

DBMS Performance Evaluation - DS- Project-I

Context

DBMS can help us manage data conveniently and significantly improve the efficiency of data retrieval.

PostgreSQL is a popular open-source RDBMS known for its robustness, advanced features, and strong compliance with SQL standards. openGauss is an enterprise-grade open-source RDBMS developed by Huawei, designed for high performance, security, and scalability in demanding business environments.

Preface

| Latest Report

[Tsawke's Blog](#)

[Github](#)

| Environment

System: Alibaba Cloud Linux 3.2104 LTS 64位

IP: 47.115.128.238

openGauss (docker)

- Version: openGauss-Docker-6.0.2-x86_64
- Port: 15432:5432
- User: omm (Operation & Maintenance Manager)
- Password: opengauss

PostgreSQL

- Version: 17.6
- Port: 5432
- User: postgres
- Password: postgres

Datasets

Clickstream

The Clickstream dataset aggregates counts of **(referrer, resource)** pairs from Wikipedia request logs, showing how readers arrive at an article and what they click next.

Run **DownloadDatasets.sh** to download all the datasets.

To import datasets to PostgreSQL, run:

```
1 | psql "postgresql://postgres:postgres@127.0.0.1:5432/project1" -f  
   './ImportDatasets.sql'
```

Results:

```
1 rows_loaded
2 -----
3      35512282
```

What are the unique advantages of a DBMS compared with data operations in files?

Purpose: Find all events that contains 'main' in 'curr'(current).

PostgreSQL

```
1 EXPLAIN (ANALYZE)
2 SELECT * FROM clickstream.events WHERE curr ILIKE '%main%';
```

```
1 QUERY PLAN
2 -----
3 Gather (cost=1000.00..547907.69 rows=124279 width=47) (actual
4 time=0.437..24740.854 rows=162828 loops=1)
5   Workers Planned: 2
6   Workers Launched: 2
7   -> Parallel Seq Scan on events (cost=0.00..534479.79 rows=51783
8     width=47) (actual time=2.401..24640.354 rows=54276 loops=3)
9     Filter: (curr ~* '%main% '::text)
10     Rows Removed by Filter: 11783151
11 Planning Time: 1.428 ms
    Execution Time: 24750.391 ms
(8 rows)
```

During 5 tests, the average result is **24662.1ms**.

C++

```
1 | g++ ./SelectAll.cpp -o SelectAll -std=c++17 -O2 && ./SelectAll
```

```
1 | Find 162828 results in total
2 | 21234.4 ms
```

During 5 tests, the average result is **20116.2ms**.

Purpose: Update every '_' in 'curr' to '^'.

PostgreSQL

```
1 | EXPLAIN (ANALYZE)
2 | UPDATE clickstream.events SET curr = REPLACE(curr, '_', '^') WHERE curr
   | LIKE '%_%';
```

```
1 | QUERY PLAN
2 | -----
3 | Update on events (cost=0.00..882401.93 rows=0 width=0) (actual
   | time=225011.949..225014.086 rows=0 loops=1)
4 |   -> Seq Scan on events (cost=0.00..882401.93 rows=35522510
   | width=38) (actual time=0.024..41868.876 rows=35512282 loops=1)
5 |     Filter: (curr ~~ '%_%'::text)
6 | Planning Time: 10.490 ms
7 | Execution Time: 225019.108 ms
8 | (5 rows)
```

The result is **225019.1ms**.

C++

```
1 | g++ ./UpdateAll.cpp -o UpdateAll -std=c++17 -O2 && ./UpdateAll
```

```
1 | Update 30866250 results in total
2 | 61515.7 ms
```

The result is **61515.7ms**.

Purpose: Find Top-K Popular Pages

PostgreSQL

```
1 | SET search_path = clickstream, public;
2 | EXPLAIN (ANALYZE)
3 | SELECT curr, SUM(n) AS clicks
4 |     FROM events
5 |     GROUP BY curr
6 |     ORDER BY clicks DESC
7 |     LIMIT 20;
```

```
1 | ...
2 | Planning Time: 9.009 ms
3 | Execution Time: 101260.848 ms
```

The result is **101260.8ms**.

C++

```
1 | g++ ./SelectTopK.cpp -o SelectTopK -std=c++17 -O2 && ./SelectTopK
```

```
1 | 69683.3 ms
```

The result is **69683.3ms**.

Conclusion

Overall, C++ streaming program **beat** the DBMS(PostgreSQL) on all tests.

And there're multiple reasons why DBMS seems to be slower than C++:

- All tasks are **one-shot full-scan or rewriting**, thus C++ can lightly and easily streams files, avoiding DBMS overheads.
- We are using **basic DBMS** without optimization like pg_trgm, B-Tree. (The efficiency won't increase too much even with pg_trgm.)

Thus, the results **doesn't negate** the DBMS strengths.

At some circumstances include **reusable queries, concurrency, strong consistency, complex joins/transactions**, e.t.c., DBMS will perform significantly **better** than data operations in files.

Which DBMS is better? PostgreSQL or openGauss, and by which standard?

Preparation

```
1 | SET max_parallel_workers_per_gather = 4;  
2 | SET work_mem = '256MB';
```

Comparison of Select

```
1 | EXPLAIN (ANALYZE)  
2 | SELECT * FROM clickstream.events WHERE curr ILIKE '%main%';
```

Results:

PostgreSQL: **24662.1ms**.

openGauss: **27012.6ms**.

Comparison of Update

```
1 | EXPLAIN (ANALYZE)
2 | SELECT * FROM clickstream.events WHERE curr ILIKE '%main%';
```

Results:

PostgreSQL: **225019.1ms**.

openGauss: **244217.9ms**.

Comparison of Table Join

```
1 | EXPLAIN (ANALYZE)
2 | SELECT a.prev AS src, b.curr AS dst, COUNT(*) AS paths
3 |     FROM clickstream.events a
4 |     JOIN clickstream.events b ON a.curr = b.prev
5 |     GROUP BY a.prev, b.curr
6 |     ORDER BY paths DESC
7 |     LIMIT 20;
```

Results:

PostgreSQL: **180567.2ms**.

openGauss: **191124.3ms**.

Comparison of Top-K Query

```
1 SET search_path = clickstream, public;  
2 EXPLAIN (ANALYZE)  
3 SELECT curr, SUM(n) AS clicks  
4     FROM events  
5     GROUP BY curr  
6     ORDER BY clicks DESC  
7     LIMIT 20;
```

Results:

PostgreSQL: **97628.7ms**.

openGauss: **115431.5ms**.

Conclusion

Overall, we can conclude that PostgreSQL and openGauss have **almost the same efficiency**. They are very **similar** and openGauss is usually **a little bit slower** than PostgreSQL. Therefore, at some specific circumstances, like **looser settings**, openGauss will perhaps **performs better**, but still **a little**.

Remarks

How to connect?

PostgreSQL

```
1 sudo -u postgres psql -p 5432 -d postgres
```

```
1 psql -h 127.0.0.1 -p 5432 -U postgres -d postgres
```



```
1 | psql "postgresql://postgres:postgres@127.0.0.1:5432/project1"
```

openGauss

```
1 | docker exec -e PGPASSWORD='opengauss' -u omm opengauss15432 \  
2 |   bash -lc "gsql -h 127.0.0.1 -p 5432 -U omm -d postgres"
```

```
1 | gsql -h 127.0.0.1 -p 15432 -d postgres -U omm
```