

第九届中国数据库技术大会 DATABASE TECHNOLOGY CONFERENCE CHINA 2018

使用动态跟踪技术SystemTap 监控MySQL、Oracle性能

吕海波(VAGE)











第九届中国数据库技术大会 DATABASE TECHNOLOGY CONFERENCE CHINA 2018

如何不

《读源码从入门到放弃》

For MySQL 5.7

吕海波(VAGE)









What was really happening

"Report Regarding Comprehensive Reorganization and Standardization of the Company Export Control Related Matters"

	中国建筑股份	WW.5			ght Codamenta
_	- systematics	and the set	- A. 19 1M		rijens i
68 (B) 4 668	关于全直整领海损在公司出口管 制制天政条件服务			の収入支付 日か会報の ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・ ・	
		それでは記憶をおけるのである。 利力の表現をお言うと		*## C084	
HERE.	有种性	SHA.	THE	200	midefield
199	STREET IN	ma.	most.		B 138
THAT.	京田市府-		PREXITO	4.3	K. HRHUTZ
ES.	间查	ANG CAF	一个	*190	es ton
0.0944	C	10	3.	2	8
					me the D
R (MA)	13.	1	1	37	Section 12







S. DEPARTMENT OF COMMERCE . BUREAU OF INDUSTRY AND SECURITY







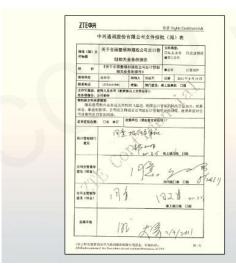




近日,美国禁止向中兴技术出口。这将进一步激发国内基础技术研究的高潮。 0S、数据库······,都是重要的基础技术。

于是,有朋友跟我说,想研究下MySQL的源码。但是,代码量太大了,不知道 从哪里入手,就放弃了。

还有朋友说,我已经编译成功了,就差看了。。。然后,然 后就没有下文了。











为什么……

因为没有建立一个正向反馈机制。

你花很多精力研究源码, 并没什么用。

你已经看了一个月代码,

再坚持一个月。

也没什么不同。

于是放弃吧。



我来过,我看过。"





如果我研究个一、两周,就能有点用处,说不定我也可以坚持下去,不断积累,集小胜为大胜,最终成为技术大牛,走向人生巅峰,迎娶白富美(这是按原来的说法,中兴被美国封锁后,现在改为"最终成为技术大牛,为科技兴国贡献力量)。

源码看个一、两周,就能有点小用途,这就是正向反馈,它会激励人不断深入下去。

如何建立正向反馈呢?

Systemtap



▶ 什么是动态跟踪技术:

又称:动态调试, Dynamic Tracing。不影响程序的情况下,动态的观察程序运行。

> 常见的动态追踪技术:

Dtrace: solaris

SystemTap: Linux

•••••

➤ Systemtap基础

基本的开发基础:变量、常量、条件、分枝、循环、子程序或子函数、数组、指针与地址。

▶ 重要概念:探针, probe

动态跟踪技术的基础,相当于数据库的触发器,探针是操作系统的触发器。我们可以利用systemtap,将我们的脚本动态的插入内存。等到一定的条件被解发,脚本就会被操作系统执行。

通常脚本都是显示个内存状态、统计时间等。



```
1 #!/usr/bin/stap
2345678
  global tm;
   probe begin {
       printf("Begin.\n");
   probe process("/data/db/mysql/lib/bin/mysqld").function
                                                ("buf page get gen") {
       tm=gettimeofday us();
10
11
12
  probe process("/data/db/mysql/lib/bin/mysqld").function
                                    ("buf page release latch").return {
       printf("Time:%d\n", gettimeofday_us()-tm );
14
15 }
```

```
//定义解释器,类似Shell脚本
  #!/usr/bin/stap
  global tm;
  probe begin {
       printf("Begin.\n");
8
  probe process("/data/db/mysql/lib/bin/mysqld").function
                                             ("buf page get gen") {
       tm=gettimeofday us();
10
11
12
13 probe process("/data/db/mysql/lib/bin/mysqld").function
                                  ("buf page release latch").return {
      printf("Time:%d\n", gettimeofday us()-tm );
14
15 }
```

```
1 #!/usr/bin/stap //定义解释器,类似Shell脚本
  global tm; //全局变量, 计录时间。
  probe begin {
      printf("Begin.\n");
8
  probe process("/data/db/mysql/lib/bin/mysqld").function
                                            ("buf page get gen") {
      tm=gettimeofday us();
10
11
12
13 probe process("/data/db/mysql/lib/bin/mysqld").function
                                 ("buf_page_release_latch").return {
      printf("Time:%d\n", gettimeofday_us()-tm );
14
15 }
```

```
1 #!/usr/bin/stap //定义解释器,类似Shell脚本
  global tm; //全局变量, 计录时间。
  probe begin {
      printf("Begin.\n");
  probe process("/data/db/mysql/lib/bin/mysqld").function
                                            ("buf page get gen") {
      tm=gettimeofday us();
10
11
12
13 probe process("/data/db/mysql/lib/bin/mysqld").function
                                 ("buf_page_release_latch").return {
      printf("Time:%d\n", gettimeofday us()-tm );
14
15 }
```

```
探针: 进入函数
 1 #!/usr/bin/stap //定义解释器,
                                   buf_page_get_gen入口点时触发。
  global tm; //全局变量, 计录时间。
                                   相当于before buf_page_get_gen触
                                   发器。
  probe begin {
      printf("Begin.\n");
  probe process("/data/db/mysql/lib/bin/mysqld").function
                                            ("buf_page_get_gen") {
      tm=gettimeofday us();
10
11
12
13 probe process("/data/db/mysql/lib/bin/mysqld").function
                                 ("buf_page_release_latch").return {
      printf("Time:%d\n", gettimeofday_us()-tm );
14
15 }
```

```
1 #!/usr/bin/stap //定义解释器,类似Shell脚本
                                      探针: 退出函数
  global tm; //全局变量, 计录时间。
                                   buf_page_release_latch入口点时触
  probe begin {
                                   发。相当于after
      printf("Begin.\n");
                                   buf_page_release_latch触发器。
8
  probe process("/data/db/mysql/lib/bin/mysqld").fur
                                            ("buf pa
                                                        Et_gen") {
      tm=gettimeofday_us();
10
11 }
12
13 probe process("/data/db/mysql/lib/bin/mysqld").function
                                 ("buf_page_release_latch").return
      printf("Time:%d\n", gettimeofday_us()-tm );
14
15 }
```

```
[root@VAGE01 ~]#
[root@VAGE01 ~]# ./lgr1.stp -x 3164
Begin.
Time:63
```

```
mysql>
mysql>
mysql>
mysql> select * from t4 where id=1;
  id | c1
   1 | BBBBBBaaaaaaabbbb
1 row in set (0.01 sec)
```

在我的测试机上,MySQL完成一次主键索引逻辑读,耗时0.063毫秒。





Oracle逻辑读方式

再来看看Oracle的逻辑读耗时吧。Oracle的逻辑读有好多种,我们试一种最普通的吧,以rowid方式访问读。这种读与MySQL中的逻辑读非常不同,但Oracle的DBA一定耳熟能详,就是先获得CBC Latch,在CBC Latch的保护下,获得块的Buffer Pin Lock。在Buffer Pin Lock保护下,读块中行数据。从获得CBC Latch始,作为逻辑读的开始。到再次获得CBC Latch,释放Buffer Pin Lock为止,为一次逻辑读的结束。

> 对应函数:

- kcbgtcr: 获得CBC Latch,在CBC Latch保护下,获得块的Buffer Pin Lock。对应MySQL的buf_page_get_gen函数。
- kcbrls:再次获得CBC Latch,并在CBC Latch保护下,释放Buffer Pin Lock。对应MySQL的buf_page_release_latch函数。



有了函数名,可以很容易写出Oracle统计逻辑读耗时的systemtap脚本。

```
1 #!/usr/bin/stap
2
3 global tm;
5 probe begin {
    printf("Begin.\n");
6
7 }
8
9 probe process("/opt/oracle/product/11204/bin/oracle").function("kcbgtcr") {
     tm = gettimeofday_us();
10
11 }
12
13 probe process("/opt/oracle/product/11204/bin/oracle").function("kcbrls").return {
     printf("Time:%d\n", gettimeofday us() - tm );
14
15 }
```





```
[root@VAGE01 ~]#
[root@VAGE01 ~]# ./lgr_o.stp -x
2630
Begin.
Time:31
```

```
SQL>
SQL> select * from t1 where
rowid='AAADs6AAEAAAACRAAB';
                                 C2
        ID C1
         1 aaaaaaaa
bbbbbbbbbbbb
SQL>
```

Oracle完成一次Rowid逻辑读,耗时0.031毫秒。

MySQL一次逻辑读用了0.065毫秒,比Oracle慢了一倍,因为MySQL的逻辑读中包含了索 引对比。

脚本很简单,可以满足我们一个小好奇心。

问题是,如何得知对应某操作的函数名呢?

如何得知MySQL逻辑读,是从buf_page_get_gen函数开始,到buf_page_release_latch函数为止?

这就需要从源码中发掘了。

按照后面我所介绍的方法, 估计找到这两个函数需要个一、两天左右时间(如果是刚开始搞,可要久一点)。

也就是说,一、两天左右的努力,可以写一个小脚本,统计某项性能指标、输出某些状态标志。让几天的努力钻研,有可用武之地。



这就是前面所说的正向反馈。读个几天代码,搞点小脚本,监控一下性能、分析一下消耗等等。再读个几天代码,继续搞小脚本......。而且,很多时候,脚本和源码,可以互相印证、互相促进,这样不断正向反馈,不断进步。

MySQL有那么多需要我们去监控性能的地方吗?

Oracle 11G有1367个等待事件和一千多个性能资料。Oracle共计用近3000个指标,反应性能、状态。其他数据库相比之下弱了很多,这给systemtap这样的动态跟踪技术,留下了巨大的空间。

阅读源码,不再是无用的屠龙之技,你可以不断在源码中发掘,利用动态跟踪技术积累脚本,既能得到成长,没准这些脚本在工作上可以帮上大忙。



下面来看"时间都去哪儿了"的脚本。

脚本很简单,就是一个统计DML执行时,各阶段时间消耗的脚本。关键还是脚本后面 隐藏的一系列方法。

写这个脚本的初衷,是为了统计一个大加载操作的时间消耗。因为前一段时间,有客户的导入数据操作,加载百万级的数据量,每天进行一次,耗时十分钟左右,DBA诊断后,根据他的环境,判断这是一个正常的耗时。感觉慢是因为没有并行,单个线程的吞吐有限。另外,Innodb引擎的事务、恢复机制,也会消耗一定时间(根据应用,表其实并不需要事务)。但客户始终觉得不应该这样慢,但又不知道时间消耗到哪儿了,于是,就有了这个systemtap脚本。



一条SQL的执行,总是分阶段的,有解析、执行、抓取等等阶段,我们来统计如下阶段的时间消耗:

- 1. 解析
- 2. 处理Binlog
- 3. 处理UNDO
- 4. 处理Redo
- 5. 总执行时间

这些阶段当然不够细,只是粗略统计下,以介绍方法为主。 掌握方法后,你可以根据需要,调节脚本,更细分、或只针对某些事件统计。

下面,我们以Binlog为例,讲解一下如何使用systemtap脚本+源码,统计DML执行期间,处理Binlog所消耗的时间。







如何开始?

Easy! 先找到SQL执行过程中,处理Binlog的函数。然后,就能像前面统计逻辑读耗时一样,统计Binlog的时间消耗。这里说的简单,是systemtap脚本简单,寻找binlog函数,就稍有点复杂了。

如何寻找处理Binlog的函数?

重要技能: 逆流而上, 顺腾摸瓜。

"流":内存数据流。

以Binlog为例,Binlog数据最终一定要写进磁盘的。写磁盘时候,一定是从某一个Buffer 写磁盘,我们称这个Buffer为Buffer Z。那么,Buffer Z的数据,又是从那里来的呢?

假设Buffer Z的数据是从Buffer Y中拷贝过去的,我们称Buffer Z的上线为Buffer Y。那么,Buffer Y的上线呢?

只要顺腾摸瓜,每找到一个上线,观察一下它的代码,看是不是我们要找的Binlog处理 函数,如此,即可。



具体过程:

要先找到写binlog到文件的函数。

这个简单,使用systemtap,显示所有名称中包含write字样的函数。脚本如下:

```
1 #!/usr/bin/stap
2
                                                       *write*: 跟踪所
  probe Ulong arg(1): 函数的第一参数。
                                                        有名称中包含
      pr 131: Binlog文件的文件号。可
                                                      write 字样的函数。
5 }
         以从/proc/进程号/fd中得到。
                                   pp(): 跟踪点的信息。
                                    ulong_arg(n): 参数
  probe process __ata/db/mysql/lib/bin,
                                                     write*") {
      if ( ulong arg(1) == 131 ) {
         printf("%s %x %x %x\n",pp(), ulong_arg(1), ulong_arg(2),
ulong_arg(3));
10
         print_ubacktrace();
         printf("----
11
                                     显示调用堆栈。即跟踪点的上
12
                                       层、上上层、.....函数。
13 }
```

```
[root@VAGE01 ~]#
[root@VAGE01 ~]# ./mio.stp -x 3164
Begin.
process("/data/db/mysql/lib/bin/mysqld").funct
ion("my_write@/home/mysql/mysql-
5.7.19/mysys/my_write.c:40") 83 38114b0 117
0x...: my_write+0x0/0x41c
0x...: inline_mysql_file_write+0x98/0x120
0x...: my_b_flush_io_cache+0x2f3/0x47f
0x...:
_ZN13MYSQL_BIN_LOG19flush_cache_to_fileEPy+0x2
7/0x8c
0x...
```

```
mysql>
mysql> insert into t1 values(43,'43aaaaabbbbb',
'AAAAAAAAA');
Query OK, 1 row affected (7 min 37.49 sec)

mysql> commit;
Query OK, 0 rows affected (0.01 sec)

mysql>
```



```
[root@VAGE01 ~]# ./mio .stp -x 3164
Begin.
process("/data/db/mysql/lib/bin/mysqld").function("my_write@/home/mysql/mysql
-5.7.19/mysys/my_write.c:40") 83 38114b0 117
0x...: my write+0x0/0x41c
0x...: inline mysql file write
                                  /0x120
                                                       my write函数:
0x...: my_b_flush_io_c
0x...: _ZN13MYSQL_BIN_ 83: 十六进制的文件号,十进制
                                                   /0x8 写Binlog的函数。
0x...: _ZN13MYSQL_BIN_ 为131。
                                                   a/0x
Ox...: _ZN13MYSQL_BIN_ 38114b0: 前文中的Buffer Z。
                      117: 十六进制的长度。一共279
0x...: _Z11mysql_parse 字节。
0x...: Z10do commandP
0x...: handle connection+0x1e0/0x2d4
0x...: pfs spawn thread+0x170/0x177
```







```
my write.c中的写I/O函数:
40 size_t my_write(File Filedes, const uchar *Buffer, size_t Count,
myf MyFlags)
41 {
   size t writtenbytes;
42
 57 for (;;)
 58
 59
       errno= 0;
 60 #ifdef WIN32
 61
       writtenbytes= my win write(Filedes, Buffer, Count);
 62 #else
       writtenbytes= write(Filedes, Buffer, Count);
 63
 64 #endif
                                         从这里,库函数write将Binlog数
128 DBUG RETURN(sum written);
                                     据写进磁盘。下面,我们开始找它
129 } /* my_write */
                                     的上线,Buffer Z。
```



▶ 准备工作

首先,要准备一条SQL,SQL一定要有标识度:

insert into t1 values(42, 'aaaaaaaa', 'AAAAAAAA');

有大量的重复数据,容易辨认。



```
[root@VAGE01 ~]# gdb -p 3164
GNU gdb (GDB) Red Hat
Enterprise Linux 7.6.1-51.el7
(gdb)
(gdb) b my_write.c:63
Breakpoint 1 at 0x1877865:
file /home/mysql/mysql-
5.7.19/mysys/my_write.c,
line 63.
(gdb) C
```

```
mysql>
mysql>
mysql>
mysql> insert into t1 values(42,
'aaaaaaaa', 'AAAAAAAA');
 ===>HANG住<===
```



```
(gdb) b my_write.c:63
Breakpoint 1 at 0x1877865:
file /home/mysql/mysql-
5.7.19/mysys/my write.c,
line 63.
(gdb)C
Continuing.
```

```
mysql>
mysql>
mysql>
mysql> insert into t1 values(42,
'aaaaaaaa', 'AAAAAAAA');
Query OK, 1 row affected (7 min
37.49 sec)
mysql>
```



```
(gdb) b my write.c:63
Breakpoint 1 at 0x1877865:
file /home/mysql/mysql-
5.7.19/mysys/my_write.c,
line 63.
(gdb)C
Continuing.
Breakpoint 1, my_write
(Filedes=131, Buffer=0x3810a40
"\362\337\275Z\"q\004",
Count=284, MyFlags=52) at
/home/mysql/mysql-
5.7.19/mysys/my write.c:63
63 writtenbytes= write(Filedes,
Buffer, Count);
(gdb)
```

```
mysql>
mysql>
mysql>
mysql> insert into t1 values(42,
'aaaaaaaa', 'AAAAAAAA');
Query OK, 1 row affected (7 min
37.49 sec)

mysql> commit;
===>HANG住<===</pre>
```







```
(gdb) b my_write.c:63
Breakpoint 1 at 0x1877865:
file /home/mysql/mysql-
5.7.19/mysys/my_write.c,
line 63.
(gdb)C
Continuing.
Breakpoint 1, my_write
(Filedes=131, Buffer=0x3810a40
"\362\337\275Z\" (\004",
Count=284, MyFlags=52) at
/home/mysql/mysql-
5.7.19/mysys/my_write.c:63
63 writtenbytes= write(Filedes,
Buffer, Count);
(gdb)
```

```
mysql>
mysql>
mysql>
mysql> insert into t1 values(42,
'aaaaaaaa', 'AAAAAAAA');
Query OK, 1 row affected (7 min 37.49 sec)

mysql> commit;
===>HANG住<===</pre>
```





```
(gdb) b my_write.c:63
Breakpoint 1 at 0x1877865:
file /home/mysql/mysql-
5.7.19/mysys/my_write.c,
line 63.
(gdb)C
Continuing.
Breakpoint 1, my_write
(Filedes=131, Buffer=0x3810a40
"\362\337\275Z\" (\004",
Count=284, MyFlags=52) at
/home/mysql/mysql-
5.7.19/mysys/my_write.c:63
63 writtenbytes= write(Filedes,
Buffer, Count);
(gdb)
```

Filedes: 文件号

Buffer: 缓存

Count: 写长度





```
(gdb) b my_write.c:63
Breakpoint 1 at 0x1877865:
file /home/mysql/mysql-
5.7.19/mysys/my_write.c,
line 63.
(gdb)C
Continuing.
Breakpoint 1, my_write
(Filedes=131, Buffer=0x3810a40
"\362\337\275Z\" (\004",
Count=284, MyFlags=52) at
/home/mysql/mysql-
5.7.19/mysys/my write.c:63
63 writtenbytes= write(Filedes,
Buffer, Count);
(gdb) p Filedes
$1 = 131
```

```
[root@VAGE01 fd]# pwd
/proc/3164/fd
[root@VAGE01 fd]# ls -lFrt|grep 131
l-wx----. 1 mysql mysql 64 Mar
30 09:27 131 ->
/data/db/mysql/binlog.000048
```





```
(gdb)C
Continuing.
Breakpoint 1, .
63 writtenbytes= write(Filedes,
Buffer, Count);
(gdb) p Filedes
$1 = 130
(gdb) p Buffer
$2 = (const uchar *) 0x3810a40
(gdb)
```

```
[root@VAGE01 fd]# pwd
/proc/7792/fd
[root@VAGE01 fd]# ls -lFrt|grep 130
l-wx----. 1 mysql mysql 64 Mar
30 09:27 131 ->
/data/db/mysql/binlog.000048
```







```
(gdb)C
Continuing.
Breakpoint 1,
    writtenbytes
63
(gdb) p Filedes
$1 = 130
                             0x61080000
(gdb) p Buffer
$2 = (const ucl
                             0x00000861
(gdb) x/80x 0x
0x3810a40:
                                                     004100
                                                                  0x000d0e00
                                                   0317473
0x3810ae0:
                                                                  0x00317402
0x3810af0:
                0x0f0f
                                                                  0xdhdff85f
                                0x3c003c04
                                                 0xf3e60700
0x3810b00:
                0xca1e
                                0x3f0004a0
                                                                  0x0000000d
                                                 0xca000000
                0x0000 200
0x3810b10:
                                0x01000000
                                                 0x03000200
                                                                  0x002bf8ff
0x3810b20:
                0x61080000
                                0x61616161
                                                 0x08616161
                                                                  0x41414141
0x3810b30:
                0x41414141
                                0xbddff2e8
                                                 0x1ca5dc41
                                                                  0xbddff2e8
0x3810b40:
                0xa0ca105a
                                0x001f0004
                                                 0x0de90000
                                                                  0x00000000
```











```
(gdb)C
Continuing.
Breakpoint 1,
                           内存地址
                                         数
                                             据
    writtenby
63
                         0x3810b20: 0x61080000
(gdb) p Filed
$1 = 130
                         0x3810b20: 0x00000861
(gdb) p Buffe
$2 = (const u)
(gdb) x/80x 0
0x3810a40:
                                                                0x000d0e00
                         0x3810b22: 0x08616161
0x3810ae0:
                                                                0x00317402
                0x0100
                                X65/40500
                                                0x0031/4/3
0x3810af0:
                0x0f0f
                                                                0xdbdff85f
                                0x3c003c04
                                                0xf3e60700
0x3810b00:
                0xca1e
                                                                0x0000000d
                                0x3f0004a0
                                                0xca000000
                0x0000
0x3810b10:
                                0x01000000
                                                0x03000200
                                                                0x002bf8ff
                0x61080000
0x3810b20:
                                0x61616161
                                                0x08616161
                                                                0x41414141
0x3810b30:
                0x41414141
                                0xbddff2e8
                                                0x1ca5dc41
                                                                0xbddff2e8
0x3810b40:
                0xa0ca105a
                                0x001f0004
                                                0x0de90000
                                                                0x00000000
```











(gdb) watch *0x3810b22 Hardware watchpoint 2: *0x3810b22 (gdb)

Watch: 监察点,观察点。

GDB中的触发器,当进程读或写指定内存地址的数据时被触发。

如前文所述,0x3810b22,是我们的"Buffer Z",只要监控Bufer Z,看"谁"向Buffer Z中写数据,就能找到Buffer Z的上线: Buffer Y。





```
(gdb) watch *0x3810b22
Hardware watchpoint 2:
*0x3810b22
(gdb) c
Continuing.
```

```
mysql> insert into t1 values(42,
'aaaaaaaa', 'AAAAAAAA');
Query OK, 1 row affected (7 min
37.49 sec)
mysql> commit;
Query OK.....
mysql> insert into t1 values(42,
'aaaaaaaa', 'AAAAAAAA');
Query OK .....
mysql>
```











```
(gdb) watch *0x3810b22
Hardware watchpoint 2:
*0x3810b22
(gdb) c
Continuing.
Hardware watchpoint 2:
*0x3810b22
Old value = 1630745612
New value = 1647588364
0x00007ffff6ae6e7b in
 _memcpy_ssse3_back () from
/lib64/libc.so.6
(gdb)
```

```
mysql> commit;
Query OK.....
mysql> insert into t1
values(41, 'aaaaaaaa', 'AAAAAAAA');
Query OK .....
mysql> commit;
 ===>HANG住<===
```









```
(gdb) watch *0x3810b22
Hardware watchpoint 2:
*0x3810b22
(gdb) c
Continuing.
Hardware watchpoint 2:
*0x3810b22
Old value = 1630745612
New value = 1647588364
0x00007ffff6ae6e7b in
 _memcpy_ssse3_back () from
/lib64/libc.so.6
(gdb)
```

原值: 6133340C 新值: 61616108

新值是08616161。









```
(gdb) watch *0x3810b22
Hardware watchpoint 2:
*0x3810h22
(gdb) c
Continuing.
Hardware watchpoint 2:
*0x3810b22
Old value = 1630745612
New value = 1647588364
0x00007ffff6ae6e7b in
 memcpy ssse3 back () from
/lib64/libc.so.6
(gdb)
```

Libc库中的函数,它是memcpy的底层函数。实现内存块拷贝。

由于这个断点还是在commit时被触发的,所以Buffer Y不会是我们要找的源头: SQL执行时处理Binlog的函数。继续找Buffer Y的上线。

首先,确定Buffer Y的准确地址,然后使用同样的方法,使用watch观察点,监控Buffer Y的地址即可。





info registers (gdb)

0x38113h1 rax rbx 0xd6b639c0

0x7fff89fff0h8 rcx

0x7ffff6ae7062 rdx

0x7fff89fff0dc rsi

rdi 0x3810B45

0x7fffd 5f2 rbp

0x7f f228 rsp

rip

还记得Buffer Z的准确地址。 和当前rdi的值相差多少? 0x

(gdb)

煨水或多存器值i中,复制内存:

memcpy(rdi, rsi, rdx)

这是因为,watch观察点,是在被观察的内存 位置被修改后触发,而不是修改前触发。观察 因此,此处的rdi,rsi严格上不是memcpy 的第一、二个参数值了。Memcpy已经改了 它们。但是memcpy虽对它们有改变,但改 变不大。

所以我们可以得出: rdi不是Buffer Z, 但它在Buffer Z附近。

同理: rsi也不是Buffer Y的准确地址, 也就是说: rdi-0x23, 等于B Buffer Y在rsi附近。

有很大的可能, rsi-0x23, 就是Buffer Y的地址。







存。

di∖

等,

的

四

(gdb) x/4x 0x7fff89fff0b9 0x7fff89fff0b9: 0x61616108 UXU_ 0x41410861 0x41414141 rdi (gdb) ·cx等, 午89fff0dc也减去0x23字节,为 数的 ~f0b9,显示此处的内存值 0x> E、四 0x7fff89fff0b9处的值,正是 08616161......。也就是说,它就是我 们要找的Buffer Z的上线,Buffer Y。 **第二个** 参数是源。 rsi:源缓存。 rdi: 目标缓存。







(gdb) x/4x 0x7fff89fff0b9

0x41410861 0x41414141

(gdb) delete 2

(gdb) watch *0x7fff89fff0b9

Hardware watchpoint 3: *0x7fff89fff0b9

(gdb)

删除刚才的watch 监察点,设置新的 监察点为 0x7fff89fff0b9。









```
(gdb) delete 2
(gdb) watch *0x7fff89fff0b9
Hardware watchpoint 3:
*0x7fff89fff0b9

(gdb) c
Continuing.
Hardware watchpoint 3:
```

Old value = 1094795528
New value = 1111638536
0x00007ffff6ae75e4 in
 __memcpy_ssse3_back () from
/lib64/libc.so.6

*0x7fff89fff0b9

```
mysql> insert into t1 values(42,
    'aaaaaaaa', 'AAAAAAAA');
Query OK .....

mysql> commit;
Query OK.....

mysql> insert into t1 values(43,
    'bbbbbbbb', 'BBBBBBBB');
===>HANG住<===</pre>
```





```
bt,显示调用栈。
调用栈是函数一层一层的调用关系。比如此处,__memcpy_ssse3_back被
my b safe write调用, my b safe write被
Log event::wrapper my b safe write调用,等等。
                         THD::binlog query的下层函数,调用memcpy,将
                         Buffer X中的数据,拷贝到Buffer Y。
(gdb) bt
                         然后又经Buffer Y拷贝到Buffer Z,又经Buffer Z写到
#0
   0x00007ffff6ae75e4 in_
   0x00000000185d155
#1
                    但THD::binlog_query不一定是DML执
         at /home/mys
                    行期间,唯一的Binlog处理函数。因
                                                |行期间,处理Binlog的|
   0x00000000017a11aa
#2
                    此,要继续"顺藤摸瓜"。
   0x00000000017bda73
#3
   0x0000000017c84c7
#8
   0x0000000017f04c4 in THD::binlog query
   0x000000001761333 in Sql cmd insert::mysql insert
#17 0x00007fffff7bc7df3 in start thread () from
#18 0x00007ffff6a933dd in clone () from /lib64/libc.so.6
(gdb)
```







```
(gdb) info registers
                 0x7fffb42a4550
                                    140736216057168
rcx
rdx
                0x7ffff6ae75d6
                                    140737332016598
                                    140736216057168
                 0x7fffb4219598
rsi
                 0x7fff89fff0b4
rdi
rbp
                 0x7fffd6b5f7f0
                                    0
                                       这一次是在rsi之前。
                 0x7fffd6b5f7c8
rsp
                                              e75e4
rip
                                    0x7ff
                0x7ffff6ae75e4
                 <__memcpy_ssse3_back+10
                                              >
(gdb)
(gdb) x/100x 0x7fffb4219508
                                            0000000000
0x7fffb4219508: 0x0210a108 0x00000000 0x000000000
0x7fffb4219578: 0x8f8f8f8f 0x8f8f8f8f 0x000033f8
                                            0x62620800
0x7fffb4219588: 0x62626262 0x42086262 0x42424242
                                            0x42424200
```

步1:

确认Buffer X的准确 地址。 可以从rsi入手,因为

rsi是Buffer X的附近。

步2:

在Buffer X处设置观察点。

步3:

再运行一条近似的SQL,等待观察点被触发。

步4:

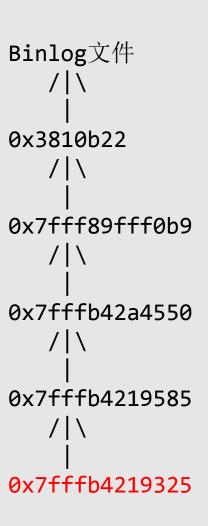
使用bt命令,观察调用栈中是否有Binlog处理函数。





从当前的调用栈来看,已经出了Binlog的范围。 从函数名上看,当前在"解析"中。 我们的源头,已经找到了: 0x7fffb4219325

```
(gdb) bt
    in memcpy ssse3 back ()
    in get_text
    in lex one token
    in MYSQLlex
#4
    in MYSQLparse
   in parse_sql
#5
    in mysql_parse
    in dispatch_command
    in do command
#8
    in handle connection
#10 in pfs_spawn_thread
#11 in start_thread ()
#12 in clone
```





〉 顺腾摸瓜结果:

除了THD::binlog_query函数,还有binlog_log_row函数,也是DML期间处理Binlog的函数。

▶ 更进一步:

我们已经知道binlog_query、binlog_log_row函数是处理Binlog的起点,有兴趣的话,你可以看看它的源码。它的代码并不长,使用gdb,单步跟踪一遍它的执行过程,看看在它们的执行过程中,数据是怎么处理的,需要哪些锁,数据最后怎么被复制到IO_CACHE中。

我们也知道,在提交时,Binlog数据用my_write函数写到磁盘文件中,在my_write上设断点,使用bt显示调用堆栈,可以很容易发现,MYSQL_BIN_LOG类的flush cache to file函数,是IO的发起函数。相关代码在binlog.cc的8850行。



我们也知道,在提交时,Binlog数据用my write函数写到磁盘文件中,在my write上 设断点,使用bt显示调用堆栈,可以很容易发现,MYSQL BIN LOG类的 flush_cache_to_file函数,是IO的发起函数。相关代码在binlog.cc的8850行。

```
Breakpoint 3, my_write (Filedes=125,
      Buffer=0x380d280 "# Time: 2018-04-18T01:38:54.290419Z\n# User@Host: root[root] @ [192.168
_sent: 0 Rows_examined: 0\nSET timestamp=1524015534;\ninsert "..., Count=245, MyFlags=20) at
               size t sum written= 0:
(qdb) c
Continuing.
Breakpoint 3, my_write (Filedes=131, Buffer=0x38112e0 "\264\241\326Z\"q\004", Count=280, MyFla
               size_t sum_written= 0;
43
(adb) bt
     my_write (Filedes=131, Buffer=0x38112e0 "\264\241\326Z\"q\004", Count=280, MyFlags=52) at 0x00000000185a1e9 in inline_mysql_file_write (src_file=0x21eeef0 "/home/mysql/mysql-5.7.1 buffer=0x38112e0 "\264\241\326Z\"q\004", count=280, flags=52) at /home/mysql/mysql-5.7.19/0x00000000185d6cc in mv b flush io cache (info=0x2d84ec8 <mysql_bin_log+840>, need_append 0x0000000017ea4e1 in MYSQL_BIN_LOG::flush_cache_to_file (this=0x2d84b80 <mysql_bin_log>,
      at /home/mysql/mysql-5.7.19/sql/binlog.cc:8850
      0x000000001/eb2b0 in MYSQL_BIN_LOG::ordered_commit (this=0x2d84b80 <mysql_bin_log>, thd=0
      at /home/mysql/mysql-5.7.19/sql/binlog.cc:9246
      0x0000000017e9959 in MYSQL_BIN_LOG::commit (this=0x2d84b80 <mysql_bin_log>, thd=0x7fffbc2
     0x000000000f1fadb in ha_commit_trans (thd=0x7fffbc296060, all=true, ignore_global_read_lo
     0x00000000163d42a in trans_commit (thd=0x7fffbc296060) at /home/mysql/mysql-5.7.19/sql/tr
0x00000000153f992 in mysql_execute_command (thd=0x7fffbc296060, first_level=true) at /home/mysql_execute_command (thd=0x7fffbc296060, first_level=true) at /home/mysql_execute_command (thd=0x7fffbc296060, first_level=true) at /home/mysql_execute_command (thd=0x7fffbc296060).
     0x000000001542ed4 in mysql_parse (thd=0x7fffbc296060, parser_state=0x7fffd6b61db0) at /ho
#10 0x0000000015386a4 in dispatch_command (thd=0x7fffbc296060, com_data=0x7fffd6b62df0, comma
#11 0x000000001537582 in do_command (thd=0x7fffbc296060) at /home/mysql/mysql-5.7.19/sql/sql_
#12 0x00000000166a016 in handle_connection (arg=0x39460f0) at /home/mysql/mysql-5.7.19/sql/co
#13 0x000000001cf5690 in pfs_spawn_thread (arg=0x3928730) at /home/mysql/mysql-5.7.19/storage
#14 0x00007fffff7bc7df3 in start_thread () from /lib64/libpthread.so.0
#15 0x00007fffff6a933dd in clone () from /lib64/libc.so.6
(qdb)
```



) 更进一步:

从MYSQL BIN LOG::flush cache to file, 到在my write中调用libc库函数 write,完成写I/O,整个流程也不复杂,使用gdb跟踪一下,留意一下数据的流动、并 发控制(锁,或Mutex锁)的申请释放,将使你对这一部分有更深入的了解。

很多朋友像了解下源码,但不知道如何入手,使用这种方法,可以很简单的找到感 兴趣的入手点。

另外,刚开始阅读源码的时候,要尽量缩小自己的目标与范围,整个解析、执行、 提交等过程是很复杂的。但提交时的I/O处理部分,目标就小多了,也简单很多。再进一 步的: "提交时Binlog的I/O处理",这个目标就更小了,代码量也更少,更容易入手。

而且,利用systemtap脚本,还可以不断写些时间统计、性能监控脚本,既有实际 意义、又可以使用自己不断成长。

再而且,如果觉得源码那里不合适了,还可以自己改改。

这就是我所说的正向反馈式的"源码研究"。



```
找到了BinLog函数,可以很马上写出统计时间的systemtap脚本:
#!/usr/bin/stap
global tm binlog, tt binlog;
probe begin {
   printf("Begin.\n");
probe process("/data/db/mysql/lib/bin/mysqld").function("binlog_log_row") {
   tm binlog = gettimeofday us();
probe
process("/data/db/mysql/lib/bin/mysqld").function("binlog_log_row").return {
   tt binlog += (gettimeofday us() - tm binlog);
   printf("binlog:%d\n", tt_binlog);
```

```
[root@VAGE01 stp]#
[root@VAGE01 stp]# ./exmp1.stp -g 3714
Begin.
binlog:80
                             mysql>
                             mysql>
                             mysql> insert into t1
                                 -> values(5b,'55bbbbbbbbbbbb', 'AAAAAAAAA');
                             Query OK, 1 row affected (0.00 sec)
```

> 各个阶段对应的函数

使用前面所述方法,我们可以找到各个阶段对应的函数,结果如下表:

解析	parse_sql	
Binlog	binlog_log_row、binlog_query	
回滚	trx_undo_report_row_operation(包含trx_undof_page_add_undo_rec_log)	
Redo	UNDO的Redo	trx_undof_page_add_undo_rec_log、
	页的Redo	<pre>page_cur_insert_rec_write_log btr_cur_update_in_place_log page_cur_delete_rec_write_log</pre>
总执行	dispatch_command	





函数有了,可以写了如下时间统计脚本:

```
#!/usr/bin/stap
global tm_binlog, tm_undo, tm_u_redo, tm_redo, tm_exec, tm_parse;
global tt_binlog, tt_u_redo, tt_undo, tt_redo, tt_parse, tt_exec;
global cs;
probe begin {
    printf("Begin.\n");
probe process("/data/db/mysql/lib/bin/mysqld").function("parse_sql") {
    if ( tid() == $1 )
        tm parse = gettimeofday us();
probe process("/data/db/mysql/lib/bin/mysqld").function("parse_sql").return {
    if ( tid() == $1 ) {
        tt_parse += (gettimeofday_us() - tm_parse);
}
probe process("/data/db/mysql/lib/bin/mysqld").function("binlog_query") {
    if ( tid() == $1 )
        tm binlog = gettimeofday us();
probe process("/data/db/mysql/lib/bin/mysqld").function("binlog_query").return {
    if ( tid() == $1 ) {
            tt_binlog += (gettimeofday_us() - tm_binlog);
}
probe process("/data/db/mysql/lib/bin/mysqld").function("binlog_log_row") {
    if ( tid() == $1 )
        tm_binlog = gettimeofday_us();
probe process("/data/db/mysql/lib/bin/mysqld").function("binlog_log_row").return {
    if ( tid() == $1 ) {
        tt_binlog += (gettimeofday_us() - tm_binlog);
probe process("/data/db/mysql/lib/bin/mysqld").function("trx_undo_report_row_operation") {
    if ( tid() == $1 )
        tm_undo = gettimeofday_us();
probe process("/data/db/mysql/lib/bin/mysqld").function("trx undof page add undo rec log") {
    if ( tid() == $1 )
        tm_u_redo = gettimeofday_us();
probe process("/data/db/mysql/lib/bin/mysqld").function("trx undof page add undo rec log").return {
```

```
#!/usr/bin/stap
global tm_binlog, tm_undo, tm_uglobal tt_binlog tt_u redo t
                                                         tm parse;
                                   在函数执行完成时,取
                                                        , t<del>t</del> exec;
ğlobal
                                得函数的耗时,并累加进全
                                局变量。
probe b
                                                                定义全局变量。
    pri
                                                                Systemtap的全局变
probe process(\data/db/mysql/lib/bin/mysqld").function("parse
                                                                统将自动赋值为0
    if ( tid() == $1
        <del>tm_parse = gett</del>imeofday us();
probe process("/data/
                                                           arse sql").return {
                      time.s(1): 定义器探针,每一秒触发一次。
    if ( tid() == $1
                      触发后输出当前的累计时间。
        tt parse += (
                      效果类似于 "vmstat 1",一秒显示一行。
probe timer.s(1)
    if ( cs \%^220 = 0 ) {
        printf("parse\text{tbin log\tundolog\t redo for undo\tredolog\t Total
Time\n");
        printf(
\n")
    printf("%d\t%d\t%d\t%d\t%d\t%d\n", tt parse, tt binlog, tt undo, tt u redo,
tt redo, tt exec );
    CS++;
```

运行脚本需要知道线程号,在performance_schema.threads视图中的PROCESSLIST_ID列、THREAD_OS_ID列,可以得到我们需要的线程ID。以下语句,得到当前线程的ID:

```
mysql>
mysql> select PROCESSLIST_ID, THREAD_OS_ID from
performance schema.threads where PROCESSLIST INFO like 'select
PROCESSLIST_ID, THREAD_OS_ID from performance_schema.threads where
PROCESSLIST_INFO%';
  PROCESSLIST_ID | THREAD_OS_ID |
                           3659
1 row in set (0.00 sec)
mysql>
```

[root@VAGE01 old]# ./tmstat.stp -x 3047 3659 Begin.

3047: mysqld进程号。 3569: 要跟踪的线程号。

```
bin log undolog redo for undo redolog Total Time
parse
0
       0
0
0
97
      25649
              38279 6641
                                   9854
    41822 88855
97
                     15717
                                   23190
                                          502118
97
     41822 88855
                     15717
                                   23190
                                           502118
^C
```

```
mysql> insert into t5(c1, c2)
  select c1, c2 from t5;
Query OK, 4096 rows affected (0.50 sec)
Records: 4096 Duplicates: 0 Warnings: 0
```

[root@VAGE01 old]# ./tmstat.stp -x 3047 3659 Begin. bin log undolog redo for undo redolog parse 脚本统计的总执行时间, 502118微秒,和MySQL自己 0 0 0 统计的时间一致, 0.50秒。 97 25649 38279 6641 9854 97 41822 88855 15717 23190 502178 41822 502118 97 88855 15717 23190

Redo For Undo的时间,占了总时间的3.1%,页Redo占了总时间的4.6%,Redo共计占7.7%

Undo占了总时间的17.7%

Binlog占了总时间的8.3%

这些为了事务、更高的可靠性,而占的总时间,达到33.7%。

```
mysql> insert into t5(c1, c2)
  select c1, c2 from t5;
Query OK, 4096 rows affected (0.50 sec)
Records: 4096 Duplicates: 0 Warnings: 0
```

开课了

我在ITPUB平台上,将陆续发布系列Oracle/MySQL调试课程,如果想进 -步学习Oracle,或学习更多调试MySQL的技巧,敬请期待。





THANKS SQL BigDate



讲师申请

联系电话(微信号): 18612470168

关注"ITPUB"更多技术干货等你来拿~

与百度外卖、京东、魅族等先后合作系列分享活动





让学习更简单

微学堂是以ChinaUnix、ITPUB所组建的微信群为载体,定期邀请嘉宾对热点话题、技术难题、新产品发布等进行移动端的在线直播活动。

截至目前,累计举办活动期数60+,参与人次40000+。

◯ ITPUB学院

ITPUB学院是盛拓传媒IT168企业事业部(ITPUB)旗下 企业级在线学习咨询平台 历经18年技术社区平台发展 汇聚5000万技术用户 紧随企业一线IT技术需求 打造全方式技术培训与技术咨询服务 提供包括企业应用方案培训咨询(包括企业内训) 个人实战技能培训(包括认证培训) 在内的全方位IT技术培训咨询服务

ITPUB学院讲师均来自于企业
一些工程师、架构师、技术经理和CTO
大会演讲专家1800+
社区版主和博客专家500+

培训特色

无限次免费播放 随时随地在线观看 碎片化时间集中学习 聚焦知识点详细解读 讲师在线答疑 强大的技术人脉圈

八大课程体系

基础架构设计与建设 大数据平台 应用架构设计与开发 系统运维与数据库 传统企业数字化转型 人工智能 区块链 移动开发与SEO



联系我们

联系人: 黄老师

电 话: 010-59127187 邮 箱: edu@itpub.net 网 址: edu.itpub.net

培训微信号: 18500940168