

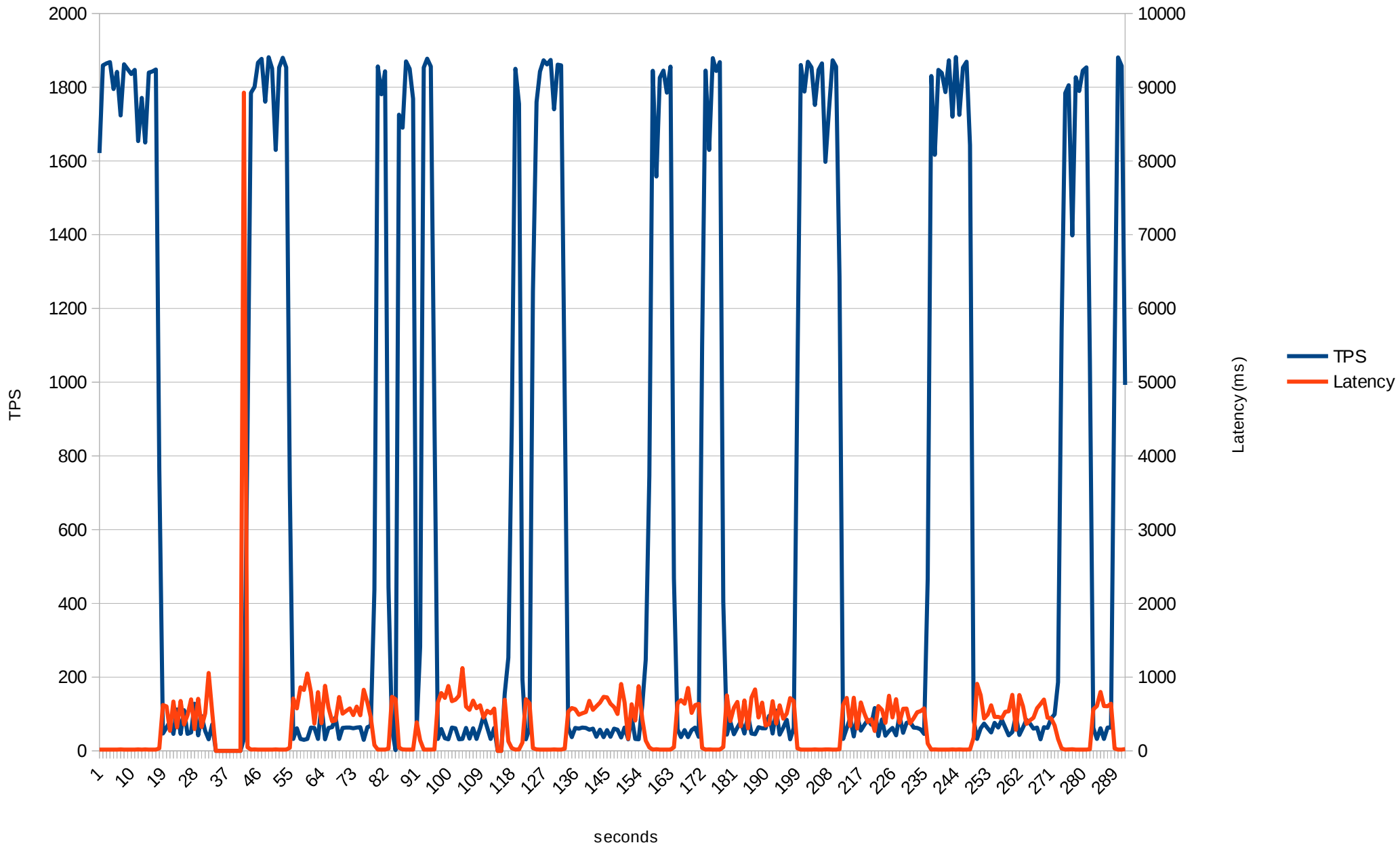
Postgres' IO Architecture, Tuning, Problems

Andres Freund
PostgreSQL Developer & Committer
Citus Data – citusdata.com - @citusdata

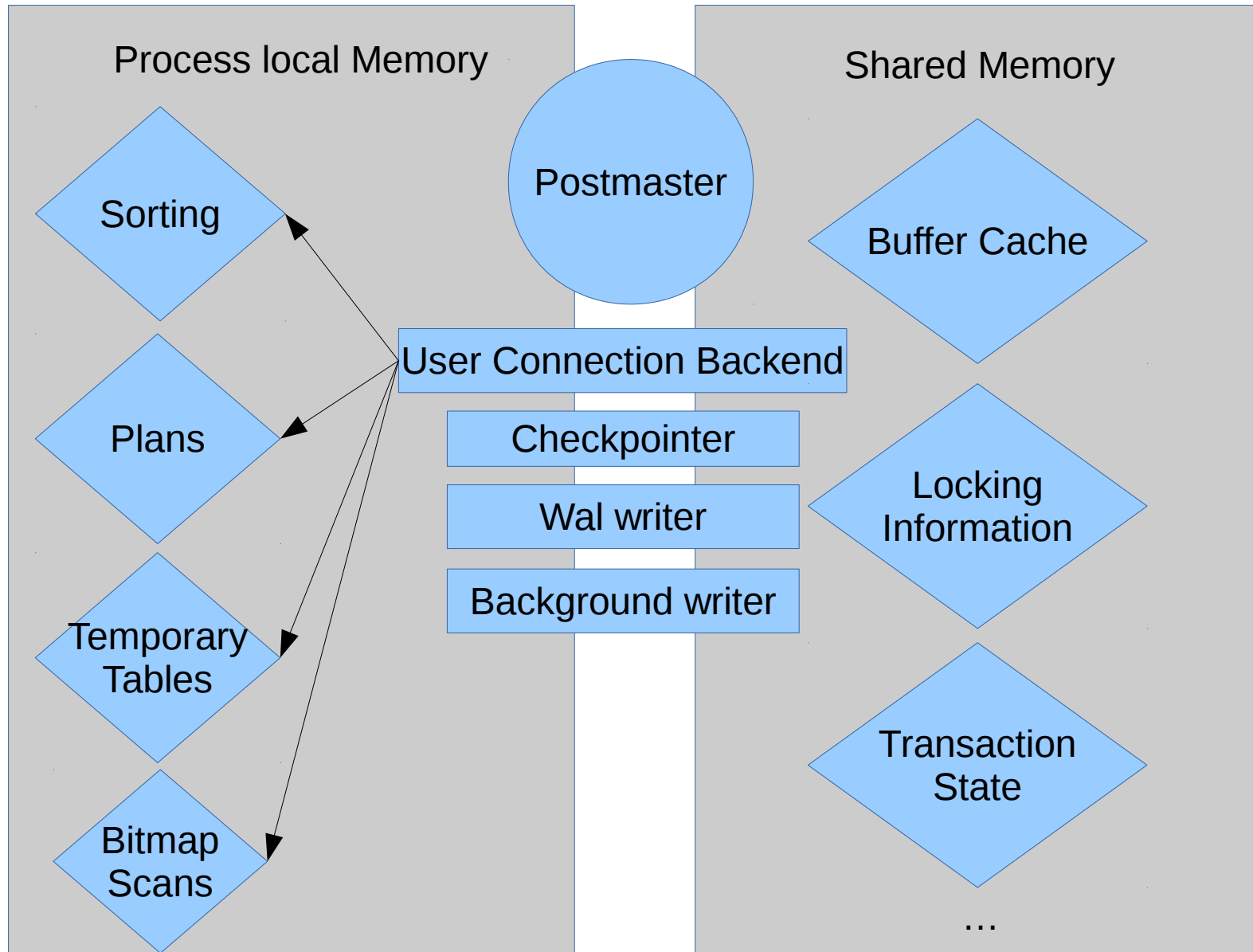
<http://anarazel.de/talks/pgdevday-prague-2016-02-18/io.pdf>

pgbench -M prepared -c 32 -j 32

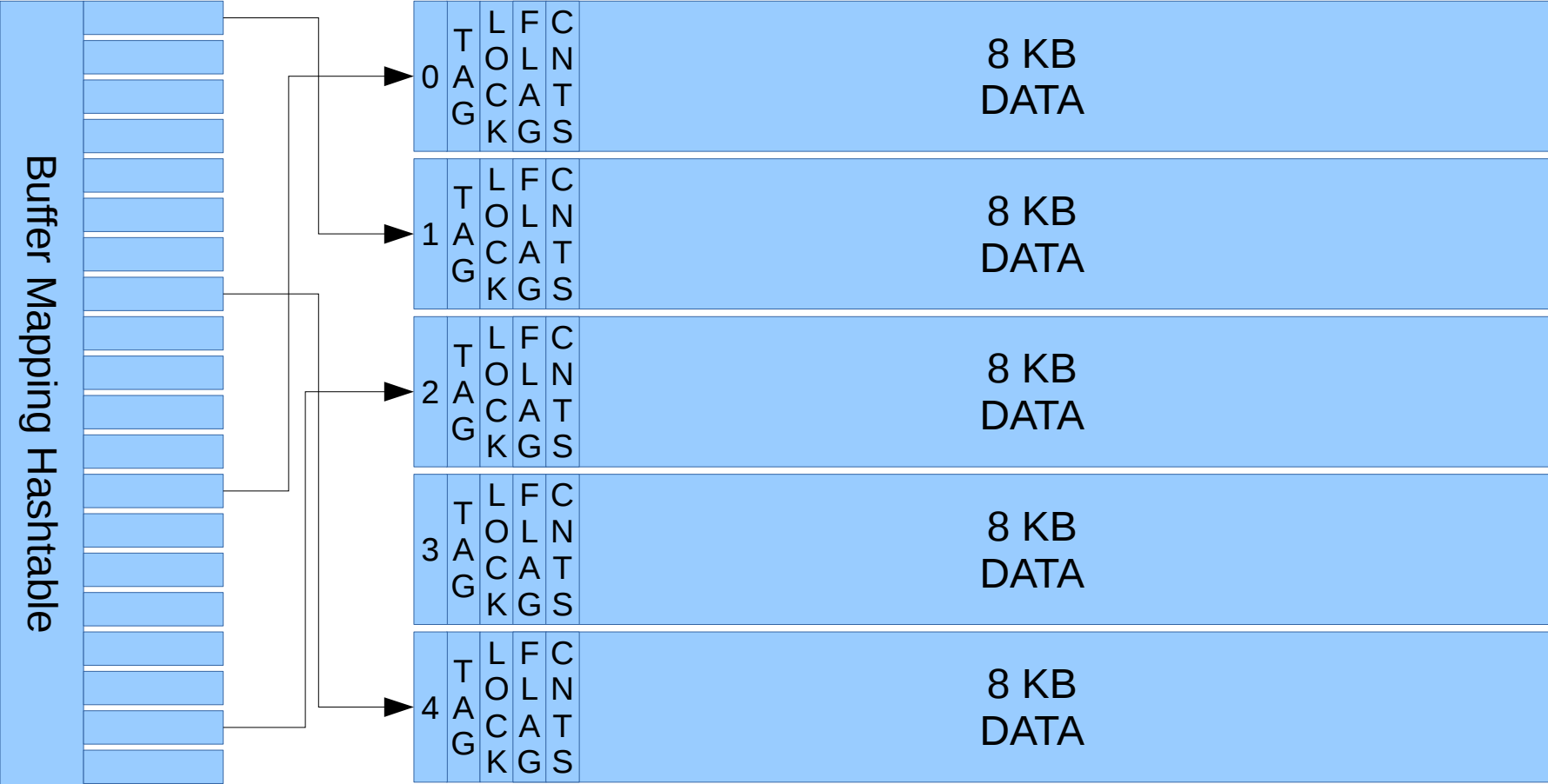
standard settings



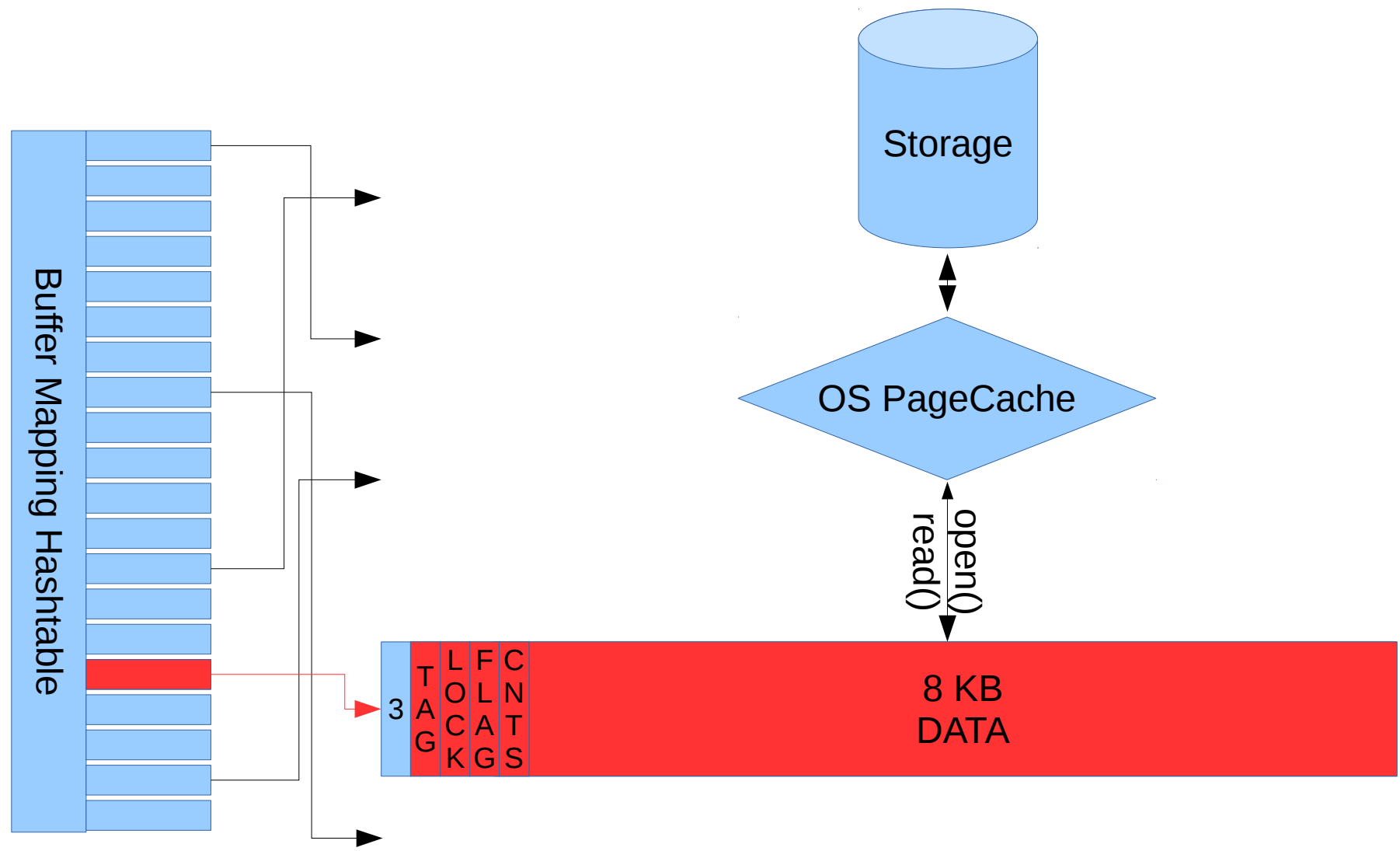
Memory Architecture



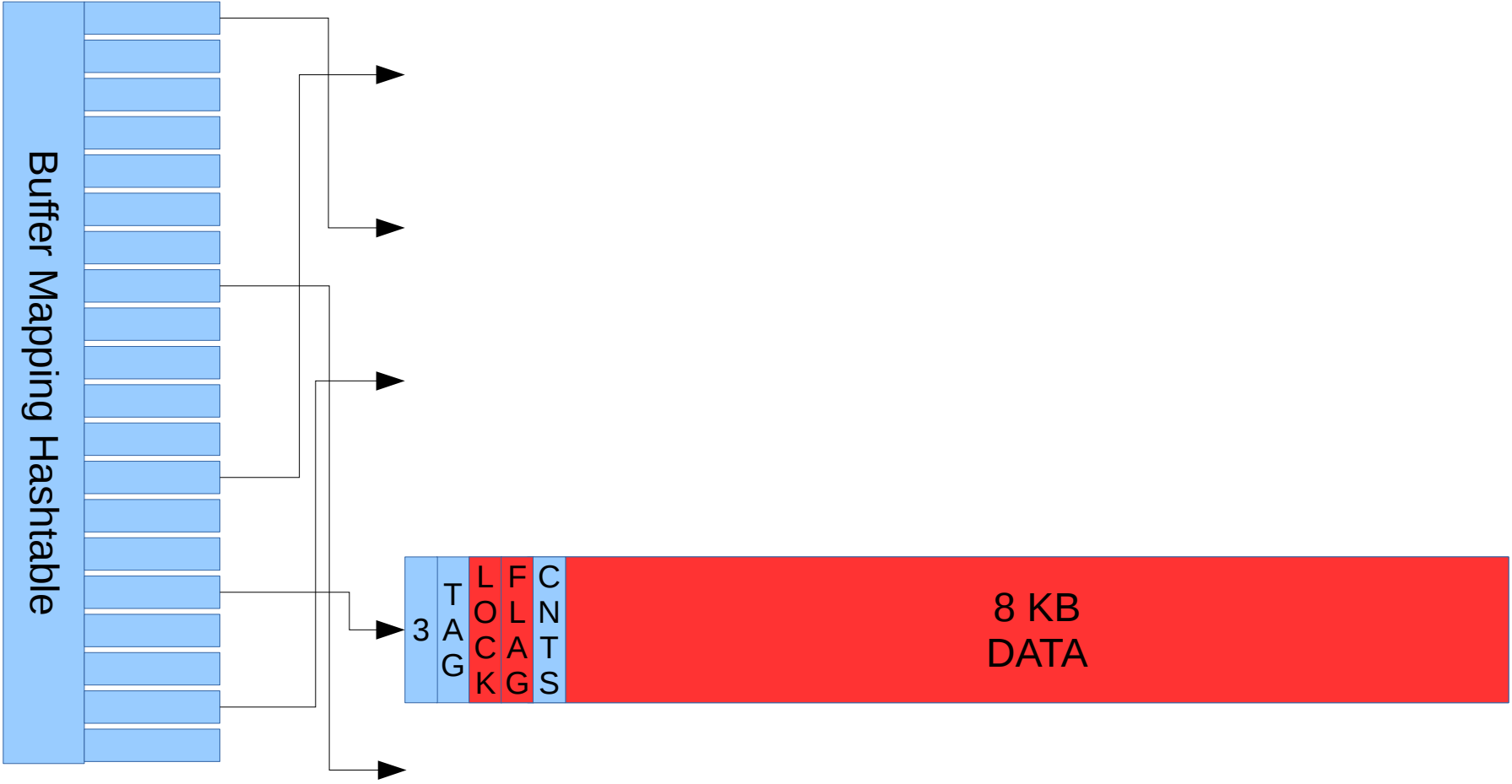
Shared Buffers



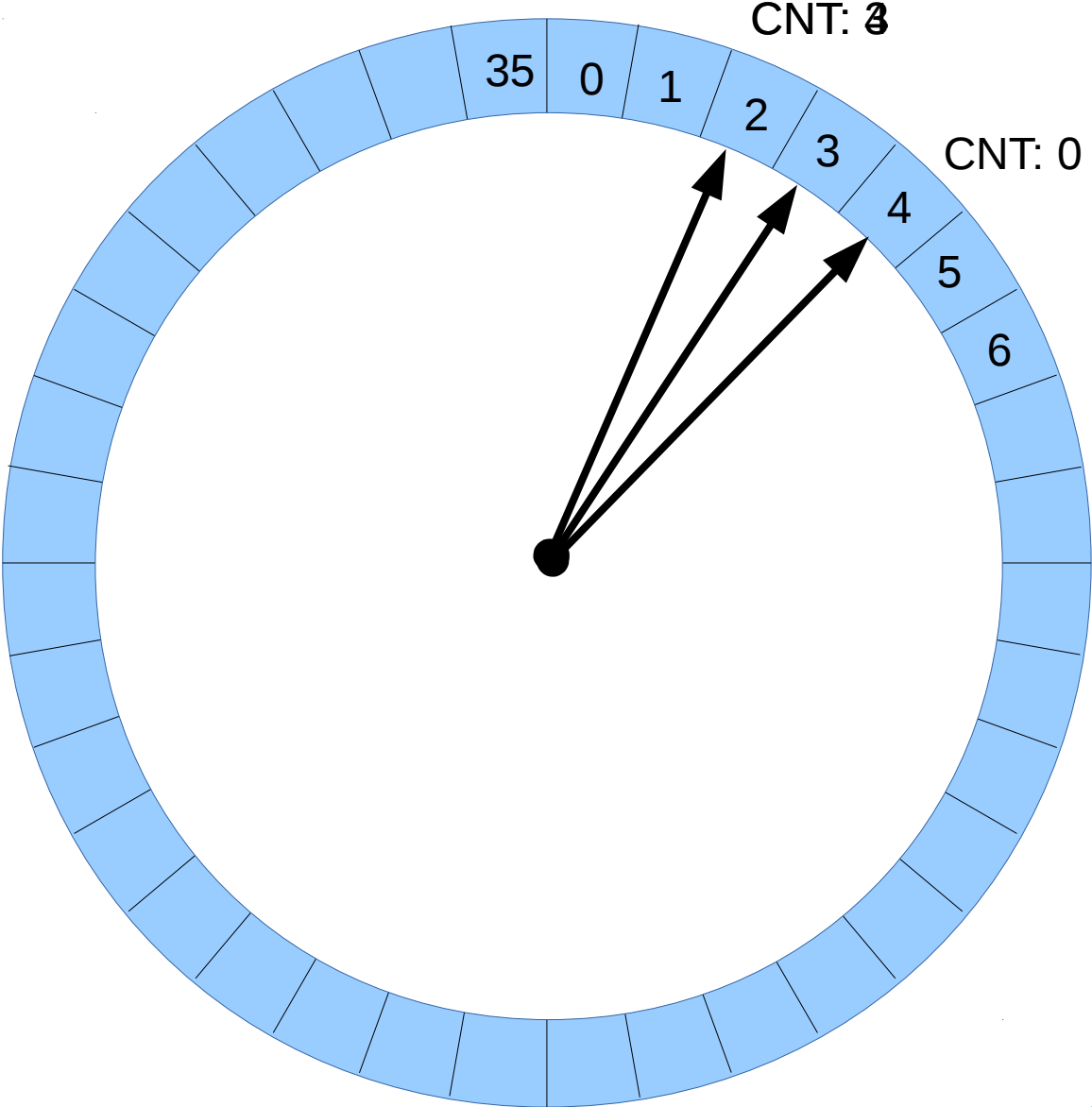
Reading Data



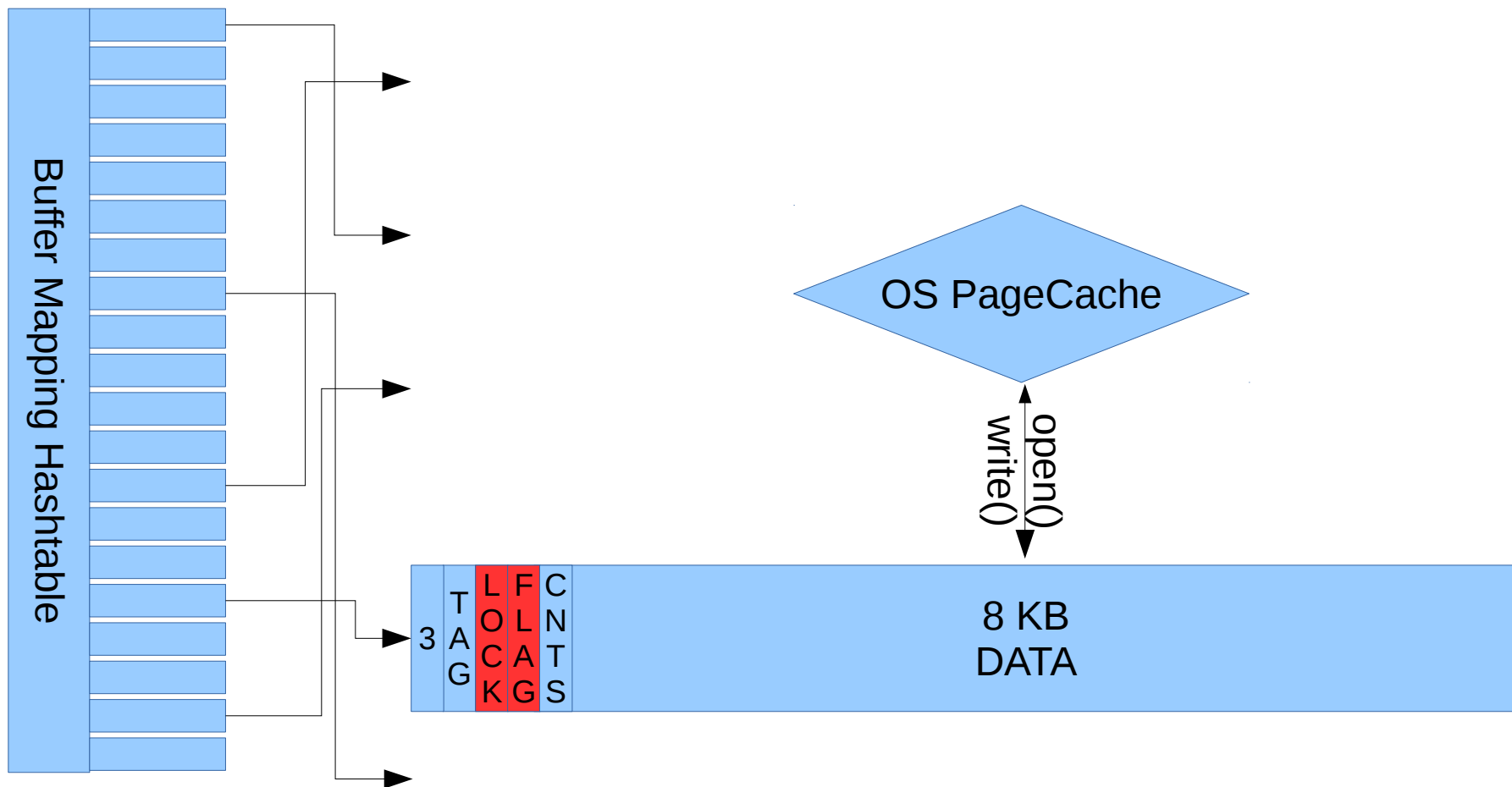
Writing Data



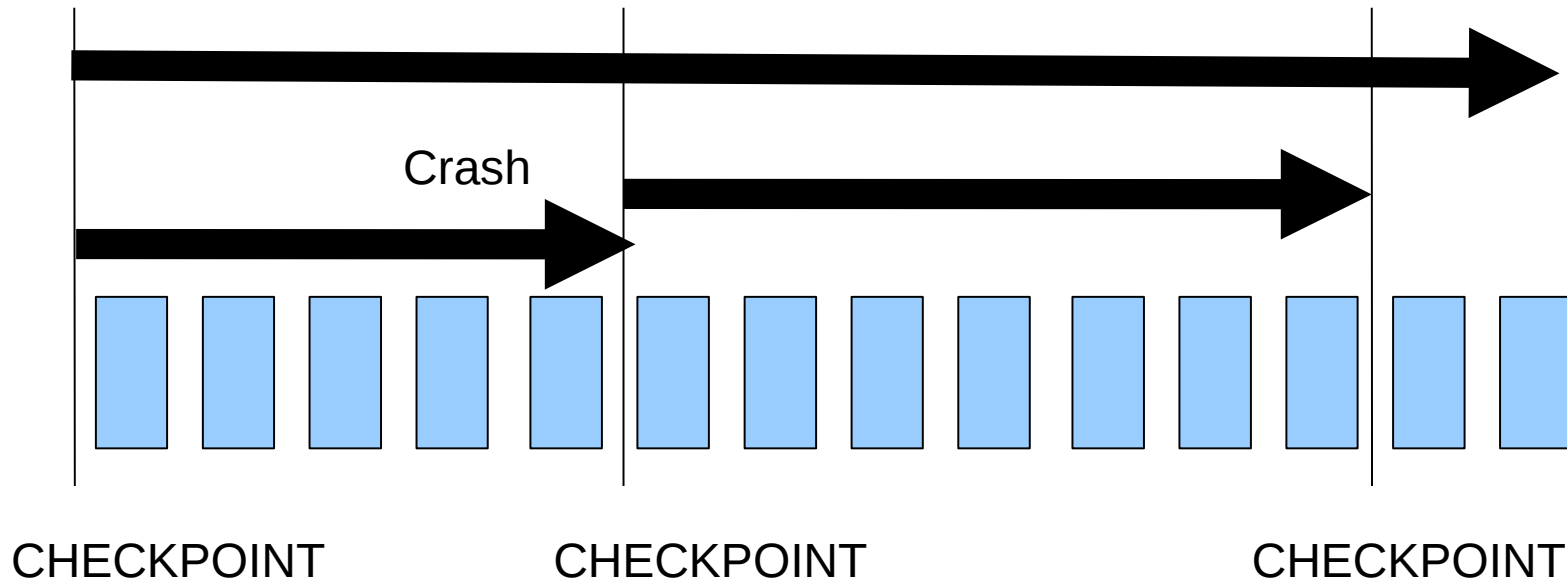
Clock-Sweep



Writing Data Out



Recovery & Checkpoints



Checkpoints

- 1) Remember current position in WAL
- 2) Do some boring things
- 3) Write out all dirty buffers
- 4) Fsync all files modified since last checkpoint
- 5) Write checkpoint WAL record, pg_control etc.
- 6) Remove old WAL

Triggering Checkpoints

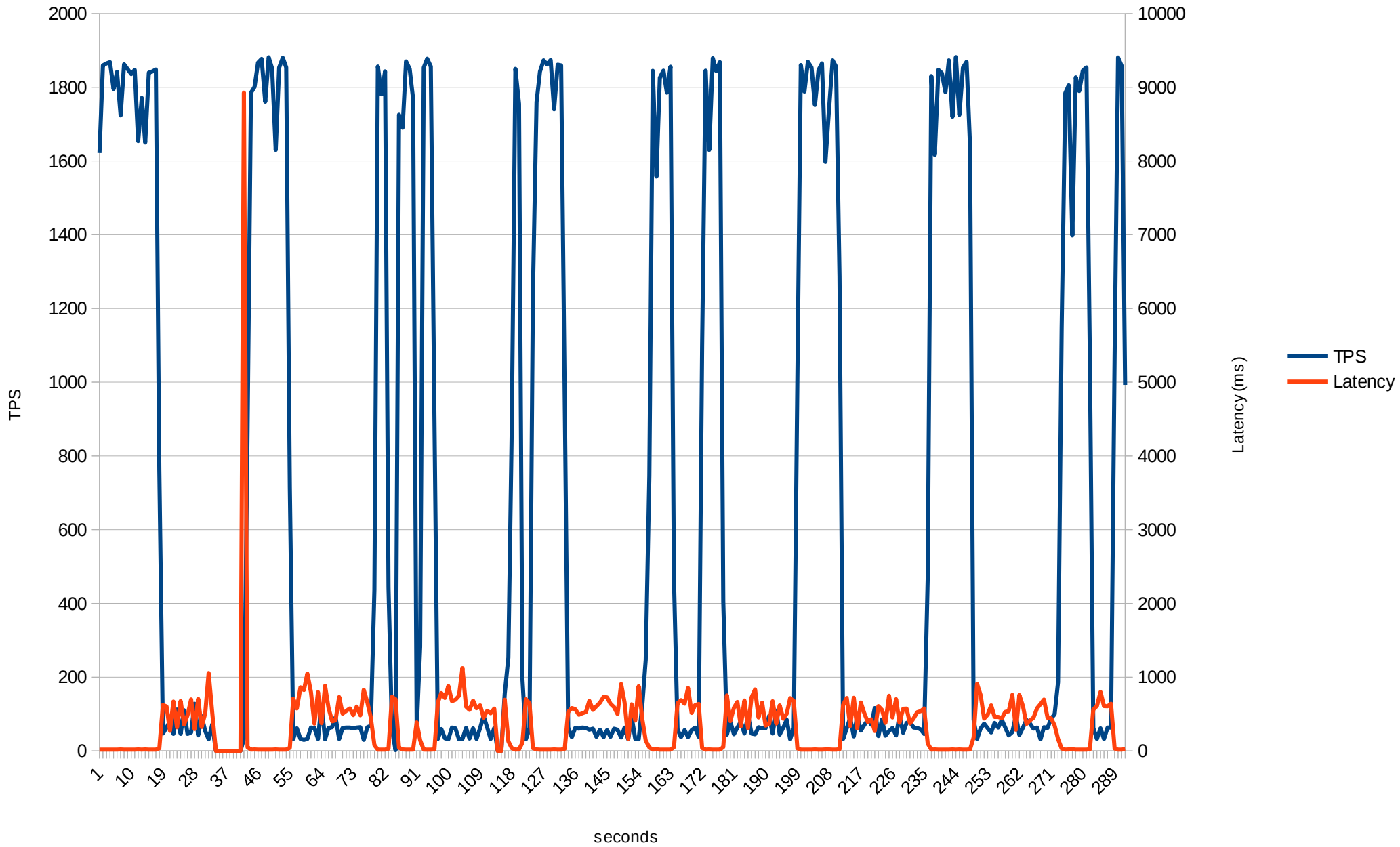
- `checkpoint_timeout = 5min`
 - LOG: checkpoint starting: time
- `checkpoint_segments = 3 / max_wal_size = 1GB`
 - LOG: checkpoint starting: xlog
 - LOG: checkpoints are occurring too frequently (2 seconds apart)
- `shutdown`
 - LOG: checkpoint starting: shutdown immediate
- manually (`CHECKPOINT;`)

Spreading Checkpoints

- `checkpoint_completion_target = 0.5`
- estimation based on
 - `checkpoint_timeout`
 - `checkpoint_segments/max_wal_size`
- Spread checkpoints over `completion_target * timeout/segments` till next checkpoint
- Try to keep pace

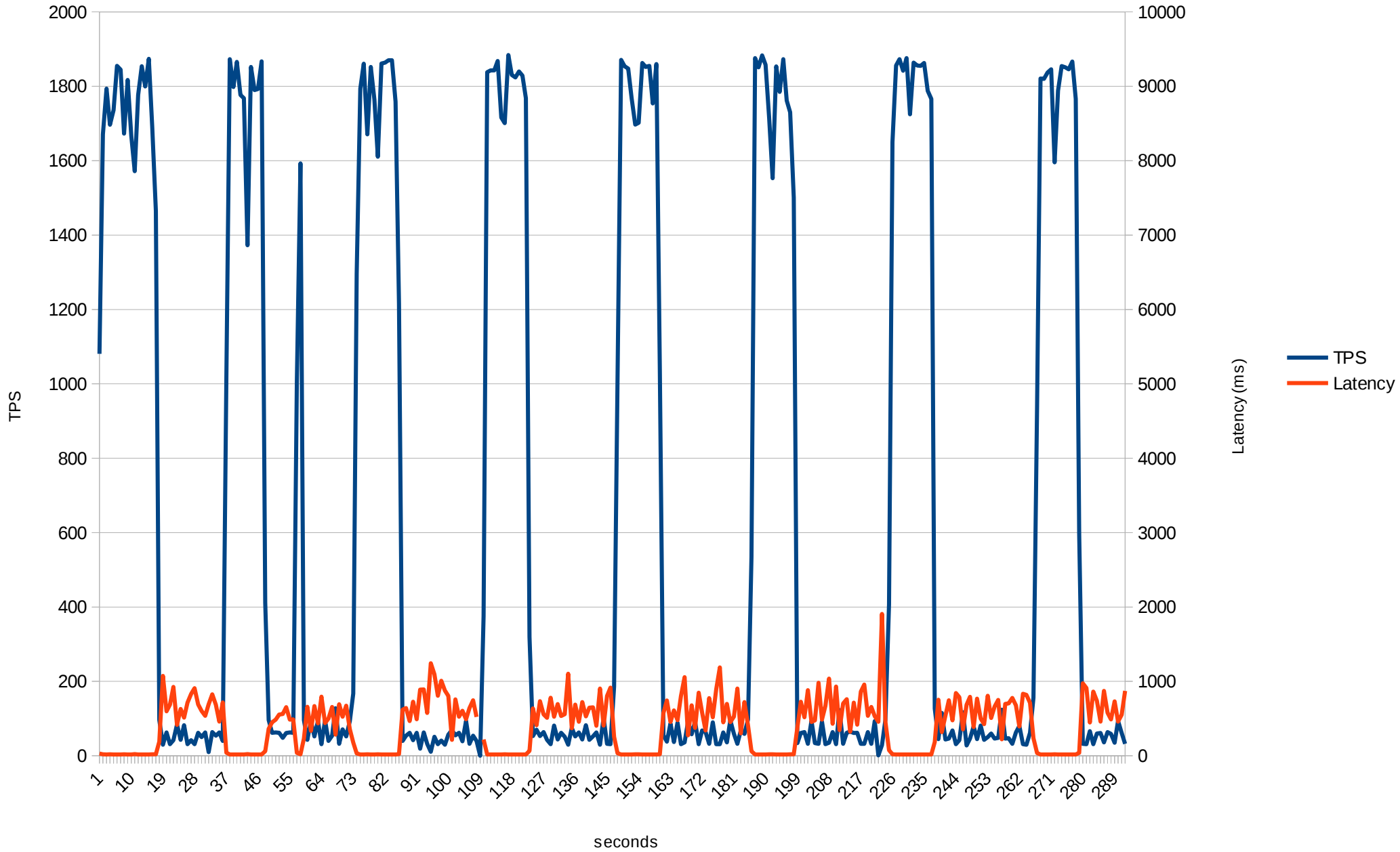
pgbench -M prepared -c 32 -j 32

standard settings



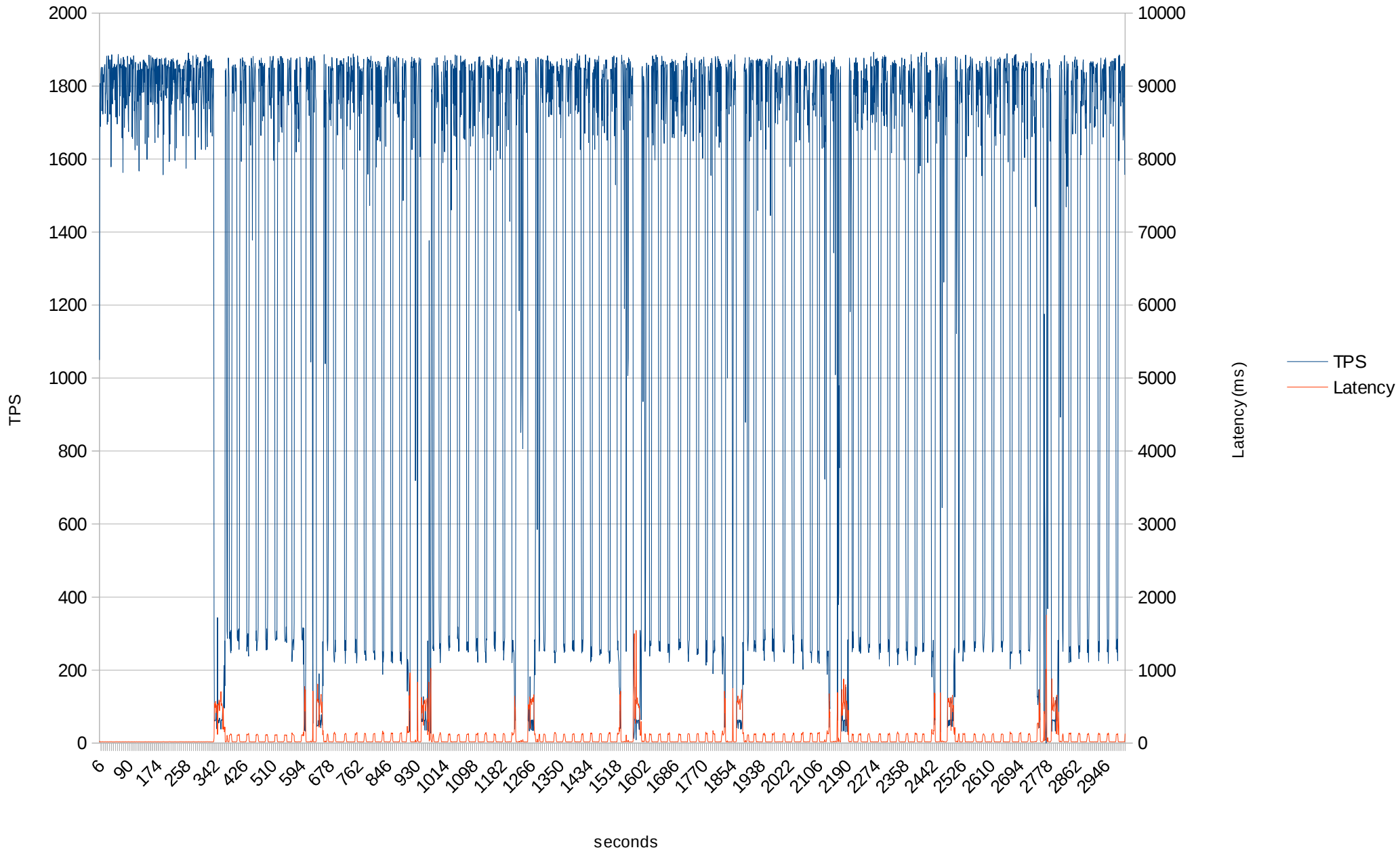
pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB

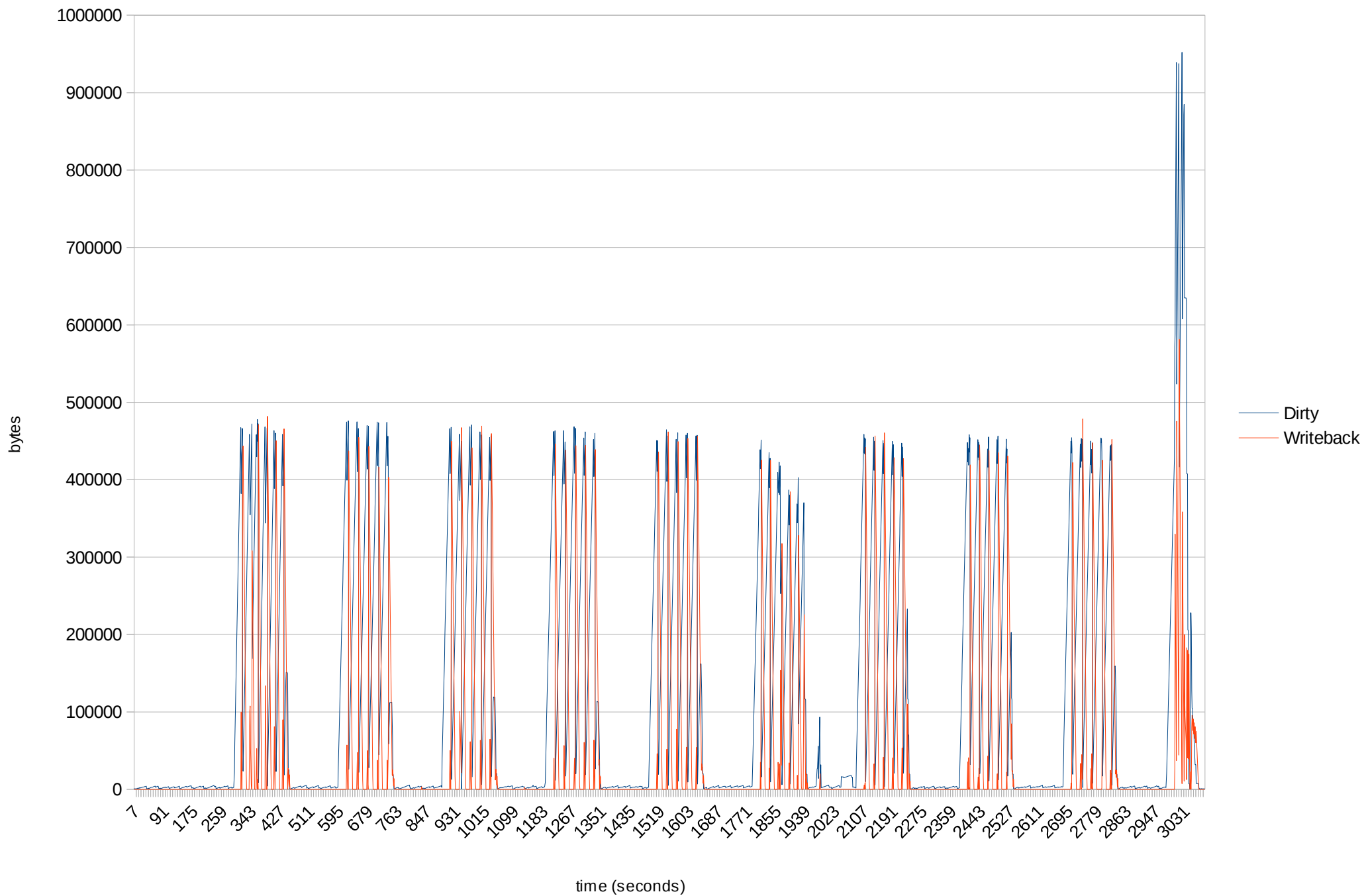


pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB, max_wal_size = 100GB

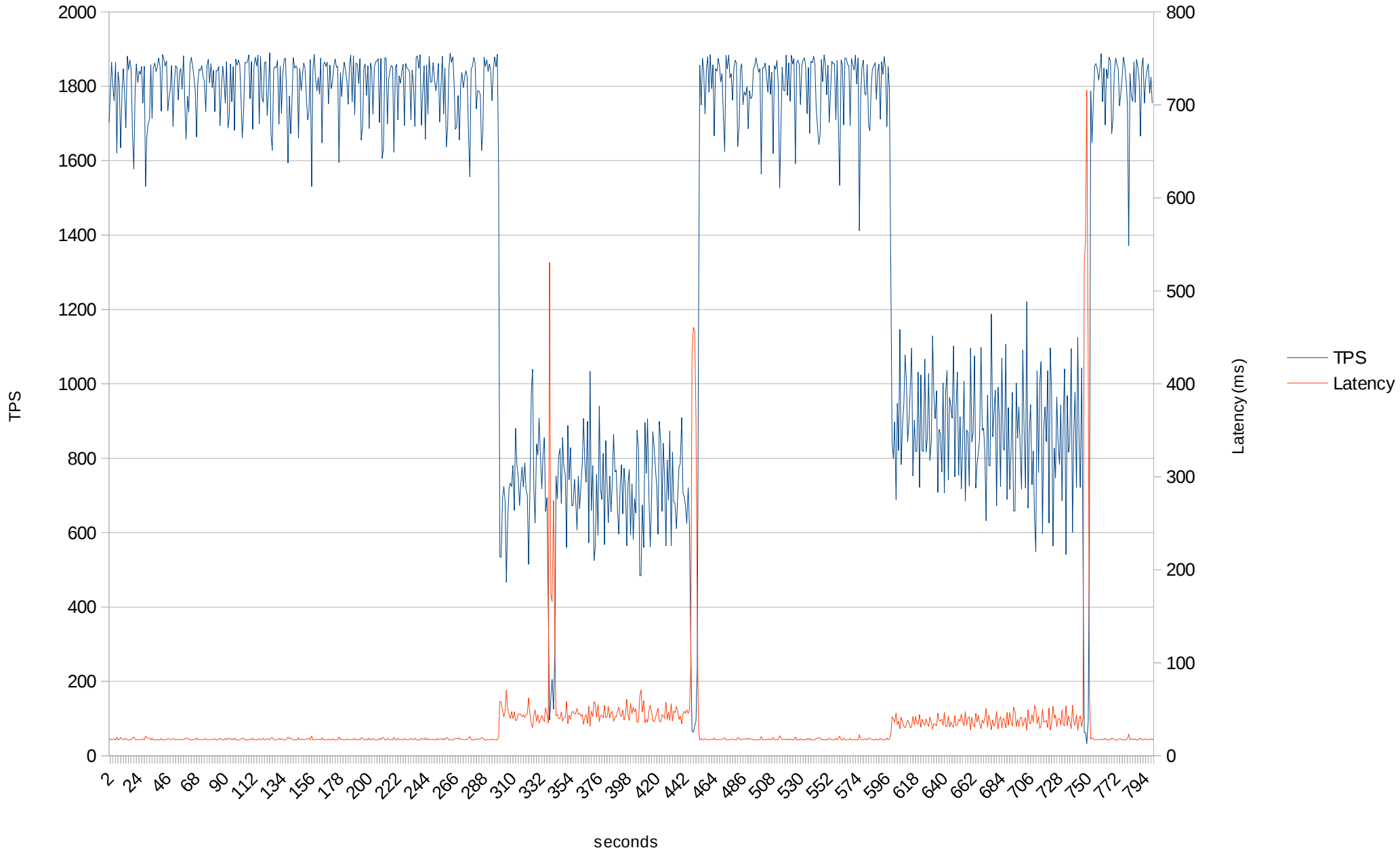


Dirty Data



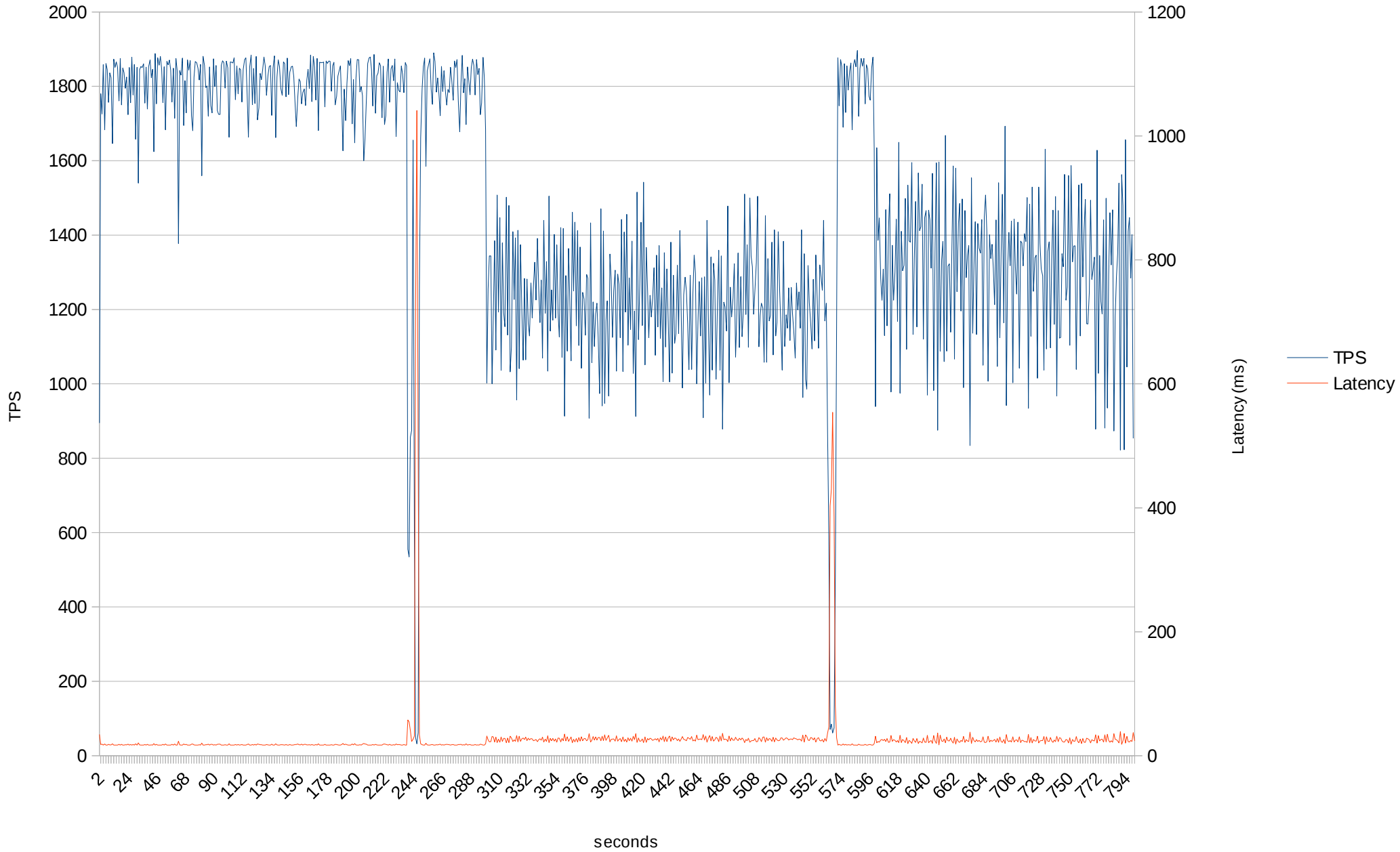
pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB, max_wal_size = 100GB, OS tuning (no dirty)



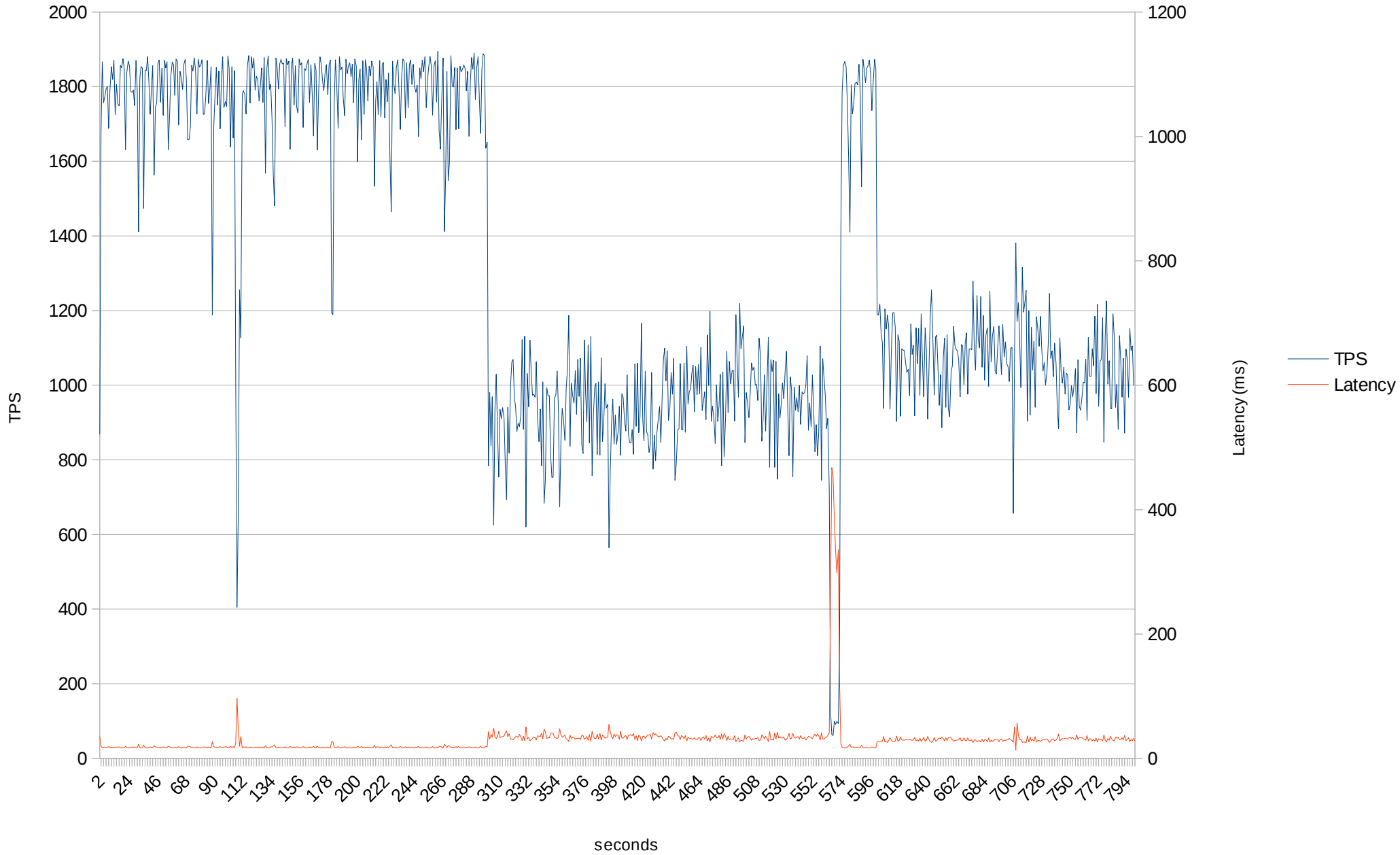
pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB, max_wal_size = 100GB, target = 0.9; OS tuning (no dirty)



pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB, max_wal_size = 100GB, target = 0.9; 9.6 flushing



Shared Buffers Tuning

- Hot data fits into shared_buffers => increase s_b
- Bulk-Writes in a bigger than shared_buffers workload => measure decreasing s_b
- Large Shared Buffers => enable huge pages
- Frequent Relation DROP/REINDEX => decrease s_b

OS Dirty Data Tuning

- `dirty_writeback_centisecs` => lower
 - how often to check for writeback
- `dirty_bytes/dirty_ratio` => lower
 - when to block writing data
- `dirty_background_bytes` => lower
 - when to write data back in the background
- Increases random writes!
- Often slows total throughput, but improves latency

WAL tuning

- Checkpoints should be triggered by time!
 - high enough `checkpoint_segments/wal_max_size`
 - Monitor!
- Except maybe at night, during batch runs or such
- Consider recovery time → less frequent checkpoints, crash recovery takes longer
- Consider full page writes → more frequent checkpoints mean much much more WAL
- separate `pg_xlog` can help a lot!

WAL Writer

- Writes WAL instead backends
- Important for `synchronous_commit = off`
- Otherwise boring

Background Writer

- Write dirty buffers before backends
- Not very good
- All random writes
- Defaults write max 4MB/s
- `bgwriter_delay` → lower, wakes up more often
- `bgwriter_lru_maxpages` → increases, writes more at once

Problem – Bad Benchmarks

- pgbench has unrealistic workload
- hard to measure regressions
- contribute!

Problem – Dirty Buffers in Kernel

- Massive Latency Spikes, up to hundreds of seconds
- Force flush using `sync_file_range()` or `msync()`
 - Decreases jitter
 - Increases randomness
- Sort checkpointed buffers
 - Decreases randomness
 - Increases Throughput
- Hopefully 9.6

Problem – Hashtable

- Can't efficiently search for the next buffer
 - need to sort for checkpoints
 - can't write combine to reduce total number of writes
- Expensive Lookups
 - Cache inefficient datastructure
- Possible Solution: Radix Tree
- Hopefully 9.7

Problem - Cache Replacement Scales Badly

- Single Lock for Clock Sweep!
 - fixed in 9.5
- Every Backend performs Clock Sweep
 - fixed in 9.6
- Algorithm is fundamentally expensive
 - UH, Oh.

Problem - Cache Replacement Replaces Badly

- Usagecount of 5 (max) reached very quickly
 - Often all buffers have 5
- Increasing max usagecount increases cost, the worst case essentially is
 $O(N_{\text{Buffer}} * \text{max_usagecount})$
- Hard to solve, patent issues

Problem: Kernel Page Cache

- Double buffering decreases effective memory utilization
- Use `O_DIRECT`?
 - Requires lots of performance work on our side
 - Considerably faster in some scenarios
 - Less Adaptive
 - Very OS specific