


ORACLE®

# MySQL Performance: Benchmarks, Tuning and “Best” Practices..

Dimitri KRAVTCHUK  
MySQL Performance Architect @Oracle





The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract.

It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

# Are you Dimitri?..

- Yes, it's me :-)
- Hello from Paris! ;-)
- Passionated by Systems and Databases Performance
- Previous 15 years @Sun Benchmark Center
- Started working on MySQL Performance since v3.23
- But during all that time just for fun only ;-)
- Since last years officially @MySQL Performance full time now
- <http://dimitrik.free.fr/blog> / @dimitrik\_fr



# Agenda

- Overview of MySQL Performance
- Workload oriented tuning and MySQL Internals
- Performance improvements in MySQL 5.6 & Benchmark results
- Pending issues..
- Progress in MySQL 5.7 Performance
- Q & A



# Why MySQL Performance ?...

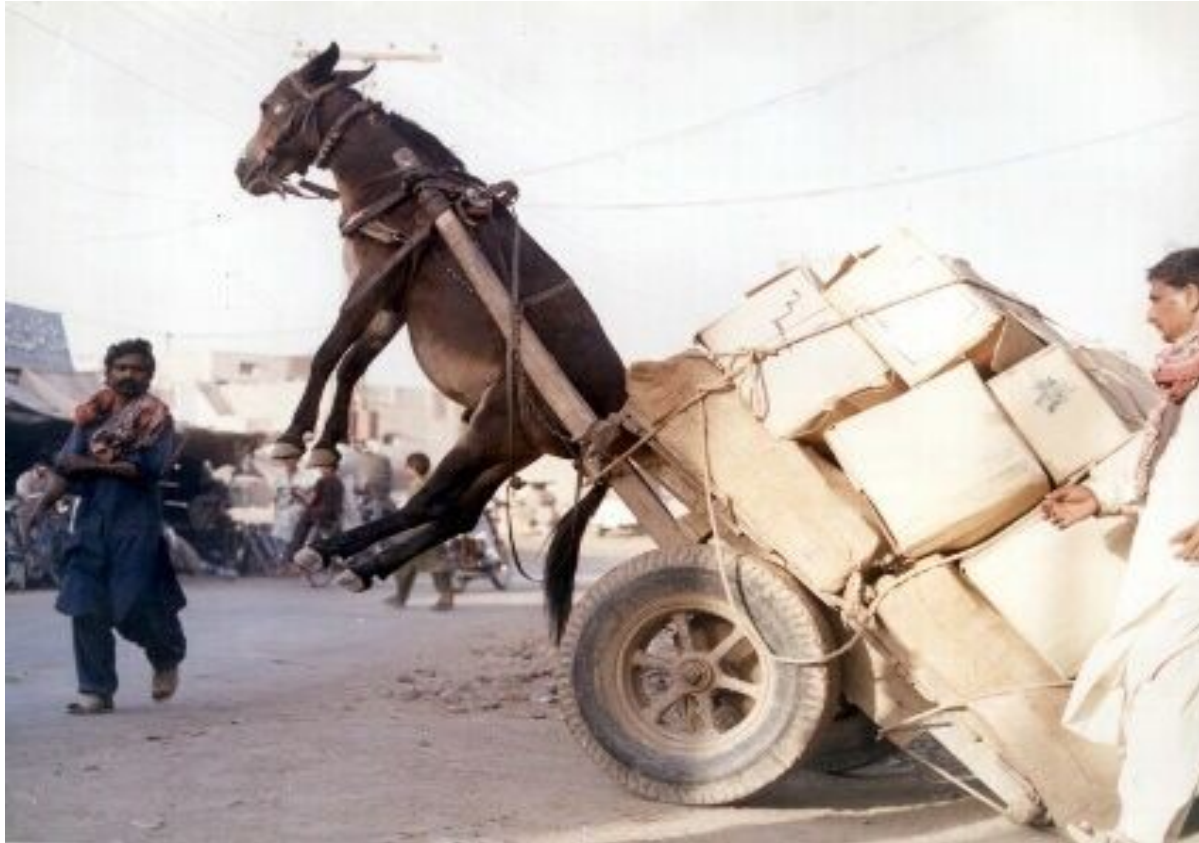
# Why benchmarking MySQL?..

- Any solution may look “good enough”...



# Why benchmarking MySQL?..

- Until it did not reach its limit..





# Why benchmarking MySQL?..

- And even improved solution may not resist to increasing load..



[www.freeuniverse4all.com](http://www.freeuniverse4all.com)

# Why benchmarking MySQL?..

- And reach a similar limit..



# Why benchmarking MySQL?..

- A good benchmark testing may help you understand ahead the resistance of your solution to incoming potential problems ;-)



# Why benchmarking MySQL?..

- But keep it in mind:
  - Even a very powerful solution but leaved in wrong hands may still be easily broken!... :-)





The Main MySQL Performance Tuning  
**#1** Best Practice is... ???..



The Main MySQL Performance Tuning  
#1 Best Practice is... ???..

**USE YOUR BRAIN !!! :-)**



The Main MySQL Performance Tuning  
#1 Best Practice is... ???..

**USE YOUR BRAIN !!! :-)**



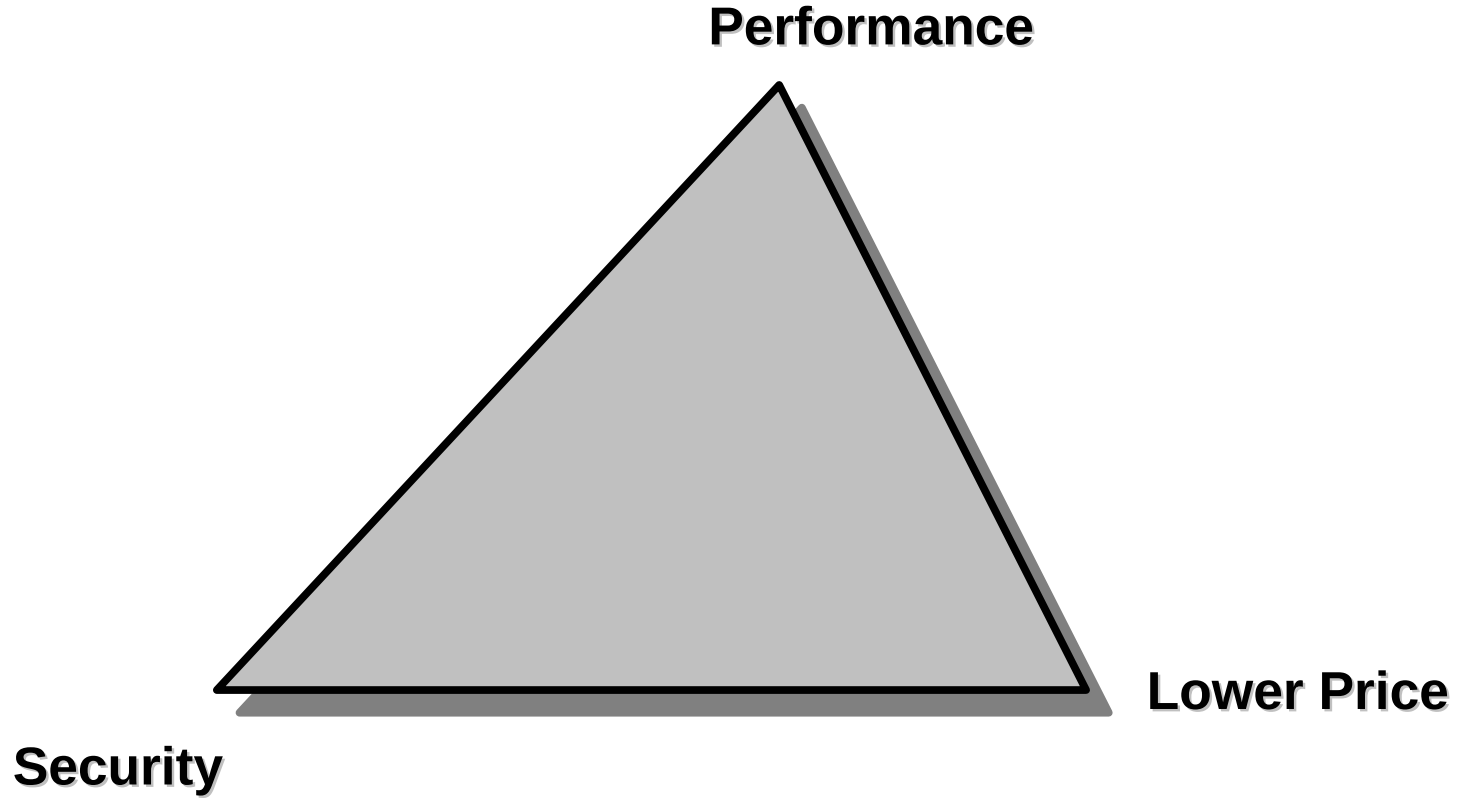
**IN FACT  
THIS IS THE  
MAIN SLIDE! ;-))**

# Think “Database Performance” from the beginning!

- **Server:**
  - Having faster CPU is still better! 32 cores is good enough ;-)
  - OS is important! - Linux, Solaris, etc.. (and Windows too!)
  - Right malloc() lib!! (Linux: jemalloc, Solaris: libumem)
- **Storage:**
  - Don't use slow disks! (except if this is a test validation goal :-))
  - SSD helping random access! (index/data) more and more cheaper
  - FS is important! - ZFS, UFS, QFS, VxFS, EXT3, EXT4, XFS, etc..
  - O\_DIRECT or not O\_DIRECT, AIO or not AIO, and be aware of bugs! ;-)
  - Do some generic I/O tests first (Sysbench, IObench, iohome, etc.)
- **Don't forget network !! :-)** (faster is better, 10Gbit is great!)



**Seek for your best option..**



# Know your platform limits / “features”..

- My backup is finished on Linux faster than on Solaris same HW
  - Be sure first there is really no more I/O activity once backup is “finished”
  - Keep in mind Linux buffering..
- Linux distro: MySQL Performance has x4 regression! Fix it!
  - How did you see it? – Our QA test is taking x4 times more time..
  - Which engine? – InnoDB..
  - What is innodb\_flush\_log\_at\_trx\_commit value? – set to 1.. why?
  - Tried innodb\_flush\_log\_at\_trx\_commit=2 ?.. – Oh! You fixed it!! Thanks!!
  - Wait, what did you “improve” recently in distro? – FS flushing, why?..
  - Well, the test in fact is proving that you did not “sync” on every fsync() before, that's all.. But now in your FS flushing you get it fixed ;-)

# Advise: Benchmark your platform / prototype!

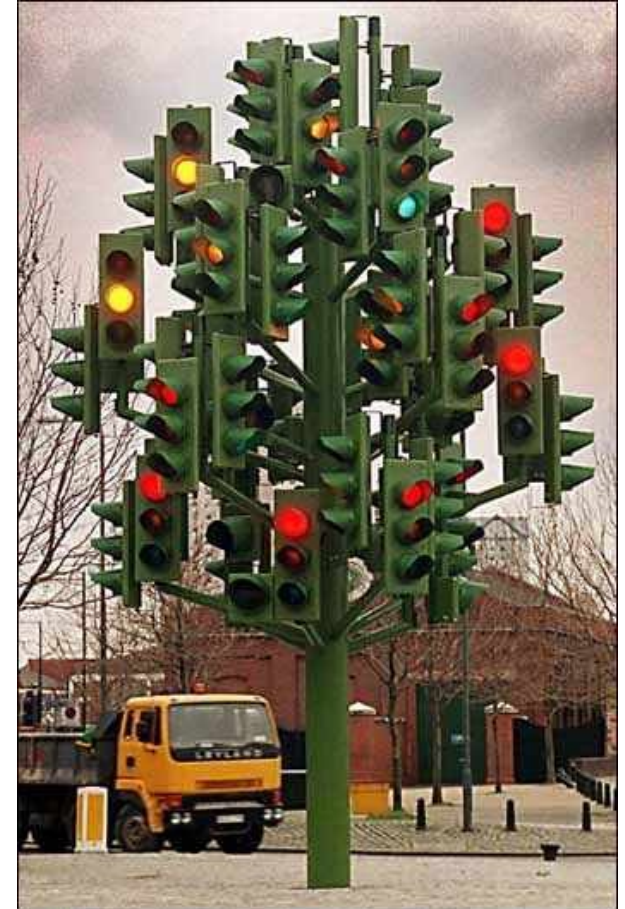
- Have a clear goal!
  - Otherwise: I've obtained all these results, and now... so what?..
- Want to simulate your production workload?..
  - Then just simulate it! (many SW available, not always OSS/free)
  - Hard to simulate? - adapt some generic tests
- Want to know capacity limits of a given platform?
  - Still try to focus on the test which are most significant for you!
- Want just to validate config settings impacts?
  - Focus on tests which are potentially depending on these settings
  - Or any, if the goal to prove there are not depending ;-)
- Well, just keep thinking about what you're doing ;-)

# Popular “Generic” Test Workloads @MySQL

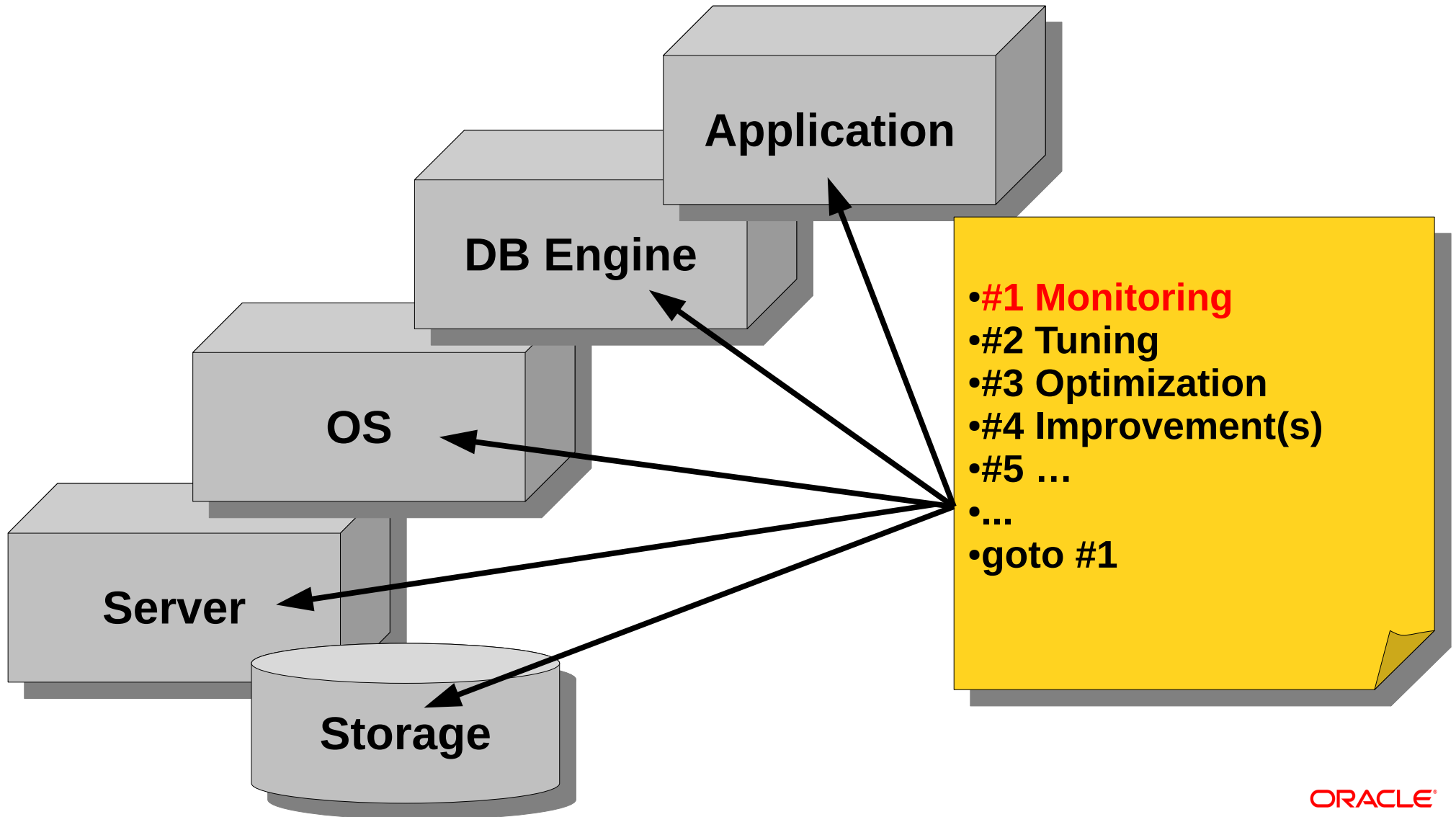
- Sysbench
  - OLTP, RO/RW, 1-table, since v0.5 N-table, lots load options, deadlocks
- DBT2 / TPCC-like
  - OLTP, RW, very complex, growing db, no options, deadlocks
  - In fact using mostly only 2 tables! (thanks Performance Schema ;-))
- dbSTRESS
  - OLTP, RO/RW, several tables, one most hot, configurable, no deadlocks
- linkbench (Facebook)
  - OLTP, RW, very intensive
- DBT3
  - DWH, complex heavy query, loved by Optimizer Team ;-)

# Test Workload

- Before to do something complex...
  - Be sure first you're comfortable with “basic” operations!
  - Single table?
  - Many tables?
  - Short queries?
  - Long queries?
- Remember: any complex load just represents a mix of simple operations..
  - So, start from as simple as possible..
  - And then increase complexity progressively..

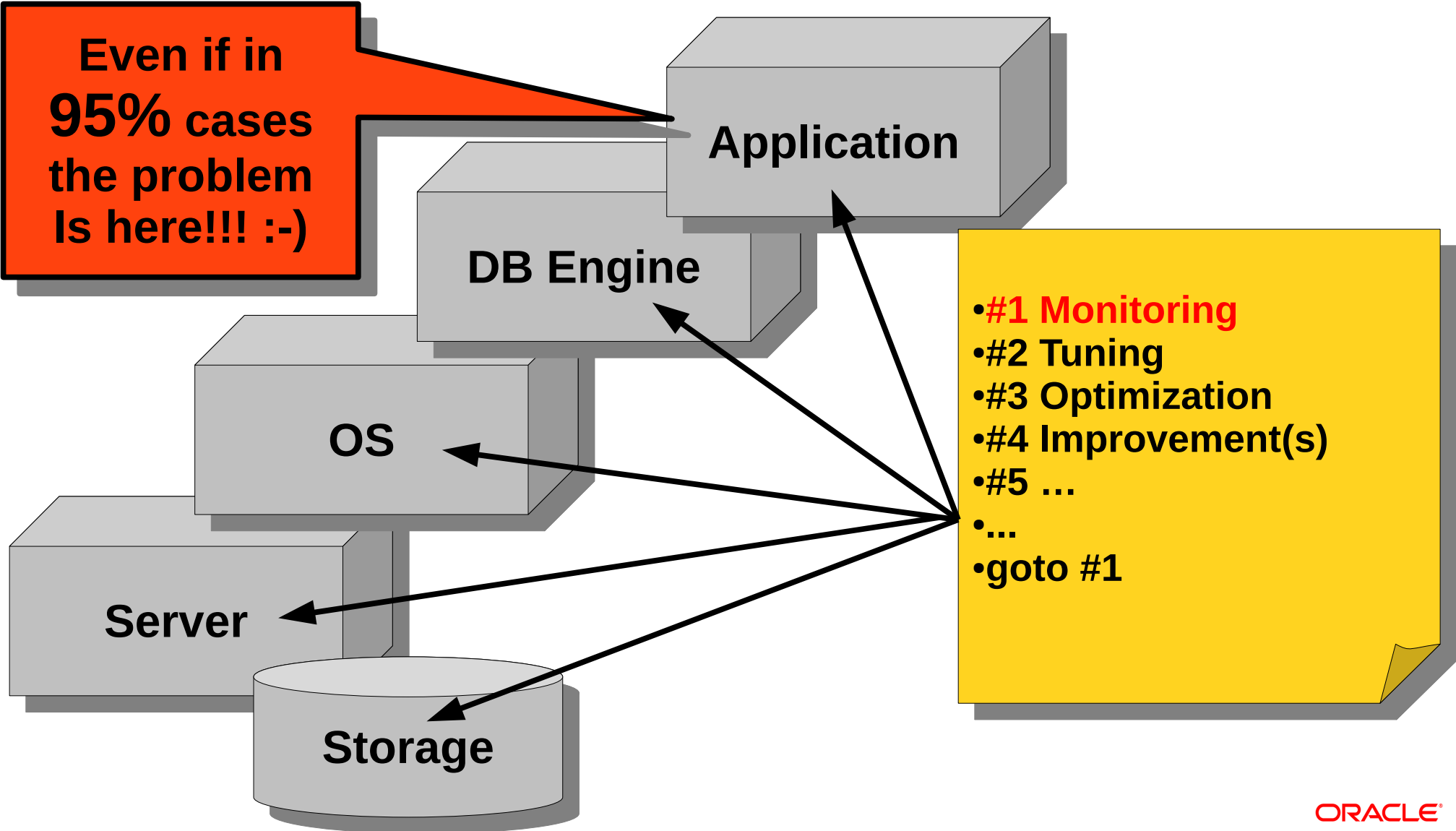


# The Infinite Loop of Database Tuning...



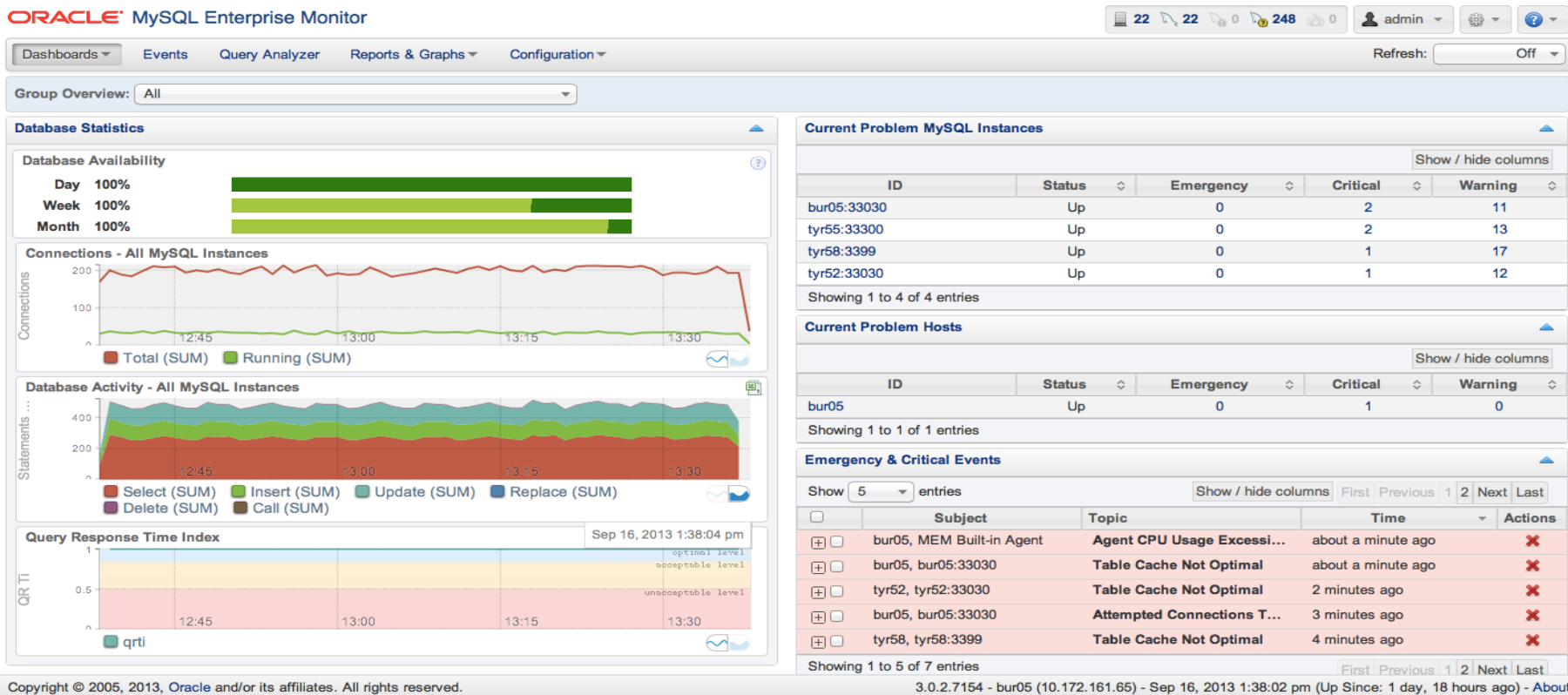
# The Infinite Loop of Database Tuning...

Even if in  
**95%** cases  
the problem  
is here!!! :-)



# MySQL Enterprise Monitor

- Fantastic tool!
  - Did you already try it?.. Did you see it live?..





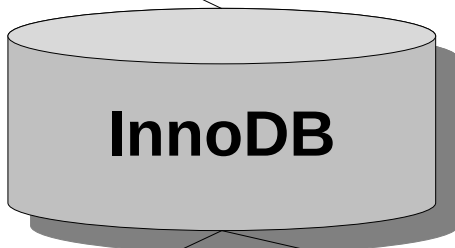
# Other Monitoring Tools

- Cacti
- Zabbix
- Nagios
- Etc.....
- ***dim\_STAT***
  - well, I'm using this one, sorry ;-)
  - all graphs within presentation were made with it
  - details are in the end of presentation..

# MySQL Performance



Internal Limits..



Configuration Settings..

Query Optimization..

Server, OS, FS

Storage  
BBU, SSD

Application Contentions..

# Basic Tuning

- **Understanding HW platform limits**
  - helps you to deploy your MySQL Server in the most optimal way..
- **Understanding MySQL Server internals**
  - helps you to configure your database settings in the most optimal way..
  - use the best adapted Storage Engine
- **Understanding of your Workload**
  - helps you to tune the whole solution in the most optimal way ;-)
  - 20% of known issues covering 80% of observed problems..
  - So, adapt some best practices from the beginning..

# Storage Engines: use InnoDB & MySQL 5.6 ;-)

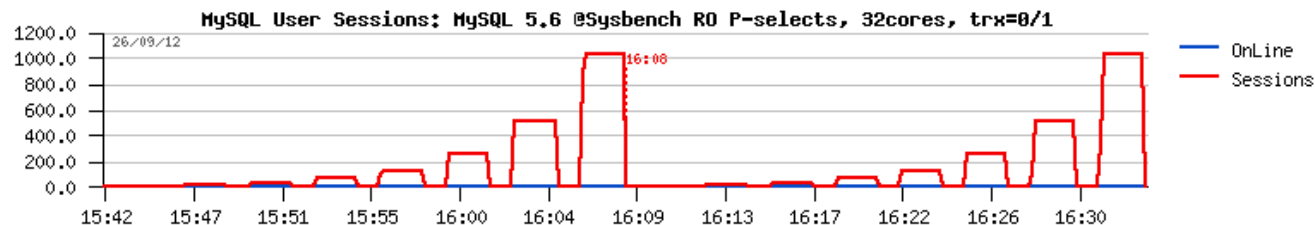
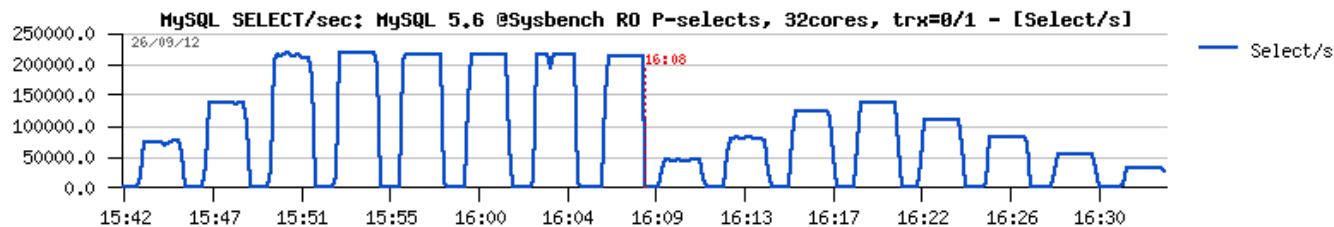
- MyISAM
  - Table level locking for everything, cache(s) for indexes only
  - Easily spending 50% of time on syscalls to read() data..
  - Many global locks, no transactions, no recovery..
- InnoDB
  - Row locking, transactions for everything, Buffer Pool
  - Double write buffer? Checksums?
  - innodb\_flush\_log\_at\_trx\_commit= 1 / 2 ??
  - **many huge improvements since MySQL 5.6!**
- Binlog / Replication
  - Sync? - **binlog group commit** is since MySQL 5.6

# Workload: Read-Only oriented

- Bigger Buffer Pool (BP) is better
  - $BP < \text{dataset} = \text{IO-bound}$
- TRX list (kernel\_mutex, since 5.6: trx\_sys mutex)
- Read view
- Auto-commit or transactions?..
  - Grouping **many** queries within a single transaction may also largely reduce MDL locking, but still keep them **short** ! (check with PFS)
- Prepared statements
  - Observed 10% performance improvement in 5.6 (while Parser time is not more than 3% according to profiler)..
- Read-Only transactions!

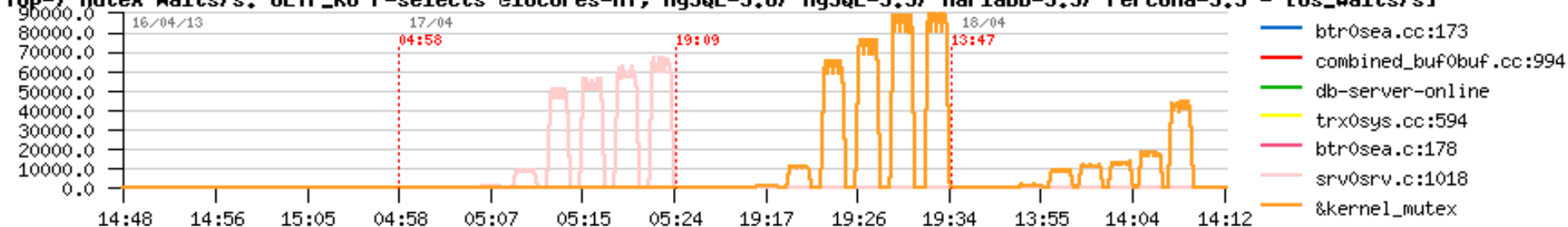
# InnoDB: Read-Only Transactions in 5.6

- Sysbench OLTP\_RO Point-Selects:
  - Concurrent user sessions: 1, 2, 4 .. 1024
  - Using of transactions in sysbench = 0 / 1

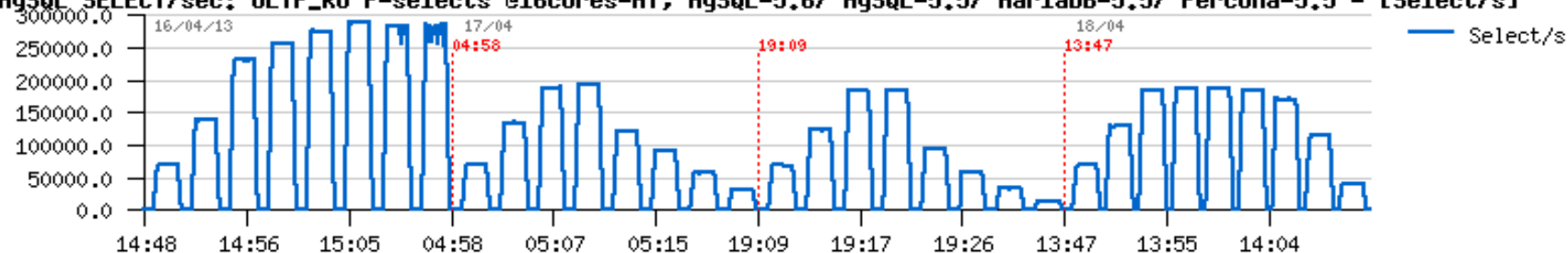


# InnoDB: Read-Only Transactions in 5.6 (Apr.2013)

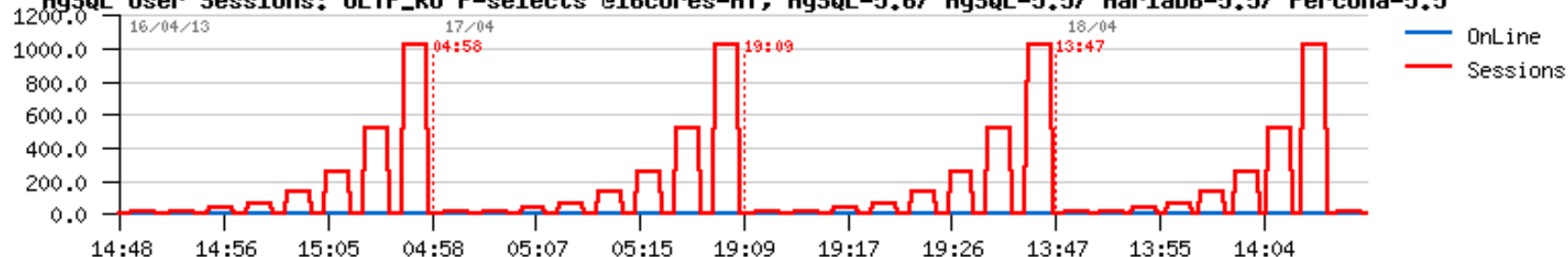
InnoDB Top-7 Mutex Waits/s: OLTP\_RO P-selects @16cores-HT, MySQL-5.6/ MySQL-5.5/ MariaDB-5.5/ Percona-5.5 - [os\_waits/s]



MySQL SELECT/sec: OLTP\_RO P-selects @16cores-HT, MySQL-5.6/ MySQL-5.5/ MariaDB-5.5/ Percona-5.5 - [Select/s]

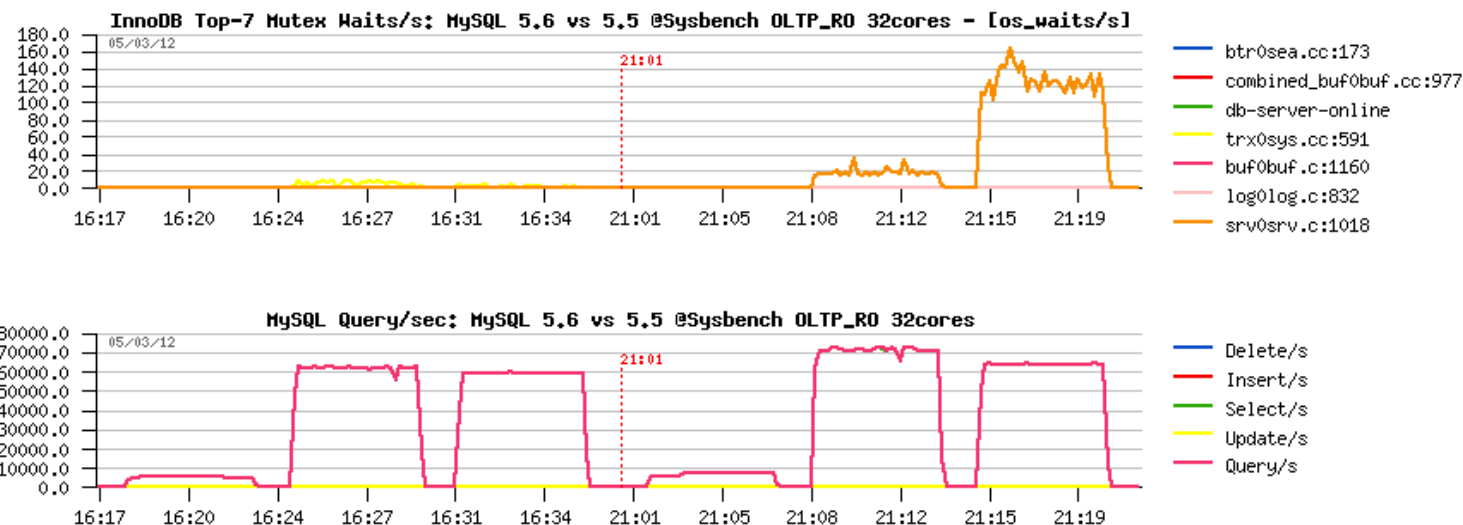


MySQL User Sessions: OLTP\_RO P-selects @16cores-HT, MySQL-5.6/ MySQL-5.5/ MariaDB-5.5/ Percona-5.5



# InnoDB : false sharing of cache-line = true killer

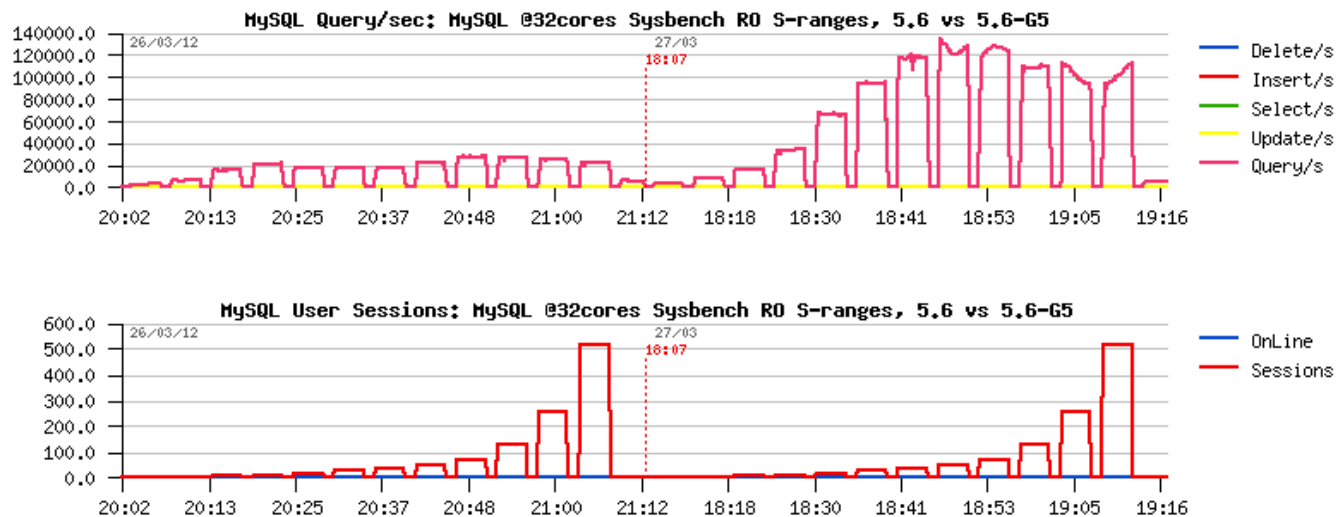
- RO or RW Workloads
  - Same symptoms in 5.5 & 5.6 : no QPS improvement between 16 and 32 user sessions:





# InnoDB : false sharing of cache-line fixed!

- RO or RW Workloads
  - “G5” patch! :-)
  - Over **x2(!)** times better on Sysbench OLTP\_RO,
  - **x6(!)** times better on SIMPLE-Ranges!
  - NOTE: the fix is not applicable on 5.5..

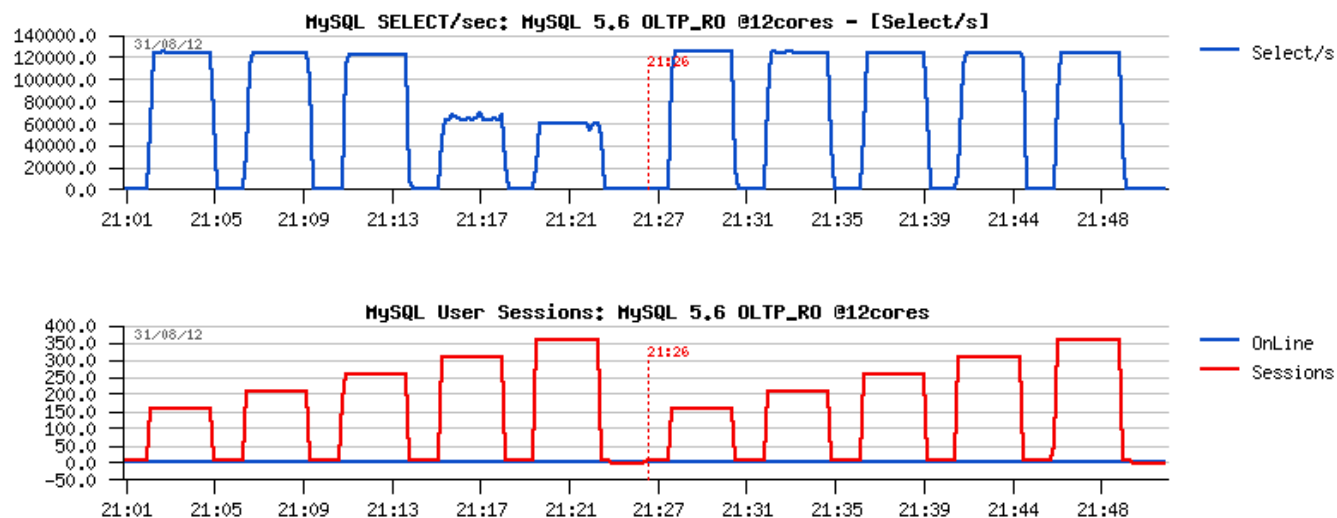


# MySQL Internals: “killer” LOCK\_open mutex

- MySQL 5.5 and before:
  - Keep “table\_open\_cache” setting big enough!
  - Monitor global status for '%opened%'
  - Once this contention become the most hot – well, time to upgrade to 5.6 ;-))
- Since MySQL 5.6:
  - Fixed: several table open cache instances
  - But it doesn't mean you can use a small “table\_open\_cache” either ;-)
  - Monitor PFS Waits!
  - Monitor “table\_open\_cache%” status variables!
  - Keep “table\_open\_cache\_instances” at least bigger than 1

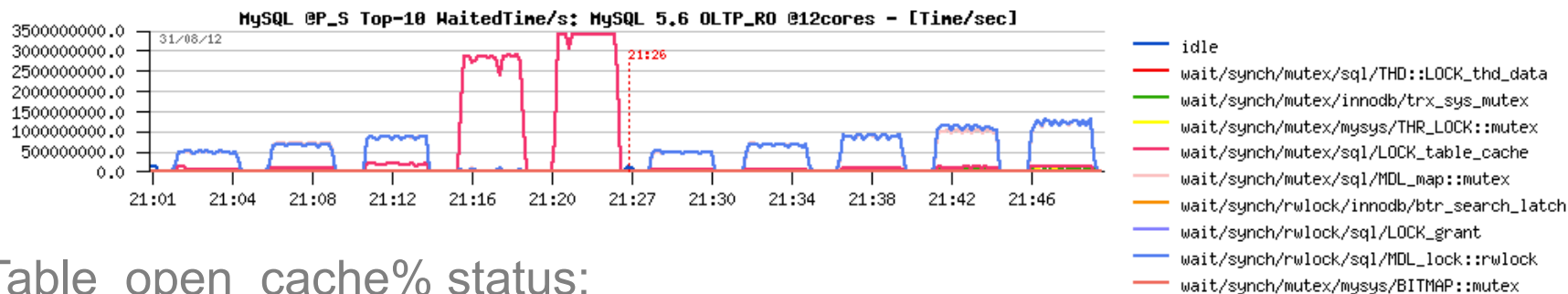
# MySQL 5.6 Internals : low table\_open\_cache

- MySQL 5.6 :
  - Not big enough “table\_open\_cache” setting

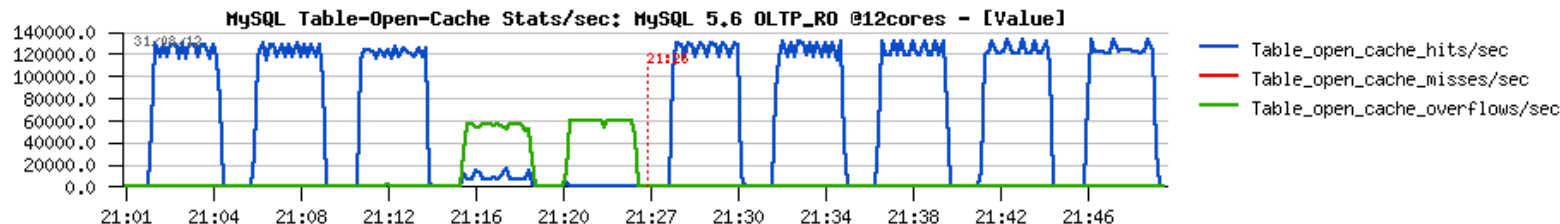


# MySQL 5.6 Internals : low table\_open\_cache (2)

- MySQL 5.6 :
  - Not big enough “table\_open\_cache” setting
  - PFS Waits monitoring: LOCK\_table\_cache become the most hot:

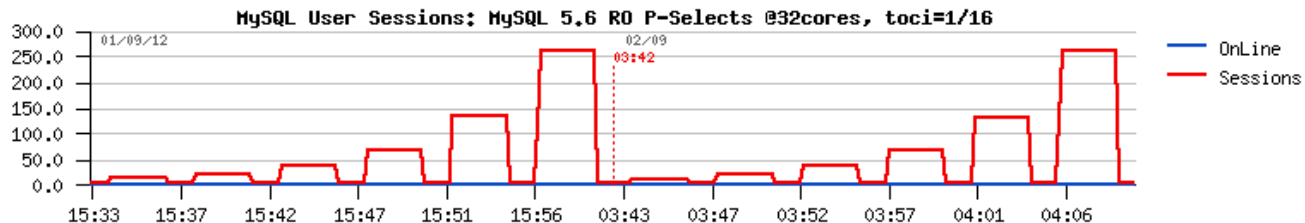
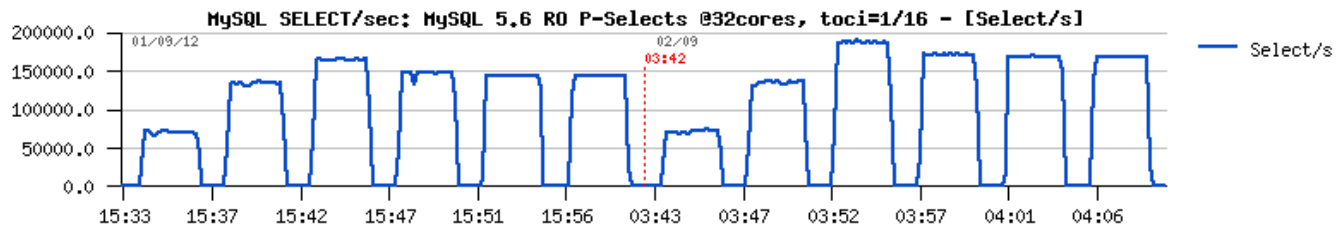
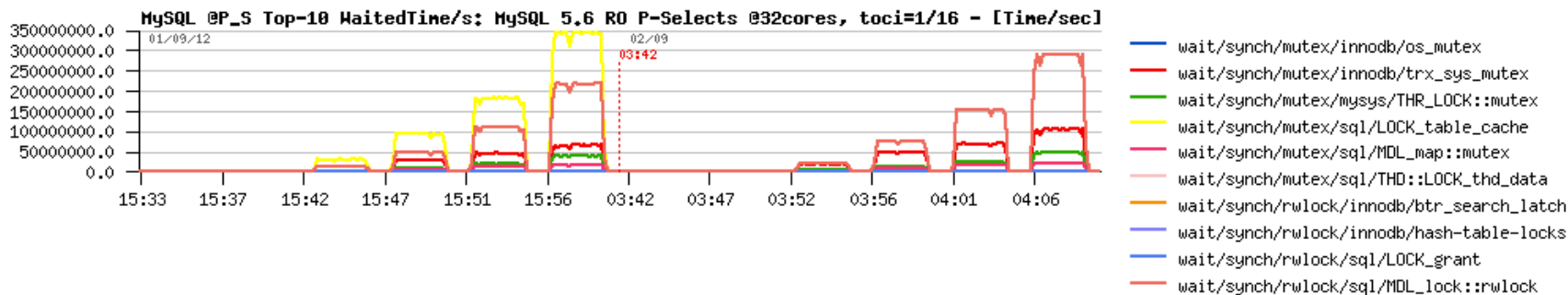


- Table\_open\_cache% status:



# MySQL 5.6 Internals : table\_open\_cache\_instances

- MySQL 5.6 :
  - When LOCK\_table\_cache wait is on top, the gain is usually well visible:



# Workload: Read+Write

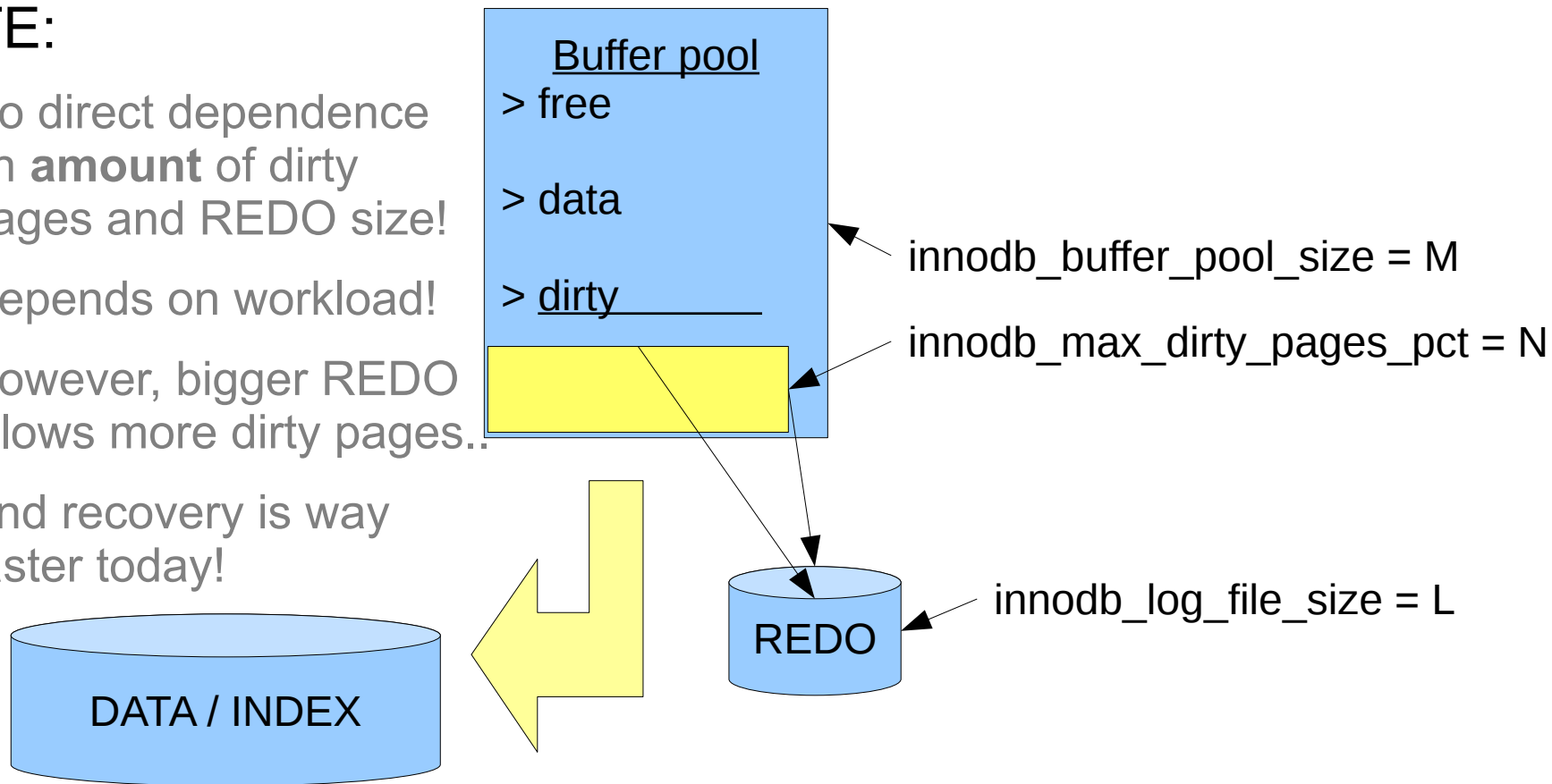
- RW activity
  - Updates only? Insert? Delete? R/W %ratio?
- Bigger Buffer Pool (BP) is still better
  - $BP < \text{dataset}$  = IO-bound Reads(!) or R+W
  - $BP > \text{dataset}$  = CPU-bound or IO-bound Writes(!)
- REDO size matters a lot! (up to 2TB in 5.6)
- Adaptive Flushing matters a lot!
- LRU flushing matters a lot as well!
- Tip: Neighbor Pages flushing = off / on

# InnoDB: Dirty pages, Flushing

- Direct dependence on REDO log size

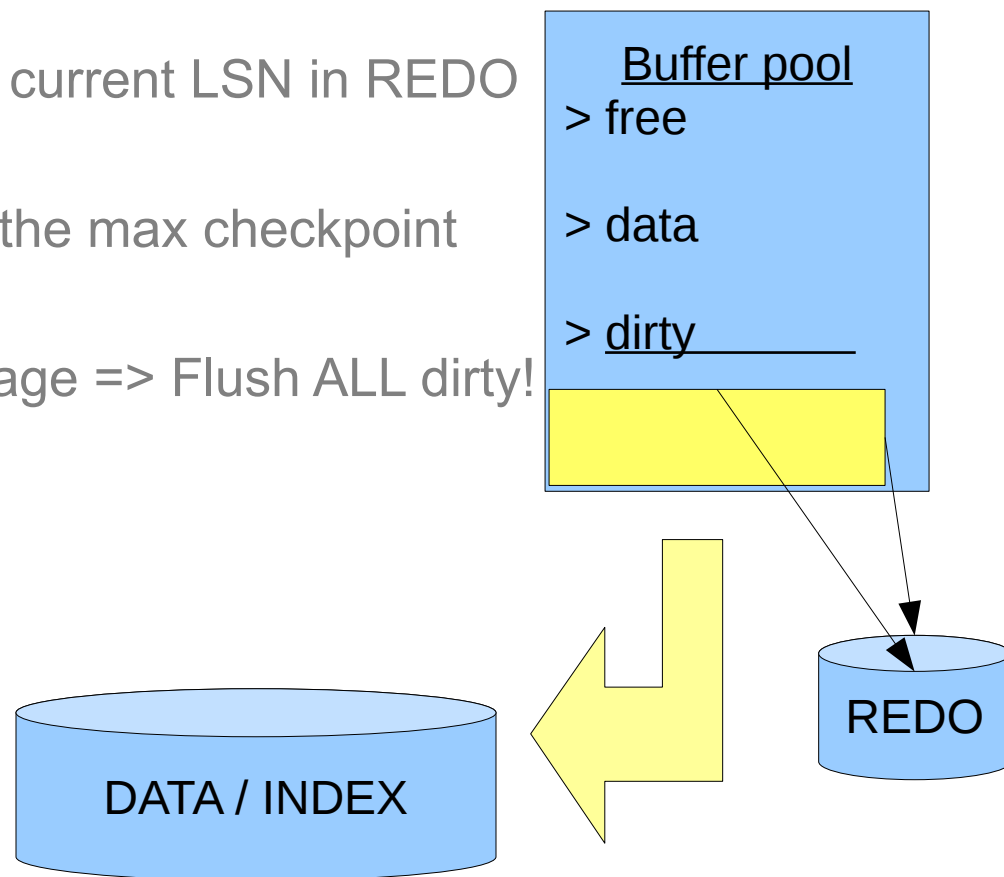
- NOTE:

- No direct dependence on **amount** of dirty pages and REDO size!
- Depends on workload!
- However, bigger REDO allows more dirty pages..
- And recovery is way faster today!



# InnoDB: REDO log constraints

- REDO log constraints: (Always monitor Checkpoint Age!!!)
  - Cyclic, need free space
  - Checkpoint age: diff between the current LSN in REDO and the oldest dirty page LSN
  - Checkpoint age cannot out-pass the max checkpoint age (redo log size)
  - If Checkpoint age  $\geq 7/8$  of Max age  $\Rightarrow$  Flush ALL dirty!
  - $\Rightarrow$  AKA “furious flushing”...
- Adaptive Flushing:
  - Keep REDO under Max age
  - Respecting IO capacity limit



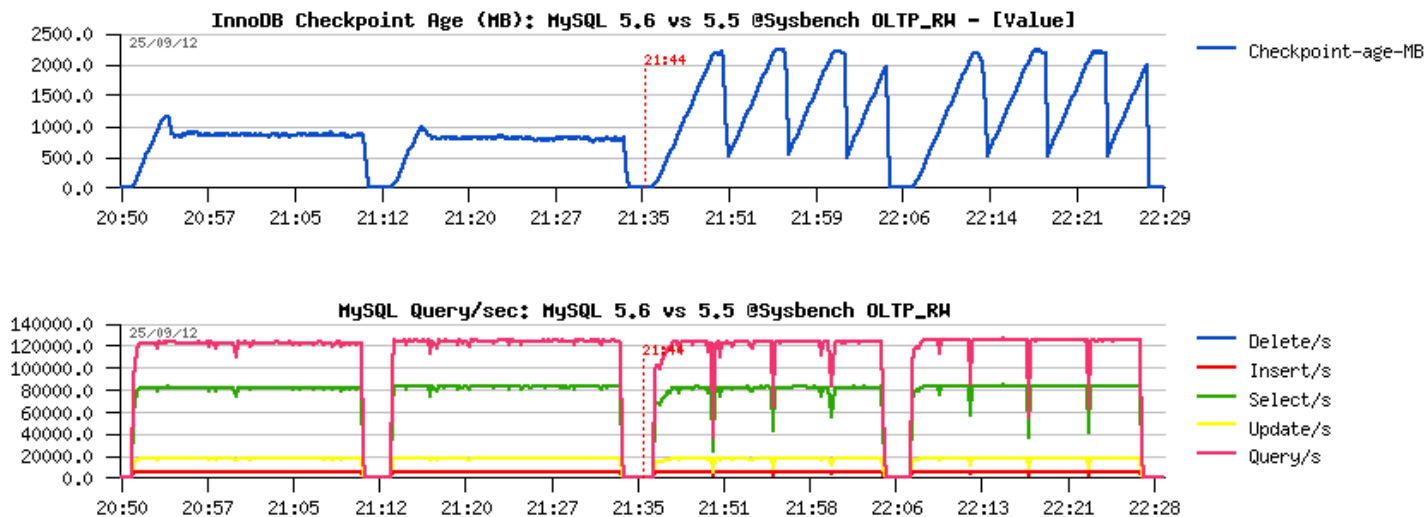


# InnoDB: Adaptive Flushing

- MySQL 5.5:
  - Estimation based
  - Sometimes works ;-)
- MySQL 5.6 :
  - Based on REDO write rate + I/O capacity Max
  - Involving batch flushing with N pages to flush (progressive, depending on REDO %free) + page age limit (according REDO rate)
- Tuning:
  - `innodb_io_capacity / innodb_io_capacity_max`
  - `innodb_adaptive_flushing_lwm / innodb_max_dirty_pages_pct_lwm`
  - ALL are dynamic!
  - Monitor Checkpoint Age..

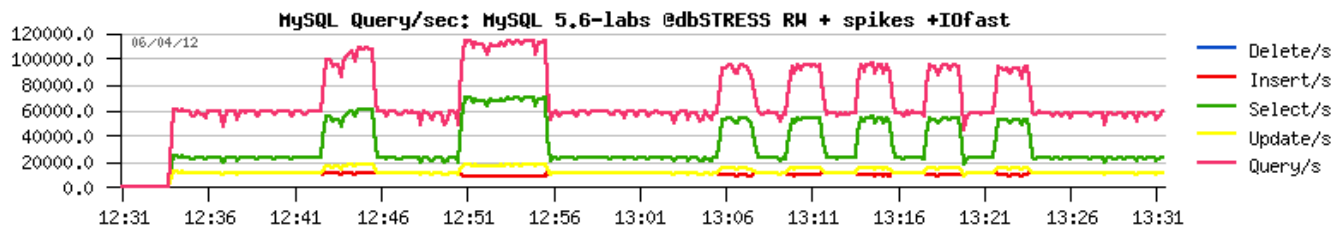
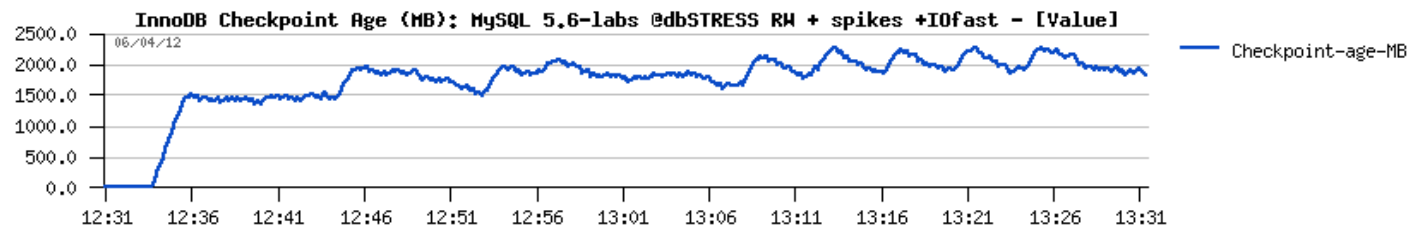
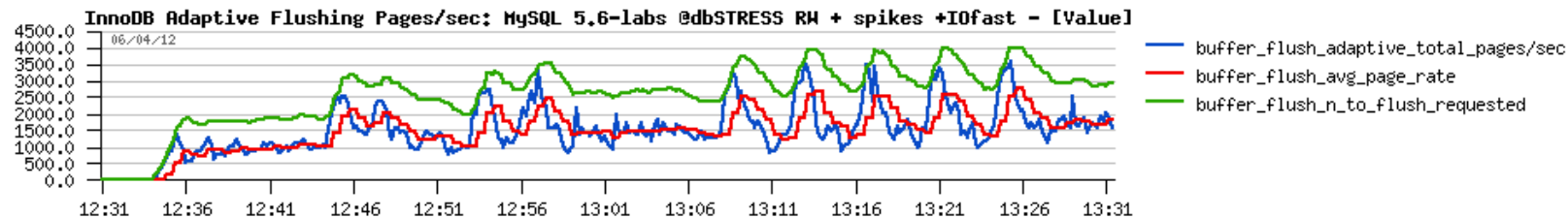
# Adaptive Flushing: MySQL 5.6 vs 5.5

- OLTP\_RW Workload:
  - Same IO capacity
  - Different logic..



# InnoDB : Resisting to activity spikes in 5.6

- dbSTRESS R+W with spikes



# User Concurrency scenarios

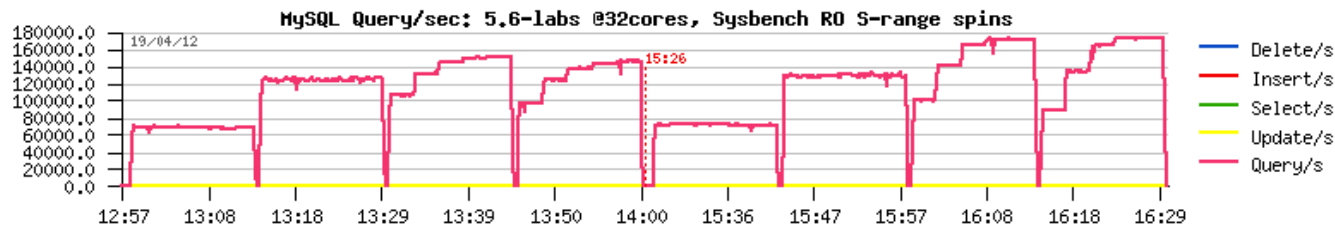
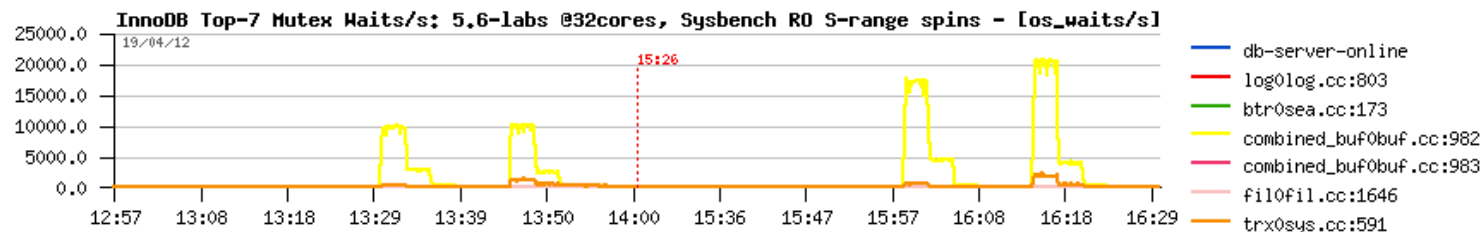
- Single user?..
  - With a bigger code path today 5.6 simply cannot be faster than 5.5
  - But then, why you're not considering Query Cache? ;-)
- More users?..
  - Up to 8-16 concurrent users all internal contention are not yet hot
  - So, 5.6 will not be better yet..
- More than 16 users?..
  - Then you'll feel a real difference, but if you have at least 16cores ;-)
  - Or if you have really a lot of concurrent users
- But don't forget other 5.6 improvements either!
  - On-line DDL, Binlog group commit, Memcached, etc..

# High Concurrency Tuning

- If bottleneck is due a concurrent access on the same data (due application design) – ask dev team to re-design ;-)
- If bottleneck is due MySQL/InnoDB internal contentions, then:
  - If you cannot avoid it, then at least don't let them grow ;-)
  - Try to increase InnoDB spin wait delay (dynamic)
  - Try innodb\_thread\_concurrency=N (dynamic)
  - CPU taskset / prcset (Linux / Solaris, both dynamic)
  - Thread Pool
  - NOTE: things with contentions may radically change since 5.7, so stay tuned ;-)

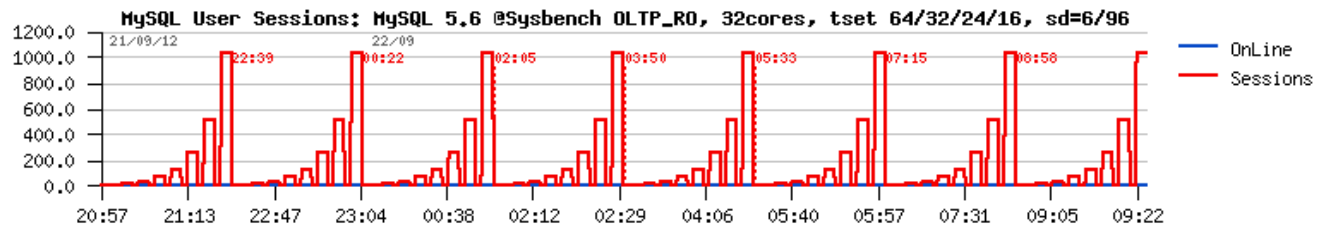
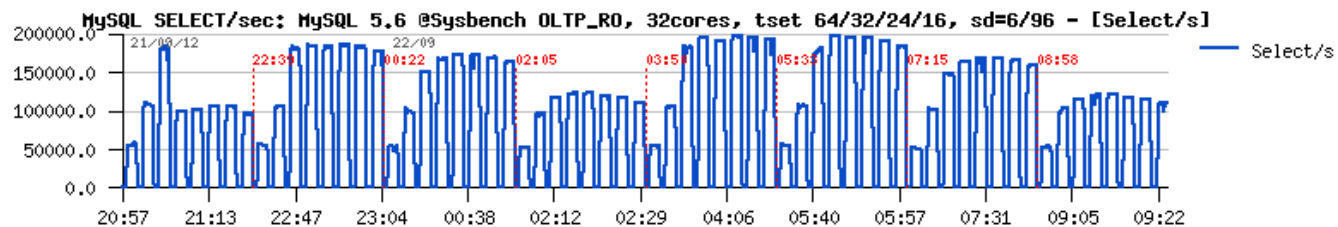
# InnoDB Spin Wait Delay

- RO/RW Workloads:
  - With more CPU cores internal contentions become more hot..
  - Bind mysqld to less cores helps, but the goal is to use more cores ;-)
  - Using innodb\_thread\_concurrency is not helping here anymore..
  - So, **innodb\_spin\_wait\_delay** is entering in the game:



# Tune InnoDB Spin Wait Delay

- Notes :
  - is the max random delay on “sleep” within a spin loop in wait for lock..
  - Ideally should be auto.. while the same tuning works for 5.5 as well ;-)
  - General rule: default is 6, may need an increase with more cores
  - Test: 32-HT/ 32/ 24/ 16cores, spin delay = 6 / 96 :



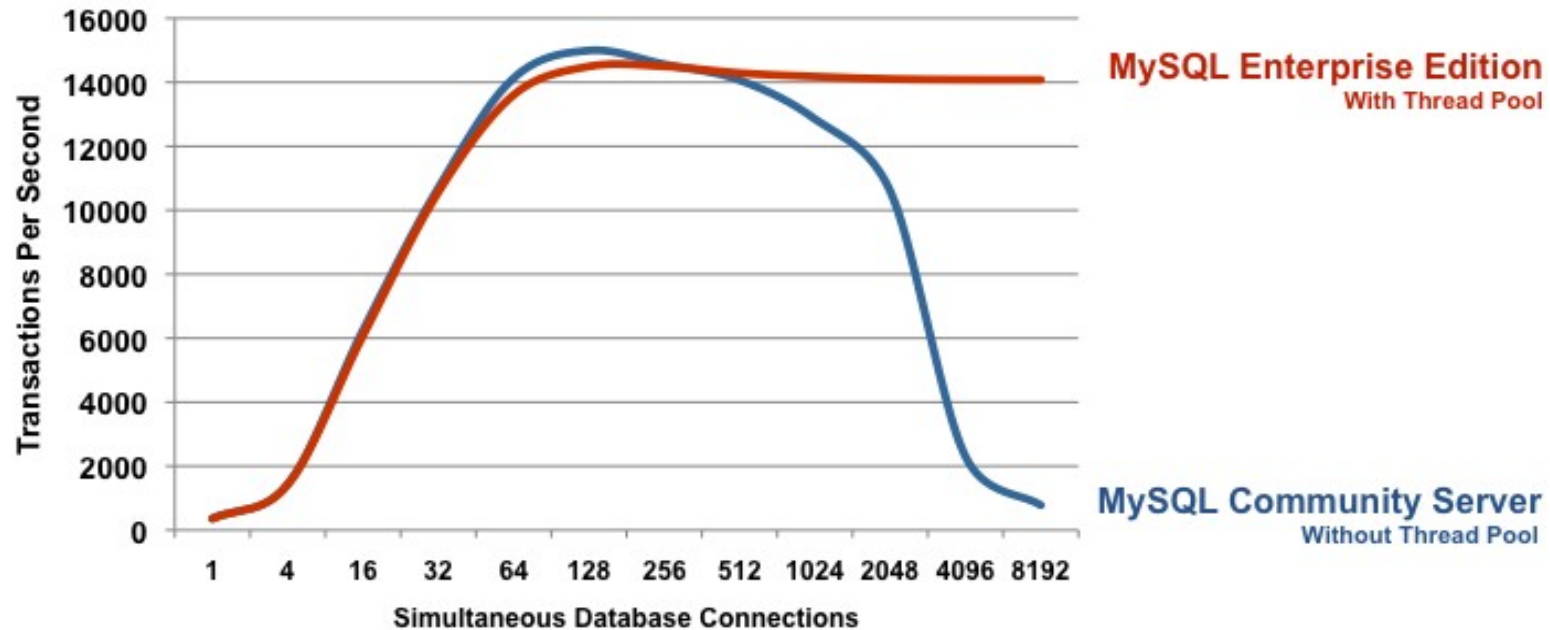
# Thread Pool @MySQL

- None of these solutions will help to increase performance!
  - it'll just help to keep the peak level constant (and you yet need to discover on which level of concurrency you're reaching your peak ;-))
- ThreadPool in MySQL 5.5 and 5.6 is aware if I/O are involved!
  - So, better than innodb thread concurrency setting or taskset
  - May still require spin wait delay tuning!
  - The must for high concurrency loads!
  - May still start to show a difference since 32-128 concurrent users! (all depends on workload)..
  - Keep in mind that OS scheduler is not aware how to manage user threads most optimally, but ThreadPool does ;-)



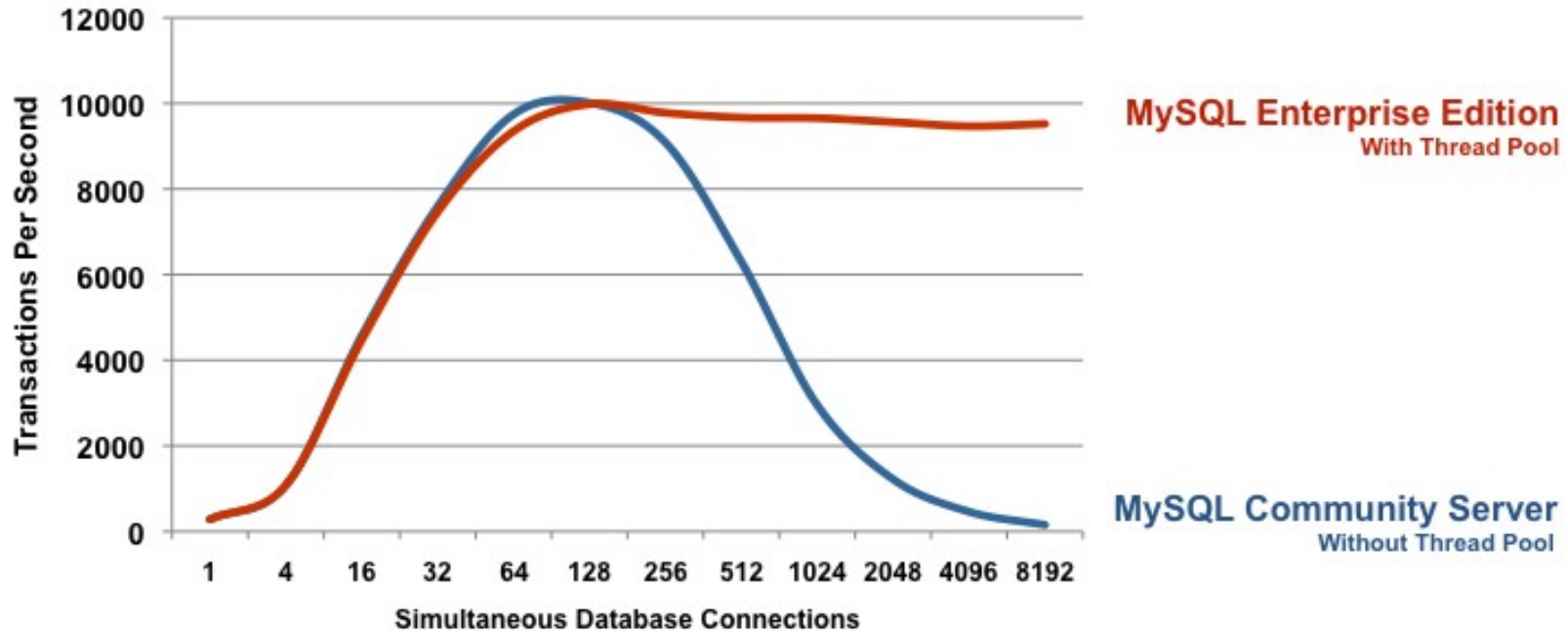
# Thread Pool in MySQL 5.6

- OLTP\_RO:



# Thread Pool in MySQL 5.6

- OLTP\_RW:



# Thread Pool in MySQL 5.7 @Heavy OLTP\_RW



# InnoDB High Concurrency: AHI

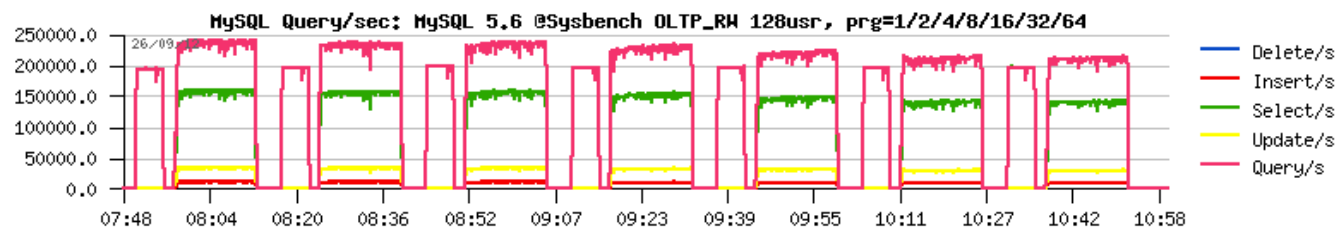
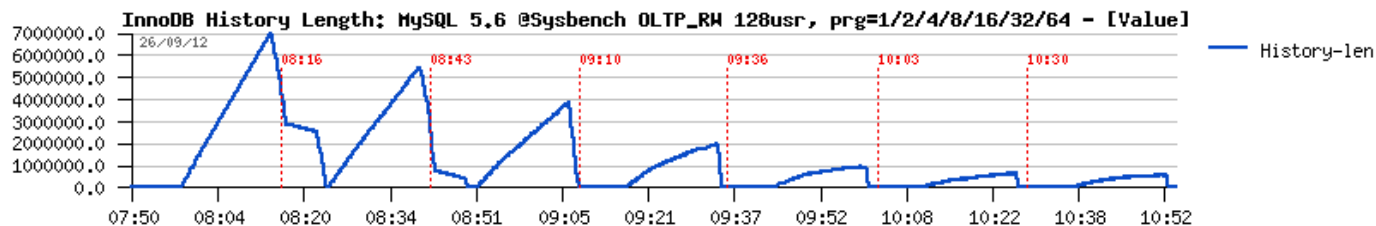
- Adaptive Hash Index (AHI)
  - Helps a lot on Read-Only workloads
  - In fact it helps always until itself become not actively modified
  - AHI contention is seen as its btr\_search\_latch RW-lock contention
  - So, on Read+Write become a huge bottleneck..
  - In many cases on RW the result is better with AHI=off..
  - NOTE: there is still a big mystery around AHI when it's having btr\_search\_latch contention even when there is no changes at all (pure RO in memory).. - expected to be fixed in 5.7 ;-)

# InnoDB Purge

- Purging (similar to Garbage Collecting)
  - Since MySQL 5.5: purge thread
  - Since MySQL 5.6: purge thread(s) (up to 32)
- Having Purge following workload is very important!
  - Ex.: On aggressive RW got 400GB of undo records within few hours(!)
  - Then it took days to reach zero in History Length..
- The main problem is the past – how to dose it?..
  - Since 5.6: with many threads, Purge become auto-stable itself
  - Still missing a dynamic config option to say how many purge threads to run in parallel right now (but it'll be fixed soon ;-))

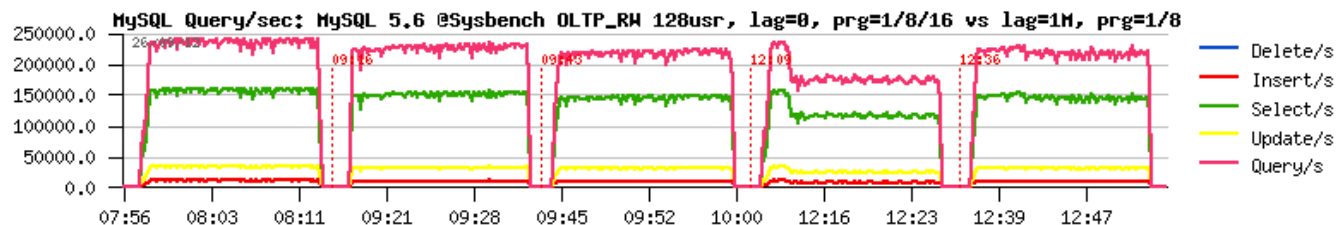
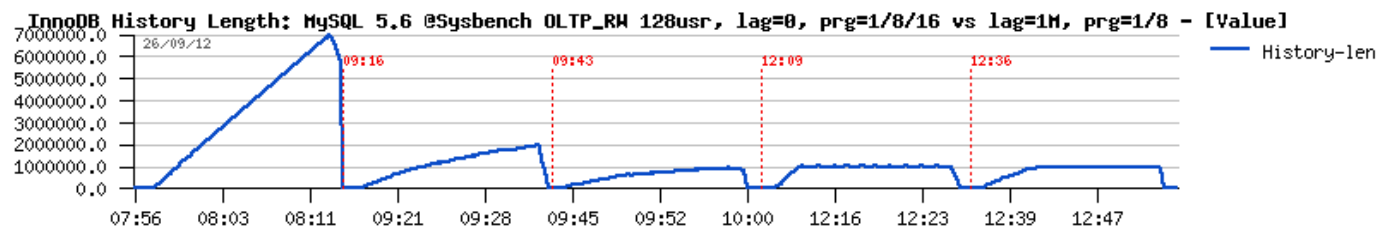
# InnoDB : Purge improvement in 5.6

- Several Purge Threads :
  - NOTE: activation is auto-magical (I'm serious ;-))



# InnoDB : Purge improvement in 5.6

- Fixed max purge lag code!
  - innodb\_max\_purge\_lag
  - innodb\_max\_purge\_lag\_delay <= configurable!
- Setting innodb\_max\_purge\_lag=1M:

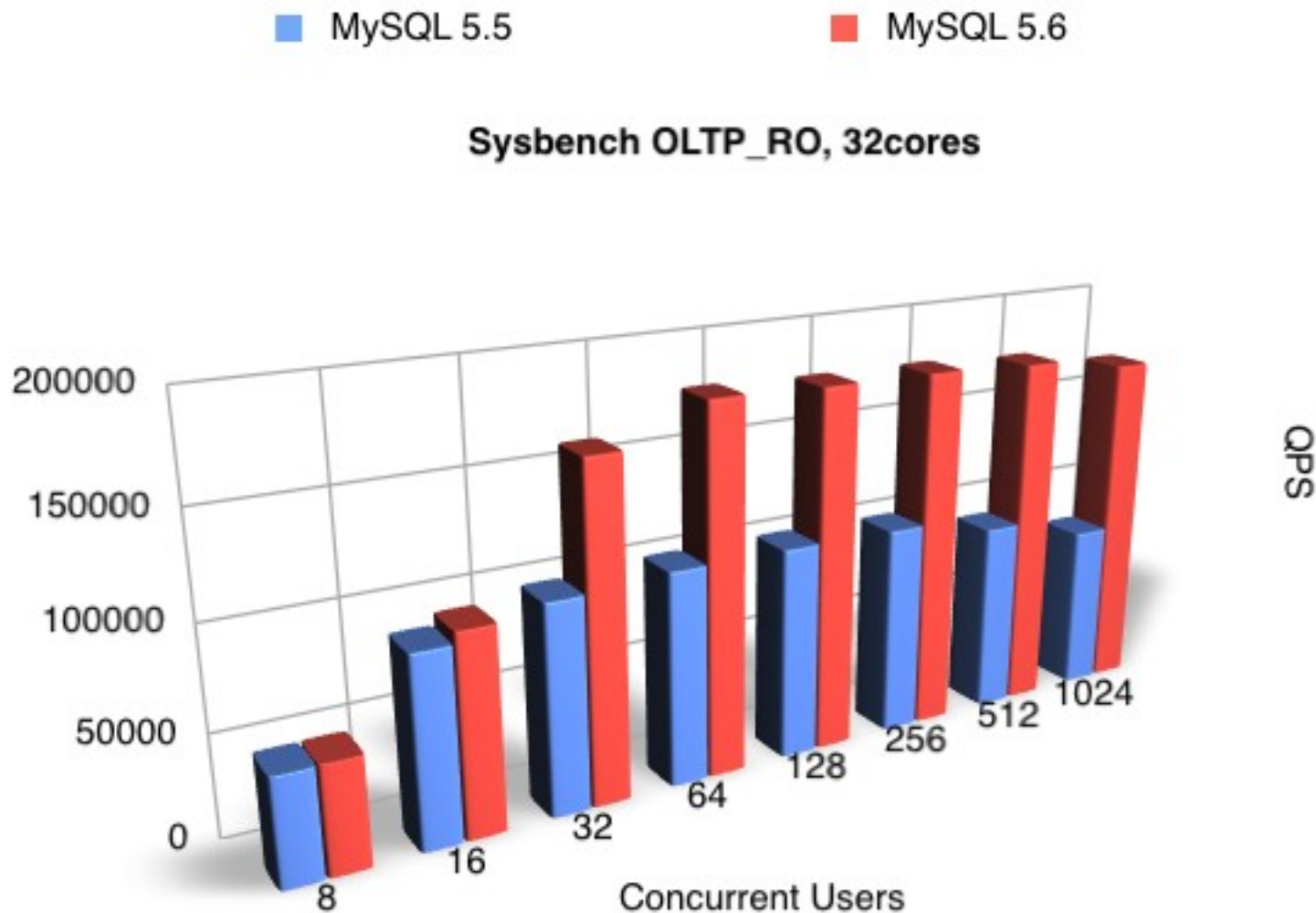


# Testing Apples-to-Apples...

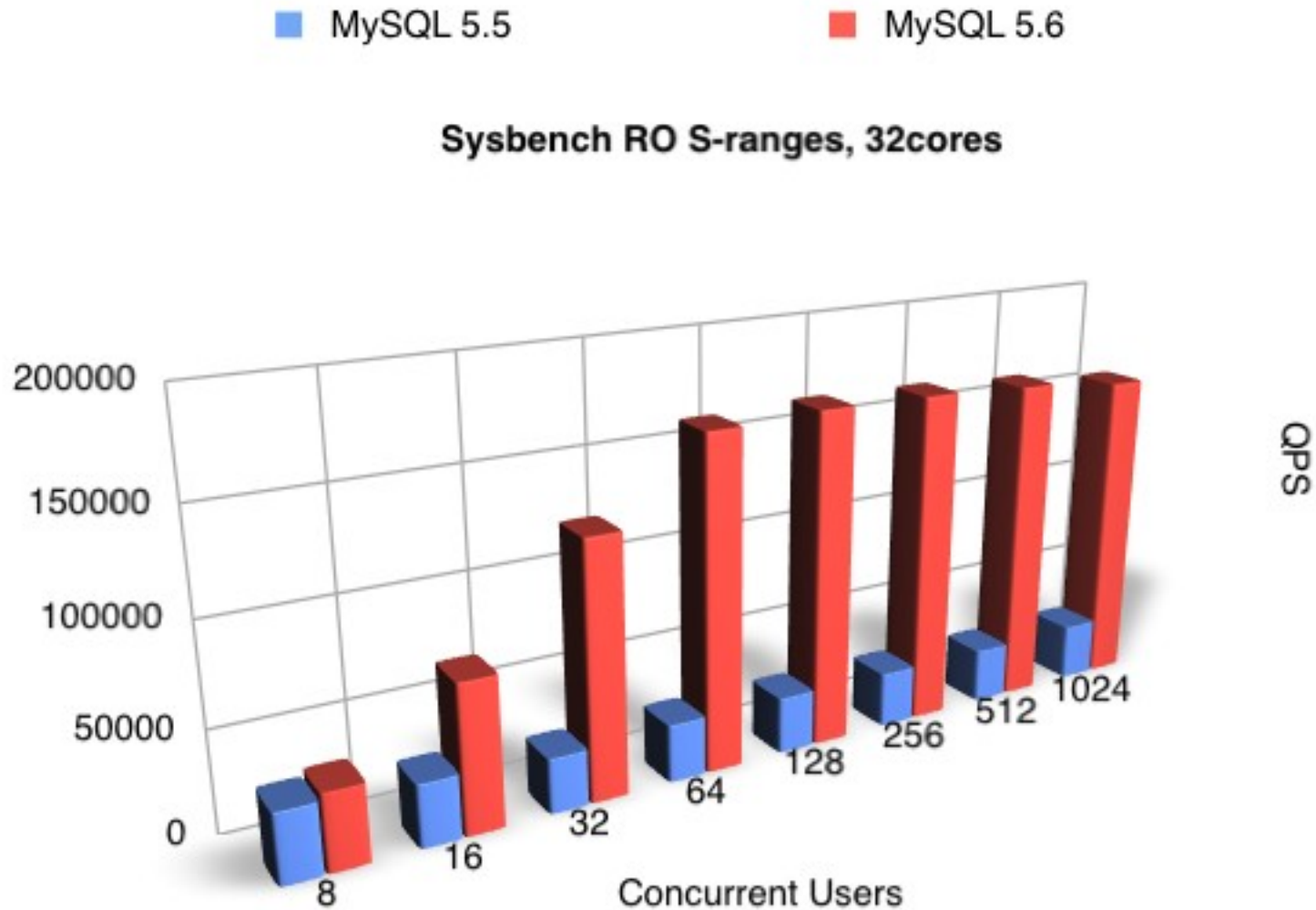
- Comparing MySQL 5.6 vs 5.5 :
  - Don't have G5: dead..
  - Don't have open table cache instances: bad..
  - Don't have improved Adaptive Flushing; bad..
  - Don't have fixed Purge & Lag: danger!..
  - Don't have binlog group commit and use binlog: dead..
  - Etc. etc. etc.
  - NOTE: some “improvement” are also fixes which are making stuff working properly, but coming with additional overhead (like Purge)..
  - NOTE: when comparing 5.6 and 5.5 keep in mind that Performance Schema is enabled by default in 5.6, and not in 5.5, so think to disable it in both (as 5.5 also has a way less PFS instrumentation)..



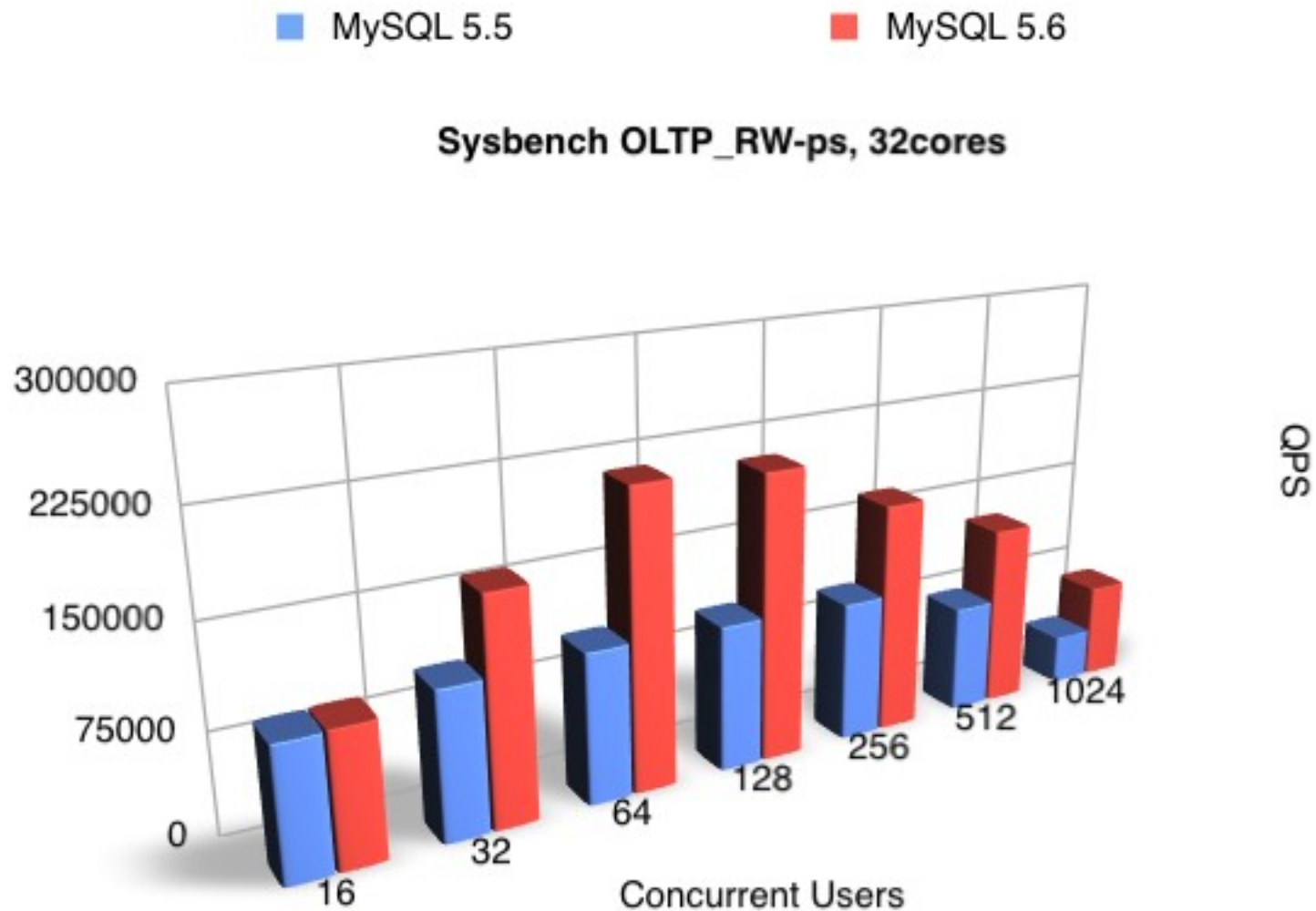
# Hope you did not miss ;-)



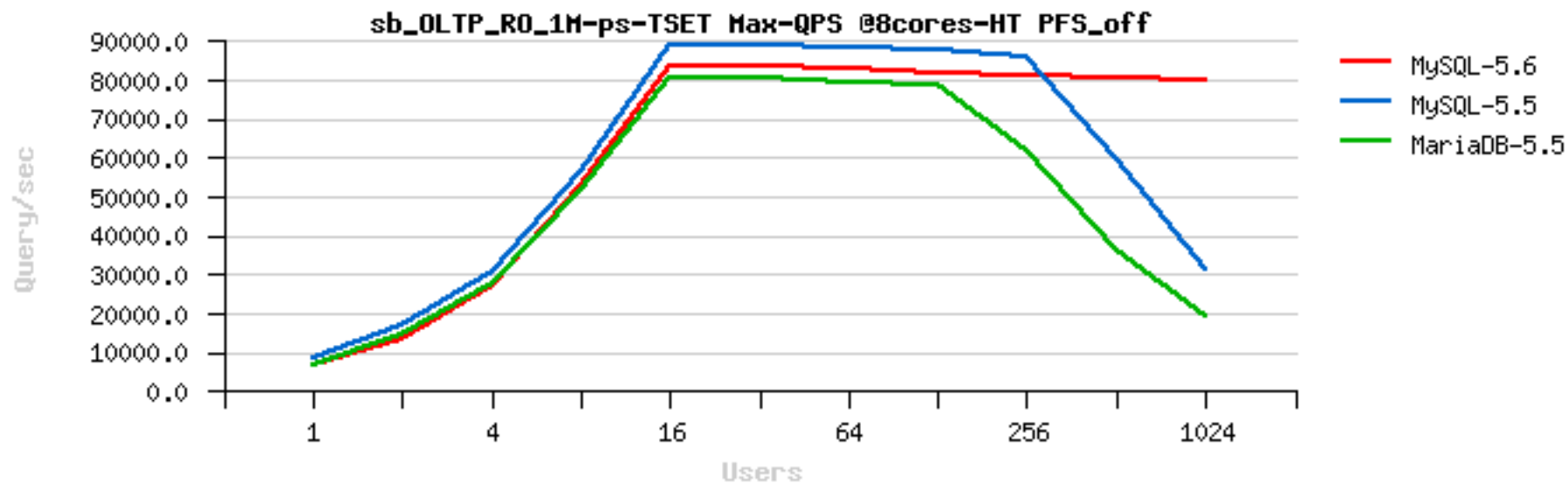
# Hope you did not miss ;-) (2)



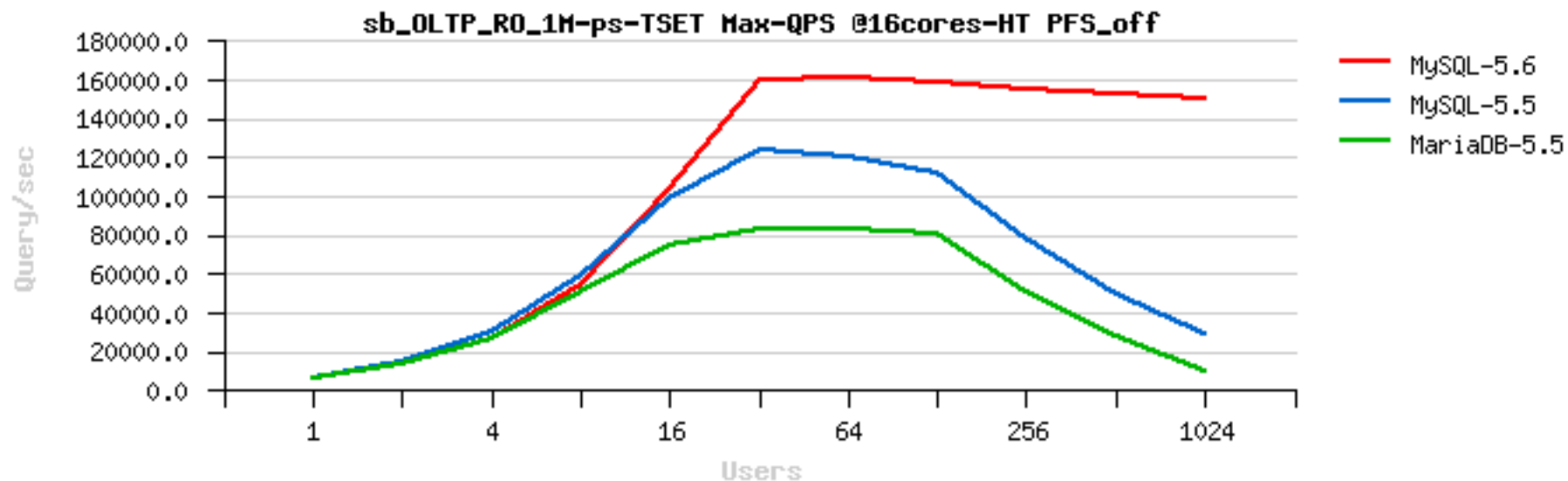
# Hope you did not miss ;-) (3)



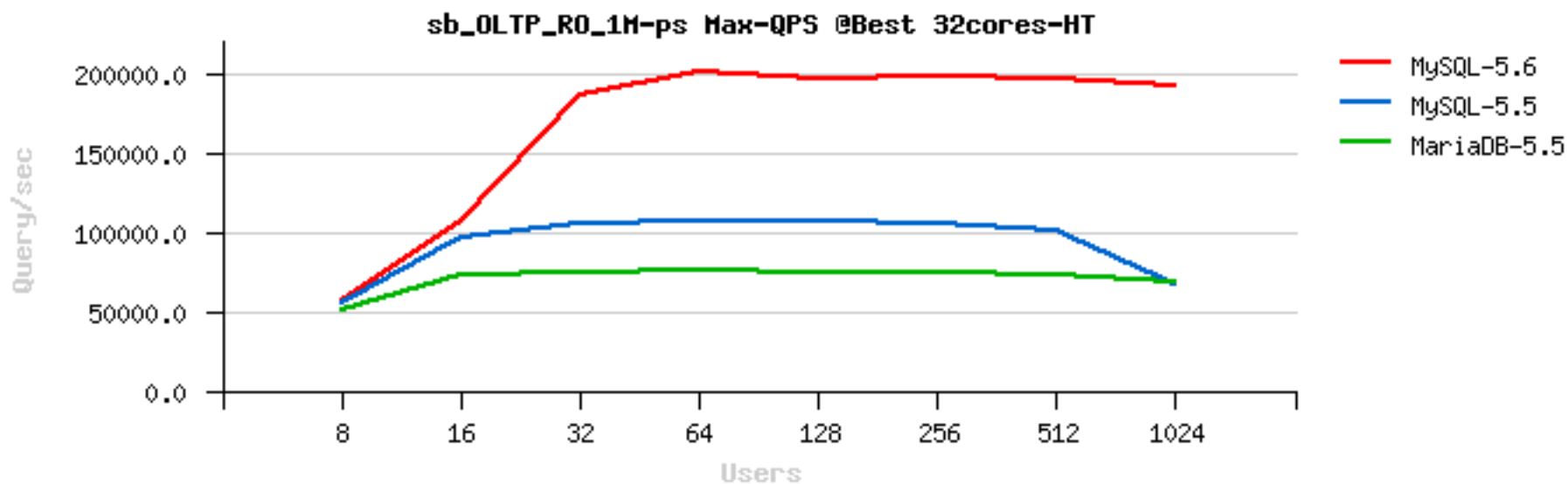
# Sysbench OLTP\_RO @8cores-HT (Apr.2013)



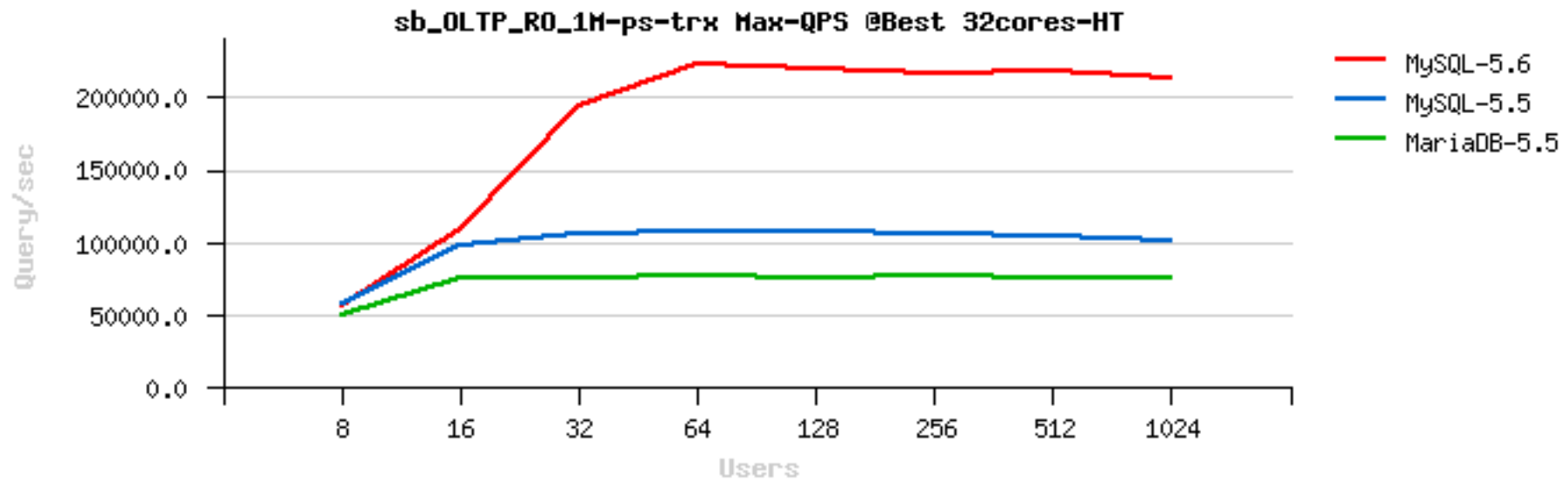
# Sysbench OLTP\_RO @16cores-HT (Apr.2013)



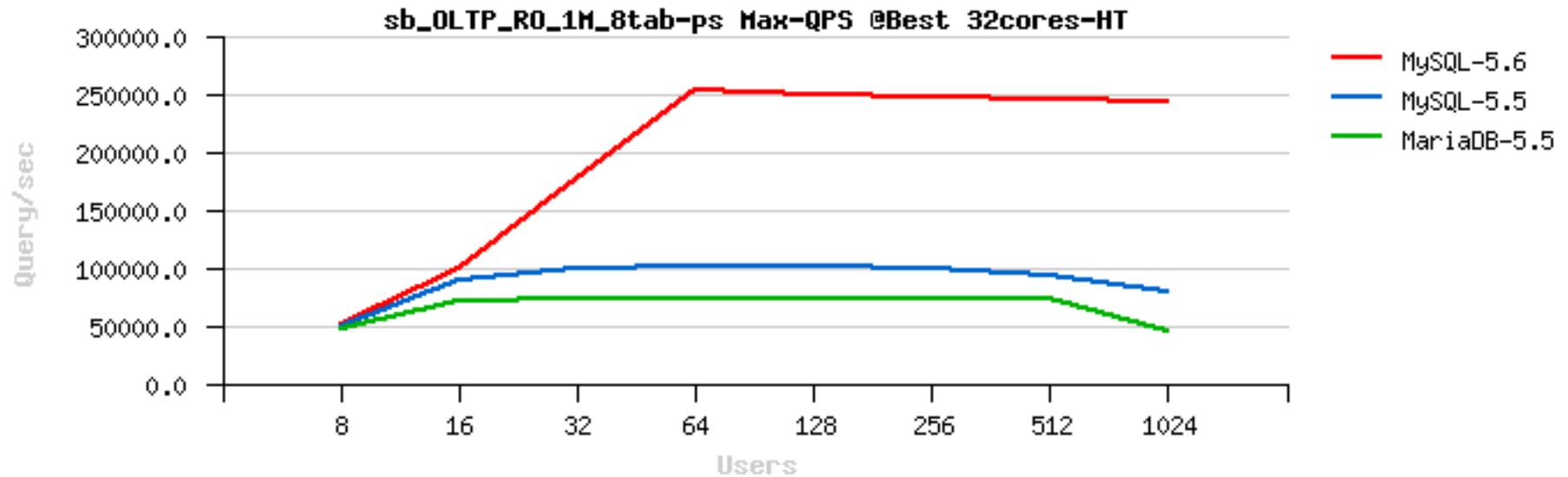
# Sysbench OLTP\_RO @32cores-HT (Apr.2013)



# Sysbench OLTP\_RO-trx @32cores-HT (Apr.2013)

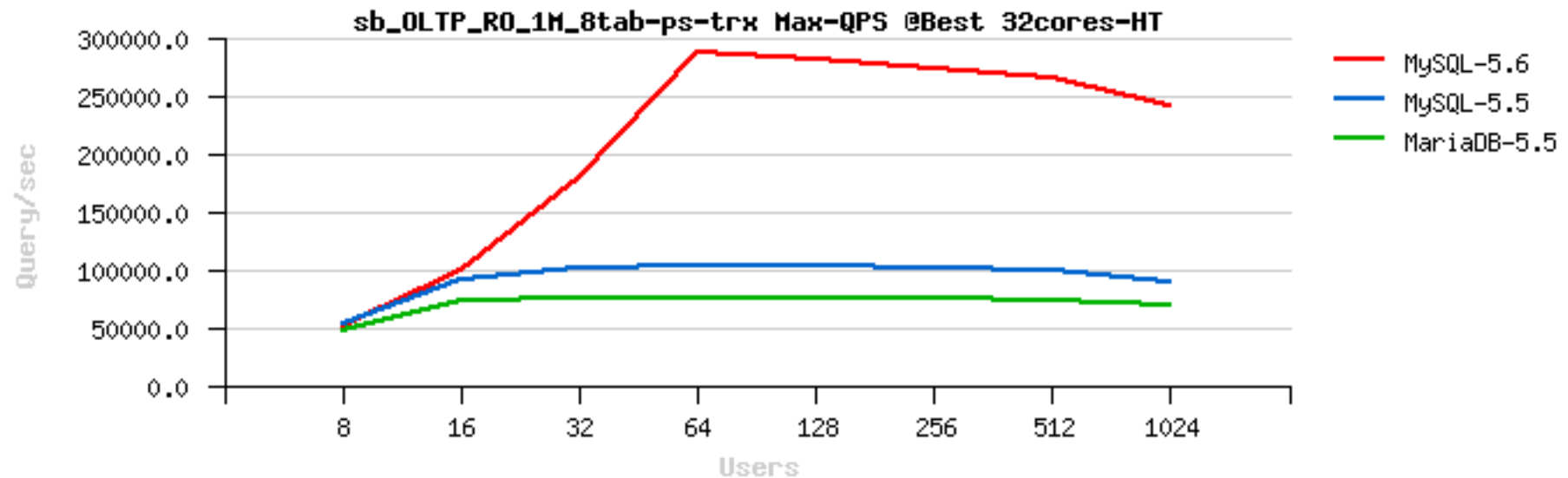


# Sysbench OLTP\_RO 8-tab @32cores-HT (Apr.2013)

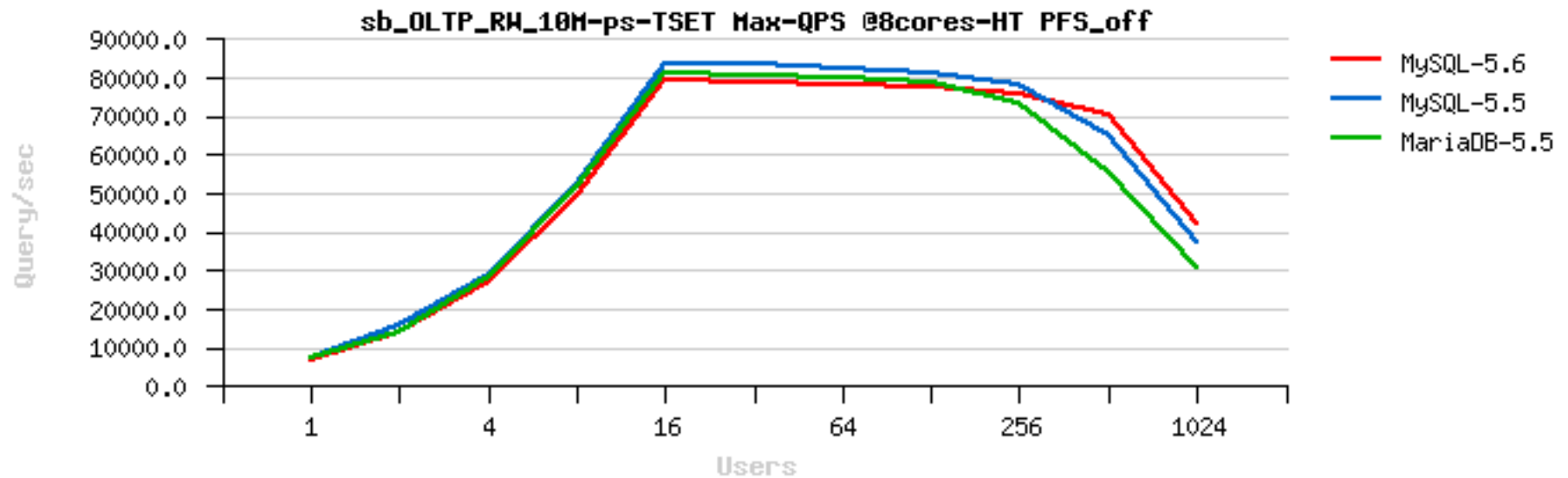




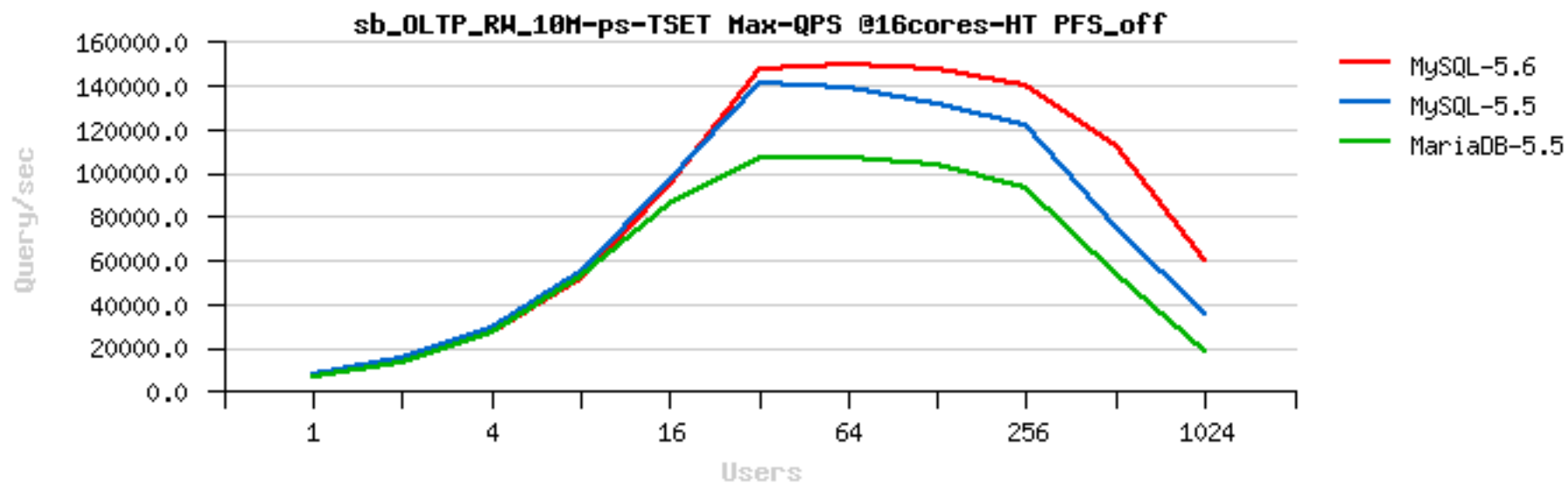
# Sysbench OLTP\_RO-trx 8-tab @32cores-HT (Apr.2013)



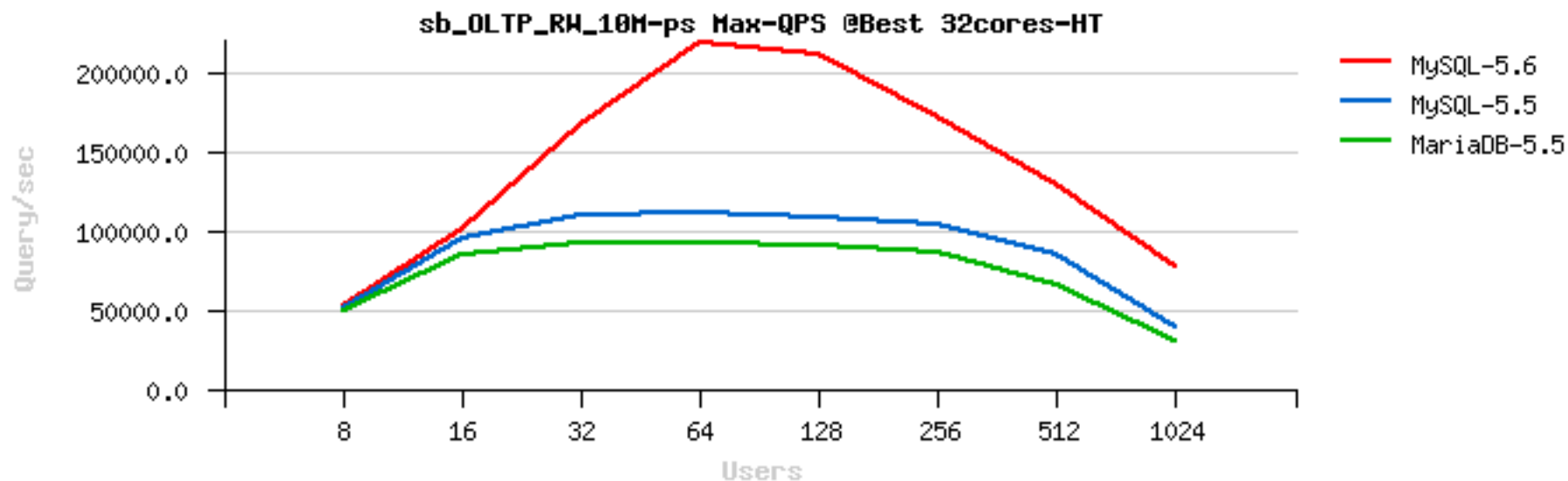
# Sysbench OLTP\_RW @8cores-HT (Apr.2013)



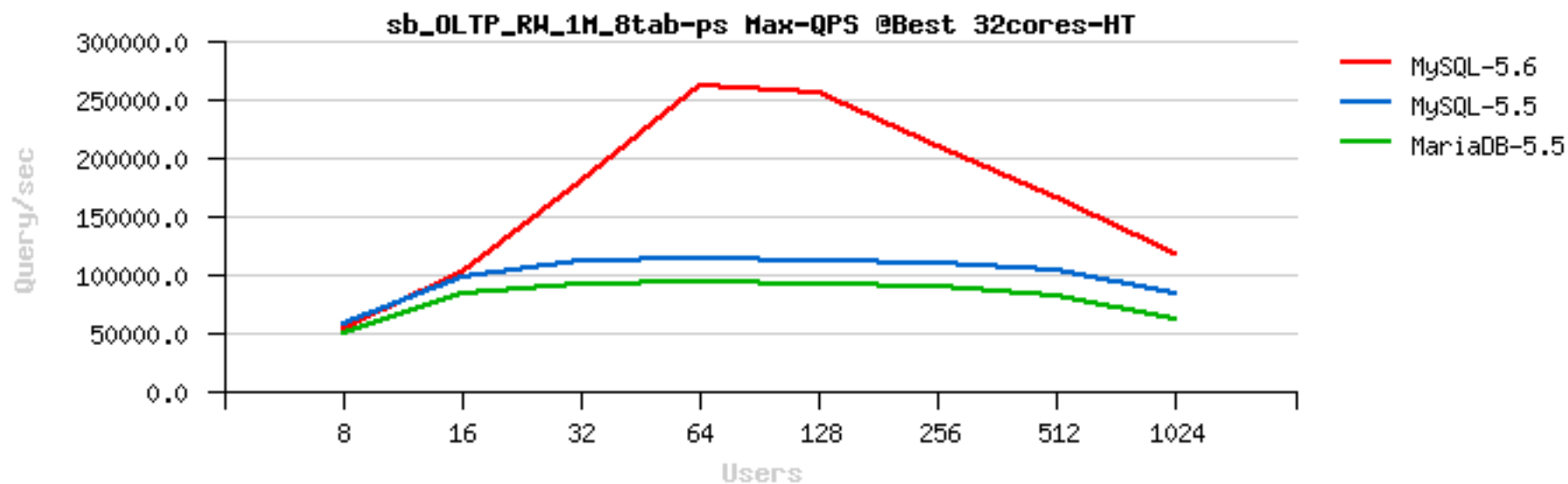
# Sysbench OLTP\_RW @16cores-HT (Apr.2013)



# Sysbench OLTP\_RW @32cores-HT (Apr.2013)



# Sysbench OLTP\_RW 8-tab @32cores-HT (Apr.2013)



# Things are changing constantly, stay tuned ;-)

- MySQL/InnoDB Scalability:
  - 2007 : up to 2CPU...
  - 2008 : up to 4CPU cores
  - 2009 : up to 16CPU cores (+Sun)
  - 2010 : up to 32CPU cores (+Oracle)
  - 2012 : up to 48CPU cores..
  - 2014 : ...?? ;-)
  - NOTE: on the same HW performance is better from version to version!
- InnoDB today:
  - At least x4-8 times better performance than 2-3 years ago ;-)
  - Capable of over ~~100K~~ ~~300K~~ **500K** QPS(!) + FTS & Memcached

## MySQL 5.6: Pending issues

- Index lock..
- Lock\_sys contention..
- Trx\_sys contention..
- MDL scalability..
- Flushing limits..
- LRU flushing..
- Design bug on block locking.. (was here from the beginning)
- Not able yet to use 100% I/O capacity on a powerful storage..
- “Mysterious” contentions on dbSTRESS..
- etc..

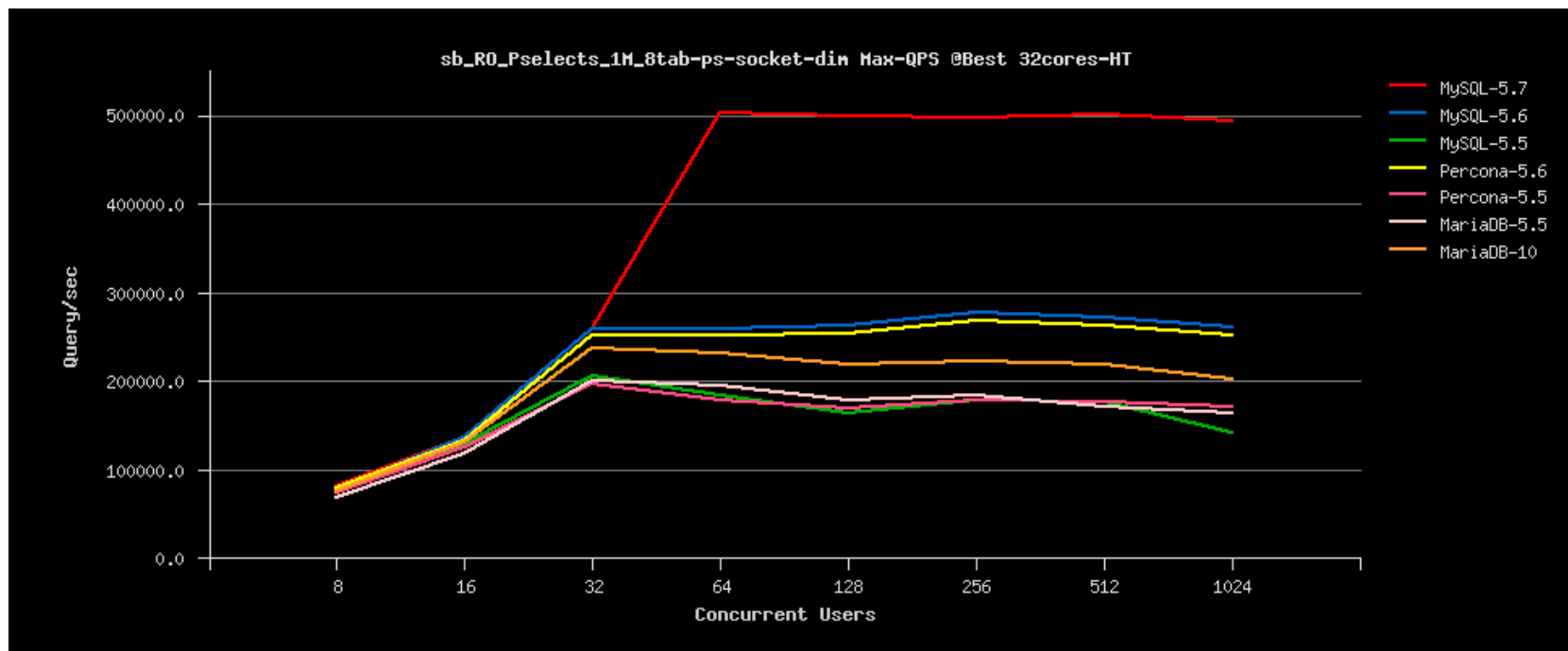
## MySQL 5.7: Work in progress.. ;-)

- Index lock.. <== **fixed** !
- Lock\_sys contention.. <== **lowered** !
- Trx\_sys contention.. <== **improved a lot !!!**
- MDL scalability.. <== in progress..
- Flushing limits.. <== in progress..
- LRU flushing.. <== in progress..
- Design bug on block locking.. (was here from the beginning)
- Not able yet to use 100% I/O capacity on a powerful storage..
- “Mysterious” contentions on dbSTRESS..
- Etc.. <== well, **ALL in progress / investigation** ;-)



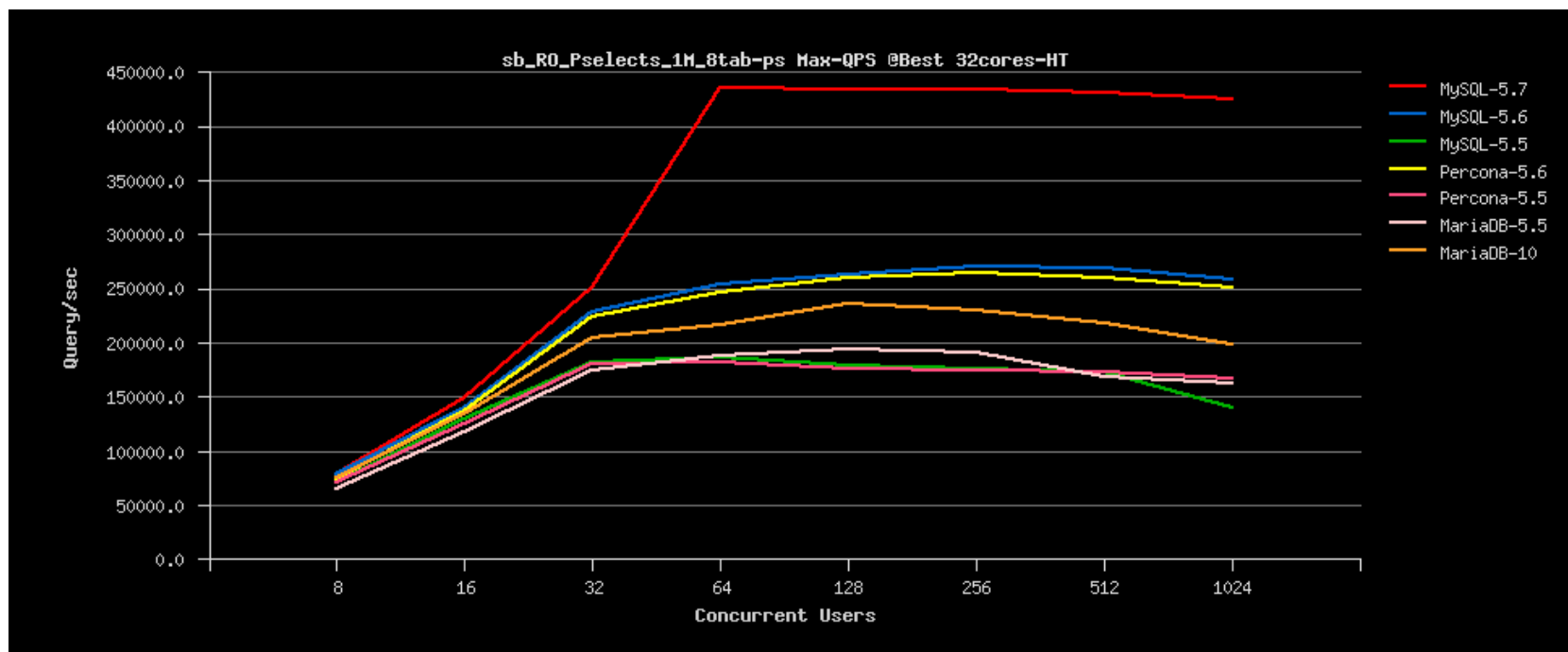
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RO Point-Selects 8-tables: **500K QPS !!!**
  - UNIX socket, sysbench 0.4.8 (older, using less CPU)



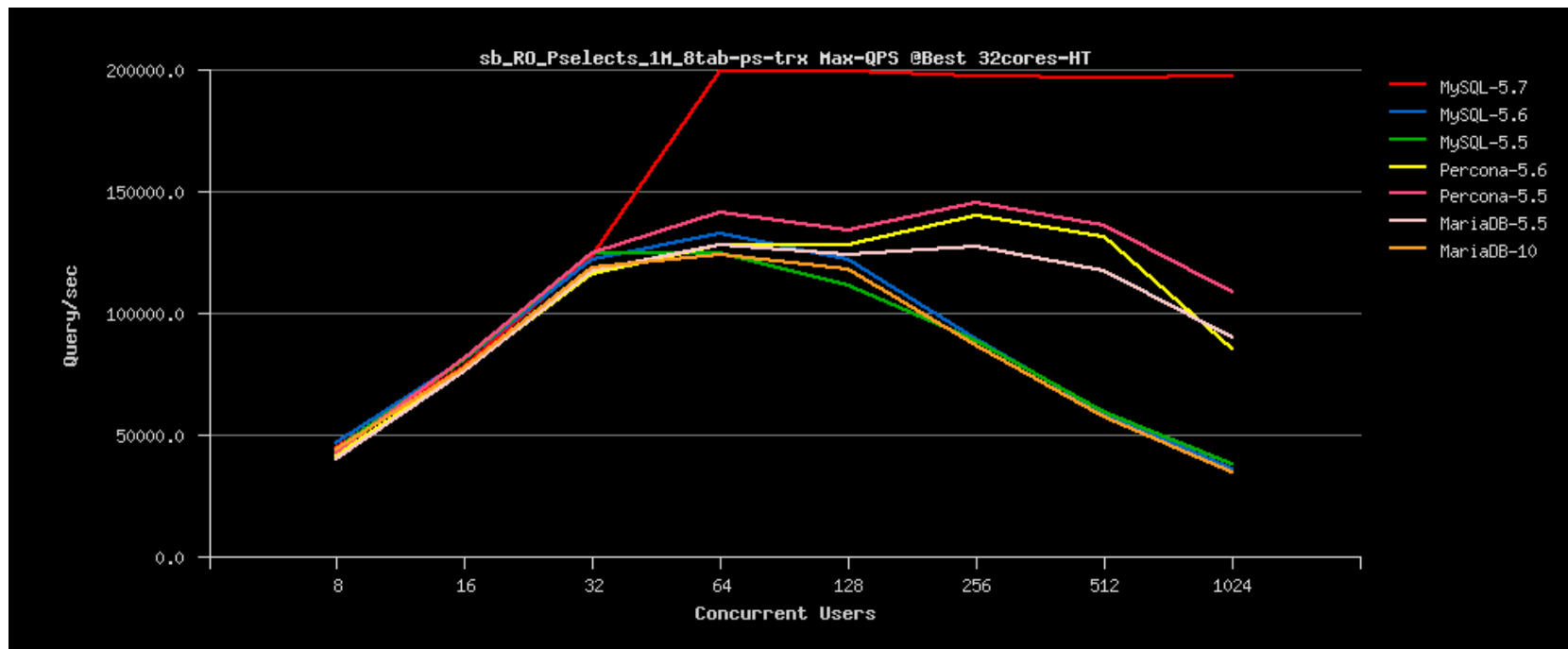
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RO Point-Selects 8-tables: **440K QPS**
  - IP port, sysbench 0.4.13 (“common”, using more CPU)



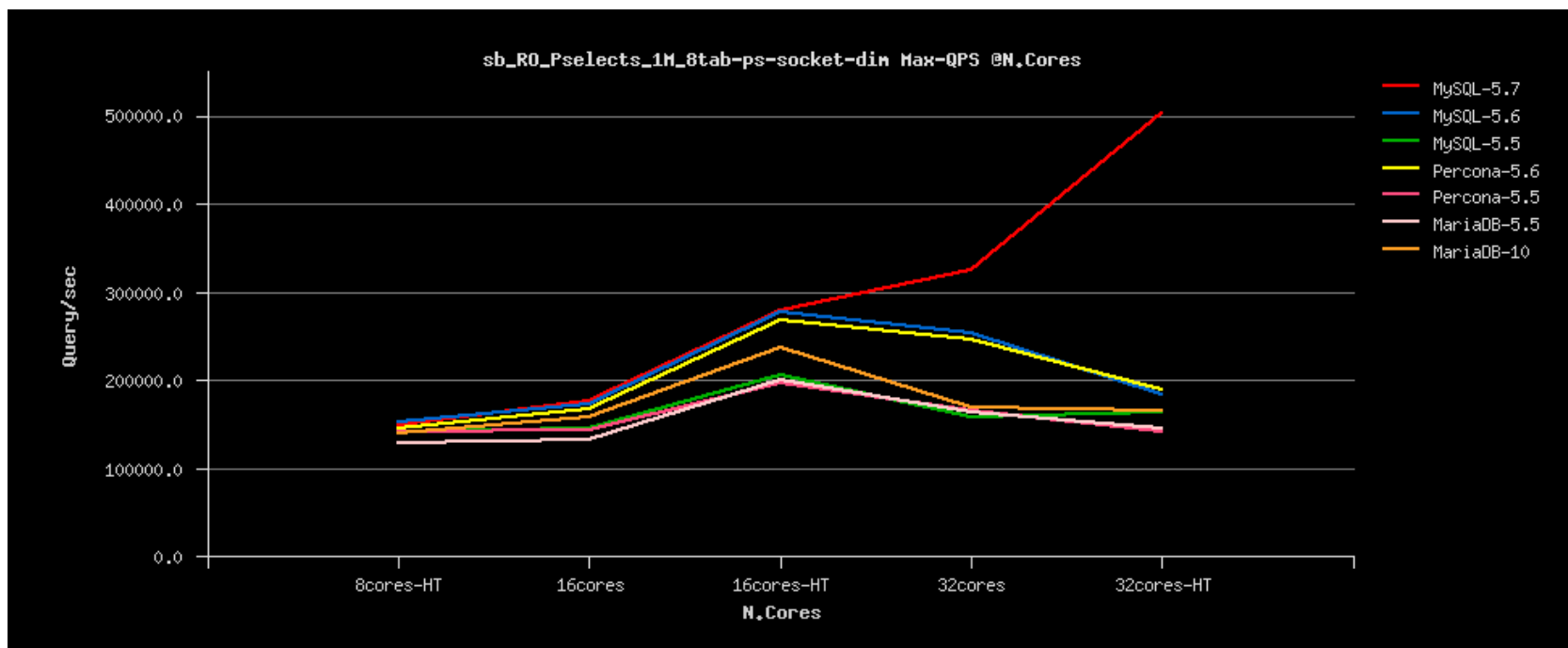
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RO Point-Selects-TRX 8-tables: **200K QPS**
  - IP port, sysbench 0.4.13 (“common”, using more CPU)



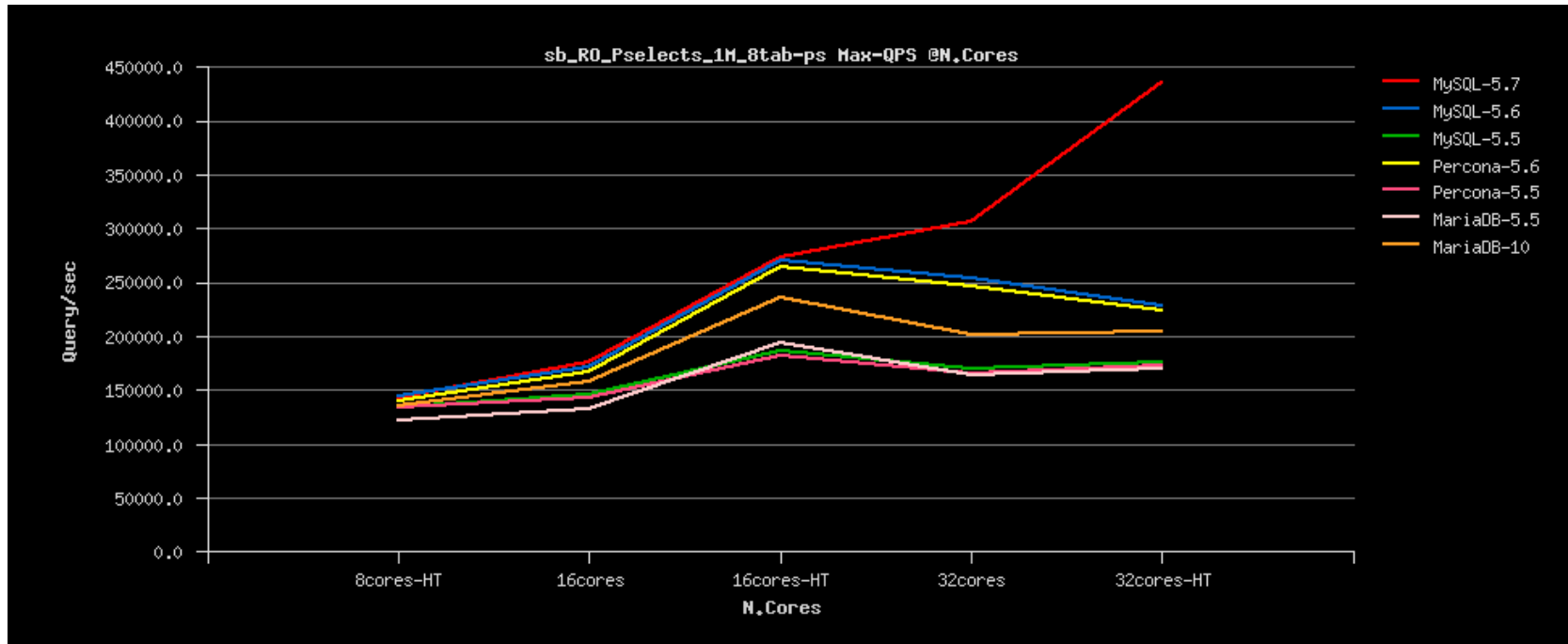
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RO Point-Selects 8-tables: Scalability..
  - UNIX socket, sysbench 0.4.8 (older, using less CPU)



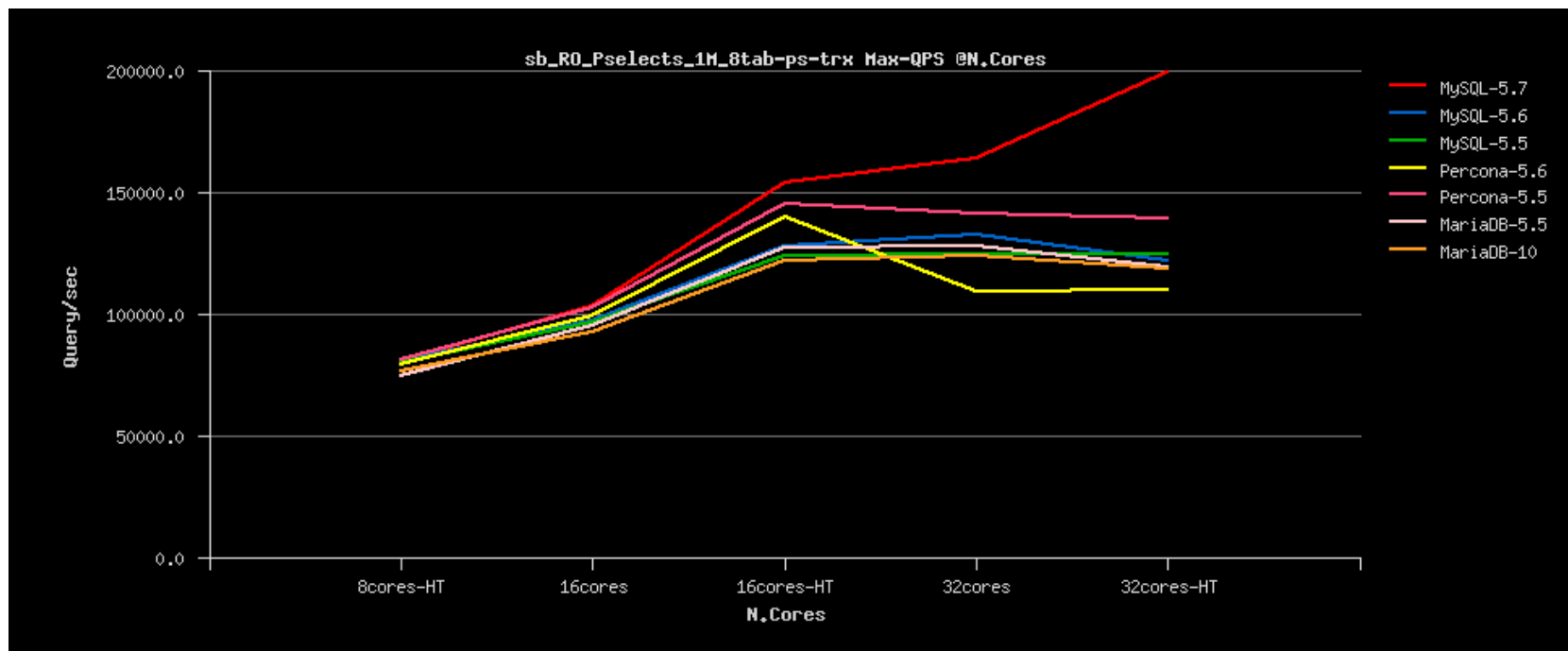
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RO Point-Selects 8-tables: Scalability..
  - IP port, sysbench 0.4.13 (using more CPU)



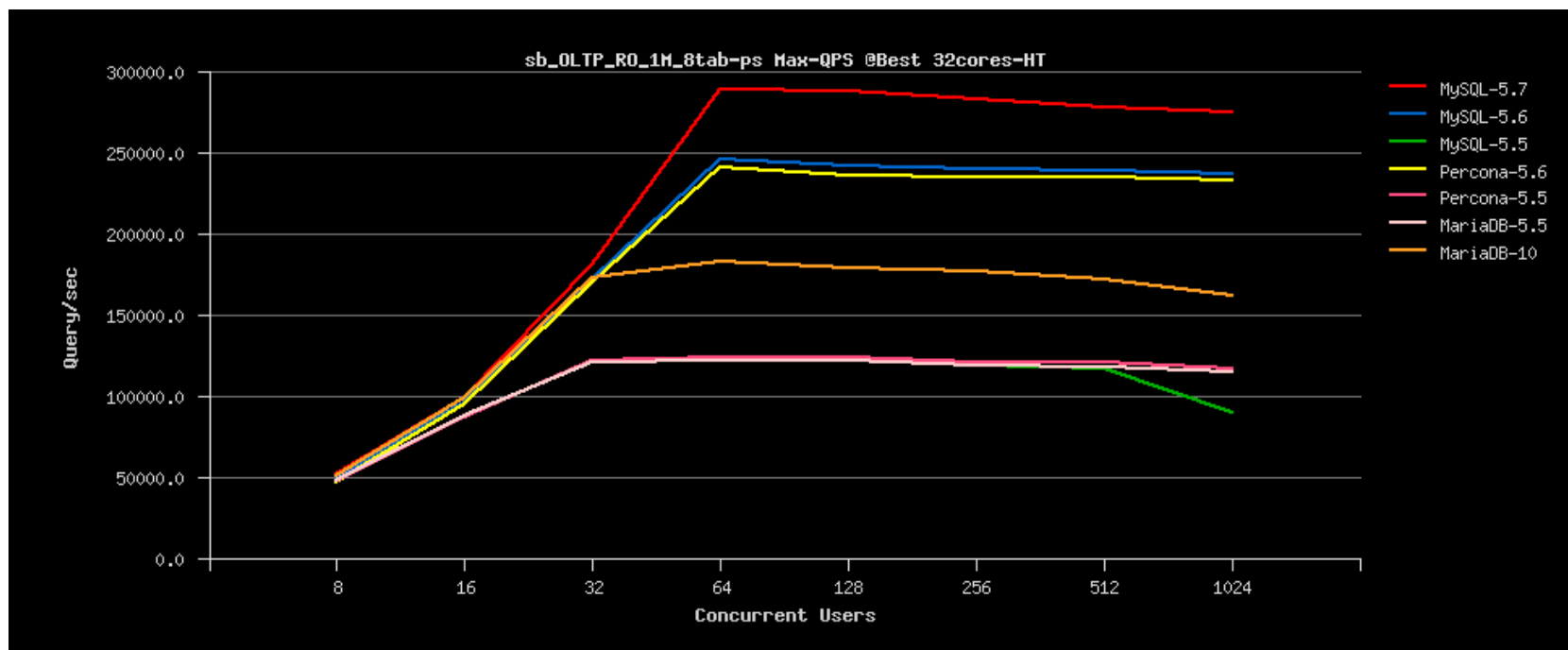
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RO Point-Selects-TRX 8-tables: Scalability..
  - IP port, sysbench 0.4.13 (using more CPU)



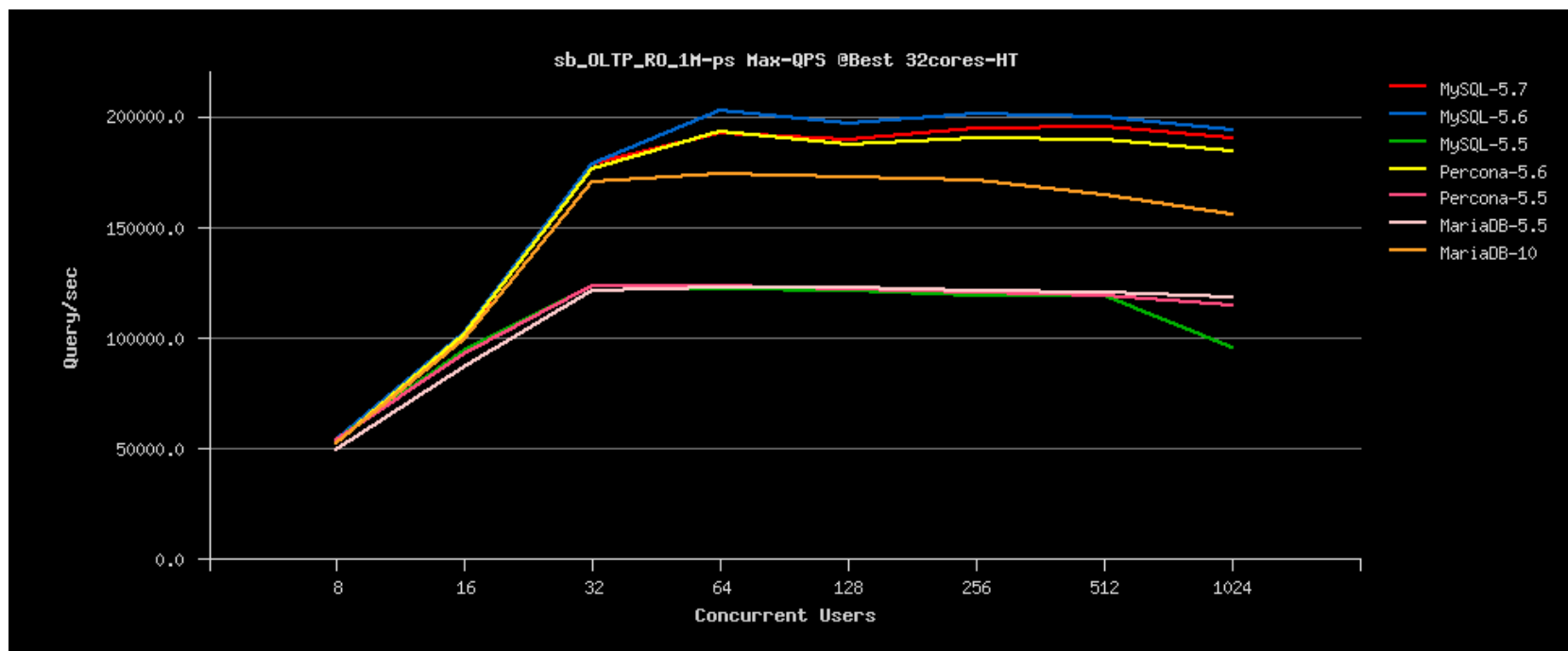
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RO 8-tables: **280K QPS**
  - IP port, sysbench 0.4.13 (“common”, using more CPU)



# MySQL 5.7: DMR2 (Sep.2013)

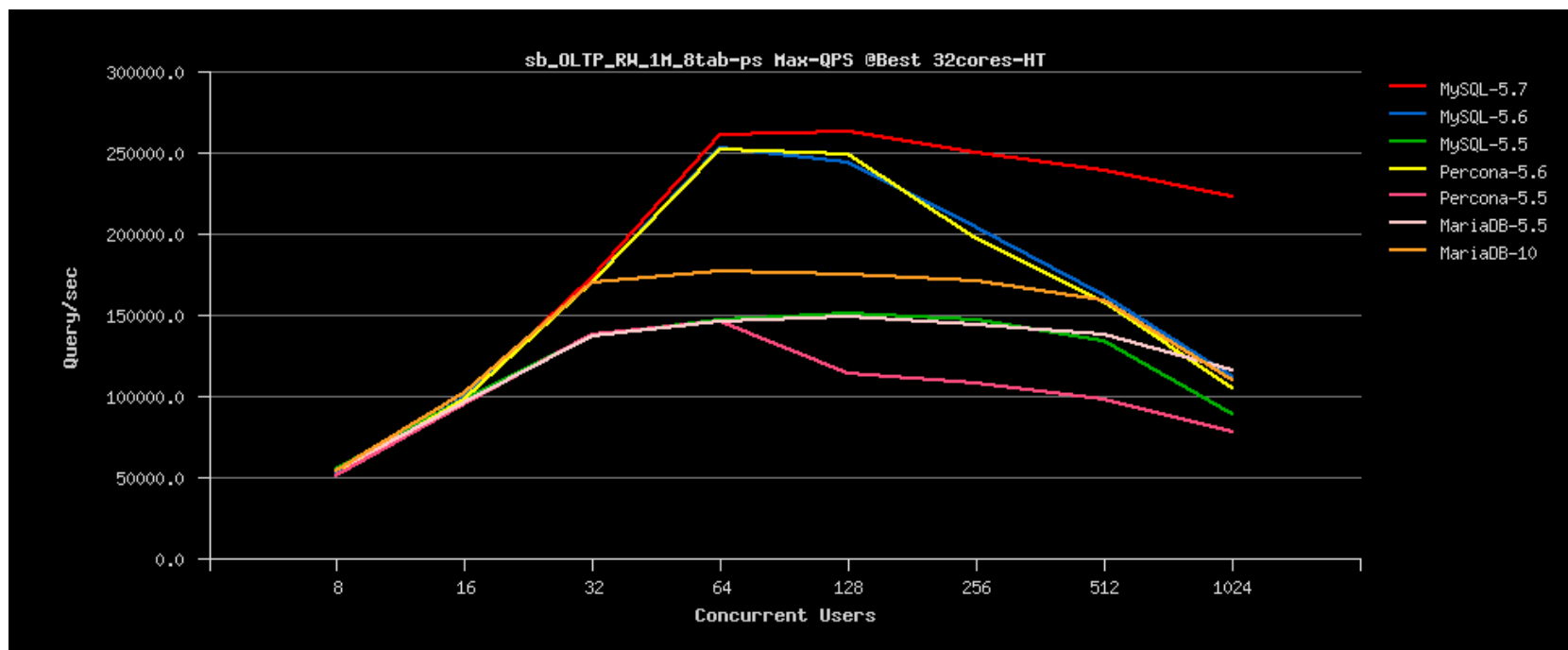
- OLTP\_RO 1-table: lower than 5.6...
  - Due higher MDL contentions, work in progress..





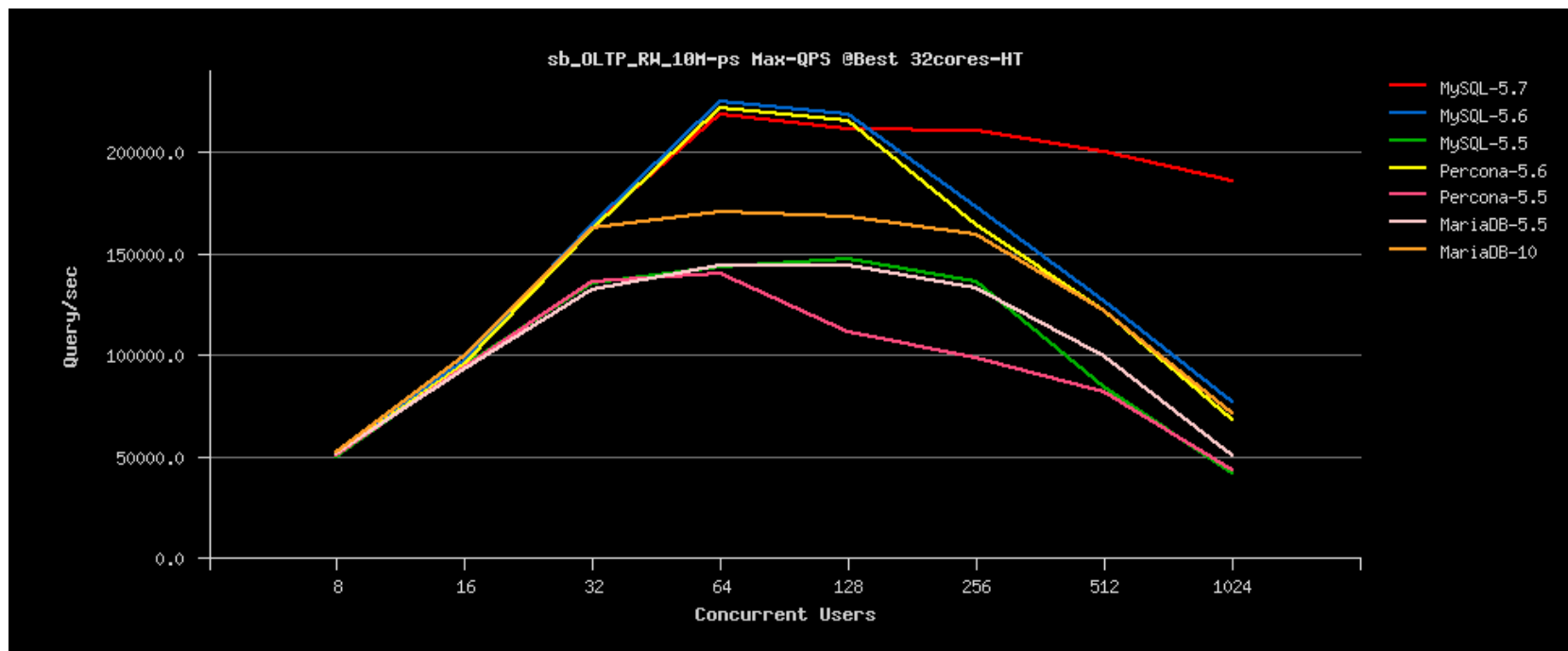
# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RW 8-tables: **265K QPS**
  - IP port, sysbench 0.4.13 (“common”, using more CPU)



# MySQL 5.7: DMR2 (Sep.2013)

- OLTP\_RW 1-table: lower than 5.6
  - MDL contention, work in progress..



## Few words about *dim\_STAT* (as you'll ask ;-))

- All graphs are built with ***dim\_STAT*** (<http://dimitrik.free.fr>)
  - All System load stats (CPU, I/O, Network, RAM, Processes,...)
  - Mainly for Solaris & Linux, but any other UNIX too :-)
  - Add-Ons for Oracle, MySQL, PostgreSQL, Java, etc.
  - MySQL Add-Ons:
    - mysqlSTAT : all available data from “show status”
    - mysqlLOAD : compact data, multi-host monitoring oriented
    - mysqlWAITS : top wait events from Performance SCHEMA
    - InnodbSTAT : most important data from “show innodb status”
    - innodbMUTEX : monitoring InnoDB mutex waits
    - innodbMETRICS : all counters from the METRICS table
  - And any other you want to add! :-)



# THANK YOU !!!

- All details about presented materials you may find on:
  - <http://dimitrik.free.fr> - dim\_STAT, Benchmark Reports
  - <http://dimitrik.free.fr/blog> - Articles about MySQL Performance