

PGConf Europe 2024

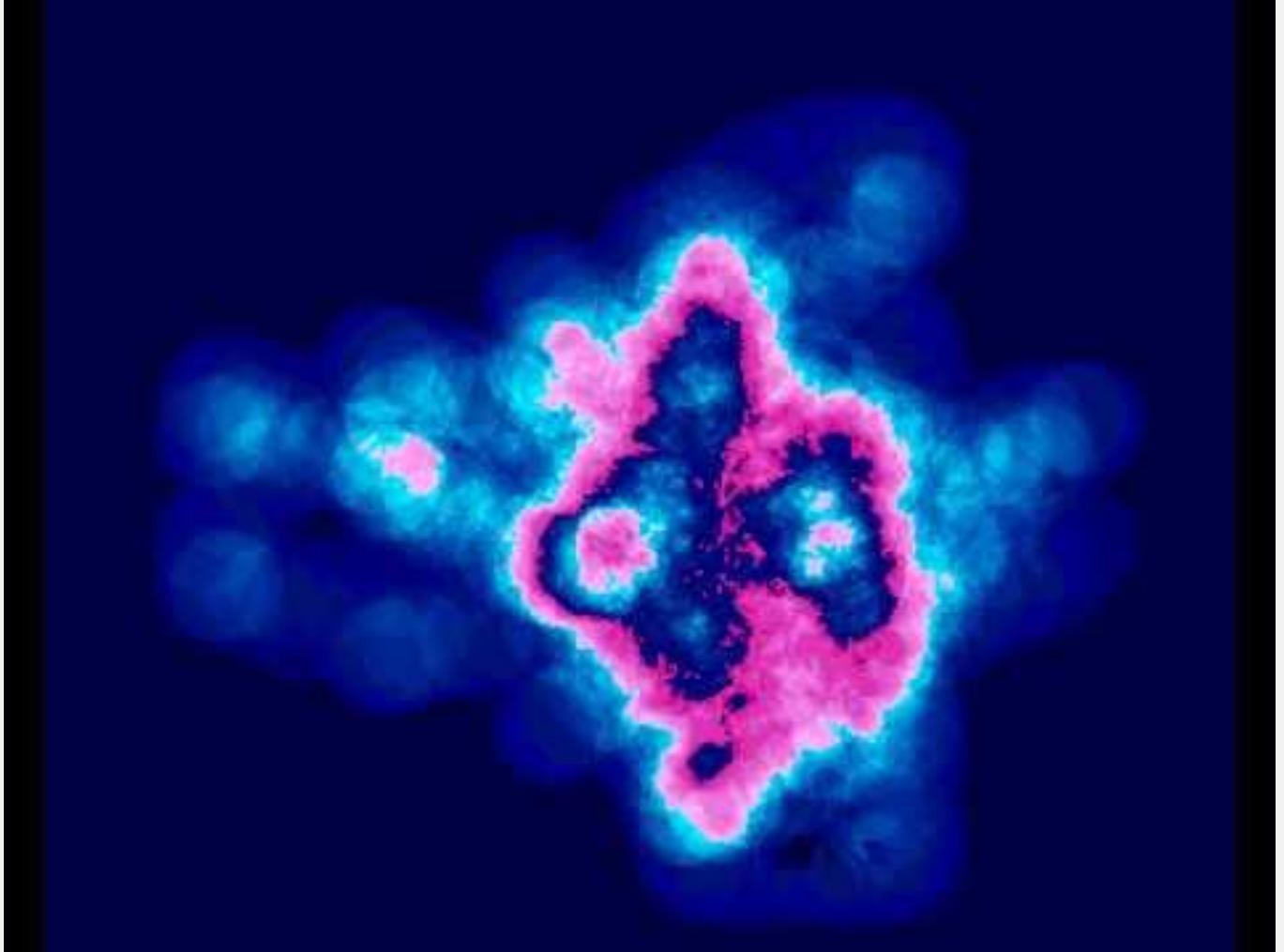
Are you collecting the right metrics?

frederic.delacourt@data-bene.io

Frédéric Delacourt

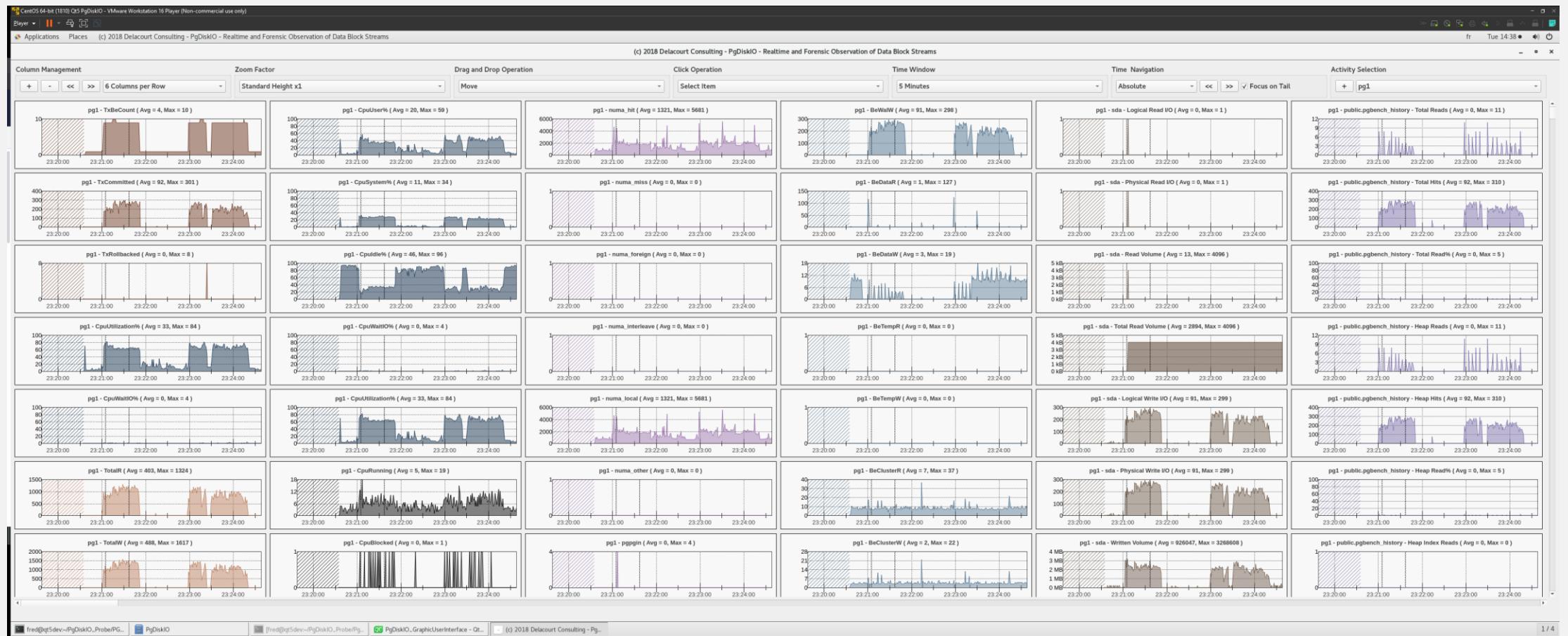


Demomaker
on Amiga 500



Frédéric Delacourt

Creator of PgDiskIO

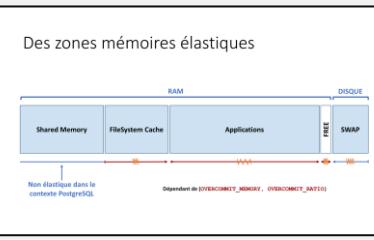
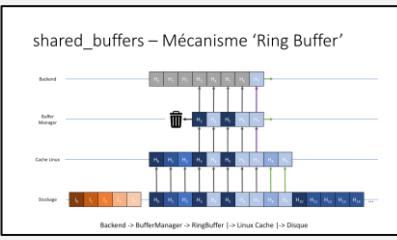
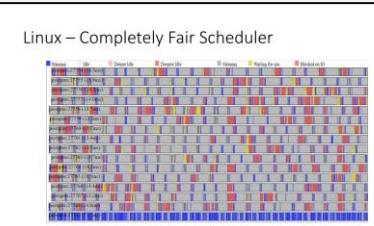


PostgreSQL

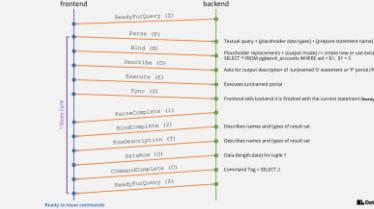
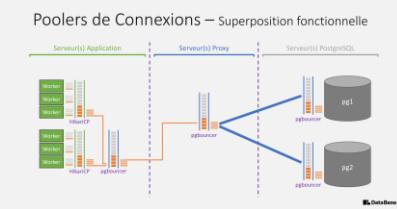
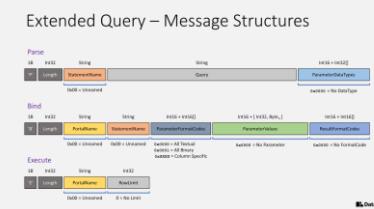
- Since 2006 (PostgreSQL 7.2) in the Telecoms
 - Database Architect (+support +dev)
- Since 2017
 - Delacourt Consulting
 - 2ndQuadrant
 - Data Bene
 - Consultant Expert PostgreSQL

Conferences

PGDay France 2021
PostgreSQL – Les règles de dimensionnement (bien connues).
Frédéric Delacourt – Data Bene
(<https://www.youtube.com/watch?v=unwvewoDmwo>)



PGDay France 2023
Tour d'horizon des Connection Poolers



PGDAY FRANCE 2022
CRÉATION D'UN NOUVEAU TYPE DE DONNÉES POUR UN CHIFFREMENT DES DONNÉES TRANSPARENT

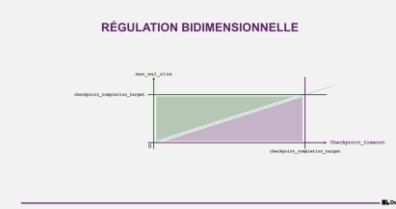
(<https://www.youtube.com/watch?v=tzvOvzPf2M4>)

RÉSULTAT – QUE POUVONS-NOUS FAIRE ?

```
SET enable_seqscan = off
EXPLAIN ANALYZE SELECT * FROM app.test WHERE encrypted = '150';
PostgreSQL: 14
-----+-----+-----+-----+-----+
| QUERY PLAN | | | | |
-----+-----+-----+-----+-----+
| Bitmap Heap Scan on test  (cost=943.81..2108.83 rows=40002 width=10) |
|   (actual time=195.0..198 rows=1 loops=1)
|     Recheck Cond: (encrypted = '150')::text
|       Heap Blocks: exact=1
|         -> Bitmap Index Scan on test_encrypted_idx
|           (rows=5002, actual time=0.184..0.185 rows=1 loops=1)
|             Index Cond: (encrypted = '150')::text
| Planning Time: 0.329 ms
| Execution Time: 0.246 ms
-----+-----+-----+-----+-----+
```

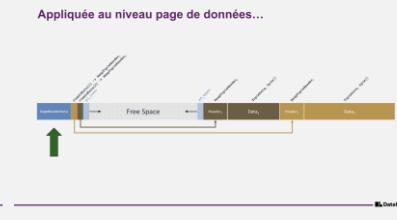
PGDAY FRANCE 2022 - LIGHTNING TALK
MÉCANISME DE RÉGULATION DU CHECKPOINT

(https://www.youtube.com/watch?v=tF_BgYSdWm)



CRÉATION D'UN NOUVEAU TYPE DE DONNÉES

```
CREATE TYPE name
    INPUT = input_function
    OUTPUT = output_function
    INTERNALLENGTH = internal_length
    LENGTH = length_function
    SEMI = semi_function
    TYPEOF = type_modifier_input_function
    TYPEDBY = type_modifier_output_function
    ANALYZE = analyze_function
    HSTORE = hstore_function
    INTERNALLENGTH = internal_length
    INTERNALALIGNMENT = internal_alignment
    INTERNALFORMAT = internal_format
    LIKE = like_type
    CATALOG = catalog
    PRECISION = precision
    DEFVAL = defVal
    COLUMNS = columns
    DELIMITER = delimiter
    COLLATE = collatable
)
```

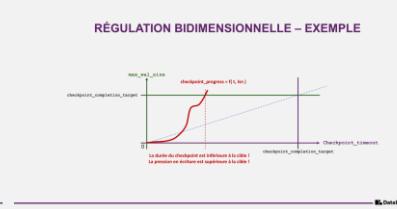


PILOTAGE DE LA RÉGULATION

```
while not processed < new_nb_rows
    wait_for_checkpoint()
    process_rows()

checkpoint_progress = max(processed, min(1, nb_nb_rows))
checkpoint_completion_target = nb_nb_rows * nb_nb_rows / nb_nb_rows

set_checkpoint_progress = correct_lbm / start_lbm * start_lbm + nb_nb_rows * nb_nb_rows / nb_nb_rows
time_checkpoint = (new_start_time - start_time) / nb_nb_rows
time_progress = (new_start_time - start_time) / nb_nb_rows * nb_nb_rows
if checkpoint_progress > new_start_time:
    sleep((checkpoint_completeness - new_start_time) * nb_nb_rows)
```



What about you?

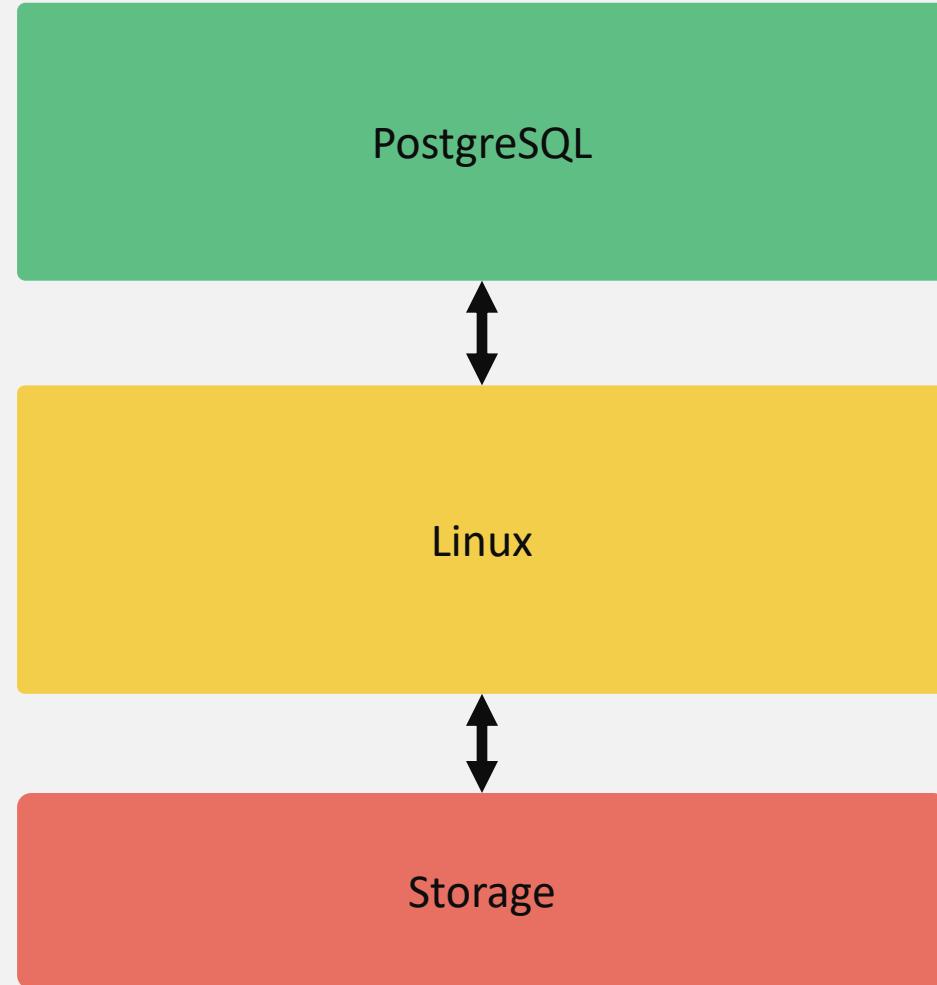
- Who are you?
- How are you related to PostgreSQL?
- What do you expect from this session?

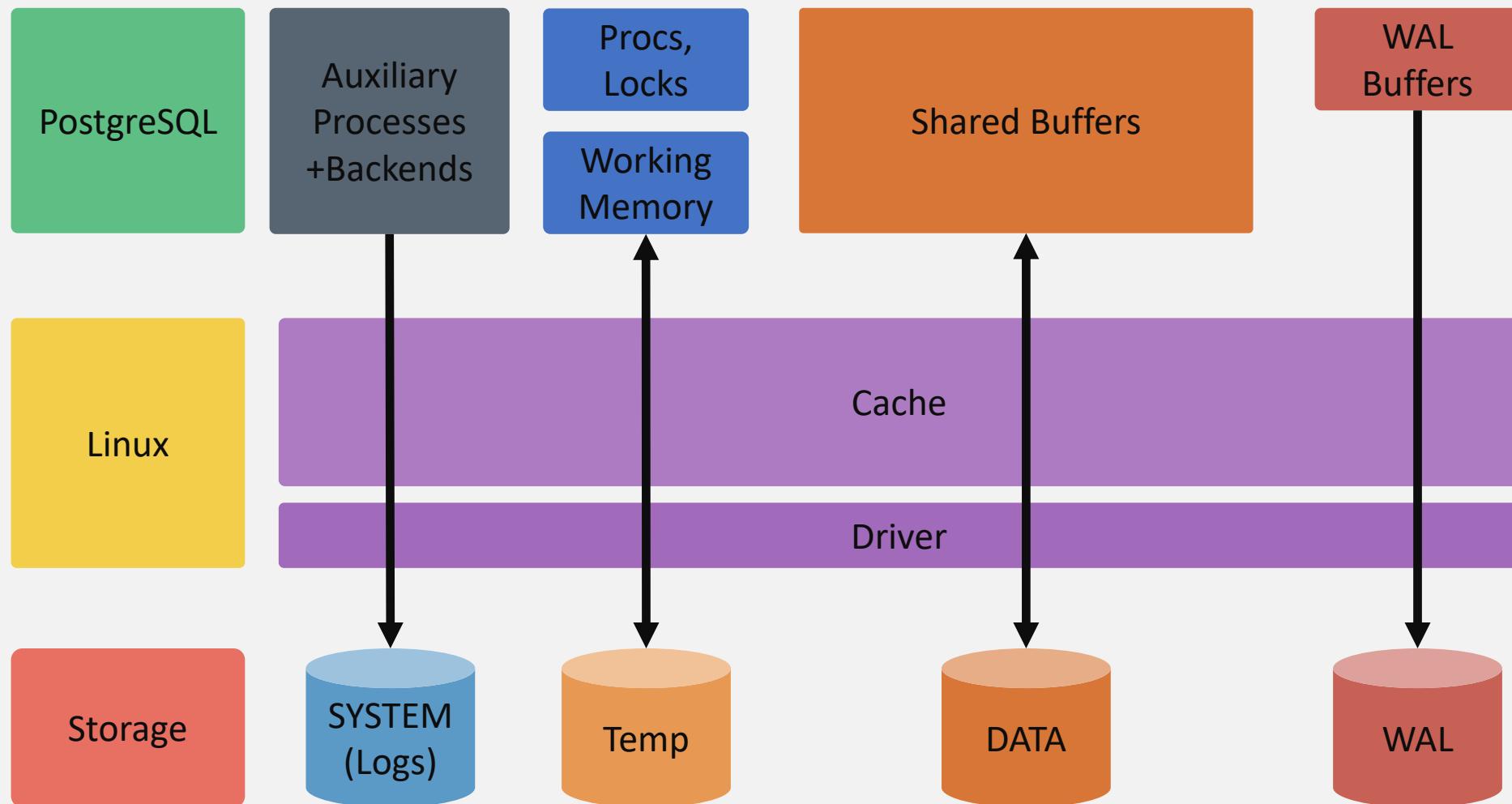
Why collecting metrics?

- Workload Analysis
 - Baseline when the system works fine
 - Allows unexpected pattern detections
- Trends
 - Capacity Planning
- Alerting
 - When things are broken
 - When things are breaking or about to break
- Detail Analyzing
 - Troubleshooting
 - Learning behaviours

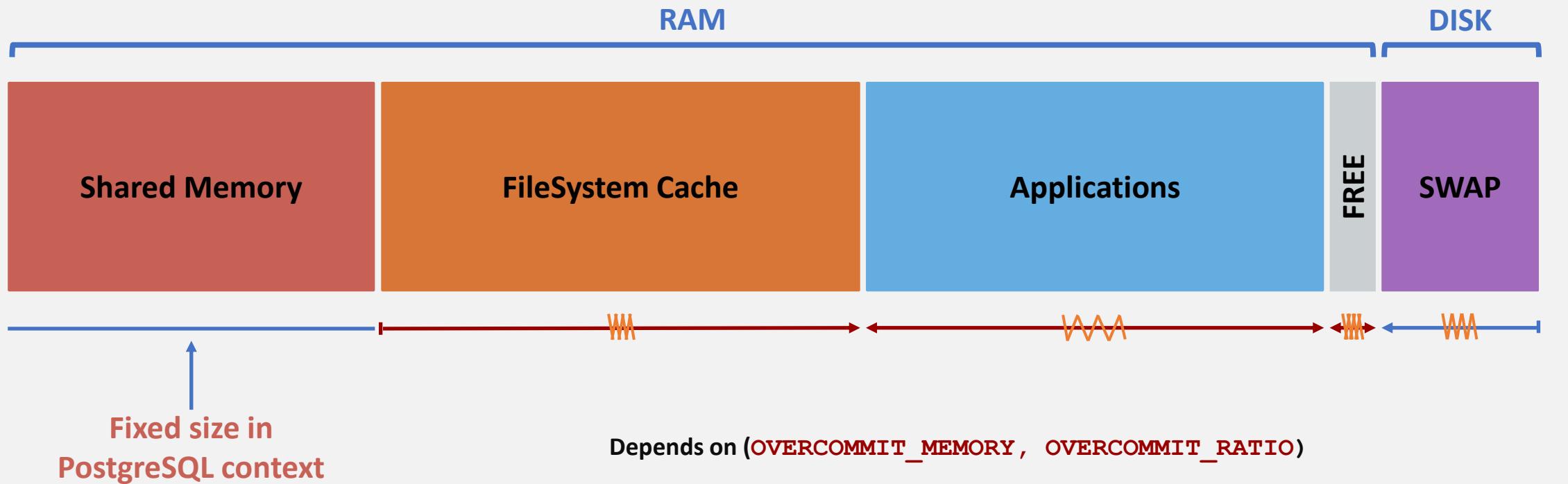
PostgreSQL Diagrams

A Graphic Description of PostgreSQL Internals.



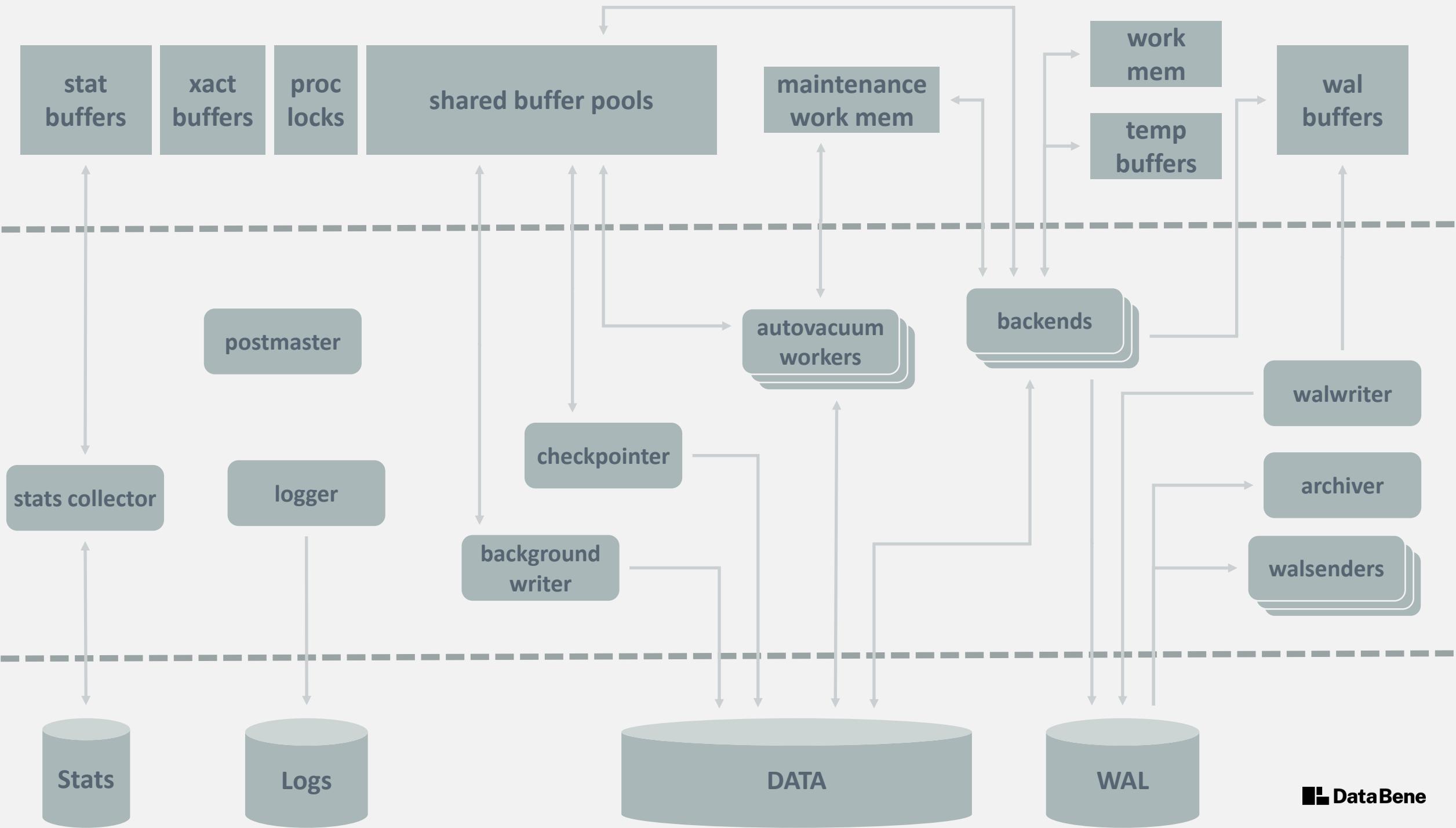


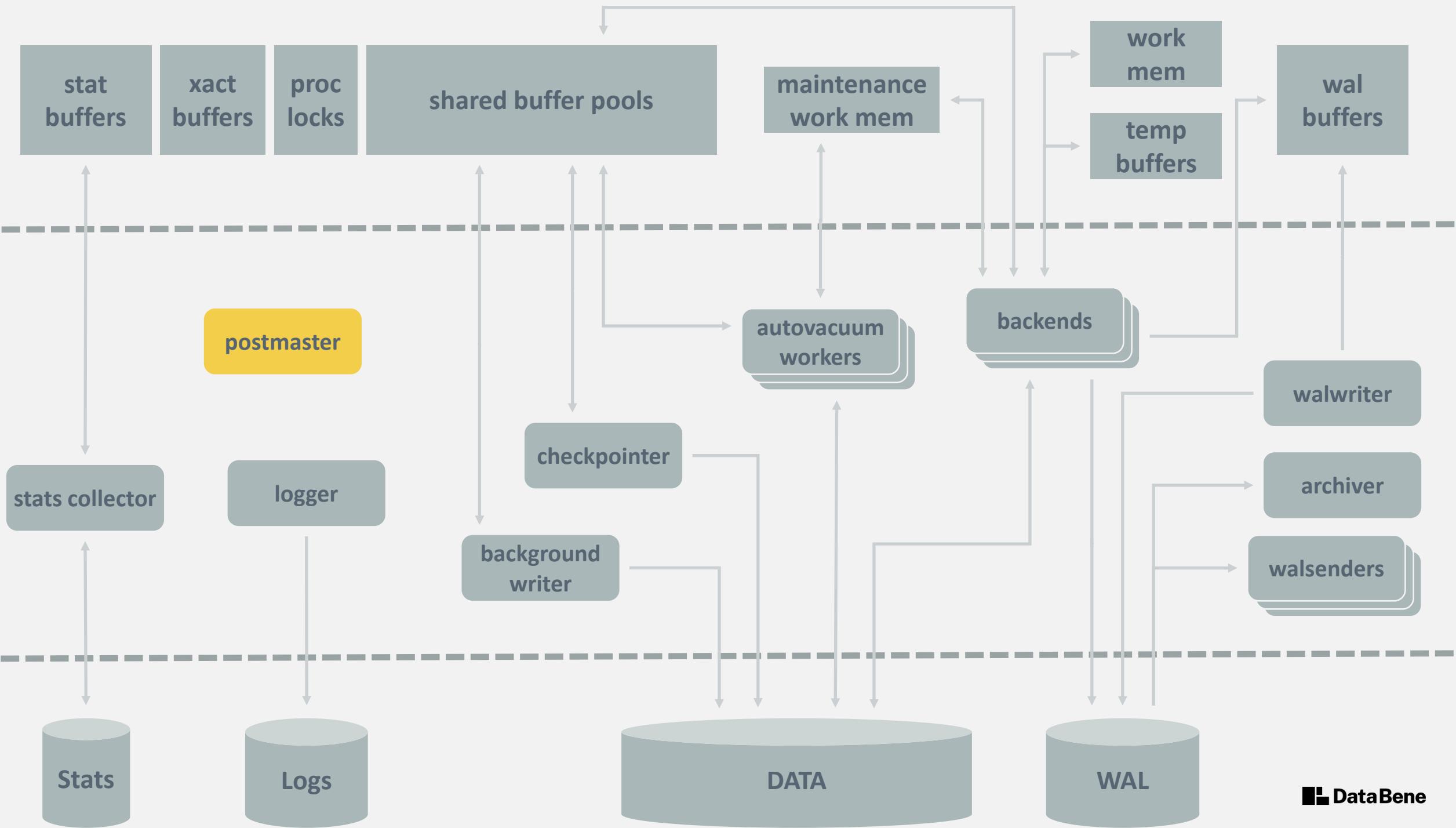
Memory Zones

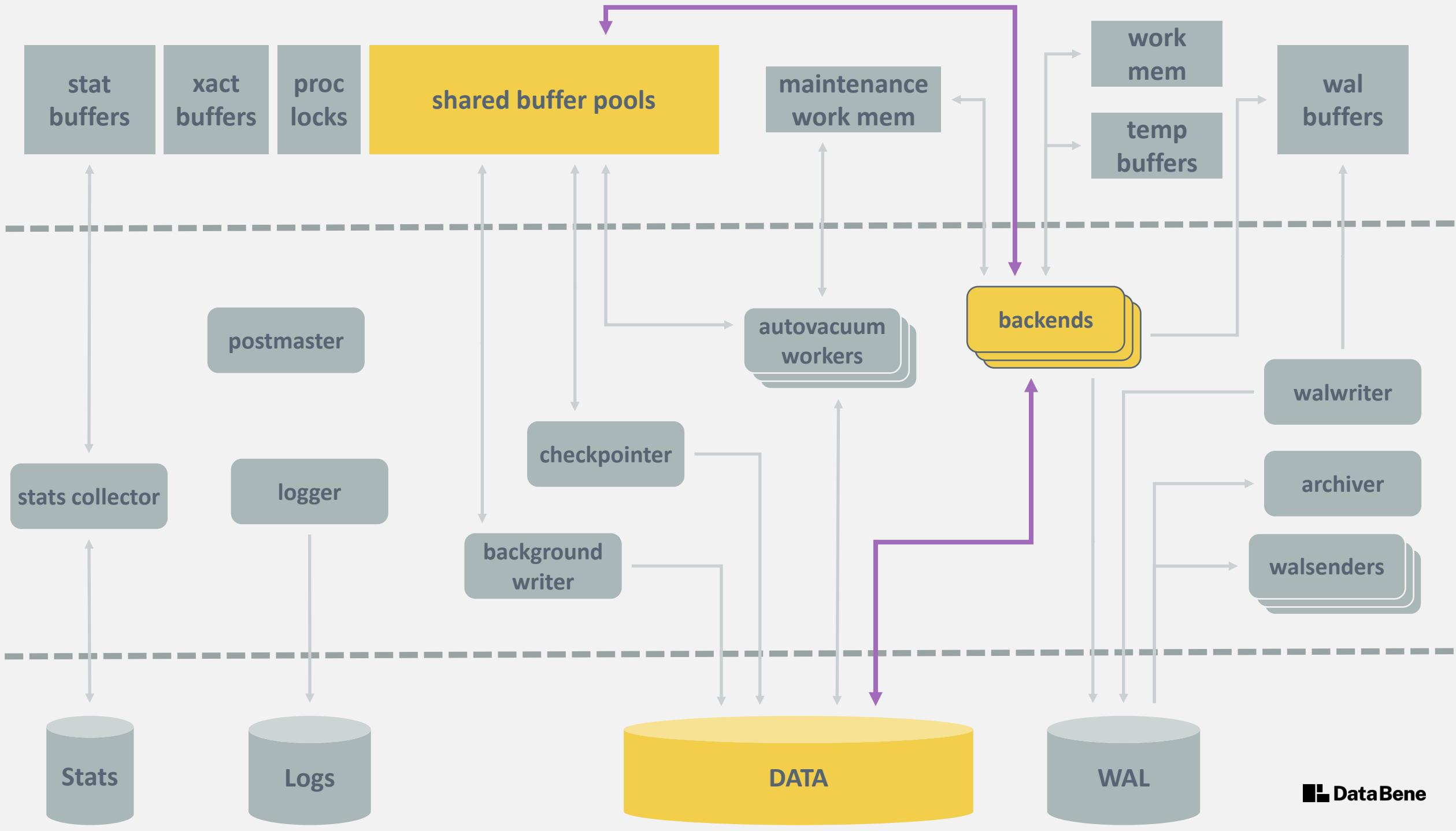


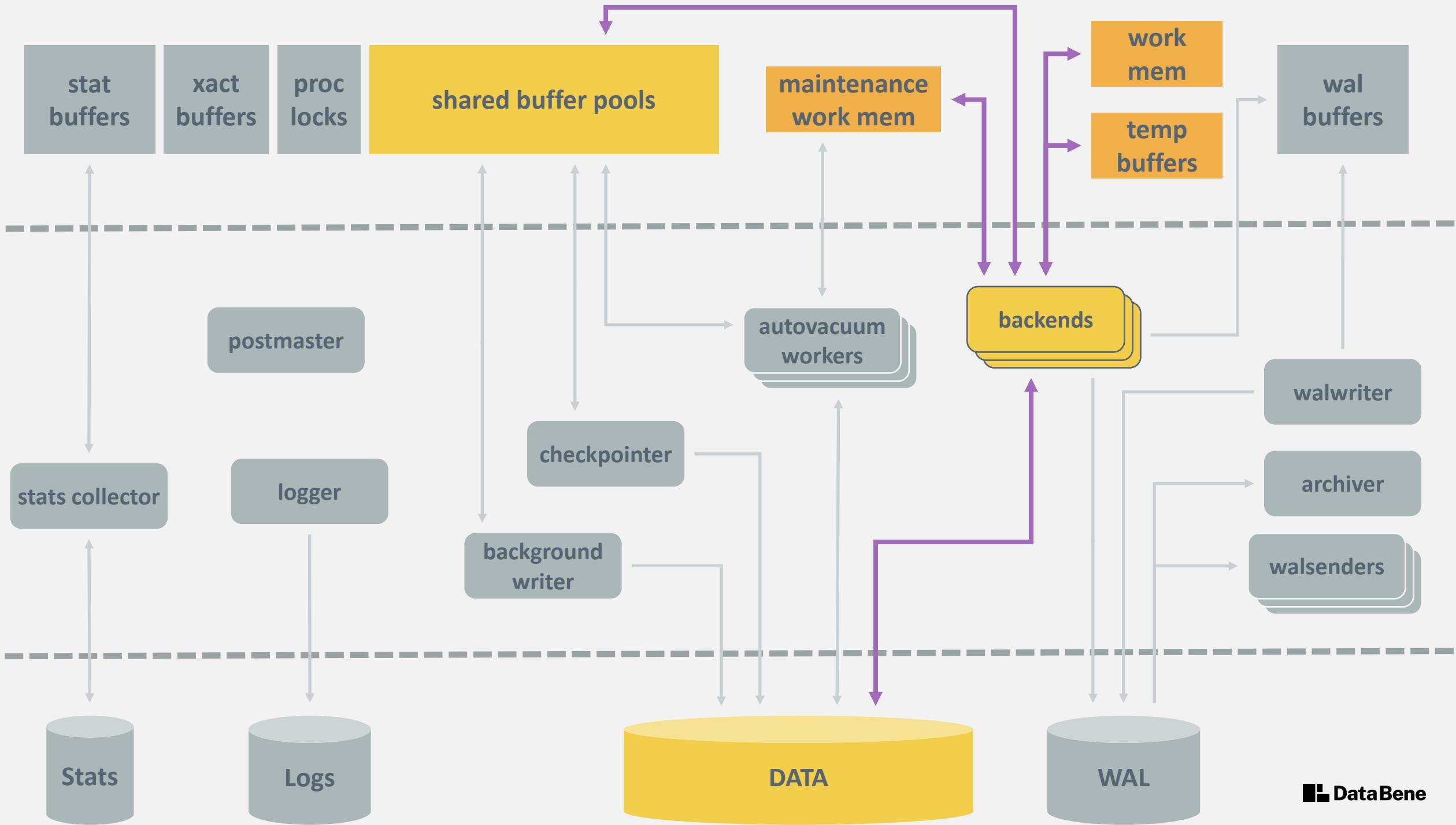
PostgreSQL Architecture

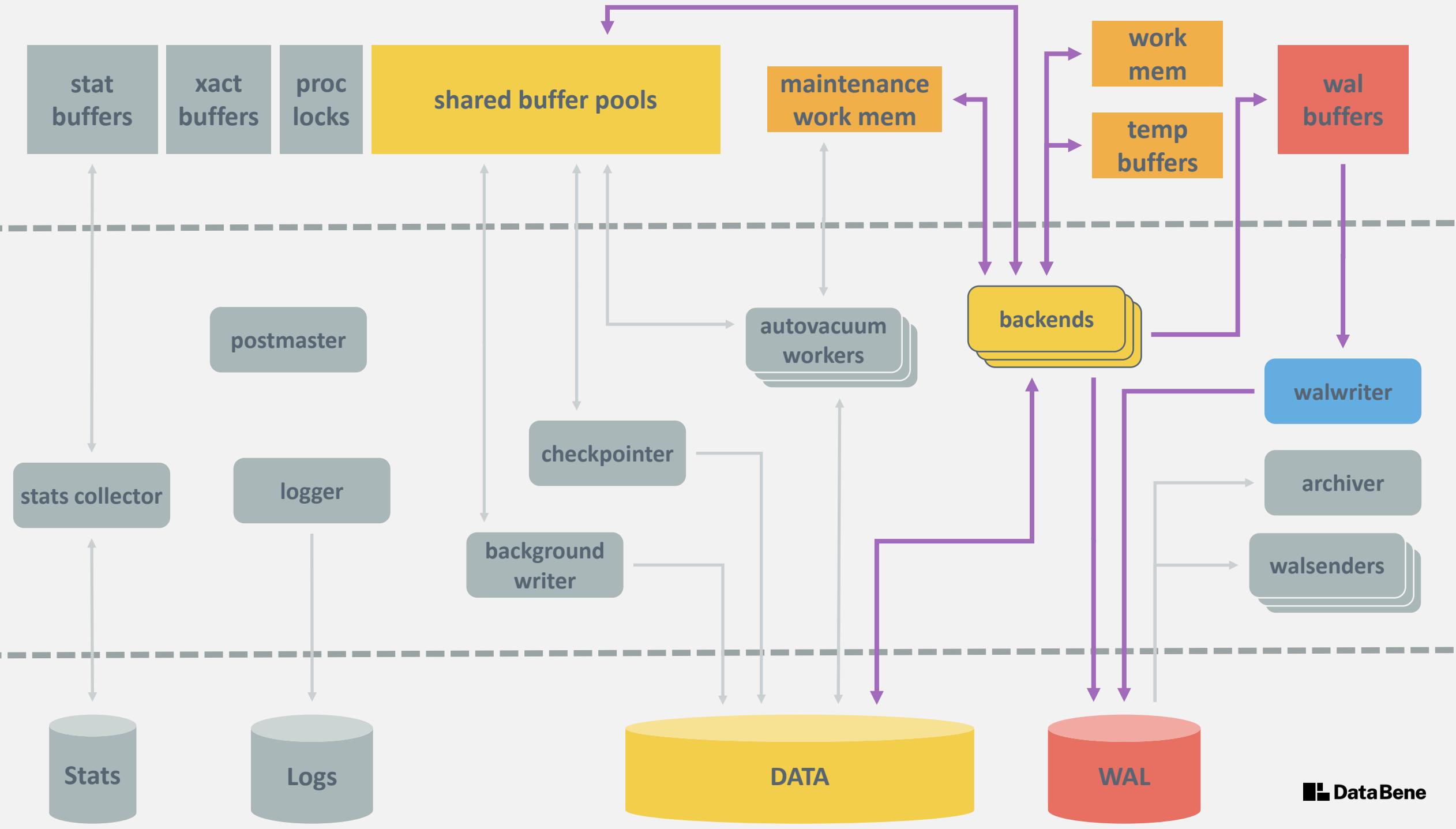
Memory – Process - Storage

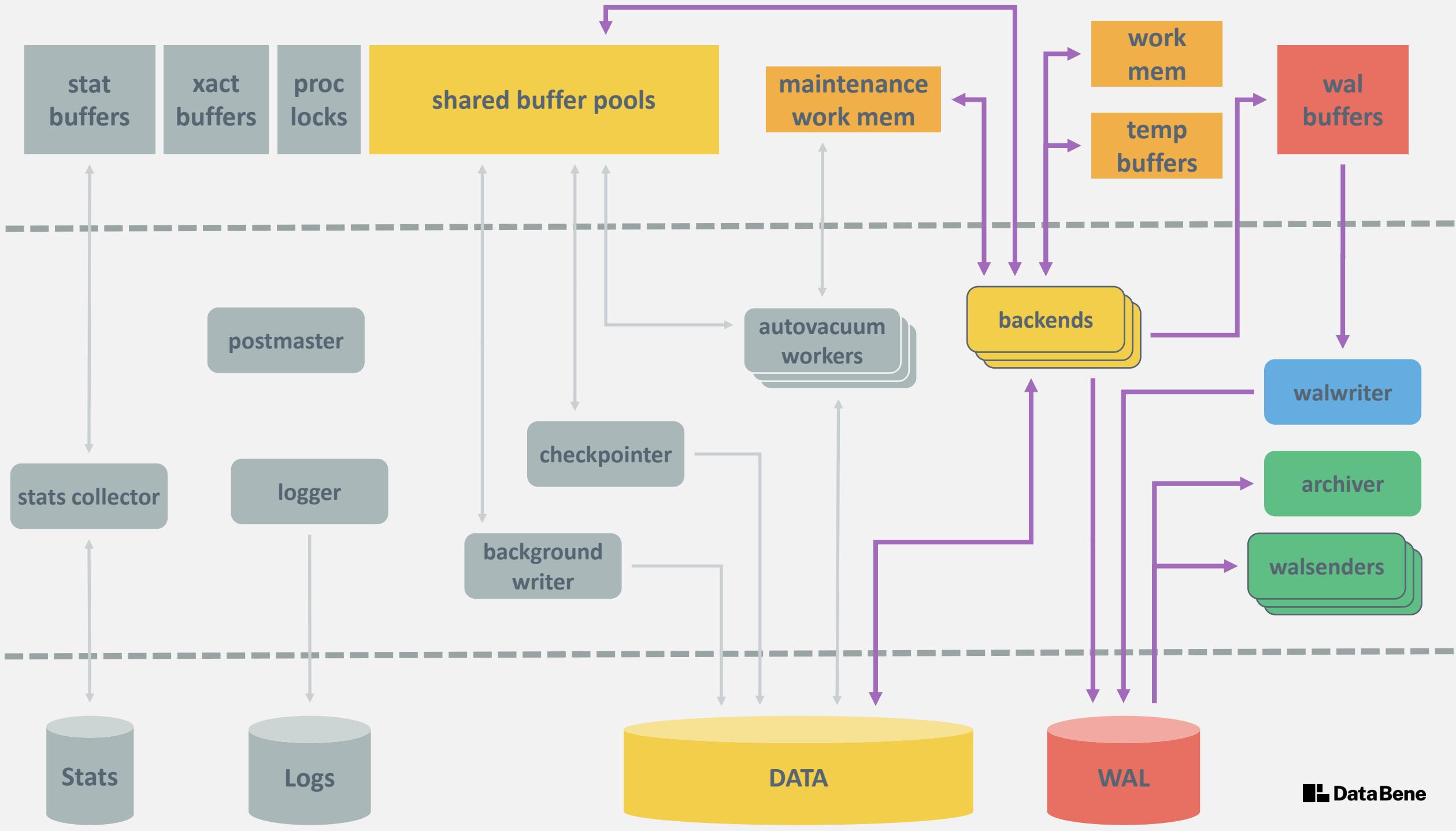


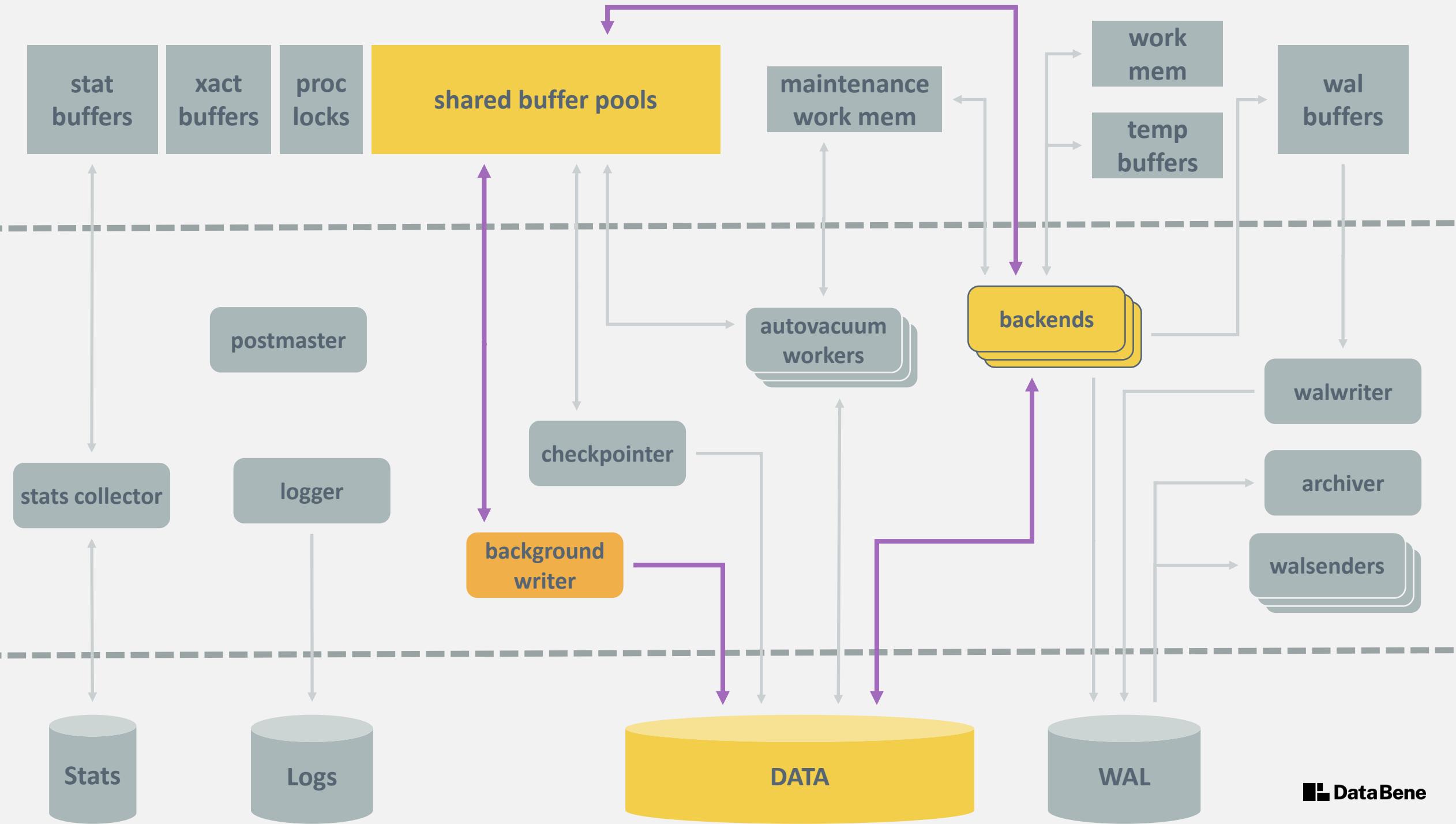


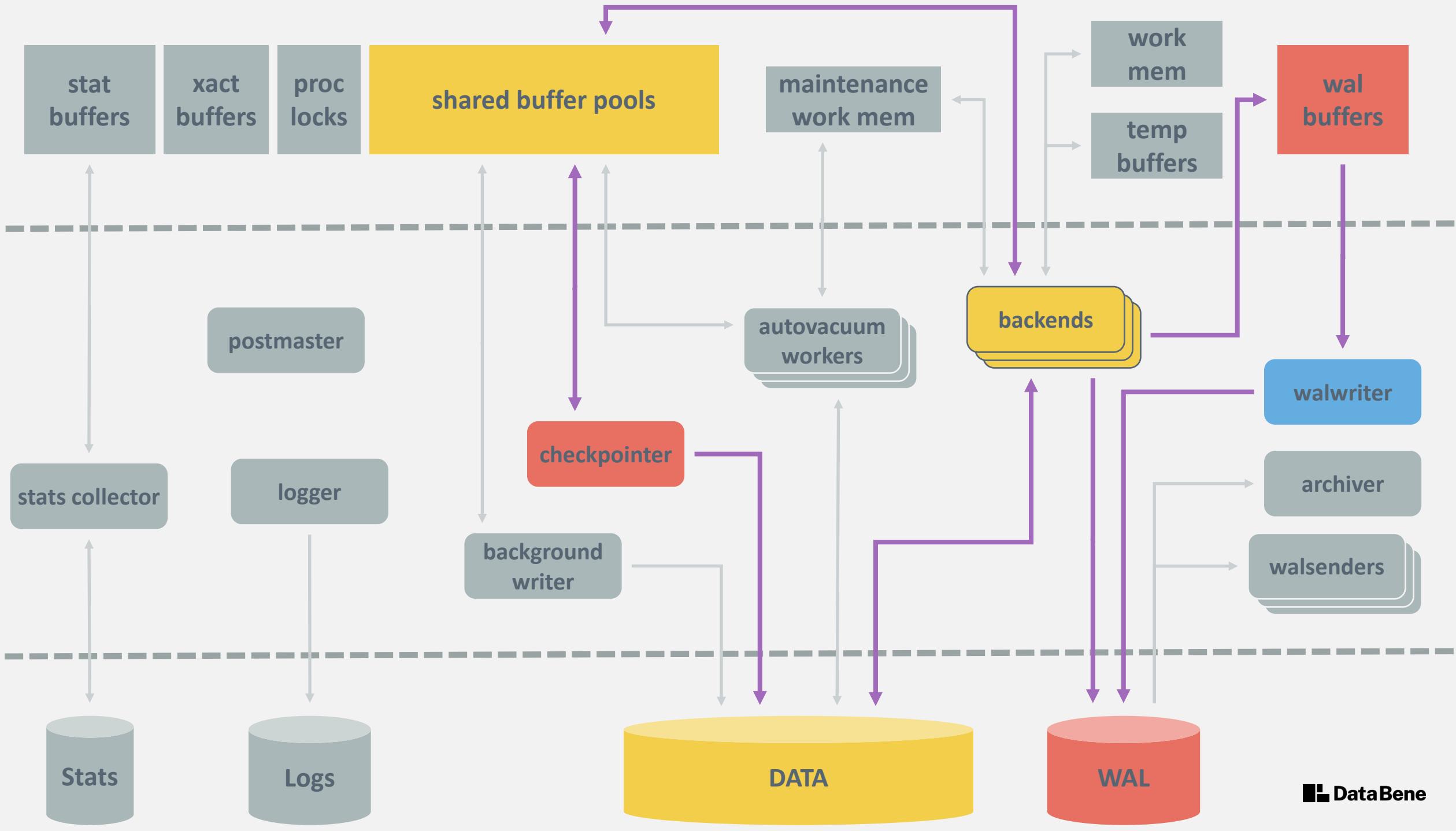


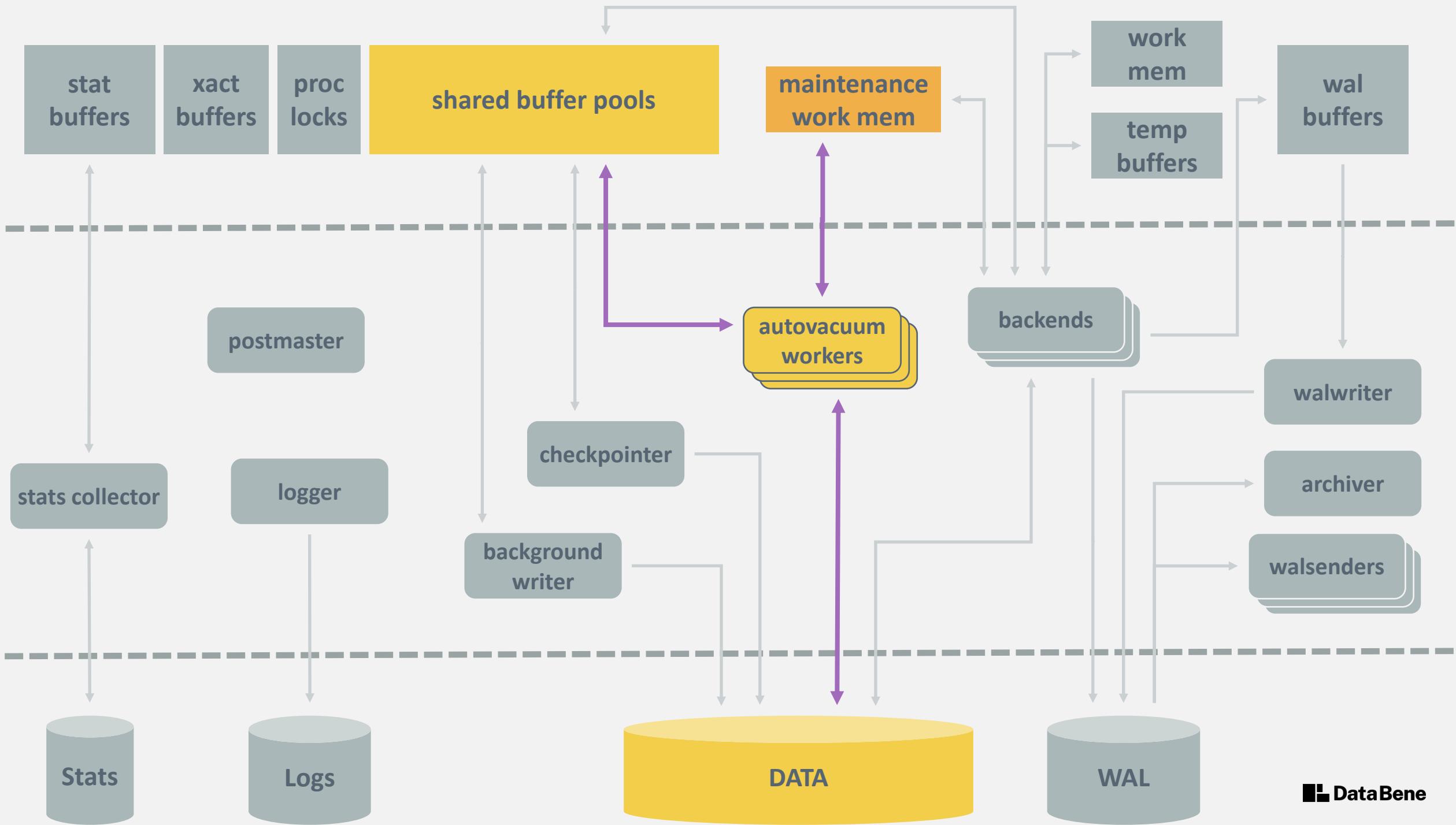


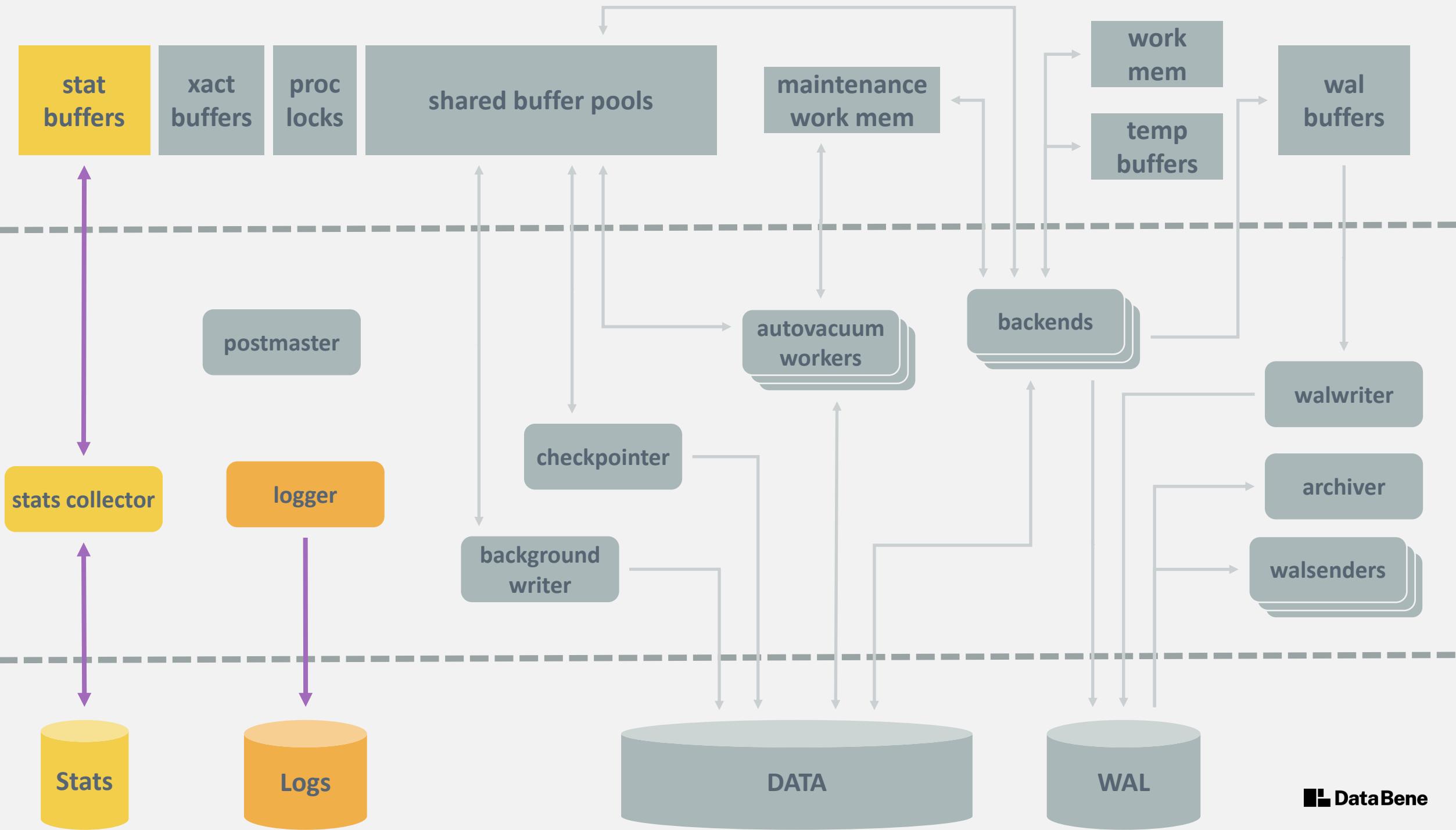


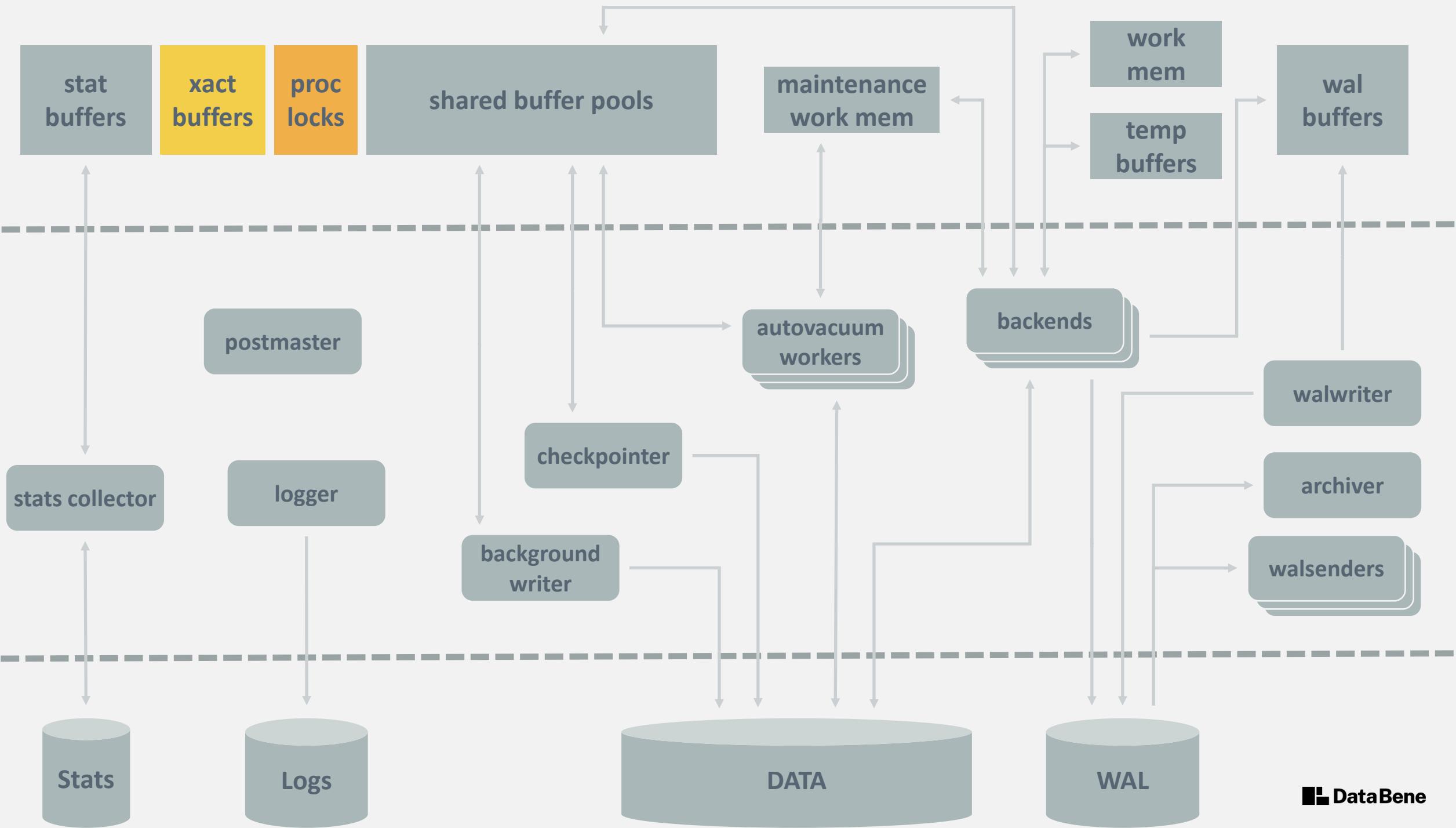






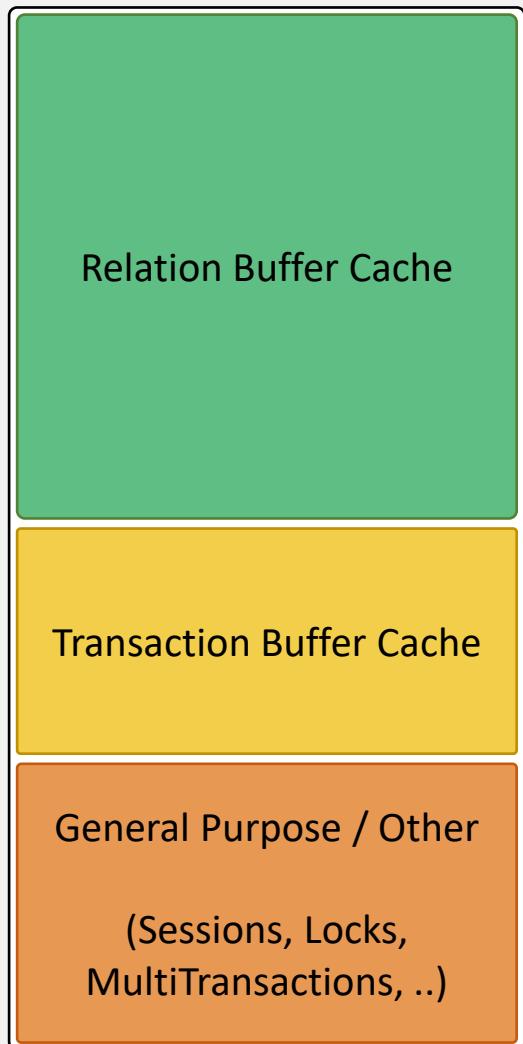




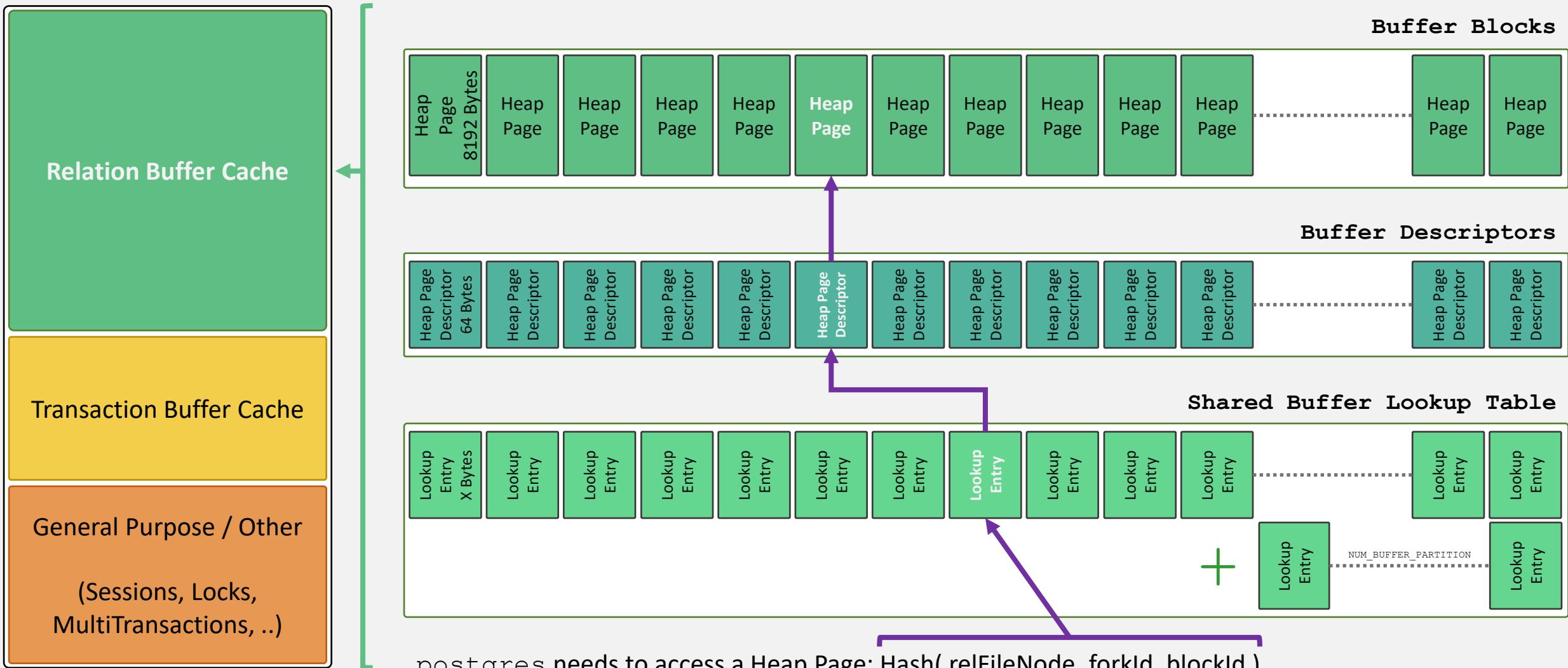


Shared Buffers

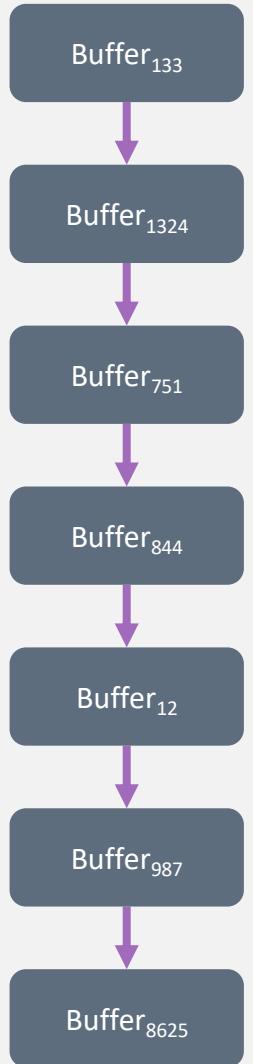
Structure – Usage – Processing



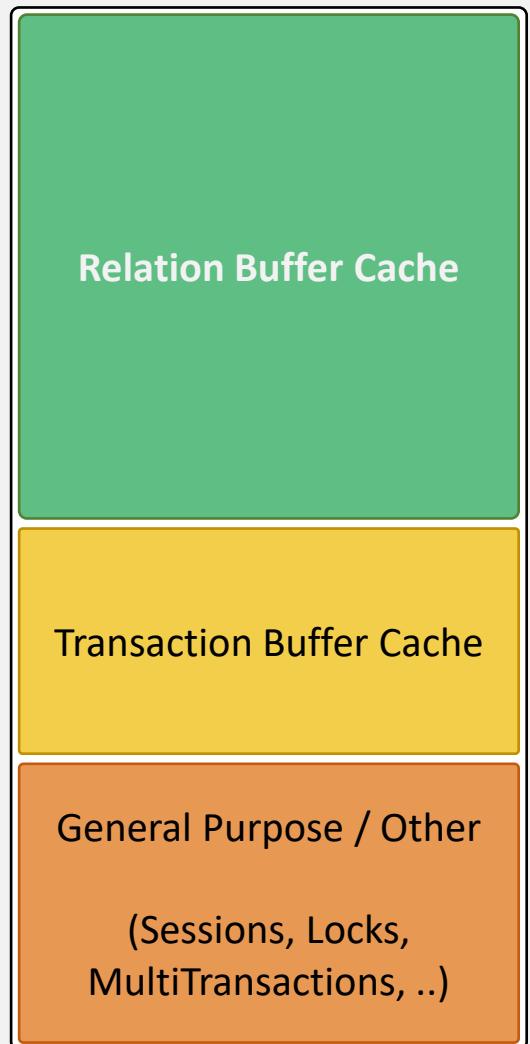
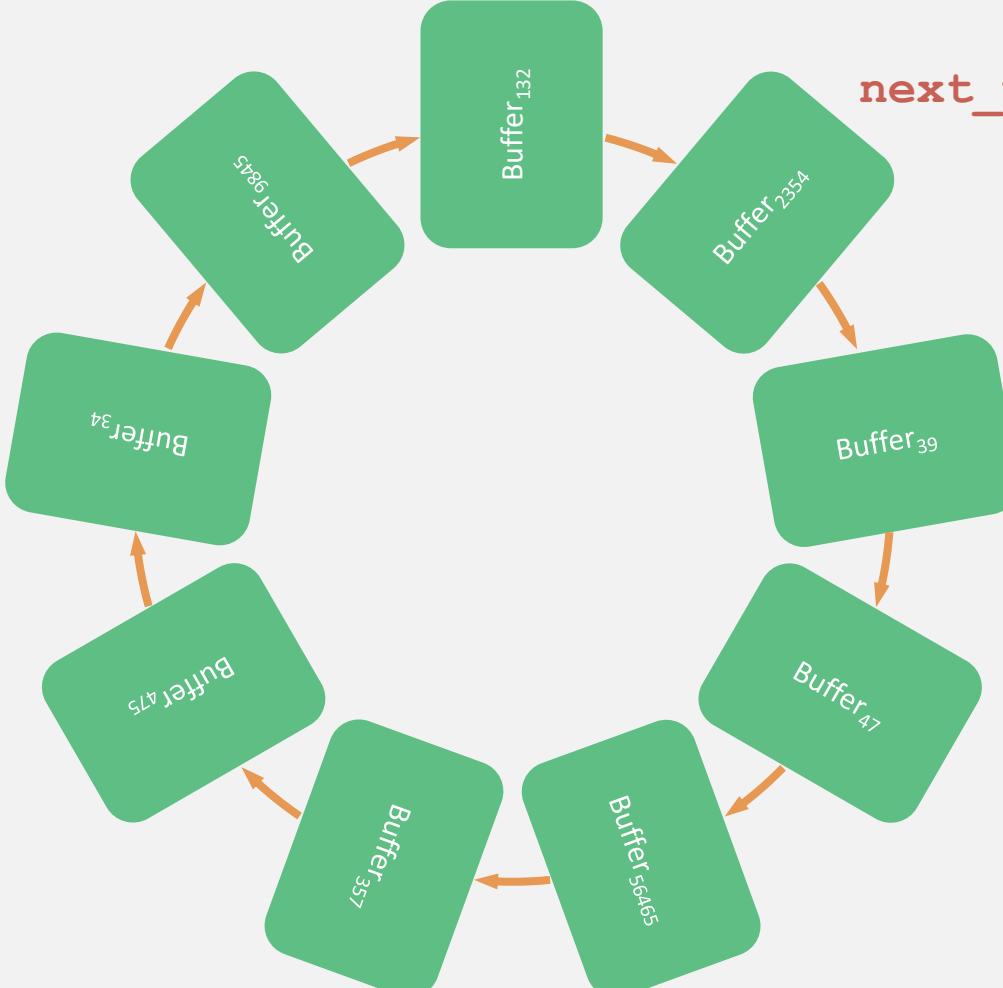
- Stored Objects
 - Tables, Indexes, Materialized Views
 - Free Space Map and Visibility Map.
- Used by Processes
 - `postgres`, `writer`, `checkpointer`.
- Modifications (INSERT, UPDATE, DELETE)
 - Are always done in the cache (`postgres`),
 - Written on Disk later (`checkpointer`, `writer`, `postgres`).

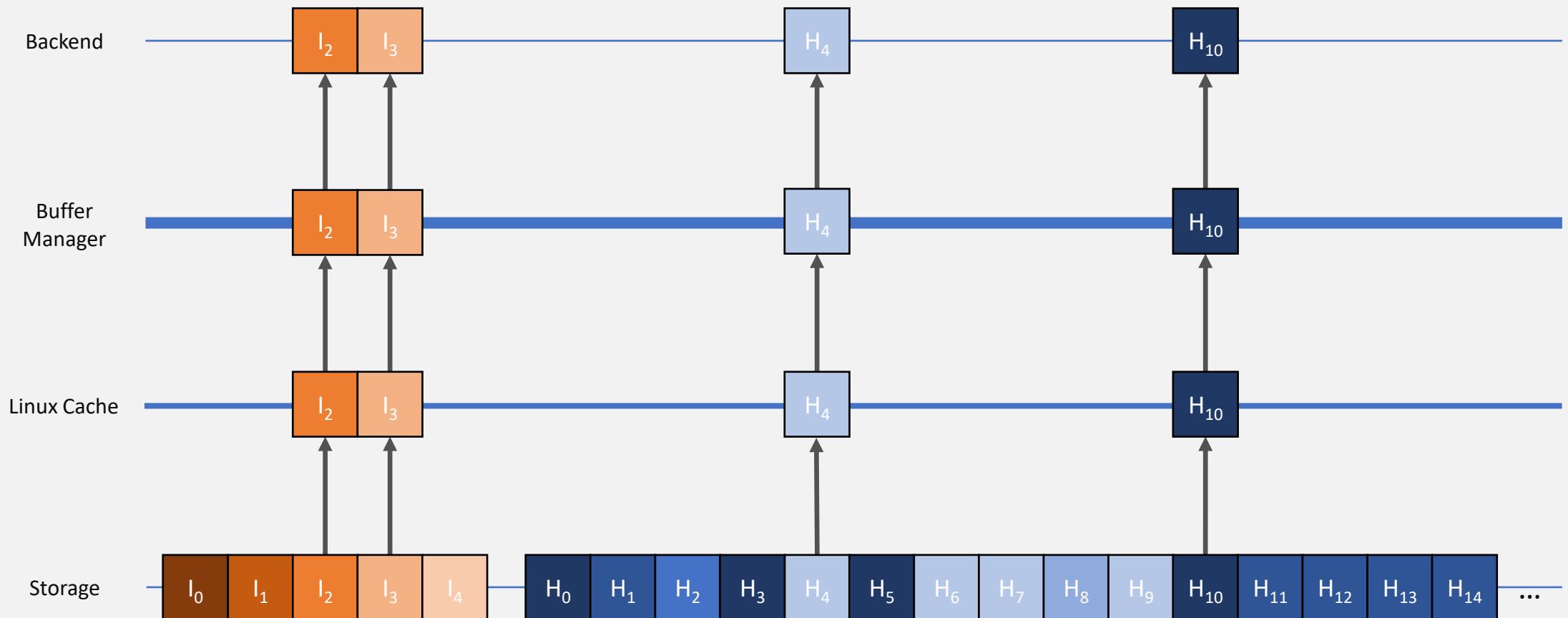


freelist

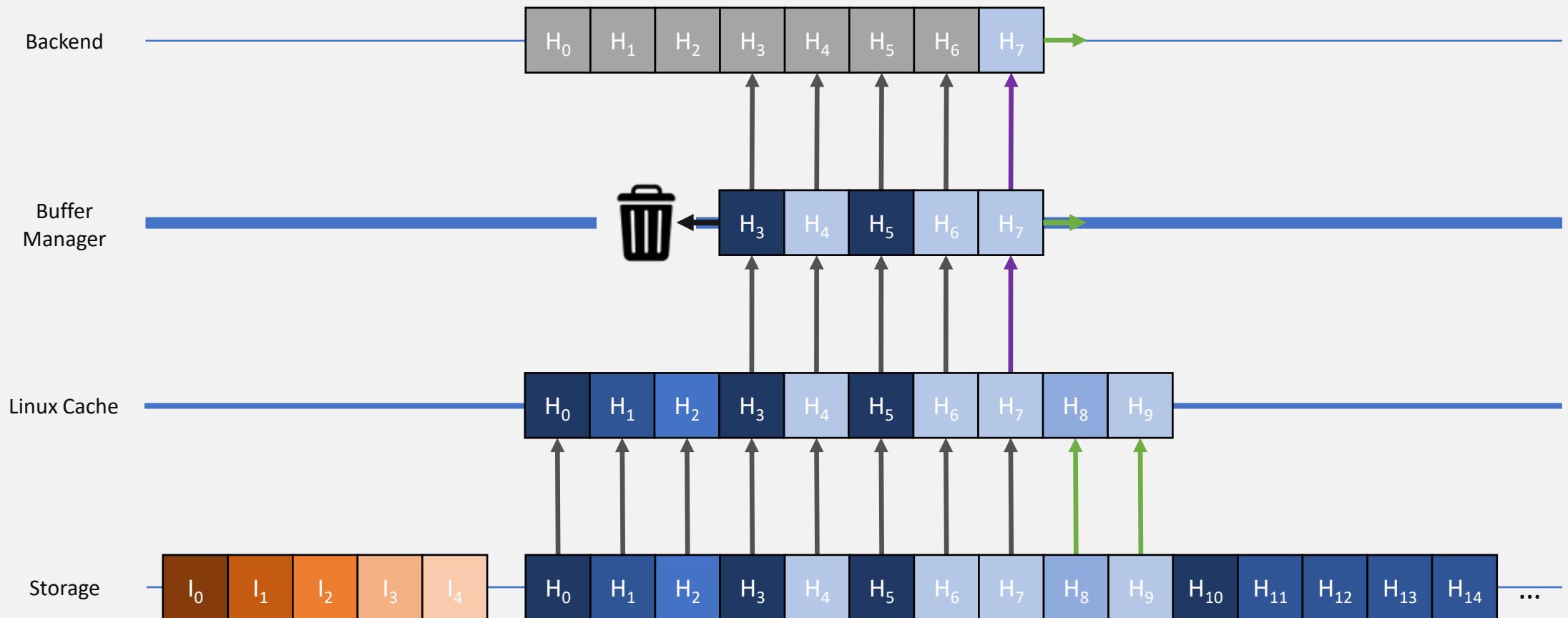


next_victim





Backend -> BufferManager -> BufferPool |-> Linux Cache |-> Storage



Backend -> BufferManager -> RingBuffer |-> Linux Cache |-> Storage

Index Level Statistics

pg_catalog.pg_stat_[all|sys|user]_indexes

Column	Type
relid	oid
indexrelid	oid
schemaname	name
relname	name
indexrelname	name
idx_scan	bigint
last_idx_scan	timestamp with time zone
idx_tup_read	bigint
idx_tup_fetch	bigint

pg_catalog.pg_stat_[all|sys|user]_indexes

Column	Type
relid	oid
indexrelid	oid
schemaname	name
relname	name
indexrelname	name
<u>idx_scan</u>	<u>bigint</u>
<u>last_idx_scan</u>	<u>timestamp with time zone</u>
idx_tup_read	bigint
idx_tup_fetch	bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

$$rate_{idx_scan} = \frac{\Delta idx_scan}{\Delta T}$$

Old event when $last_idx_{scan} \leq now() - k\Delta T; k \in [0.5 ; 3[$

pg_catalog.pg_stat_[all|sys|user]_indexes

Column	Type
relid	oid
indexrelid	oid
schemaname	name
relname	name
indexrelname	name
idx_scan	bigint
last_idx_scan	timestamp with time zone
idx_tup_read	bigint
idx_tup_fetch	bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

- **idx_tup_read** is the number of returned entries by the index scan
- **idx_tup_fetch** is the number of entries which produced reads in heap
- **idx_tup_fetch** \leq **idx_tup_read**

pg_catalog.pg_stat_[all|sys|user]_indexes

Column	Type
relid	oid
indexrelid	oid
schemaname	name
relname	name
indexrelname	name

idx_scan | bigint
last_idx_scan | timestamp with time zone
idx_tup_read | bigint
idx_tup_fetch | bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

$$ratio_{fetch} = \frac{\Delta idx_tup_fetch}{\Delta idx_tup_read} \quad (\text{index_onlyscan})$$

$$\text{avg}_{index_by_scan} = \frac{\Delta idx_tup_read}{\Delta idx_scan} \quad (\text{for watch/baseline})$$

pg_catalog.pg_statio_[all|sys|user]_indexes

Column	Type
relid	oid
indexrelid	oid
schemaname	name
relname	name
indexrelname	name
idx_blkss_read	bigrnt
idx_blkss_hit	bigrnt

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

$$\text{cache}_{\text{idx_hit}} = \frac{\Delta idx_blkss_hit}{\Delta idx_blkss_hit + \Delta idx_blkss_read} \geq 0.98$$

$$\text{cache}_{\text{idx_miss}} = \frac{\Delta idx_blkss_read}{\Delta idx_blkss_hit + \Delta idx_blkss_read}$$

Table Level Statistics

pg_catalog.pg_stat_[all|sys|user]_tables

Column	Type	Column	Type
relid	oid	n_live_tup	bigint
schemaname	name	n_dead_tup	bigint
relname	name	n_mod_since_analyze	bigint
seq_scan	bigint	n_ins_since_vacuum	bigint
last_seq_scan	timestamp with time zone	last_vacuum	timestamp with time zone
seq_tup_read	bigint	last_autovacuum	timestamp with time zone
idx_scan	bigint	last_analyze	timestamp with time zone
last_idx_scan	timestamp with time zone	last_autoanalyze	timestamp with time zone
idx_tup_fetch	bigint	vacuum_count	bigint
n_tup_ins	bigint	autovacuum_count	bigint
n_tup_upd	bigint	analyze_count	bigint
n_tup_del	bigint	autoanalyze_count	bigint
n_tup_hot_upd	bigint		
n_tup_newpage_upd	bigint		

pg_catalog.pg_stat_[all|sys|user]_tables

Column	Type	
relid	oid	$rate_x = \frac{\Delta x}{\Delta T}$; $\Delta T \in (1 - 15 \text{ mins})$
schemaname	name	
relname	name	
seq_scan	bigrnt	
last_seq_scan	timestamp with time zone	$\text{avg tuples by scan} = \frac{\Delta \text{seq_tup_read}}{\Delta \text{seq_scan}}$ (for watch/baseline)
seq_tup_read	bigrnt	Old event when $\text{last_*_scan} \leq \text{now}() - k\Delta T$; $k \in [0.5 ; 3[$
idx_scan	bigrnt	
last_idx_scan	timestamp with time zone	
idx_tup_fetch	bigrnt	

pg_catalog.pg_stat_[all|sys|user]_tables

Column	Type
relid	oid
schemaname	name
relname	name
seq_scan	bigint
last_seq_scan	timestamp with time zone
seq_tup_read	bigint
idx_scan	bigint
last_idx_scan	timestamp with time zone
idx_tup_fetch	bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

$$\text{avg}_{\text{tuples by scan}} = \frac{\Delta \text{idx_tup_fetch}}{\Delta \text{idx_scan}} \text{ (for watch/baseline)}$$

Old event when $\text{last_*_scan} \leq \text{now}() - k\Delta T$; $k \in [0.5 ; 3[$

pg_catalog.pg_stat_[all|sys|user]_tables

Column	Type
relid	oid
schemaname	name
relname	name
n_tup_ins	bigint
n_tup_upd	bigint
n_tup_del	bigint
n_tup_hot_upd	bigint
n_tup_newpage_upd	bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

$$ratio_{hot} = \frac{\Delta n_tup_hot_upd}{\Delta n_tup_upd}$$

$$ratio_{newpage} = \frac{\Delta n_tup_newpage_upd}{\Delta n_tup_upd}$$

$$ratio_{rw} = \frac{\Delta seq_tup_read + \Delta idx_tup_fetch}{\Delta n_tup_ins + \Delta n_tup_upd + \Delta n_tup_del}$$

pg_catalog.pg_stat_[all|sys|user]_tables

Column	Type
relid	oid
schemaname	name
relname	name

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

<u>n_mod_since_analyze</u>	bigint
<u>last_analyze</u>	timestamp with time zone
<u>last_autoanalyze</u>	timestamp with time zone
analyze_count	bigint
autoanalyze_count	bigint

Make sure table received (auto)ANALYZE

$$n_mod_since_analyze \geq analyze_threshold + analyze_scale_factor * pg_class.reltuples$$

pg_catalog.pg_stat_[all|sys|user]_tables

Column	Type
relid	oid
schemaname	name
relname	name
n_live_tup	bigint
n_dead_tup	bigint
n_ins_since_vacuum	bigint
!last_vacuum	timestamp with time zone
!last_autovacuum	timestamp with time zone
vacuum_count	bigint
autovacuum_count	bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

Make sure table received (auto)VACUUM

*n_ins_since_vacuum ≥ vacuum_insert_threshold + vacuum_insert_scale_factor * pg_class.reltuples*
*n_dead_tup ≥ vacuum_threshold + vacuum_scale_factor * pg_class.reltuples*

pg_catalog.pg_statio_[all|sys|user]_tables

Column	Type
relid	oid
schemaname	name
relname	name
<code>heap_blks_read</code>	<code>bigint</code>
<code>heap_blks_hit</code>	<code>bigint</code>
<code>idx_blks_read</code>	<code>bigint</code>
<code>idx_blks_hit</code>	<code>bigint</code>
<code>toast_blks_read</code>	<code>bigint</code>
<code>toast_blks_hit</code>	<code>bigint</code>
<code>tidx_blks_read</code>	<code>bigint</code>
<code>tidx_blks_hit</code>	<code>bigint</code>

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 15 \text{ mins})$$

$$cache_{hit} = \frac{\Delta x_{hit}}{\Delta x_{hit} + \Delta x_{read}}$$

$$cache_{miss} = \frac{\Delta x_{read}}{\Delta x_{hit} + \Delta x_{read}}$$

Database Level Statistics

pg_catalog.pg_stat_database

Column	Type	Column	Type
datid	oid	deadlocks	bigint
datname	name	checksum_failures	bigint
numbackends	integer	checksum_last_failure	timestamp with time zone
xact_commit	bigint	blk_read_time	double precision
xact_rollback	bigint	blk_write_time	double precision
blks_read	bigint	session_time	double precision
blks_hit	bigint	active_time	double precision
tup_returned	bigint	idle_in_transaction_time	double precision
tup_fetched	bigint	sessions	bigint
tup_inserted	bigint	sessions_abandoned	bigint
tup_updated	bigint	sessions_fatal	bigint
tup_deleted	bigint	sessions_killed	bigint
conflicts	bigint	stats_reset	timestamp with time zone
temp_files	bigint		
temp_bytes	bigint		

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
numbackends	integer

The number of currently connected backends to this database!

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
stats_reset	timestamp with time zone



The cumulated statistic counters started accumulating since **stats_reset**.

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
xact_commit	bigrnt
xact_rollback	bigrnt

Number of COMMITTED transactions since `stats_reset`.
SELECT are accounted!

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
xact_commit	bigrnt
xact_rollback	bigrnt

$$tpS_{committed} = \frac{\Delta xact_commit}{\Delta T}$$

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
xact_commit	bignum
xact_rollback	bignum

$$tpS_{committed} = \frac{\Delta xact_commit}{\Delta T}$$

- Should we alert when $tpS_{committed}$ is outside upper and lower boundaries?
- Absolute boundaries
 - $\leq 0.5x$ or $\geq 2x$ window averaged
 - Compared to expected timely activity

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
xact_commit	bignum
xact_rollback	bignum

$$tpS_{committed} = \frac{\Delta xact_commit}{\Delta T}$$

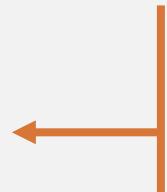
Should we up/down-scale the system?

pg_catalog.pg_stat_database

Column	Type	
datid	oid	
datname	name	
xact_commit	bigint	
xact_rollback	bigint	Should be kept as low as possible!

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
xact_commit	bigint
xact_rollback	bigint



Too much rollbacks means:

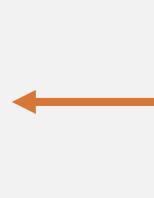
- concurrency issues (or bad design)
- producing dead tuples!

pg_catalog.pg_stat_database

Column	Type	
datid	oid	
datname	name	
xact_commit	bigint	
xact_rollback	bigint	 $tps_{rolled} = \frac{\Delta xact_rollback}{\Delta T}$

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
xact_commit	bigrnt
xact_rollback	bigrnt



$$tps_{committed} = \frac{\Delta xact_commit}{\Delta T}$$

$$tps_{rollacked} = \frac{\Delta xact_rollback}{\Delta T}$$

$$\text{ratio}_{rollacked} = \frac{\Delta xact_rollback}{\Delta xact_committed}$$

Should we alert when
 $\text{ratio}_{rollacked}$ is higher than a threshold?

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
xact_commit	bigrnt
xact_rollback	bigrnt

$tpS_{committed} = \frac{\Delta xact_commit}{\Delta T}$

$tpS_{rollacked} = \frac{\Delta xact_rollback}{\Delta T}$

$ratio_{rollacked} = \frac{\Delta xact_rollback}{\Delta xact_committed}$

You just degraded performance and power efficiency by $ratio_{rollacked}$ of your PostgreSQL server over the last observed period.

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
tup_returned	bigrnt
tup_fetched	bigrnt
tup_inserted	bigrnt
tup_updated	bigrnt
tup_deleted	bigrnt

tup_returned ← Number of visible tuples manipulated in SELECT, UPDATE, DELETE, VACUUM... excluding INSERT

pg_catalog.pg_stat_database

Column	Type	
datid	oid	
datname	name	
tup_returned	bigint	
tup_fetched	bigint	
tup_inserted	bigint	
tup_updated	bigint	
tup_deleted	bigint	

Number of tuples fetched
from Index Entries [to Heap Tuples]

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
tup_returned	bigrnt
tup_fetched	bigrnt
tup_inserted	bigrnt
tup_updated	bigrnt
tup_deleted	bigrnt

$$\text{index_scan_ratio} = \frac{\Delta \text{tup_fetched}}{\Delta \text{tup_returned}}$$

pg_catalog.pg_stat_database

Column	Type	
datid	oid	
datname	name	
tup_returned	bigint	
tup_fetched	bigint	
tup_inserted	bigint	
tup_updated	bigint	
tup_deleted	bigint	

Data Manipulation Language

- INSERT
- UPDATE
- DELETE

Produce dead tuples (**no tracking at db level**)

- INSERT + ROLLBACK | UPDATE + ROLLBACK
- UPDATE + COMMIT | DELETE + COMMIT

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
tup_returned	bigint
tup_fetched	bigint
tup_inserted	bigint
tup_updated	bigint
tup_deleted	bigint

$$rate_x = \frac{\Delta x}{\Delta T}$$

$$rate_{DML} = \frac{\Delta inserted + \Delta updated + \Delta deleted}{\Delta T}$$

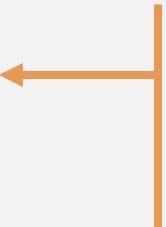
$$ratio_{rw} = \frac{rate_{tup_returned}}{rate_{DML}}$$

Do you find it useful?

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
blks_read	bigint
blks_hit	bigint
temp_files	bigint
temp_bytes	bigint
blk_read_time	double precision
blk_write_time	double precision

Missed data pages from Shared Buffers
PostgreSQL Cache Miss



pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
blks_read	bigint
blks_hit	bigint
temp_files	bigint
temp_bytes	bigint
blk_read_time	double precision
blk_write_time	double precision

Accessed data pages from Shared Buffers
PostgreSQL Cache Hit

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
blks_read	bigrnt
blks_hit	bigrnt
temp_files	bigrnt
temp_bytes	bigrnt
blk_read_time	double precision
blk_write_time	double precision

$$rate_x = \frac{\Delta x}{\Delta T}$$

$$cache_{hit} = \frac{\Delta blks_hit}{\Delta blk_hit + \Delta blks_read}$$

$$cache_{miss} = \frac{\Delta blks_read}{\Delta blk_hit + \Delta blks_read}$$

- Usually $cache_{hit} < 0.7$ is correlated with performance issues.
- Low $cache_{hit}$ ratio is expected when workload consists in accessing uniform distributed tuples over large tables.

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
blks_read	bigint
blks_hit	bigint
temp_files	bigint
temp_bytes	bigint
blk_read_time	double precision
blk_write_time	double precision

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (5, 15, 30, 60 \text{ mins}, 1 \text{ day})$$

$$avgfilesize = \frac{\Delta temp_bytes}{\Delta temp_files} \quad <- \text{not that useful}$$

More information in log files (pid<->temp_files)

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
blks_read	bigint
blks_hit	bigint
temp_files	bigint
temp_bytes	bigint
blk_read_time	double precision
blk_write_time	double precision

$$rate_x = \frac{\Delta x}{\Delta T} ; \text{ for which } \Delta T ?$$

To be compared with what ?
Useful or useless ?

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
deadlocks	bigint
conflicts	bigint
checksum_failures	bigint
checksum_last_failure	timestamp with time zone

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (60 \text{ mins}, 1 \text{ day})?$$

Catches bugs not performance issues
Use log content!

Alert on detection!

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
deadlocks	bigint
conflicts	bigint
checksum_failures	bigint
checksum_last_failure	timestamp with time zone

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (30s, 1 - 60 \text{ mins}, 1 \text{ day})?$$

Number of cancelled queries due to conflicts with recovery (standby only)

pg_catalog.pg_stat_database_conflicts

Column	Type
datid	oid
datname	name
confl_tablespace	bigint
confl_lock	bigint
confl_snapshot	bigint
confl_bufferpin	bigint
confl_deadlock	bigint
confl_active_logicals_slot	bigint

Only on Standby

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (30s, 1 - 60 \text{ mins}, 1 \text{ day})?$$

Causes of cancelled queries due to conflicts with recovery

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
deadlocks	bigint
conflicts	bigint
checksum_failures	bigint
checksum_last_failure	timestamp with time zone

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 - 60 \text{ mins})$$

ALERT!!!

Be proactive on this, check regularly.

pg_catalog.pg_stat_database

Column	Type
datid	oid
datname	name
session_time	double precision
active_time	double precision
idle_in_transaction_time	double precision
sessions	bigint
sessions_abandoned	bigint
sessions_fatal	bigint
sessions_killed	bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (5 - 60 \text{ mins}, 1 \text{ day})$$

Updated on « pg_stat_activity.state_change »

$$ratio_{active} = \frac{\Delta active_time}{\Delta session_time}$$

$$rate_{idle_trans} = \frac{\Delta idle_in_transaction_time}{\Delta session_time}$$

Cluster Level Statistics

Standalone

pg_catalog.pg_stat_io

Column	Type	Column	Type
backend_type	text	evictions	bigint
object	text	reuses	bigint
context	text	fsyncs	bigint
reads	bigint	fsync_time	double precision
read_time	double precision	stats_reset	timestamp with time zone
writes	bigint		
write_time	double precision		
writebacks	bigint		
writeback_time	double precision		
extends	bigint		
extend_time	double precision		
op_bytes	bigint		
hits	bigint		

pg_catalog.pg_stat_io

Column	Type
<code>backend_type</code>	<code>text</code>
<code>stats_reset</code>	<code>timestamp with time zone</code>

autovacuum launcher,
autovacuum worker,
logical replication launcher,
logical replication worker,
client backend,
parallel worker,
background writer,
checkpointer,
archiver,
standalone backend,
startup,
walreceiver,
walsender,
walwriter
walsummarizer

pg_catalog.pg_stat_io

Column	Type	
backend_type	text	
object	text	<ul style="list-style-type: none">• relation• temp relation
context	text	<ul style="list-style-type: none">• normal• vacuum• bulkread• bulkwrite• temporary relation

pg_catalog.pg_stat_io

Column	Type	
backend_type	text	
object	text	
context	text	
<u>op_bytes</u>	<u>bigint</u>	IO Operation Size (BLOCKSIZE = 8192)

pg_catalog.pg_stat_io

Column	Type
backend_type	text
object	text
context	text
reads	bigrnt
writes	bigrnt
writebacks	bigrnt
writeback_time	double precision
extends	bigrnt
extend_time	double precision
hits	bigrnt

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1s - 15 \text{ mins})$$

$$\text{cache}_{\text{hit_ratio}} = \frac{\Delta \text{hits}}{\Delta \text{hits} + \Delta \text{reads}}$$

$$\text{cache}_{\text{pending_writes}} = \Delta \text{writes} - \Delta \text{writebacks}$$

$$\text{avg}_{\text{writeback_time}} = \frac{\Delta \text{writeback_time}}{\Delta \text{writebacks}}$$

$$\text{avg}_{\text{extend_time}} = \frac{\Delta \text{extend_time}}{\Delta \text{extends}}$$

pg_catalog.pg_stat_io

Column	Type
backend_type	text
object	text
context	text
evictions	bigrnt
reuses	bigrnt
fsyncs	bigrnt
fsync_time	double precision

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1s - 15 \text{ mins})$$

reuses is related to BAS Ring Buffers

$$\text{avg}_{\text{fsync_time}} = \frac{\Delta \text{fsync_time}}{\Delta \text{fsyncs}} \leq \sim 100ms$$

pg_catalog.pg_stat_wal

Column	Type
wal_records	bigint
wal_fpi	bigint
wal_bytes	numeric
wal_buffers_full	bigint
wal_write	bigint
wal_sync	bigint
wal_write_time	double precision
wal_sync_time	double precision
stats_reset	timestamp with time zone

pg_catalog.pg_stat_wal

Column	Type
wal_records	bigint
wal_fpi	bigint
wal_bytes	numeric
wal_buffers_full	bigint
wal_write	bigint
wal_write_time	double precision
wal_sync	bigint
wal_sync_time	double precision
stats_reset	timestamp with time zone

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1s - 5 \text{ mins})$$

$$\Delta \text{Volume}_{\text{fpi}} = \Delta \text{wal_fpi} * \text{wal_block_size}$$

$$\Delta \text{Volume}_{\text{rec}} = \Delta \text{wal_bytes} - \Delta \text{Volume}_{\text{fpi}}$$

$$\text{avg}_{\text{wal_write}} = \frac{\Delta \text{wal_write_time}}{\Delta \text{wal_write}}$$

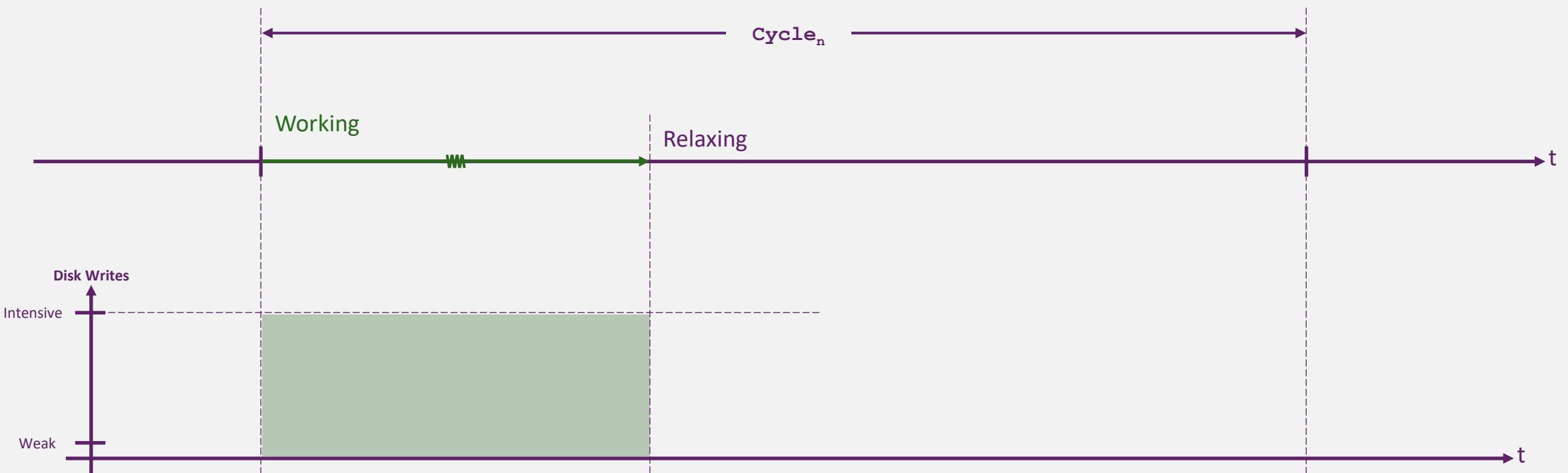
$$\text{avg}_{\text{wal_sync}} = \frac{\Delta \text{wal_sync_time}}{\Delta \text{wal_sync}} \leq 100\mu\text{s} (\text{pg_test_sync})$$

$$\text{ratio}_{\text{fullwrite}} = \frac{\Delta \text{wal_buffer_full} * \text{wal_buffers} * \text{wal_block_size}}{\Delta \text{wal_bytes}}$$

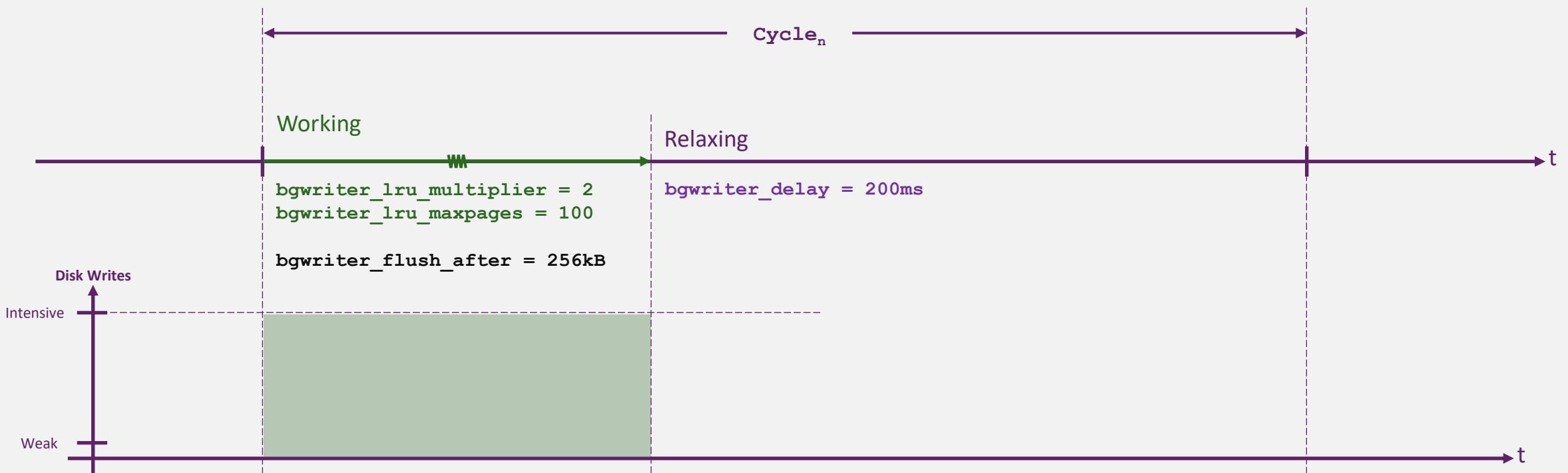
Background Writer

Structure – Usage – Processing

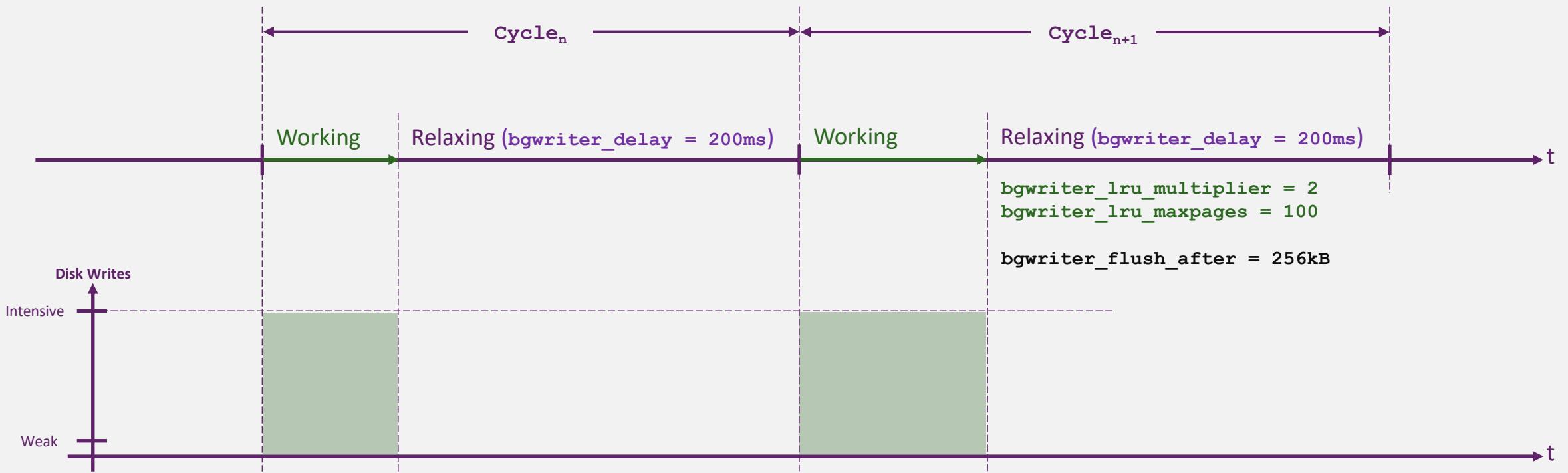
ACTIVITY and IO PRESSURE



ACTIVITY and IO PRESSURE



ACTIVITY and IO PRESSURE



pg_catalog.pg_stat_bgwriter

Column (pg16)	Type	Column (pg17)	Type
checkpoints_timed	bigint	buffers_clean	bigint
checkpoints_req	bigint	maxwritten_clean	bigint
checkpoint_write_time	double precision	buffers_alloc	bigint
checkpoint_sync_time	double precision	stats_reset	timestamp with time zone
buffers_checkpoint	bigint	Column (pg17)	Type
buffers_clean	bigint	num_timed	bigint
maxwritten_clean	bigint	num_requested	bigint
buffers_backend	bigint	restartpoints_timed	bigint
buffers_backend_fsync	bigint	restartpoints_req	bigint
buffers_alloc	bigint	restartpoints_done	bigint
stats_reset	timestamp with time zone	write_time	double precision
		sync_time	double precision
		buffers_written	bigint
		stats_reset	timestamp with time zone

pg_catalog.pg_stat_bgwriter

Column	Type
<code>buffers_clean</code>	<code>bigint</code>
<code>maxwritten_clean</code>	<code>bigint</code>
<code>buffers_alloc</code>	<code>bigint</code>
<code>stats_reset</code>	<code>timestamp with time zone</code>

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (\textcolor{red}{1s} - 15 \text{ mins})$$

- $maxwritten_clean \times bgwriter_lru_maxpages \leq buffers_clean$
- pg16
 - Minimize `pg_stat_bgwriter.backend_buffers`
- Pg17
 - Minimize `pg_stat_io.evictions` for « client backend »

Checkpoint(er)

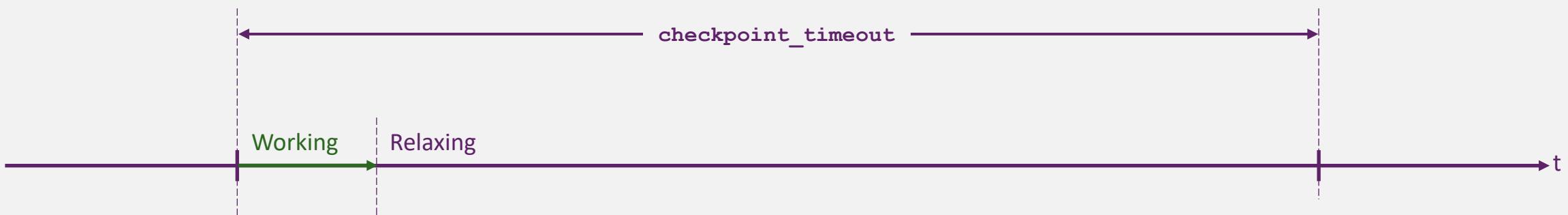
Understanding Performance Key Points

CHECKPOINT KEY CONCEPTS

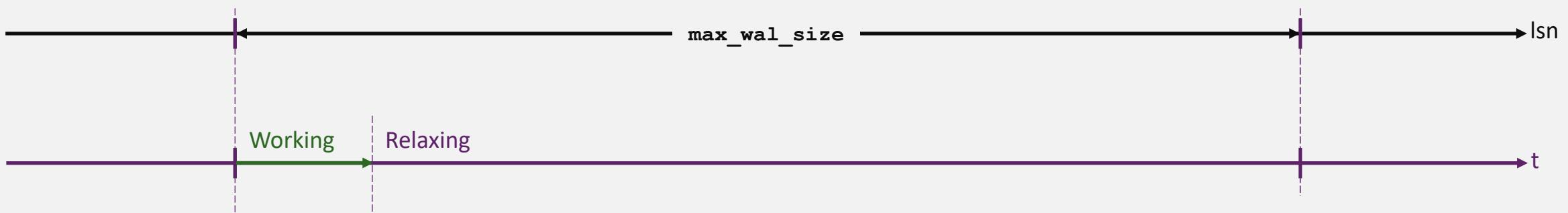
`checkpoint_timeout - max_wal_size - checkpoint_completion_target`
`checkpoint_warning - min_wal_size - checkpoint_flush_after`

CHECKPOINT TRIGGERING

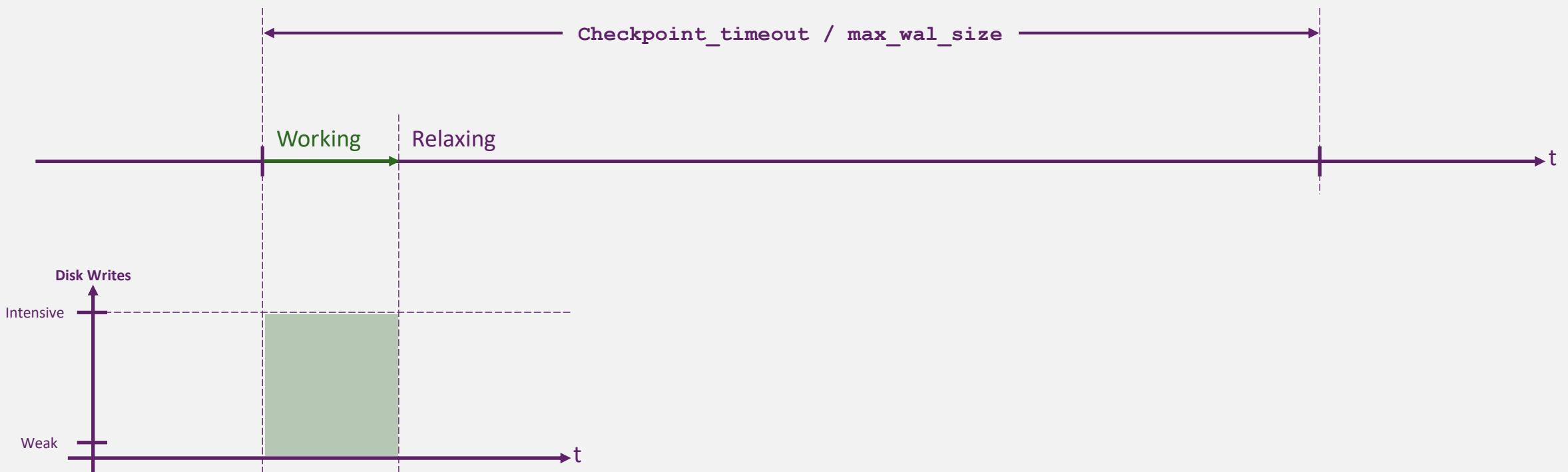
TEMPORAL TRIGGERING



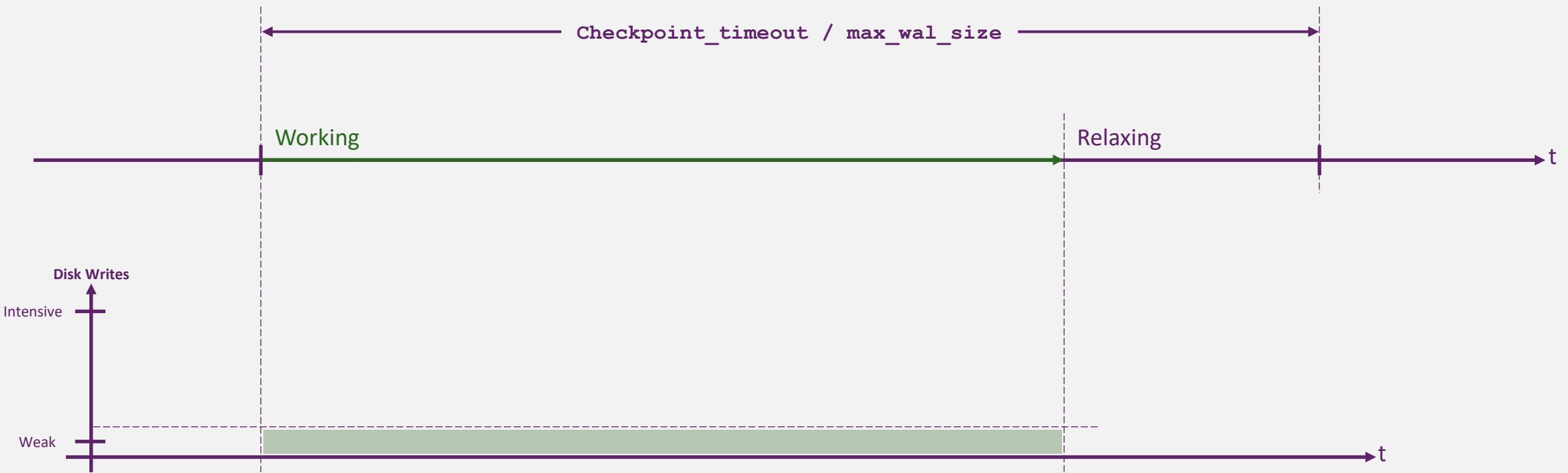
SPATIAL TRIGGERING



IO PRESSURE (intensive)

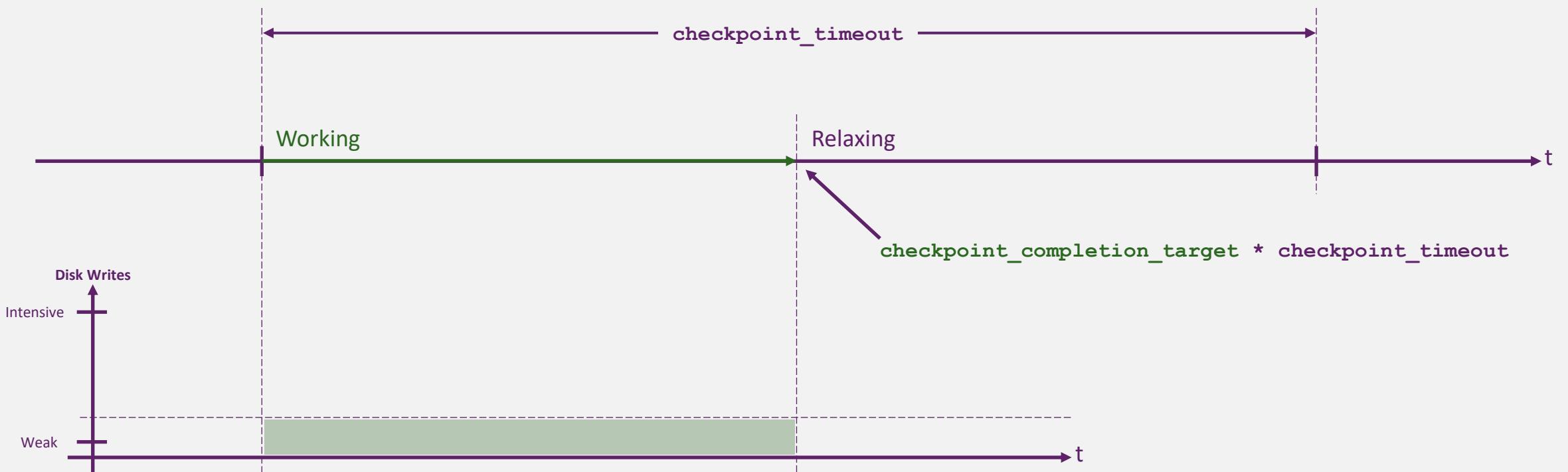


IO PRESSURE (spread)

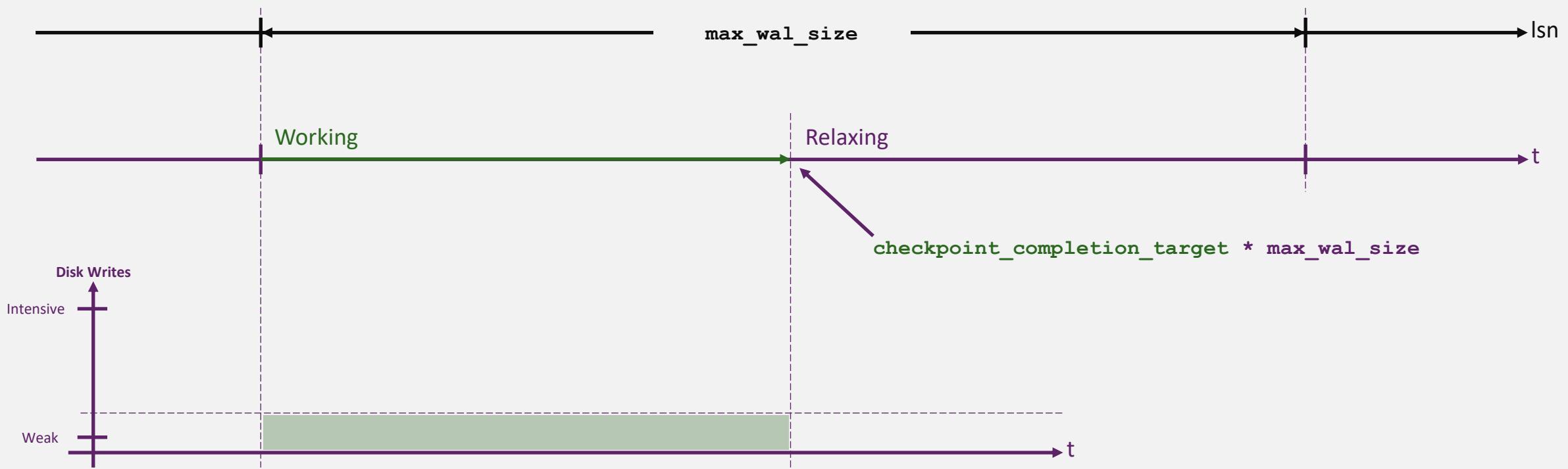


IO PRESSURE SPREADING FEEDBACK LOOP

TEMPORAL SPREADING FEEDBACK LOOP



SPATIAL SPREADING FEEDBACK LOOP



FEEDBACK LOOP

```
while num_processed < num_to_scan

    WriteDirtyPageToDisk()
    num_processed++

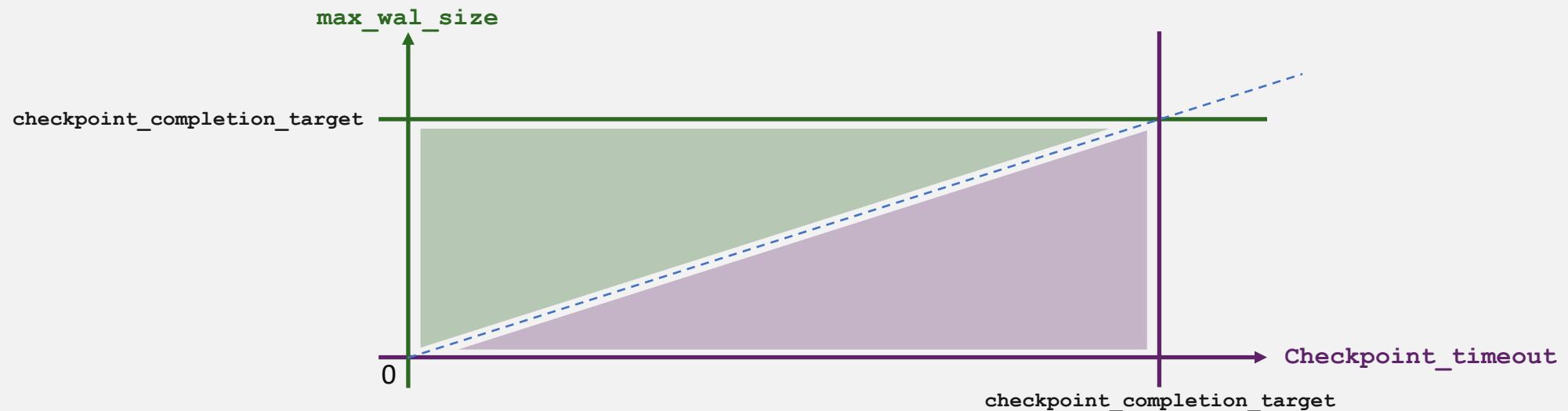
    checkpoint_progress = num_processed / num_to_scan
    checkpoint_progress = checkpoint_progress * checkpoint_completion_target

    wal_progress = (current LSN - start LSN) / max_wal_size
    if checkpoint_progress < wal_progress then continue

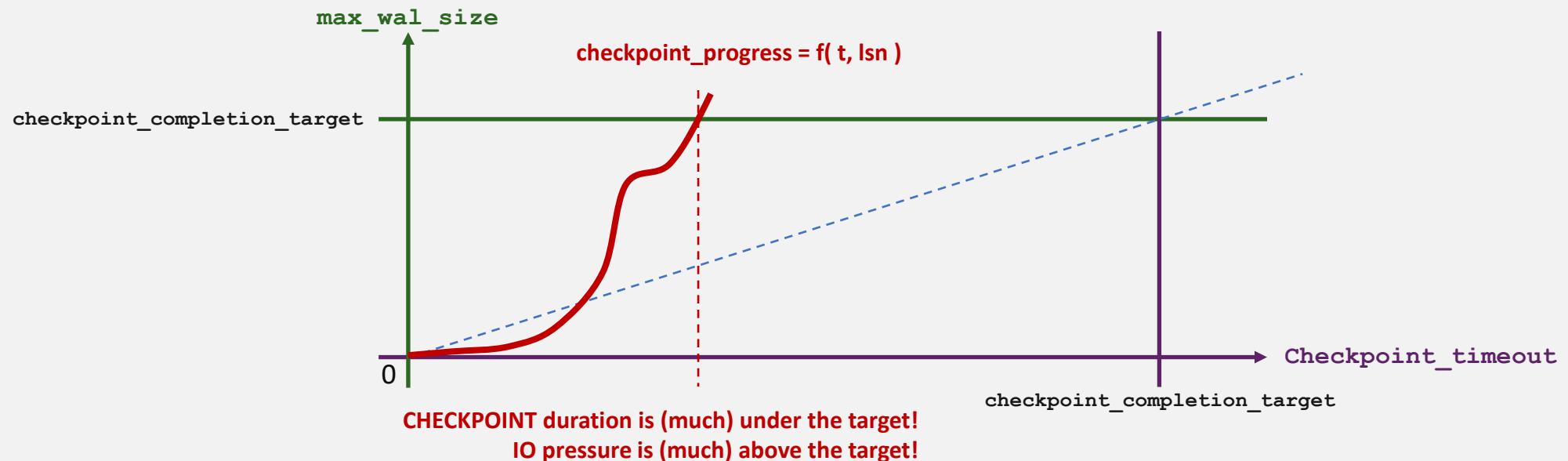
    time_progress = (now - start_time) / checkpoint_timeout
    if checkpoint_progress < time_progress then continue

    sleep( CheckpointWriteDelay = 100ms )
```

BIDIMENSIONNAL FEEDBACK LOOP DIAGRAM



BIDIMENSIONNAL FEEDBACK LOOP – EXAMPLE



pg_catalog.pg_stat_checkpointer

Column	Type
num_timed	bigint
num_requested	bigint
restartpoints_timed	bigint
restartpoints_req	bigint
restartpoints_done	bigint
write_time	double precision
sync_time	double precision
buffers_written	bigint

$$rate_x = \frac{\Delta x}{\Delta T} ; \Delta T \in (1 \text{ mins})$$

- sync_time must be kept low
 $\Delta sync_time \leq \sim 100ms$
- Check IO spreading
$$\text{avg}_{\text{iops}} = \frac{\Delta buffer_written}{\Delta write_time} \leq IOPressure$$
- Check completion_target is respected
$$\frac{avg_w(write_time)}{checkpoint_timeout} = checkpoint_completion_target$$

Logging

- Do not forgot to logs some events
 - logging_collector = on
 - log_destination = 'stderr, csv' (csv can be injected to a database and then queried)
 - log_temp_files = 1
 - log_lock_waits = on
 - log_line_prefix = '%m [%p] %a %u@%d %a' as a minimum
 - log_checkpoints = on
 - log_autovacuum_min_duration = 60000 (ms)
 - log_min_duration_statement = the value which make sense for you
- Mine your log every day
 - pgbadger can help with that

Hey, I did not cover the following yet:

- pg_stat_slru
- pg_stat_replication
 - `SELECT pg_wal_lsn_diff(pg_current_wal_lsn(), sent_lsn) AS sent_lag_bytes FROM pg_stat_replication;`
- pg_class size catalog tables
 - `SELECT relname, relpages, pg_size.pretty(8192::bigint * relpages) FROM pg_class WHERE relnamespace = 'pg_catalog'::regnamespace ORDER BY 2 DESC;`
- pg_index is not valid (`pg_index indisvalid`)
 - `select * from pg_index where not indisvalid;`
- pg_stat_archiver
 - `SELECT now() - COALESCE(last_failed_time, '1970-01-01'::date) <= interval '5 mins' AS recent_failure FROM pg_stat_archiver ;`
- pg_stat_activity (wait_event_type, wait_event, state, long running query/snapshot) with pg_locks (granted?)

Linux

procFS, cgroups, perf

ProcFS - CPU

```
$ cat /proc/stat
```

#	user	nice	system	idle	iowait	irq	softirq	steal	guest	guest_nice
cpu	31597599	2158402	9278540	1083188543	1623186	0	303365	0	0	0
cpu0	6700402	488305	2278642	272134488	392242	0	19813	0	0	0
cpu1	7706591	563056	2072923	271348166	413324	0	18152	0	0	0
cpu2	8151274	502605	2405923	270298148	422380	0	173093	0	0	0
cpu3	9039330	604434	2521050	269407739	395239	0	92305	0	0	0

$$CPU\ Utilization = \frac{user + nice + system + stealth + guest}{user + nice + system + stealth + guest + idle [+ iowait]} < 0.85$$

ProcFS - CPU

```
$ cat /proc/stat
...
intr 697368832 48 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 25 23437850 15669310 2 11670157 0 0 0 ... 0 0 0
ctxt 426424223
btime 1726754569
processes 1486375
procs_running 1
procs_blocked 0
#
#           TIMER      NET_TX     NET_RX      BLOCK      TASKLET    SCHED      HRTIMER   RCU
softirq 837985299 6 108387462 15434898 35102963 11475396 0 4390        475399864 48060        192132260
```

ProcFS - Network

```
$ cat /proc/net/dev
```

Inter- Receive	face	bytes	packets	errs	drop	fifo	frame	compressed	multicast
	lo:	19449	119	0	0	0	0	0	0
	eno0:	10296287653	26877099	0	0	0	0	0	1416671
	docker0:	684107064	12947192	0	0	0	0	0	0
	veth303a4b4:	250762391	3804884	0	0	0	0	0	0
	veth37d40dd:	251701833	3817537	0	0	0	0	0	0
	veth4c6976c:	33518004	489729	0	0	0	0	0	0
	veth75f0dc9:	252377309	3828110	0	0	0	0	0	0
	veth49d99b2:	24114346	312563	0	0	0	0	0	0
	veth4fd0c00:	23808783	309302	0	0	0	0	0	0
	veth2a4269b:	23354884	302527	0	0	0	0	0	0

ProcFS - Disk

```
$ cat /proc/diskstats
8 0 sda 1274558 151495 85887753 4817963 9019869 7340199 238013064 28267395 0 24084484 48700750 0 0 0 0 1615784 15615391
8 1 sda1 236 0 11618 368 0 0 0 0 660 368 0 0 0 0 0 0
8 2 sda2 437 211 40076 1732 257 114 85912 2763 0 3452 4495 0 0 0 0 0 0
```

ProcFS - Diskstat

```
== =====  
1 major number  
2 minor number  
3 device name  
4 reads completed successfully  
5 reads merged  
6 sectors read  
7 time spent reading (ms)  
8 writes completed  
9 writes merged  
10 sectors written  
11 time spent writing (ms)  
12 I/Os currently in progress  
13 time spent doing I/Os (ms)  
14 weighted time spent doing I/Os (ms)  
== =====
```

Kernel 4.18+ appends four more fields for discard

tracking putting the total at 18:

```
== =====  
15 discards completed successfully  
16 discards merged  
17 sectors discarded  
18 time spent discarding  
== =====
```

Kernel 5.5+ appends two more fields for flush requests:

```
== =====  
19 flush requests completed successfully  
20 time spent flushing  
== =====
```

ProcFS – Diskstat – Service Time

$$\frac{\delta[io_{ticks}]}{(\delta[read_{IOs} + read_{merges} + write_{IOs} + write_{merges}])}$$

ServiceTime is given in milliseconds.

Rewritten with field numbers:

$$\frac{\delta[f13]}{(\delta[f4 + f5 + f8 + f9])}$$

ProcFS – Diskstat – Queue Time

$$\frac{\delta[time_{in_queue}]}{(\delta[read_{IOs} + read_{merges} + write_{IOs} + write_{merges}] + \delta[in_flight])} - servicetime$$

QueueTime is given in milliseconds.

Rewritten with field numbers:

$$\frac{\delta[f14]}{(\delta[f4 + f5 + f8 + f9] + \delta[f12])} - \frac{\delta[f13]}{(\delta[f4 + f5 + f8 + f9])}$$

cgroups

- cgroups are used everywhere and automatically
- /sys/fs/cgroup
 - blkio
 - cpu | cpu,cpuacct | cpuacct | cpuset
 - devices
 - freezer
 - hugetlb | memory
 - misc
 - net_cls | net_cls,net_prio | net_prio
 - perf_event
 - pids
 - rdma
 - systemd
 - unified

CPU consumption of my PostgreSQL service

```
$ cat /sys/fs/cgroup/cpuacct/system.slice/ \
    system-postgresql.slice/cpuacct.usage_all
cpu user system
0 1841620000000 3844232000000
1 7006000000000 1351788000000
2 1627472000000 3376616000000
3 1489560000000 3066904000000
```

Unit: 10ms (1/100th s) => direct percent of CPU consumption.

Compute rate as usual $rate_{cpu\%} = \frac{\Delta user + \Delta system}{\Delta T}$

IODisk comsumption for my docker

```
cat /sys/fs/cgroup/blkio/docker/0f040c...8c7bc0e3\  
/blkio.throttle.io_service_bytes  
8:0 Read 3897999360  
8:0 Write 2056773632  
8:0 Sync 4718170112  
8:0 Async 1236602880  
8:0 Discard 0  
8:0 Total 5954772992  
Total 5954772992
```

Perf

- Very interesting tool you have to try at least once.
- Only for troubleshooting because it is kind of manual.
- Usually perf.data file has to be processed on the recorded server.
 - Processing data file can be consuming

Perf – System Probe

```
sudo perf list
```

Gives a lot of available hardware and system metrics!

Perf - Example

```
sudo \
perf stat -e
block:block_rq_*,syscalls:sys_enter_write,syscalls:sys_enter_fsync \
-a -r 5 -- psql -q -p 5434 -U postgres postgres \
-c "drop table if exists x; create table x as select a FROM
generate_series(1,1000000) a;";
```

Perf – Completely Fair Scheduler

To investigate

```
perf sched record sleep 0.1
```

```
perf sched latency -s max
```

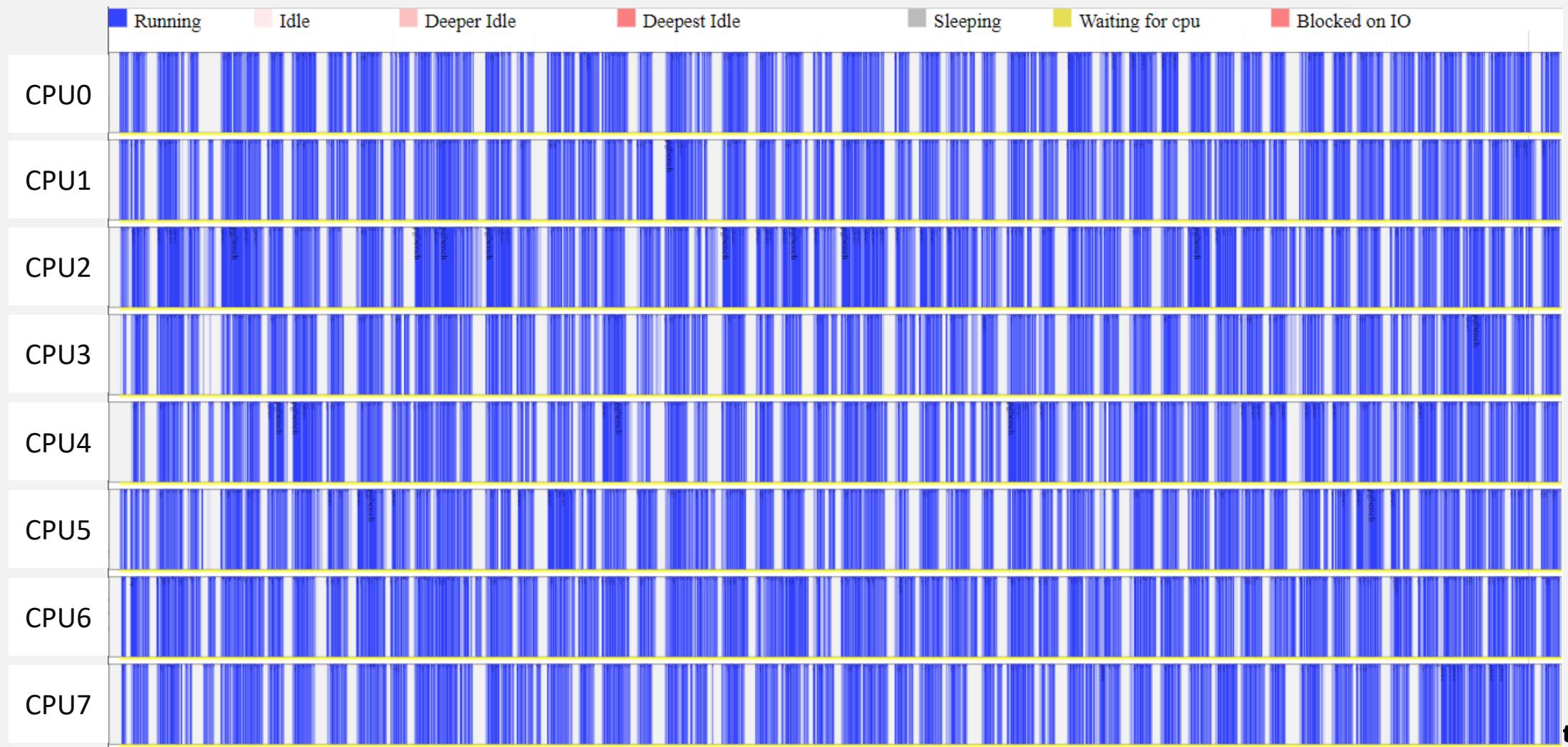
```
perf sched map
```

```
perf timechart
```

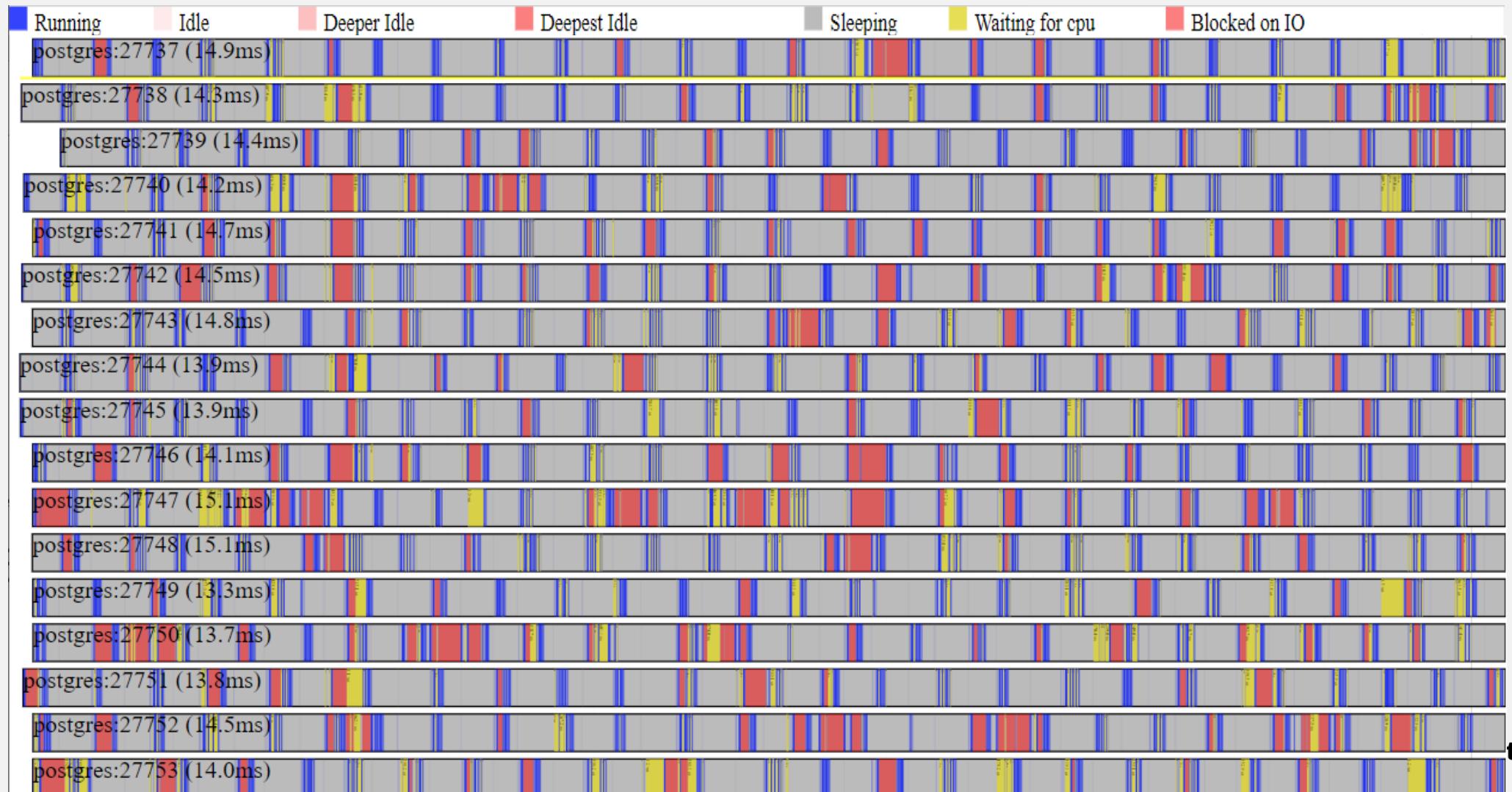
[Announce] 'perf sched': Utility to capture, measure and analyze scheduler latencies and behavior

<https://lwn.net/Articles/353295/>

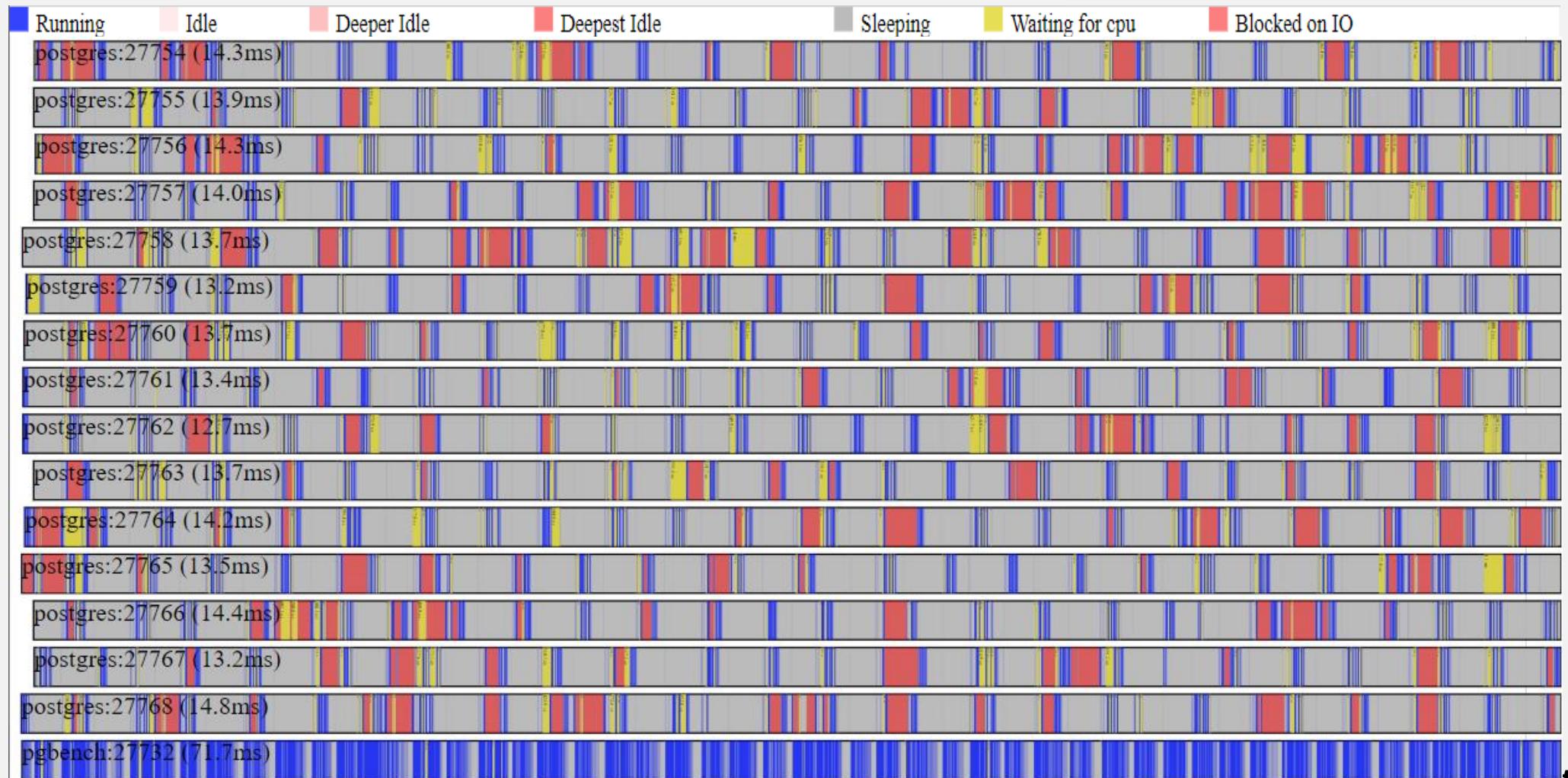
Linux – Completely Fair Scheduler



Linux – Completely Fair Scheduler



Linux – Completely Fair Scheduler



Thank You

Without Heap Only Tuples

Structure – Usage – Processing

PageHeader_k

Free Space

Index Points to



Line Pointer₀ (LP_NORMAL)

PageHeader_k

Insert Tuple_A

Free Space

Tuple_A xmin=53, xmax= , ctid=(k, 0)



LP_NORMAL

LP_REDIRECT

LP_DEAD

LP_UNUSED

Unvisible

 DataBene

Index Points to



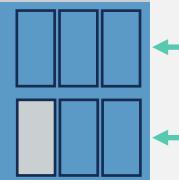
Line Pointer₀ (LP_NORMAL)
Line Pointer₁ (LP_NORMAL)

Index Points to



Free Space

Tuple_B xmin=54, xmax= , ctid=(k,1)
Tuple_A xmin=53, xmax=54, ctid=(k,1)



Insert Tuple_A
UPDATE_{A->B}

LP_NORMAL

LP_REDIRECT

LP_DEAD

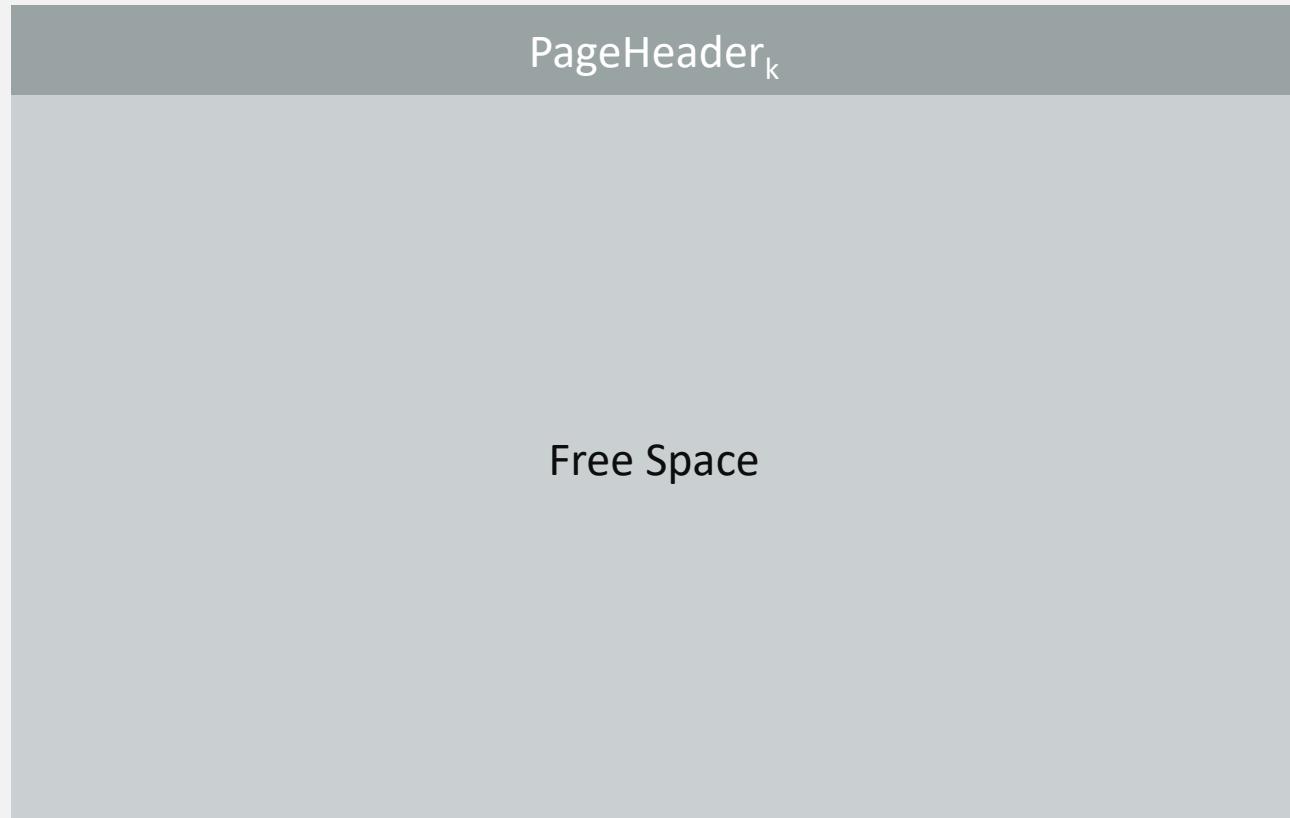
LP_UNUSED

Unvisible

 DataBene

Heap Only Tuples

Structure – Usage – Processing



LP_NORMAL

LP_REDIRECT

LP_DEAD

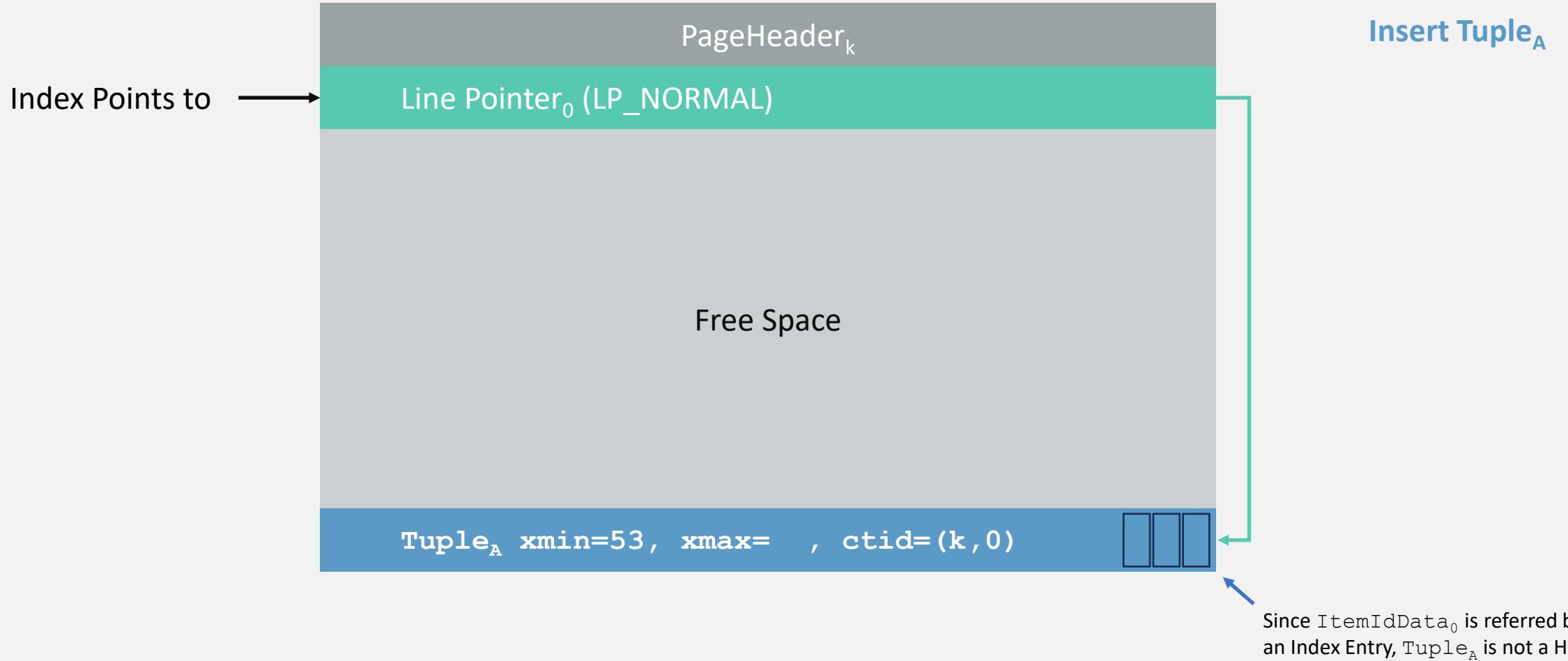
LP_UNUSED

Unvisible

HEAP_ONLY_UPDATE

HEAP_ONLY_TUPLE

 DataBene



LP_NORMAL

LP_REDIRECT

LP_DEAD

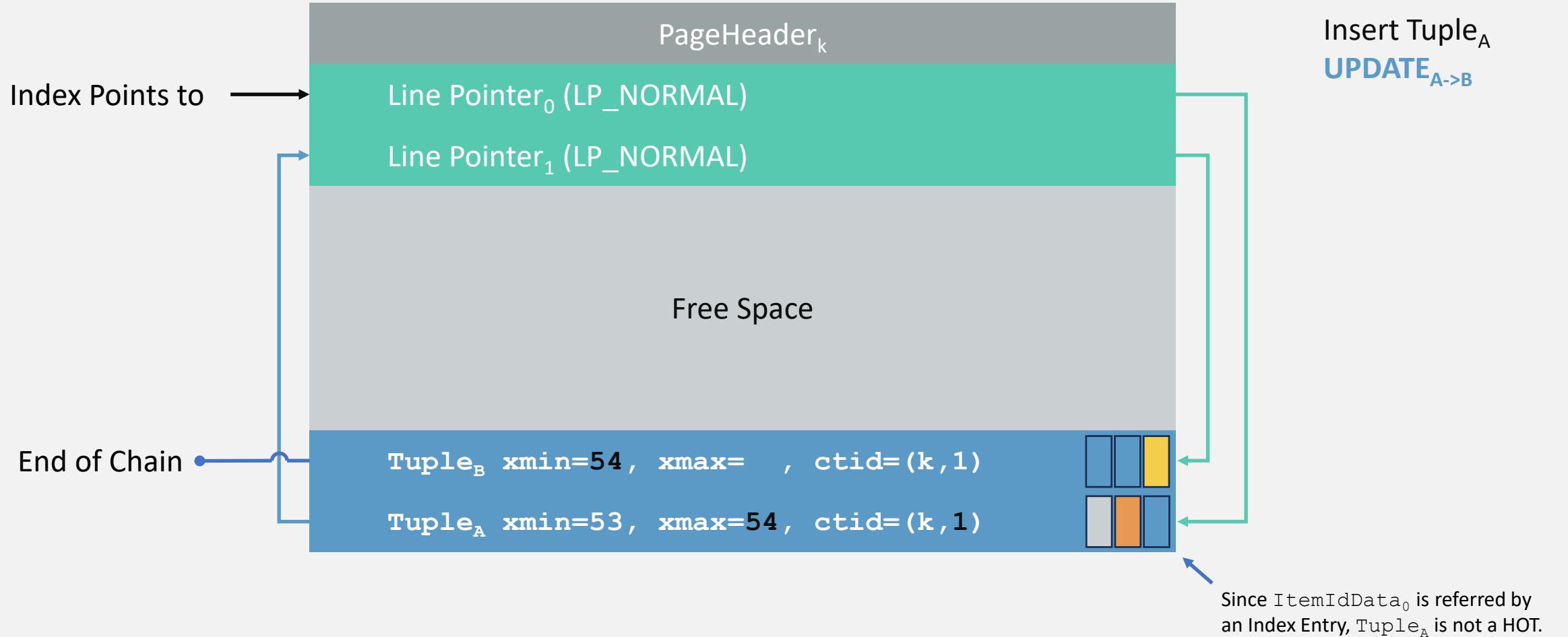
LP_UNUSED

Unvisible

HEAP_ONLY_UPDATE

HEAP_ONLY_TUPLE

DataBene



LP_NORMAL

LP_REDIRECT

LP_DEAD

LP_UNUSED

Unvisible

HEAP_ONLY_UPDATE

HEAP_ONLY_TUPLE

DataBene

