

Vacuum through Pictures



About 2ndQuadrant

- Founded by an Enterprise Architect / Postgres Developer
- Contributed to features to help make Postgres Enterprise Ready
 - Backup and Recovery
 - Point in Time Recovery
 - Streaming Replication
 - Logical Replication
 - Stored Procedures with Transaction control
 - Performance improvements to partitioning
- Funded by support of Postgres Server





Why this talk

- Years of talking to smart people about vacuum
- Years talking to customers who are confused about vacuum
- Years talking to customers who were in pain
- Needed pictures to understand vacuum concepts
- Need more pictures around Postgres



Best Practices : Vacuum

- Understand what it is and why it is necessary
- Design for Vacuum
- Monitor Vacuum
 - Tools
 - Queries
 - Extensions
- Tune for Vacuum



Agenda

- Why Postgres Vacuum
- What it does
- How does it do what it does
- Your options for tuning it
- Best practices



Future Talks

- Freezing through pictures
- HOT updates and Fill through pictures
- Monitoring Vacuum

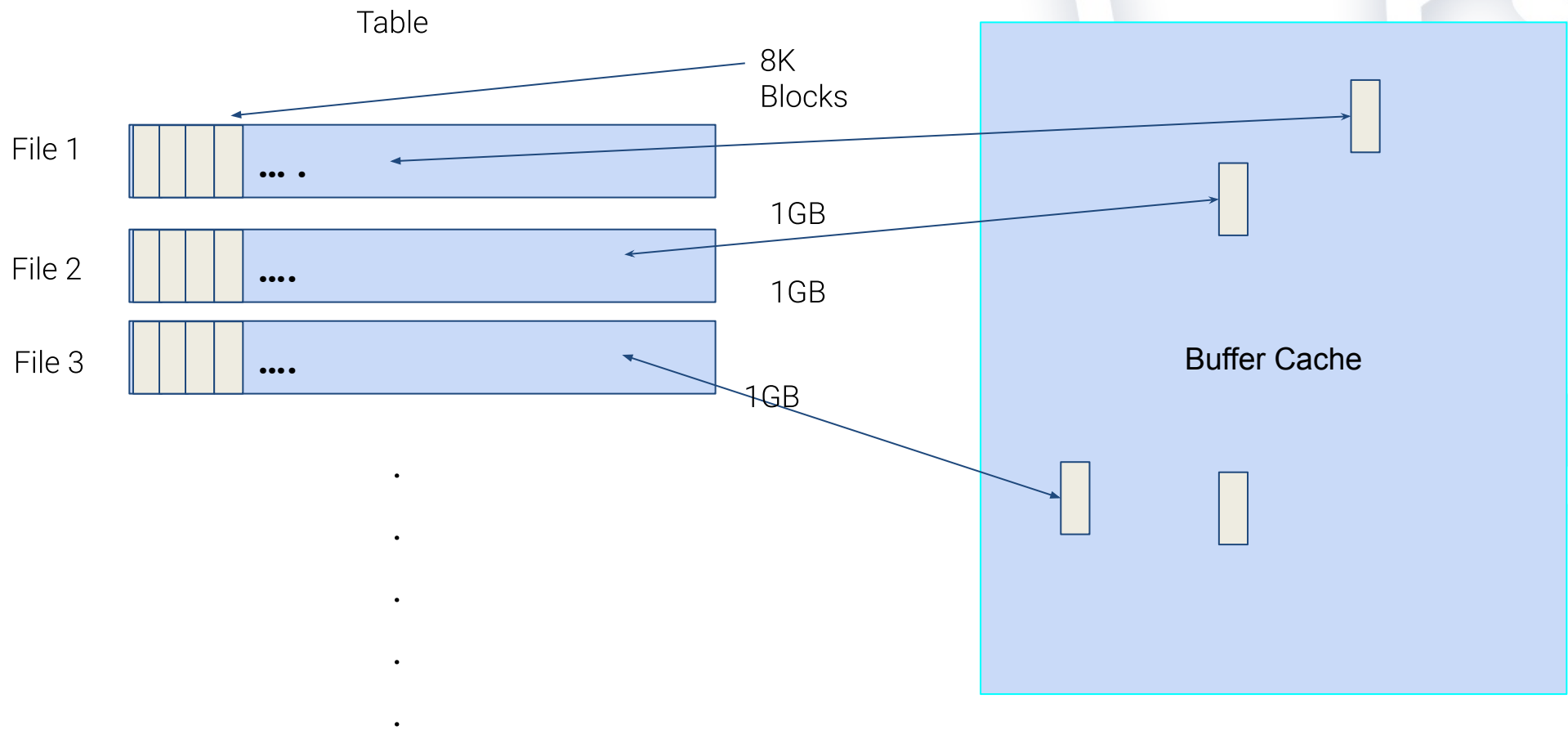


Postgres Internals Basics

- Postgres Tables reside on disk as files 1 GB in size
- Files consists of 8K (usually) pages or blocks
- A page is the unit of transfer from disk to memory (shared buffers)

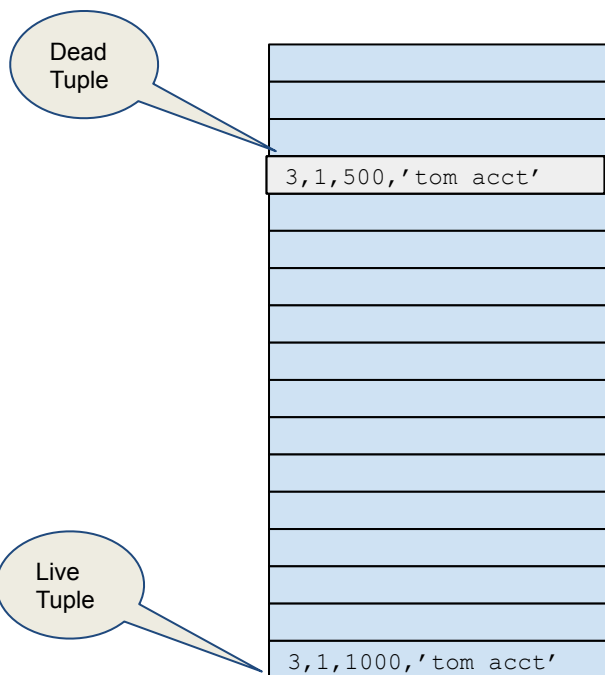


Tables, Files, Blocks and Memory





About MVCC



Trans.1 (starts at time 1)

Begin

Select bal from Account
where aid =3;

Commit

Trans. 2 (starts at time 2)

Begin
Update Accounts set bal =
1000
where aid = 3;
Commit

Trans. 2 (commits at time 3)

```
(1 row)

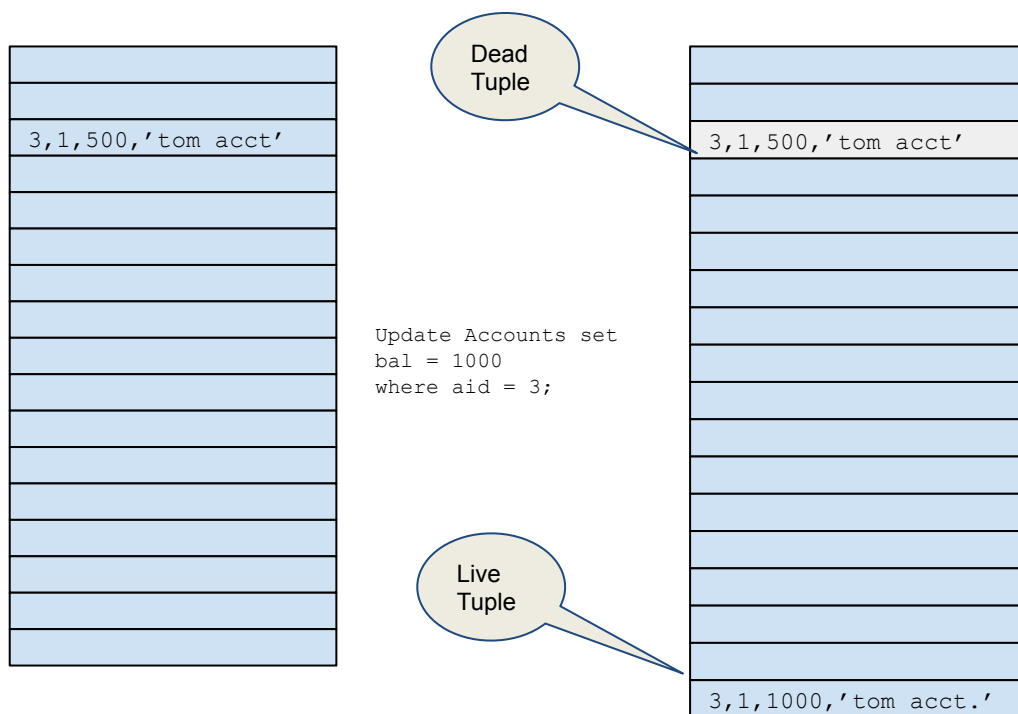
pgbench=# \d pgbench_accounts;
          Table "public.pgbench_accounts"
   Column   |      Type      | Collation | Nullable | Default
-----+-----+-----+-----+-----
 aid        | integer         |           | not null |
 bid        | integer         |           |          |
 abalance   | integer         |           |          |
 filler     | character(84)   |           |          |
Indexes:
    "pgbench_accounts_pkey" PRIMARY KEY, btree (aid)

pgbench=#
```

Different Transactions need to see different versions of a row depending on a variety of factors.



Something that happens regularly



```
Account (aid int, bid int, bal int, comment char(50))
```

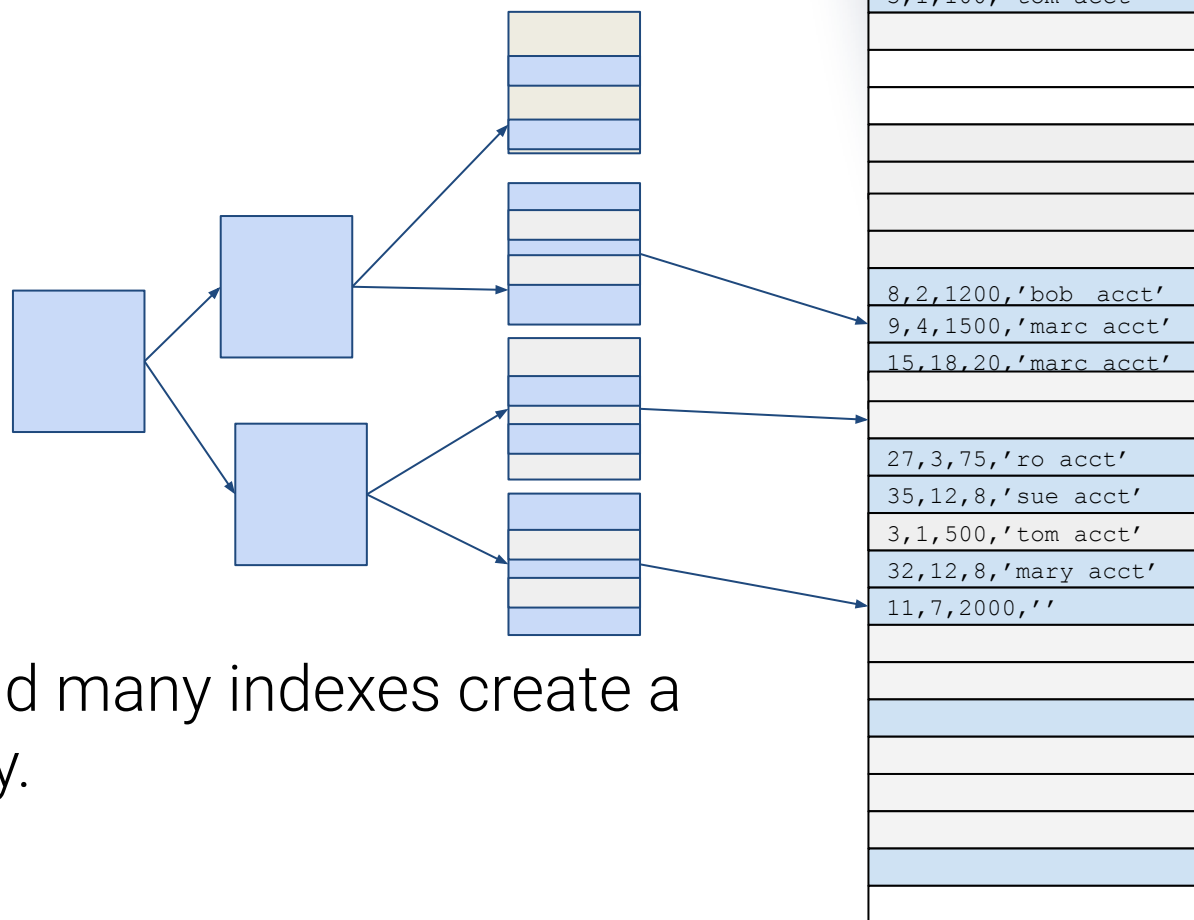
“The application was running great ... but then performance went sideways”

update	...
delete	...
update	...
update	...
update	...
delete	...
update	...
delete	...
update	...
update	...
delete	...
update	...
update	...
update	...
delete	...
update	...
delete	...
update	...
update	...
update	...
update	...
delete	...
update	...
delete	...
delete	...
update	...

[illegible]



Not just the table .. but the indexes



Many tables and many indexes create a problem quickly.

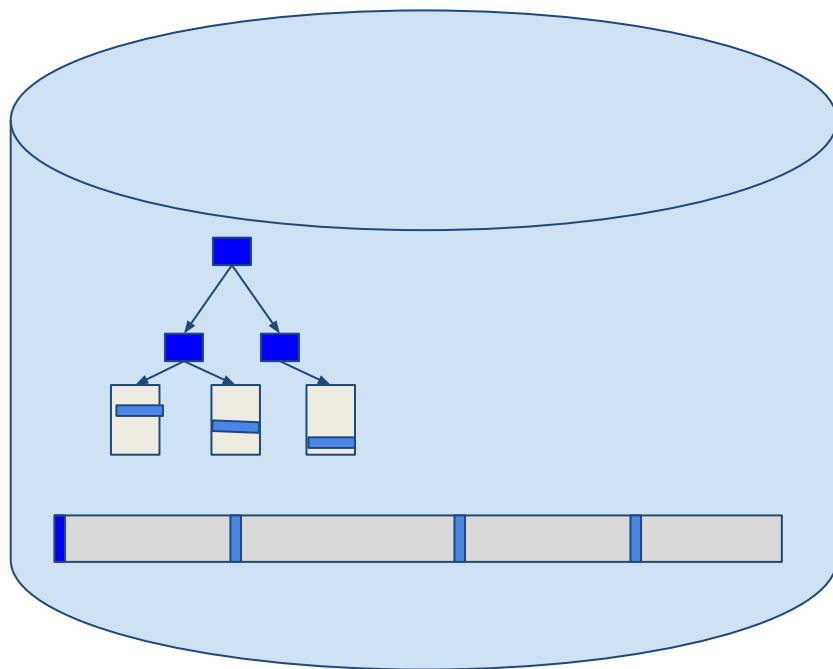


Why bloat is a problem

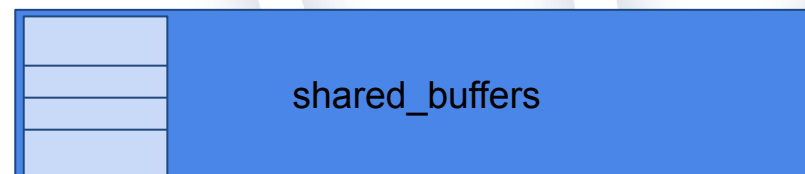
- Database occupies more disk
- Less actual data found in the cache
- Eventually you need to pay a huge vacuum price
- Pay a little frequently or Pay a lot later



Why excessive bloat is a problem



```
select * from pgbench_accounts;
```



or

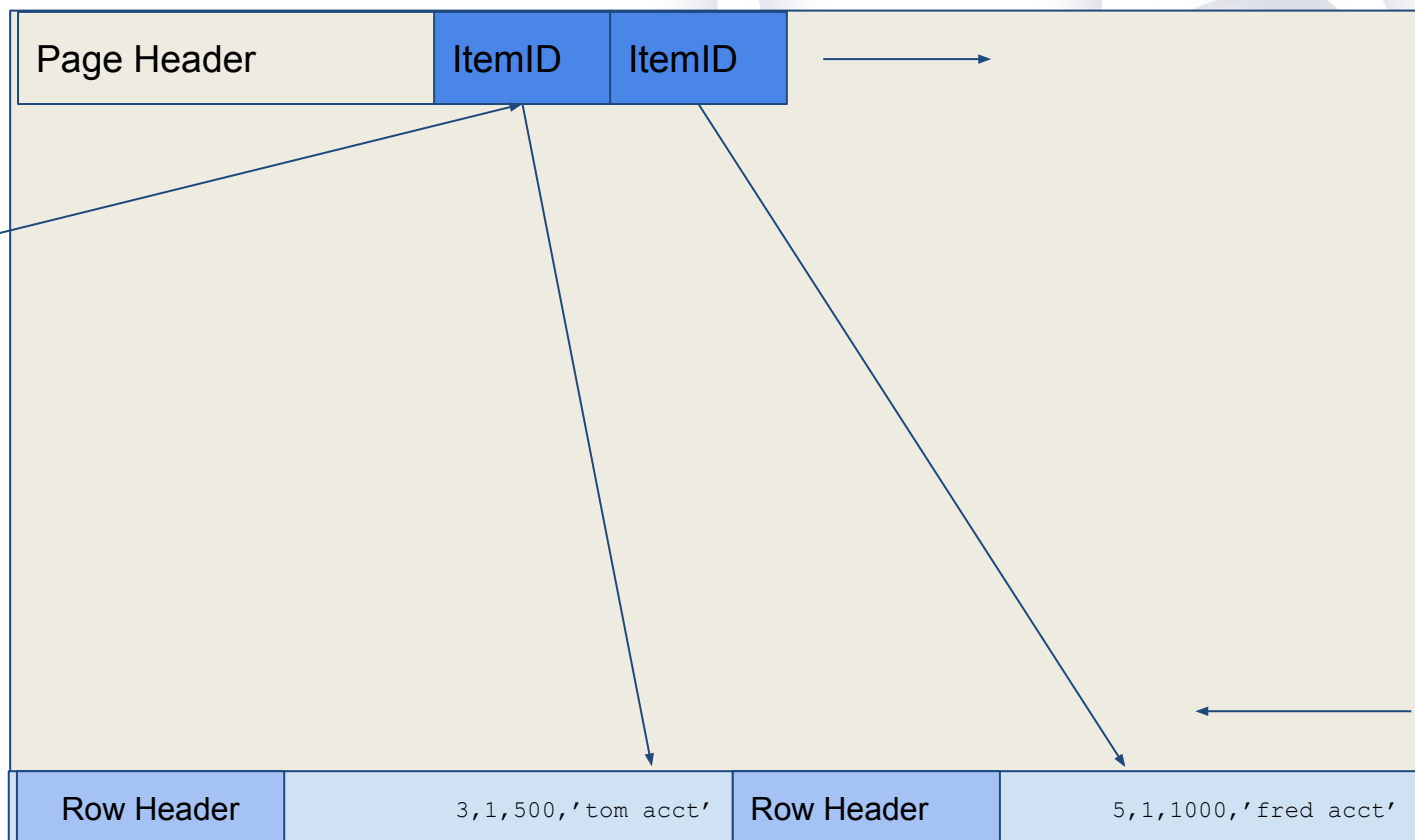
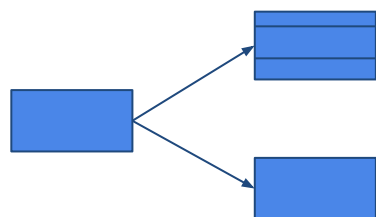




Some Elements of a Page and a Tuple

Typically 8K in size

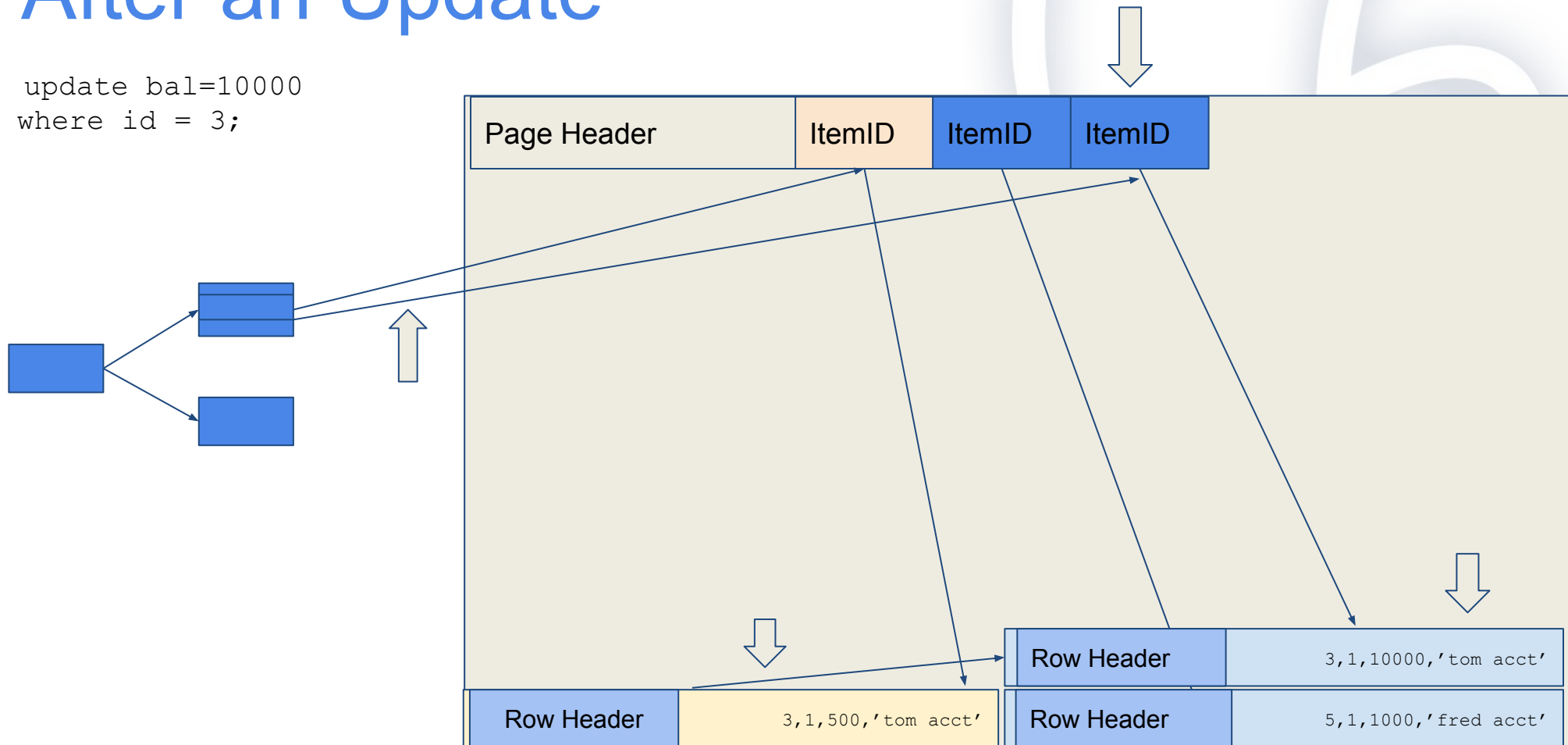
Indexes point to Item Pointers
and not to Tuples



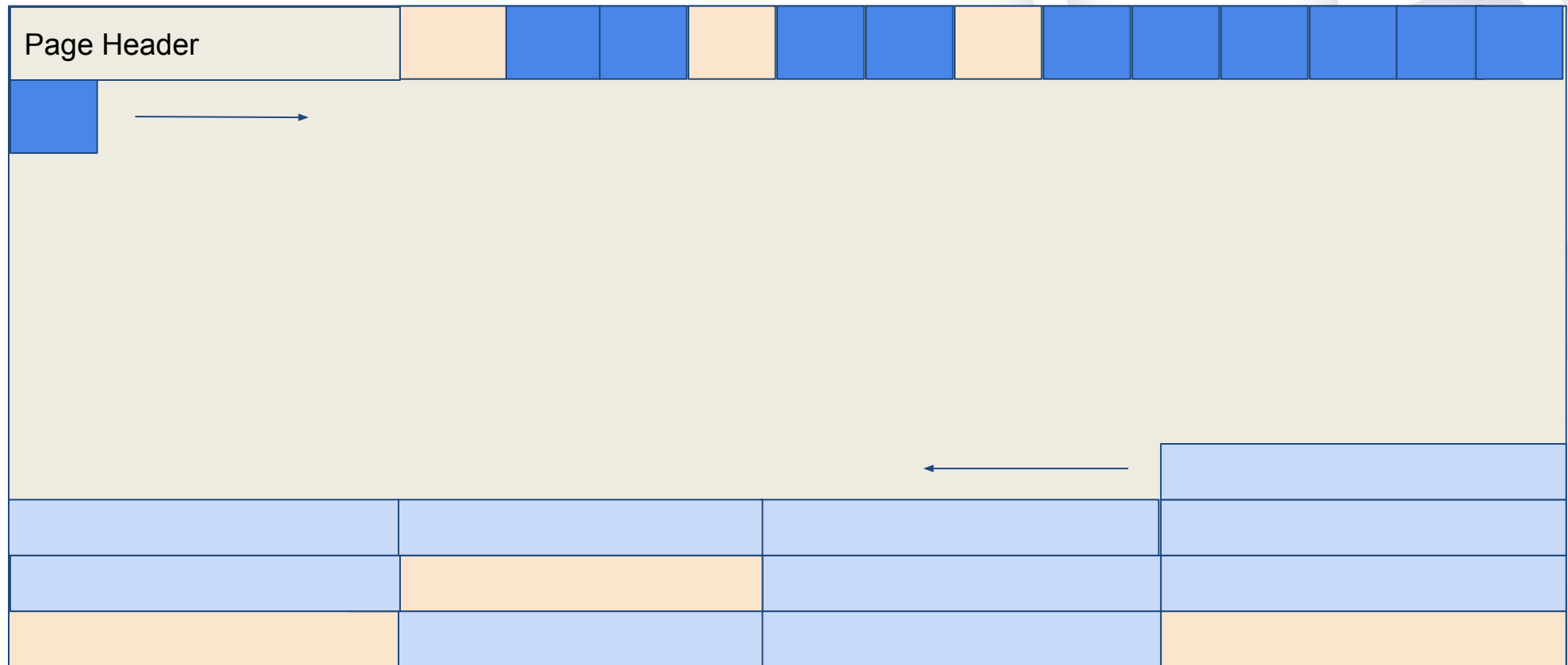


After an Update

```
update bal=10000  
where id = 3;
```



Over time



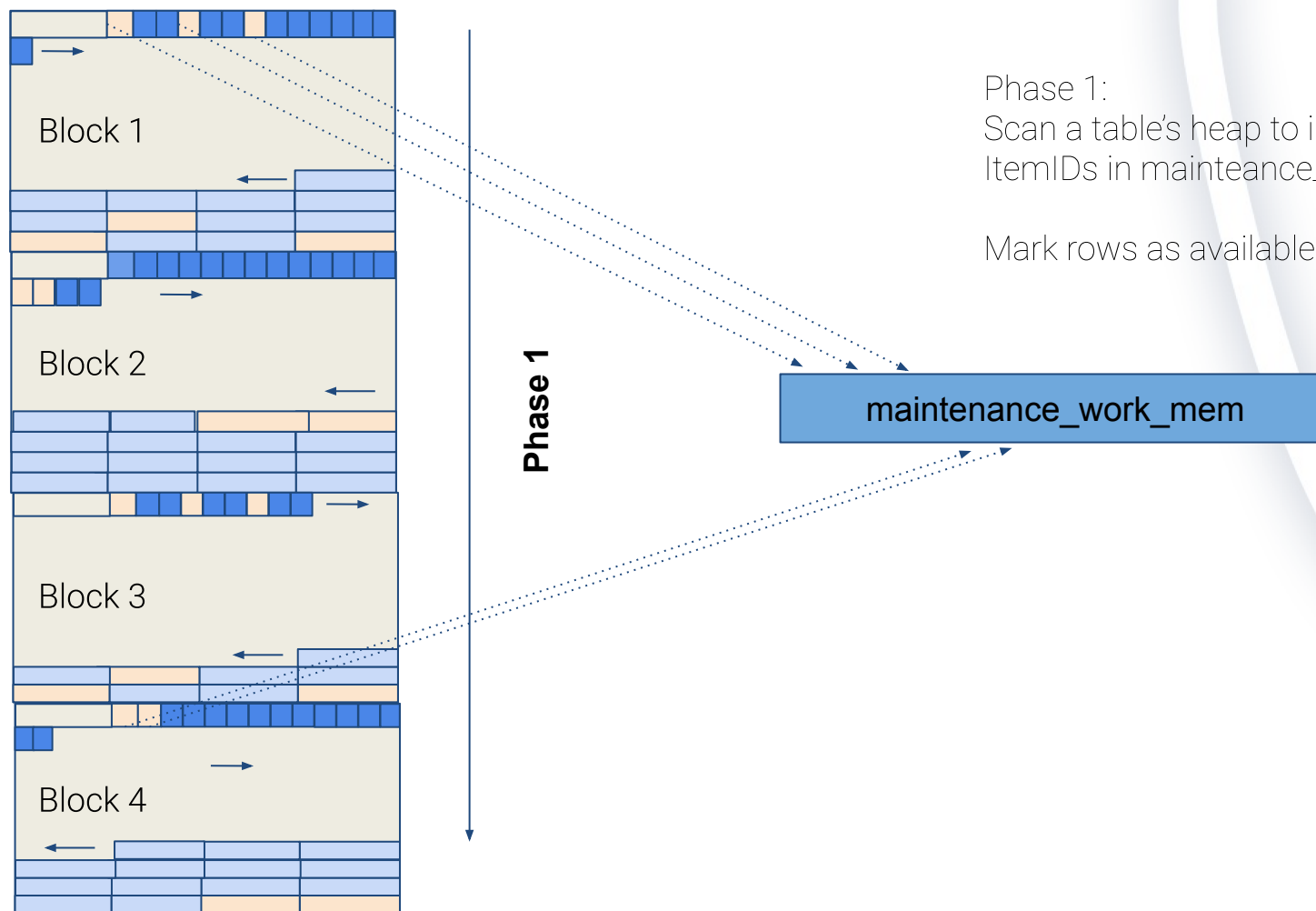


Vacuum is a 3 phase process

1. Scan the table to identify dead rows and ItemIDs that can be reclaimed
 - a. Place the ItemIDs in buffer (maintenance_work_mem in size)
 - b. Make row spaces reusable
 - c. Reorganize pages
 - d. Return space to OS if possible
2. Scan indexes to remove ItemIDs of dead rows reusable
3. After scanning indexes mark ItemIDs as re-suable in Table



A three phase process: Phase 1

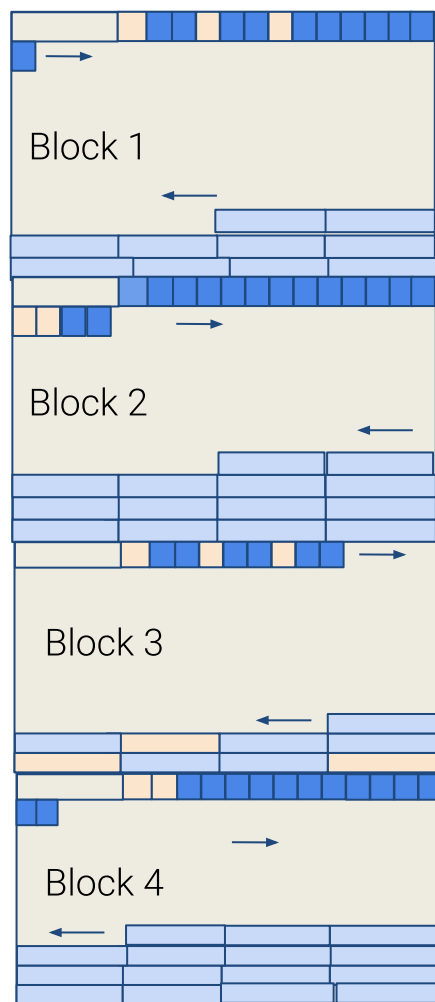


Phase 1:
Scan a table's heap to identify dead tuples, place
ItemIDs in maintenance_work_mem buffer.

Mark rows as available for future updates



A three phase process: Phase 1



Phase 1:

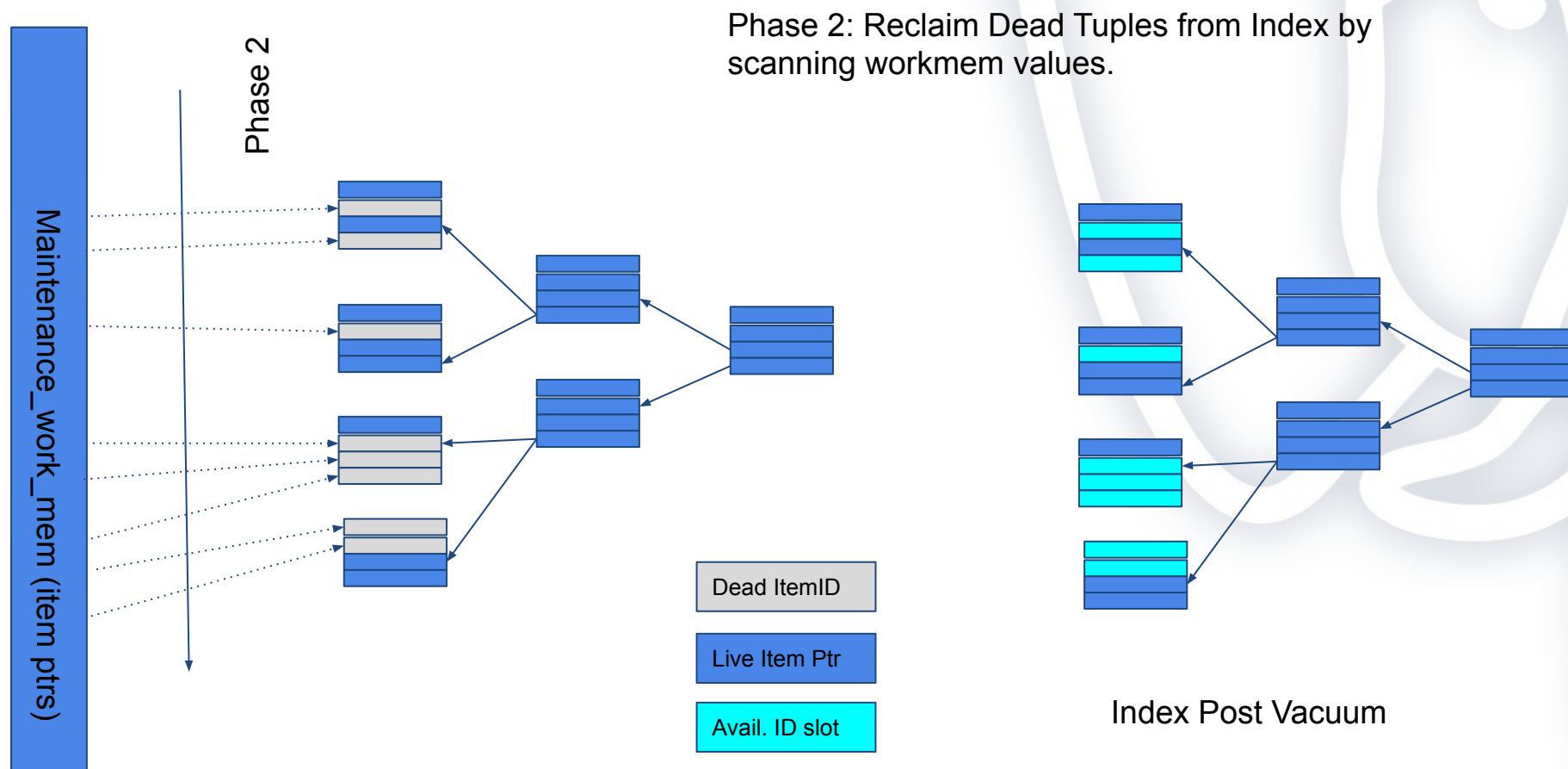
Reorganize Prunes / pages making room for additional tuples.

Updates free space map

Note: Block 3 still has tuples that are visible to some transactions.

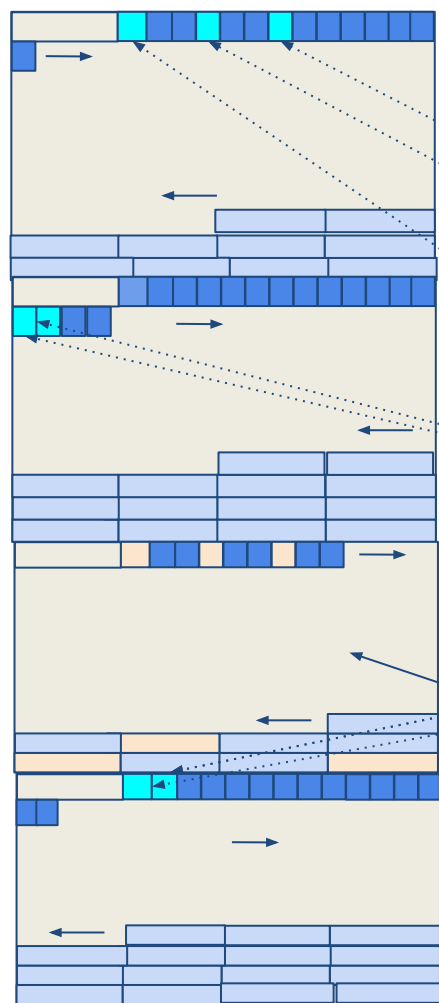


A three phase process: Phase 2





A three phase process: Phase 3



Phase 3:

Make ItemIDs space available for reuse for future tuples on page.

Phase 3

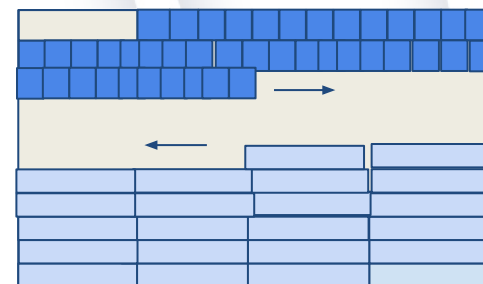
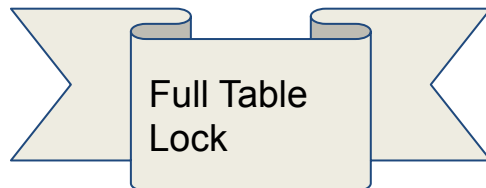
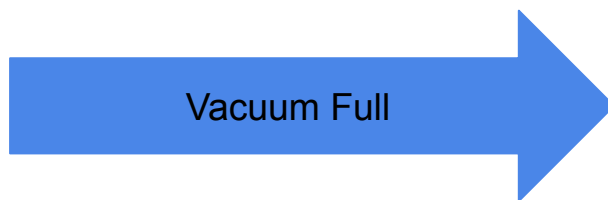
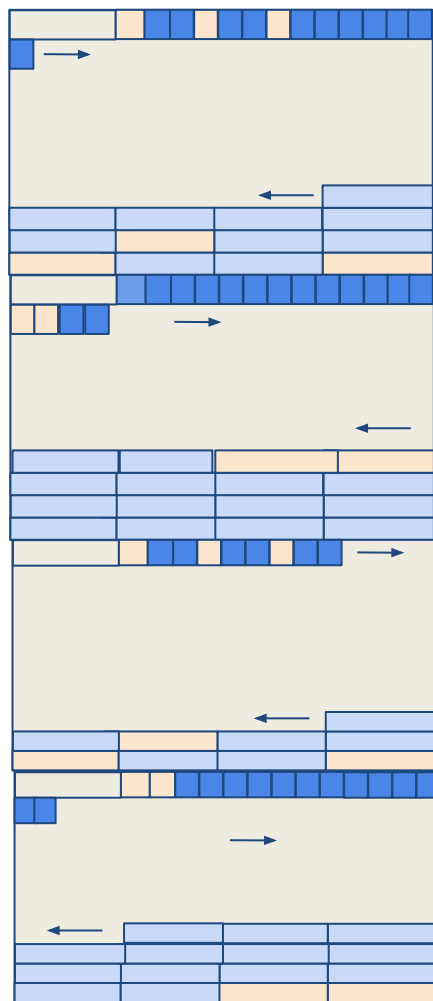
maintenance_work_mem

Page contained tuples that were visible to some transactions at time of vacuum.



Vacuum Full

- 1) Return space to operating system
- 2) Requires full table lock
- 3) Beware of dependent tables
- 4) It is just a rewrite of the table.



Space returned to the Operating System



Running Vacuum

- Raise default `maintenance_work_mem`
 - Default is generally too low
- Remove unused indexes
 - Check the server stats for index usage
 - Be sure to check standbys as well
- Let vacuum operations finish if you can
 - Minimal benefit from aborting a vacuum



Running Vacuum

1. Manually via vacuum command
 - a. Frequently run via cron jobs
 - b. Can vacuum specific tables or the database
2. Via auto vacuum workers
 - a. Run when certain events occur (thresholds crossed)
 - b. The preferred way of doing things ... if can



Autovacuum Workers

Vacuum
launcher
Process

```
[tomkincaid@tommac ~]$ psql -c 'select pid, (now() - pg_stat_activity.query_start) AS duration
M pg_stat_activity;' ^C
[tomkincaid@tommac ~]$ ps -ef | grep vacuum
postgres  7141   7020  0 Jun12 ?        00:00:13 postgres: autovacuum launcher process
tomkinc+  51889  85485  0 12:21 pts/7    00:00:00 grep --color=auto vacuum
tomkinc+ 109079   8796  0 Jun17 pts/3    00:00:00 vim autovacuum.c
[tomkincaid@tommac ~]$
```

Buffer
(maint_work_mem
in size)

Worker
Process

Accounts

Branches

Tellers

Buffer
(maint_work_mem
in size).

Worker
Process



Tuning Autovacuum

- When to start a vacuum on a Table?
- How long to vacuum before yielding?
- How long to yield?



Should auto-vacuum vacuum a table?

A table needs to be vacuumed if the number of dead tuples exceeds a threshold. This threshold is calculated as

$$\text{threshold} = \text{vac_base_thresh} + \text{vac_scale_factor} * \text{reltuples}$$

Percent dead_tuples vs. total_tuples + threshold

```
postgresql.conf
autovacuum_vacuum_scale_factor = 0.2
autovacuum_vacuum_threshold = 5
```

34 Total Tuples

19 Dead Tuples

$$5 + 0.2 * 34$$

Vacuum Threshold is: 12

19 is greater than 12

Autovacuum will run this table

[illegible]



Throttling Autovacuum

`autovacuum_cost_limit`

Work measured in cache hit and cache miss pages. Default is 200. Very low.

`autovacuum_cost_delay`

Sleep specified number of ms (default is very high).

`autovacuum_cost_limit`



Autovacuum and Workers

Enable auto_vacuum postgresql.conf

You configure the number of `auto_vacuum` workers in `postgresql.conf` setting `autovacuum max workers`.

When `autovacuum_naptime` expires vacuum determines if a worker should be launched based on state of tables and `postgresql.conf` settings.

[illegible]



Best Practices Thresholds

- Remember design first if you can
 - Often times it is too late
- Decide on how much bloat you are going to tolerate
 - A rule of thumb is 20%
- Almost always adjust default values to be more aggressive
 - Raise `vacuum_cost_limit`
 - Lower `vacuum_cost_delay`
 - Lower `vacuum_scale_factor`
- Attempt to avoid table specific vacuum thresholds
- Remember `vacuum_cost_limit` shared across all workers



Best Practices Thresholds

- Remember design first if you can
 - Often times it is too late
- Decide on how much bloat you are going to tolerate
 - A rule of thumb is 20%
- Almost always adjust default values to be more aggressive
 - Raise `vacuum_cost_limit`
 - Lower `vacuum_cost_delay`
 - Lower `vacuum_scale_factor`
- Attempt to avoid table specific vacuum thresholds



Some design principles around vacuum

- If possible, avoid having large update intensive tables
- Don't spuriously add indexes
- Consider a fill factor when creating table
 - More details coming in future talk
- Consider using Table Partitioning
- Will you have regular maintenance windows?
 - How frequent, how long will they be?
- Avoid long running transactions
- DDL operations



Best Practices Monitoring

- Table Bloat
- CPU Usage
- Disk I/O Rates
- Frozen Tuple Age and XID consumption
 - To be covered in a future talk
- Time since last analyze
 - Future talk
- Long Running Transactions



Vacuum Summary

- Understand
- Design
- Monitor
- Tune



Thank you!

tom.kincaid@2ndquadrant.com

andrew.dunstan@2ndquadrant.com

info@2ndquadrant.com

www.2ndquadrant.com



Thank you!

tom.kincaid@2ndquadrant.com

andrew.dunstan@2ndquadrant.com

info@2ndquadrant.com

www.2ndquadrant.com



A three phase process: Phase 1

9,12,33,'steve acct'
9,12,8,'steve acct'
9,12,75,'steve acct'
3,1,100,'tom acct'
9,12,00,'steve acct'
3,1,101,'tom acct'
3,1,100,'frank acct'
3,1,100,'joe acct'
3,1,4,'tom acct'
3,1,75,'tom acct'
3,1,21,'tom acct'
8,2,1200,'bob acct'
9,4,1500,'marc acct'
15,18,20,'marc acct'
3,1,21,'tom a'
3,1,21,'tom acct'
27,3,75,'ro acct'
35,12,8,'sue acct'
3,1,500,'tom acct'
32,12,8,'mary acct'
11,7,2000,''
3,1,21,'liz acct'
3,1,21,'frnk acct'
3,1,233,'tom acct'
3,1,21,'tom acct'
3,1,17,'tom acct'
3,1,12,'tom acct'
3,1,54,'tom acct'
3,1,75,'tom acct'

Phase 1

Phase 1:

Scan a table's heap to identify dead tuples, place item_ptrs in maintenance_work_mem buffer.

Mark rows as available for future updates

maintenance_work_mem



A three phase process: Phase 1

9,12,33,'steve acct'
9,12,8,'steve acct'
3,1,100,'tom acct'
8,2,1200,'bob acct'
9,4,1500,'marc acct'
15,18,20,'marc acct'
27,3,75,'ro acct'
35,12,8,'sue acct'
32,12,8,'mary acct'
11,7,2000,''
3,1,21,'liz acct'
3,1,21,'frnk acct'
3,1,233,'tom acct'
3,1,54,'tom acct'

Phase 1



Live Tuples



Rows than can be used for inserts and updates



Rows that could not be reclaimed (still visible to running transactions).



Vacuum Full

- 1) Return space to operating system
- 2) Requires full table lock
- 3) Beware of dependent tables
- 4) It is just a rewrite of the table.

Vacuum Full

Full Table
Lock

9,12,8,'steve acct'
3,1,100,'tom acct'
8,2,1200,'bob acct'
9,4,1500,'marc acct'
15,18,20,'marc acct'
27,3,75,'ro acct'
35,12,8,'sue acct'
32,12,8,'mary acct'
11,7,2000,''



A three phase process: Phase 3

Phase 3: Make the item ptr slots available for reuse.

