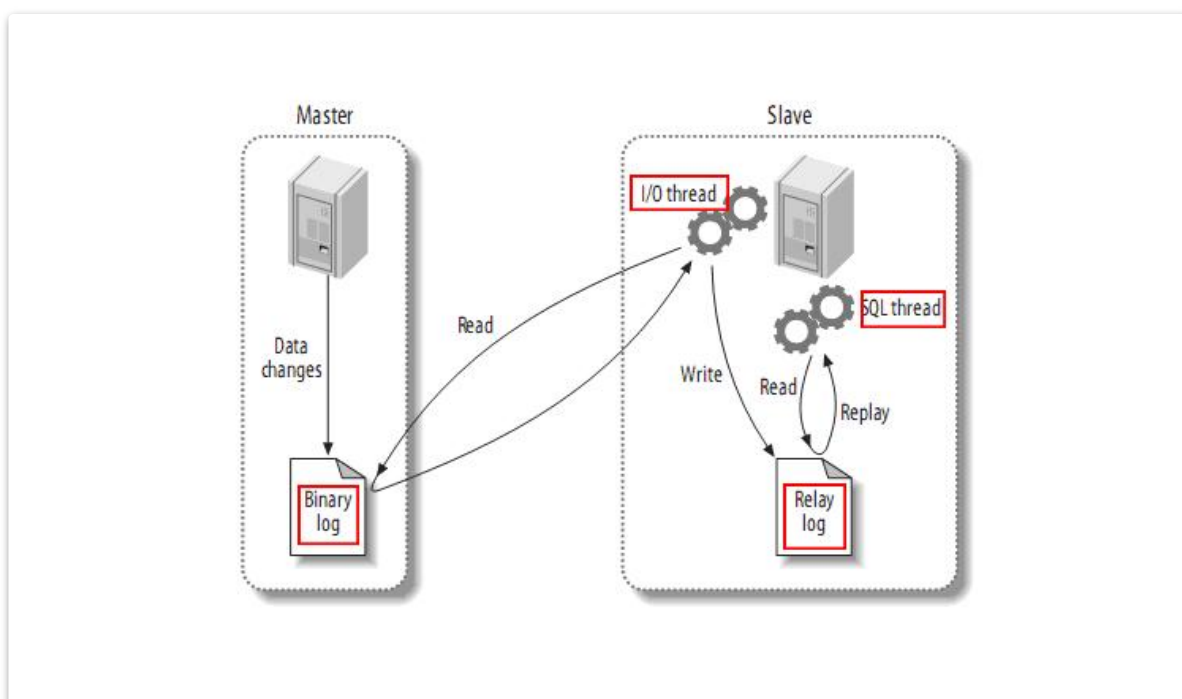
当MySQL主服务宕机后，可以由一台从服务切换成为主服务，继续提供数据读写功能。

对于高可用架构，主从数据的同步也只是实现故障转移的一个前提条件，要实现MySQL主从切换，还需要依靠一些其他的中间件来实现。比如MMM、MHA、MGR。

在一般项目中，如果数据库的访问压力没有那么大，那读写分离不一定是必须要做的，但是，主从架构和高可用架构则是必须要搭建的。

2、同步的原理

MySQL服务的主从架构一般都是通过binlog日志文件来进行的。即在主服务上打开binlog记录每一步的数据库操作，然后从服务上会有一个IO线程，负责跟主服务建立一个TCP连接，请求主服务将binlog传输过来。这时，主库上会有一个IO dump线程，负责通过这个TCP连接把Binlog日志传输给从库的IO线程。接着从服务的IO线程会把读取到的binlog日志数据写入自己的relay日志文件中。然后从服务上另外一个SQL线程会读取relay日志里的内容，进行操作重演，达到还原数据的目的。我们通常对MySQL做的读写分离配置就必须基于主从架构来搭建。



MySQL的binlog不光可以用于主从同步，还可以用于缓存数据同步等场景。

例如Canal，可以模拟一个slave节点，向MySQL发起binlog同步，然后将数据落地到Redis、Kafka等其他组件，实现数据实时流转。

搭建主从集群时，有两个必要的要求：

- 双方MySQL必须版本一致。至少需要主服务的版本低于从服务
- 两节点间的时间需要同步。

3、搭建主从集群

3.1 配置master主服务器

首先，配置主节点的mysql配置文件： /etc/my.cnf 这一步需要对master进行配置，主要是需要打开binlog日志，以及指定serverId。我们打开MySQL主服务的my.cnf文件，在文件中一行server-id以及一个关闭域名解析的配置。然后重启服务。

```
1  [mysqld]
2  server-id=47
3  #开启binlog
4  log_bin=master-bin
5  log_bin-index=master-bin.index
6  skip-name-resolve
7  # 设置连接端口
8  port=3306
9  # 设置mysql的安装目录
10 basedir=/usr/local/mysql
11 # 设置mysql数据库的数据的存放目录
12 datadir=/usr/local/mysql/mysql-files
13 # 允许最大连接数
14 max_connections=200
15 # 允许连接失败的次数。
16 max_connect_errors=10
17 # 服务端使用的字符集默认为UTF8
18 character-set-server=utf8
19 # 创建新表时将使用的默认存储引擎
20 default-storage-engine=INNODB
21 # 默认使用“mysql_native_password”插件认证
22 #mysql_native_password
23 default_authentication_plugin=mysql_native_password
```

配置说明：主要需要修改的是以下几个属性：

server-id：服务节点的唯一标识。需要给集群中的每个服务分配一个单独的ID。

log_bin：打开Binlog日志记录，并指定文件名。

log_bin-index：Binlog日志文件

重启MySQL服务， `service mysqld restart`

然后，我们需要给root用户分配一个replication slave的权限。

```
1 #登录主数据库
2 mysql -u root -p
3 GRANT REPLICATION SLAVE ON *.* TO 'root'@'%';
4 flush privileges;
5 #查看主节点同步状态:
6 show master status;
```

在实际生产环境中，通常不会直接使用root用户，而会创建一个拥有全部权限的用户来负责主从同步。

```
mysql> show master status;
```

File	Position	Binlog_Do_DB	Binlog_Ignore_DB	Executed_Gtid_Set
master-bin.000004	156			

```
1 row in set (0.00 sec)
```

这个指令结果中的File和Position记录的是当前日志的binlog文件以及文件中的索引。

而后面的Binlog_Do_DB和Binlog_Ignore_DB这两个字段是表示需要记录binlog文件的库以及不需要记录binlog文件的库。目前我们没有进行配置，就表示是针对全库记录日志。这两个字段如何进行配置，会在后面进行介绍。

开启binlog后，数据库中的所有操作都会被记录到datadir当中，以一组轮询文件的方式循环记录。而指令查到的File和Position就是当前日志的文件和位置。而在后面配置从服务时，就需要通过这个File和Position通知从服务从哪个地方开始记录binLog。

```

[root@worker1 mysql-files]# ll
总用量 179660
-rw-r----- 1 mysql mysql      56 12月 30 14:24 auto.cnf
-rw----- 1 mysql mysql     1676 12月 30 14:24 ca-key.pem
-rw-r--r-- 1 mysql mysql     1112 12月 30 14:24 ca.pem
-rw-r--r-- 1 mysql mysql     1112 12月 30 14:24 client-cert.pem
-rw----- 1 mysql mysql     1676 12月 30 14:24 client-key.pem
-rw-r----- 1 mysql mysql    196608 12月 30 14:33 #ib_16384_0.dblwr
-rw-r----- 1 mysql mysql    8585216 12月 30 14:24 #ib_16384_1.dblwr
-rw-r----- 1 mysql mysql      3768 12月 30 14:32 ib_buffer_pool
-rw-r----- 1 mysql mysql    12582912 12月 30 14:33 ibdata1
-rw-r----- 1 mysql mysql    50331648 12月 30 14:33 ib_logfile0
-rw-r----- 1 mysql mysql    50331648 12月 30 14:24 ib_logfile1
-rw-r----- 1 mysql mysql    12582912 12月 30 14:32 ibtmp1
drwxr-x--- 2 mysql mysql      187 12月 30 14:32 #innodb temp
-rw-r----- 1 mysql mysql      179 12月 30 14:32 master-bin.000001
-rw-r----- 1 mysql mysql     1829 12月 30 14:33 master-bin.000002
-rw-r----- 1 mysql mysql       40 12月 30 14:32 master-bin.index
drwxr-x--- 2 mysql mysql      143 12月 30 14:24 mysql
-rw-r----- 1 mysql mysql    28311552 12月 30 14:33 mysql.ibd
drwxr-x--- 2 mysql mysql     8192 12月 30 14:24 performance_schema
-rw----- 1 mysql mysql     1676 12月 30 14:24 private_key.pem
-rw-r--r-- 1 mysql mysql       452 12月 30 14:24 public_key.pem
-rw-r--r-- 1 mysql mysql     1112 12月 30 14:24 server-cert.pem
-rw----- 1 mysql mysql     1680 12月 30 14:24 server-key.pem
drwxr-x--- 2 mysql mysql       28 12月 30 14:24 sys
-rw-r----- 1 mysql mysql   10485760 12月 30 14:33 undo_001
-rw-r----- 1 mysql mysql   10485760 12月 30 14:33 undo_002
-rw-r----- 1 mysql mysql      2200 12月 30 14:32 worker1.err
-rw-r----- 1 mysql mysql        5 12月 30 14:32 worker1.pid
[root@worker1 mysql-files]# pwd
/usr/local/mysql/mysql-files

```

3.2 配置slave从服务

下一步，我们来配置从服务mysqls。我们打开mysqls的配置文件my.cnf，修改配置文件：

```

1  [mysqld]
2  #主库和从库需要不一致
3  server-id=48
4  #打开MySQL中继日志
5  relay-log-index=slave-relay-bin.index
6  relay-log=slave-relay-bin
7  #打开从服务二进制日志
8  log-bin=mysql-bin
9  #使得更新的数据写进二进制日志中
10 log-slave-updates=1
11 # 设置3306端口
12 port=3306
13 # 设置mysql的安装目录
14 basedir=/usr/local/mysql
15 # 设置mysql数据库的数据的存放目录
16 datadir=/usr/local/mysql/mysql-files
17 # 允许最大连接数
18 max_connections=200
19 # 允许连接失败的次数。
20 max_connect_errors=10
21 # 服务端使用的字符集默认为UTF8
22 character-set-server=utf8

```

```
23 # 创建新表时将使用的默认存储引擎
24 default-storage-engine=INNODB
25 # 默认使用“mysql_native_password”插件认证
26 #mysql_native_password
27 default_authentication_plugin=mysql_native_password
```

配置说明：主要需要关注的几个属性：

server-id：服务节点的唯一标识

relay-log：打开从服务的relay-log日志。

log-bin：打开从服务的bin-log日志记录。

然后我们启动mysqls的服务，并设置他的主节点同步状态。

```
1 #登录从服务
2 mysql -u root -p;
3 #设置同步主节点:
4 CHANGE MASTER TO
5 MASTER_HOST='192.168.232.128',
6 MASTER_PORT=3306,
7 MASTER_USER='root',
8 MASTER_PASSWORD='root',
9 MASTER_LOG_FILE='master-bin.000004',
10 MASTER_LOG_POS=156
11 GET_MASTER_PUBLIC_KEY=1;
12 #开启slave
13 start slave;
14 #查看主从同步状态
15 show slave status;
16 或者用 show slave status \G; 这样查看比较简洁
```

注意，CHANGE MASTER指令中需要指定的MASTER_LOG_FILE和MASTER_LOG_POS必须与主服务中查到的保持一致。

并且后续如果要检查主从架构是否成功，也可以通过检查主服务与从服务之间的File和Position这两个属性是否一致来确定。


```
mysql> show slave status \G;
***** 1. row *****
      Slave_IO_State: Waiting for master to send event
      Master_Host: 127.0.0.1
      Master_User: root
      Master_Port: 3307
      Connect_Retry: 60
      Master_Log_File: master-bin.000004
      Read_Master_Log_Pos: 156
      Relay_Log_File: slave-relay-bin.000006
      Relay_Log_Pos: 373
      Relay_Master_Log_File: master-bin.000004
      Slave_IO_Running: Yes
      Slave_SQL_Running: Yes
      Replicate_Do_DB:
      Replicate_Ignore_DB:
      Replicate_Do_Table:
      Replicate_Ignore_Table:
      Replicate_Wild_Do_Table:
      Replicate_Wild_Ignore_Table:
      Last_Errno: 0
      Last_Error:
      Skip_Counter: 0
```

我们重点关注其中红色方框的两个属性，与主节点保持一致，就表示这个主从同步搭建是成功的。

从这个指令的结果能够看到，有很多Replicate_开头的属性，这些属性指定了两个服务之间要同步哪些数据库、哪些表的配置。只是在我们这个示例中全都没有进行配置，就标识是全库进行同步。后面我们会补充如何配置需要同步的库和表。

3.3 主从集群测试

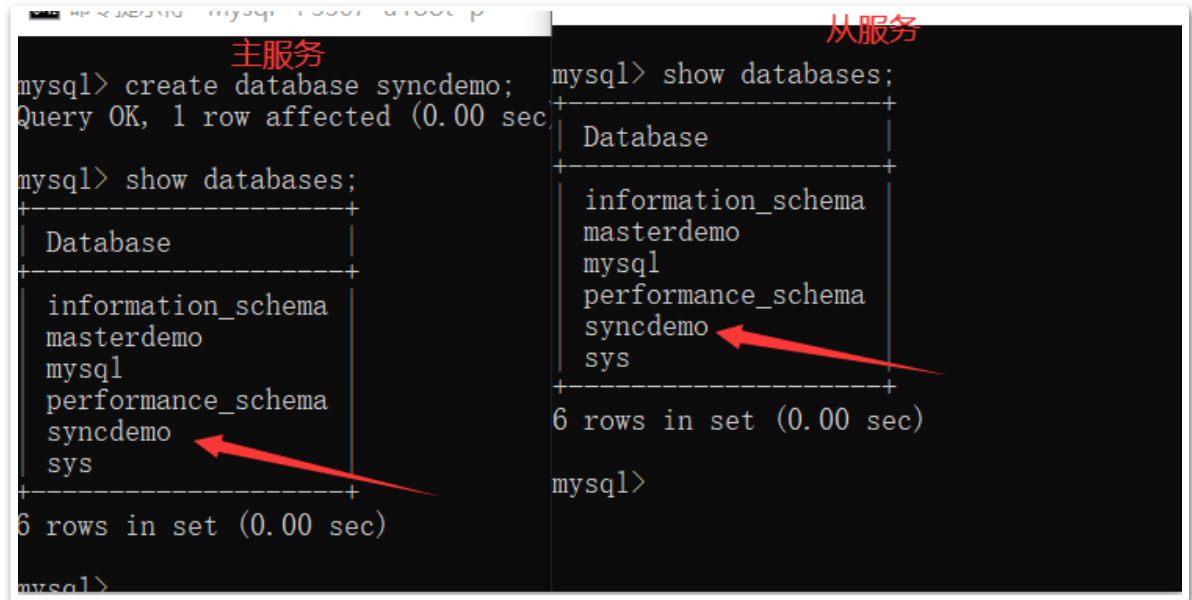
测试时，我们先用showdatabases，查看下两个MySQL服务中的数据库情况

主服务	从服务
<pre>mysql> show databases; +-----+ Database +-----+ information_schema masterdemo mysql performance_schema sys +-----+ 5 rows in set (0.00 sec)</pre>	<pre>mysql> show databases; +-----+ Database +-----+ information_schema masterdemo mysql performance_schema sys +-----+ 5 rows in set (0.01 sec)</pre>

然后我们在主服务器上创建一个数据库

```
1 | mysql> create database syncdemo;
2 | Query OK, 1 row affected (0.00 sec)
```

然后我们再用show databases，来看下这个syncdemo的数据库是不是已经同步到了从服务。



接下来我们继续在syncdemo这个数据库中创建一个表，并插入一条数据。

```
1 | mysql> use syncdemo;
2 | Database changed
3 | mysql> create table demoTable(id int not null);
4 | Query OK, 0 rows affected (0.02 sec)
5 |
6 | mysql> insert into demoTable value(1);
7 | Query OK, 1 row affected (0.01 sec)
```

然后我们也同样到主服务与从服务上都来查一下这个demoTable是否同步到了从服务。

主服务	从服务
<pre>syncdemo sys 6 rows in set (0.00 sec) mysql> use syncdemo; Database changed mysql> create table demoTable(id int); Query OK, 0 rows affected (0.02 sec) mysql> insert into demoTable value (1); Query OK, 1 row affected (0.01 sec) mysql> show tables; +-----+ Tables_in_syncdemo +-----+ demotable +-----+ 1 row in set (0.00 sec) mysql> select * from demotable; +----+ id +----+ 1 +----+ 1 row in set (0.00 sec) mysql></pre>	<pre>Database +-----+ information_schema masterdemo mysql performance_schema syncdemo sys +-----+ 6 rows in set (0.00 sec) mysql> use syncdemo; Database changed mysql> show tables; +-----+ Tables_in_syncdemo +-----+ demotable +-----+ 1 row in set (0.00 sec) mysql> select * from demotable; +----+ id +----+ 1 +----+ 1 row in set (0.00 sec) mysql></pre>

主服务上创建的
demoTable表已经同步到
了从服务

从上面的实验过程看到，我们在主服务中进行的数据操作，就都已经同步到了从服务上。这样，我们一个主从集群就搭建完成了。

另外，这个主从架构是有可能失败的，如果在slave从服务上查看slave状态，发现Slave_SQL_Running=no，就表示主从同步失败了。这有可能是因为在从数据库上进行了写操作，与同步过来的SQL操作冲突了，也有可能是slave从服务重启后有事务回滚了。

如果是因为slave从服务事务回滚的原因，可以按照以下方式重启主从同步：

```
1 | mysql> stop slave ;
2 | mysql> set GLOBAL SQL_SLAVE_SKIP_COUNTER=1;
3 | mysql> start slave ;
```

而另一种解决方式就是重新记录主节点的binlog文件消息

```
1 | mysql> stop slave ;
2 | mysql> change master to .....
3 | mysql> start slave ;
```

但是这种方式要注意binlog的文件和位置，如果修改后和之前的同步接不上，那就会丢失部分数据。所以不太常用。

3.4 集群搭建扩展：

在完成这个基本的MySQL主从集群后，我们还可以进行后续的实验：

1、全库同步与部分同步

之前提到，我们目前配置的主从同步是针对全库配置的，而实际环境中，一般并不需要针对全库做备份，而只需要对一些特别重要的库或者表来进行同步。那如何针对库和表做同步配置呢？

首先在Master端：在my.cnf中，可以通过以下这些属性指定需要针对哪些库或者哪些表记录binlog

```
1  #需要同步的二进制数据库名
2  binlog-do-db=masterdemo
3  #只保留7天的二进制日志，以防磁盘被日志占满(可选)
4  expire-logs-days = 7
5  #不备份的数据库
6  binlog-ignore-db=information_schema
7  binlog-ignore-db=performance_schema
8  binlog-ignore-db=sys
```

然后在Slave端：在my.cnf中，需要配置备份库与主服务的库的对应关系。

```
1  #如果slave库名称与master库名相同，使用本配置
2  replicate-do-db = masterdemo
3  #如果master库名[mastdemo]与slave库名[mastdemo01]不同，使用以下配置[需要做映射]
4  replicate-rewrite-db = masterdemo -> masterdemo01
5  #如果不是要全部同步[默认全部同步]，则指定需要同步的表
6  replicate-wild-do-table=masterdemo01.t_dict
7  replicate-wild-do-table=masterdemo01.t_num
```

配置完成了之后，在show master status指令中，就可以看到Binlog_Do_DB和Binlog_Ignore_DB两个参数的作用了。

2、读写分离配置

我们要注意，目前我们的这个MySQL主从集群是单向的，也就是只能从主服务同步到从服务，而从服务的数据表更是无法同步到主服务的。

```
mysql> 主服务
mysql> 插入记录为3的数据
mysql> insert into demotable value(3);
Query OK, 1 row affected (0.01 sec)

mysql> select * from demotable;
+----+
| id |
+----+
| 1  |
| 3  |
+----+
2 rows in set (0.00 sec)

mysql>

mysql> 从服务
mysql>
mysql> insert into demotable value(2);
Query OK, 1 row affected (0.00 sec)

mysql> select * from demotable;
+----+
| id |
+----+
| 1  |
| 2  |
| 3  |
+----+
3 rows in set (0.00 sec)

mysql>
```

所以，在这种架构下，为了保证数据一致，通常会需要保证数据只在主服务上写，而从服务只进行数据读取。这个功能，就是大名鼎鼎的读写分离。但是这里要注意下，mysql主从本身是无法提供读写分离的服务的，需要由业务自己来实现。这也是我们后面要学的ShardingSphere的一个重要功能。

到这里可以看到，在MySQL主从架构中，是需要严格限制从服务的数据写入的，一旦从服务有数据写入，就会造成数据不一致。并且从服务在执行事务期间还很容易造成数据同步失败。

如果需要限制用户写数据，我们可以在从服务中将`read_only`参数的值设为1(`set global read_only=1;`)。这样就可以限制用户写入数据。但是这个属性有两个需要注意的地方：

- 1、`read_only=1`设置的只读模式，不会影响slave同步复制的功能。所以在MySQL slave库中设定了`read_only=1`后，通过"`show slave status\G`"命令查看slave状态，可以看到slave仍然会读取master上的日志，并且在slave库中应用日志，保证主从数据库同步一致；
- 2、`read_only=1`设置的只读模式，限定的是普通用户进行数据修改的操作，但不会限定具有super权限的用户的数据修改操作。在MySQL中设置`read_only=1`后，普通的应用用户进行insert、update、delete等会产生数据变化的DML操作时，都会报出数据库处于只读模式不能发生数据变化的错误，但具有super权限的用户，例如在本地或远程通过root用户登录到数据库，还是可以进行数据变化的DML操作；如果需要限定super权限的用户写数据，可以设置`super_read_only=0`。另外 **如果想连super权限用户的写操作也禁止，就使用"`flush tables with read lock;`"，这样设置也会阻止主从同步复制！**

3、其他集群方式

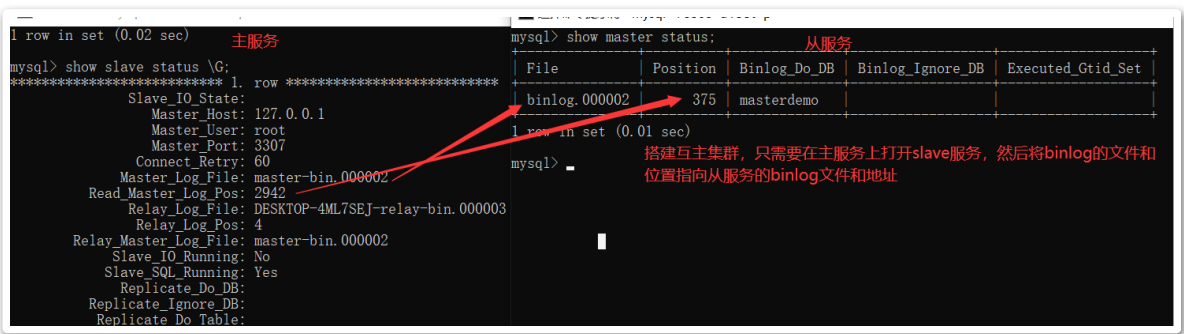
我们到这里搭建出了一个一主一从的MySQL主从同步集群，具有了数据同步的基础功能。而在生产环境中，通常会以此为基础，根据业务情况以及负载情况，搭建更大更复杂的集群。

例如为了进一步提高整个集群的读能力，可以扩展出一主多从。而为了减轻主节点进行数据同步的压力，可以继续扩展出多级从的主从集群。

为了提高整个集群的高可用能力，可以扩展出多主的集群。

我们也可以扩展出互为主从的互主集群甚至是环形的主从集群，实现MySQL多活部署。

搭建互主集群只需要按照上面的方式，在主服务上打开一个slave进程，并且指向slave节点的binlog当前文件地址和位置。



```
mysql> show slave status \G;
1 row in set (0.02 sec)

Slave_IO_State:
Master_Host: 127.0.0.1
Master_User: root
Master_Port: 3307
Connect_Retry: 60
Master_Log_File: master-bin.000002
Read_Master_Log_Pos: 2942
Relay_Log_File: DESKTOP-4ML7SEJ-relay-bin.000003
Relay_Log_Pos: 4
Relay_Master_Log_File: master-bin.000002
Slave_IO_Running: No
Slave_SQL_Running: Yes
Replicate_Do_DB:
Replicate_Ignore_DB:
Replicate_Do_Table:

mysql> show master status;
1 row in set (0.01 sec)

File: binlog.000002
Position: 375
Binlog_Do_DB: masterdemo
Binlog_Ignore_DB:
Executed_Gtid_Set:
```

搭建互主集群，只需要在主服务上打开slave服务，然后将binlog的文件和位置指向从服务的binlog文件和地址

我们这里是使用的最为传统的Binlog方式搭建集群，是基于日志记录点的方式来进行主从同步的。在这个实验中，Executed_Grid_Set一列，实际上就是另外一种搭建主从同步的方式，即GTID搭建方式。GTID的本质也是基于Binlog来实现的主从同步，只是他会基于一个全局的事务ID来标识同步进度。这个GTID全局事务ID是一个全局唯一、并且趋势递增的分布式ID策略。我们这里就不再去搭建了。

3.5、GTID同步集群

上面我们搭建的集群方式，是基于Binlog日志记录点的方式来搭建的，这也是最为传统的MySQL集群搭建方式。而在这个实验中，可以看到有一个Executed_Grid_Set列，暂时还没有用上。实际上，这就是另外一种搭建主从同步的方式，即GTID搭建方式。这种模式是从MySQL5.6版本引入的。

GTID的本质也是基于Binlog来实现主从同步，只是他会基于一个全局的事务ID来标识同步进度。GTID即全局事务ID，全局唯一并且趋势递增，他可以保证为每一个在主节点上提交的事务在复制集群中可以生成一个唯一的ID。

在基于GTID的复制中，首先从服务器会告诉主服务器已经在从服务器执行完了哪些事务的GTID值，然后主库会有把所有没有在从库上执行的事务，发送到从库上进行执行，并且使用GTID的复制可以保证同一个事务只在指定的从库上执行一次，这样可以避免由于偏移量的问题造成数据不一致。

他的搭建方式跟我们上面的主从架构整体搭建方式差不多。只是需要在my.cnf中修改一些配置。

在主节点上：

```
1 | gtid_mode=on
2 | enforce_gtid_consistency=on
3 | log_bin=on
4 | server_id=单独设置一个
5 | binlog_format=row
```

在从节点上：

```
1 | gtid_mode=on
2 | enforce_gtid_consistency=on
3 | log_slave_updates=1
4 | server_id=单独设置一个
```

然后分别重启主服务和从服务，就可以开启GTID同步复制方式。

4、集群扩容

我们现在已经搭建成功了一主一从的MySQL集群架构，那要扩展到一主多从的集群架构，其实就比较简单了，只需要增加一个binlog复制就行了。

但是如果我们的集群是已经运行过一段时间，这时候如果要扩展新的从节点就有一个问题，之前的数据没办法从binlog来恢复了。这时候在扩展新的slave节点时，就需要增加一个数据复制的操作。

MySQL的数据备份恢复操作相对比较简单，可以通过SQL语句直接来完成。具体操作可以使用mysql的bin目录下的mysqldump工具。


```
1 | mysqldump -u root -p --all-databases > backup.sql
2 | #输入密码
```

通过这个指令，就可以将整个数据库的所有数据导出成backup.sql，然后把这个backup.sql分发到新的MySQL服务器上，并执行下面的指令将数据全部导入到新的MySQL服务中。

```
1 | mysql -u root -p < backup.sql
2 | #输入密码
```

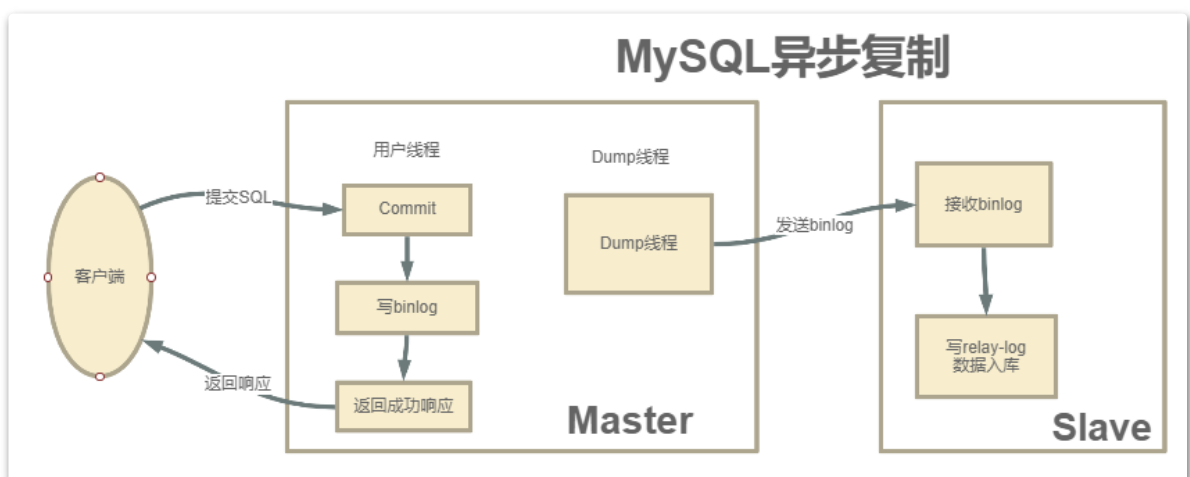
这样新的MySQL服务就已经有了所有的历史数据，然后就可以再按照上面的步骤，配置Slave从服务的数据同步了。

5、半同步复制

1、理解半同步复制

到现在为止，我们已经可以搭建MySQL的主从集群，互主集群，但是我们这个集群有一个隐患，就是有可能会丢数据。这是为什么呢？这要从MySQL主从数据复制分析起。

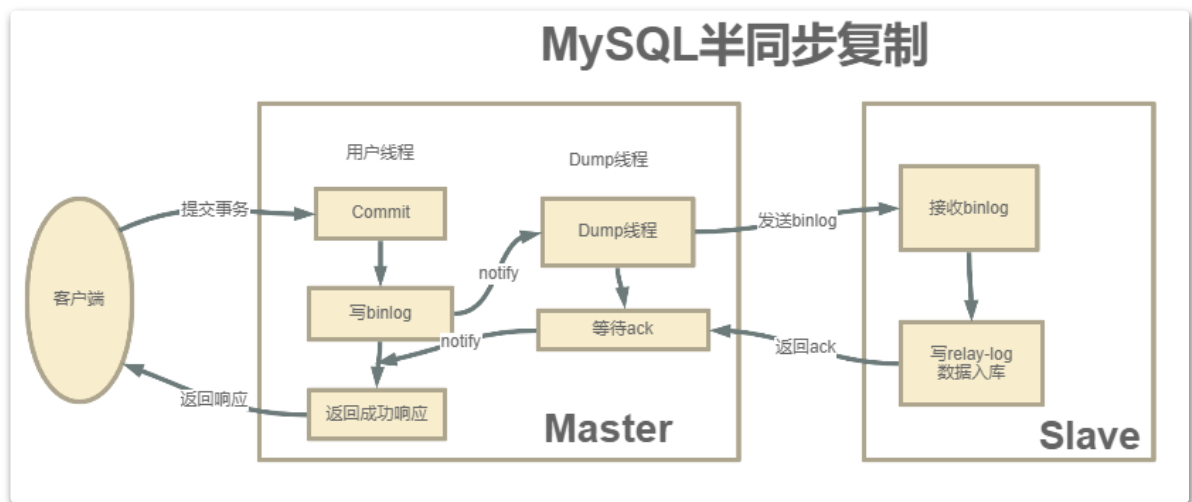
MySQL主从集群默认采用的是一种异步复制的机制。主服务在执行用户提交的事务后，写入binlog日志，然后就给客户端返回一个成功的响应了。而binlog会由一个dump线程异步发送给Slave从服务。



由于这个发送binlog的过程是异步的。主服务在向客户端反馈执行结果时，是不知道binlog是否同步成功了的。这时候如果主服务宕机了，而从服务还没有备份到新执行的binlog，那就有可能会丢数据。

那怎么解决这个问题呢，这就要靠MySQL的半同步复制机制来保证数据安全。

半同步复制机制是一种介于异步复制和全同步复制之前的机制。主库在执行完客户端提交的事务后，并不是立即返回客户端响应，而是等待至少一个从库接收并写到relay log中，才会返回给客户端。MySQL在等待确认时，默认会等10秒，如果超过10秒没有收到ack，就会降级成为异步复制。



这种半同步复制相比异步复制，能够有效的提高数据的安全性。但是这种安全性也不是绝对的，他只保证事务提交后的binlog至少传输到了一个从库，并且并不保证从库应用这个事务的binlog是成功的。另一方面，半同步复制机制也会造成一定程度的延迟，这个延迟时间最少是一个TCP/IP请求往返的时间。整个服务的性能是会有所下降的。而当从服务出现问题时，主服务需要等待的时间就会更长，要等到从服务的服务恢复或者请求超时才能给用户响应。

2、搭建半同步复制集群

半同步复制需要基于特定的扩展模块来实现。而mysql从5.5版本开始，往上的版本都默认自带了这个模块。这个模块包含在mysql安装目录下的lib/plugin目录下的semisync_master.so和semisync_slave.so两个文件中。需要在主服务上安装semisync_master模块，在从服务上安装semisync_slave模块。

```
[root@worker1 plugin]# pwd
/usr/local/mysql/lib/plugin
[root@worker1 plugin]# ll
总用量 90952
-rwxr-xr-x 1 7161 31415 166168 3月 26 2020 adt_null.so
-rwxr-xr-x 1 7161 31415 639696 3月 26 2020 authentication_ldap_sasl_client.so
-rwxr-xr-x 1 7161 31415 113784 3月 26 2020 auth_socket.so
-rwxr-xr-x 1 7161 31415 271816 3月 26 2020 component_audit_api_message_emit.so
-rwxr-xr-x 1 7161 31415 192920 3月 26 2020 component_log_filter_dragnet.so
-rwxr-xr-x 1 7161 31415 36968 3月 26 2020 component_log_sink_json.so
-rwxr-xr-x 1 7161 31415 146024 3月 26 2020 component_log_sink_syseventlog.so
-rwxr-xr-x 1 7161 31415 566016 3月 26 2020 component_mysqlbackup.so
-rwxr-xr-x 1 7161 31415 510288 3月 26 2020 component_validate_password.so
-rwxr-xr-x 1 7161 31415 1311200 3月 26 2020 connection_control.so
-rwxr-xr-x 1 7161 31415 3183232 3月 26 2020 ddl_rewriter.so
drwxr-xr-x 2 7161 31415 4096 3月 26 2020 debug
-rwxr-xr-x 1 7161 31415 60459336 3月 26 2020 group_replication.so
-rwxr-xr-x 1 7161 31415 633896 3月 26 2020 ha_example.so
-rwxr-xr-x 1 7161 31415 1303856 3月 26 2020 ha_mock.so
-rwxr-xr-x 1 7161 31415 1229744 3月 26 2020 innodb_engine.so
-rwxr-xr-x 1 7161 31415 2648040 3月 26 2020 keyring_file.so
-rwxr-xr-x 1 7161 31415 518872 3月 26 2020 keyring_udf.so
-rwxr-xr-x 1 7161 31415 1208848 3月 26 2020 libmemcached.so
-rwxr-xr-x 1 7161 31415 9064128 3月 26 2020 libpluginmecab.so
-rwxr-xr-x 1 7161 31415 91000 3月 26 2020 locking_service.so
-rwxr-xr-x 1 7161 31415 115064 3月 26 2020 mypluglib.so
-rwxr-xr-x 1 7161 31415 2675360 3月 26 2020 mysql_clone.so
-rwxr-xr-x 1 7161 31415 112048 3月 26 2020 mysql_no_login.so
-rwxr-xr-x 1 7161 31415 114176 3月 26 2020 rewrite_example.so
-rwxr-xr-x 1 7161 31415 1570552 3月 26 2020 rewriter.so
-rwxr-xr-x 1 7161 31415 1646120 3月 26 2020 semisync_master.so
-rwxr-xr-x 1 7161 31415 750408 3月 26 2020 semisync_slave.so
-rwxr-xr-x 1 7161 31415 438736 3月 26 2020 validate_password.so
-rwxr-xr-x 1 7161 31415 1340600 3月 26 2020 version_token.so
```

首先我们登陆主服务，安装semisync_master模块：

```
1 mysql> install plugin rpl_semi_sync_master soname 'semisync_master.so';
2 Query OK, 0 rows affected (0.01 sec)
3
4 mysql> show global variables like 'rpl_semi%';
5 +-----+-----+
6 | Variable_name | Value |
7 +-----+-----+
8 | rpl_semi_sync_master_enabled | OFF |
9 | rpl_semi_sync_master_timeout | 10000 |
10 | rpl_semi_sync_master_trace_level | 32 |
11 | rpl_semi_sync_master_wait_for_slave_count | 1 |
12 | rpl_semi_sync_master_wait_no_slave | ON |
13 | rpl_semi_sync_master_wait_point | AFTER_SYNC |
14 +-----+-----+
15 6 rows in set, 1 warning (0.02 sec)
16
17 mysql> set global rpl_semi_sync_master_enabled=ON;
18 Query OK, 0 rows affected (0.00 sec)
```

这三行指令中，第一行是通过扩展库来安装半同步复制模块，需要指定扩展库的文件名。

第二行查看系统全局参数，`rpl_semi_sync_master_timeout`就是半同步复制时等待应答的最长等待时间，默认是10秒，可以根据情况自行调整。

第三行则是打开半同步复制的开关。

在第二行查看系统参数时，最后的一个参数
`rpl_semi_sync_master_wait_point`其实表示一种半同步复制的方式。

半同步复制有两种方式，一种是我们现在看到的这种默认的
`AFTER_SYNC`方式。这种方式下，主库把日志写入binlog，并且复制给从库，然后开始等待从库的响应。从库返回成功后，主库再提交事务，接着给客户端返回一个成功响应。

而另一种方式是叫做`AFTER_COMMIT`方式。他不是默认的。这种方式，在主库写入binlog后，等待binlog复制到从库，主库就提交自己的本地事务，再等待从库返回给自己一个成功响应，然后主库再给客户端返回响应。

然后我们登陆从服务，安装`semisync_slave`模块

```
1 mysql> install plugin rpl_semi_sync_slave soname 'semisync_slave.so';
2 Query OK, 0 rows affected (0.01 sec)
3
4 mysql> show global variables like 'rpl_semi%';
5 +-----+-----+
6 | Variable_name          | Value |
7 +-----+-----+
8 | rpl_semi_sync_slave_enabled | OFF   |
9 | rpl_semi_sync_slave_trace_level | 32    |
10 +-----+-----+
11 2 rows in set, 1 warning (0.01 sec)
12
13 mysql> set global rpl_semi_sync_slave_enabled = on;
14 Query OK, 0 rows affected (0.00 sec)
15
16 mysql> show global variables like 'rpl_semi%';
17 +-----+-----+
18 | Variable_name          | Value |
19 +-----+-----+
20 | rpl_semi_sync_slave_enabled | ON    |
21 | rpl_semi_sync_slave_trace_level | 32    |
22 +-----+-----+
23 2 rows in set, 1 warning (0.00 sec)
```

```
24
25 mysql> stop slave;
26 Query OK, 0 rows affected (0.01 sec)
27
28 mysql> start slave;
29 Query OK, 0 rows affected (0.01 sec)
```

slave端的安装过程基本差不多，不过要注意下安装完slave端的半同步插件后，需要重启下slave服务。

6、主从架构的数据延迟问题：

在我们搭建的这个主从集群中，有一个比较隐藏的问题，就是这样的主从复制之间会有延迟。这在做了读写分离后，会更容易体现出来。即数据往主服务写，而读数据在从服务读。这时候这个主从复制延迟就有可能造成刚插入了数据但是查不到。当然，这在我们目前的这个集群中是很难出现的，但是在大型集群中会很容易出现。

出现这个问题的根本在于：面向业务的主服务数据都是多线程并发写入的，而从服务是单个线程慢慢拉取binlog，这中间就会有效率差。所以解决这个问题的关键是要让从服务也用多线程并行复制binlog数据。

MySQL自5.7版本后就已经支持并行复制了。可以在从服务上设置slave_parallel_workers为一个大于0的数，然后把slave_parallel_type参数设置为LOGICAL_CLOCK，这就可以了。

四、MySQL的高可用方案--了解

我们之前的MySQL服务集群，都是使用MySQL自身的功能来搭建的集群。但是这样的集群，不具备高可用的功能。即如果是MySQL主服务挂了，从服务是没办法自动切换成主服务的。而如果要实现MySQL的高可用，需要借助一些第三方工具来实现。

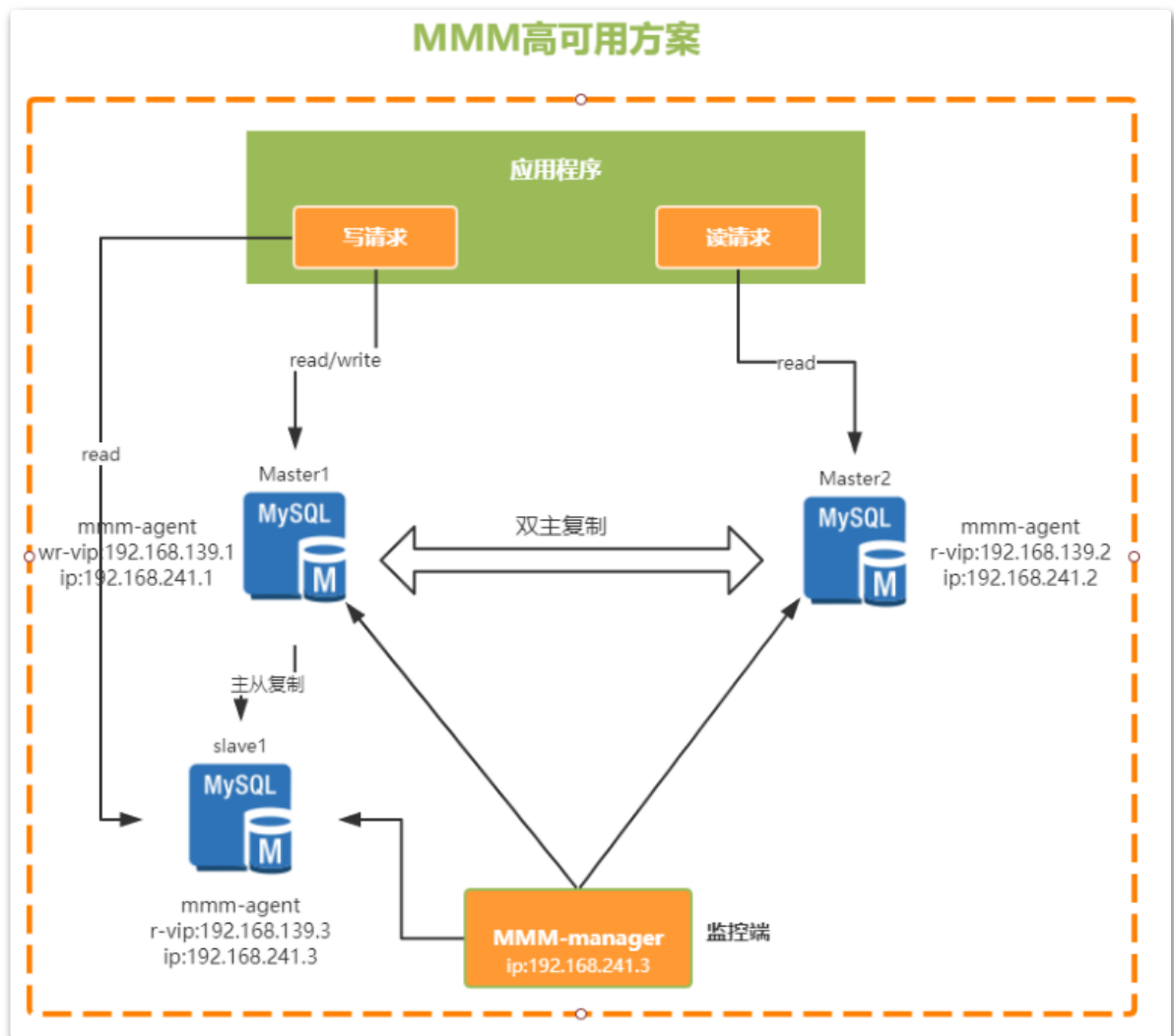
常见的MySQL集群方案有三种: MMM、MHA、MGR。这三种高可用框架都有一些共同点：

- 对主从复制集群中的Master节点进行监控
- 自动的对Master进行迁移，通过VIP。
- 重新配置集群中的其它slave对新的Master进行同步

1、MMM

MMM(Master-Master replication manager for Mysql, Mysql主主复制管理器)是一套由Perl语言实现的脚本程序，可以对mysql集群进行监控和故障迁移。他需要两个Master，同一时间只有一个Master对外提供服务，可以说是主备模式。

他是通过一个VIP(虚拟IP)的机制来保证集群的高可用。整个集群中，在主节点上会通过一个VIP地址来提供数据读写服务，而当出现故障时，VIP就会从原来的主节点漂移到其他节点，由其他节点提供服务。



优点：

- 提供了读写VIP的配置，使读写请求都可以达到高可用
- 工具包相对比较完善，不需要额外的开发脚本
- 完成故障转移之后可以对MySQL集群进行高可用监控

缺点：

- 故障简单粗暴，容易丢失事务，建议采用半同步复制方式，减少失败的概率

- 目前MMM社区已经缺少维护，不支持基于GTID的复制

适用场景：

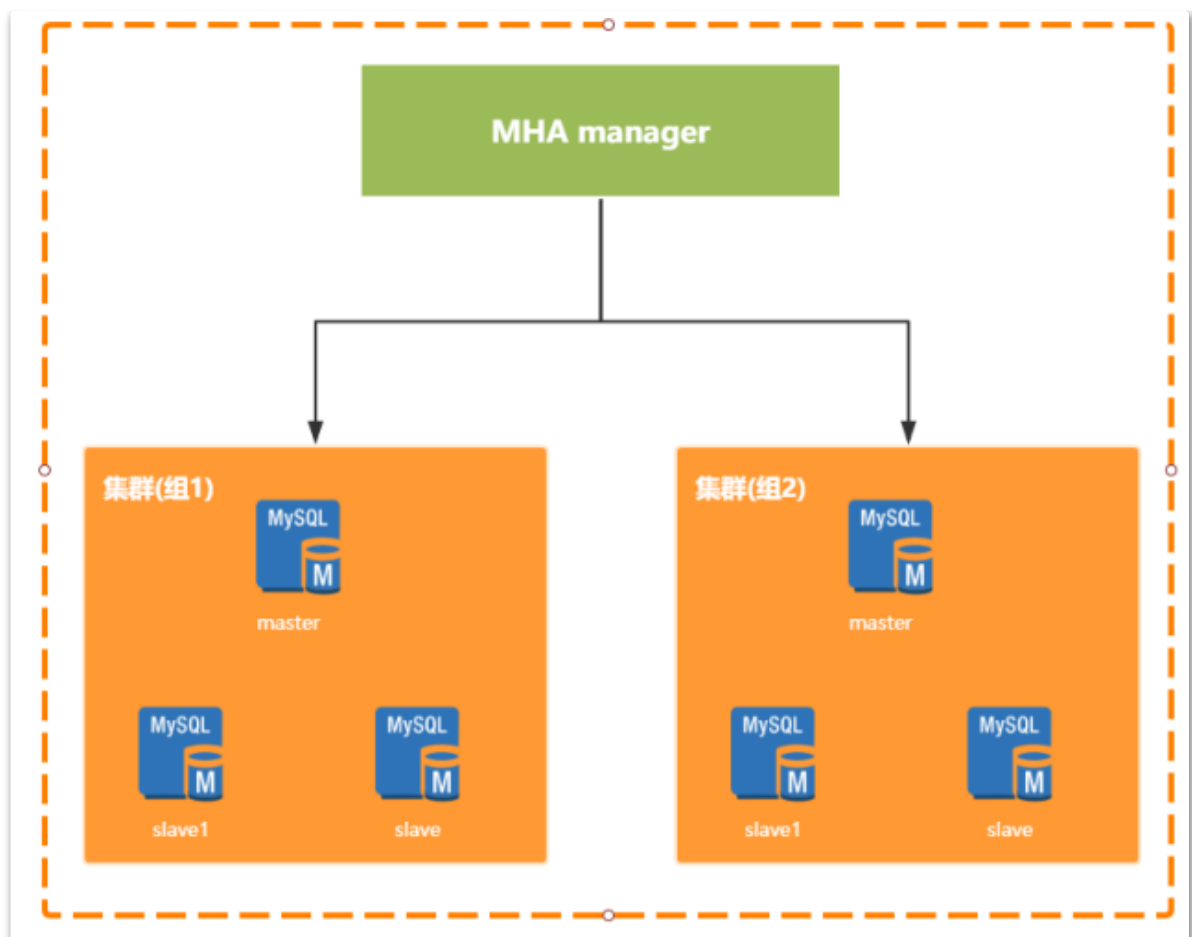
- 读写都需要高可用的
- 基于日志点的复制方式

2、MHA

Master High Availability Manager and Tools for MySQL。是由日本人开发的一个基于Perl脚本写的工具。这个工具专门用于监控主库的状态，当发现master节点故障时，会提升其中拥有新数据的slave节点成为新的master节点，在此期间，MHA会通过其他从节点获取额外的信息来避免数据一致性方面的问题。MHA还提供了mater节点的在线切换功能，即按需切换master-slave节点。MHA能够在30秒内实现故障切换，并能在故障切换过程中，最大程度的保证数据一致性。在淘宝内部，也有一个相似的TMHA产品。

MHA是需要单独部署的，分为Manager节点和Node节点，两种节点。其中Manager节点一般是单独部署的一台机器。而Node节点一般是部署在每台MySQL机器上的。Node节点得通过解析各个MySQL的日志来进行一些操作。

Manager节点会通过探测集群里的Node节点去判断各个Node所在机器上的MySQL运行是否正常，如果发现某个Master故障了，就直接把他的一个Slave提升为Master，然后让其他Slave都挂到新的Master上去，完全透明。



优点：

- MHA除了支持日志点的复制还支持GTID的方式
- 同MMM相比，MHA会尝试从旧的Master中恢复旧的二进制日志，只是未必每次都能成功。如果希望更少的数据丢失场景，建议使用MHA架构。

缺点：

MHA需要自行开发VIP转移脚本。

MHA只监控Master的状态，未监控Slave的状态

3、MGR

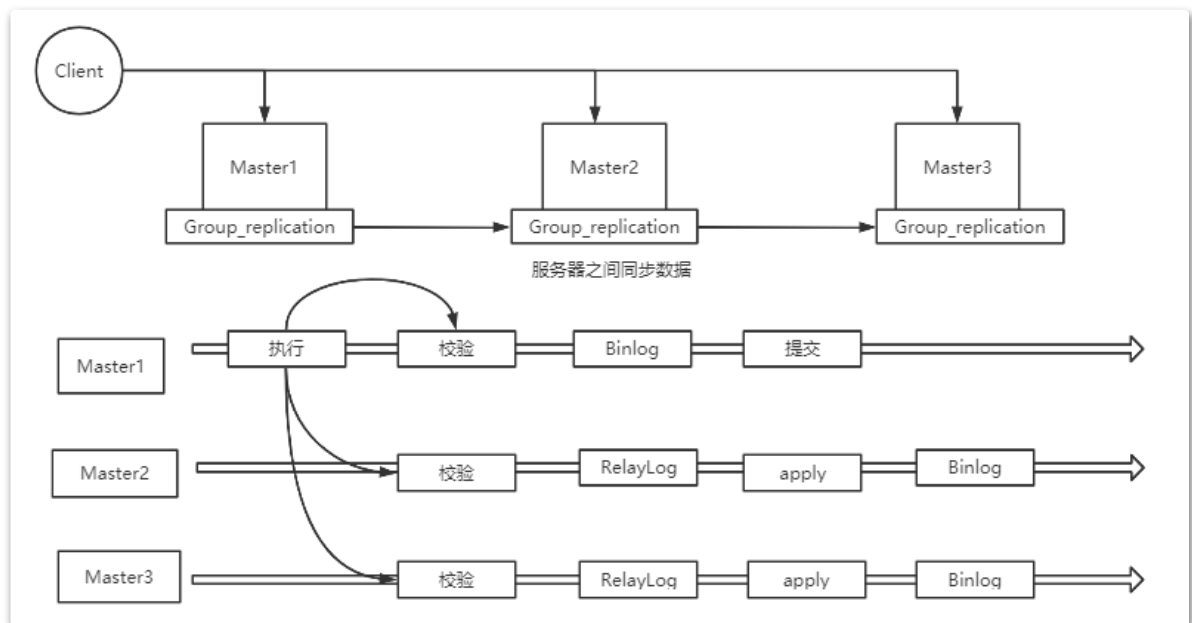
MGR: MySQL Group Replication。是MySQL官方在5.7.17版本正式推出的一种组复制机制。主要是解决传统异步复制和半同步复制的数据一致性问题。

由若干个节点共同组成一个复制组，一个事务提交后，必须经过超过半数节点的决议并通过后，才可以提交。引入组复制，主要是为了解决传统异步复制和半同步复制可能产生数据不一致的问题。MGR依靠分布式一致性协议(Paxos协议的一个变体)，实现了分布式下数据的最终一致性，提供了真正的数据高可用方案(方案落地后

是否可靠还有待商榷)。

支持多主模式，但官方推荐单主模式：

- 多主模式下，客户端可以随机向MySQL节点写入数据
- 单主模式下，MGR集群会选出primary节点负责写请求，primary节点与其它节点都可以进行读请求处理。



优点：

- 基本无延迟，延迟比异步的小很多
- 支持多写模式，但是目前还不是很成熟
- 数据的强一致性，可以保证数据事务不丢失

缺点：

- 仅支持innodb，且每个表必须提供主键。
- 只能用在GTID模式下，且日志格式为row格式。

适用的业务场景：

- 对主从延迟比较敏感
- 希望对对写服务提供高可用，又不想安装第三方软件
- 数据强一致的场景

五、分库分表

前面我们做的一大段实验，目的是为了大家能够理解MySQL的主从集群。而主从集群的作用，在我们开发角度更大的是作为读写分离的支持，也是我们后面学习ShardingSphere的重点。我们这一部分就来介绍下分库分表。

分库分表就是业务系统将数据写请求分发到master节点，而读请求分发到slave节点的一种方案，可以大大提高整个数据库集群的性能。但是要注意，分库分表的一整套逻辑全部是由客户端自行实现的。而对于MySQL集群，数据主从同步是实现读写分离的一个必要前提条件。

1、分库分表有什么用

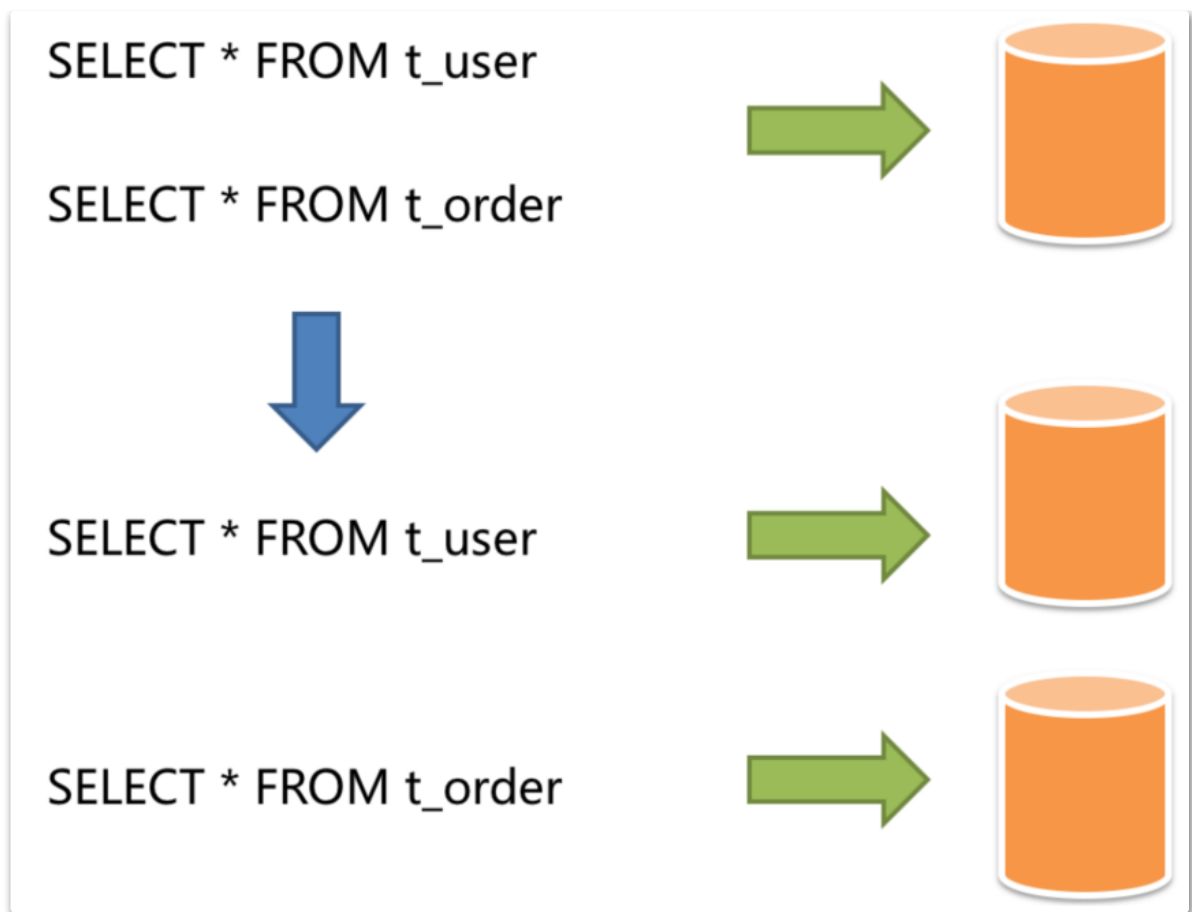
分库分表就是为了解决由于数据量过大而导致数据库性能降低的问题，将原来独立的数据库拆分成若干数据库组成，将数据大表拆分成若干数据表组成，使得单一数据库、单一数据表的数据量变小，从而达到提升数据库性能的目的。

例如：微服务架构中，每个服务都分配一个独立的数据库，这就是分库。而对一些业务日志表，按月拆分成不同的表，这就是分表。

2、分库分表的方式

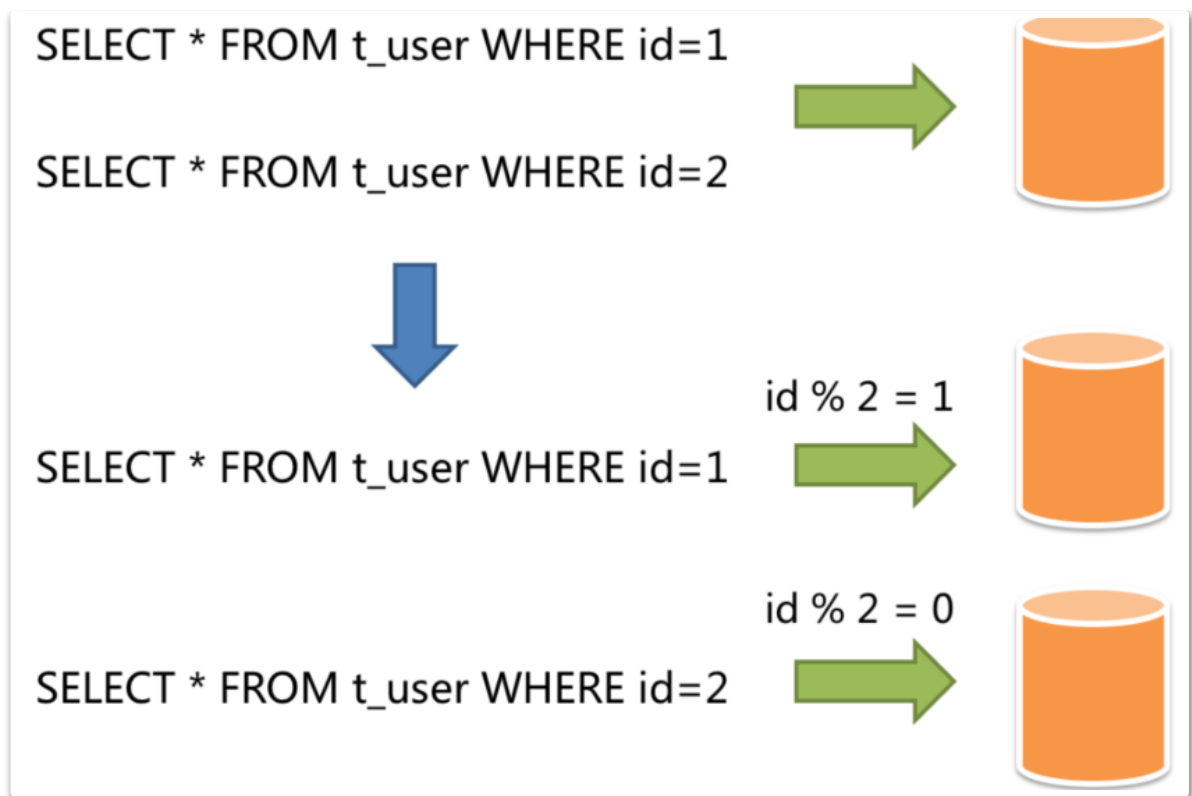
分库分表包含分库和分表 两个部分，而这两个部分可以统称为数据分片，其目的都是将数据拆分成不同的存储单元。另外，从分拆的角度上，可以分为垂直分片和水平分片。

- 垂直分片：按照业务来对数据进行分片，又称为纵向分片。他的核心理念就是转库专用。在拆分之前，一个数据库由多个数据表组成，每个表对应不同的业务。而拆分之后，则是按照业务将表进行归类，分布到不同的数据库或表中，从而将压力分散至不同的数据库或表。例如，下图将用户表和订单表垂直分片到不同的数据库：



垂直分片往往需要对架构和设计进行调整。通常来讲，是来不及应对业务需求快速变化的。而且，他也无法真正的解决单点数据库的性能瓶颈。垂直分片可以缓解数据量和访问量带来的问题，但无法根治。如果垂直分片之后，表中的数据量依然超过单节点所能承载的阈值，则需要水平分片来进一步处理。

- 水平分片：又称横向分片。相对于垂直分片，它不再将数据根据业务逻辑分类，而是通过某个字段(或某几个字段)，根据某种规则将数据分散至多个库或表中，每个分片仅包含数据的一部分。例如，像下图根据主键机构分片。



常用的分片策略有：

取余\取模：优点 均匀存放数据，缺点 扩容非常麻烦

按照范围分片：比较好扩容，数据分布不够均匀

按照时间分片：比较容易将热点数据区分出来。

按照枚举值分片：例如按地区分片

按照目标字段前缀指定进行分区：自定义业务规则分片

水平分片从理论上突破了单机数据量处理的瓶颈，并且扩展相对自由，是分库分表的标准解决方案。

一般来说，在系统设计阶段就应该根据业务耦合松紧来确定垂直分库，垂直分表方案，在数据量及访问压力不是特别大的情况，首先考虑缓存、读写分离、索引技术等方案。若数据量极大，且持续增长，再考虑水平分库水平分表方案

3、分库分表的缺点

虽然数据分片解决了性能、可用性以及单点备份恢复等问题，但是分布式的架构在获得收益的同时，也引入了非常多新的问题。

- 事务一致性问题

原本单机数据库有很好的事务机制能够帮我们保证数据一致性。但是分库分表后，由于数据分布在不同库甚至不同服务器，不可避免会带来分布式事务问题。

- 跨节点关联查询问题

在没有分库时，我们可以进行很容易的进行跨表的关联查询。但是在分库后，表被分散到了不同的数据库，就无法进行关联查询了。

这时就需要将关联查询拆分成多次查询，然后将获得的结果进行拼装。

- 跨节点分页、排序函数

跨节点多库进行查询时，limit分页、order by排序等问题，就变得比较复杂了。需要先在不同的分片节点中将数据进行排序并返回，然后将不同分片返回的结果集进行汇总和再次排序。

这时非常容易出现内存崩溃的问题。

- 主键避重问题

在分库分表环境中，由于表中数据同时存在不同数据库中，主键值平时使用的自增长将无用武之地，某个分区数据库生成的ID无法保证全局唯一。因此需要单独设计全局主键，以避免跨库主键重复问题。

- 公共表处理

实际的应用场景中，参数表、数据字典表等都是数据量较小，变动少，而且属于高频联合查询的依赖表。这一类表一般就需要在每个数据库中都保存一份，并且所有对公共表的操作都要分发到所有的分库去执行。

- 运维工作量

面对散乱的分库分表之后的数据，应用开发工程师和数据库管理员对数据库的操作都变得非常繁重。对于每一次数据读写操作，他们都需要知道要往哪个具体的数据库的分表去操作，这也是其中重要的挑战之一。

4、什么时候需要分库分表？

在阿里巴巴公布的开发手册中，建议MySQL单表记录如果达到500W这个级别，或者单表容量达到2GB，一般就建议进行分库分表。而考虑到分库分表需要对数据进行再平衡，所以如果要使用分库分表，就要在系统设计之初就详细考虑好分库分表的方案，这里要分两种情况。

一般对于用户数据这一类后期增长比较缓慢的数据，一般可以按照三年左右的业务量来预估使用人数，按照标准预设好分库分表的方案。

而对于业务数据这一类增长快速且稳定的数据，一般则需要按照预估量的两倍左右预设分库分表方案。并且由于分库分表的后期扩容是非常麻烦的，所以在进行分库分表时，尽量根据情况，多分一些表。最好是计算一下数据增量，永远不用增加更多的表。

另外，在设计分库分表方案时，要尽量兼顾业务场景和数据分布。在支持业务场景的前提下，尽量保证数据能够分得更均匀。

最后，一旦用到了分库分表，就会表现为对数据查询业务的灵活性有一定的影响，例如如果按userId进行分片，那按age来进行查询，就必然会增加很多麻烦。如果再进行排序、分页、聚合等操作，很容易就扛不住了。这时候，都要尽量在分库分表的同时，再补充设计一个降级方案，例如将数据转存一份到ES，ES可以实现更灵活的大数据聚合查询。

5、常见的分库分表组件

由于分库分表之后，数据被分散在不同的数据库、服务器。因此，对数据的操作也就无法通过常规方式完成，并且它还带来了一系列的问题。好在，这些问题不是所有都需要我们在应用层面上解决，市面上有很多中间件可供我们选择，我们来了解一下它。

- shardingSphere 官网地址：<https://shardingsphere.apache.org/document/current/cn/overview/>

Sharding-JDBC是当当网研发的开源分布式数据库中间件，他是一套开源的分布式数据库中间件解决方案组成的生态圈，它由Sharding-JDBC、Sharding-Proxy和Sharding-Sidecar（计划中）这3款相互独立的产品组成。他们均提供标准化的数据分片、分布式事务和数据库治理功能，可适用于如Java同构、异构语言、容器、云原生等各种多样化的应用场景。

这也是我们这次学习的重点。

- mycat 官网地址: <http://www.mycat.org.cn/>

基于阿里开源的Cobar产品而研发, Cobar的稳定性、可靠性、优秀的架构和性能以及众多成熟的使用案例使得MYCAT一开始就拥有一个很好的起点, 站在巨人的肩膀上, 我们能看到更远。业界优秀的开源项目和创新思路被广泛融入到MYCAT的基因中, 使得MYCAT在很多方面都领先于目前其他一些同类的开源项目, 甚至超越某些商业产品。

MyCAT虽然是从阿里的技术体系中出来的, 但是跟阿里其实没什么关系。

- DBLE 官网地址: <https://opensource.actionsky.com/>

该网站包含几个重要产品。其中分布式中间件可以认为是MyCAT的一个增强版, 专注于MySQL的集群化管理。另外还有数据传输组件和分布式事务框架组件可供选择。

有道云笔记分享链接:

文档: VIP01-MySQL主从架构及读写分离实战.md

链接: <http://note.youdao.com/noteshare?id=df15b9a4c400815870b14a58c25dbb8f&sub=27F8AA04B3254C51A63932A208D6B224>