




Strace, Perf & Gdb – three friends for the DBA

Kamil Stawiarski (@ora600pl) 

DB FILE SCATTERED READ – MBRC AND ASYNC IO ON LINUX

Recently one of my students asked me if there is any correlation between the AU size in ASM diskgroup and parameter `db_file_multiblock_read_count`. I made some tests but results are very strange.

I have 4 same tables, located in 4 different diskgroups:

```
1
2     1  select t.table_name,  g.name, t.blocks, g.ALLOCATION_UNIT_SIZE
3     2  from dba_tables t, dba_data_files d, v$asm_diskgroup g
4     3  where t.tablespace_name=d.tablespace_name
5     4  and   g.name=substr(d.file_name,2,instr(d.file_name,'/')-2)
6     5  and   t.owner='SH'
7     6* and   t.table_name like 'SALES%'
8 SQL> /
9
10  TABLE_NAME NAME                                BLOCKS ALLOCATION_UNIT_SIZE
11  -----
12  SALES16M    DATA16M                                93551   16777216
13  SALES32M    DATA32M                                93551   33554432
14  SALES1M     DATA1M                                  93551   1048576
15  SALES4M     DATA4M                                  93551   4194304
```

All tables are smaller then a hidden parameter "`_small_table_threshold`" so table access full on all tables will be using db file scattered read.

```
1
2 SQL> ;
3     1  select p.kspinm, v.kspstvl
4     2  from x$kspin p, x$kspsv v
5     3  where p.indx=v.indx
6     4* and p.kspinm='_small_table_threshold'
7 SQL> /
8
9  KSPINM                                KSPSTVL
10  -----
11  _small_table_threshold                101527
```

Let's start the main test. Parameter `db_file_multiblock_read_count` is set to the default value

```
1
2 [oracle@rico ~]$ sqlplus sh/sh
3
4 SQL*Plus: Release 12.1.0.2.0 Production on Fri Jul 17 10:30:36 2015
5
6 Copyright (c) 1982, 2014, Oracle. All rights reserved.
7
8 Last Successful login time: Wed Jul 01 2015 11:16:33 +02:00
9
10 Connected to:
11 Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit Production
12 With the Partitioning, Automatic Storage Management, OLAP, Advanced Analytics
13 and Real Application Testing options
14
15 SQL> sho parameters db_file_multiblock
16
17 NAME                                TYPE      VALUE
18 -----
19 db_file_multiblock_read_count        integer   128
20 SQL>
```

I will identify the server process, correlated with my session

```

1
2  SQL> ;
3      1  select spid
4      2  from v$process p, v$session s
5      3  where p.addr=s.paddr
6      4* and s.username='SH'
7  SQL> /
8
9  SPID
10 -----
11 10430

```

Now I will use strace on this server process:

```

1 [root@rico ~]# strace -c -p 10430
2 Process 10430 attached - interrupt to quit

```

I will gather results for the following queries:

```

1
2  SQL> select count(1) from SALES16M;
3
4      COUNT(1)
5  -----
6      19295703
7
8  SQL> select count(1) from SALES32M;
9
10     COUNT(1)
11  -----
12     19295703
13
14  SQL> select count(1) from SALES1M;
15
16     COUNT(1)
17  -----
18     19295703
19
20  SQL> select count(1) from SALES4M;
21
22     COUNT(1)
23  -----
24     19295703

```

Results for SALES16M

```

1 [root@rico ~]# strace -c -p 10430
2 Process 10430 attached - interrupt to quit
3 ^CProcess 10430 detached
4 % time      seconds  usecs/call   calls   errors syscall
5 -----
6 100.00     0.053000      71         750
7 0.00       0.000000       0          31
8 0.00       0.000000       0           2
9 0.00       0.000000       0          16
10 0.00       0.000000       0           9
11 0.00       0.000000       0           1
12 0.00       0.000000       0           1
13 0.00       0.000000       0          28
14 0.00       0.000000       0           3
15 0.00       0.000000       0           2

```

```

14      0.00      0.000000          0          1          brk
15      0.00      0.000000          0          1          uname
16      0.00      0.000000          0         10          fcntl
17      0.00      0.000000          0          1          getrlimit
18      0.00      0.000000          0        164          getrusage
19      0.00      0.000000          0         28          times
20      0.00      0.000000          0          1          getuid
21      0.00      0.000000          0          1          io_getevents
22      0.00      0.000000          0          1          io_submit
-----
23      100.00    0.053000                    1051          7 total
24
25
26

```

Results for SALES32M

```

1
2      [root@rico ~]# strace -c -p 10430
3      Process 10430 attached - interrupt to quit
4      ^CProcess 10430 detached
5      % time      seconds  usecs/call      calls      errors syscall
6      -----
7      98.16      0.053462          72          747          pread
8      1.84      0.001000          500          2          read
9      0.00      0.000000          0           2          write
10     0.00      0.000000          0           1          open
11     0.00      0.000000          0           1          close
12     0.00      0.000000          0           1          mmap
13     0.00      0.000000          0           3          fcntl
14     0.00      0.000000          0           1          getrlimit
15     0.00      0.000000          0        164          getrusage
16     0.00      0.000000          0         28          times
17     -----
18     100.00    0.054462                    950          total
19
20

```

Results for SALES1M

```

1
2      [root@rico ~]# strace -c -p 10430
3      Process 10430 attached - interrupt to quit
4      ^CProcess 10430 detached
5      % time      seconds  usecs/call      calls      errors syscall
6      -----
7      48.58      0.015000          26          583          io_getevents
8      39.99      0.012349          75          164          pread
9      11.43      0.003528          6           583          io_submit
10     0.00      0.000000          0           3          read
11     0.00      0.000000          0           2          write
12     0.00      0.000000          0           2          open
13     0.00      0.000000          0           2          close
14     0.00      0.000000          0           1          semctl
15     0.00      0.000000          0           3          fcntl
16     0.00      0.000000          0           1          getrlimit
17     0.00      0.000000          0        164          getrusage
18     0.00      0.000000          0         29          times
19     0.00      0.000000          0           3          semtimedop
20     -----
21     100.00    0.030877                    1540          total
22
23

```

Results for SALES4M

```

1      [root@rico ~]# strace -c -p 10430
2      Process 10430 attached - interrupt to quit

```

```

3 ^CProcess 10430 detached
4 % time      seconds  usecs/call      calls      errors  syscall
-----
5 95.16      0.043000      65             664
6 2.62       0.001186      14             83
7 2.21       0.001000      12             83
8 0.00       0.000000      0              2
9 0.00       0.000000      0              2
10 0.00       0.000000      0              21
11 0.00       0.000000      0              12
-----
12 100.00     0.045186      867            total
13
14

```

As we can see, the ASYNC IO represented by syscalls: *io_submit* and *io_getevents* are mostly used while scanning table located in diskgroup with 1MB AU.

At first I thought that this is correlated with *db_file_multiblock_read_count*. The default value (128) multiplied by 8k (block size for my database) is exactly 1M. So I made a few more tests with altered values for this parameter. First results were promising – I have managed to get ASYNC IO for all of the tables.

But then I made the final (a little stupid :)) test. After creating a session I have altered value of the *db_file_multiblock_read_count* to 128. Yes: 128. So the same as default:

```

1
2 [oracle@rico ~]$ sqlplus sh/sh
3
4 SQL*Plus: Release 12.1.0.2.0 Production on Fri Jul 17 10:49:25 2015
5
6 Copyright (c) 1982, 2014, Oracle. All rights reserved.
7
8 Last Successful login time: Fri Jul 17 2015 10:48:43 +02:00
9
10 Connected to:
11 Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit Production
12 With the Partitioning, Automatic Storage Management, OLAP, Advanced Analytics
13 and Real Application Testing options
14
15 SQL> sho parameters db_file_multiblock
16
17 NAME                                TYPE      VALUE
18 -----
19 db_file_multiblock_read_count        integer   128
20
21 SQL> alter session set db_file_multiblock_read_count=128;
22
23 Session altered.
24
25

```

Let's check the results now:

SALES16M

```

1 [root@rico ~]# strace -c -p 11208
2 Process 11208 attached - interrupt to quit
3 ^CProcess 11208 detached
4 % time      seconds  usecs/call      calls      errors  syscall
-----
5 53.96      0.007032      9             747
6 46.04      0.006000      8             747
7 0.00       0.000000      0             31
8 0.00       0.000000      0             2
9 0.00       0.000000      0             16
10 0.00       0.000000      0             9
11 0.00       0.000000      0             1
12 0.00       0.000000      0             1
13
14

```

```

11      0.00      0.000000          0      28      lseek
12      0.00      0.000000          0       8      mmap
13      0.00      0.000000          0       1      munmap
14      0.00      0.000000          0       1      brk
15      0.00      0.000000          0       4      pread
16      0.00      0.000000          0       1      uname
17      0.00      0.000000          0      10      fcntl
18      0.00      0.000000          0       1      getrlimit
19      0.00      0.000000          0      164     getrusage
20      0.00      0.000000          0      34      times
21      0.00      0.000000          0       1      getuid
-----
22 100.00      0.013032          1807      7 total
23
24
25
26

```

SALES32M

```

1
2 [root@rico ~]# strace -c -p 11208
3 Process 11208 attached - interrupt to quit
4 ^CProcess 11208 detached
5 % time      seconds  usecs/call      calls      errors syscall
6 -----
7 54.18      0.008278          11      746      io_submit
8 45.82      0.007000          9      746      io_getevents
9 0.00      0.000000          0       2      read
10 0.00      0.000000          0       2      write
11 0.00      0.000000          0       1      open
12 0.00      0.000000          0       1      close
13 0.00      0.000000          0       1      pread
14 0.00      0.000000          0       3      fcntl
15 0.00      0.000000          0       1      getrlimit
16 0.00      0.000000          0      164     getrusage
17 0.00      0.000000          0      28      times
18 -----
19 100.00      0.015278          1695      total

```

SALES1M

```

1
2 [root@rico ~]# strace -c -p 11208
3 Process 11208 attached - interrupt to quit
4 ^CProcess 11208 detached
5 % time      seconds  usecs/call      calls      errors syscall
6 -----
7 65.36      0.010000          13      749      io_getevents
8 34.64      0.005300          7      746      io_submit
9 0.00      0.000000          0       3      read
10 0.00      0.000000          0       2      write
11 0.00      0.000000          0       2      open
12 0.00      0.000000          0       2      close
13 0.00      0.000000          0       1      pread
14 0.00      0.000000          0       1      semctl
15 0.00      0.000000          0       3      fcntl
16 0.00      0.000000          0       1      getrlimit
17 0.00      0.000000          0      164     getrusage
18 0.00      0.000000          0      28      times
19 0.00      0.000000          0       2      semtimedop
20 -----
21 100.00      0.015300          1704      total

```

SALES4M

```
1
2 [root@rico ~]# strace -c -p 11208
3 Process 11208 attached - interrupt to quit
4 ^CProcess 11208 detached
5 % time      seconds  usecs/call   calls   errors syscall
6 -----
7 50.10      0.006000         8     746   io_getevents
8 49.90      0.005975         8     746   io_submit
9 0.00       0.000000         0         2     read
10 0.00       0.000000         0         2     write
11 0.00       0.000000         0         1     pread
12 0.00       0.000000         0        21     getrusage
13 0.00       0.000000         0        12     times
14 -----
15 100.00     0.011975        1530   total
```

So EVERYTHING is being read with ASYNC IO!!! And this not the craziest thing – the other sessions are also using ASYNC IO for those tables now. Without altering the parameter. I have tried to flush `buffer_cache` and `shared_pool`, but with no success. The situation goes back to square one after restarting the instance.

It occurs that it takes only one session to set the parameter to any value – even the same value as it already was by default. And when the first query will use db file scattered read in async mode – all the other sessions, querying tables on ASM will also use `io_getevents` and `io_submit` instead of `pread`. (!!!)

```
1 [oracle@rico ~]$ sqlplus sh/sh
2
3 SQL*Plus: Release 12.1.0.2.0 Production on Fri Jul 17 10:59:22 2015
4
5 Copyright (c) 1982, 2014, Oracle. All rights reserved.
6
7 Last Successful login time: Fri Jul 17 2015 10:56:53 +02:00
8
9 Connected to:
10 Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit Production
11 With the Partitioning, Automatic Storage Management, OLAP, Advanced Analytics
12 and Real Application Testing options
13
14 SQL> alter session set db_file_multiblock_read_count=128;
15
16 Session altered.
17
18 SQL> select count(1) from sales16m;
19
20      COUNT(1)
21 -----
22      19295703
23
24 SQL> Disconnected from Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 -
25 With the Partitioning, Automatic Storage Management, OLAP, Advanced Analytics
26 and Real Application Testing options
27 [oracle@rico ~]$ sqlplus sh/sh
28
29 SQL*Plus: Release 12.1.0.2.0 Production on Fri Jul 17 10:59:51 2015
30
31 Copyright (c) 1982, 2014, Oracle. All rights reserved.
32
33 Last Successful login time: Fri Jul 17 2015 10:59:22 +02:00
34
35 Connected to:
36 Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit Production
37 With the Partitioning, Automatic Storage Management, OLAP, Advanced Analytics
38 and Real Application Testing options
```

```

34
35 SQL> select count(1) from SALES4M;
36
37     COUNT(1)
38     -----
39     19295703
40
41 SQL> Disconnected from Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 -
42 With the Partitioning, Automatic Storage Management, OLAP, Advanced Analytics
43 and Real Application Testing options
44 [oracle@rico ~]$ sqlplus sh/sh
45
46 SQL*Plus: Release 12.1.0.2.0 Production on Fri Jul 17 11:00:43 2015
47
48 Copyright (c) 1982, 2014, Oracle. All rights reserved.
49
50 Last Successful login time: Fri Jul 17 2015 10:59:51 +02:00
51
52 Connected to:
53 Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit Production
54 With the Partitioning, Automatic Storage Management, OLAP, Advanced Analytics
55 and Real Application Testing options
56
57 SQL> select count(1) from sales32m;
58
59     COUNT(1)
60     -----
61     19295703
62
63
64
65
66

```

The appropriate *strace* results are:

```

1 [root@rico ~]# strace -c -p 11782
2 Process 11782 attached - interrupt to quit
3 ^CProcess 11782 detached
4 % time      seconds  usecs/call   calls   errors syscall
5 -----
6 53.97      0.006000      8         747         io_getevents
7 45.60      0.005069      7         747         io_submit
8 0.43       0.000048      1         34         times
9 0.00       0.000000      0         31         read
10 0.00       0.000000      0         2         write
11 0.00       0.000000      0         16         7 open
12 0.00       0.000000      0         9         close
13 0.00       0.000000      0         1         stat
14 0.00       0.000000      0         1         fstat
15 0.00       0.000000      0         28         lseek
16 0.00       0.000000      0         7         mmap
17 0.00       0.000000      0         1         munmap
18 0.00       0.000000      0         1         brk
19 0.00       0.000000      0         4         pread
20 0.00       0.000000      0         1         uname
21 0.00       0.000000      0         10         fcntl
22 0.00       0.000000      0         1         getrlimit
23 0.00       0.000000      0         164        getrusage
24 0.00       0.000000      0         1         getuid
25 -----
26 100.00     0.011117             1806         7 total
27 [root@rico ~]# strace -c -p 11793
28 Process 11793 attached - interrupt to quit
29 ^CProcess 11793 detached
30 % time      seconds  usecs/call   calls   errors syscall

```



```

26 -----
27 51.71 0.007000 10 731 io_getevents
28 48.29 0.006536 9 731 io_submit
29 0.00 0.000000 0 32 read
30 0.00 0.000000 0 2 write
31 0.00 0.000000 0 17 7 open
32 0.00 0.000000 0 10 close
33 0.00 0.000000 0 1 stat
34 0.00 0.000000 0 1 fstat
35 0.00 0.000000 0 28 lseek
36 0.00 0.000000 0 9 mmap
37 0.00 0.000000 0 1 munmap
38 0.00 0.000000 0 1 brk
39 0.00 0.000000 0 17 pread
40 0.00 0.000000 0 1 uname
41 0.00 0.000000 0 1 semctl
42 0.00 0.000000 0 10 fcntl
43 0.00 0.000000 0 1 getrlimit
44 0.00 0.000000 0 209 getrusage
45 0.00 0.000000 0 50 times
46 0.00 0.000000 0 1 getuid
47 0.00 0.000000 0 2 semtimedop
48 -----
49 100.00 0.013536 1856 7 total
50 [root@rico ~]# strace -c -p 11823
51 Process 11823 attached - interrupt to quit
52 ^CProcess 11823 detached
53 % time seconds usecs/call calls errors syscall
54 -----
55 51.65 0.005342 7 731 io_submit
56 48.35 0.005000 7 731 io_getevents
57 0.00 0.000000 0 31 read
58 0.00 0.000000 0 2 write
59 0.00 0.000000 0 16 7 open
60 0.00 0.000000 0 9 close
61 0.00 0.000000 0 1 stat
62 0.00 0.000000 0 1 fstat
63 0.00 0.000000 0 28 lseek
64 0.00 0.000000 0 9 mmap
65 0.00 0.000000 0 1 munmap
66 0.00 0.000000 0 1 brk
67 0.00 0.000000 0 17 pread
68 0.00 0.000000 0 1 uname
69 0.00 0.000000 0 10 fcntl
70 0.00 0.000000 0 1 getrlimit
71 0.00 0.000000 0 209 getrusage
72 0.00 0.000000 0 50 times
73 0.00 0.000000 0 1 getuid
74 -----
75 100.00 0.010342 1850 7 total
76
77
78
79
80

```

Results from *perf* seams to confirm this strange behaviour.

Before setting the parameter:

```
1 [root@rico ~]# perf record -p 13521
```

```

2 ^C[ perf record: Woken up 1 times to write data ]
3 [ perf record: Captured and wrote 0.061 MB perf.data (~2644 samples) ]
4
5 [root@rico ~]# perf report | grep aio
6 0.06% oracle_13521_or [kernel.kallsyms] [k] do_aio_read
7 [root@rico ~]# perf report | grep pread
8 0.06% oracle_13521_or libpthread-2.12.so [.] __pread_nocancel
9

```

After setting the parameter:

```

1 [root@rico ~]# perf record -p 13573
2 ^C[ perf record: Woken up 1 times to write data ]
3 [ perf record: Captured and wrote 0.061 MB perf.data (~2680 samples) ]
4
5 [root@rico ~]# perf report | grep pread
6 [root@rico ~]# perf report | grep aio
7 0.18% oracle_13573_or libaio.so.1.0.1 [.] 0x00000000000000615
8 0.06% oracle_13573_or oracle [.] ksfd aio
9 0.06% oracle_13573_or [kernel.kallsyms] [k] aio_put_req
10 0.06% oracle_13573_or [kernel.kallsyms] [k] aio_run_ioctx

```

I made tests on: Oracle 12.1.0.2 EE (OEL 6.6) and Oracle 11.2.0.3 EE (OEL 6.3)

So the conclusion is: WTF?!

ORACLE 12C: PRAGMA UDF – THE TRUTH

Let's check, why PRAGMA UDF makes execution faster, then regular function.
C function responsible for calling a PL/SQL code from SQL is called "plsqli_run":

```

1 [oracle@rico bin]$ nm oracle | grep plsqli_run
2 0000000000c00000 T plsqli_run

```

Now let's try to create a new HR session and create a simple function which will multiple values by 2:

```

1
2 SQL> conn hr/hr
3 Connected.
4 SQL> get by2
5 1 create or replace
6 2 function by2(x number) return number is
7 3 begin
8 4 return x*2;
9 5* end;
10 SQL> /
11
12 Function created.

```

From other terminal I will connect to the HR session with GDB:

```

1 [root@rico ~]# gdb -p 10330
2 GNU gdb (GDB) Red Hat Enterprise Linux (7.2-83.el6)
3 Copyright (C) 2010 Free Software Foundation, Inc.
4 License GPLv3+: GNU GPL version 3 or later <http://gnu.org/licenses/gpl.html>
5 This is free software: you are free to change and redistribute it.
6 There is NO WARRANTY, to the extent permitted by law. Type "show copying"
7 and "show warranty" for details.
8 This GDB was configured as "x86_64-redhat-linux-gnu".

```

Database Whisperers sp. z o. o. sp. k.

Ul. Płocka 5a
01-231 Warszawa
NIP: 527274487
www.ora-600.pl

```
8 For bug reporting instructions, please see:
9 <http://www.gnu.org/software/gdb/bugs/>.
10 Attaching to process 10330
11 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/bin/oracle...(no debug
12 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libodm12.so...(no
13 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libodm12.so
14 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libcell12.so...don
15 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libcell12.so
16 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libskgxp12.so...(n
17 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libskgxp12.so
18 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libskjcx12.so...do
19 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libskjcx12.so
20 Reading symbols from /lib64/librt.so.1...Reading symbols from /usr/lib/debug/lib64/
21 done.
22 Loaded symbols for /lib64/librt.so.1
23 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libmql1.so...(no d
24 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libmql1.so
25 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libipcl.so...(no d
26 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libipcl.so
27 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libclsra12.so...do
28 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libclsra12.so
29 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libdbcfg12.so...(n
30 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libdbcfg12.so
31 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libhasgen12.so...d
32 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libhasgen12.so
33 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libskgxn2.so...(no
34 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libskgxn2.so
35 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libocr12.so...done
36 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libocr12.so
37 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libocrb12.so...don
38 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libocrb12.so
39 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libocrut12.so...d
40 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libocrut12.so
41 Reading symbols from /lib64/libaio.so.1...Reading symbols from /usr/lib/debug/lib64/
42 done.
43 Loaded symbols for /lib64/libaio.so.1
44 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libons.so...(no de
45 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libons.so
46 Reading symbols from /lib64/libdl.so.2...Reading symbols from /usr/lib/debug/lib64/
47 done.
48 Loaded symbols for /lib64/libdl.so.2
49 Reading symbols from /lib64/libm.so.6...Reading symbols from /usr/lib/debug/lib64/l
50 done.
51 Loaded symbols for /lib64/libm.so.6
52 Reading symbols from /lib64/libpthread.so.0...Reading symbols from /usr/lib/debug/l
53 [Thread debugging using libthread_db enabled]
54 done.
55 Loaded symbols for /lib64/libpthread.so.0
56 Reading symbols from /lib64/libnsl.so.1...Reading symbols from /usr/lib/debug/lib64
57 done.
58 Loaded symbols for /lib64/libnsl.so.1
59 Reading symbols from /lib64/libc.so.6...Reading symbols from /usr/lib/debug/lib64/l
60 done.
61 Loaded symbols for /lib64/libc.so.6
62 Reading symbols from /lib64/ld-linux-x86-64.so.2...Reading symbols from /usr/lib/de
63 done.
64 Loaded symbols for /lib64/ld-linux-x86-64.so.2
65 Reading symbols from /usr/lib64/libnuma.so.1...Reading symbols from /usr/lib/debug/
66 done.
67 Loaded symbols for /usr/lib64/libnuma.so.1
68 Reading symbols from /lib64/libnss_files.so.2...Reading symbols from /usr/lib/debug
69 done.
70 Loaded symbols for /lib64/libnss_files.so.2
71 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libnque12.so...(no
72 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libnque12.so
73 Reading symbols from /u01/app/oracle/product/12.1.0/dbhome_1/lib/libshpkavx12.so...
74 Loaded symbols for /u01/app/oracle/product/12.1.0/dbhome_1/lib/libshpkavx12.so
75
76 warning: no loadable sections found in added symbol-file system-supplied DSO at 0x7
77 0x000000390c00e7c0 in __read_nocancel () at ../sysdeps/unix/syscall-template.S:82
78 82 T_PSEUDO (SYSCALL_SYMBOL, SYSCALL_NAME, SYSCALL_NARGS)
79 (gdb)
80
81
```

Database Whisperers sp. z o. o. sp. k.

Ul. Płocka 5a
01-231 Warszawa
NIP: 527274487
www.ora-600.pl

69
70
71
72
73
74
75
76
77
78
79
80

and I will set the breakpoint for plsql_run function

```
1
2 (gdb) b plsql_run
3 Breakpoint 1 at 0xce1c090
4 (gdb) info breakpoint
5 Num      Type          Disp Enb Address          What
6 1        breakpoint    keep y  0x000000000ce1c090 <plsql_run>
7 (gdb) command 1
8 Type commands for breakpoint(s) 1, one per line.
9 End with a line saying just "end".
10 >c
11 >end
12 (gdb) info b
13 Num      Type          Disp Enb Address          What
14 1        breakpoint    keep y  0x000000000ce1c090 <plsql_run>
15 c
16 (gdb) set pagination off
17 (gdb) c
Continuing.
```

In the HR session, I will use my function on the table "EMPLOYEES".

```
1 SQL> select by2(salary) from employees;
2 (...output removed...)
3 107 rows selected.
4
5 SQL>
```

In the gdb we can see, that the breakpoint was hit 107 times – this is related to the number of rows in the table.

And what will happen if we add "PRAGMA UDF" syntax to our little function?

```
1 SQL> ed
2 Wrote file afiedt.buf
3
4 1 create or replace
5 2 function by2(x number) return number is
6 3 pragma udf;
7 4 begin
8 5 return x*2;
9 6* end;
10 SQL> /
11
12 Function created.
13
14 SQL> select by2(salary) from employees;
15 (...output removed...)
16 107 rows selected.
```

```
16 SQL>
17
18
```

The gdb shows a funny thing:

```
1 (gdb) info b
2 Num      Type          Disp Enb Address          What
3 2        breakpoint    keep y 0x000000000ce1c090 <plsqli_run>
4          breakpoint already hit 9 times
```

So why, the function was called only 9 times?

The tkprof shows something that can be interesting:

```
1
2 select by2(salary)
3 from
4 employees
5
6 call      count      cpu      elapsed      disk      query      current      rows
7 -----
8 Parse     1          0.00     0.00         0          0          0            0
9 Execute   1          0.00     0.00         0          0          0            0
10 Fetch    9          0.00     0.00         0          15         0           107
11 -----
12 total    11         0.00     0.01         0          15         0           107
```

There was 9 fetches! Is it possible that the function was called once for each fetch? We can make sure by changing the arraysize in SQL*Plus

```
1 SQL> set arraysize 66
2 SQL> select by2(salary) from employees;
3 (...output removed...)
4 107 rows selected.
5
6 SQL>
```

The gdb shows three plsqli_run executions:

```
1 (gdb) info b
2 Num      Type          Disp Enb Address          What
3 3        breakpoint    keep y 0x000000000ce1c090 <plsqli_run>
4          breakpoint already hit 3 times
5          c
```

And the tkprof also:

```
1
2 select by2(salary)
3 from
4 employees
5
6 call      count      cpu      elapsed      disk      query      current      rows
7 -----
8 Parse     1          0.00     0.00         0          0          0            0
9 Execute   1          0.00     0.00         0          0          0            0
10 Fetch    3          0.00     0.00         0          9          0           107
11 -----
12 total     5          0.00     0.00         0          9          0           107
```

So the fetch size can have actual impact on PRAGMA UDF functionality. Let's see what will happen from inside of the PL/SQL – without throwing the results to the client:

```
1
2  SQL> ed
3  Wrote file afiedt.buf
4
5      1 declare
6      2 type t_nums is table of number index by pls_integer;
7      3 v_nums t_nums;
8      4 begin
9      5 select by2(salary) bulk collect into v_nums
10     6 from employees;
11     7* end;
12 SQL> /
13
14 PL/SQL procedure successfully completed.
```

gdb shows two plsql_run hits – one for anonymous block and one for the "BY2" function.

CONTEXT SWITCH – PL/SQL CURSOR LOOPS AND FETCHSIZE VS OPIFCH2

Those context switches can be a real pain in the ass – there is a great article by Frits Hoogland about context switching from SQL to PL/SQL – you can read it here: <https://fritshoogland.wordpress.com/2016/01/23/plsql-context-switch/>

You should also read a great blog post by Martin Widlake: Pragma UDF – Some Current Limitations <https://mwidlake.wordpress.com/2015/11/11/pragma-udf-some-current-limitations/>

I also wrote a few words about reducing context switching with PRAGMA UDF: <http://blog.ora-600.pl/2015/10/29/oracle-12c-pragma-udf-the-truth/>

To sum up – since there is no wait event for context switch, you can use for example perf, to verify if the function "plsql_run" has a lot of executions on your system, which can indicate that you have the context switching problem.

But this method can be used only to reveal problems with PL/SQL functions being called from SQL.

In a lots of cases there is a problem with developers not using BULK COLLECT syntax. From Oracle 10g onwards, there is a performance feature in PL/SQL that transforms regular cursor loops to BULK COLLECTs.

Let's check the fetch size of this optimization – the function responsible for fetching the row is called opifch2.

How do we know that? Well, one of my favourite examples of showing the context switching (or fetching) problem is comparing the execution times of those two anonymous blocks:

```
1  SQL> declare
2      x number;
3  begin
4      for i in 1..1000000 loop
5          select 1 into x from dual;
6      end loop;
7  end;
8  /
```

```

7   PL/SQL procedure successfully completed.
8
9   Elapsed: 00:00:20.08
10
11  SQL> ed
12  Wrote file afiedt.buf
13
14  1  declare
15  2    x number;
16  3  begin
17  4    for i in 1..1000000 loop
18  5      x:=1;
19  6    end loop;
20  7* end;
21  SQL> /
22
23  PL/SQL procedure successfully completed.
24
25  Elapsed: 00:00:00.01
26
27

```

With perf we can see the top function calls while executing the first anonymous block:

```

1
2  # Overhead          Command          Shared Object          Symbol
3  # .....
4  #
5  20.51% oracle_3902_orc [vdso]                [.] 0x0000000000000bb0
6  3.28%  oracle_3902_orc oracle                [.] opiexe
7  2.46%  oracle_3902_orc [kernel.kallsyms]    [k] __audit_syscall_entry
8  2.21%  oracle_3902_orc [kernel.kallsyms]    [k] system_call_after_swapgs
9  1.83%  oracle_3902_orc oracle                [.] opifch2

```

The function opiexe is responsible for executing some stuff – in this example: the select statement for 1 000 000 times.

The opifch2 is the function which we will be examining, since this is the function responsible for fetching.

After some investigation, I've noticed that the number of rows fetched by this function is



stored in RCX registry – let's try to prove it

I will use a simple "bulk collect limit" syntax to change the number of rows being fetched. To display the RCX value I will use GDB with breakpoint at opifch2 function.

Preparing GDB session:

- get SPID info about your session:

```

1  SQL> get myspid
2  1  select spid
3  2  from v$session s, v$process p
4  3  where s.paddr=p.addr
5  4* and  s.sid=sys_context('userenv','sid')
6  SQL> /
7
8  SPID
9  -----
10 5543

```

9
10

- start GDB and prepare a breakpoint:

```
1
2 [root@rico ~]# gdb -p 5543
3 (gdb) b opifch2
4 Breakpoint 1 at 0xcb0acb0
5 (gdb) command 1
6 Type commands for breakpoint(s) 1, one per line.
7 End with a line saying just "end".
8 >p/d $rcx
9 >c
10 >end
11 (gdb) set pagination off
12 (gdb) c
Continuing.
```

Now we will execute anonymous PL/SQL block and check our RCX value for LIMIT 25, 59, 123

- LIMIT 25

PL/SQL block:

```
1
2
3 SQL> ;
4   1 declare
5   2   cursor c_sql is
6   3   select salary from employees;
7   4   type t_sql is table of c_sql%ROWTYPE index by pls_integer;
8   5   v_sql t_sql;
9   6   v_x   number;
10  7   begin
11  8   open c_sql;
12  9   loop
13 10   fetch c_sql bulk collect into v_sql limit 25;
14 11   exit when c_sql%NOTFOUND;
15 12   for i in v_sql.first..v_sql.last loop
16 13   v_x:=v_x+v_sql(i).salary;
17 14   end loop;
18 15   end loop;
19 16   close c_sql;
20 17* end;
21 SQL> /
```

GDB:

```
1 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
2 $1 = 140159623635016
3
4 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
5 $2 = 25
6
7 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
8 $3 = 25
9
10 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
11 $4 = 25
12
13 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
14 $5 = 25
```



```
14 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
15 $6 = 25
16
17
```

- LIMIT 59

PL/SQL block:

```
1
2
3 SQL> ed
4 Wrote file afiedt.buf
5
6 1 declare
7 2 cursor c_sql is
8 3 select salary from employees;
9 4 type t_sql is table of c_sql%ROWTYPE index by pls_integer;
10 5 v_sql t_sql;
11 6 v_x number;
12 7 begin
13 8 open c_sql;
14 9 loop
15 10 fetch c_sql bulk collect into v_sql limit 59;
16 11 exit when c_sql%NOTFOUND;
17 12 for i in v_sql.first..v_sql.last loop
18 13 v_x:=v_x+v_sql(i).salary;
19 14 end loop;
20 15 end loop;
21 16 close c_sql;
22 17* end;
23 SQL> /
24
25 PL/SQL procedure successfully completed.
26
27
28
```

GDB:

```
1 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
2 $7 = 59
3
4 Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
5 $8 = 59
```

- LIMIT 123

PL/SQL block:

```
1 SQL> ed
2 Wrote file afiedt.buf
3
4 1 declare
5 2 cursor c_sql is
6 3 select salary from employees;
7 4 type t_sql is table of c_sql%ROWTYPE index by pls_integer;
8 5 v_sql t_sql;
9 6 v_x number;
10 7 begin
11 8 open c_sql;
12 9 loop
13 10 fetch c_sql bulk collect into v_sql limit 123;
14 11 exit when c_sql%NOTFOUND;
15 12 for i in v_sql.first..v_sql.last loop
16 13 v_x:=v_x+v_sql(i).salary;
17 14 end loop;
18 15 end loop;
19 16 close c_sql;
```

```

16      17* end;
17      SQL> /
18
19      PL/SQL procedure successfully completed.
20
21
22
23

```

GDB:

```

1      Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
2      $9 = 123

```

Great! So now we have proved, that the number of rows fetched by opifch2 is stored at RCX.

Now we can check how many rows are being fetched by standard cursor loop optimization:

PL/SQL block:

```

1
2
3      SQL> ;
4      1 declare
5      2   cursor c_sql is
6      3   select salary from employees;
7      4   type t_sql is table of c_sql%ROWTYPE index by pls_integer;
8      5   v_sql t_sql;
9      6   v_x   number;
10     7   begin
11     8   for i in c_sql loop
12     9   v_x:=v_x+i.salary ;
13    10   end loop;
14    11* end;
15    SQL> /
16
17    PL/SQL procedure successfully completed.
18
19

```

GDB:

```

1      Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
2      $10 = 100
3
4      Breakpoint 1, 0x00000000cb0acb0 in opifch2 ()
5      $11 = 100

```

I made tests for different table sizes and each time it was exactly 100 rows – so this is fixed value and it doesn't depend on statistics.

I think that main problem is, that you can't figure it out based on regular wait events or statistics (or at least I don't know any).

If you check active sessions with a lot of cursor loops being executed on big tables, you will see something like this:

```

1      Oracle 12c - orc 19:49:45 up: 6.4h, 1 ins, 2 sn, 1 us, 2.1G mt, 66.0% db
2      ID %CPU LOAD %DCU AAS ASC ASI ASW AST IOPS %FR PGA UTPS UCPS SSRT %DBT
3      1 16 0 15 0.7 1 1 0 1 5 7 82M 0 2 41u 100
4
5      EVENT (RT) TOT WAITS TIME(s) AVG_MS PCT WAIT_CLASS

```

```

5      DB CPU                                7          100
6      direct path read                      37          0          0          0      User I/O
7
8      ID      SID      SPID  USR  PROG S   PGA  SQLID/BLOCKER  OPN  E/T  STA  STE  EVENT/*LA  W/T
9      1       48      5543 HR   sqlp D 3.1M 6ucj0mtv5bf9b SEL 18s ACT CPU direct pa 3m
10
11      Enter sql_id: 6ucj0mtv5bf9b
12
13      PLAN_TABLE_OUTPUT
14      -----
15      SQL_ID: 6ucj0mtv5bf9b, child number 0
16      -----
17      SELECT AMOUNT_SOLD FROM SH.SALES2
18
19      Plan hash value: 497282378
20
21      -----
22      |   Id | Operation                               | Name          | Rows | Cost | Stale |
23      |-----|-----|-----|-----|-----|-----|
24      |    0 | SELECT STATEMENT                         |                |      |    13k |      |
25      |    1 | TABLE ACCESS FULL                       | SALES2        | 10M  |    13k | YES  |
26      |-----|-----|-----|-----|-----|-----|
27

```

Fortunately you can use operating system tools like perf – you should check out this amazing article by Luca Canali about Linux Perf Probes for Oracle Tracing:<http://externaltable.blogspot.com/2016/02/linux-perf-probes-for-oracle-tracing.html>

Based on Luca's findings we can create a perf probe for opifch2 function and check out fetch size for all server processes on our operating system:

```

1
2      [root@rico ~]# perf probe -x /u01/app/oracle/product/12.1.0/dbhome_1/bin/oracle opi
3      Added new event:
4          probe_oracle:opifch2 (on 0xc70acb0 with fetchsize=%cx)
5
6      You can now use it in all perf tools, such as:
7
8          perf record -e probe_oracle:opifch2 -aR sleep 1
9      [root@rico ~]# perf record -e probe_oracle:opifch2 -aR
10      ^C[ perf record: Woken up 149 times to write data ]
11      [ perf record: Captured and wrote 37.395 MB perf.data (~1633825 samples) ]
12
13      [root@rico ~]# perf script > opifch2.out

```

After recording events for some time we can process the output with a very simple AWK oneliner:

```

1      [root@rico ~]# awk '/oracle_/ {fsize=strtonum("0x"substr($NF,11));fsize_tab[$1 " was " f
2      " times)}' opifch2.out | sort -n -k7
3      Fetchsize for oracle_5857_orc was 2 (executed 1 times)
4      Fetchsize for oracle_5857_orc was 2147483647 (executed 1 times)
5      Fetchsize for oracle_5857_orc was 100 (executed 2 times)
6      Fetchsize for oracle_5857_orc was 140171241660488 (executed 2 times)
7      Fetchsize for oracle_5857_orc was 1 (executed 4 times)
8      Fetchsize for oracle_5857_orc was 227183808 (executed 4 times)
9      Fetchsize for oracle_5857_orc was 25 (executed 5 times)
10     Fetchsize for oracle_5755_orc was 2 (executed 22 times)
11     Fetchsize for oracle_5755_orc was 140058172228192 (executed 36 times)
12     Fetchsize for oracle_5543_orc was 100 (executed 441045 times)

```

The high number of executions for fetchsize 100 suggests that your system is using a lot of cursor loops with auto optimization – but without actual BULK COLLECT syntax.

Please notice, that this method can reveal also wrong fetch sizes from any other client – like JDBC or .NET.