

# Oracle and storage IOs, explanations and experience at CERN CHEP 2009 [paper 28]

Eric Grancher  
eric.grancher@cern.ch  
CERN IT



- Logical and Physical IO
- Measuring IO
- Exadata
- SSD
- Conclusions
- References

- Even so memory access is fast compared to disk access, LIO are actually expensive
- LIO cost latching and CPU
- Tuning using LIO reduction as a reference is advised
  
- See “Why You Should Focus on LIOs Instead of PIOs” Carry Millsap

- One has to measure “where the PIO are performed” and “how long they take / how many per second are performed”
- Oracle instrumentation and counters provide us the necessary information, raw and aggregated
- Counters:
  - Aggr file: V\$FILESTAT / DBA\_HIST\_FILESTATXS (\*), V\$FILE\_HISTOGRAM
  - Aggr session: V\$SESS\_IO
  - Aggr system-wide: V\$SYSSTAT

- Individual wait events:

- 10046 event or EXECUTE

```
DBMS_MONITOR.SESSION_TRACE_ENABLE(83,5,  
TRUE, FALSE); then EXECUTE
```

```
DBMS_MONITOR.SESSION_TRACE_DISABLE(83,5);
```

- Trace file contains lines like:

```
WAIT #5: nam='db file sequential read'  
  ela=6784 file#=6 block#=467667 blocks=1  
  obj#=73442 tim=1490530491532
```

- Session wait: V\$SESSION\_WAIT (V\$SESSION 10.1+)

- Aggregated wait events:

- Aggr session: V\$SESSION\_EVENT

- Aggr system-wide: V\$SYSTEM\_EVENT

## – Statspack/AWR(\*) reports

```
SQL> execute dbms_workload_repository.create_snapshot;
...
SQL> @?/rdbms/admin/awrrpt
```

Top 5 Timed Foreground Events

~~~~~

| Event                        | Waits  | Time(s) | Avg wait (ms) | % D time | Wait Class |
|------------------------------|--------|---------|---------------|----------|------------|
| db file sequential read      | 17,049 | 122     | <b>7</b>      | 94.2     | User I/O   |
| DB CPU                       |        | 7       |               | 5.3      |            |
| log file sync                | 4      | 2       | 570           | 1.8      | Commit     |
| db file scattered read       | 21     | 0       | 1             | .0       | User I/O   |
| control file sequential read | 694    | 0       | 0             | .0       | System I/O |

```
^LWait Event Histogram
DB/Inst: ORCL/orcl Snaps: 367-368
-> Units for Total Waits column: K is 1000, M is 1000000, G is 1000000000
-> % of Waits: value of .0 indicates value was <.05%. Value of null is truly 0
-> % of Waits: column heading of <=1s is truly <1024ms, >1s is truly >=1024ms
-> Ordered by Event (idle events last)
```

| Event                      | Total Waits | % of Waits |      |      |      |       |       |      |     |
|----------------------------|-------------|------------|------|------|------|-------|-------|------|-----|
|                            |             | <1ms       | <2ms | <4ms | <8ms | <16ms | <32ms | <=1s | >1s |
| SQL*Net message to client  | 1           | 100.0      |      |      |      |       |       |      |     |
| control file parallel writ | 40          |            |      |      |      |       | 55.0  | 45.0 |     |
| control file sequential re | 758         | 100.0      |      |      |      |       |       |      |     |
| db file parallel write     | 111         |            |      | 1.8  | 3.6  | 27.9  | 31.5  | 35.1 |     |
| db file scattered read     | 22          | 95.5       |      |      |      |       | 4.5   |      |     |
| db file sequential read    | 16K         | 3.5        | .0   | 4.5  | 58.1 | 32.9  | .9    | .1   |     |

# How to measure IO (4/4)

^LTablespace IO Stats DB/Inst: ORCL/orcl Snaps: 367-368  
-> ordered by IOs (Reads + Writes) desc

| Tablespace | Av Reads | Av Reads/s | Av Rd(ms)  | Av Blks/Rd | Writes | Av Writes/s | Buffer Waits | Av Buf Wt (ms) |
|------------|----------|------------|------------|------------|--------|-------------|--------------|----------------|
| TESTTBS    | 16,323   | 131        | <b>7.5</b> | 1.0        | 0      | 0           | 0            | 0.0            |
| SYSAUX     | 614      | 5          | 0.0        | 1.0        | 121    | 1           | 0            | 0.0            |
| SYSTEM     | 221      | 2          | 0.0        | 1.6        | 7      | 0           | 0            | 0.0            |
| UNDOTBS1   | 1        | 0          | 0.0        | 1.0        | 17     | 0           | 0            | 0.0            |

^LFile IO Stats DB/Inst: ORCL/orcl Snaps: 367-368  
-> ordered by Tablespace, File

| Tablespace | Av Reads | Av Reads/s | Av Rd(ms)  | Av Blks/Rd | Writes | Av Writes/s | Buffer Waits | Av Buf Wt (ms) | Filename                                             |
|------------|----------|------------|------------|------------|--------|-------------|--------------|----------------|------------------------------------------------------|
| SYSAUX     | 614      | 5          | 0.0        | 1.0        | 121    | 1           | 0            | 0.0            | /export/home/oracle/app/oracle/oradata/orcl/sysaux01 |
| SYSTEM     | 221      | 2          | 0.0        | 1.6        | 7      | 0           | 0            | 0.0            | /export/home/oracle/app/oracle/oradata/orcl/system01 |
| TESTTBS    | 8,001    | 64         | <b>7.5</b> | 1.0        | 0      | 0           | 0            | 0.0            | /ORA/orcl_testtbs.dbf                                |
| TESTTBS    | 8,322    | 67         | 7.5        | 1.0        | 0      | 0           | 0            | 0.0            | /ORA2/orcl_testtbs2.dbf                              |
| UNDOTBS1   | 1        | 0          | 0.0        | 1.0        | 17     | 0           | 0            | 0.0            | /export/home/oracle/app/oracle/oradata/orcl/undotbs0 |

- Using ASH(\*)
  - Sampling of session information every 1s
  - Not biased (just time sampling), so reliable source of information
  - Obviously not all information is recorded so some might be missed
- Can be accessed via
  - @ashrpt / @ashrpti
  - v\$active\_session\_history /  
DBA\_HIST\_ACTIVE\_SESS\_HISTORY(\*)



# ASH and IO (2/2)

```
SQL> EXECUTE DBMS_MONITOR.SESSION_TRACE_ENABLE(69,17062, TRUE, FALSE);
```

PL/SQL procedure successfully completed.

```
SQL> select to_char(sample_time,'HH24MISS') ts,seq#,p1,p2,time_waited from v$active_session_history where SESSION_ID=
69 and session_serial#=17062
 2 and SESSION_STATE = 'WAITING' and event='db file sequential read' and sample_time>sysdate -5/24/3600
 3 order by sample_time;
```

| TS     | SEQ#  | P1 | P2            | TIME_WAITED  |
|--------|-------|----|---------------|--------------|
| 001557 | 45565 | 6  | 449426        | 5355         |
| 001558 | 45716 | 6  | <b>179376</b> | <b>10118</b> |
| 001559 | 45862 | 6  | 702316        | 7886         |
| 001600 | 46014 | 7  | 91988         | 5286         |
| 001601 | 46167 | 7  | 424665        | 7594         |
| 001602 | 46288 | 6  | <b>124184</b> | <b>0</b>     |

```
SQL> EXECUTE DBMS_MONITOR.SESSION_TRACE_DISABLE(69,17062);
```

PL/SQL procedure successfully completed.

```
-bash-3.00$ grep -n 124184 orcl_ora_15854.trc
676:WAIT #2: nam='db file sequential read' ela= 5355 file#=6 block#=449426 blocks=1 obj#=73442 tim=2707602560910
[...]
829:WAIT #2: nam='db file sequential read' ela= 10118 file#=6 block#=179376 blocks=1 obj#=73442 tim=2707603572300
[...]
977:WAIT #2: nam='db file sequential read' ela= 7886 file#=6 block#=702316 blocks=1 obj#=73442 tim=2707604583489
[...]
1131:WAIT #2: nam='db file sequential read' ela= 5286 file#=7 block#=91988 blocks=1 obj#=73442 tim=2707605593626
[...]
1286:WAIT #2: nam='db file sequential read' ela= 7594 file#=7 block#=424665 blocks=1 obj#=73442 tim=2707606607137
[...]
1409:WAIT #2: nam='db file sequential read' ela= 8861 file#=6 block#=124184 blocks=1 obj#=73442 tim=2707607617211
```

- How to identify which segments are accessed most often from a given session? (ashrpti can do it as well)
- Ultimate information is in a 10046 trace
- Extract necessary information, load into t(p1,p2)

```
> grep "db file sequential read" accmeas2_j004_32116.trc | head -2
WAIT #12: nam='db file sequential read' ela= 11175 file#=13 block#=200041
blocks=1 obj#=67575 tim=1193690114589134
WAIT #12: nam='db file sequential read' ela= 9454 file#=6 block#=587915
blocks=1 obj#=67577 tim=1193690114672648
```

```
accmeas_2 bdump > grep "db file sequential read" accmeas2_j004_32116.trc
| head -2 | awk '{print $9"="$10}' | awk -F= '{print $2","$4}'
13,200041
6,587915
```

```
SQL> select distinct
e.owner,e.segment_name,e.PARTITION_NAME,(e.bytes/1024/1024) size_MB from
t, dba_extents e where e.file_id=t.p1 and t.p2 between e.block_id and
e.block_id+e.blocks order by e.owner,e.segment_name,e.PARTITION_NAME;
```

- Take information from v\$active\_session\_history

```
create table t as select p1,p2 from v$active_session_history h where  
h.module like 'DATA_LOAD%' and h.action like 'COLLECT_DN%' and  
h.event ='db file sequential read' and h.sample_time>sysdate-4/24;
```

```
SQL> select distinct  
e.owner,e.segment_name,e.PARTITION_NAME, (e.bytes/1024/1024) size_MB from  
t, dba_extents e where e.file_id=t.p1 and t.p2 between e.block_id and  
e.block_id+e.blocks order by e.owner,e.segment_name,e.PARTITION_NAME;
```

# ashrpti and DB objects

```
CSR:SQL> select user_id from dba_users where username='CINBAD';
```

```
USER_ID
```

```
-----
```

**55**

```
CSR:SQL> select event,user_id,session_id,session_serial#,to_char(SAMPLE_TIME,'YYYYMMDD-  
HH24MISS') from v$active_session_history where SAMPLE_TIME>sysdate-2/24 and user_id=55 and  
event is not null order by sample_time;
```

```
EVENT                                USER_ID  SESSION_ID  SESSION_SERIAL#  TO_CHAR(SAMPLE_  
-----
```

|                         |           |            |       |                 |
|-------------------------|-----------|------------|-------|-----------------|
| [...]                   |           |            |       |                 |
| db file sequential read | 55        | 530        | 28609 | 20081130-225208 |
| db file sequential read | <b>55</b> | <b>530</b> | 28609 | 20081130-230708 |

```
CSR:SQL> @ashrpti
```

```
[...]
```

```
Specify SESSION_ID (eg: from V$SESSION.SID) report target:
```

```
Defaults to NULL:
```

```
Enter value for target_session_id: 530
```

```
[...]
```

## Top DB Objects

- With respect to Application, Cluster, User I/O and buffer busy waits only.

| Object ID | % Activity | Event                   | % Event | Object Name (Type)           | Tablespace |
|-----------|------------|-------------------------|---------|------------------------------|------------|
| 101841    | 2.38       | db file sequential read | 2.38    | CINBAD.SFL_IP_PK (INDEX)     | INDX01     |
| 101847    | 2.38       | db file sequential read | 2.38    | CINBAD.SFL_IP_UDP_PK (INDEX) | INDX01     |

- First identify how the IO is performed:
  - DstackProf (Tanel Poder)
  - strace (Linux) / truss (Solaris)
  - Dtruss
  - DTrace (example later)

# DStackProf example

```
-bash-3.00$ ./dstackprof.sh 11073
```

```
DStackProf v1.02 by Tanel Poder ( http://www.tanelpoder.com )  
Sampling pid 11073 for 5 seconds with stack depth of 100 frames...
```

```
[...]
```

```
780 samples with stack below
```

```
-----  
libc.so.1`_pread
```

```
skgfqio
```

```
ksfd_skgfqio
```

```
ksfd_io
```

```
ksfdread1
```

```
kcfxbd
```

```
kcbzib
```

```
kcbgtcr
```

```
operations used by capabilities such as direct load, has clusters , etc.
```

```
ktrget2
```

```
kdsgrp
```

```
existing row data
```

```
qetlbr
```

```
qertbFetchByRowID qertb - table row source
```

```
qerjotRowProc
```

```
qerjo - row source: join
```

```
kdstf0000001000kmP
```

```
kdsttgr
```

```
kds: operations on data such as retrieving a row and updating
```

```
existing row data
```

```
qertbFetch
```

```
qertb - table row source
```

```
qerjotFetch
```

```
qerjo - row source: join
```

```
qergsFetch
```

```
qergs - group by sort row source
```

```
opifch2
```

```
Kpoal8 / opiodr / ttcpip/ opitsk / opiino / opiodr / opidrv / sou2o / a.out`main / a.out`_start
```

**note 175982.1**

ksfd: support for various kernel associated capabilities manages and coordinates operations on the control file(s)

kcb: manages Oracle's buffer cache operation as well as operations used by capabilities such as direct load, has clusters , etc.

ktr - kernel transaction read consistency

kds: operations on data such as retrieving a row and updating

qertb - table row source

qerjo - row source: join

kds: operations on data such as retrieving a row and updating

qertb - table row source

qerjo - row source: join

qergs - group by sort row source

- BTW you can get the status of execution this way, instead of looking at note 175982.1, you can use Tanel Poder's `os_explain.sh`

```
-bash-3.00$ pstack 11073| ./os_explain.sh
kpoal8
SELECT FETCH:
GROUP BY SORT: Fetch
NESTED LOOP JOIN: Fetch
TABLE ACCESS: Fetch
kdsttgr
kdstf0000001000kmP
NESTED LOOP JOIN: RowProc
TABLE ACCESS: FetchByRowID
getlbr
kdsgrp
ktrget2
kcbgtcr
kcbzib
kcfrbd
ksfdread1
ksfd_io
ksfd_skgfqio
skgfqio
_pread
```

```
SQL> select * from table(dbms_xplan.display_cursor
('8st7asbzt4jz7',null,'BASIC'));
[...]
```

```
PLAN_TABLE_OUTPUT
```

|  |   |  |                            |  |            |  |
|--|---|--|----------------------------|--|------------|--|
|  | 1 |  | SORT AGGREGATE             |  |            |  |
|  | 2 |  | NESTED LOOPS               |  |            |  |
|  | 3 |  | TABLE ACCESS FULL          |  | PROBETEST1 |  |
|  | 4 |  | TABLE ACCESS BY USER ROWID |  | TESTTABLE1 |  |

- You can measure (with the least overhead), selecting only the syscalls that you need
- For example, pread

```
-bash-3.00$ truss -t pread -Dp 17924
/1:      0.0065 pread(258, "06A2\0\001CA9EE1\0 !B886".., 8192, 0x153DC2000) = 8192
/1:      0.0075 pread(257, "06A2\0\0018CFEE4\0 !C004".., 8192, 0x19FDC8000) = 8192
/1:      0.0078 pread(258, "06A2\0\001C4CEE9\0 !92AA".., 8192, 0x99DD2000) = 8192
/1:      0.0103 pread(257, "06A2\0\00188 S F\0 !A6C9".., 8192, 0x10A68C000) = 8192
/1:      0.0072 pread(257, "06A2\0\0018E kD7\0 !CFC2".., 8192, 0x1CD7AE000) = 8192
```

```
-bash-3.00$ truss -t pread -Dp 15854 2>&1 | awk '{s+=$2; if (NR%1000==0) {print NR " " s " " s/NR}}'
```

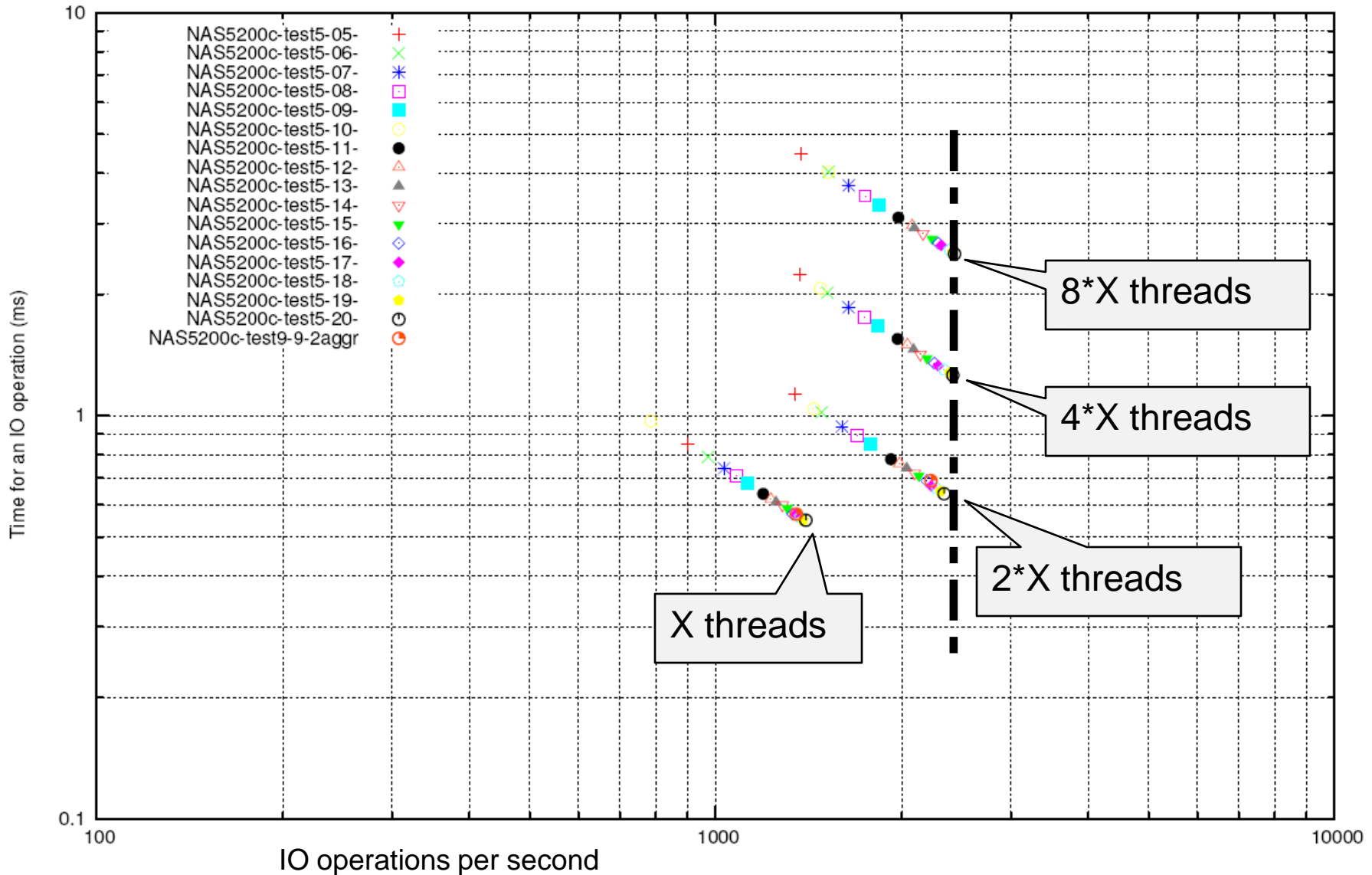
|      |         |            |
|------|---------|------------|
| 1000 | 7.6375  | 0.0076375  |
| 2000 | 15.1071 | 0.00755355 |
| 3000 | 22.4648 | 0.00748827 |



- Each (spinning) disk is capable of ~ 100 to 300 IO operations per second depending on the speed and controller capabilities
- Putting many requests at the same time from the Oracle layer, makes as if IO takes longer to be serviced

# Overload at disk driver level / system level (2/2)

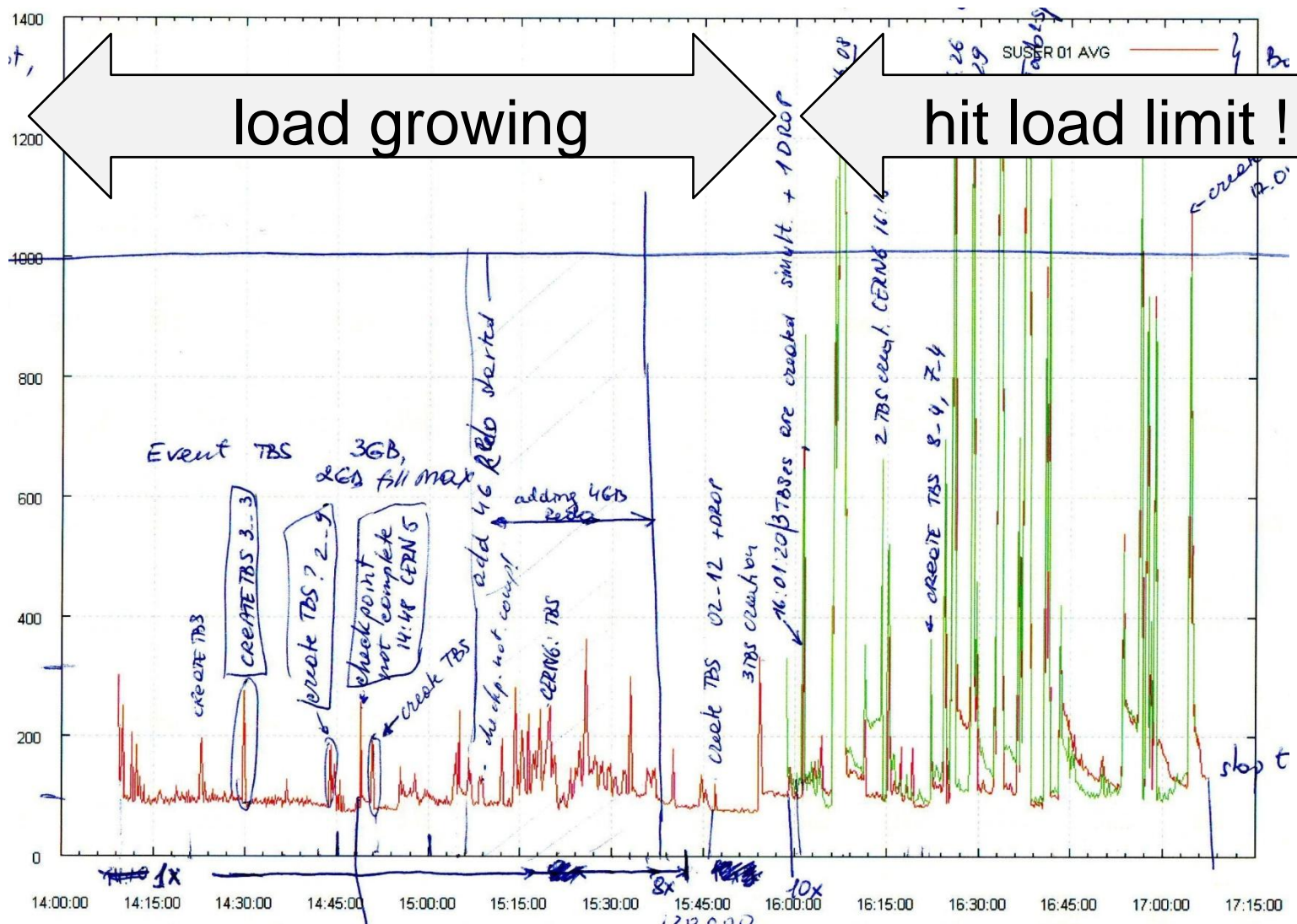
Time of an IO versus IO operations per second  
Oracle database 8kB random block read operations:



- Observed many times: “the storage is slow” (and storage administrators/specialists say “storage is fine / not loaded”)
- Typically happens that observed (from Oracle rdbms point of view) IO wait times are long if CPU load is high
- Instrumentation / on-off cpu

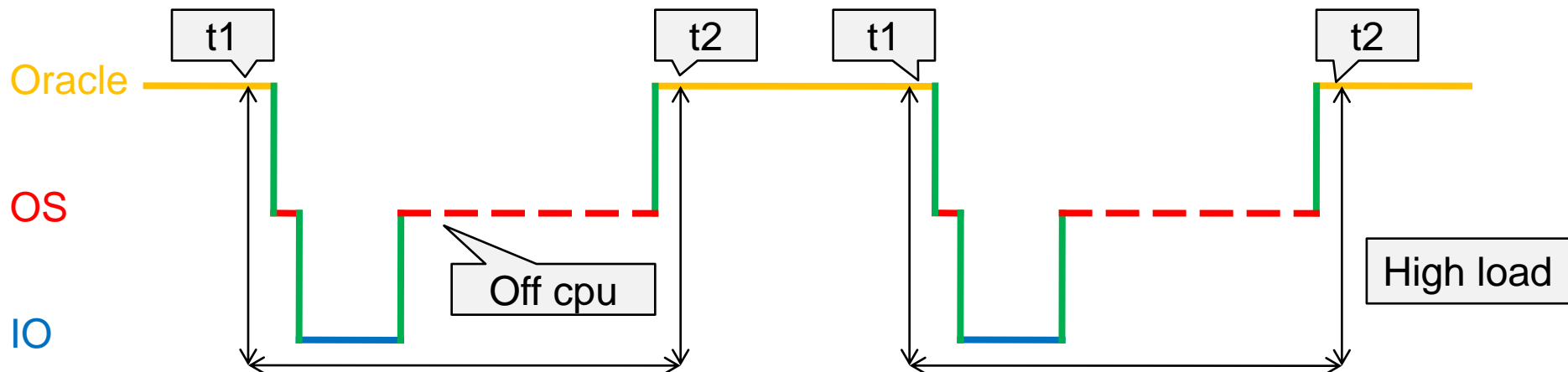
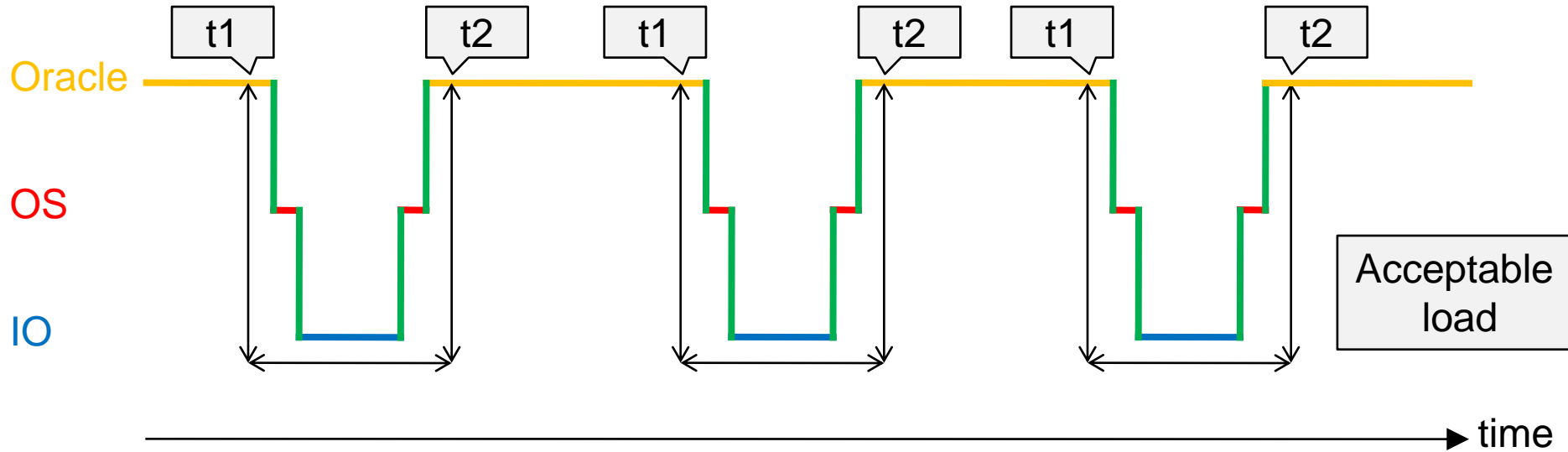
# Overload at CPU level (2/) example

Insertion time (ms), has to be less than 1000ms



15k... 30k ... 60k... 90k... 120k ...135k... || 150k (insertions per second)

# OS level / high-load



- Dtrace (Solaris) can be used at OS level to get (detailed) information at OS level

```
syscall::pread:entry
/pid == $target && self->traceme == 0 /
{
    self->traceme = 1;
    self->on = timestamp;
    self->off= timestamp;
    self->io_start=timestamp;
}
```

```
syscall::pread:entry
/self->traceme == 1 /
{
    self->io_start=timestamp;
}
```

```
syscall::pread:return
/self->traceme == 1 /
{
    @avgs["avg_io"] = avg(timestamp-self->io_start);
    @[tid,"time_io"] = quantize(timestamp-self->io_start);
    @counts["count_io"] = count();
}
```

```
sched:::on-cpu
/pid == $target && self->traceme == 1 /
{
    self->on = timestamp;
    @[tid,"off-cpu"] = quantize(self->on - self->off);
    @totals["total_cpu_off"] = sum(self->on - self->off);
    @avgs["avg_cpu_off"] = avg (self->on - self->off);
    @counts["count_cpu_on"] = count();
}
sched:::off-cpu
/self->traceme == 1/
{
    self->off= timestamp;
    @totals["total_cpu_on"] = sum(self->off - self->on);
    @avgs["avg_cpu_on"] = avg(self->off - self->on);
    @[tid,"on-cpu"] = quantize(self->off - self->on);
    @counts["count_cpu_off"] = count();
}

tick-1sec
/i++ >= 5/
{
    exit(0);
}
```

# Dtrace, “normal load”

```
-bash-3.00$ sudo ./cpu.d4 -p 15854
dtrace: script './cpu.d4' matched 7 probes
CPU      ID          FUNCTION:NAME
  3  52078          :tick-1sec

avg_cpu_on          169114
avg_cpu_off        6768876
avg_io             6850397
```

[...]

```
1  off-cpu
   value  ----- Distribution ----- count
   524288 |                                0
  1048576 |                                2
  2097152 |@@@@                               86
  4194304 |@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@ 577
  8388608 |@@@@@@@@@@@@@                       189
 16777216 |                                2
 33554432 |                                0
```

[...]

```
count_cpu_on          856
count_io              856
count_cpu_off         857
total_cpu_on          144931300
total_cpu_off        5794158700
```



# Dtrace, "high load"

```
-bash-3.00$ sudo ./cpu.d4 -p 15854
dtrace: script './cpu.d4' matched 7 probes
CPU      ID          FUNCTION:NAME
  2  52078          :tick-1sec

avg_cpu_on                210391
avg_cpu_off             10409057
avg_io                    10889597
[...]
```

| 1 | off-cpu  |                          | count |
|---|----------|--------------------------|-------|
|   | value    | ----- Distribution ----- |       |
|   | 8192     |                          | 0     |
|   | 16384    |                          | 4     |
|   | 32768    | @                        | 11    |
|   | 65536    |                          | 2     |
|   | 131072   |                          | 0     |
|   | 262144   |                          | 0     |
|   | 524288   |                          | 0     |
|   | 1048576  |                          | 0     |
|   | 2097152  | @                        | 15    |
|   | 4194304  | @@@@@@@@@@@@@@@@         | 177   |
|   | 8388608  | @@@@@@@@@@@@@@@@         | 249   |
|   | 16777216 | @@@                      | 41    |
|   | 33554432 |                          | 4     |
|   | 67108864 |                          | 0     |

```
[...]
```

|               |            |
|---------------|------------|
| count_io      | 486        |
| count_cpu_on  | 503        |
| count_cpu_off | 504        |
| total_cpu_on  | 106037500  |
| total_cpu_off | 5235756100 |

- Exadata has a number of offload features, most published about are row selection and column selection
- Some of our workloads are data insertion intensive, for these the tablespace creation is/can be a problem
- Additional load, additional IO head moves, additional bandwidth usage on the connection server→storage
- Exadata has file creation offloading
- Tested recently with 4 Exadata cells storage. Tests done with Anton Topurov / Ela Gajewska-Dendek

# Swingbench in action

Applications Places System oracle Thu Aug 28, 7:54 PM

### SwingBench 2.3.0.381 (SWPVSS1)

Time Remaining : 0:00:00

Users: 15

Transactions per Minute: 486

Transactions per Second: 10

CPU: 0

Disk Activity: 0

| Property           | Value                               |
|--------------------|-------------------------------------|
| Benchmark Name     | "PVSS Benchmark"                    |
| Connect String     | SWPVSS1                             |
| Coordinator        |                                     |
| Driver Type        | Oracle10g Type II jdbc driver (oci) |
| Maximum Think Time | 0                                   |
| Minimum Think Time | 0                                   |
| Query Time Out     | 600                                 |
| User Count         | 15                                  |
| User Name          | SUSER01                             |

### SwingBench 2.3.0.381 (SWPVSS2)

Time Remaining : 0:00:00

Users: 15

Transactions per Minute: 420

Transactions per Second: 6

CPU: 0

Disk Activity: 0

| Property           | Value                               |
|--------------------|-------------------------------------|
| Benchmark Name     | "PVSS Benchmark"                    |
| Connect String     | SWPVSS2                             |
| Coordinator        |                                     |
| Driver Type        | Oracle10g Type II jdbc driver (oci) |
| Maximum Think Time | 0                                   |
| Minimum Think Time | 0                                   |
| Query Time Out     | 600                                 |
| User Count         | 15                                  |
| User Name          | SUSER03                             |

### SwingBench 2.3.0.381 (SWPVSS1)

Time Remaining : 0:00:00

Users: 15

Transactions per Minute: 488

Transactions per Second: 13

CPU: 0

Disk Activity: 0

| Property           | Value                               |
|--------------------|-------------------------------------|
| Benchmark Name     | "PVSS Benchmark"                    |
| Connect String     | SWPVSS1                             |
| Coordinator        |                                     |
| Driver Type        | Oracle10g Type II jdbc driver (oci) |
| Maximum Think Time | 0                                   |
| Minimum Think Time | 0                                   |
| Query Time Out     | 600                                 |
| User Count         | 15                                  |
| User Name          | SUSER02                             |

### SwingBench 2.3.0.381 (SWPVSS2)

Time Remaining : 0:00:00

Users: 15

Transactions per Minute: 432

Transactions per Second: 9

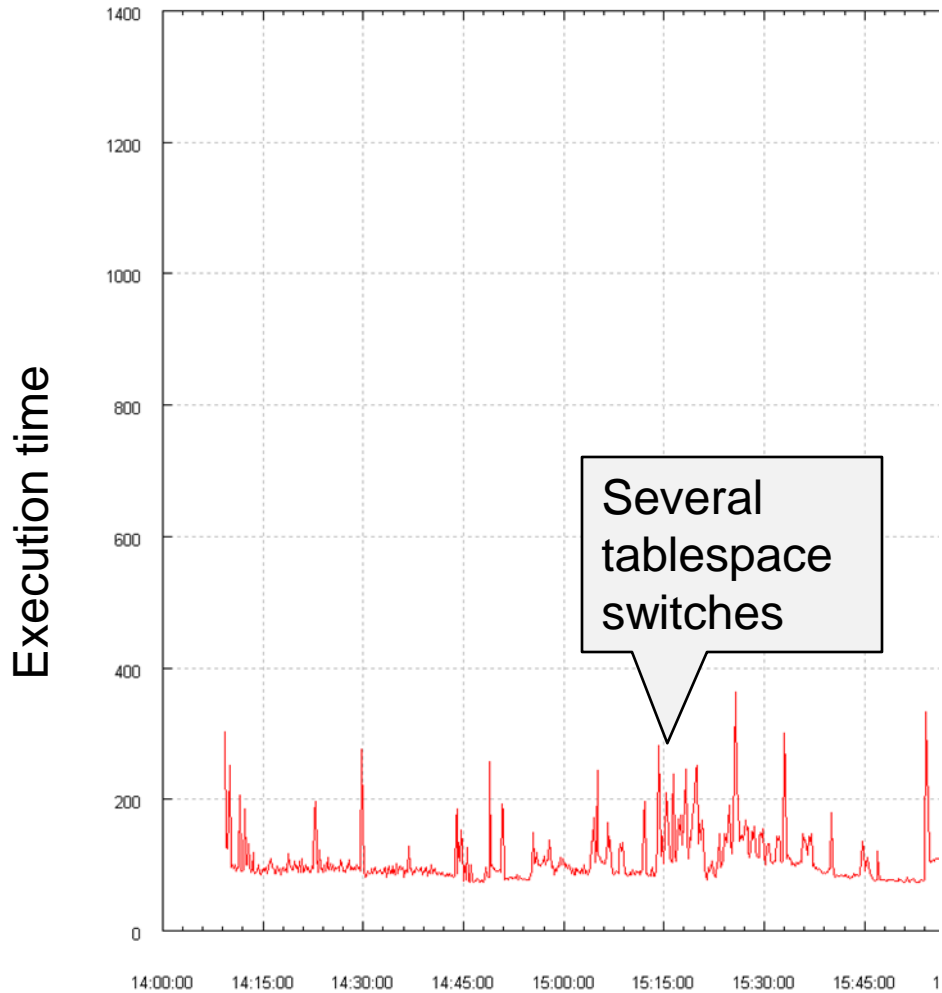
CPU: 0

Disk Activity: 0

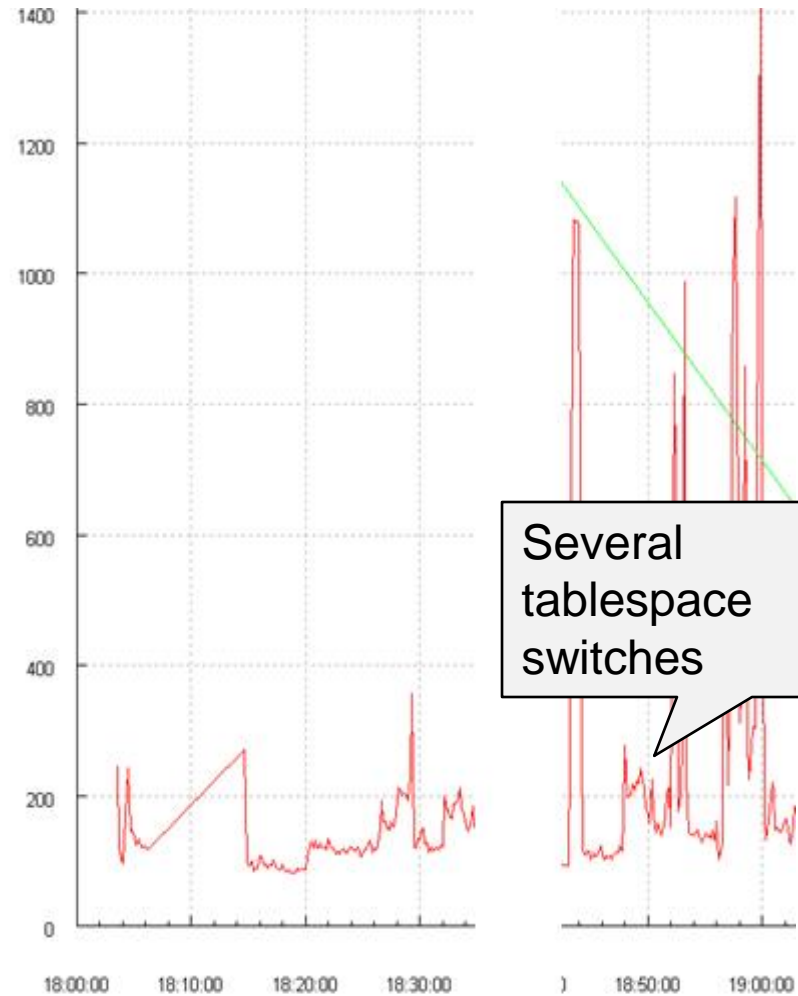
| Property           | Value                               |
|--------------------|-------------------------------------|
| Benchmark Name     | "PVSS Benchmark"                    |
| Connect String     | SWPVSS2                             |
| Coordinator        |                                     |
| Driver Type        | Oracle10g Type II jdbc driver (oci) |
| Maximum Think Time | 0                                   |
| Minimum Think Time | 0                                   |
| Query Time Out     | 600                                 |
| User Count         | 15                                  |
| User Name          | SUSER04                             |

[oracle@sof:~/sw...] SwingBench 2.3... SwingBench 2.3.0... SwingBench 2.3... [root@sof:~] SwingBench 2.3.0...

# Exadata (2/2)



`_cell_fcrc=true`

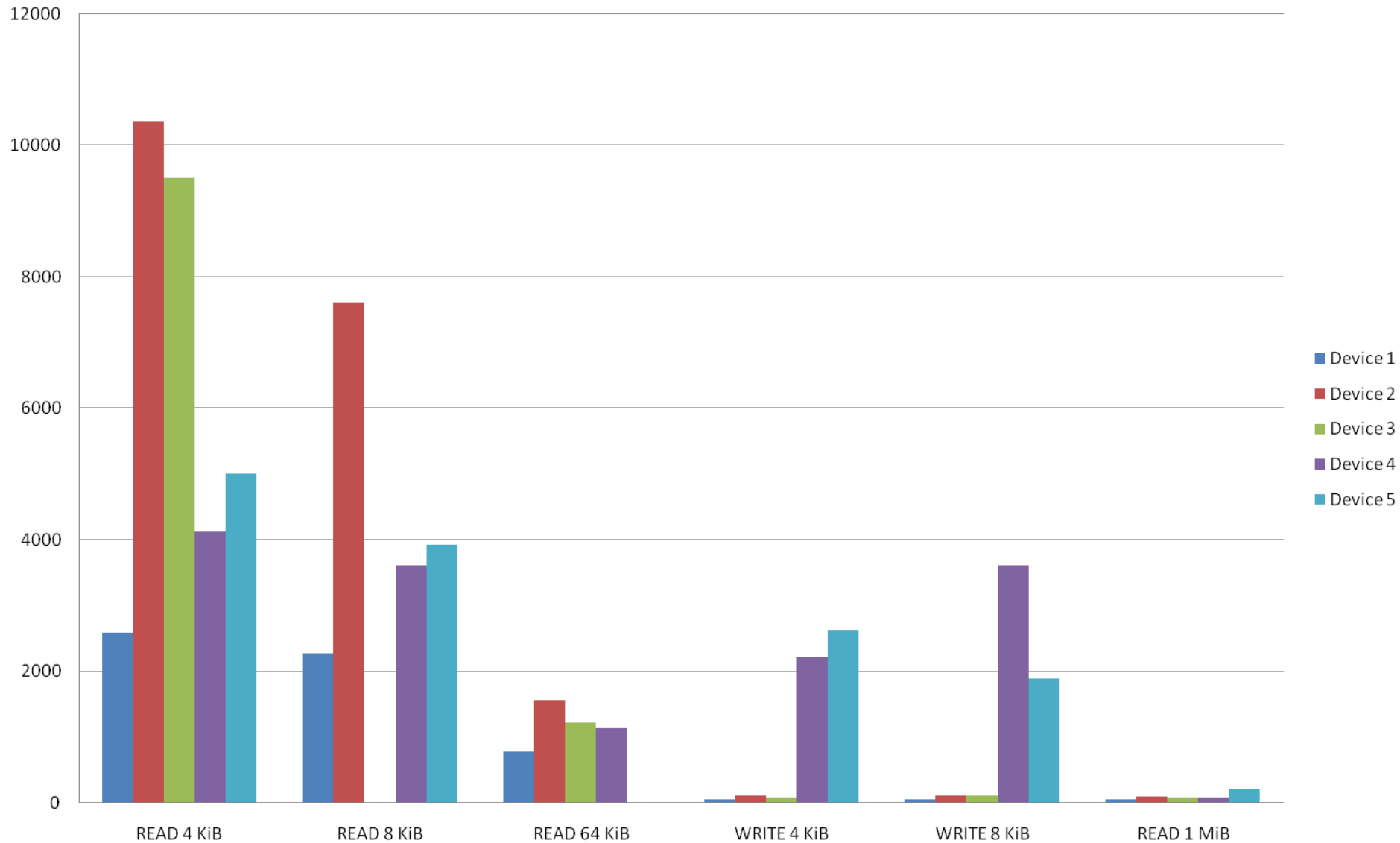


`_cell_fcrc=false`

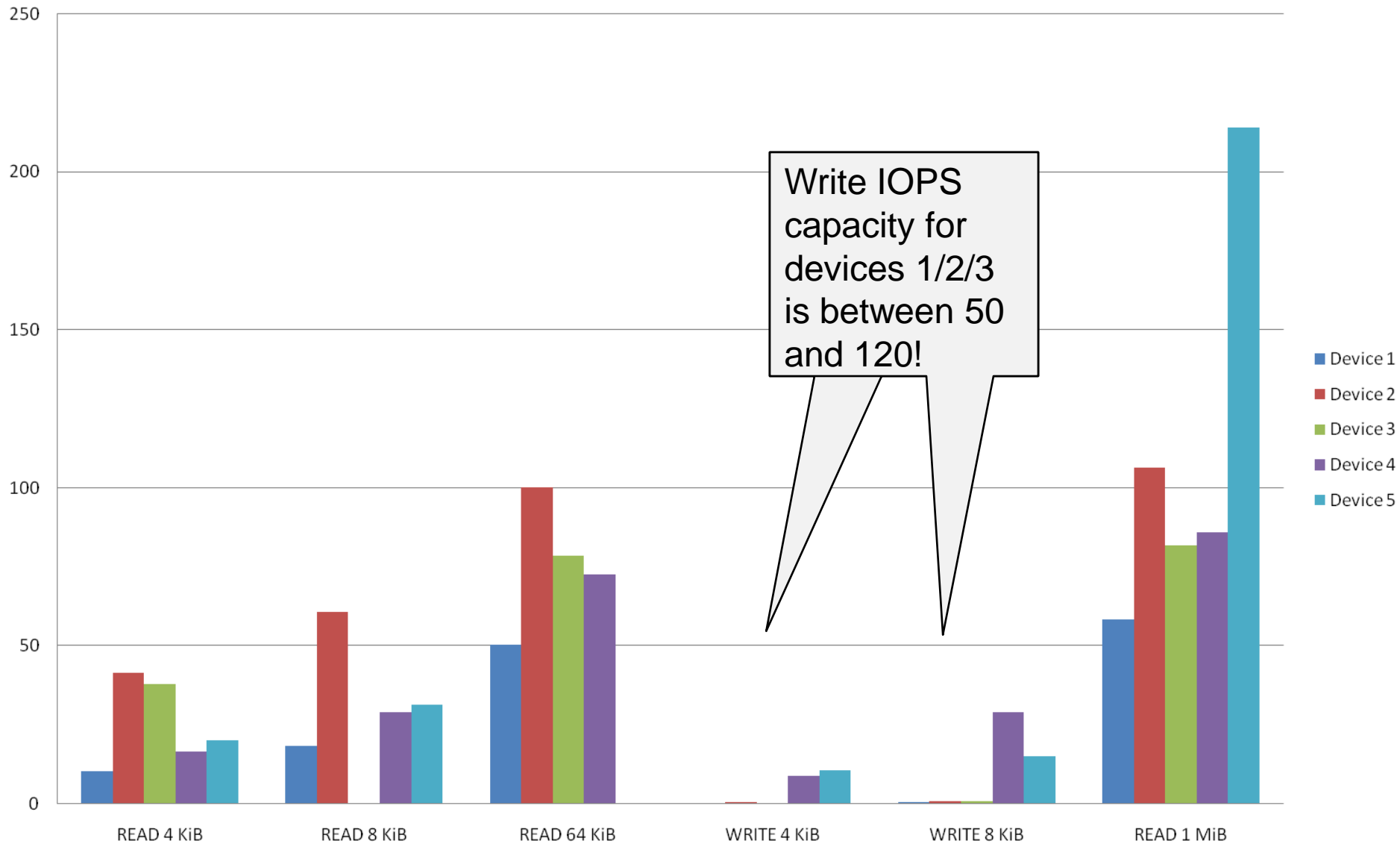
- Solid-State Drive, based on flash, means many different things
- **Single Level Cell (more expensive, said to be more reliable / faster) / Multiple Level Cell**
- Competition in the consumer market is shown on the bandwidth...
- Tests done thanks to Peter Kelemen / CERN – Linux (some done only by him)

- Here are results for 5 different types / models
- Large variety, even the “single cell” SSDs
- (as expected) The biggest difference is with the writing IO operations per second

## SSD IO operations per second



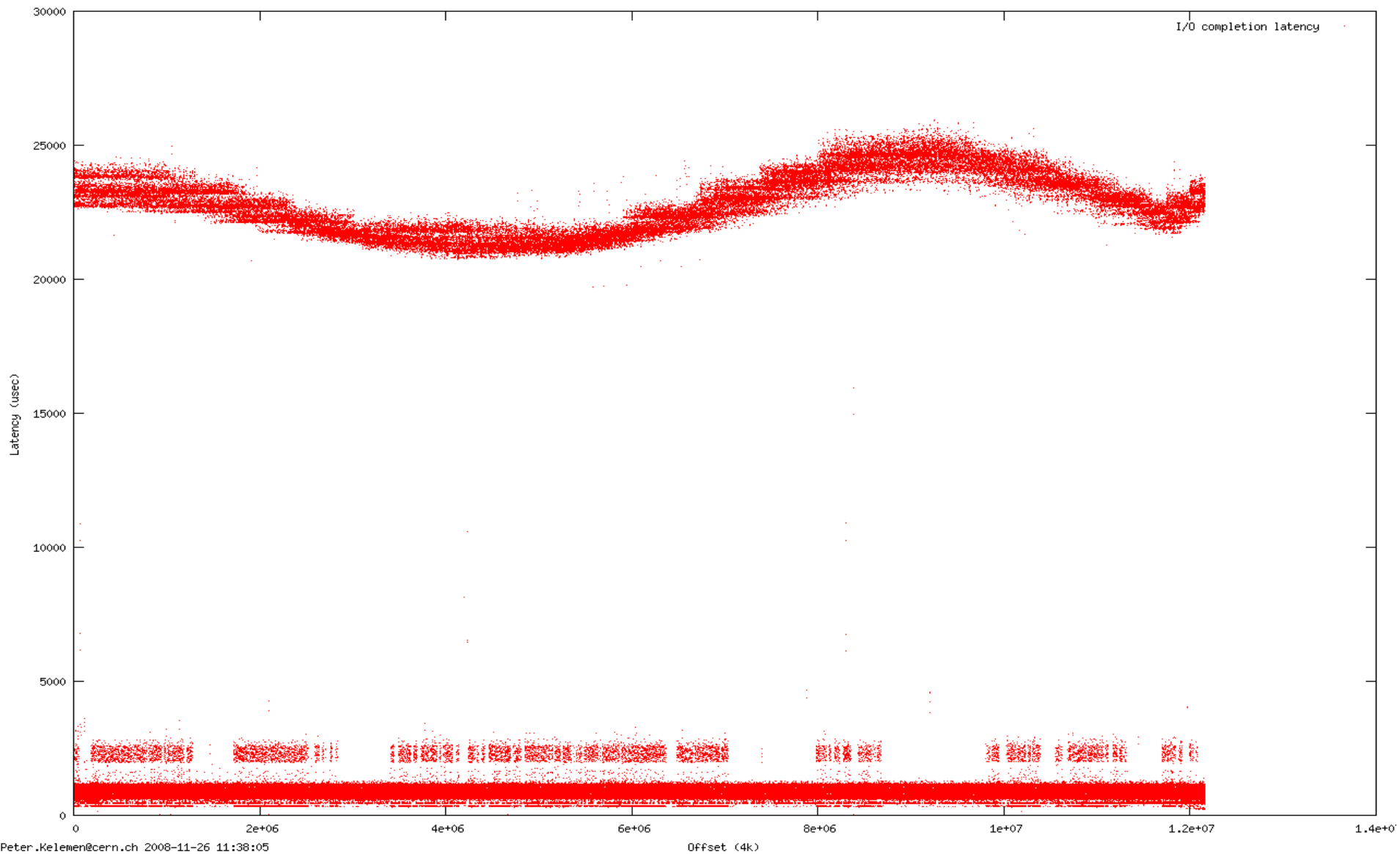
### SSD bandwidth (MiB/s)





- “expensive” and “small” (50GB), complex, very promising
- For random read small IO operations (4kiB or 8kiB), we measure ~4000 to 5000 IOPS (compare to 26 disks)
- For random write operations (4kiB or 8kiB), we measure 2000 to 3000+ write IOPS (compare to 13 disks)
- But for some of the 8K offsets the I/O completion latency is 10× the more common 0.2 ms
- “Wear-levelling/erasure block artefact”? Reported to the vendor

# SSD (6/6)



- New tools like Dtrace change the way we can track IO operations
- Overload in IO and CPU can not be seen from Oracle IO views
- Exadata offloading operations can be interesting (and promising)
- Flash SSD are coming, a lot of differences between them. Writing is the issue (and is a driving price factor). Not applicable for everything. Not to be used for everything for now (as write cache? Oracle redo logs) They change the way IO are perceived. Test/test/test!

- Why You Should Focus on LIOs Instead of PIOs  
Author, Cary Millsap <http://www.hotsos.com/e-library/abstract.php?id=7>
- Tanel Poder DStackProf  
<http://tanelpoder.otepad.com/script:dstackprof.sh>
- Metalink Note 175982.1
- Tanel Poder os\_explain.sh  
[http://www.tanelpoder.com/files/scripts/os\\_explain](http://www.tanelpoder.com/files/scripts/os_explain)

# Q&A