

Vienna 2025

Peter Hofer











11

Freedom has many difficulties & democracy is not perfect, but we have never had to put a wall up to keep our people in

John F. Kennedy, 1963



11

Freedom has many difficulties & democracy is not perfect, but we have never had to put a wall up to keep our people in

"

matches well to software

John F. Kennedy, 1963



Closed Source What does that really mean?



Bugs and errors numerous and disruptive



Knowledge Base What can one actually know?



The love of detail Purpose vs. paycheck





Bugs and errors

And how to deal with it



An example

```
1.DBMS_AUDIT_MGMT.SET_LAST_ARCHIVE_TIMESTAMP(
2.     audit_trail_type => DBMS_AUDIT_MGMT.AUDIT_TRAIL_UNIFIED,
3.     last_archive_time => TO_TIMESTAMP(:ts, 'YYYY-MM-DD HH24:MI:SS.FF')
4.);
5.DBMS_AUDIT_MGMT.CLEAN_AUDIT_TRAIL(
6.     audit_trail_type => DBMS_AUDIT_MGMT.AUDIT_TRAIL_UNIFIED,
7.     use_last_arch_timestamp => true
8.);
```

- So far, 4 different error messages...
- Identical systems with identical loads
- Occasionally, ancient data remains





Knowledge Base

What can one actually know?



"Let's ask the guru"

PostgreSQL World

- "I'll take a quick look"
- "That makes..."
- "Ah, I know that, wait a moment..."
- ... I probably shouldn't have asked ...

Commercial world

- "I don't know."
- "I guess that ..."
- "... never seen this before..."



11

The realization

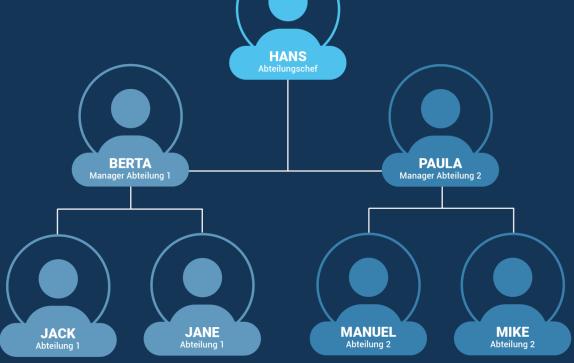
How was the colleague supposed to

know that, anyway?



What does that mean in real life?

- "Not knowing" = Research
- "Not knowing" = Delay
- "Not knowing" = Risk
- "Not knowing" = Potentially wrong
- "Not knowing" = Frustration
- "Not knowing" = Loss of trust







It can be done better

What does open source really mean?





```
What does this?
```

```
1.test=# SHOW
  effective_cache_size;
2. effective_cache_size
3.------
4. 4GB
5.(1 row)
```



Let's take a closer look...

```
1.hs@system:~/src/postgresql-17.3/src/backend$ grep -r -n -I -l effective_cache_size
    *
2.access/gist/gistbuild.c
3.optimizer/path/costsize.c
4.utils/misc/postgresql.conf.sample
5.utils/misc/guc_tables.c
```



GIST? What's happening there?

```
731 | /* subtree must fit in cache (with safety factor of 4) */
732 | if (subtreesize > effective_cache_size / 4)
733 | break;
```



Optimizer? Sounds interesting...

```
1.* costsize.c
            Routines to compute (and set) relation sizes and path costs
2. *
3....
4. *
          seq page cost
                                  Cost of a sequential page fetch
5. *
         random_page_cost
                                  Cost of a non-sequential page fetch
6. *
         cpu_tuple_cost
                                  Cost of typical CPU time to process a tuple
7. *
         cpu index tuple cost Cost of typical CPU time to process an index tuple
8. *
                                  Cost of CPU time to execute an operator or function
         cpu operator cost
9. *
          parallel tuple cost Cost of CPU time to pass a tuple from worker to leader backend
10. *
           parallel_setup_cost Cost of setting up shared memory for parallelism
11....
12. * We also use a rough estimate "effective cache size" of the number of
13. * disk pages in Postgres + OS-level disk cache. (We can't simply use
14. * NBuffers for this purpose because that would ignore the effects of
15. * the kernel's disk cache.)
16. *
17. * Obviously, taking constants for these values is an oversimplification,
18. * but it's tough enough to get any useful estimates even at this level of
19. * detail. Note that all of these parameters are user-settable, in case
20. * the default values are drastically off for a particular platform.
```



Optimizer? Sounds interesting...

```
1. * index_pages_fetched
            Estimate the number of pages actually fetched after accounting for
2. *
3. *
           cache effects.
4. *
5. * We use an approximation proposed by Mackert and Lohman, "Index Scans
6. * Using a Finite LRU Buffer: A Validated I/O Model", ACM Transactions
7. * on Database Systems, Vol. 14, No. 3, September 1989, Pages 401-424.
8. * The Mackert and Lohman approximation is that the number of pages fetched is
9. *
          PF =
10. *
                  min(2TNs/(2T+Ns), T)
                                                           when T <= b
                   2TNs/(2T+Ns)
                                                                   when T > b and Ns <= 2Tb/(2T-
11. *
 b)
                   b + (Ns - 2Tb/(2T-b))*(T-b)/T when T > b and Ns > 2Tb/(2T-b)
12. *
13. * where
       T = # pages in table, N = # tuples in table
14. *
        s = selectivity = fraction of table to be scanned, b = # buffer pages available
15. *
16. *
17. * We assume that effective_cache_size is the total number of buffer pages
18. * available for the whole query, and pro-rate that space across all the
19. * tables in the guery and the index currently under consideration. (This
20. * ignores space needed for other indexes used by the query, but since we
21. * don't know which indexes will get used, we can't estimate that very well;
22. * and in any case counting all the tables may well be an overestimate, since
23. * depending on the join plan not all the tables may be scanned concurrently.)
```



Optimizer? Sounds interesting...

```
1.
          * Estimate number of main-table pages fetched, and compute I/O cost.
2.
          * When the index ordering is uncorrelated with the table ordering,
          * we use an approximation proposed by Mackert and Lohman (see
5.
          * index_pages_fetched() for details) to compute the number of pages
6.
          * fetched, and then charge spc random page cost per page fetched.
7.
8.
          * When the index ordering is exactly correlated with the table ordering
9.
           * (just after a CLUSTER, for example), the number of pages fetched should
10.
           * be exactly selectivity * table_size.
11.
```

- Yes, the example is simple
- Yes, you still need to know something about the topic
- But: Users and support have a chance!





Love for details

To improve life for everyone



Example: JSON support

PostgreSQL

- early adopter,
- JSON in 2012
- JSONB in 2014.

Oracle

- JSON in 2014,
- native binary JSON in 2021

Oracle entered the JSON world in 2014, but its native binary JSON type only arrived in 2021 – seven years after PostgreSQL's JSONB



Recently: more than 1000 words

Note: nobody does such a thing





In a nutshell

a few thoughts

Any questions?

Ask anything





Our Partners at PGDay Austria































Your Pathway to Verified PostgreSQL Skills

Scan for Updates



oapg-edu.org