

WiredTiger scalability (and RocksDB, too)

Benchmarks and observations

Vadim Tkachenko
Percona, CTO, Performance Enthusiast |



Is performance important?

How the database can be more efficient

Getting more with less

In the end, it's about money:

- Can we use less hardware/cloud instances?

Benchmark

Sysbench for MongoDB

<https://github.com/Percona-Lab/sysbench/tree/dev-mongodb-support-1.0>

What is sysbench?

Micro-synthetic
benchmark

Primitive
operations by
primary key
(collection _id)

Easy to setup and
run

These operations
are building blocks
for more complex
queries

Sysbench operations

Read-Only

- Point lookup
- Range select
- Range aggregation
- Distinct over range

Read-Write

- Read-only plus:
 - Modify non-indexed field
 - Modify indexed field
 - Insert/delete row
 - Transactional for MySQL
 - Find-and-modify for MongoDB to guarantee atomic operation

Sysbench cons

Does not represent any real workload

But keep in mind:

It runs elementary operations that are building blocks for more complex operations

If there is problem in sysbench workload – most likely it will be a problem in real life

The reverse is not necessarily true...

Benchmark setup: hardware

Supermicro blade servers

CPU: Intel(R) Xeon(R) CPU E5-2683 v3 @ 2.00GHz

Total RAM: 256GB

Client and Server identical, connected via 10GB network

Storage, SSD:

- Samsung 850 Pro, SATA – consumer grade, \$500/TB
- Samsung SM863, SATA – server grade, \$600/TB
- Samsung PM1725, NVMe, – high end, \$1,300/TB
 - Big progress from \$10,000+/TB, ~5 years ago

Benchmark setup: configuration

8 collections x 60mln documents each (~100GB of uncompressed data)

Threads from 1 to 1000

Storage engines: WiredTiger and MongoRocks (RocksDB)

Cachesize: from 4GB to 200GB

- Cgroup limits mongod memory with 2x cachesize
- Both WiredTiger and MongoRocks keep cached data in OS cache
 - Need special attention, as otherwise all data (100GB) gets cached in OS

Engines specifics

WiredTiger

- B-Tree based
 - Fast reads
 - Slow random writes
- LSM engine is not documented

MongoRocks

- Based on RocksDB – LSM tree
 - Fast random writes
 - Slower reads

Ways to scale

Multiple user threads

Increase memory

Storage with better performance

Multiple servers (scale out) – not in this talk

Benchmark methodology

Measurements with 10sec intervals

Throughput and response time for the 10sec intervals

It helps to see trends and stalls

Benchmark length – 5-10 mins

Threads in [1-1000] interval

Cachesize in [4-200]GB interval

A lot of data points

A single number is not interesting

The trend is interesting

Ways to represent data

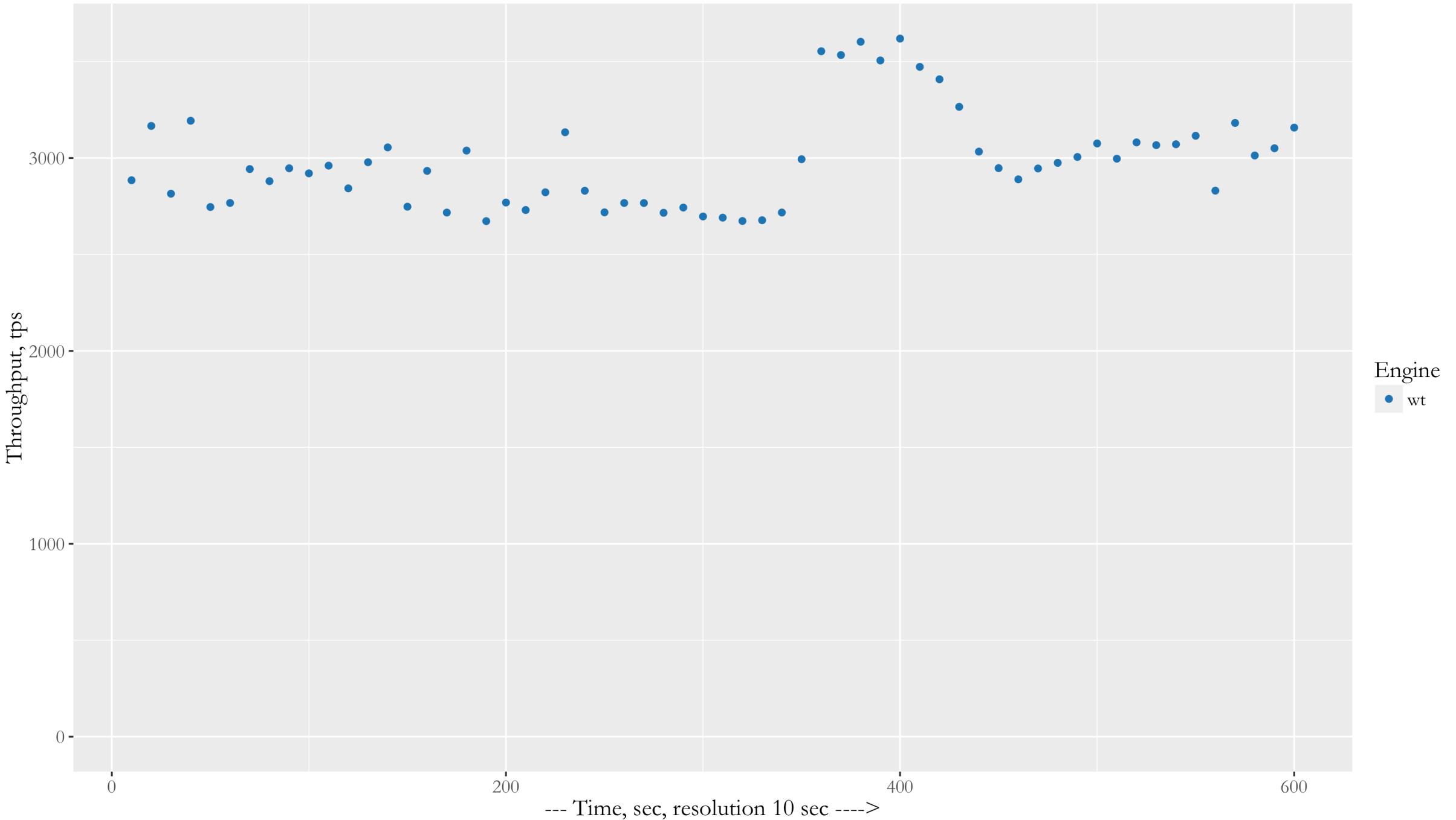
Scatterplots

Trend lines

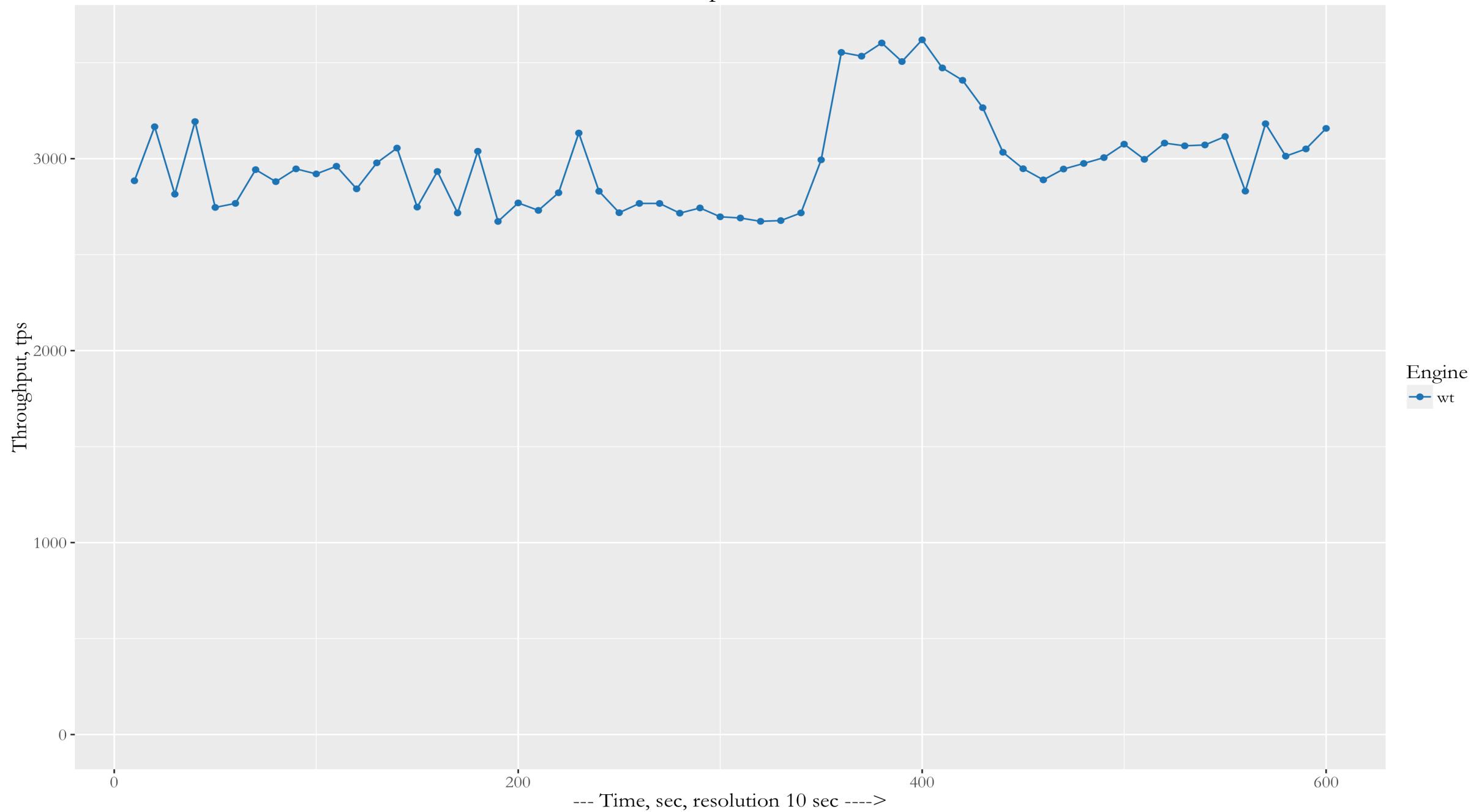
Jitters

Boxplots

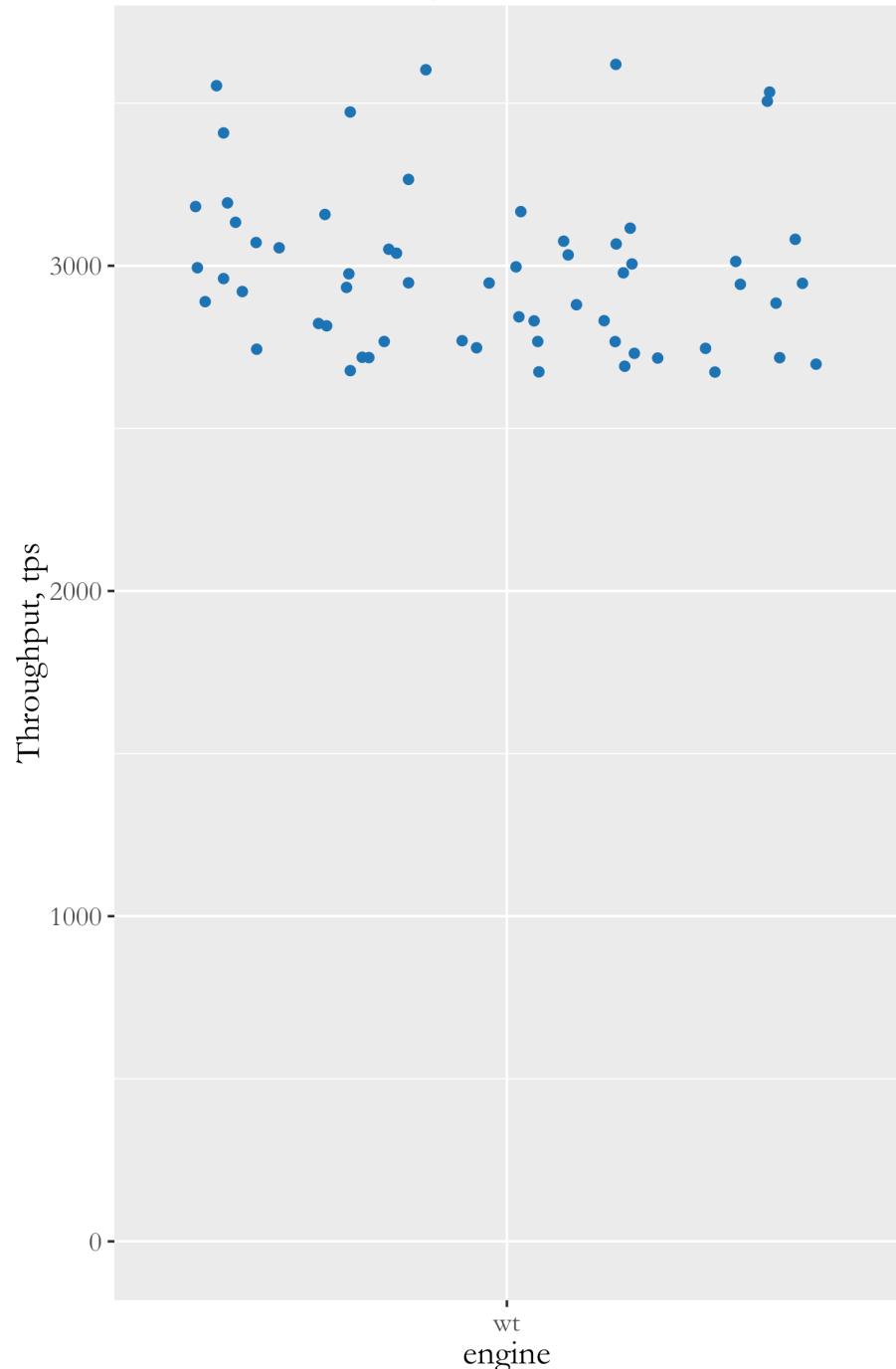
scatterplot



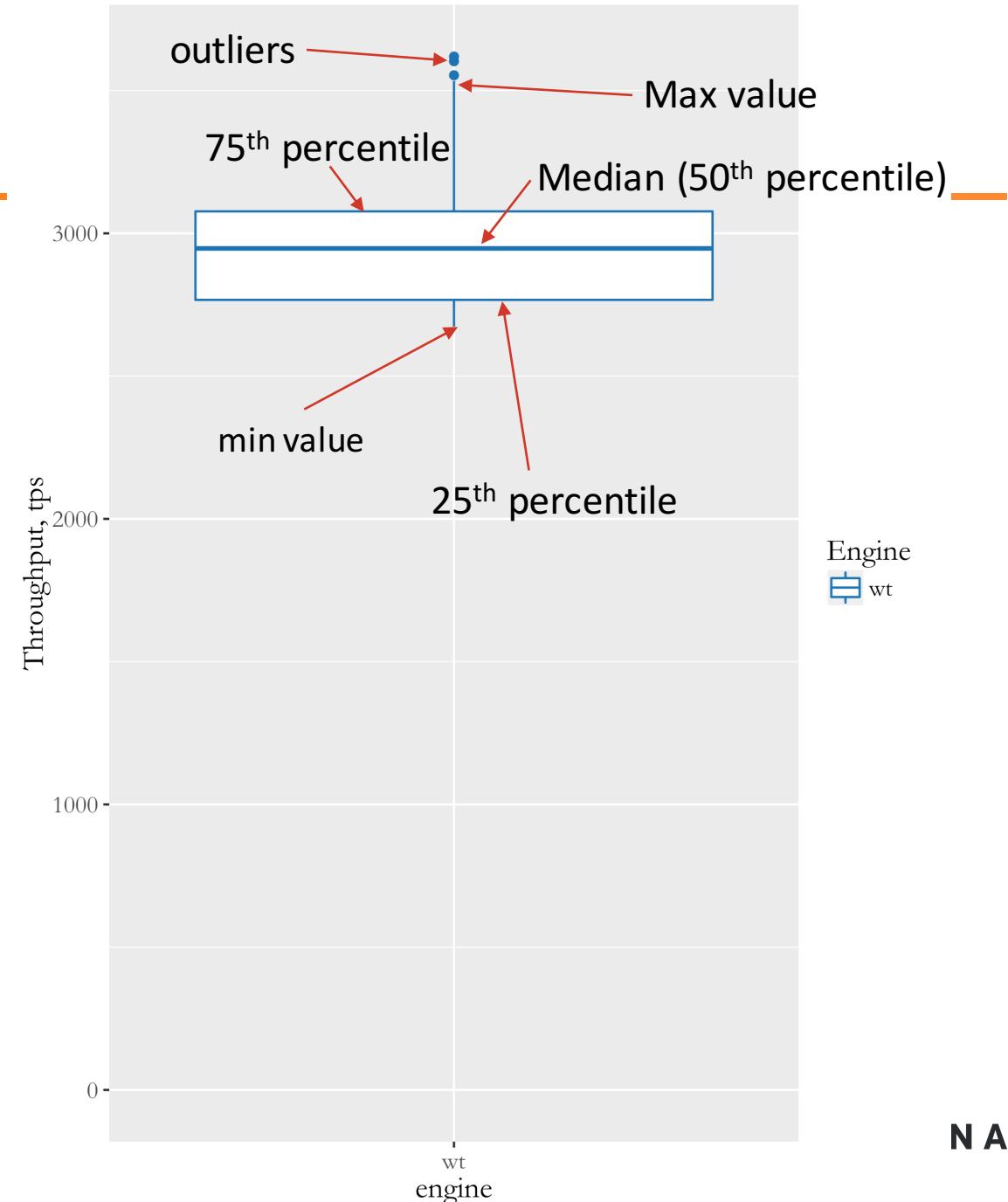
scatterplot + line



jitter plot



boxplot



Results

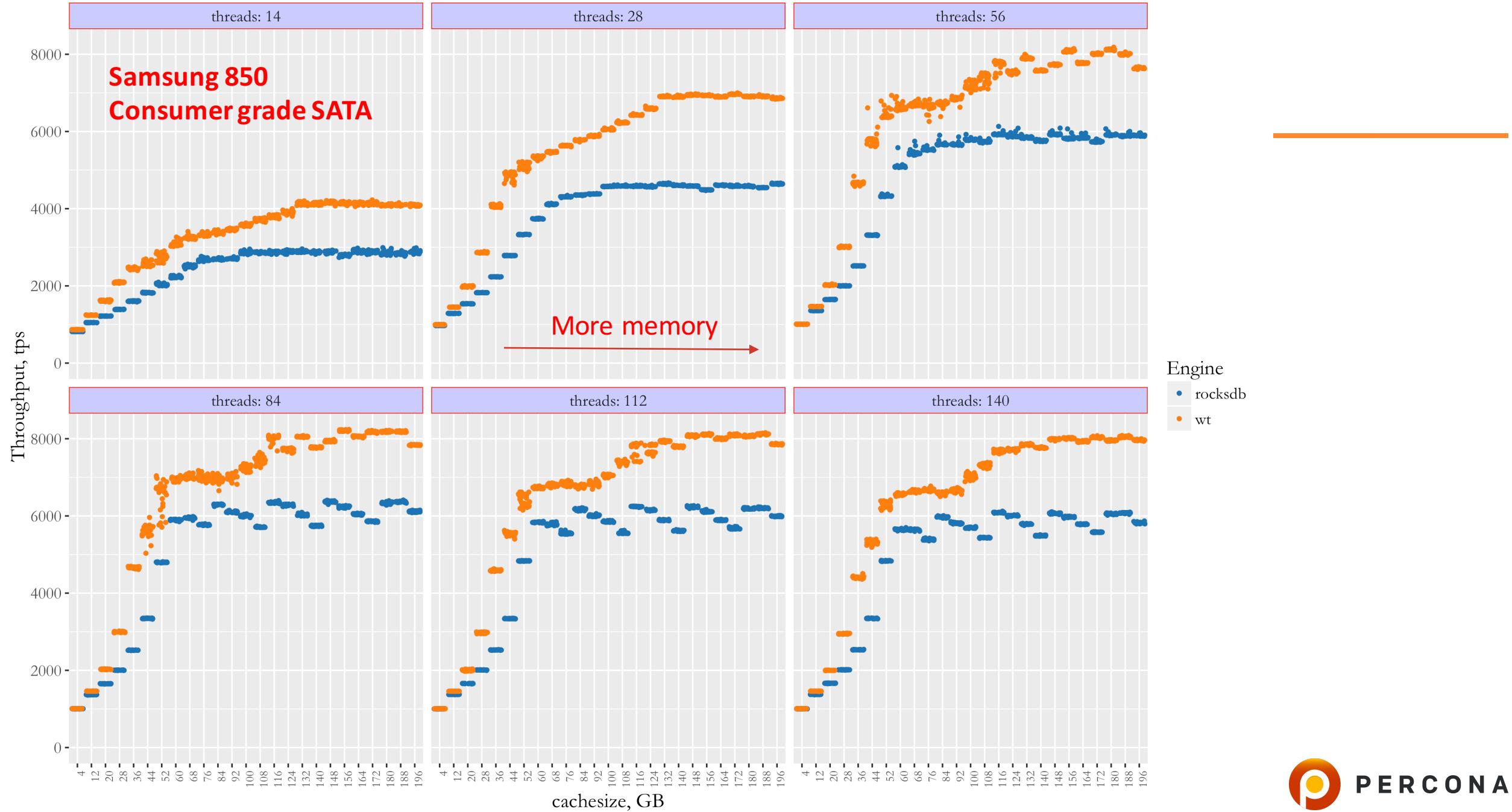


PERCONA

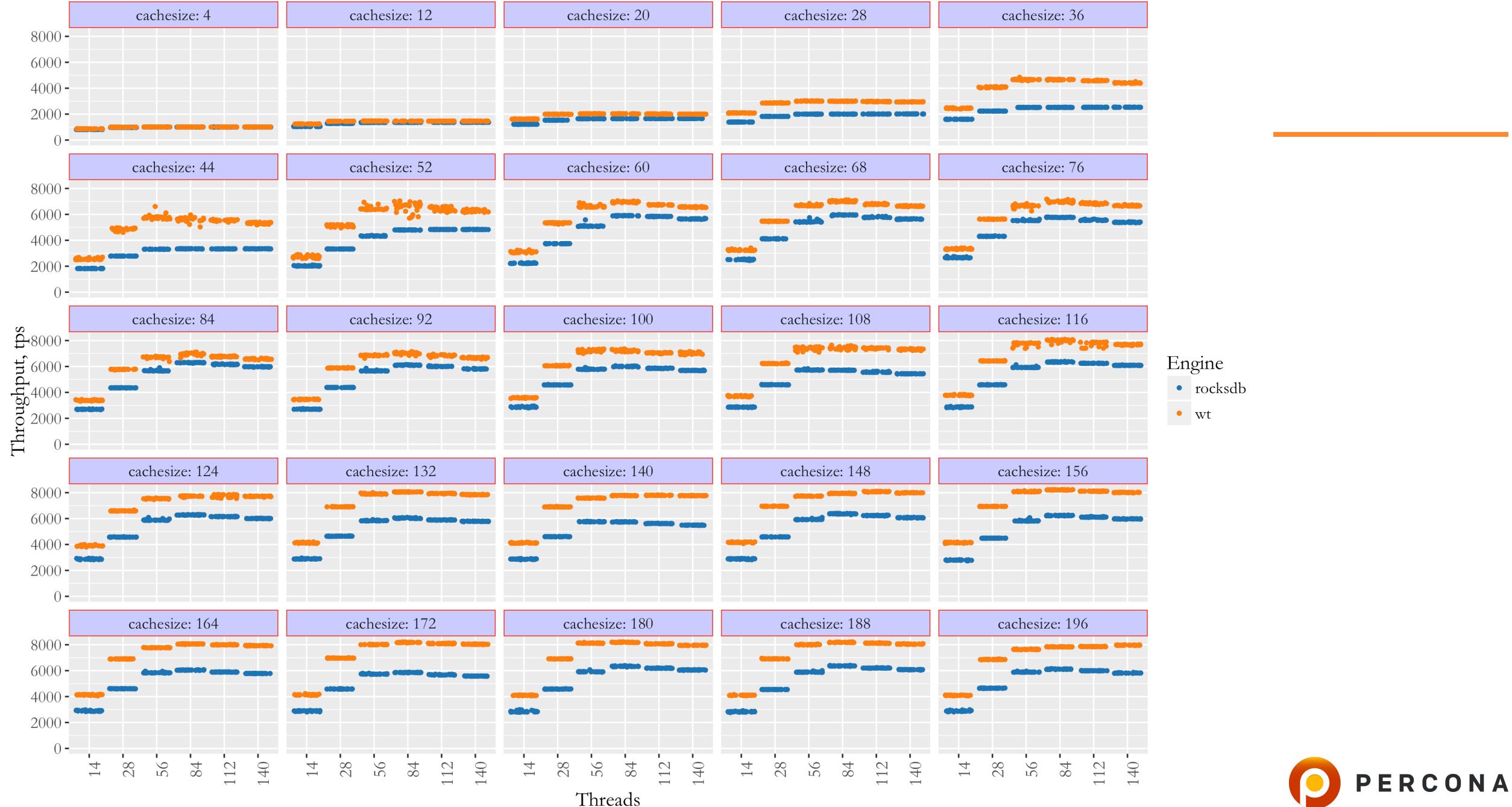
Read-only

- Different cachesizes
 - Small cachesize shows ability to handle growing data
 - Pressure on IO-reads
 - Big cachesize shows the efficiency of internal structures
- Read-only use cases:
 - Reporting
 - Read-intensive components (blog, wiki etc)
- Problems in read-only scalability signal design problems
 - e.g., mutex abuse to access a shared structure

[network] sysbench OLTP READ-ONLY samsung 850



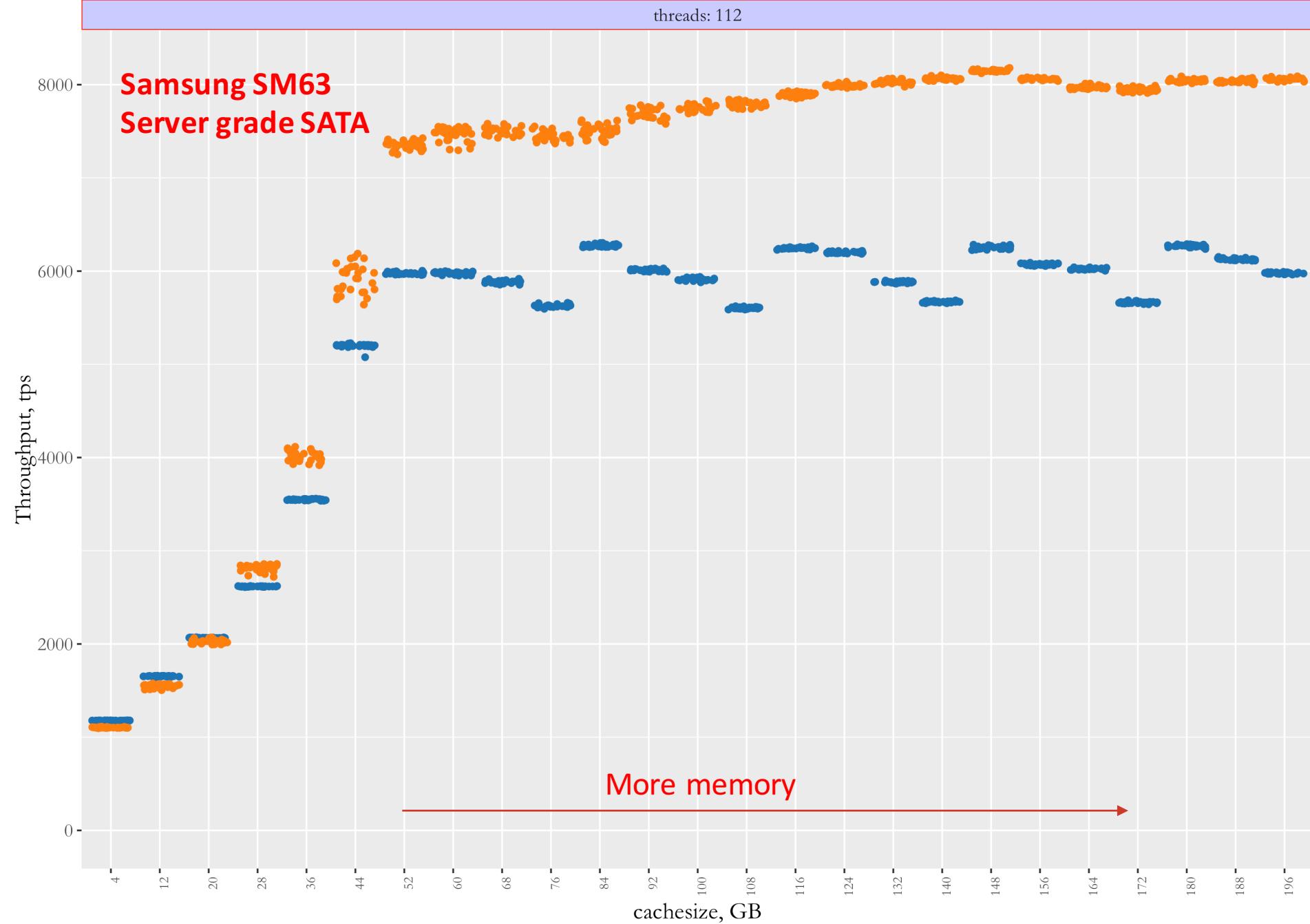
sysbench OLTP READ-ONLY samsung 850



sysbench OLTP READ-ONLY / samsung SM863

threads: 112

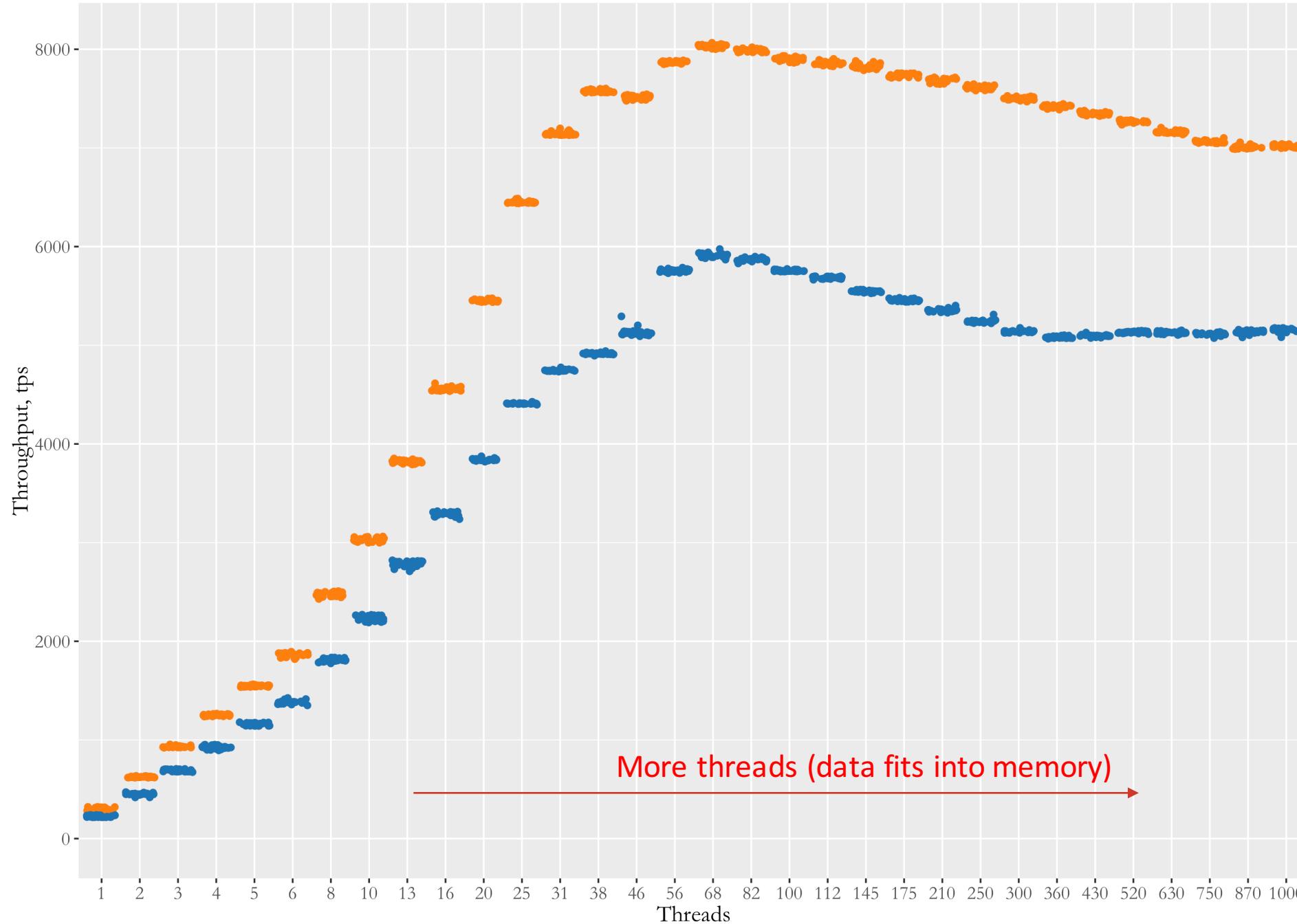
Samsung SM63
Server grade SATA



Multi-threaded scalability Read-Only

- Data fits into memory (200GB cachesize)

[network] sysbench OLTP READ-ONLY in-memory SM863



Read-only observation

- Both engines scale reasonable on read-only
- WiredTiger benefits from more memory, scales better with increasing threads
- WiredTiger seems suitable for read-intensive workloads, a lot of cached data

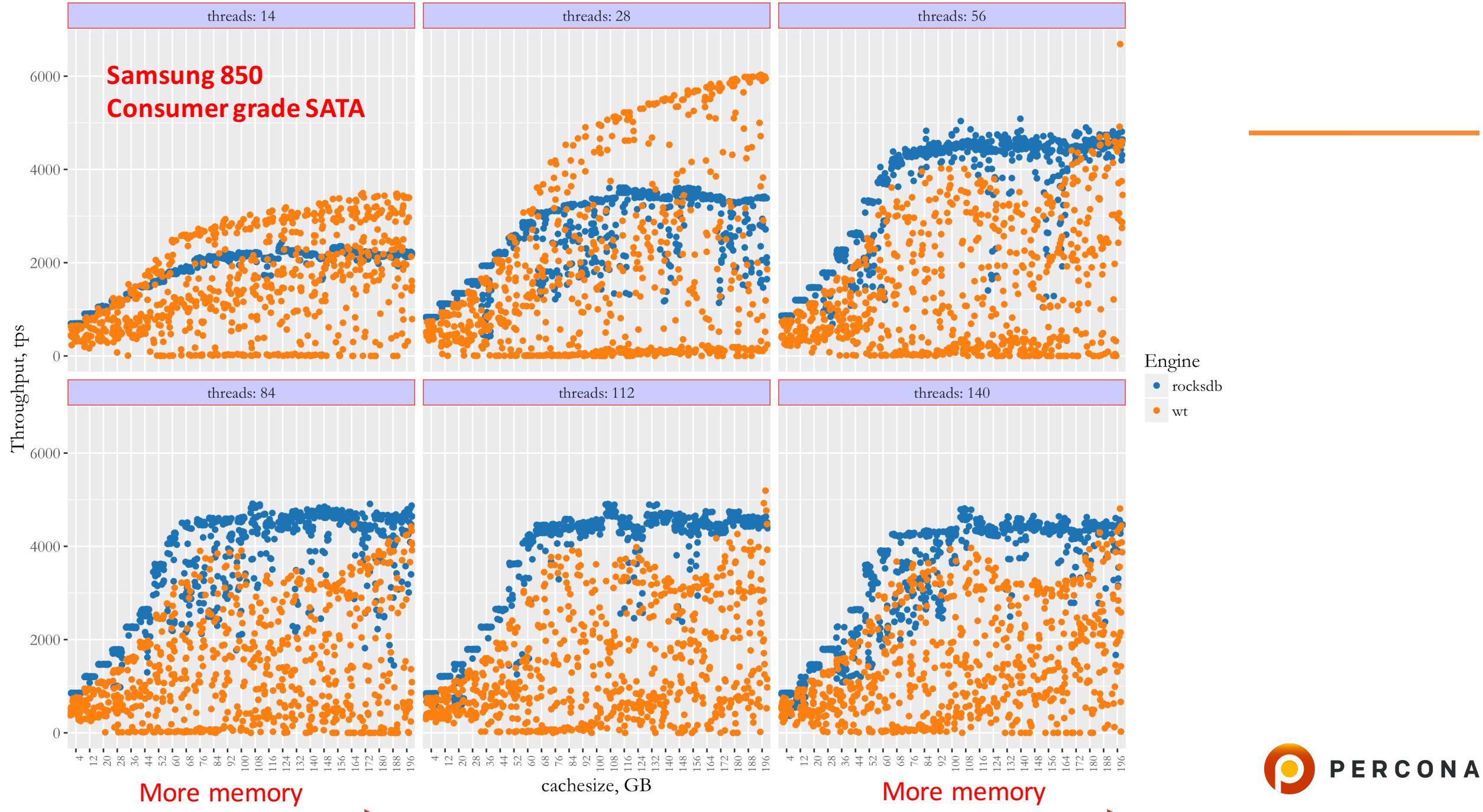
Read-Write results

Memory varies



PERCONA

[network] sysbench OLTP READ-WRITE samsung 850



[network] sysbench OLTP READ-WRITE samsung 850



[network] sysbench OLTP READ-WRITE samsung 863



[network] sysbench OLTP READ-WRITE samsung 863



Read-Write

Threads scalability – in-memory workload

Cachesize 200GB

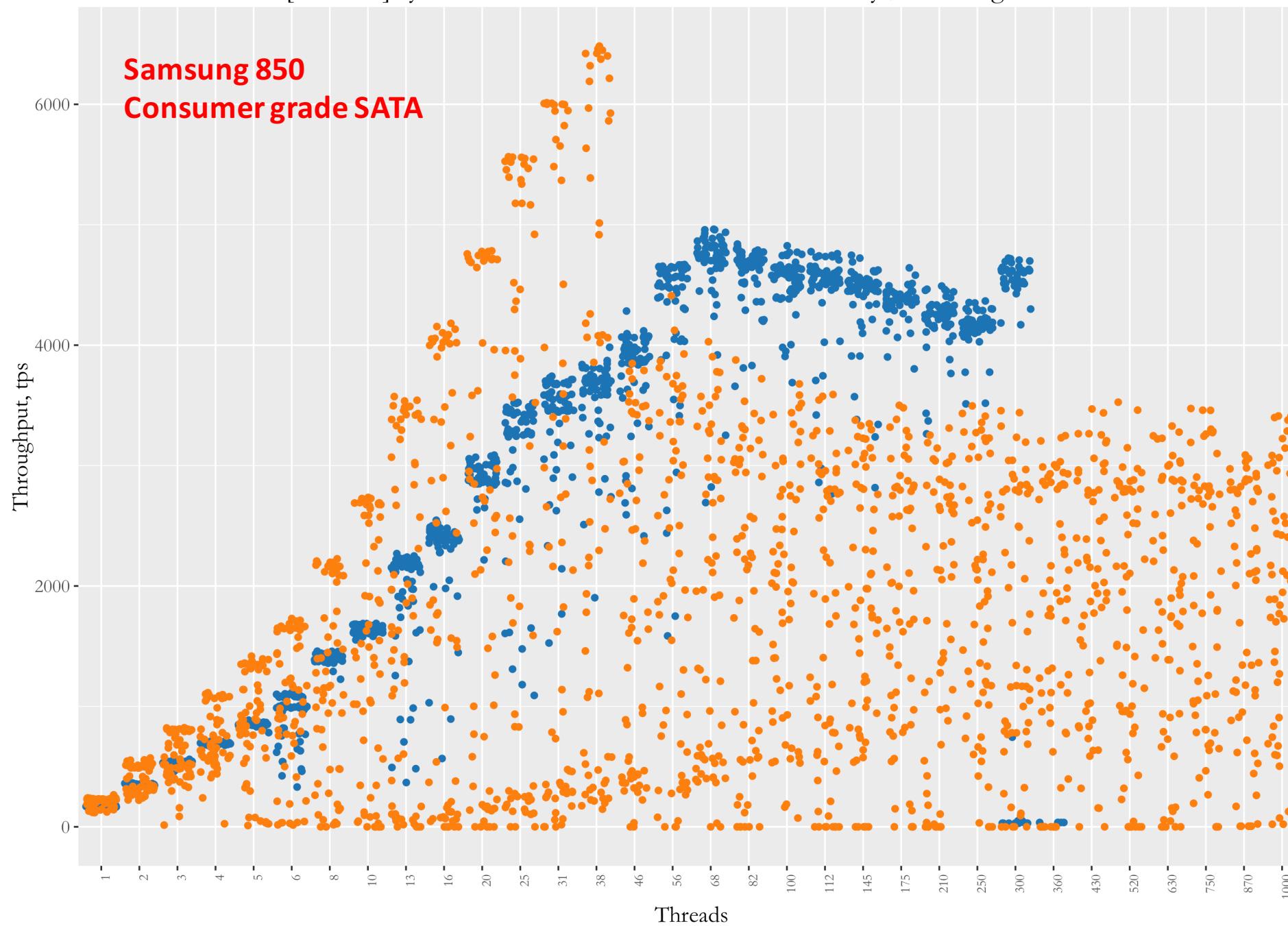


PERCONA

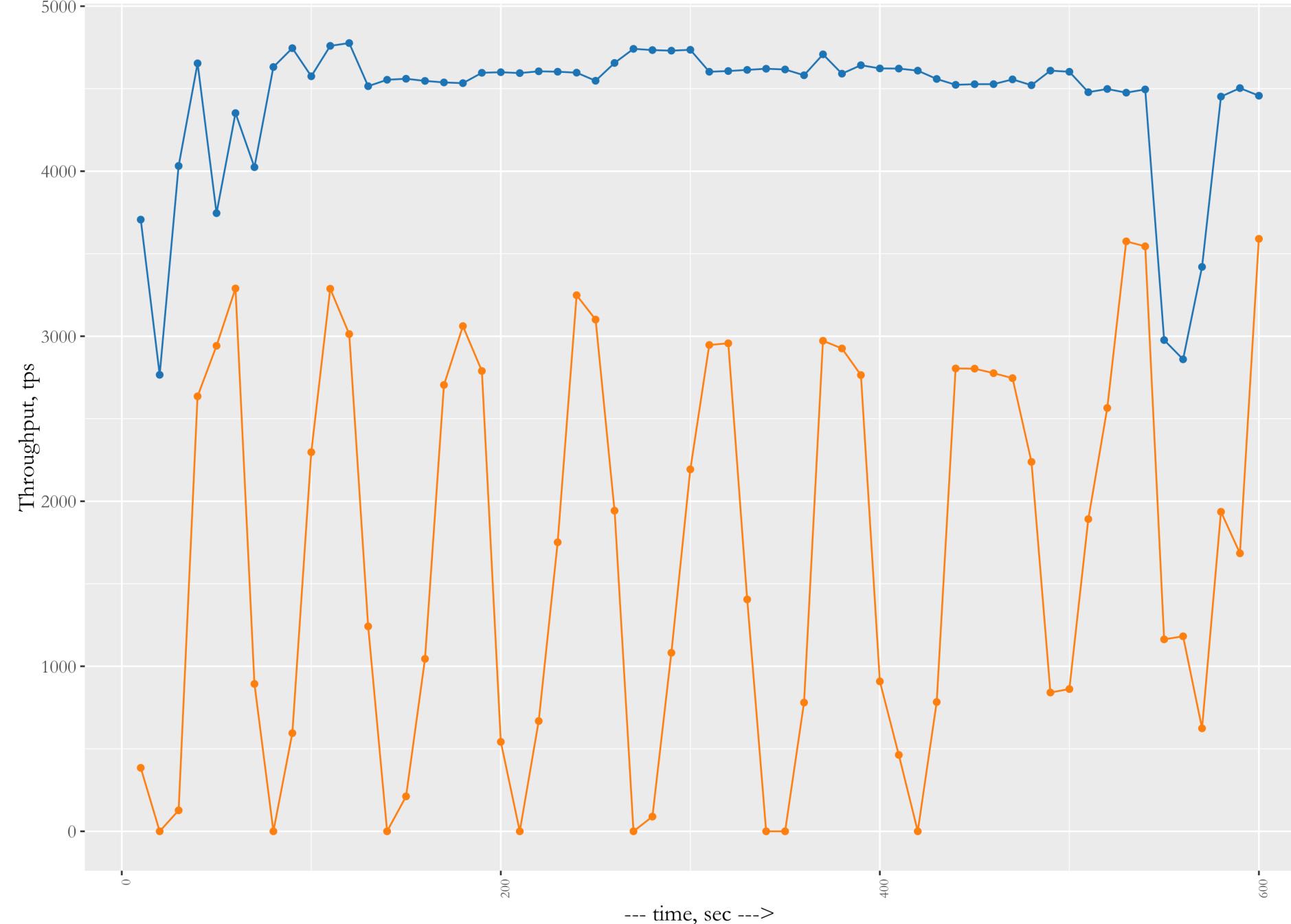
[network] sysbench OLTP READ-WRITE in memory / samsung 850 Pro

Samsung 850

Consumer grade SATA



sysbench OLTP READ-WRITE in memory / 112 threads /samsung 850 Pro

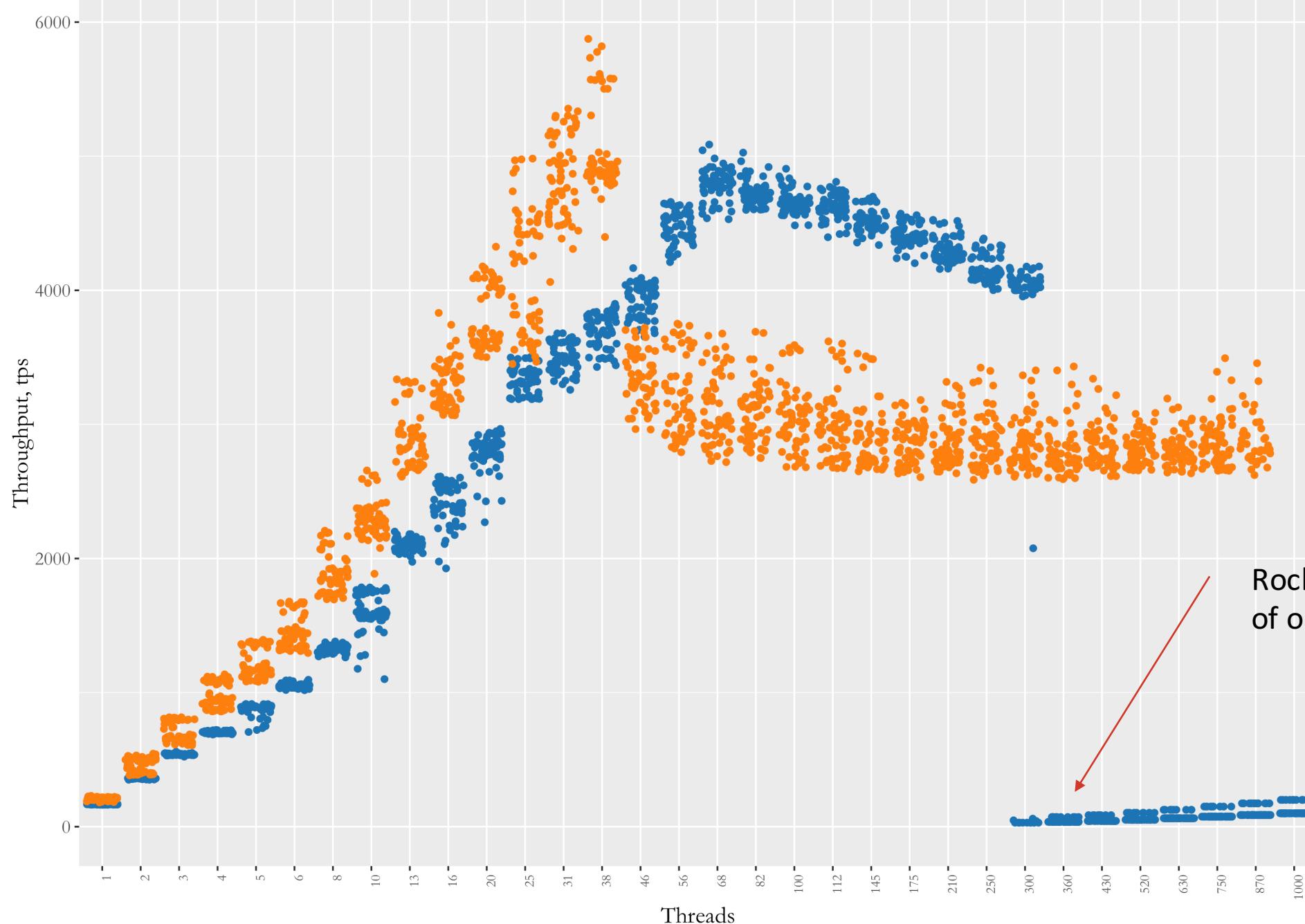


Engine
rocksdb
wt

Let's talk about checkpointing

- Writes buffered in memory (not immediately landed on storage)
- To deal with crashes:
 - Write-Ahead-Log (WAL)
 - Redo-logs
 - Journals
- To limit size of log files in WiredTiger:
 - Time-based checkpointing (60 sec default)
 - syncdelay option
 - Or 2GB of log files (journals), whatever happens sooner
 - I am not sure if 2GB is configurable

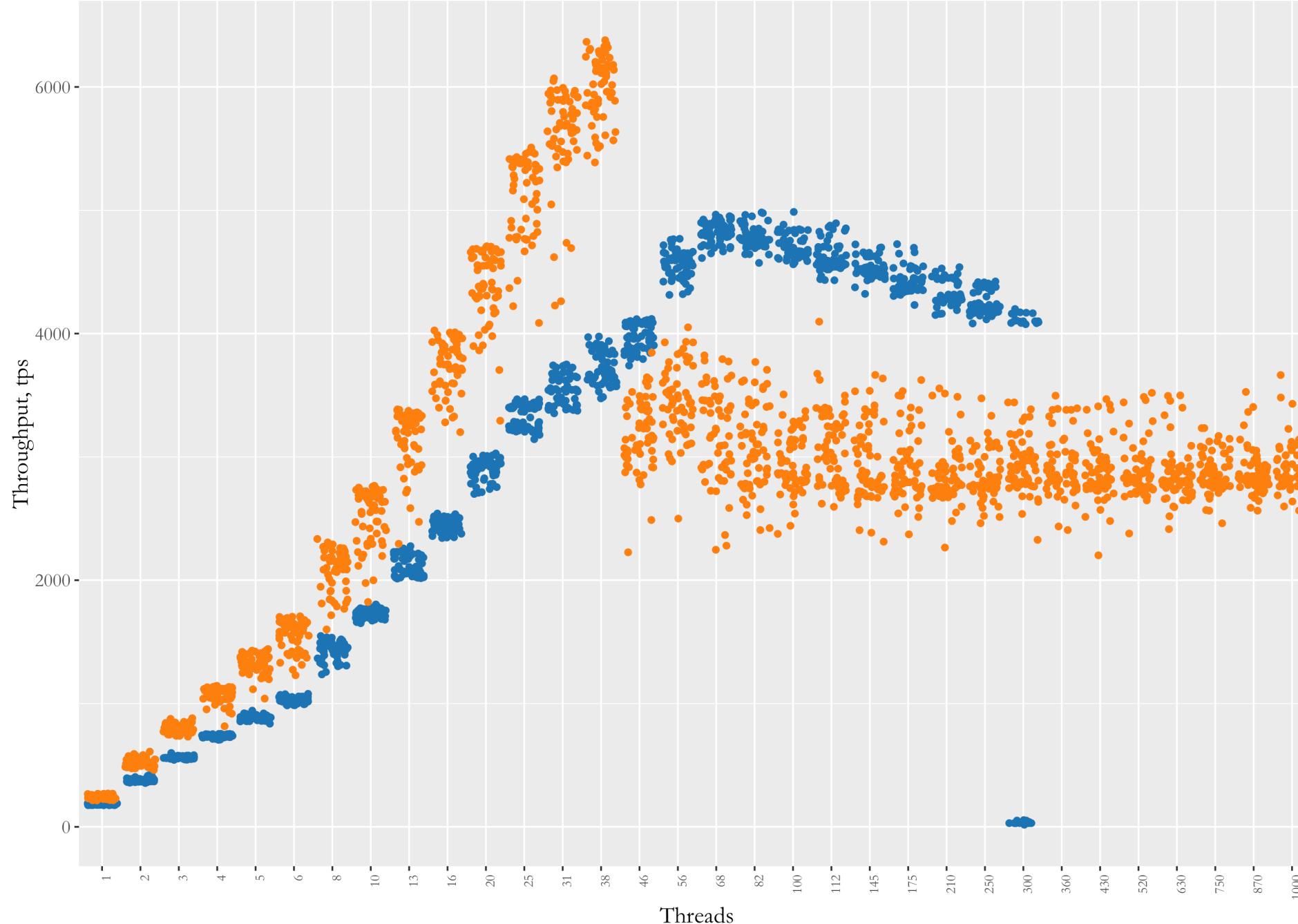
[network] sysbench OLTP READ-WRITE in memory / samsung SM863



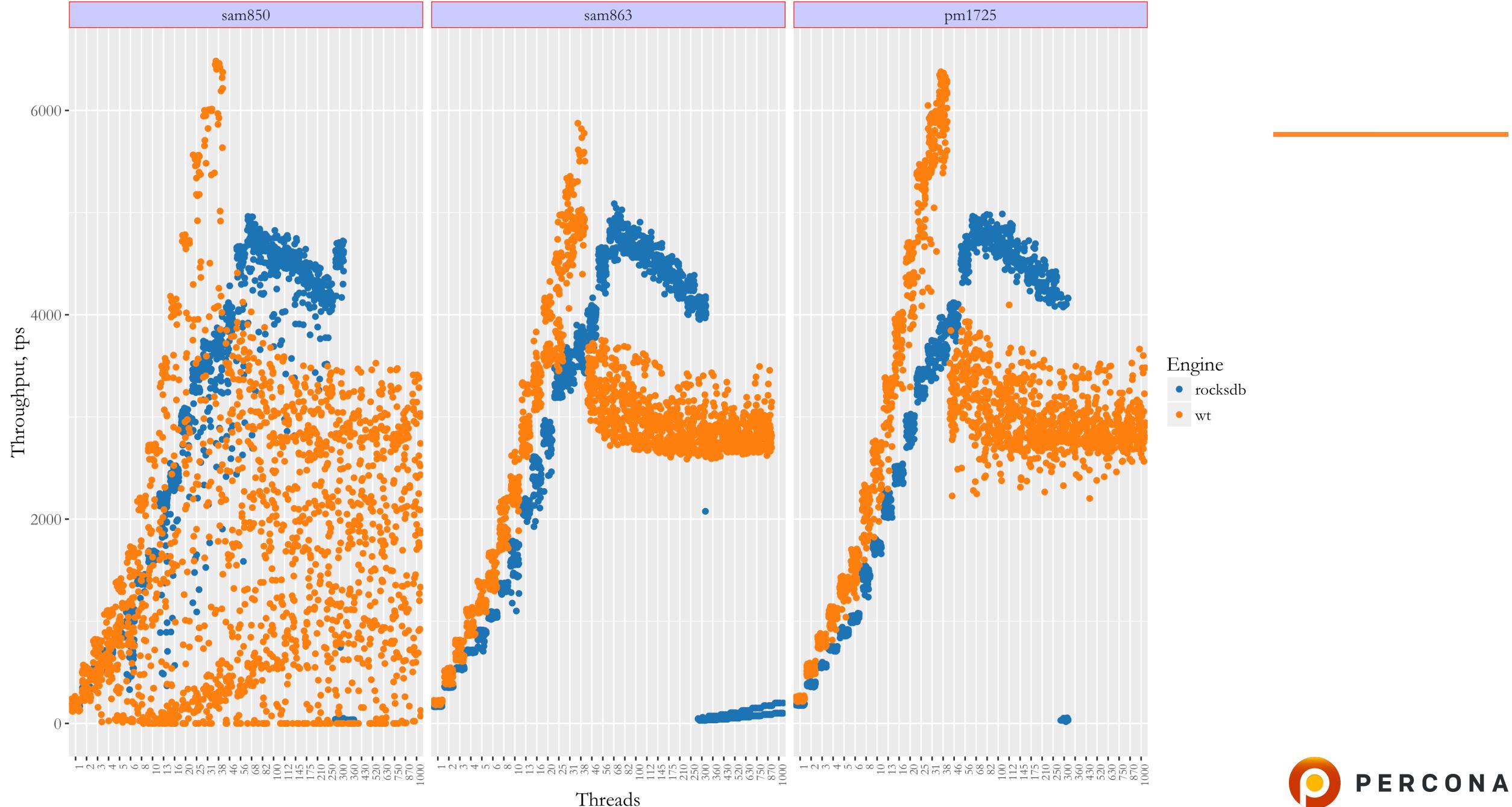
Let's talk about open files in RocksDB

- RocksDB stores data in chunk files, fixed sizes
- It can result in large number of open files at give point of time
 - Be sure to increase the limit of open files ulimit -n
 - Not quite predictable how much will be required

[network] sysbench OLTP READ-WRITE in memory / samsung PM1725



sysbench OLTP READ-WRITE in memory



Conclusions

- When data fits into memory:
 - RocksDB works well even on cheap storage
 - WiredTiger still requires storage that handles big volume writes well

Read-Write

Threads scalability – IO-bound workload

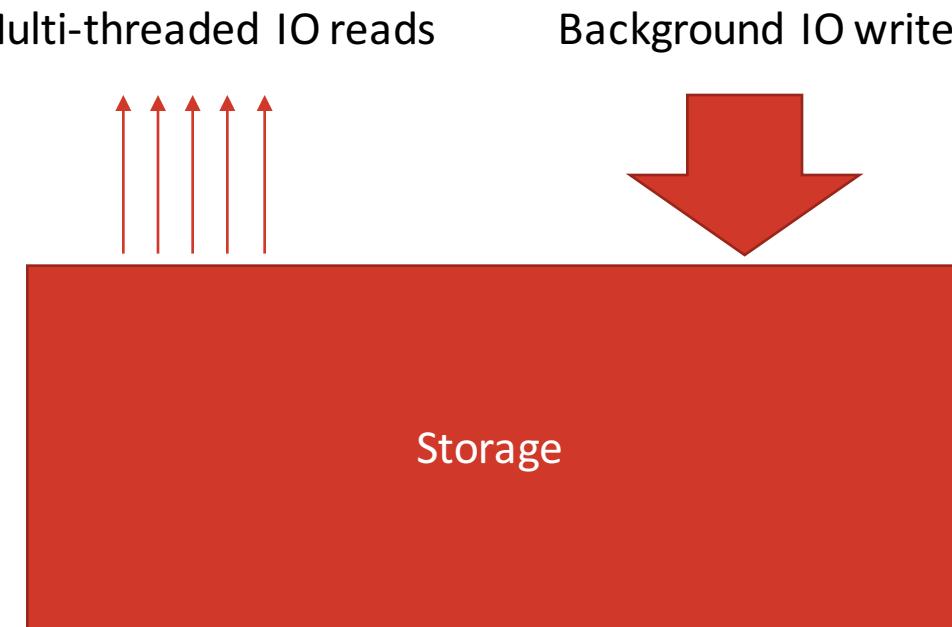
Cachesize **20GB**



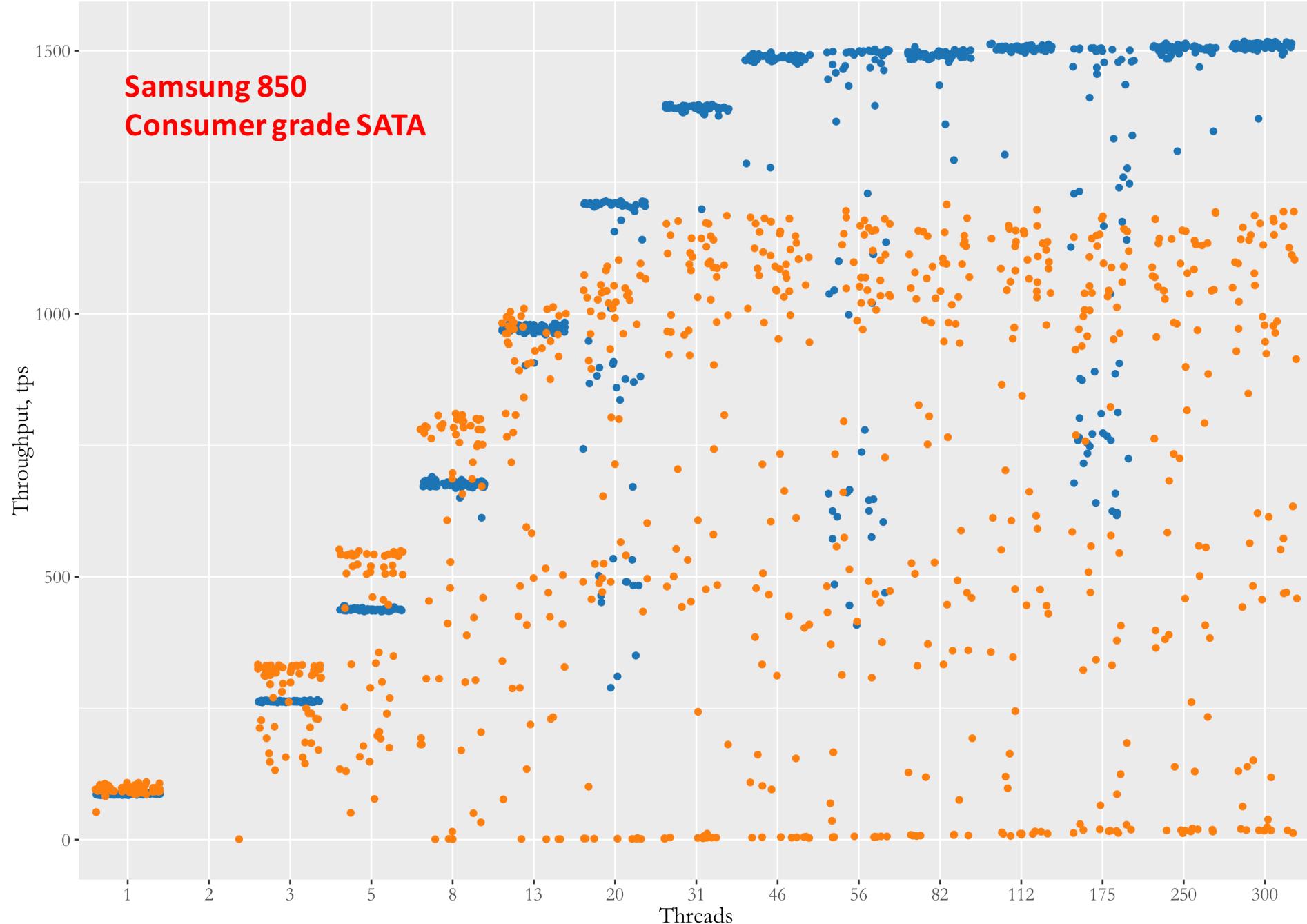
PERCONA

What to expect from this workload

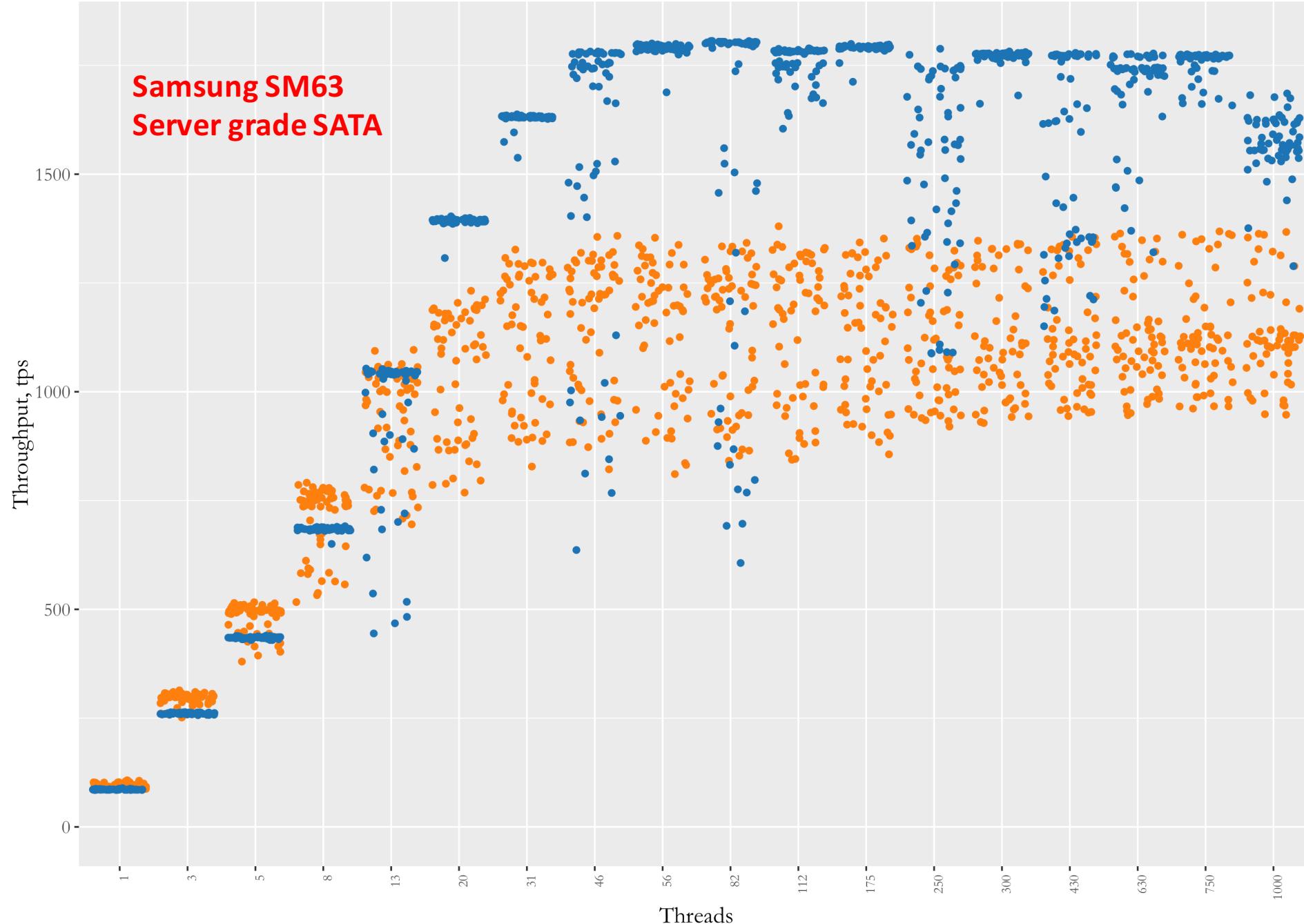
- Intensive background writes
- Response time to user query depends on IO read response time



sysbench OLTP READ-WRITE io-bound / samsung 850 Pro



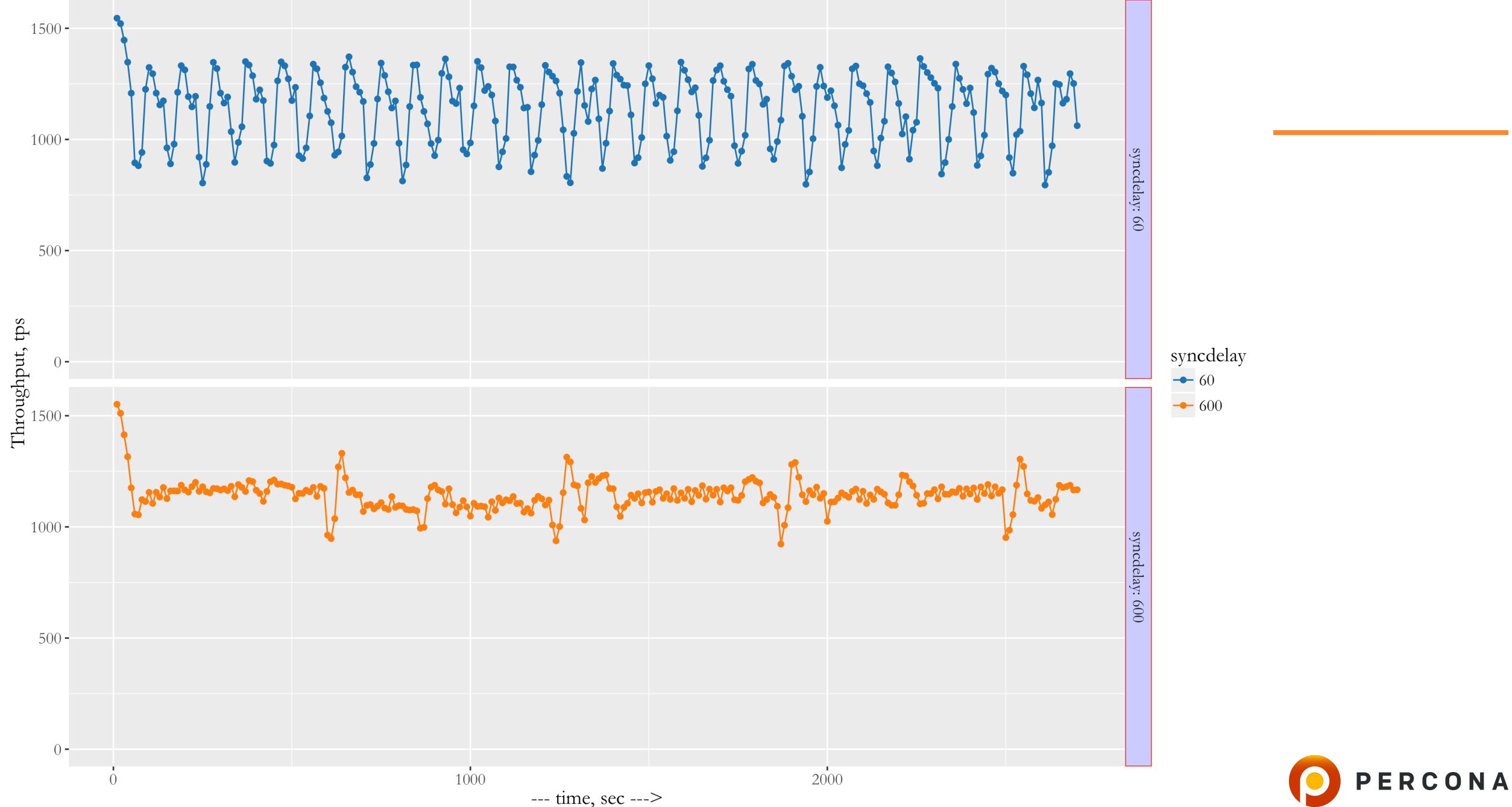
sysbench OLTP READ-WRITE iobound / samsung SM863



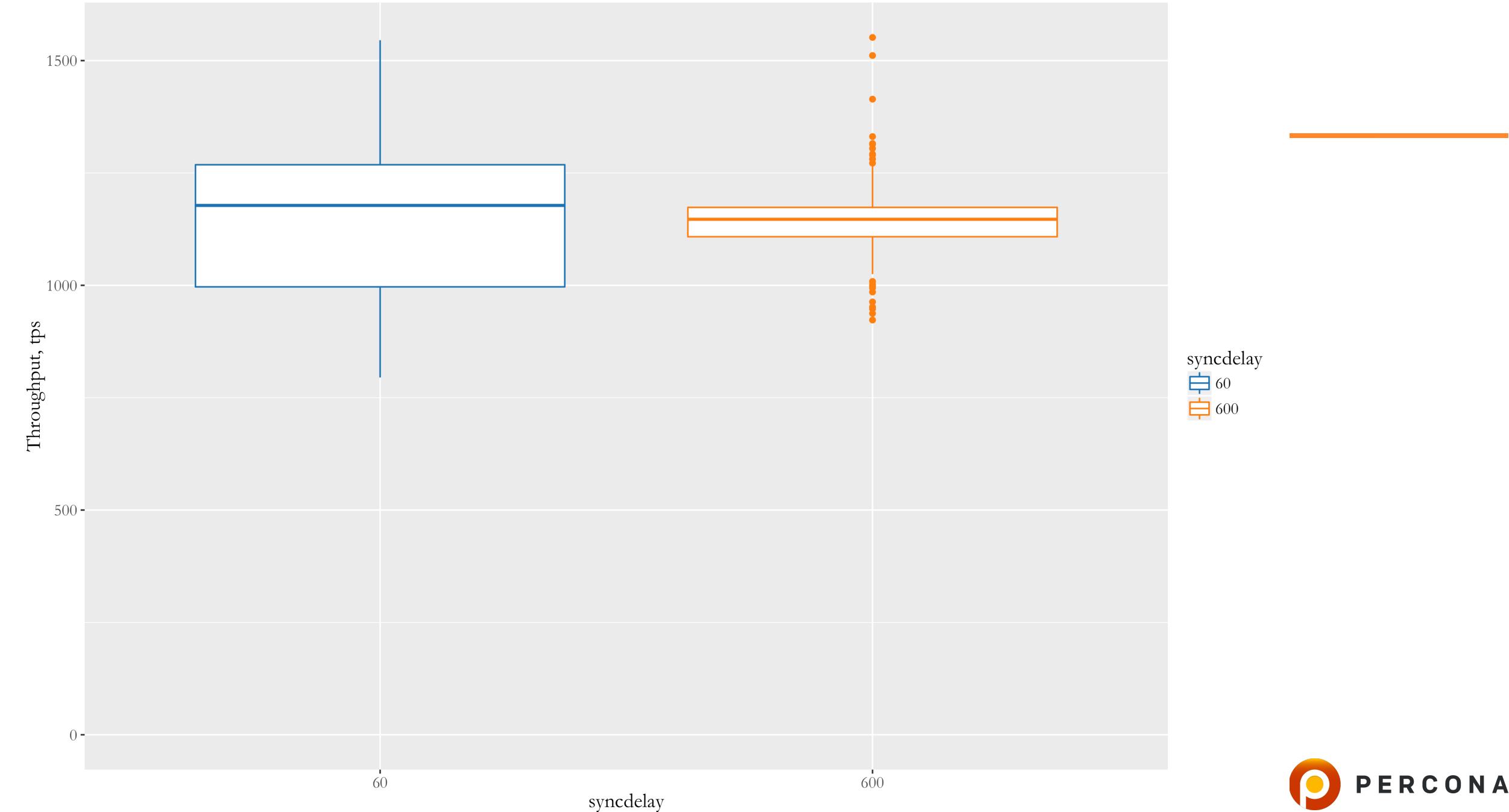
sysbench OLTP READ-WRITE iobound / samsung SM863 / 112 threads



sysbench OLTP READ-WRITE iobound / samsung SM863



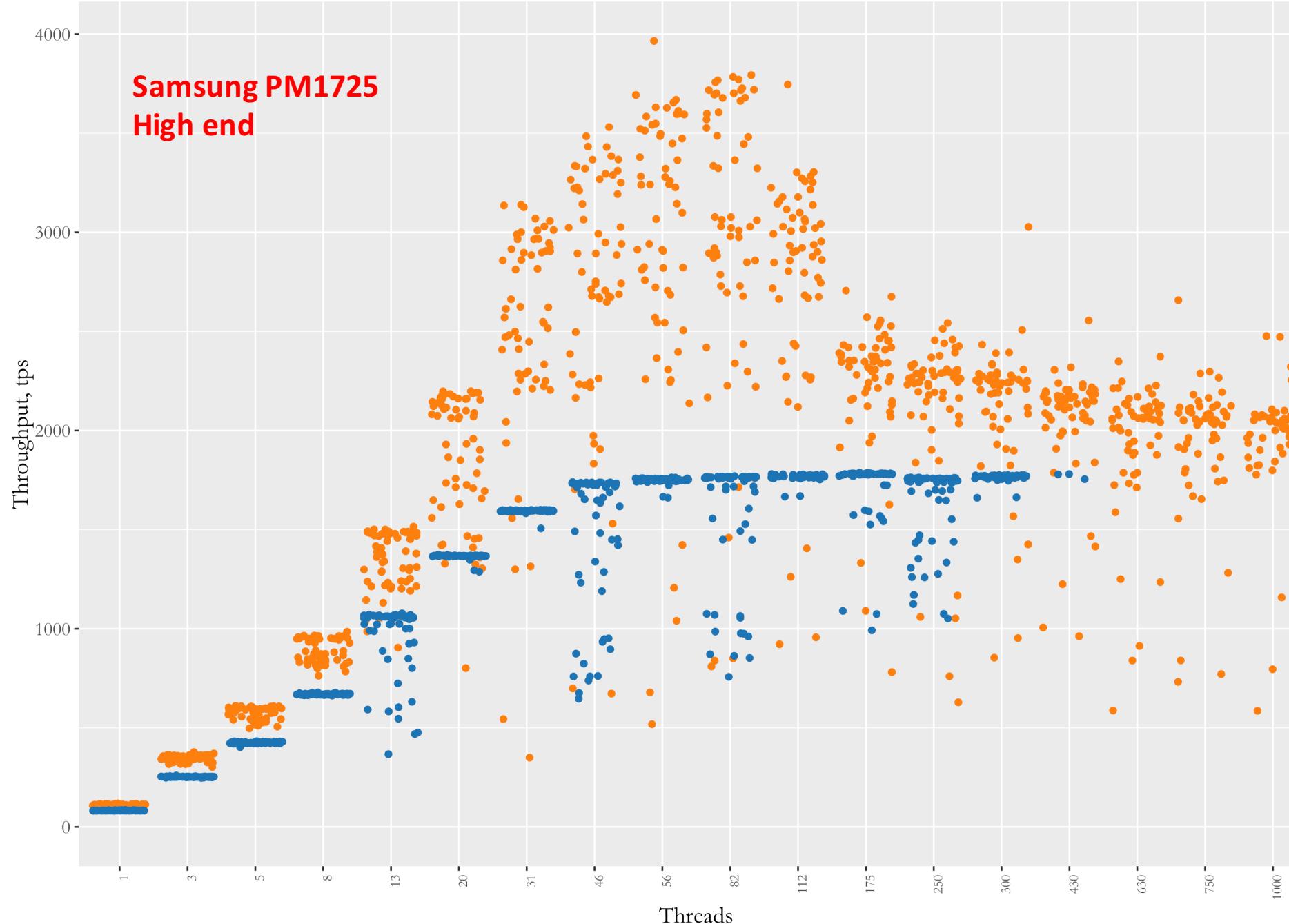
sysbench OLTP READ-WRITE iobound / samsung SM863



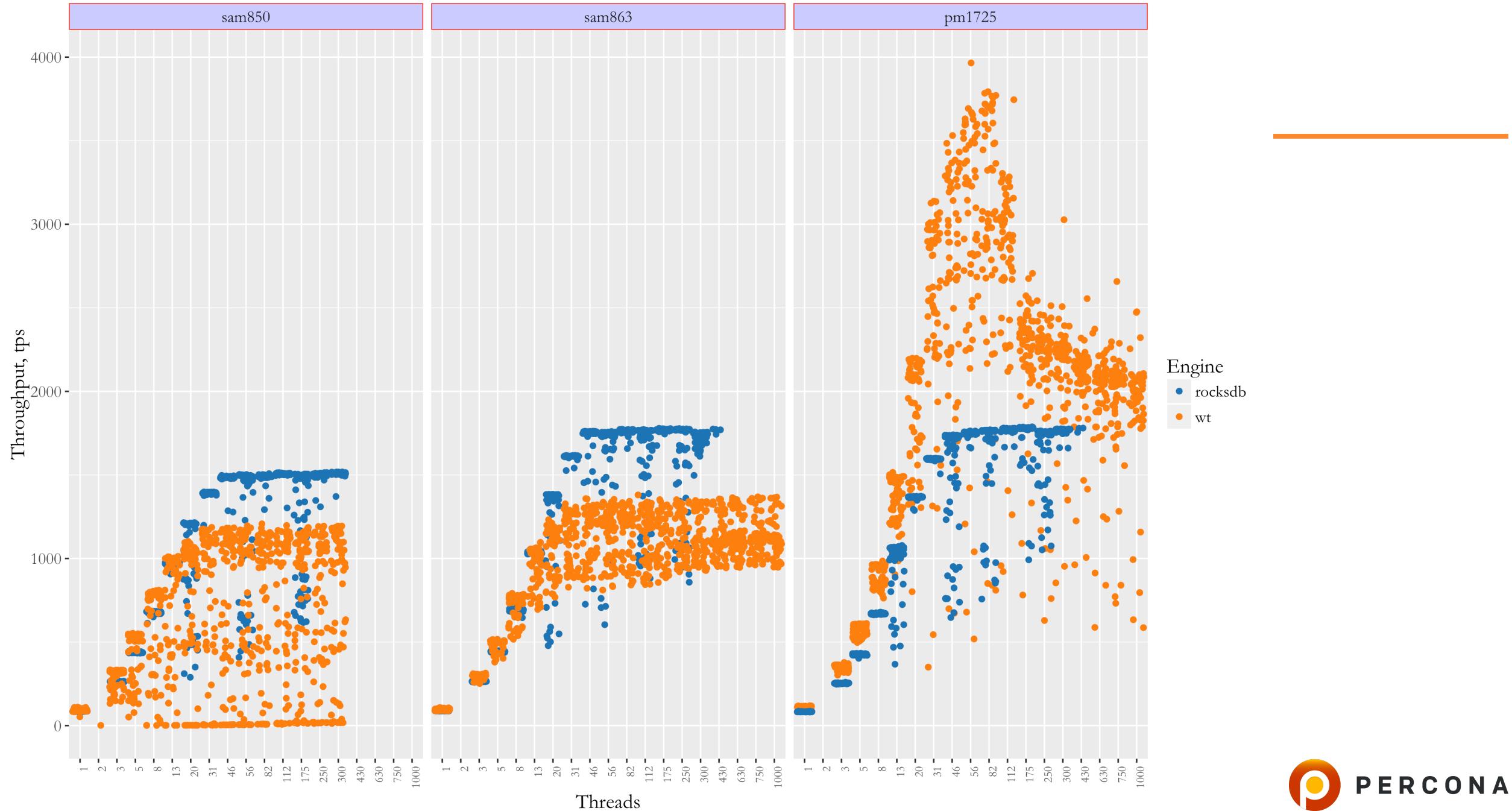
Crash Recovery time

- Syncdelay=60
 - To start mongod after a crash: ~25 mins
- Syncdelay=600
 - ~40 mins

sysbench OLTP READ-WRITE io-bound / samsung PM1725



sysbench OLTP READ-WRITE iobound



Conclusions

- RocksDB works well with cheap storage
- WiredTiger benefits from faster storage
- To improve WiredTiger performance:
 - Typical recommendations for B-Tree based engines
 - More memory if data does not fit into memory
 - Storage with better response times
 - Playing with checkpoint intervals may help with variability, but keep in mind crash recovery times
- Send questions/feedback: @VadimTk, vadim@percona.com