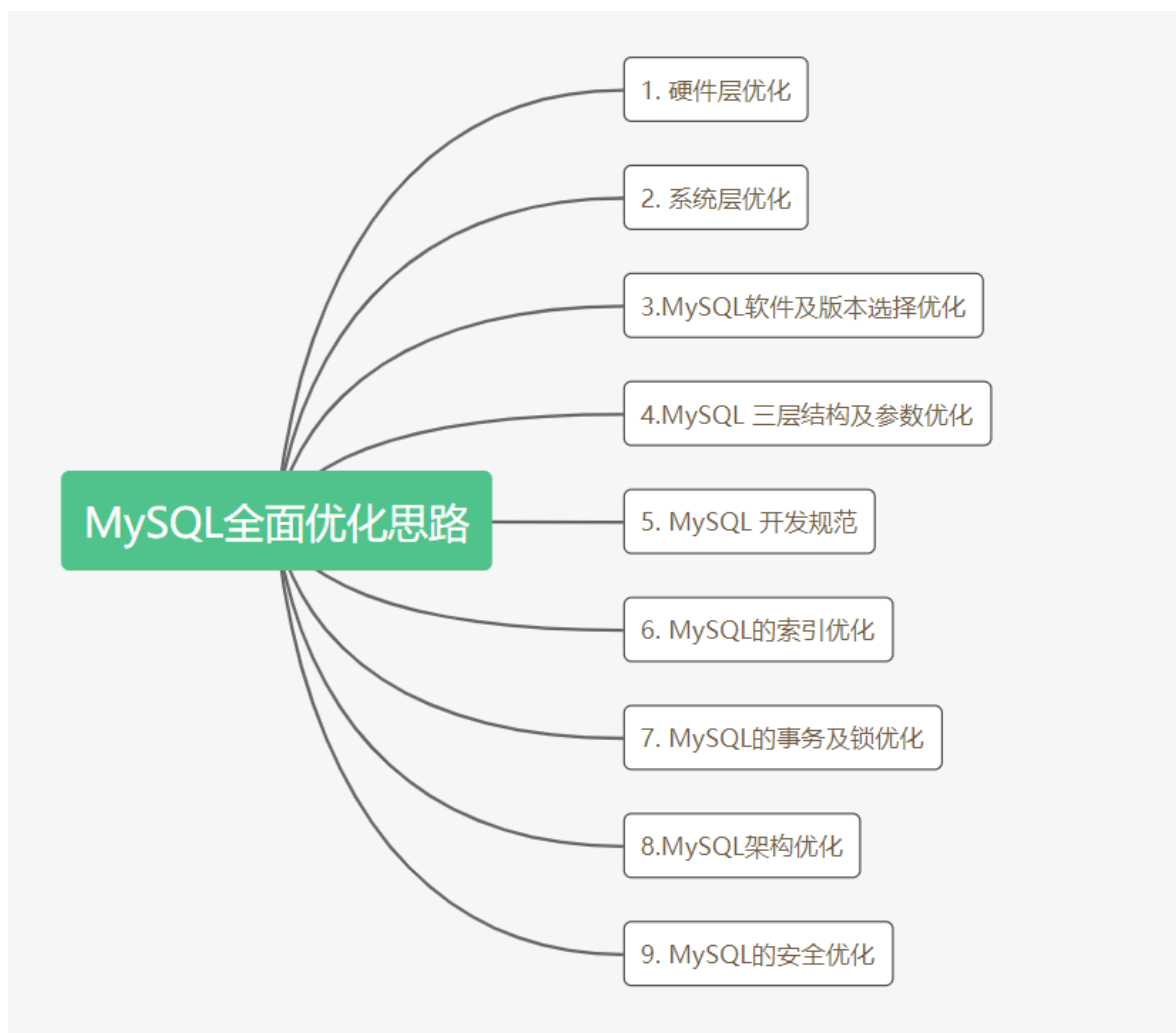


# MySQL全面优化

## 0.优化思路



## 1.硬件层面优化

### 1.0 硬件选配

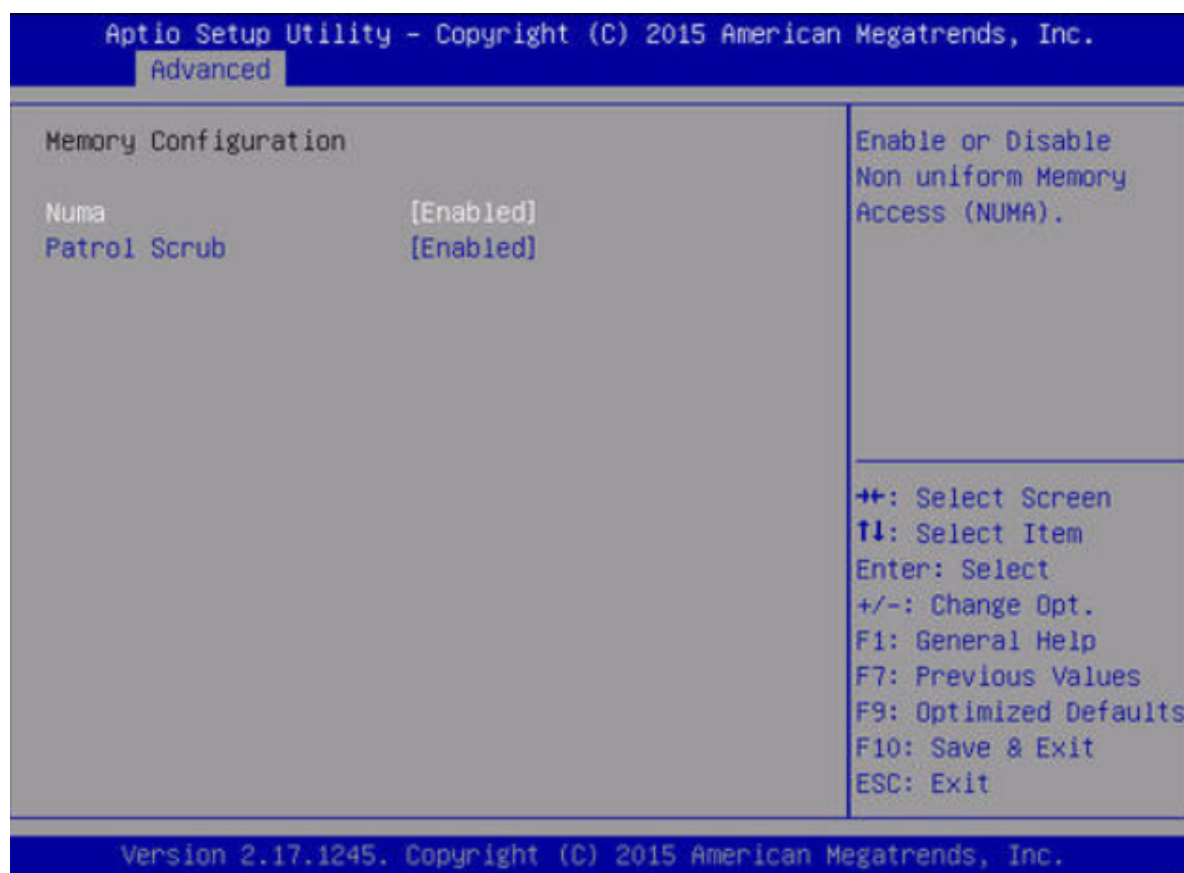
- 1 DELL、HP、IBM、华为、浪潮。
- 2 CPU: I、E
- 3 内存: ECC
- 4 IO : SAS 、 pci-e SSD 、 Nvme flash
- 5 raid卡: Raid10
- 6 网卡: 单卡单口
- 7
- 8 云服务器: ECS 、 RDS 、 PolarDB、TDSQL

### 1.1 关闭NUMA

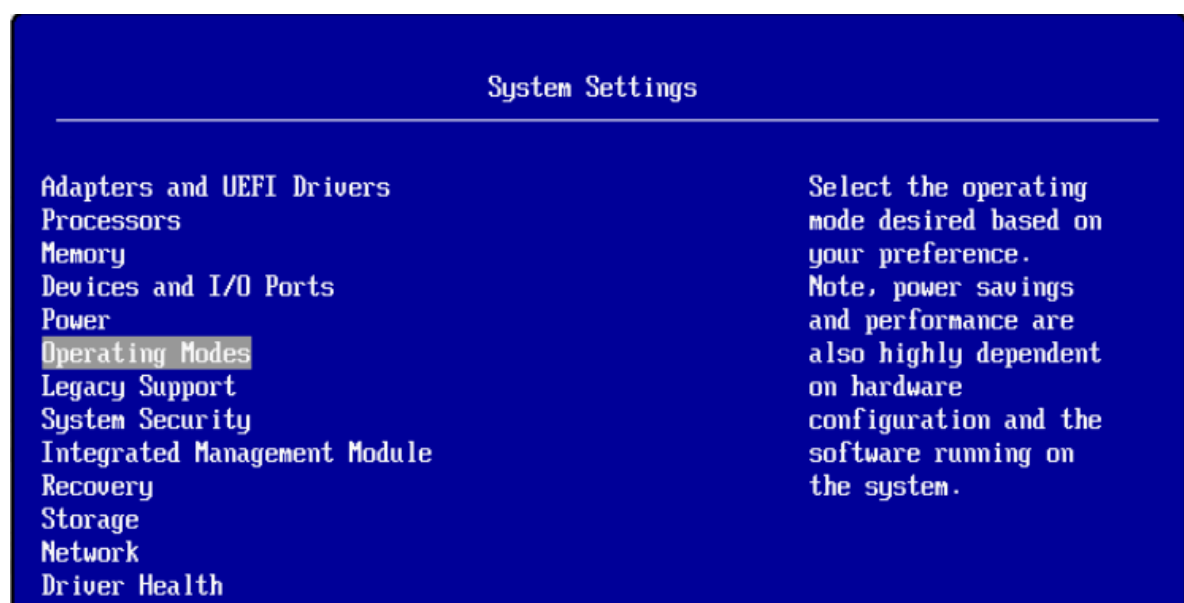
```

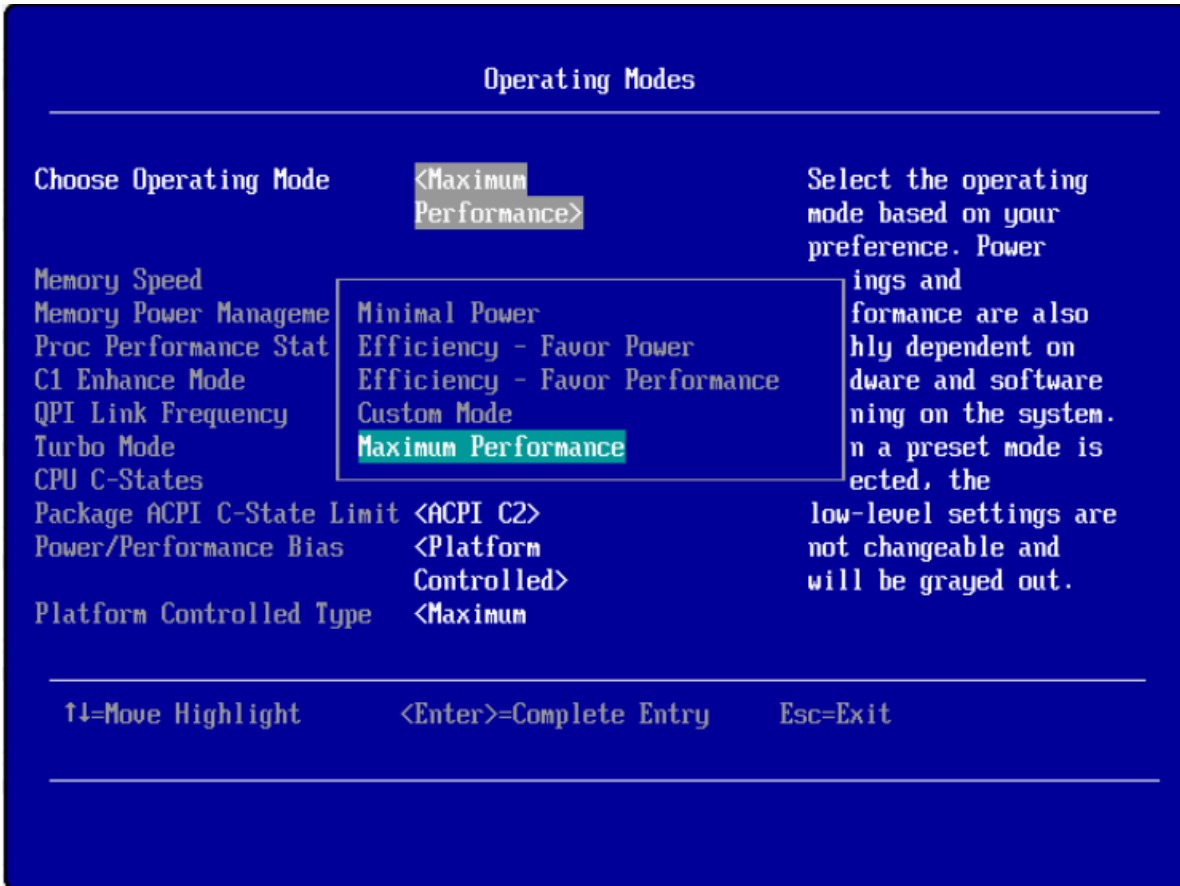
1  a. bios级别:
2  在bios层面numa关闭时, 无论os层面的numa是否打开, 都不会影响性能。
3
4  # numactl --hardware
5  available: 1 nodes (0)          #如果是2或多个nodes就说明numa没关掉
6
7  b. os grub级别:
8  vi /boot/grub/grub.conf
9  #/* Copyright 2010, Oracle. All rights reserved. */
10
11 default=0
12 timeout=5
13 hiddenmenu
14 foreground=000000
15 background=ffffff
16 splashimage=(hd0,0)/boot/grub/oracle.xpm.gz
17
18 title Trying_C0D0_as_HD0
19 root (hd0,0)
20 kernel /boot/vmlinuz-2.6.18-128.1.16.0.1.el5 root=LABEL=DBSYS ro
21 bootarea=dbsys rhgb quiet console=ttyS0,115200n8 console=tty1
22 crashkernel=128M@16M numa=off
23 initrd /boot/initrd-2.6.18-128.1.16.0.1.el5.img
24
25 在os层numa关闭时, 打开bios层的numa会影响性能, QPS会下降15-30%;
26
27 c. 数据库级别:
28
29 mysql> show variables like '%numa%';
30 +-----+-----+
31 | variable_name          | value |
32 +-----+-----+
33 | innodb_numa_interleave | OFF   |
34 +-----+-----+
35
36 或者:
37 vi /etc/init.d/mysqld
38 找到如下行
39 # Give extra arguments to mysqld with the my.cnf file. This script
40 # may be overwritten at next upgrade.
41 $bindir/mysqld_safe --datadir="$datadir" --pid-file="$mysqld_pid_file_path"
42 $other_args >/dev/null &
43
44 wait_for_pid created "$!" "$mysqld_pid_file_path"; return_value=$?
45 将$bindir/mysqld_safe --datadir="$datadir"这一行修改为:
46
47 /usr/bin/numactl --interleave all $bindir/mysqld_safe --datadir="$datadir"
48 --pid-file="$mysqld_pid_file_path" $other_args >/dev/null &
49 wait_for_pid created "$!" "$mysqld_pid_file_path"; return_value=$?
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99

```



## 1.2 开启CPU高性能模式





### 1.3 阵列卡配置建议

- 1 | raid10(推荐)
- 2 | SSD或者PCI-E或者Flash
- 3 | 强制回写 (Force writeBack)
- 4 | BBU 电池： 如果没电会有较大性能影响、定期充放电，如果UPS、多路电源、发电机。可以关闭。
- 5 | 关闭预读
- 6 | 有可能的话开启Cache(如果UPS、多路电源、发电机。)

■ Enter The Volume Attribute	
Volume Name	FA100 -VOL#000
Member Disks	4
Volume Raid Level	Raid 5
Max Capacity Allowed	750 GB
Select Volume Capacity	750 GB
Volume Initialization Mode	Foreground Initialization
Volume Stripe Size	64 KBytes
Volume Cache Mode	Write Back
Tagged Command Queuing	Enabled
Channel: LUN Base: LUN	0 : 0 : 0
Volumes To Be Created	1
<input type="checkbox"/> Confirm The Operation	
Submit Reset	

### 1.4 关闭THP

- 1 | vi /etc/rc.local

```

2  在文件末尾添加如下指令：
3  if test -f /sys/kernel/mm/transparent_hugepage/enabled; then
4      echo never > /sys/kernel/mm/transparent_hugepage/enabled
5  fi
6  if test -f /sys/kernel/mm/transparent_hugepage/defrag; then
7      echo never > /sys/kernel/mm/transparent_hugepage/defrag
8  fi
9
10 [root@master ~]# cat /sys/kernel/mm/transparent_hugepage/enabled
11 always madvise [never]
12 [root@master ~]# cat /sys/kernel/mm/transparent_hugepage/defrag
13 always madvise [never]
14

```

## 1.5 网卡绑定

- 1 bonding技术，业务数据库服务器都要配置bonding继续。建议是主备模式。
- 2 交换机一定要堆叠。

## 1.6 存储多路径

- 1 使用独立存储设备的话，需要配置多路径。
- 2 linux 自带 : multipath
- 3 厂商提供 :

# 2.系统层面优化

- 1 a. 更改文件句柄和进程数
- 2 内核优化 /etc/sysctl.conf
- 3 vm.swappiness <= 5 (也可以设置为0)
- 4 vm.dirty\_ratio <= 20
- 5 vm.dirty\_background\_ratio <= 10
- 6 net.ipv4.tcp\_max\_syn\_backlog = 819200
- 7 net.core.netdev\_max\_backlog = 400000
- 8 net.core.somaxconn = 4096
- 9 net.ipv4.tcp\_tw\_reuse=1
- 10 net.ipv4.tcp\_tw\_recycle=0
- 11
- 12 limits.conf
- 13 nofile 63000
- 14
- 15 b. 防火墙
- 16 禁用selinux : /etc/sysconfig/selinux 更改SELINUX=disabled.
- 17 iptables如果不使用可以关闭。可是需要打开MySQL需要的端口号
- 18
- 19 c. 文件系统优化
- 20 推荐使用XFS文件系统
- 21 MySQL数据分区独立 , 例如挂载点为: /data
- 22 mount参数 defaults, noatime, nodiratime, nobarrier 如/etc/fstab:
- 23 /dev/sdb /data xfs
- 24 defaults,noatime,nodiratime,nobarrier 1 2

```
25 | d. 不使用LVM
26 |
27 | e. io调度
28 | SAS : deadline
29 | SSD&PCI-E: noop
30 |
31 |
```

## 3. 数据库版本选择

- 1、稳定版：选择开源的社区版的稳定版GA版本。
- 2、选择mysql数据库GA版本发布后6个月-12个月的GA双数版本，大约在15-20个小版本左右。
- 3、要选择前后几个月没有大的BUG修复的版本，而不是大量修复BUG的集中版本。
- 4、要考虑开发人员开发程序使用的版本是否兼容你选的版本。
- 5、作为内部开发测试数据库环境，跑大概3-6个月的时间。
- 6、优先企业非核心业务采用新版本的数据库GA版本软件。
- 7、向DBA高手请教，或者在技术氛围好的群里和大家一起交流，使用真正的高手们用过的好用的GA版本产品。
- 8
- 9 最终建议： 8.0.20是一个不错的版本选择。向后可以选择双数版。

## 4.数据库三层结构及核心参数优化

### 4.1 连接层

```
1 | max_connections=单节点不高于3000
2 | max_connect_errors=大一点。
3 | wait_timeout=600
4 | interactive_wait_timeout=600
5 | net_read_timeout
6 | net_write_timeout
7 | max_allowed_packet
```

### 4.2 Server层

```
1 | sql_safe_updates           =1
2 | slow_query_log              =ON
3 | slow_query_log_file        =/xxx
4 | long_query_time             =1
5 | log_queries_not_using_indexes =ON
6 | log_throttle_queries_not_using_indexes = 10
7 | sort_buffer\join_buffer\read_buffer\read_rnd_buffer, 建议不超过8M
8 | tmp_table、heap_table,建议不要超过128M
9 | sql_mode,建议保持默认。
10 | max_execution_time, 建议跑批量是设置较大。
11 | lock_wait_timeout, 建议设置在60秒以内
12 | lower_case_table_names      =1
13 | thread_cache_size           =64
14 | character_set_server        =utf8或者utf8mb4
15 | log_timestamps               =SYSTEM
16 | init_connect                 ="set names utf8"
```

```

17 event_scheduler                =OFF
18 secure-file-priv                =/xxx
19 expire_logs_days                =10
20 sync_binlog                      =1
21 log-bin                         =/opt/log/mysql/blog/mysql-bin
22 log-bin-index                   =/opt/log/mysql/blog/mysql-bin.index
23 max_binlog_size                  =500M
24 binlog_format                   =ROW
25 max_binlog_cache_size            =2G
26 max_binlog_stmt_cache_size      =2G
27

```

## 4.3 存储引擎层

```

1 transaction-isolation           ="READ-COMMITTED"
2 innodb_data_home_dir            =/xxx
3 innodb_log_group_home_dir       =/xxx
4 innodb_log_file_size            =2048M
5 innodb_log_files_in_group       =3
6 innodb_flush_log_at_trx_commit  =2
7 innodb_flush_method             =O_DIRECT/fsync
8 innodb_io_capacity              =1000
9 innodb_io_capacity_max          =4000
10 innodb_buffer_pool_size         =64G
11 innodb_buffer_pool_instances   =4
12 innodb_log_buffer_size          =64M
13 innodb_max_dirty_pages_pct      =85
14 innodb_lock_wait_timeout        =10
15 innodb_open_files               =63000
16 innodb_page_cleaners            =4
17 innodb_sort_buffer_size         =64M
18 innodb_print_all_deadlocks      =1
19 innodb_rollback_on_timeout      =ON
20 innodb_deadlock_detect          =ON

```

## 4.4 复制

```

1 relay_log                       =/opt/log/mysql/blog/relay
2 relay_log_index                 =/opt/log/mysql/blog/relay.index
3 max_relay_log_size              =500M
4 relay_log_purge                 =ON
5 relay_log_recovery              =ON
6 rpl_semi_sync_master_enabled    =ON
7 rpl_semi_sync_master_timeout    =1000
8 rpl_semi_sync_master_trace_level =32
9 rpl_semi_sync_master_wait_for_slave_count =1
10 rpl_semi_sync_master_wait_no_slave =ON
11 rpl_semi_sync_master_wait_point =AFTER_SYNC
12 rpl_semi_sync_slave_enabled     =ON
13 rpl_semi_sync_slave_trace_level =32
14 binlog_group_commit_sync_delay   =1
15 binlog_group_commit_sync_no_delay_count =1000
16 gtid_mode                       =ON

```

```

17 enforce_gtid_consistency      =ON
18 master_verify_checksum       =ON
19 sync_master_info              =1
20 skip_slave_start              =1
21 #read_only                    =ON
22 #super_read_only              =ON
23 log_slave_updates             =ON
24 server_id                     =2330602
25 report_host                   =xxxx
26 report_port                   =3306
27 slave_parallel_type           =LOGICAL_CLOCK
28 slave_parallel_workers        =4
29 master_info_repository       =TABLE
30 relay_log_info_repository     =TABLE

```

## 4.5 其它

```

1 客户端配置：
2 [mysql]
3 no-auto-rehash
4 pager less

```

# 5.开发规范

## 5.1 字段规范

- 1 1. 每个表建议在30个字段以内。
- 2 2. 需要存储emoji字符的，则选择utf8mb4字符集。
- 3 3. 机密数据，加密后存储。
- 4 4. 整型数据，默认加上UNSIGNED。
- 5 5. 存储IPV4地址建议用INT UNSIGNED，查询时再利用INET\_ATON()、INET\_NTOA()函数转换。
- 6 6. 如果遇到BLOB、TEXT大字段单独存储表或者附件形式存储。
- 7 7. 选择尽可能小的数据类型，用于节省磁盘和内存空间。
- 8 8. 存储浮点数，可以放大倍数存储。
- 9 9. 每个表必须有主键，INT/BIGINT并且自增做为主键，分布式架构使用sequence序列生成器保存。
- 10 10. 每个列使用not null，或增加默认值。

## 5.2 SQL语句规范

- 1 ### 1. 去掉不必要的括号
- 2 如： ((a AND b) AND c OR (((a AND b) AND (c AND d))))
- 3 修改成 (a AND b AND c) OR (a AND b AND c AND d)
- 4 ### 2. 去掉重叠条件
- 5 如： (a<b AND b=c) AND a=5
- 6 修改成 b>5 AND b=c AND a=5
- 7 如： (B>=5 AND B=5) OR (B=6 AND 5=5) OR (B=7 AND 5=6)
- 8 修改成 B=5 OR B=6
- 9 ### 3. 避免使用not in、not exists、<>、like %%
- 10 ### 4. 多表连接，小表驱动大表
- 11 ### 5. 减少临时表应用，优化order by、group by、union、distinct、join等
- 12 ### 6. 减少语句查询范围，精确查询条件
- 13 ### 7. 多条件，符合联合索引最左原则



```
14  ### 8. 查询条件减少使用函数、拼接字符等条件、条件隐式转换
15  ### 9. union all 替代 union
16  ### 10.减少having子句使用
17  ### 11.如非必须不使用 for update语句
18  ### 12.update和delete, 开启安全更新参数
19  ### 13.减少inset ... select语句应用
20  ### 14.使用load 替代insert录入大数据
21  ### 15.导入大量数据时, 可以禁用索引、增大缓冲区、增大redo文件和buffer、关闭
    autocommit、RC级别可以提高效率
22  ### 16.优化limit, 最好业务逻辑中先获取主键ID, 再基于ID进行查询
23      limit 5000000,10
24  ### 17. DDL执行前要审核
25  ### 18. 多表连接语句执行前要看执行计划
26
```

## 6.索引优化

```
1  1. 非唯一索引按照“i_字段名称_字段名称[_字段名]”进行命名。
2  2. 唯一索引按照“u_字段名称_字段名称[_字段名]”进行命名。
3  3. 索引名称使用小写。
4  4. 索引中的字段数不超过5个。
5  5. 唯一键由3个以下字段组成, 并且字段都是整形时, 使用唯一键作为主键。
6  6. 没有唯一键或者唯一键不符合5中的条件时, 使用自增id作为主键。
7  7. 唯一键不和主键重复。
8  8. 索引选择度高的列作为联合索引最左条件
9  9. ORDER BY, GROUP BY, DISTINCT的字段需要添加在索引的后面。
10 10. 单张表的索引数量控制在5个以内, 若单张表多个字段在查询需求上都要单独用到索引, 需要经过
    DBA评估。查询性能问题无法解决的, 应从产品设计上进行重构。
11 11. 使用EXPLAIN判断SQL语句是否合理使用索引, 尽量避免extra列出现: Using File Sort,
    Using Temporary。
12 12. UPDATE、DELETE语句需要根据WHERE条件添加索引。
13 13. 对长度大于50的VARCHAR字段建立索引时, 按需求恰当的使用前缀索引, 或使用其他方法。
14 14. 下面的表增加一列url_crc32, 然后对url_crc32建立索引, 减少索引字段的长度, 提高效率。
15 CREATE TABLE all_url(ID INT UNSIGNED NOT NULL PRIMARY KEY AUTO_INCREMENT,
16 url VARCHAR(255) NOT NULL DEFAULT 0,
17 url_crc32 INT UNSIGNED NOT NULL DEFAULT 0,
18 index idx_url(url_crc32));
19 15. 合理创建联合索引(避免冗余), (a,b,c) 相当于 (a)、(a,b)、(a,b,c)。
20 16. 合理利用覆盖索引, 减少回表。
21 17. 减少冗余索引和使用率较低的索引
```

## 7.锁优化

### 7.1 latch 闕锁

#### a. 介绍

```
1  Latch用于管理对共享内存资源的并发访问, 例如, 操作缓冲池汇总的LRU列表, 删除、添加、移动
    LRU列表中的元素, 为了保证一致性, 必须有锁的介入, 这就是latch锁。
```

#### b.latch 和 lock的区别

	lock	latch
对象	SQL操作流程	线程
保护	数据库对象（库、表、行、索引、表空间、数据页等）	所有共享内存数据结构
生命周期	整个操作周期	临界资源
锁定模式	MDL、table、Record、Gap、Nextlock、意向	rw-latch、mutex

## c.查看latch争用的类型

```

1  mysql> show engine innodb mutex;
2  +-----+-----+-----+-----+
3  | Type   | Name                               | Status   |
4  +-----+-----+-----+-----+
5  | InnoDB | rwlock: dict0dict.cc:2687         | waits=1  |
6  | InnoDB | rwlock: dict0dict.cc:1184         | waits=13 |
7  | InnoDB | rwlock: log0log.cc:844            | waits=35 |
8  | InnoDB | sum rwlock: buf0buf.cc:1457       | waits=4  |
9  +-----+-----+-----+-----+
10
11  可以在latch争用较为严重情况下，定位到源码的位置点，从而获得到底什么原因导致争用。
12
13  也可以在此时通过以下工具分析堆栈信息：
14  pstack -p `pidof mysqld` >/tmp/aa.txt
15  pt-pmp /tmp/aa.txt|more

```

## d.什么时候发生争用

```

1  1) a 访问x内存链表
2  2) b 排队等待x解锁，占了cpu，但是cpu发现你在等待，所以cpu将b踢出
3  3) 访问锁链的时间，就是找数据的时间。
4  4) b知道很快所以，b不去排队，这时去spin也就是空转cpu，然后再去看一下内存数据结构，a是否已
   解锁
5  5) b转了一圈后，在b spin的时间段的时间中，c进来了，连续多次的spin后，产生了os waits
6  6) 操作系统将b从cpu中踢出
7
8  latch争用的表面现象：latch争用会表现为cpu繁忙，IO很闲，没有做实际的事情。
9

```

## e.如何监控是否latch争用较为严重

```

1  -----
2  SEMAPHORES
3  -----
4  OS WAIT ARRAY INFO: reservation count 13
5  OS WAIT ARRAY INFO: signal count 13
6  RW-shared spins 0, rounds 0, OS waits 0
7  RW-excl spins 2, rounds 60, OS waits 2
8  RW-sx spins 2, rounds 60, OS waits 2
9  Spin rounds per wait: 0.00 RW-shared, 30.00 RW-excl, 30.00 RW-sx
10
11  rounds: 意思是每次询问旋转的次数
12  os waits: 表示sleep，当突然增长比较快的时候，说明latch争用比较严重

```

```
13 | rw-shared spin 的次数
14 | rw-excl   spin的次数
```

## f. latch争用发生的原因

```
1 | 1、内存访问太频繁（不停地找）
2 | 2、list链太长（链上挂10000个块，被持有的几率太大）
3 |
```

## g.如何降低latch争用

```
1 | 如果出现latch争用比较严重
2 | 1.优化大sql，降低对内存读的数量—效果比较明显
3 | 2.增加instances的数量
```

## 7.2 全局锁 Global Read lock

### a. 介绍

```
1 | 全局读锁。
2 | 加锁方法： FTWRL, flush tables with read lock.
3 | 解锁方法： unlock tables;
4 | 出现场景：
5 |     mysqldump --master-data
6 |     xtrabackup（8.0之前早期版本）等备份时。
7 | 属于类型： MDL（metadata lock）层面锁
8 | 影响情况： 加锁期间，阻塞所有事务写入，阻塞所有已有事务commit。
9 | MDL，等待时间受 lock_wait_timeout=31536000
```

### b. 检测方法

```
1 | UPDATE performance_schema.setup_instruments
2 | SET ENABLED = 'YES', TIMED = 'YES'
3 | WHERE NAME = 'wait/lock/metadata/sql/mdl';
4 |
5 | mysql> select * from performance_schema.metadata_locks;
6 |
7 | mysql> select OBJECT_SCHEMA ,OBJECT_NAME
8 | ,LOCK_TYPE,LOCK_DURATION,LOCK_STATUS ,OWNER_THREAD_ID,OWNER_EVENT_ID from
9 | performance_schema.metadata_locks;
10 |
11 | mysql> show processlist;
12 |
13 | mysql> select * from sys.schema_table_lock_waits;
```

### c. 一个经典故障：5.7 xtrabackup/mysqldump备份时数据库出现hang状态，所有查询都不能进行

```
1 | session1: 模拟一个大的查询或事务
2 | mysql> select *,sleep(100) from city where id<10 limit 1 ;
3 |
```

```

4 session2: 模拟备份时的FTWRL
5 mysql> flush tables with read lock;
6 -- 此时发现命令被阻塞
7
8 session3: 发起正常查询, 发现被阻塞
9 mysql> select * from world.city where id=1;
10
11 结论: 备份时, 一定要选择业务不繁忙期间, 否则有可能会阻塞正常业务。
12
13
14 案例2:
15 5.7 innobackupex备份全库, 进程死了, mysql里就是全库读锁, 后边insert 全阻塞了
16

```

## 7.3 Table lock

### a. 介绍

```

1 表锁。
2 加锁方式:
3     lock table read. 所有会话只读。属于MDL锁。
4     lock table write. 当前持有会话可以RW, 其他会被阻塞。属于MDL锁
5     select for update ;
6     select for share ;
7 解锁方式:
8 unlock tables;

```

### b. 检测方式

```

1 [mysqld]
2 performance-schema-instrument='wait/lock/metadata/sql/mdl=ON'
3 mysql> select * from performance_schema.metadata_locks;
4

```

## 7.4 MDL锁

### a. 介绍

```

1 Metadata lock .元数据锁。
2 作用范围: global 、 commit、 tablespace、 schema、 table等
3 默认timeout时间: lock_wait_timeout
4 mysql> select @@lock_wait_timeout;
5 +-----+
6 | @@lock_wait_timeout |
7 +-----+
8 |           31536000 |
9 +-----+

```

### b. 监控

```

1  [mysqld]
2  performance-schema-instrument='wait/lock/metadata/sql/mdl=ON'
3  mysql> select * from performance_schema.metadata_locks;
4
5  找到
6  OWNER_THREAD_ID: 62
7  mysql> select * from threads where thread_id='62'\G
8
9  PROCESSLIST_ID: 21
10 kill 21;
11
12

```

## 7.5 autoinc\_lock

```

1  自增锁。
2  通过参数: innodb_autoinc_lock_mode=0,1,2设定
3  0: 表锁, 每次插入都请求表锁, 效率低。
4  1: mutex方式, 预计插入多少行, 预申请自增序列。如果出现load或者insert select方式会退化为0。
5  2: 强制使用mutex方式。并发插入可以更高效。
6
7  作用:
8  The default innodb_autoinc_lock_mode setting is now 2 (interleaved).
   Interleaved lock mode permits the execution of multi-row inserts in parallel,
   which improves concurrency and scalability.

```

## 7.6 row lock

### a. 介绍

```

1  record lock 、gap、next lock
2  都是基于索引加锁, 与事务隔离级别有关。

```

### b. 监控及分析

```

1  show status like 'innodb_row_lock%'
2  select * from information_schema.innodb_trx;
3  select * from sys.innodb_lock_waits;
4  select * from performance_schema.threads;
5  select * from performance_schema.events_statements_current;
6  select * from performance_schema.events_statements_history;

```

### c. 优化方向

```

1 | 1. 优化索引
2 | 2. 减少事务的更新范围
3 | 3. RC
4 | 4. 拆分语句:
5 | 例如: update t1 set num=num+10 where k1 <100; k1 是辅助索引,record lock gap
   | next
6 |       改为:
7 |       select id from t1 where k1 <100; ---> id: 20,30,50
8 |       update t1 set num=num+10 where id in (20,30,50);

```

## 7.7 死锁

### a.介绍

```

1 | dead lock 多个并发事务之间发生交叉资源依赖时，会出现。

```

### b. 监控及分析

```

1 | show engine innodb status \G
2 | innodb_print_all_deadlocks =1

```

### c. 经典死锁案例延时及解析

```

mysql> select * from t1;
+-----+-----+
| id | k1 |
+-----+-----+
| 1 | a |
| 2 | b |
| 3 | c |
| 4 | d |
| 5 | e |
| 6 | f |
+-----+-----+
6 rows in set (0.00 sec)

mysql> █

```

```

mysql> desc t1;
+-----+-----+-----+-----+-----+-----+
| Field | Type          | Null | Key | Default | Extra |
+-----+-----+-----+-----+-----+-----+
| id    | int           | NO   | PRI | NULL    |       |
| k1    | varchar(20)   | YES  |     | NULL    |       |
+-----+-----+-----+-----+-----+-----+
2 rows in set (0.00 sec)

```

```
mysql> use test;
Database changed
mysql> begin;
Query OK, 0 rows affected (0.00 sec)

mysql> delete from t1 where id=6;
Query OK, 1 row affected (0.00 sec)

mysql> delete from t1 where id=5;
Query OK, 1 row affected (0.00 sec)

mysql> delete from t1 where id=4;
Query OK, 1 row affected (0.00 sec)

mysql> delete from t1 where id=3;
Query OK, 1 row affected (3.46 sec)

mysql>

mysql> use test;
Database changed
mysql> begin;
Query OK, 0 rows affected (0.00 sec)

mysql> delete from t1 where id=1;
Query OK, 1 row affected (0.00 sec)

mysql> delete from t1 where id=2;
Query OK, 1 row affected (0.00 sec)

mysql> delete from t1 where id=3;
Query OK, 1 row affected (0.00 sec)

mysql> delete from t1 where id=4;
ERROR 1213 (40001): Deadlock found when trying to get lock; try restarting transaction
mysql>
```

## 8.架构优化

- 1
- 2 高可用架构:
- 3     MHA+ProxySQL+GTID
- 4     MGR\InnoDB Cluster
- 5     PXC
- 6 读写分离:
- 7     ProxySQL、MySQL-router
- 8 NoSQL:
- 9     Redis+sentinel, Redis Cluster
- 10    MongoDB RS/MongoDB SHARDING Cluster
- 11    ES

## 9.安全优化

- 1 1、 使用普通nologin用户管理MySQL
- 2 2、 合理授权用户、密码复杂度及最小权限、系统表保证只有管理员用户可访问。
- 3 3、 删除数据库匿名用户
- 4 4、 锁定非活动用户
- 5 5、 MySQL尽量不暴露互联网,需要暴露互联网用户需要设置明确白名单、替换MySQL默认端口号、使用ssl连接
- 6 6、 优化业务代码,防止SQL注入。

## 10.常用工具介绍

一、 PT (percona-toolkits) 工具的应用:

### 1. pt工具安装

```
[root@master ~]# yum install -y percona-toolkit-3.1.0-2.el7.x86_64.rpm
```

### 2. 常用工具使用介绍

#### 2.1 pt-archiver 归档表

场景:

面试题: 亿级的大表, delete批量删除100w左右数据。

面试题: 定期按照时间范围, 进行归档表。

# 重要参数

--limit 100      每次取100行数据用pt-archive处理  
--txn-size 100    设置100行为一个事务提交一次,  
--where 'id<3000' 设置操作条件  
--progress 5000   每处理5000行输出一次处理信息  
--statistics      输出执行过程及最后的操作统计。(只要不加上--quiet, 默认情况下pt-archive都会输出执行过程的)  
--charset=UTF8    指定字符集为UTF8—这个最后加上不然可能出现乱码。  
--bulk-delete     批量删除source上的旧数据(例如每次1000行的批量删除操作)

注意: 需要归档表中至少有一个索引,做好是where条件列有索引

使用案例:

## 1.归档到数据库

```
db01 [test]>create table test1 like t100w;  
pt-archiver --source h=10.0.0.51,D=test,t=t100w,u=oldguo,p=123 --dest  
h=10.0.0.51,D=test,t=test1,u=oldguo,p=123 --where 'id<10000' --no-check-charset --no-delete --  
limit=1000 --commit-each --progress 1000 --statistics
```

## 2.只清理数据

```
pt-archiver --source h=10.0.0.51,D=test,t=t100w,u=oldguo,p=123 --where 'id<10000' --purge --  
limit=1 --no-check-charset
```

## 3.只把数据导出到外部文件, 但是不删除源表里的数据

```
pt-archiver --source h=10.0.0.51,D=world,t=city,u=root,p=123 --where '1=1' --no-check-charset --  
no-delete --file="/tmp/archiver.dat"
```

## 2.2 pt-osc

场景:

修改表结构、索引创建删除  
不能加快速度, 但能减少业务影响(锁)。

面试题:

pt-osc工作流程:

- 1、检查更改表是否有主键或唯一索引, 是否有触发器
- 2、检查修改表的表结构, 创建一个临时表, 在新表上执行ALTER TABLE语句  
create table bak like t1;  
alter table bak add telnum char(11) not null;
- 3、在源表上创建三个触发器分别对于INSERT UPDATE DELETE操作  
create trigger  
a  
b  
c
- 4、从源表拷贝数据到临时表, 在拷贝过程中, 对源表的更新操作会写入到新建表中  
insert into bak select \* from t1
- 5、将临时表和源表rename (需要元数据修改锁, 需要短时间锁表)
- 6、删除源表和触发器, 完成表结构的修改。

pt-osc工具限制

- 1、源表必须有主键或唯一索引, 如果没有工具将停止工作
- 2、如果线上的复制环境过滤器操作过于复杂, 工具将无法工作
- 3、如果开启复制延迟检查, 但主从延迟时, 工具将暂停数据拷贝工作



- 4、如果开启主服务器负载检查，但主服务器负载较高时，工具将暂停操作
- 5、当表使用外键时，如果未使用--alter-foreign-keys-method参数，工具将无法执行
- 6、只支持InnoDB存储引擎表，且要求服务器上有该表1倍以上的空闲空间。

pt-osc之alter语句限制

- 1、不需要包含alter table关键字，可以包含多个修改操作，使用逗号分开，如"drop column c1, add column c2 int"
- 2、不支持rename语句来对表进行重命名操作
- 3、不支持对索引进行重命名操作
- 4、如果删除外键，需要对外键名加下划线，如删除外键fk\_uid, 修改语句为"DROP FOREIGN KEY \_fk\_uid"

pt-osc之命令模板

## --execute表示执行

---

## --dry-run表示只进行模拟测试

---

## 表名只能使用参数t来设置，没有长参数

---

```
pt-online-schema-change \  
--host="127.0.0.1" \  
--port=3358 \  
--user="root" \  
--password="root@root" \  
--charset="utf8" \  
--max-lag=10 \  
--check-salve-lag='xxx.xxx.xxx.xxx' \  
--recursion-method="hosts" \  
--check-interval=2 \  
--database="testdb1" \  
t="tb001" \  
--alter="add column c4 int" \  
--execute
```

例子：

```
pt-online-schema-change --user=oldguo --password=123 --host=10.0.0.51 --alter "add column  
state int not null default 1" D=test,t=t100w --print --execute  
pt-online-schema-change --user=oldguo --password=123 --host=10.0.0.51 --alter "add index  
idx(num)" D=test,t=t100w --print --execute
```

### 2.3 pt-table-checksum

场景： 校验主从数据一致性

#### 2.3.1 创建数据库

Create database pt CHARACTER SET utf8;

创建用户checksum并授权

```
GRANT ALL ON . TO 'checksum'@'10.0.0.0%' IDENTIFIED BY 'checksum';  
flush privileges;
```

### 2.3.2 参数:

--[no]check-replication-filters: 是否检查复制的过滤器, 默认是yes, 建议启用不检查模式。  
--databases | -d: 指定需要被检查的数据库, 多个库之间可以用逗号分隔。  
--[no]check-binlog-format: 是否检查binlog文件的格式, 默认值yes。建议开启不检查。因为在默认的row格式下会出错。  
--replicate: 把checksum的信息写入到指定表中。  
--replicate-check-only: 只显示不同步信息

```
pt-table-checksum --nocheck-replication-filters --no-check-binlog-format --replicate=pt.checksums  
--create-replicate-table --databases=test --tables=t1  
h=10.0.0.51,u=checksum,p=checksum,P=3306
```

```
#!/bin/bash
```

```
date >> /root/db/checksum.log  
pt-table-checksum --nocheck-binlog-format --nocheck-plan --nocheck-replication-filters --  
replicate=pt.checksums --set-vars innodb_lock_wait_timeout=120 --databases test --tables t1 -  
u'checksum' -p'checksum' -h'10.0.0.51' >> /tmp/checksum.log  
date >> /root/db/checksum.log
```

### 2.4 pt-table-sync

#### 主要参数介绍

--replicate: 指定通过pt-table-checksum得到的表。  
--databases: 指定执行同步的数据库。  
--tables: 指定执行同步的表, 多个用逗号隔开。  
--sync-to-master: 指定一个DSN, 即从的IP, 他会通过show processlist或show slave status 去自动的找主。  
h: 服务器地址, 命令里有2个ip, 第一次出现的是Master的地址, 第2次是Slave的地址。  
u: 帐号。  
p: 密码。  
--print: 打印, 但不执行命令。  
--execute: 执行命令。

```
pt-table-sync --replicate=pt.checksums --databases test --tables t1  
h=10.0.0.51,u=checksum,p=checksum,P=3306 h=10.0.0.52,u=checksum,p=checksum,P=3306 --  
print
```

```
pt-table-sync --replicate=pt.checksums --databases test --tables t1  
h=10.0.0.51,u=checksum,p=checksum,P=3306 h=10.0.0.52,u=checksum,p=checksum,P=3306 --  
execute
```

### 2.5 pt-duplicate-key-checker

作用: 检查数据库重复索引

```
pt-duplicate-key-checker --database=test h='10.0.0.51' --user=oldguo --password=123
```

### 2.6 pt-kill 语句

场景: 无法正常kill的连接。

#### 常用参数说明

--daemonize 放在后台以守护进程的形式运行;  
--interval 多久运行一次, 单位可以是s,m,h, d等默认是s -不加这个默认是5秒  
--victims 默认是oldest,只杀最古老的查询。这是防止被查杀是不是真的长时间运行的查询, 他们只是长期等待 这种种匹配按时间查询, 杀死一个时间最高值。  
--all 杀掉所有满足的线程  
--kill-query 只杀掉连接执行的语句, 但是线程不会被终止  
--print 打印满足条件的语句  
--busy-time 批次查询已运行的时间超过这个时间的线程;

--idle-time 杀掉sleep 空闲了多少时间的连接线程，必须在--match-command sleep时才有效—也就是匹配使用 -- -match-command 匹配相关的语句。

---ignore-command 忽略相关的匹配。这两个搭配使用一定是ignore-command在前 match-command在后，

--match-db cdelzone 匹配哪个库

command有：Query、Sleep、Binlog Dump、Connect、Delayed insert、Execute、Fetch、Init DB、Kill、Prepare、Processlist、Quit、Reset stmt、Table Dump

例子：

## 杀掉空闲链接sleep 5秒的 SQL 并把日志放到/home/pt-kill.log文件中

```
/usr/bin/pt-kill --user=用户名 --password=密码 --match-command Sleep --idle-time 5 --victim all --interval 5 --kill --daemonize -S /tmp/mysql.sock --pid=/tmp/ptkill.pid --print --log=/tmp/pt-kill.log &
```

## 查询SELECT 超过1分钟

```
/usr/bin/pt-kill --user=用户名 --password=密码 --busy-time 60 --match-info "SELECT|select" --victim all --interval 5 --kill --daemonize -S -S /tmp/mysql.sock --pid=/tmp/ptkill.pid --print --log=/tmp/pt-kill.log &
```

## Kill掉 select IFNULL.\*语句开头的SQL

```
pt-kill --user=用户名 --password=密码 --victims all --busy-time=0 --match-info="select IFNULL.*" --interval 1 -S /tmp/mysqld.sock --kill --daemonize --pid=/tmp/ptkill.pid --print --log=/tmp/pt-kill.log &
```

## kill掉state Locked

```
/usr/bin/pt-kill --user=用户名 --password=密码 --victims all --match-state='Locked' --victim all --interval 5 --kill --daemonize -S /tmp/mysqld.sock --pid=/tmp/ptkill.pid --print --log=/tmp/pt-kill.log &
```

## kill掉 a库， web为10.0.0.11的连接

```
pt-kill --user=用户名 --password=密码 --victims all --match-db='a' --match-host='10.0.0.11' --kill --daemonize --interval 10 -S /tmp/mysqld.sock --pid=/tmp/ptkill.pid --print --log=/tmp/pt-kill.log &
```

## 指定哪个用户kill

```
pt-kill --user=用户名 --password=密码 --victims all --match-user='root' --kill --daemonize --interval 10 -S /home/zb/data/my6006/socket/mysqld.sock --pid=/tmp/ptkill.pid --print --log=/home/pt-kill.log &
```

## kill掉 command query | Execute

```
pt-kill --user=用户名 --password=密码 --victims all --match-command="query|Execute" --interval 5 --kill --daemonize -S /tmp/mysqld.sock --pid=/tmp/ptkill.pid --print --log=/home/pt-kill.log &
```

7. 显示主从结构：pt-slave-find

```
[root@db01 tmp]# pt-slave-find -h10.0.0.51 -P3306 -uchecksum -pchecksum
10.0.0.51
Version      5.7.28-log
Server ID    51
```

```
Uptime      27:57 (started 2020-05-15T13:24:15)
Replication  Is not a slave, has 1 slaves connected, is not read_only
Filters
Binary logging ROW
Slave status
Slave mode   STRICT
Auto-increment increment 1, offset 1
InnoDB version 5.7.28
+- 10.0.0.52
Version      5.7.28-log
Server ID    52
Uptime      28:18 (started 2020-05-15T13:23:54)
Replication  Is a slave, has 0 slaves connected, is not read_only
Filters
Binary logging ROW
Slave status 0 seconds behind, running, no errors
Slave mode   STRICT
Auto-increment increment 1, offset 1
InnoDB version 5.7.28
[root@db01 tmp]#
```

## 8. 监控主从延时

# pt-heartbeat

---

主库:

```
pt-heartbeat --user=oldguo --ask-pass --host=10.0.0.51 --port=3306 --create-table -D test --
interval=1 --update --replace --daemonize
```

从库:

```
pt-heartbeat --user=oldguo --ask-pass --host=10.0.0.52 --port=3306 -D test --table=heartbeat --
monitor
```

## 9. pt-show-grants

---

作用: 用户和权限信息迁移。

```
pt-show-grants -h10.0.0.51 -P3306 -uchecksum -pchecksum
```

```
-- Grants dumped by pt-show-grants
-- Dumped from server 10.0.0.51 via TCP/IP, MySQL 5.7.28-log at 2020-05-15 17:11:06
-- Grants for 'checksum'@'10.0.0.%'
CREATE USER IF NOT EXISTS 'checksum'@'10.0.0.%';
ALTER USER 'checksum'@'10.0.0.%' IDENTIFIED WITH 'mysql_native_password' AS
'*E5E390AF1BDF241B51D9C0DBBEA262CC9407A2DF' REQUIRE NONE PASSWORD EXPIRE
DEFAULT ACCOUNT UNLOCK;
GRANT ALL PRIVILEGES ON . TO 'checksum'@'10.0.0.%';
-- Grants for 'mysql.session'@'localhost'
CREATE USER IF NOT EXISTS 'mysql.session'@'localhost';
ALTER USER 'mysql.session'@'localhost' IDENTIFIED WITH 'mysql_native_password' AS
```

```

'THISISNOTAVALIDPASSWORDTHATCANBEUSEDHERE' REQUIRE NONE PASSWORD EXPIRE DEFAULT
ACCOUNT LOCK;
GRANT SELECT ON mysql.user TO 'mysql.session'@'localhost';
GRANT SELECT ON performance_schema TO 'mysql.session'@'localhost';
GRANT SUPER ON . TO 'mysql.session'@'localhost';
-- Grants for 'mysql.sys'@'localhost'
CREATE USER IF NOT EXISTS 'mysql.sys'@'localhost';
ALTER USER 'mysql.sys'@'localhost' IDENTIFIED WITH 'mysql_native_password' AS
'THISISNOTAVALIDPASSWORDTHATCANBEUSEDHERE' REQUIRE NONE PASSWORD EXPIRE DEFAULT
ACCOUNT LOCK;
GRANT SELECT ON sys.sys_config TO 'mysql.sys'@'localhost';
GRANT TRIGGER ON sys TO 'mysql.sys'@'localhost';
GRANT USAGE ON . TO 'mysql.sys'@'localhost';
-- Grants for 'repl'@'10.0.0.%'
CREATE USER IF NOT EXISTS 'repl'@'10.0.0.%';
ALTER USER 'repl'@'10.0.0.%' IDENTIFIED WITH 'mysql_native_password' AS
'*23AE809DDACAF96AF0FD78ED04B6A265E05AA257' REQUIRE NONE PASSWORD EXPIRE
DEFAULT ACCOUNT UNLOCK;
GRANT REPLICATION SLAVE ON . TO 'repl'@'10.0.0.%';
-- Grants for 'root'@'10.0.0.%'
CREATE USER IF NOT EXISTS 'root'@'10.0.0.%';
ALTER USER 'root'@'10.0.0.%' IDENTIFIED WITH 'mysql_native_password' AS
'*23AE809DDACAF96AF0FD78ED04B6A265E05AA257' REQUIRE NONE PASSWORD EXPIRE
DEFAULT ACCOUNT UNLOCK;
GRANT ALL PRIVILEGES ON . TO 'root'@'10.0.0.%';
-- Grants for 'root'@'localhost'
CREATE USER IF NOT EXISTS 'root'@'localhost';
ALTER USER 'root'@'localhost' IDENTIFIED WITH 'mysql_native_password' AS
'*23AE809DDACAF96AF0FD78ED04B6A265E05AA257' REQUIRE NONE PASSWORD EXPIRE
DEFAULT ACCOUNT UNLOCK;
GRANT ALL PRIVILEGES ON . TO 'root'@'localhost' WITH GRANT OPTION;
GRANT PROXY ON "" TO 'root'@'localhost' WITH GRANT OPTION;

```

## 11.IS、PS、SYS用法介绍

---