

Improving Postgres' Buffer Manager

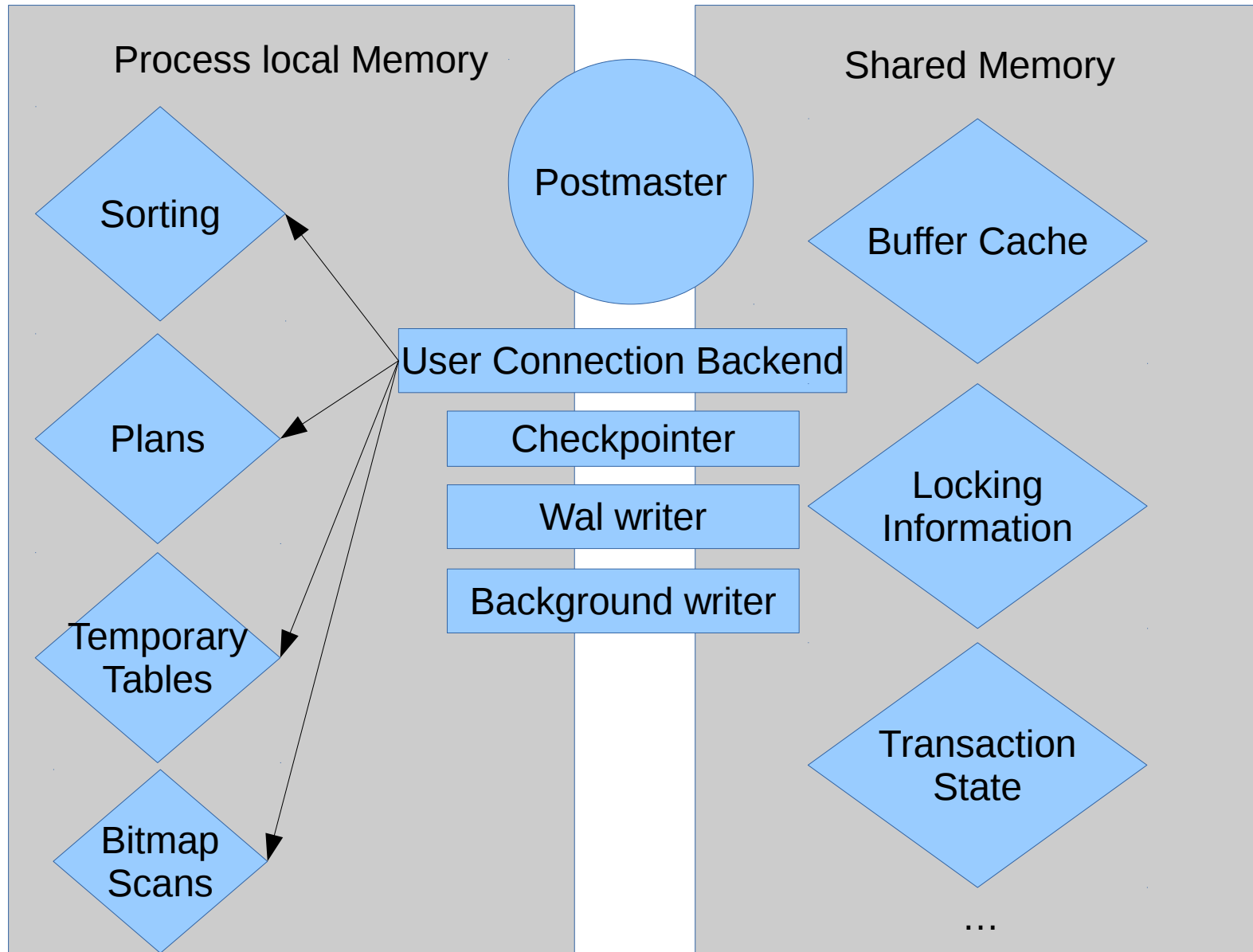
Andres Freund

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

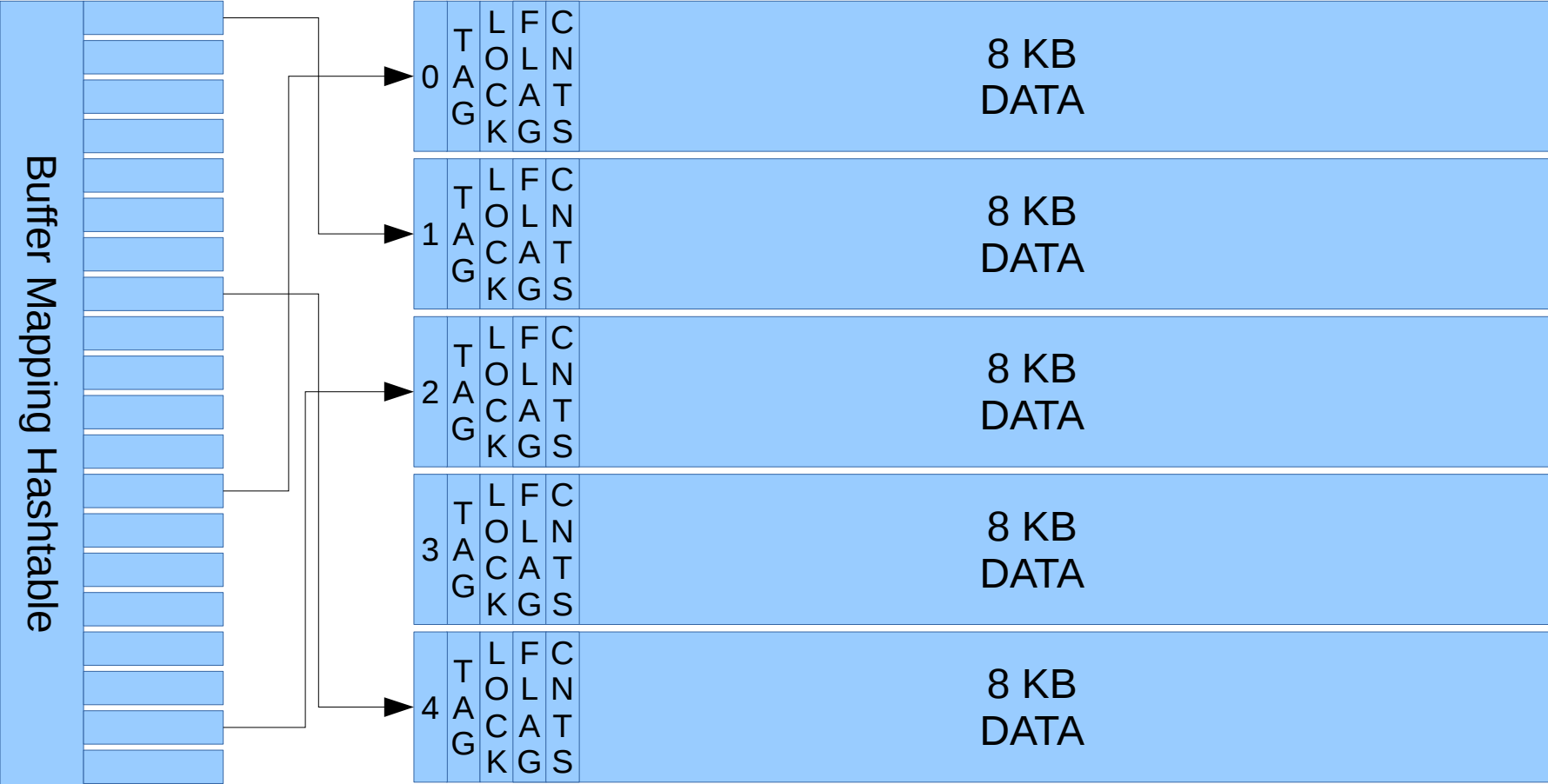
<http://anarazel.de/talks/pgconf-ru-2016-02-04/io.pdf>



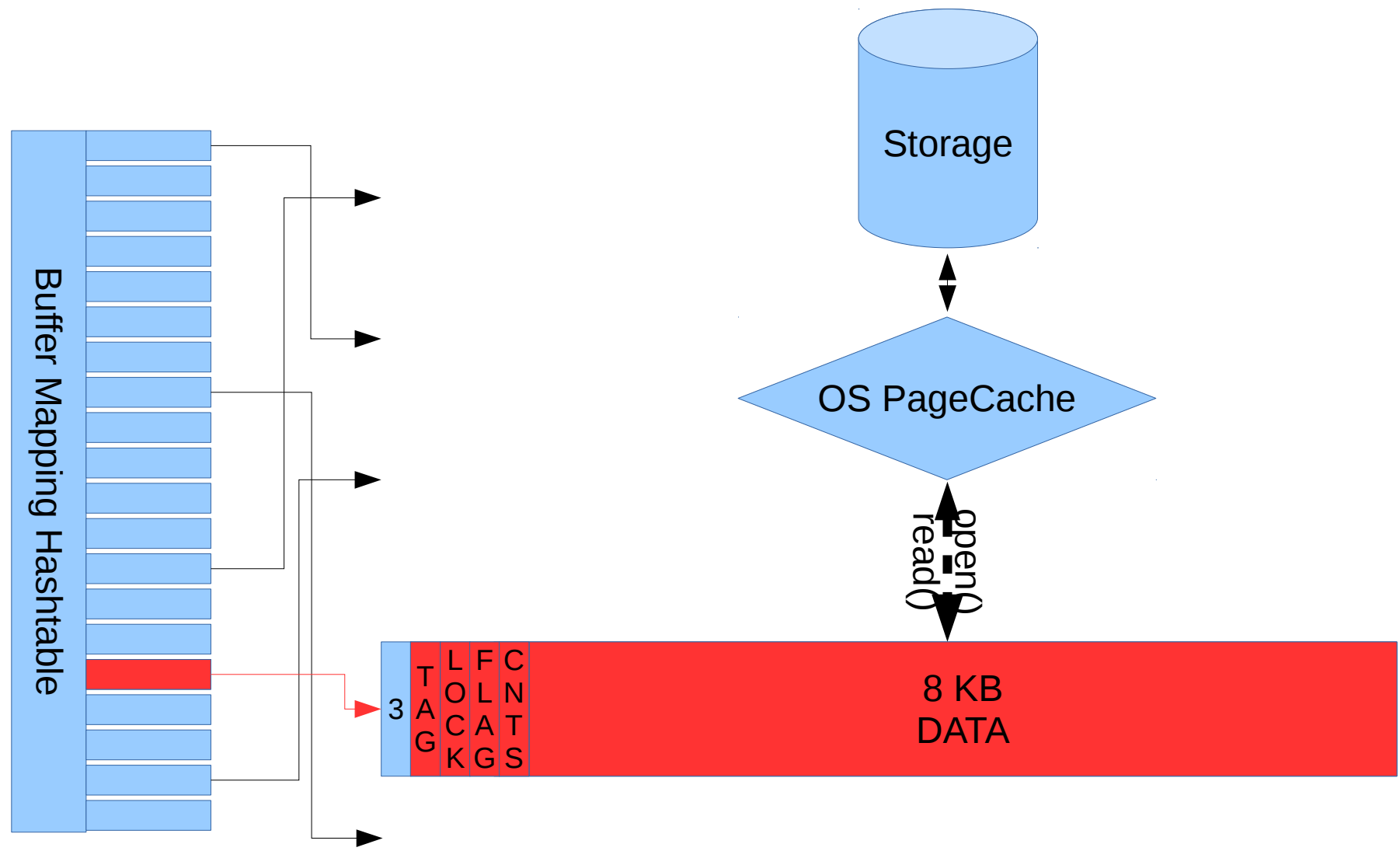
Memory Architecture



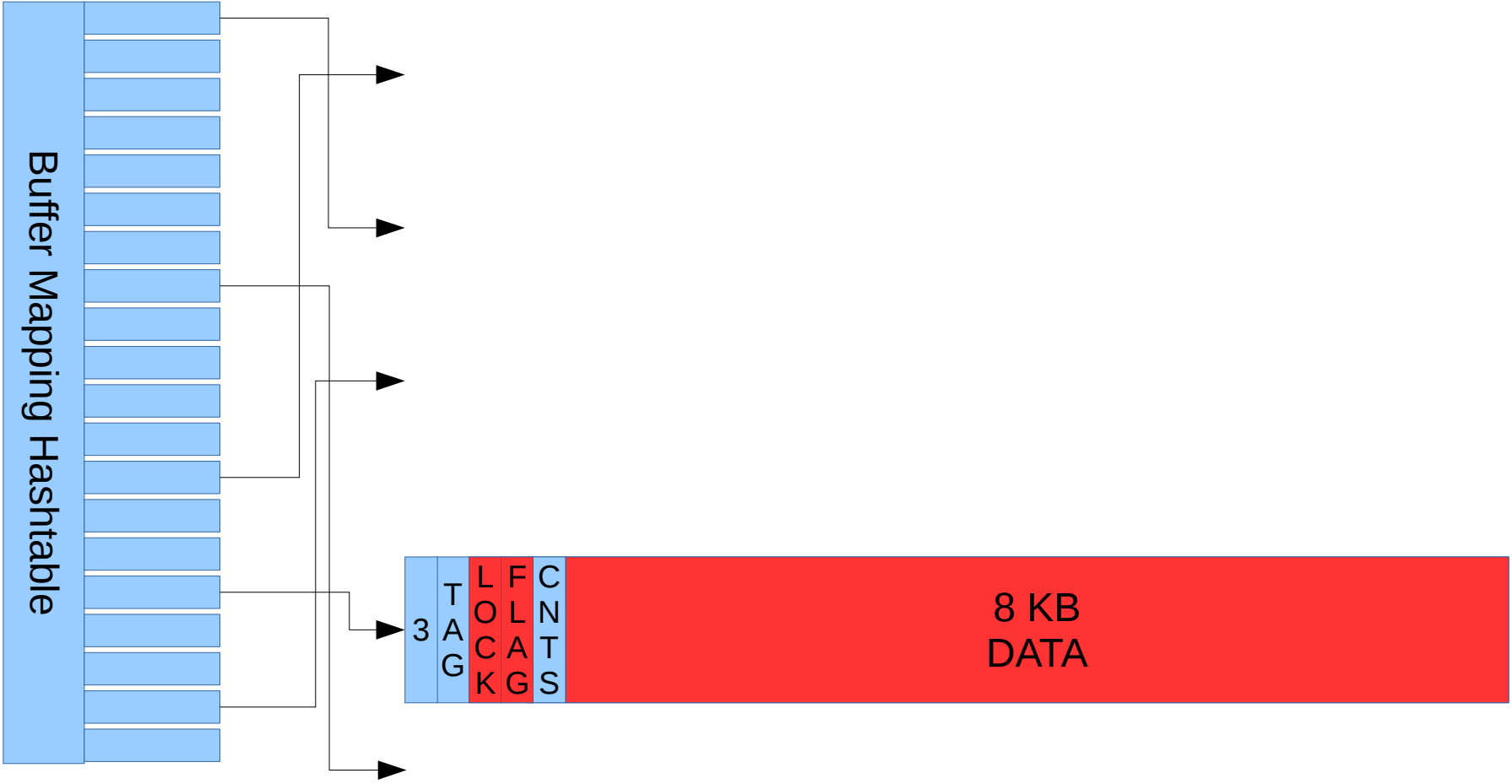
Shared Buffers



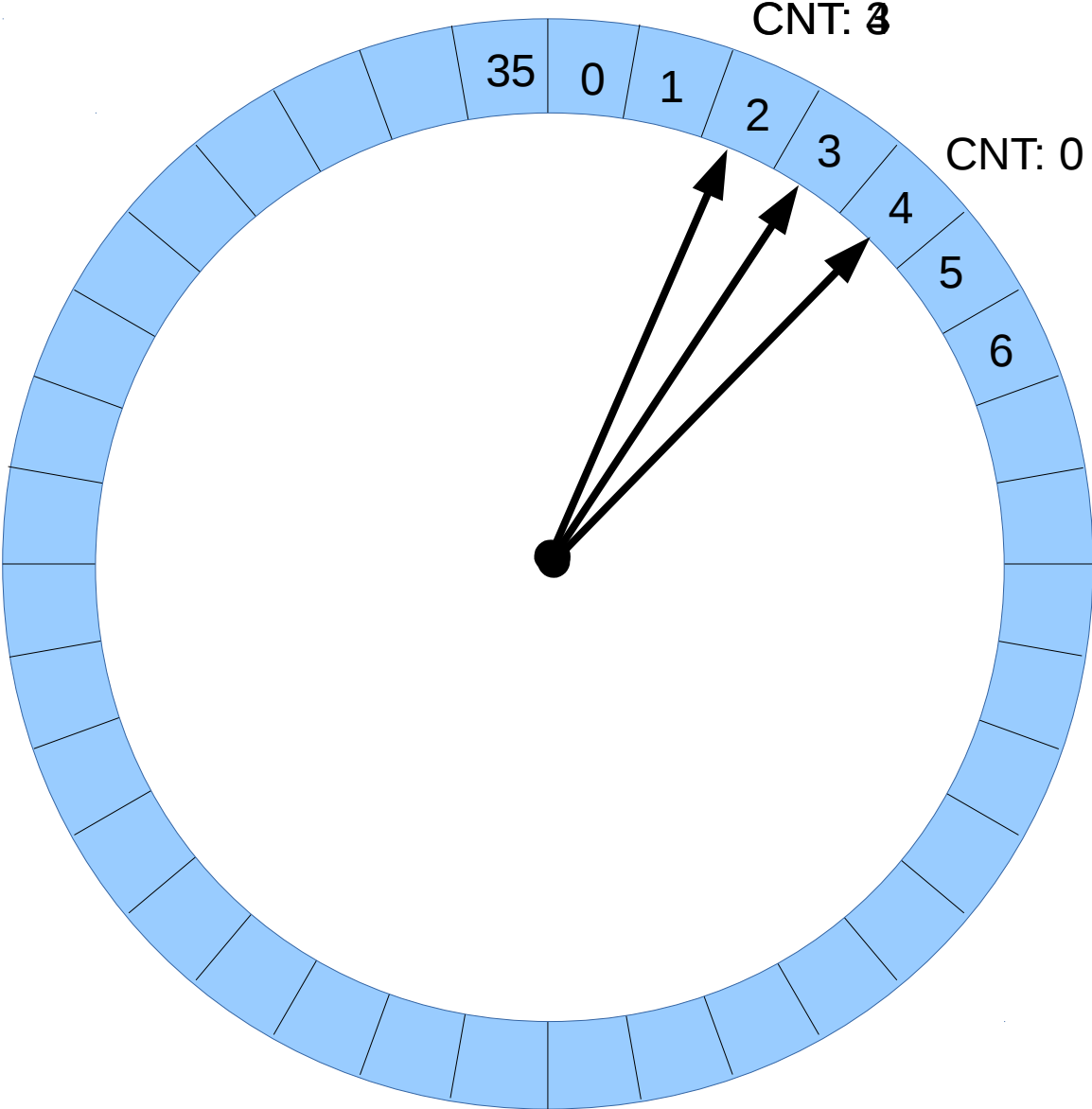
Reading Data



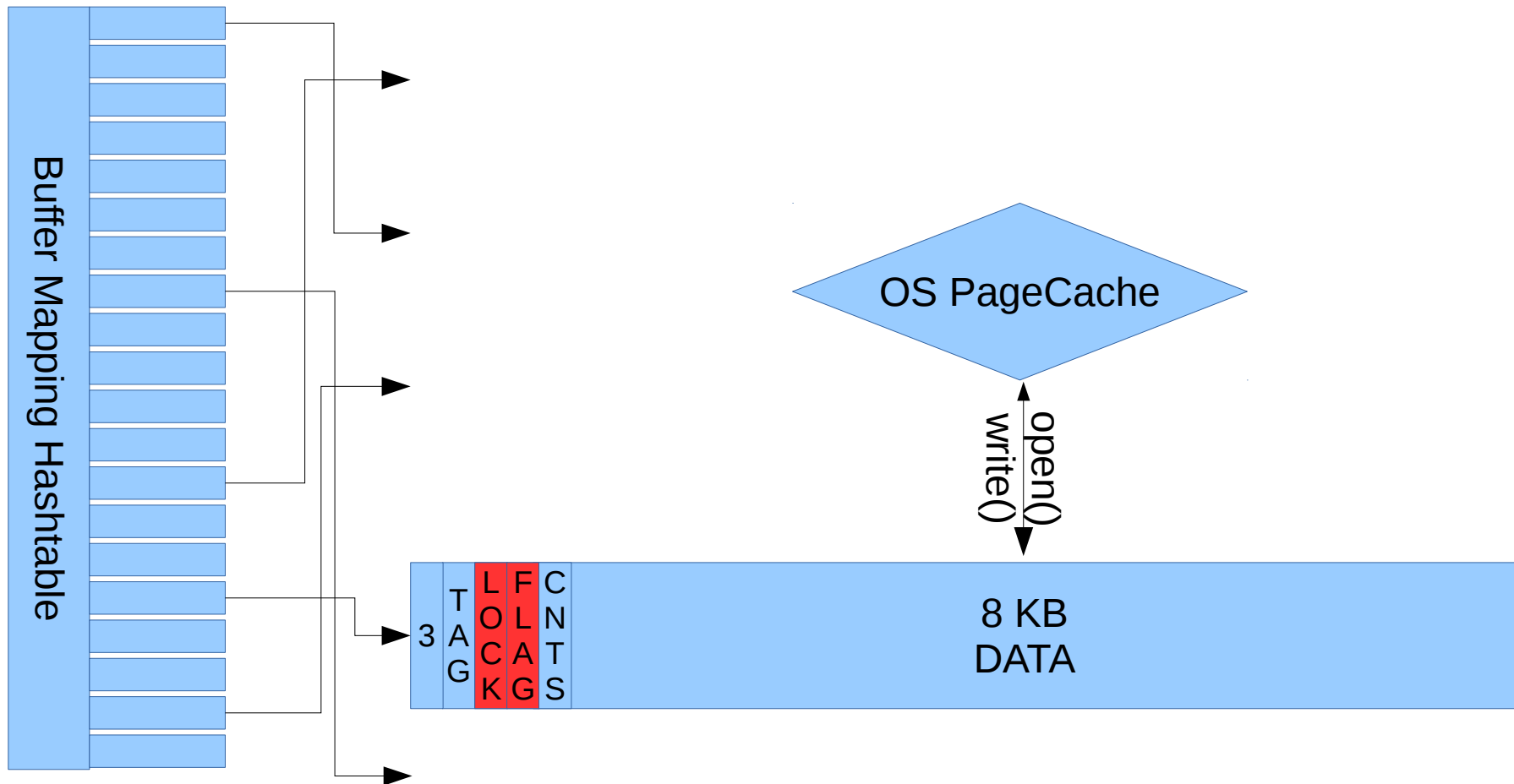
Writing Data



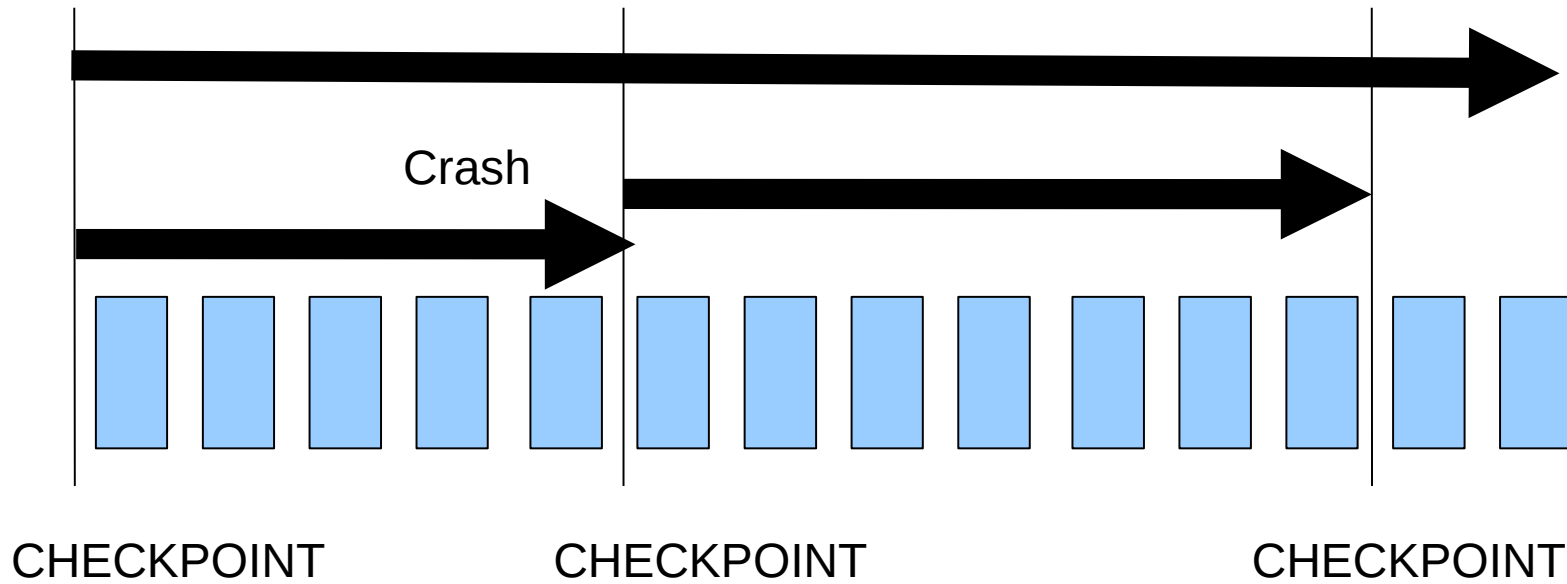
Clock-Sweep



Writing Data Out



Recovery & Checkpoints

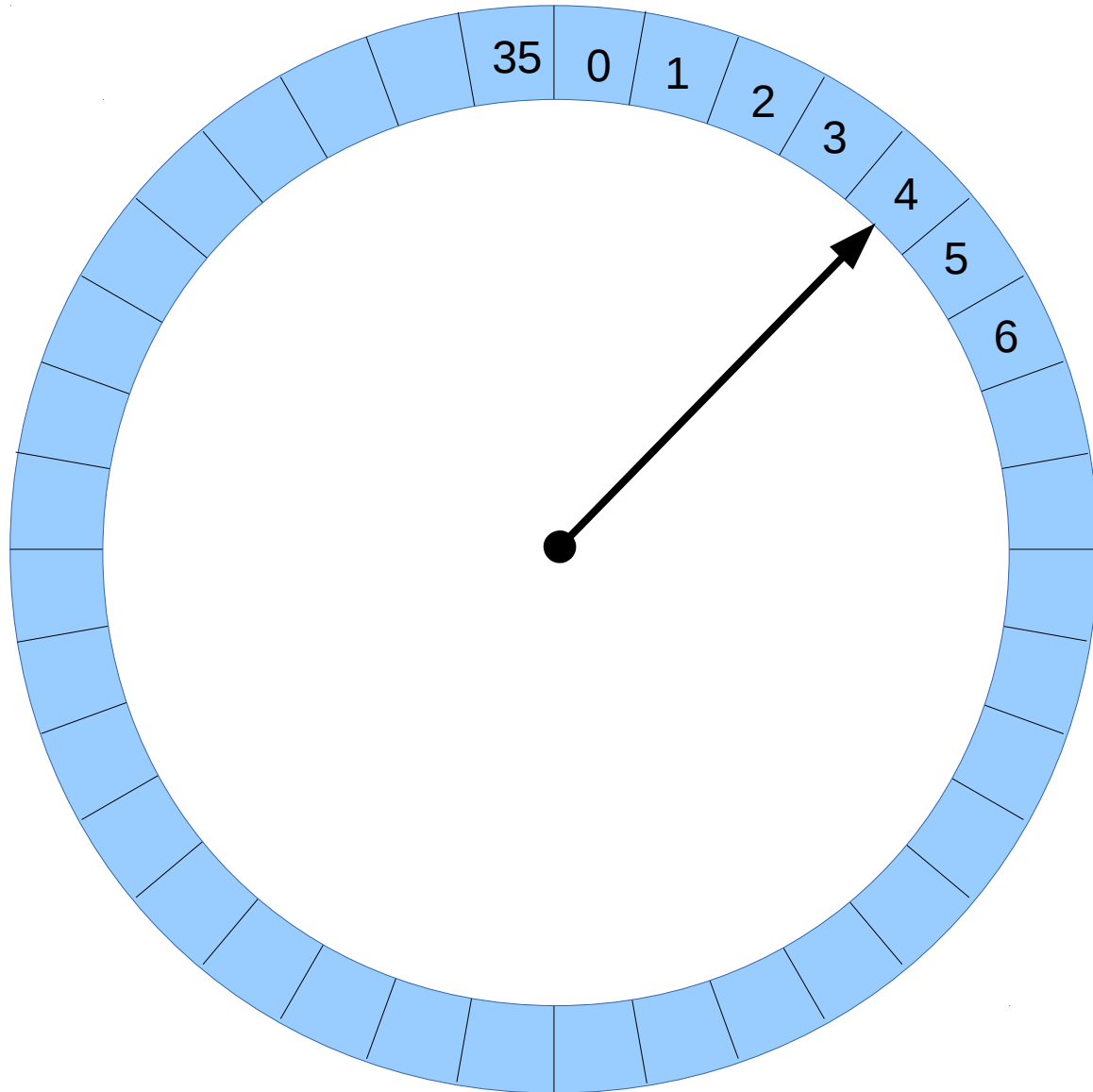


Checkpoints

- 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

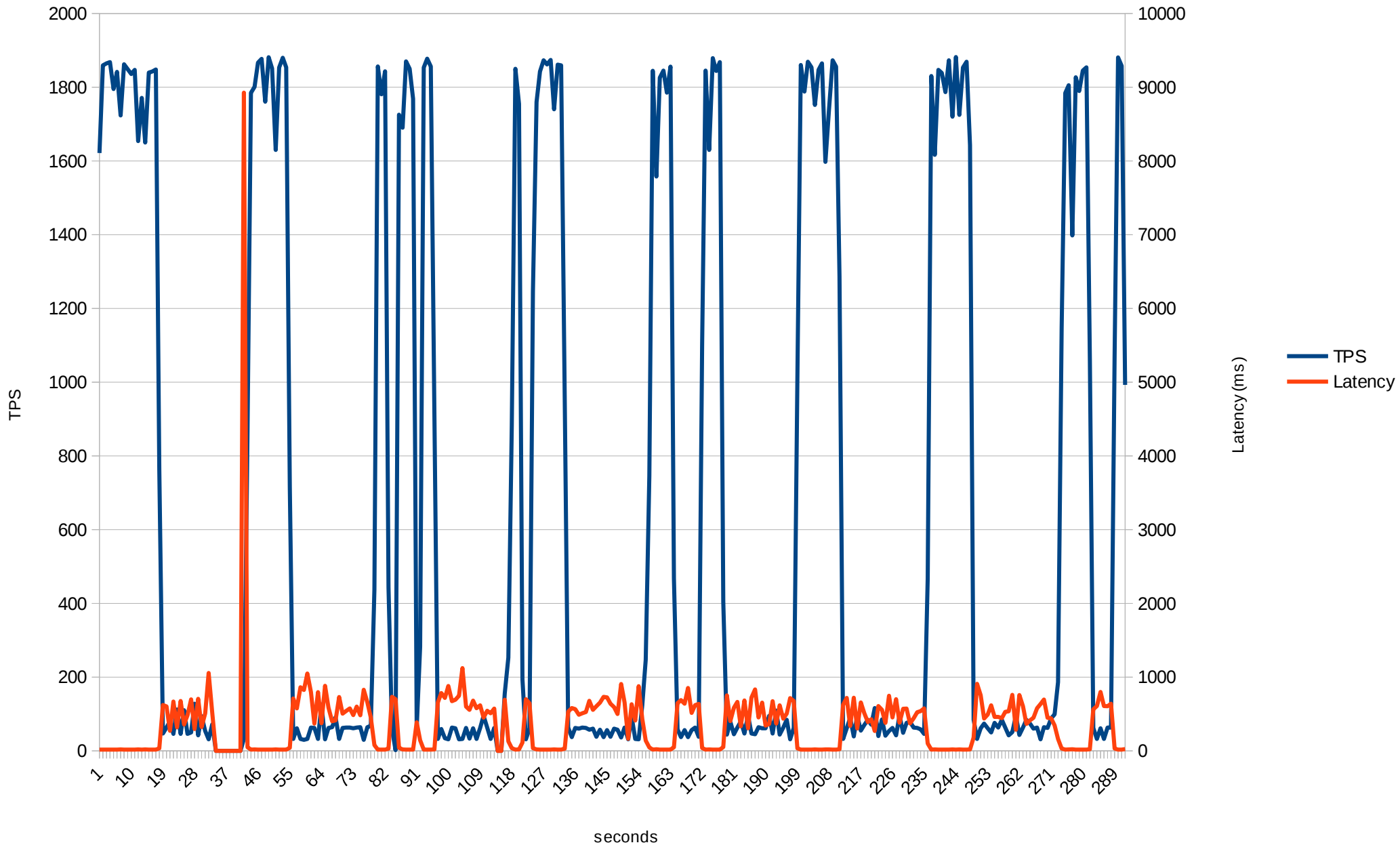


Checkpoint Writeout



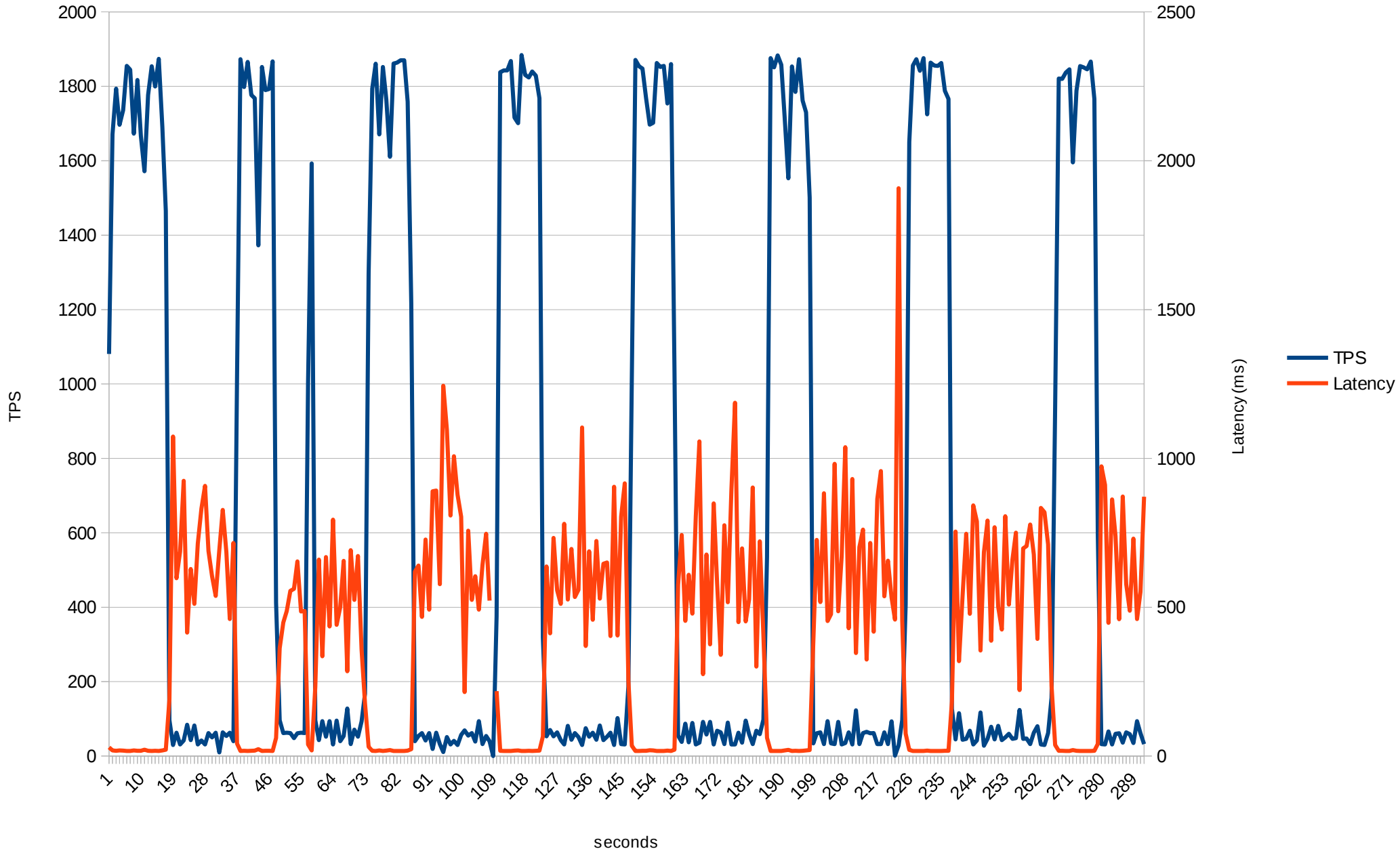
pgbench -M prepared -c 32 -j 32

standard settings



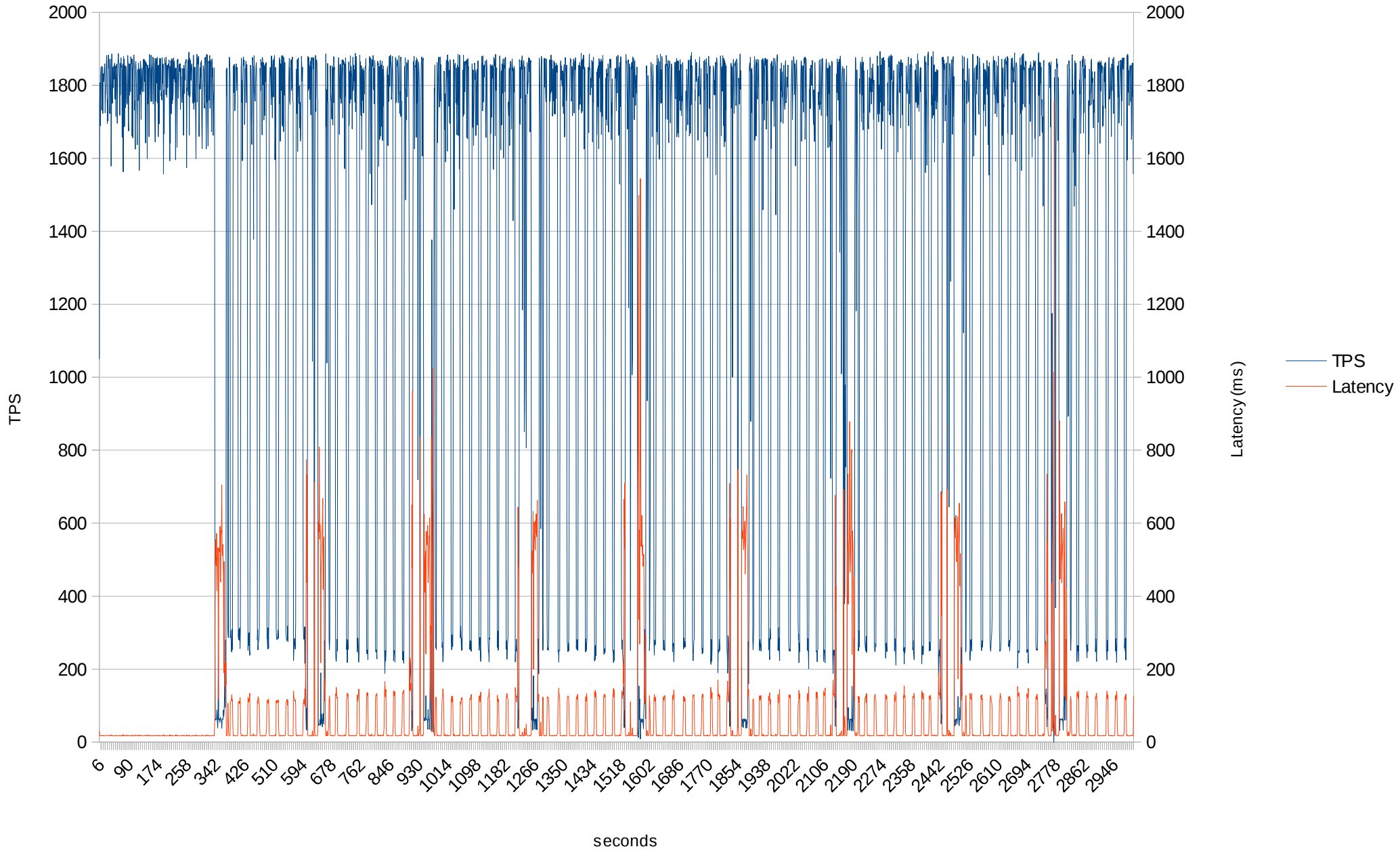
pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB



pgbench -M prepared -c 32 -j 32

shared_buffers = 16GB, max_wal_size = 100GB



Dirty Data



Problem – Dirty Buffers in Kernel

- Massive Latency Spikes, up to hundreds of seconds
- No actually efficient merging of IO requests
- latency spikes every `dirty_writeback_centisecs`
- spikes when reaching `dirty_{background_,}ratio`
- latency spikes after checkpoint's `fsync()`



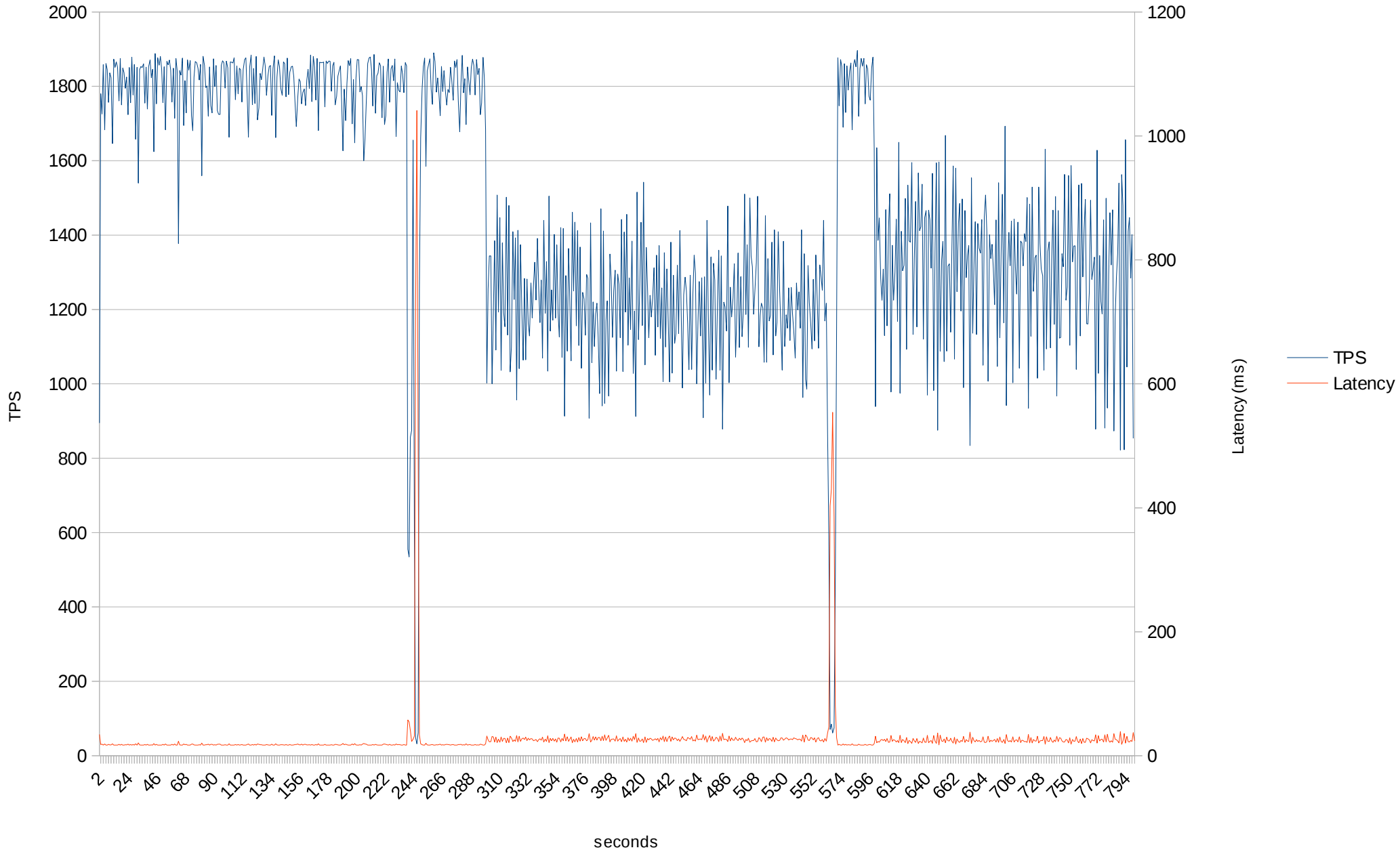
OS Dirty Data Tuning

- `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!
- Systemwide. Ooops.



pgbench -M prepared -c 32 -j 32

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



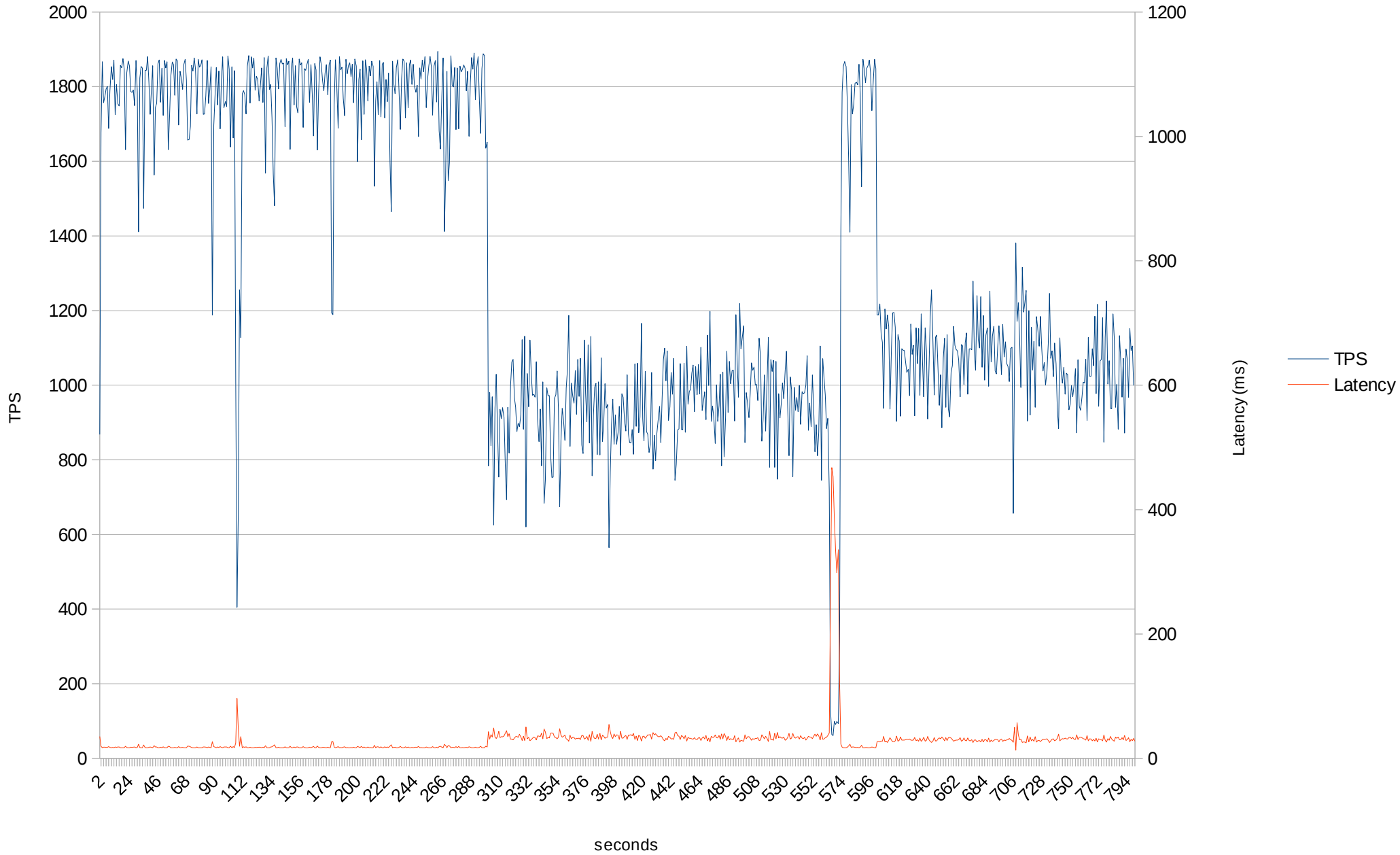
Dirty Buffers in Kernel

- Force flush using `sync_file_range()` or `msync()`
 - Decreases jitter
 - Increases randomness
 - Flushes need to happen in checkpoint, bgwriter, backends
- Sort to-be-checkpointed buffers
 - Decreases randomness
 - Increases Throughput
- Hopefully 9.6



pgbench -M prepared -c 32 -j 32

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



Problem – Hashtable

- Expensive Lookups
 - Wide keys (20 bytes)
 - Cache inefficient datastructure (spatial locality)
- Can't efficiently search for the next buffer
 - need to sort for checkpoints
 - can't write combine to reduce total number of writes
- Dropping relations very expensive
- Possible Solution: Tree of Radix Trees
- Hopefully 9.7



```
typedef struct BufferTag
{
    struct RelFileNode
    {
        Oid          spcNode;      /* tablespace */
        Oid          dbNode;      /* database */
        Oid          relNode;     /* relation */
    } rnode;                     /* physical relation identifier */

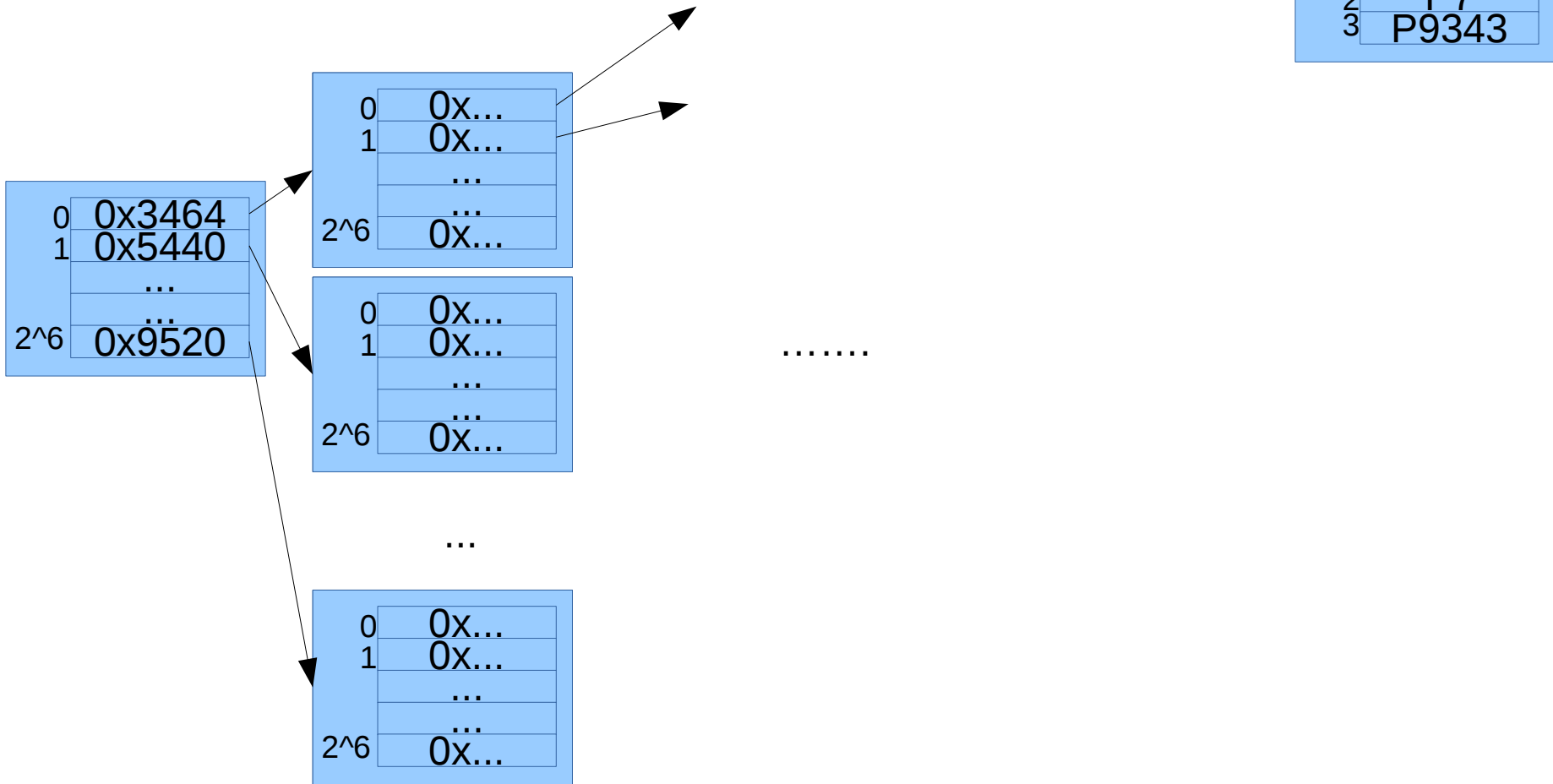
    ForkNumber      forkNum;

    BlockNumber     blockNum;     /* blknum relative to begin of reln */
} BufferTag;
```



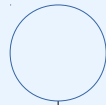
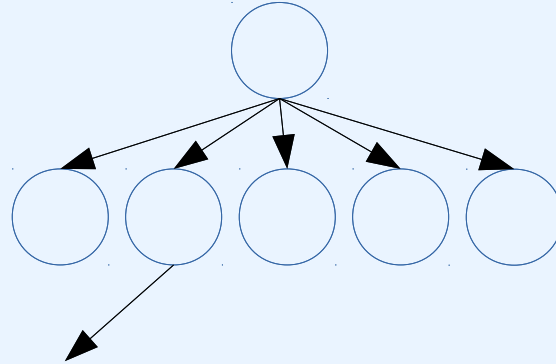
Radix Tree "Linux Style"

1 2 3 4 5 6
6b 6b 6b 6b 6b 2b



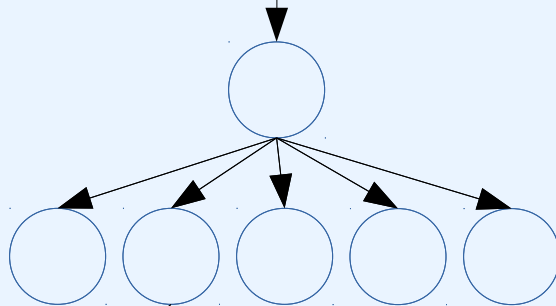
Tree of Trees

Upper Tree (tablespace, database, relfilenode, fork)



cached in SmgrRelation / permanent datastructure?

Lower Tree (block number)



Buffer



Problem - Cache Replacement Scales Badly

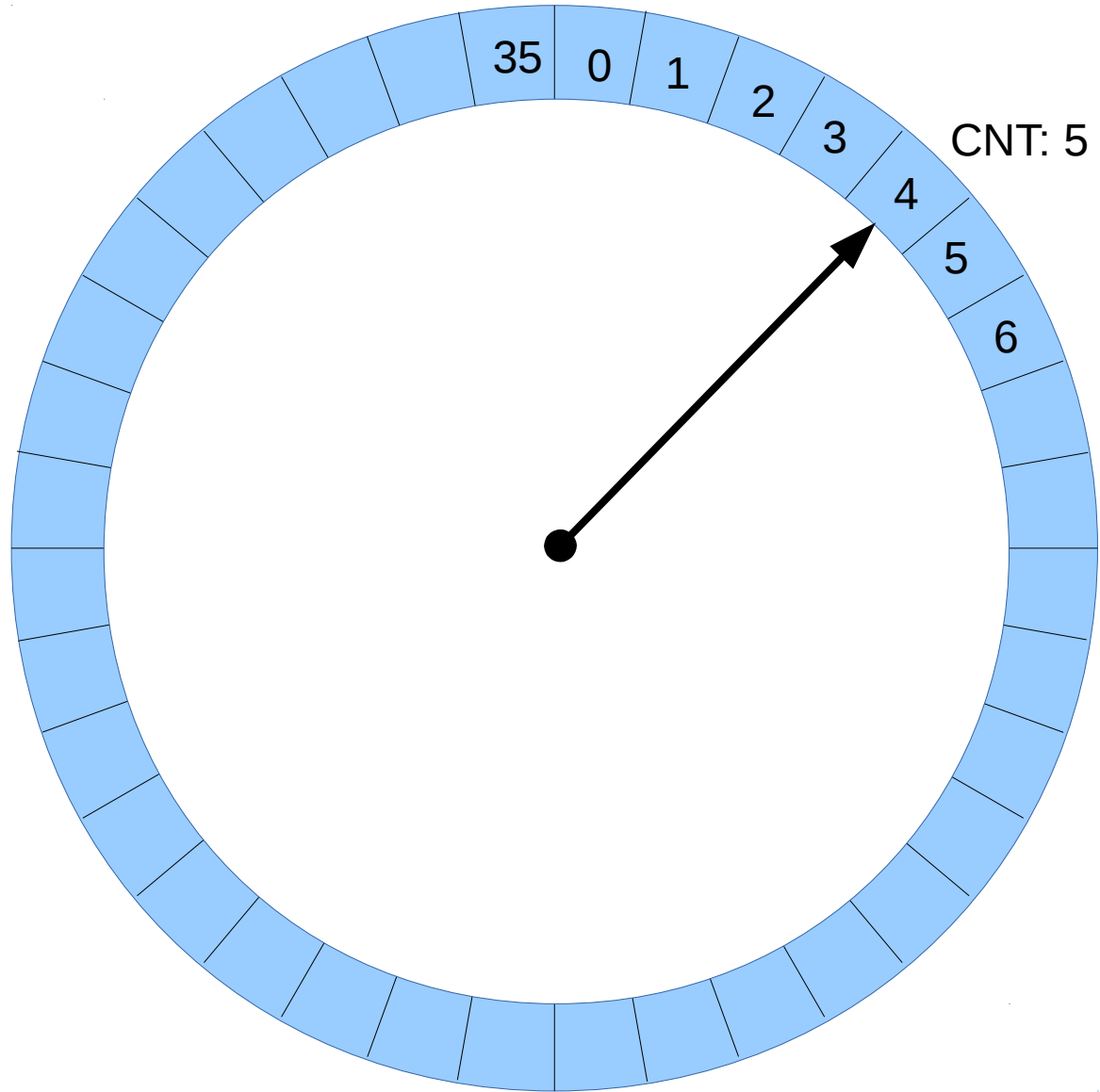
- Single Lock for Clock Sweep!
 - fixed in 9.5
- Every Backend performs Clock Sweep
 - can be moved to separate process (patches exist)
- Algorithm is fundamentally expensive
 - UH, Oh.
 - Worst case essentially is having to touch $N\text{Buffers} * 5 \text{ Buffers}$



Problem - Cache Replacement Replaces Badly

- Usagecount of 5 (max) reached very quickly
 - Often all buffers have 5
 - only works well if replacement rate is higher than average usage rate
 - very expensive form of random replacement
- Increasing max usagecount increases cost, the worst case essentially is
 $O(N_{\text{Buffer}} * \text{max_usagecount})$
- Hard to solve





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 (to be fast)



Improving Postgres' Buffer Manager

Andres Freund

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

<http://anarazel.de/talks/pgconf-ru-2016-02-04/io.pdf>

