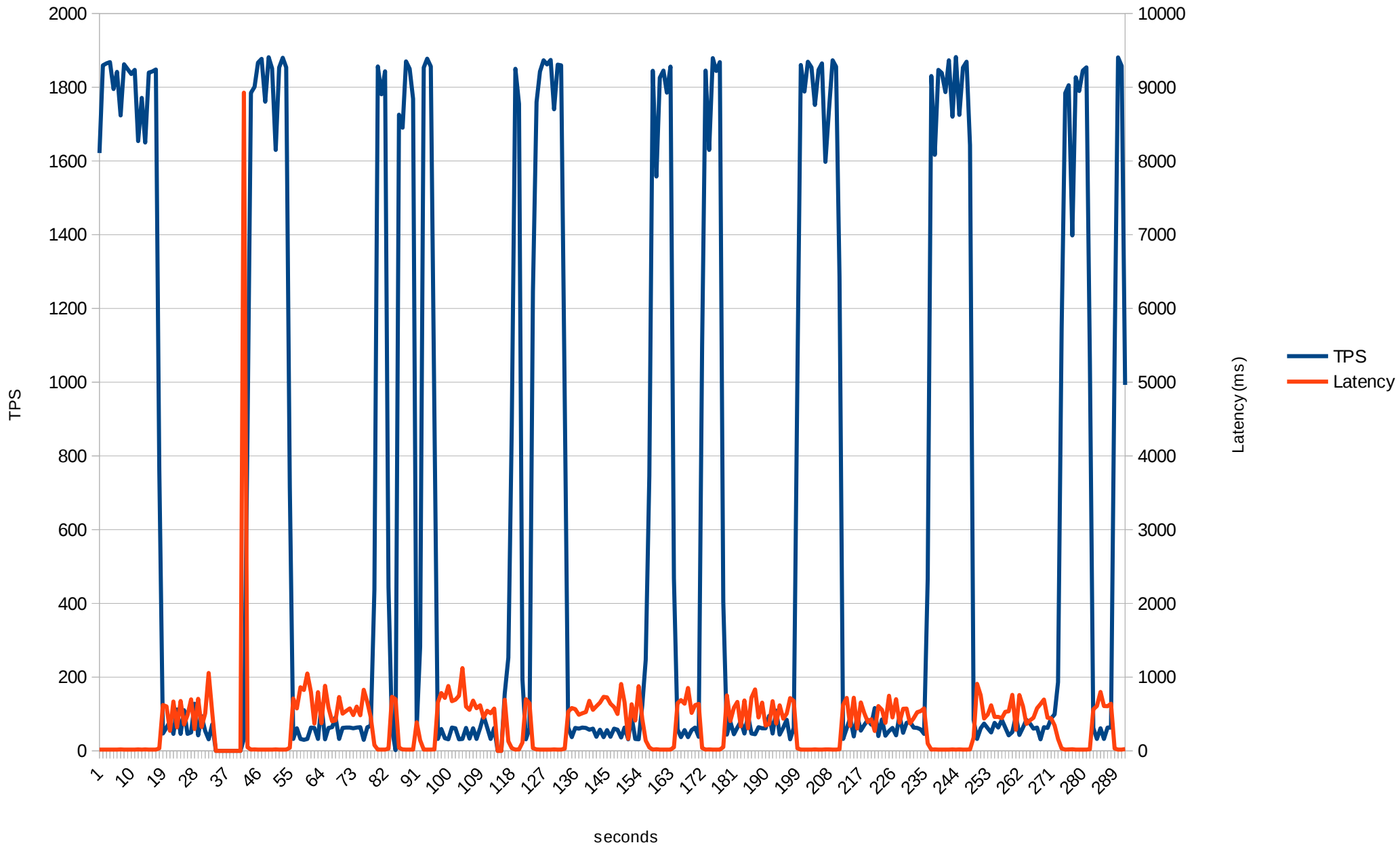# Postgres' IO
# Architecture, Tuning, Problems

## Andres Freund
## PostgreSQL Developer & Committer
## Citus Data – citusdata.com - @citusdata
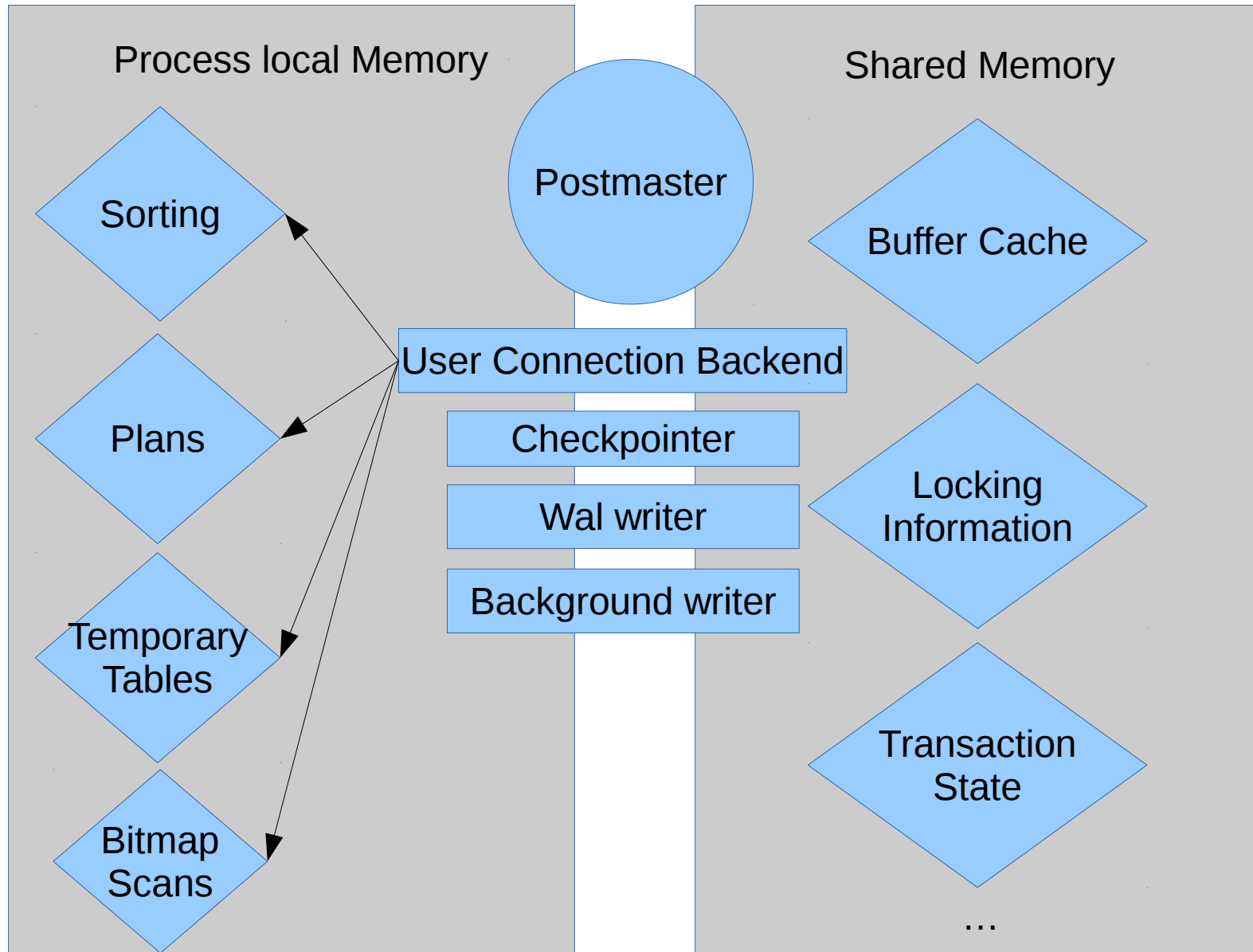
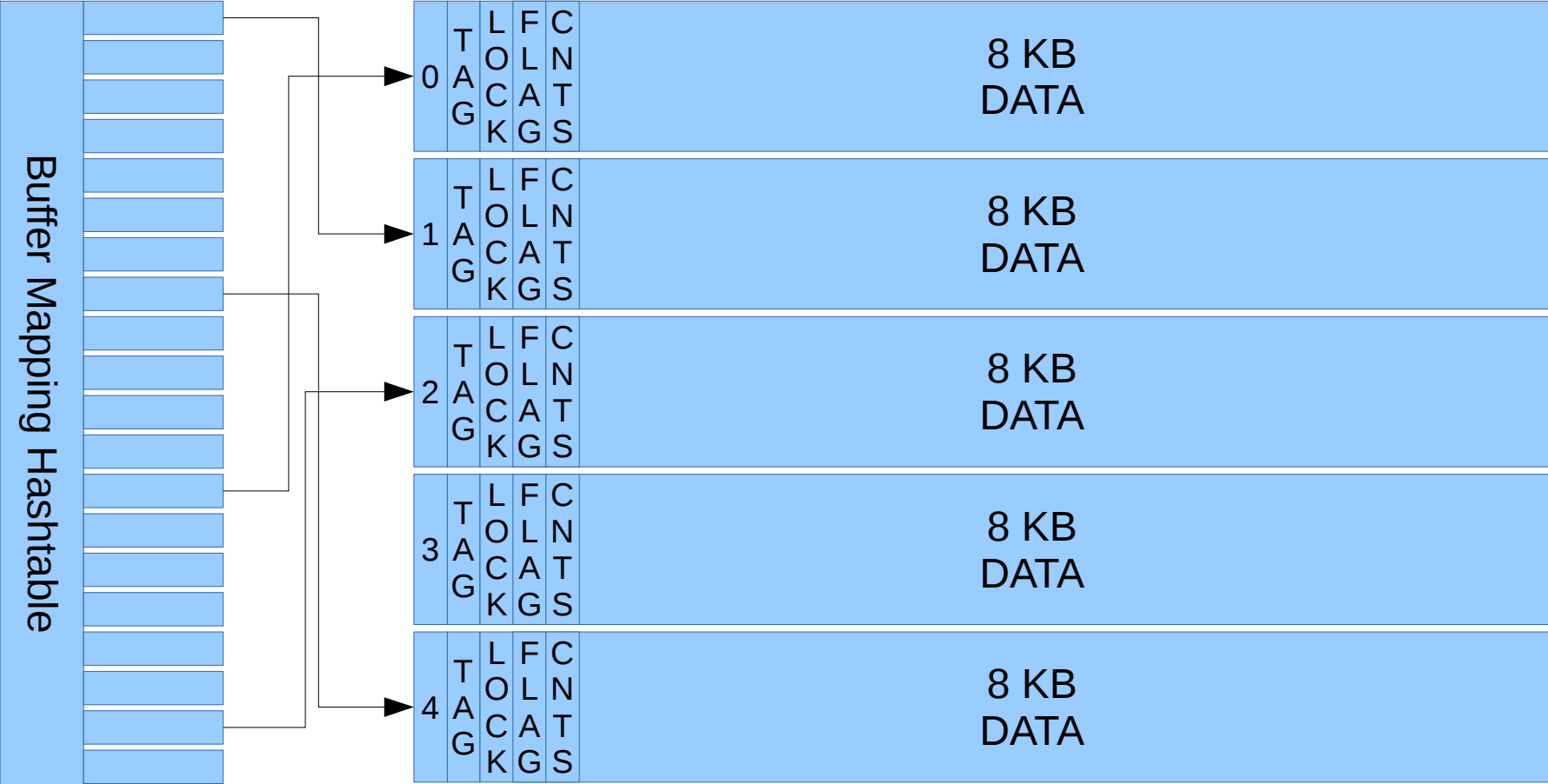http://anarazel.de/talks/berlin-meetup-2016-01-26/io.pdf

**citusdata**

# pgbench -M prepared -c 32 -j 32

## standard settings



citusdata

# Memory Architecture

## Process local Memory

Sorting

Plans

Temporary
Tables

Bitmap
Scans

Postmaster

User Connection Backend

Checkpointer

Wal writer

Background writer

## Shared Memory

Buffer Cache

Locking
Information

Transaction
State

…

citusdata

# Shared Buffers

# Reading Data



Buffer Mapping Hashtable

Storage

OS PageCache

open()
read()

3 | TAG | LOCK | FLAG | CNTS | 8 KB DATA

citusdata

# Clock-Sweep



CNT: 4 3

CNT: 0

35 0 1 2 3 4 5 6

citusdata

# Writing Data Out



Buffer Mapping Hashtable

3 | TAG | LOCK | FLAG | CNTS | 8 KB DATA

OS PageCache

open()
write()

citusdata

# Need for WAL logging

| | TAG | LOCK | FLAG | CNTS | |
|---|---|---|---|---|---|
| 0 | TAG | LOCK | FLAG | CNTS | 8 KB DATA |
| 1 | TAG | LOCK | FLAG | CNTS | 8 KB DATA |
| 2 | TAG | LOCK | FLAG | CNTS | 8 KB DATA |
| 3 | TAG | LOCK | FLAG | CNTS | 8 KB DATA |
| 4 | TAG | LOCK | FLAG | CNTS | 8 KB DATA |

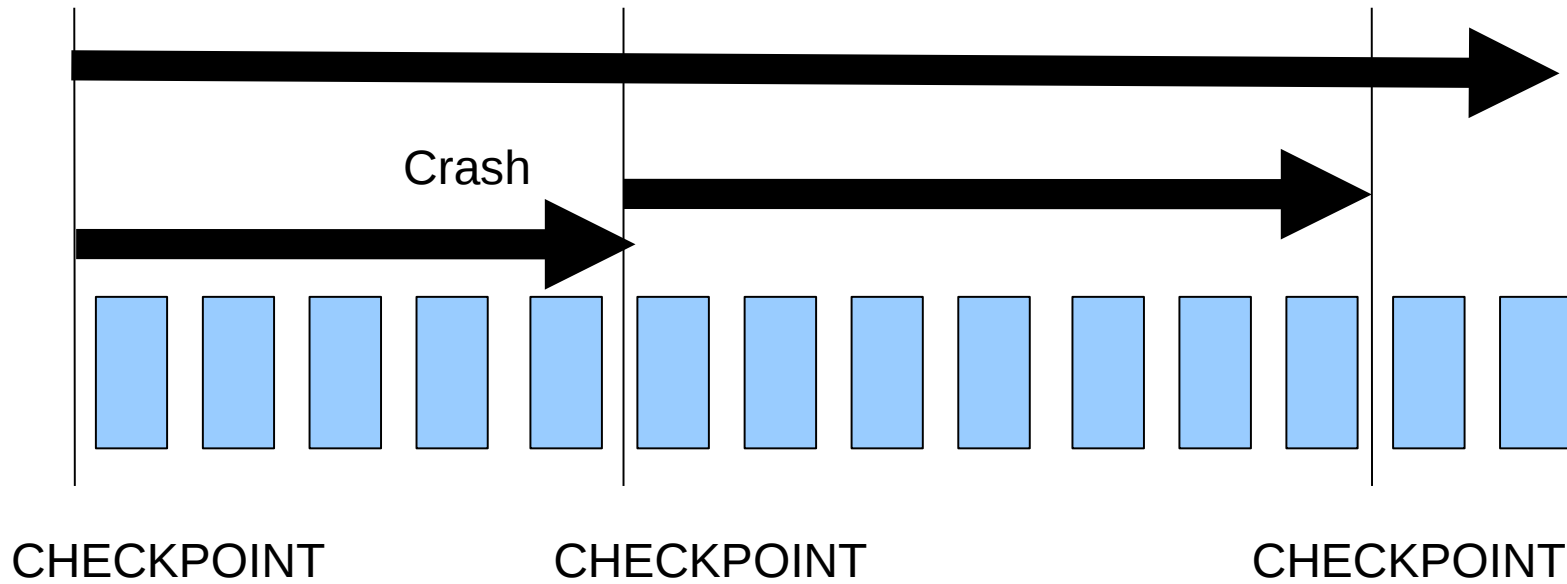citusdata

# WAL Logging

BEGIN;
INSERT INTO mytbl VALUES (…);
COMMIT;


rmgr: Heap          tx:  344576442, lsn: 317/A2333ED8, desc: INSERT off 3 blkref
< actually modify heap >
rmgr: Btree         tx:  344576442, lsn: 317/A2333F18, desc: INSERT_LEAF off 2
< actually modify btree >
rmgr: Transaction tx:  344576442, lsn: 317/A2333F58, desc COMMIT 2016-01-25 15:17:30
< actually modify transaction >

# 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.
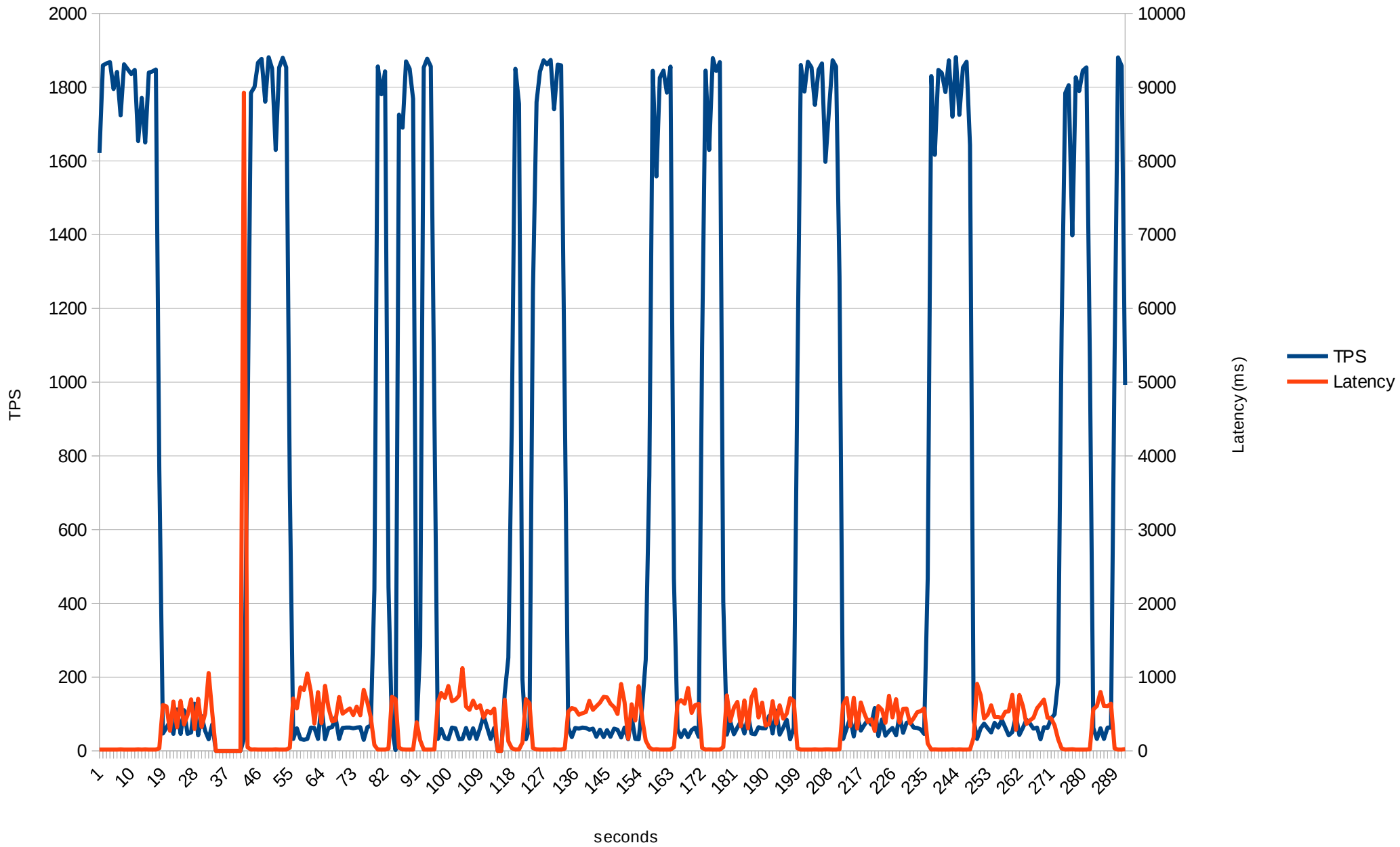
6)Remove old WAL

# 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;)

citusdata

# Spreading Checkpoints

- checkpoint_completion_target = 0.5

- estimation based on
  - checkpoint_timeout
  - checkpoint_segments/max_wal_size

- Start the next checkpoint after checkpoint_completion_target * above_param has passed

- Try to keep pace

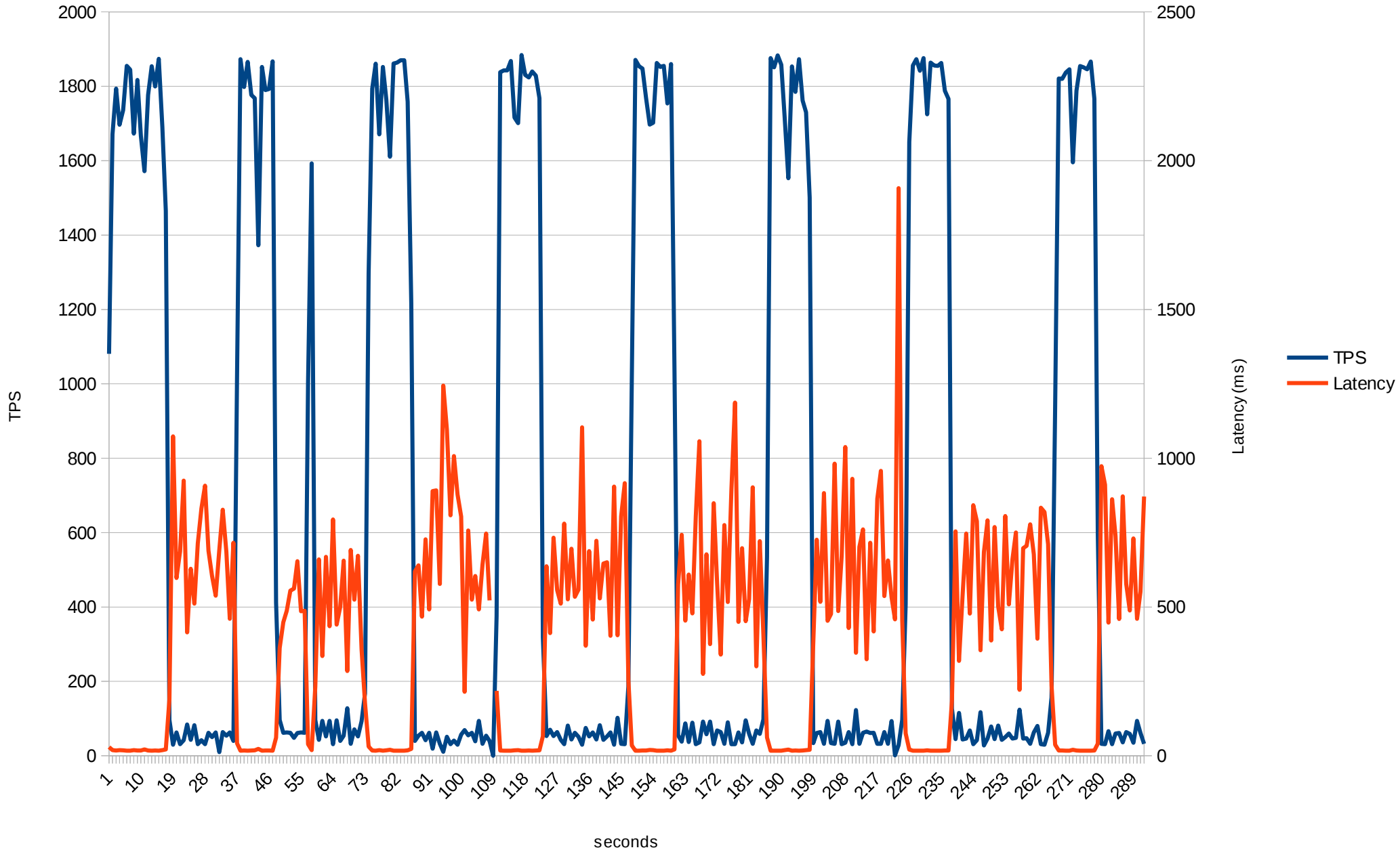citusdata

# pgbench -M prepared -c 32 -j 32

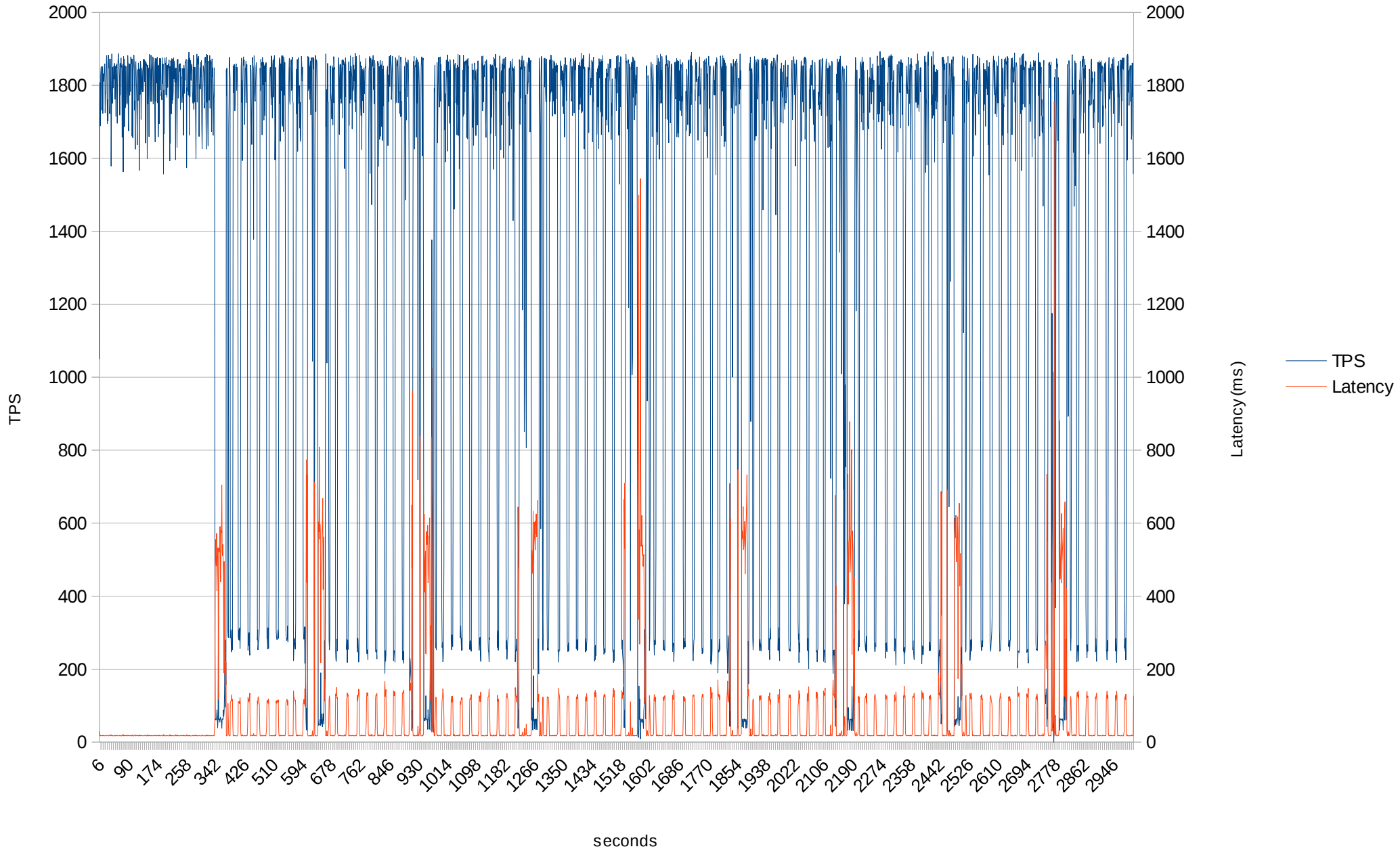## standard settings



citusdata

pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB

pgbench -M prepared -c 32 -j 32
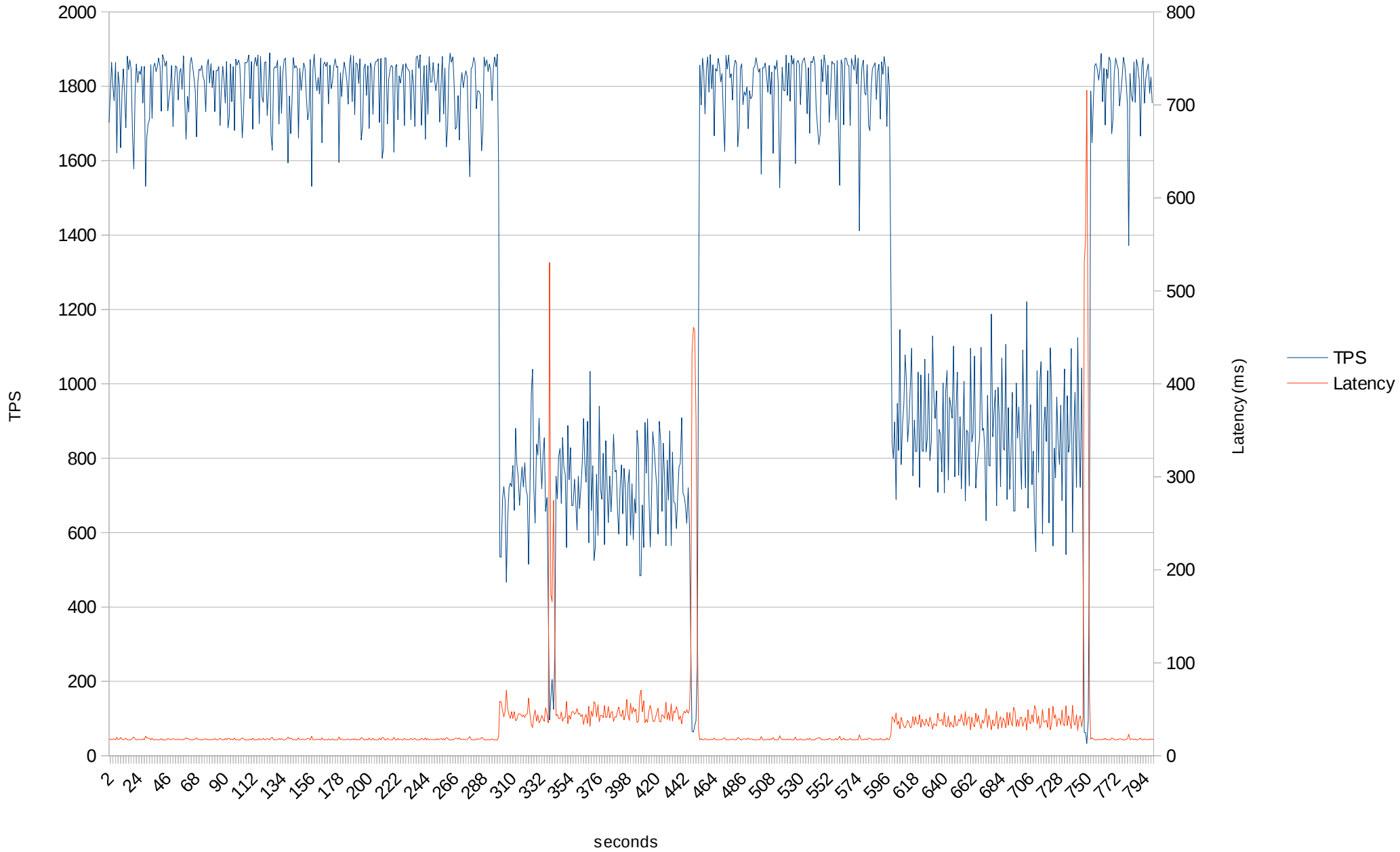
shared_buffers = 16GB, max_wal_size = 100GB

citusdata

Dirty Data

time (seconds)

citusdata

pgbench -M prepared -c 32 -j 32

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

citusdata

# 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

citusdata

# 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!

citusdata

# 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

**citusdata**

# Problem – Bad Benchmarks

- pgbench has unrealistic workload

- hard to measure regressions

- contribute!

citusdata

# 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

citusdata

# 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

citusdata

# 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.

citusdata

# 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(NBuffer * max_usagecount)
- Hard to solve, patent issues

citusdata

# 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