

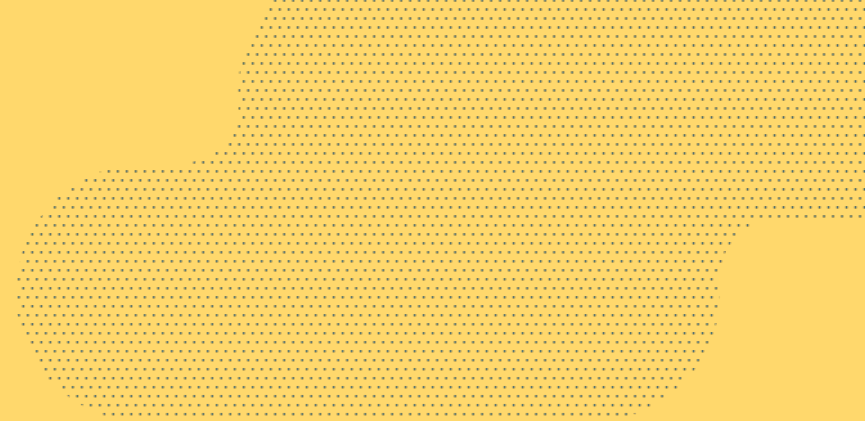


How `pg_query` rewrites, deparses & formats any valid Postgres query

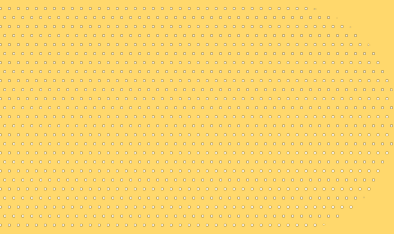
* Presented by Keiko Oda & Lukas Fittl

Agenda for today

- What is pg_query?
- Query Rewriting in Practice
- Formatting any valid Postgres Query



What is pg_query?



What is pg_query?



Postgres parser as a standalone library

First launched with Postgres 10, continuously supported new releases since PG 13.

libpg_query

C library for accessing the PostgreSQL parser outside of the server.

This library uses the actual PostgreSQL server source to parse SQL queries and return the internal PostgreSQL parse tree.

```
SELECT * FROM mytable
```

```
SelectStmt
├─ targetList
│  └─ ResTarget
│     └─ ColumnRef
│        └─ A_Star
├─ fromClause
```

```
DELETE FROM abc WHERE id = 1
```

```
DeleteStmt
├─ relation
│  └─ RangeVar
│     └─ relname = abc
├─ whereClause
│  └─ A_Expr
```

What is pg_query?



Postgres parser as a standalone library

```
#include <pg_query.h>
#include <stdio.h>

int main() {
    PgQueryParseResult result;

    result = pg_query_parse("SELECT 1");

    printf("%s\n", result.parse_tree);

    pg_query_free_parse_result(result);

    return 0;
}
```

What is pg_query?

Postgres parser as a standalone library

```
{
  "version": 170007,
  "stmts": [
    {
      "stmt": {
        "SelectStmt": {
          "targetList": [
            {
              "ResTarget": {
                "val": {
                  "A_Const": {
                    "val": {
                      "Integer": {
                        "ival": 1
                      }
                    },
                    "location": 7
                  }
                },
                "location": 7
              }
            }
          ],
          "limitOption": "LIMIT_OPTION_DEFAULT",
          "op": "SETOP_NONE"
        }
      }
    }
  ]
}
```

What is pg_query?



Or in Ruby:

```
require 'pg_query'  
PgQuery.parse("SELECT 1")
```

Or in Rust:

```
use pg_query::parse;  
pg_query::parse("SELECT * FROM contacts");
```

Or in Python (via pglast binding):

```
from pglast import parse_sql  
root = parse_sql('select 1')
```

Why pg_query?



Applications need to parse SQL with full Postgres compatibility:

At **pganalyze** we needed to parse + fingerprint Postgres queries from `pg_stat_statements` and the logs. Other projects had other needs:

Postgres Language Server: Autocomplete, syntax checks & more in your editor

SQLlint: Simple SQL linter supporting ANSI and PostgreSQL syntaxes

pgspot: Spot vulnerabilities in postgres SQL scripts

pgDog: Proxy that can rewrite queries/route them to different shards

Fun Fact: The **DuckDB** parser originated from `pg_query`, before they forked it to expand their syntax.

pg_query uses Protobufs to have language-specific structs that can be more easily accessed:

```
q = PgQuery.parse('SELECT * FROM x')
q.tree
#=> <PgQuery::ParseResult: version: 180004, stmts: [
  <PgQuery::RawStmt: stmt: <PgQuery::Node: select_stmt:
    <PgQuery::SelectStmt: distinct_clause: [],
      target_list: [
        <PgQuery::Node: res_target: <PgQuery::ResTarget: name: "", indirection: [],
          val: <PgQuery::Node: column_ref: <PgQuery::ColumnRef: fields: [<PgQuery::Node:
a_star: <PgQuery::A_Star: >>], location: 7>>, location: 7>>],
        from_clause: [<PgQuery::Node: range_var: <PgQuery::RangeVar: catalogname: "",
          schemaname: "", relname: "x", inh: true, relpersistence: "p", location: 14>>],
          ...
x.tree.stmts[0].stmt.select_stmt.from_clause
#=> [<PgQuery::Node: range_var: <PgQuery::RangeVar: catalogname: "", schemaname: "",
  relname: "x", inh: true, relpersistence: "p", location: 14>>]
```

pg_query has a deparser!



Deparsing means we turn a parse tree (AST) into SQL again

Feature complete, meaning it can turn 100% of the Postgres regression tests back into SQL, including utility statements.

`SELECT * FROM x;` ➡ `.. <PgQuery::RangeVar: relname: "x",
inh: true, relpersistence: "p"> ..` ➡ `SELECT * FROM x;`

Separate C file from the main pg_query code. It could be imported into Postgres if there was an in-core use case for raw parse tree deparsing.

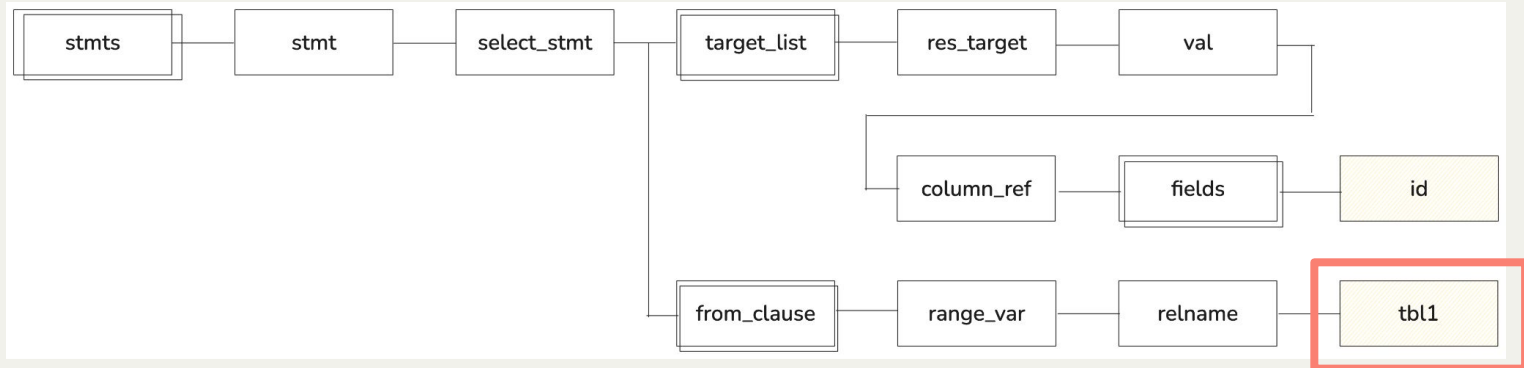
Query Rewriting in Practice

Query rewrite 101: Update Table Name

- **Goal:** Update the table name from "tbl1" to "tbl2"
- Use "**PgQuery.parse**" to turn SQL into a parse tree, then drill down to the field

```
require "pg_query"
parsed_query = PgQuery.parse("SELECT id FROM tbl1")
parsed_query.tree.stmts[0].stmt.select_stmt.from_clause[0].range_var.relname
#=> "tbl1"
```

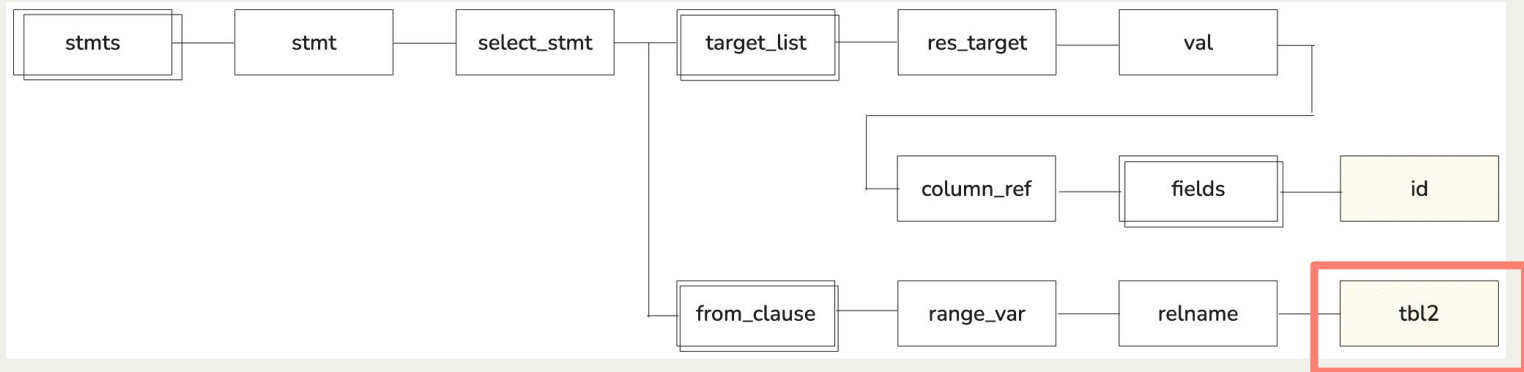
example uses the Ruby bindings



Query rewrite 101: Update Table Name

- Assign a new value to the field, then call "**#deparse**" to get SQL back
- The full pattern: **parse** → **modify** → **deparse**

```
parsed_query #⇒ input: "SELECT id FROM tbl1"  
parsed_query.tree.stmts[0].stmt.select_stmt.from_clause[0].range_var.relname = "tbl2"  
parsed_query.deparse  
#⇒ "SELECT id FROM tbl2"
```



Rewrite OR to UNION to enable index scans across JOINS



- OR on columns from different tables → JOIN runs first, then filter

```
SELECT o.id, o.total, c.email
FROM orders o JOIN customers c
  ON c.id = o.customer_id
WHERE
  o.flagged IS TRUE
OR
  c.is_blocked IS TRUE

# orders: 1M rows
(~1k rows with flagged: true)

# customers: 100k rows
(~500 rows with is_blocked: true)
```

Plan

- 1 Hash Join
Hash Cond: (o.customer_id = c.id)
Inner Unique: true
Join Filter: (o.flagged OR c.is_blocked)
Join Type: Inner
Rows Removed by Join Filter: 994294
- 2 Seq Scan on orders
- 3 Hash
- 4 Seq Scan on customers

Rewrite OR to UNION to enable index scans across JOINS



- By rewriting OR to UNION, orders and customers can each use its own index

```
SELECT o.id, o.total, c.email
FROM orders o JOIN customers c
  ON c.id = o.customer_id
WHERE o.flagged IS TRUE

UNION

SELECT o.id, o.total, c.email
FROM orders o JOIN customers c
  ON c.id = o.customer_id
WHERE c.is_blocked IS TRUE
```

Plan

```
1   Aggregate
2   Append
3     Nested Loop
4       Index Scan (Forward) on orders
5       Index Scan (Forward) on customers
6     Nested Loop
7       Index Scan (Forward) on customers
8       Index Scan (Forward) on orders
```

Rewrite OR to UNION: Plan Comparison

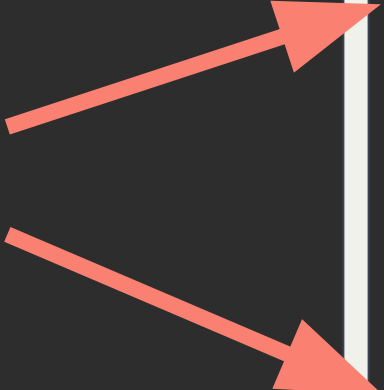


Plan A	Plan B	Plan A	Plan B
Using OR - Parameter Set 1	Using UNION - Parameter Set 1	Runtime	Runtime
-> Hash Join		177.76ms	
-> Seq Scan on orders		51.28ms	
-> Hash		32.84ms	
-> Seq Scan on customers		16.24ms	
	-> Aggregate		5.07ms
	-> Append		0.64ms
	-> Nested Loop		1.05ms
	-> Index Scan ⁴ on orders		5.17ms
	-> Index Scan ² on custome...		8.36ms
	-> Nested Loop		1.22ms
	-> Index Scan ¹ on custome...		0.74ms
	-> Index Scan ³ on orders		12.93ms

278.34ms 35.74ms

Rewrite OR to UNION: Query Comparison

```
SELECT o.id, o.total, c.email
FROM orders o JOIN customers c
  ON c.id = o.customer_id
WHERE
  o.flagged IS TRUE
OR
  c.is_blocked IS TRUE
```

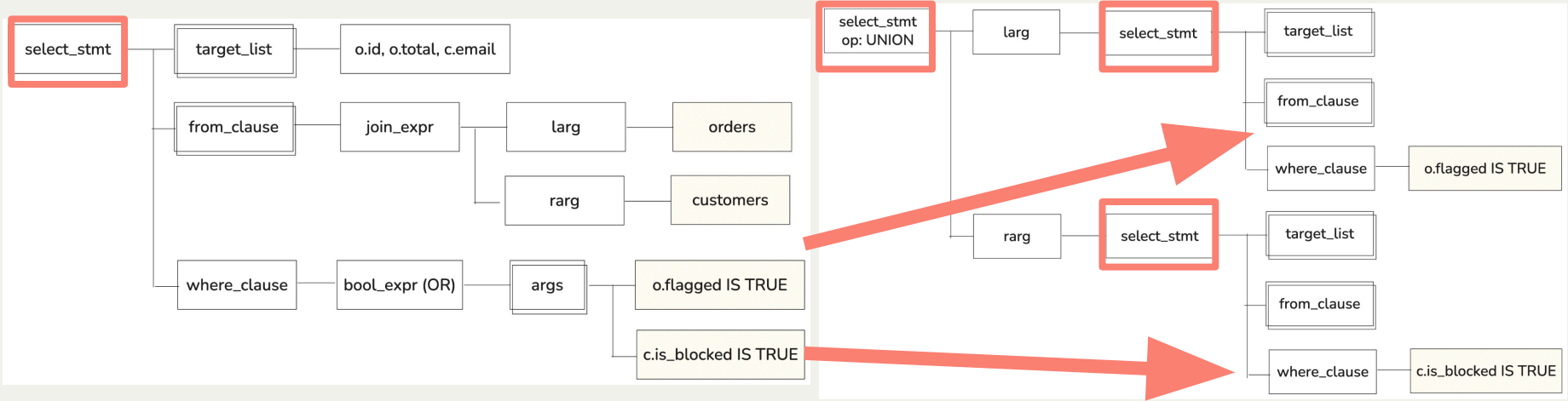


```
SELECT o.id, o.total, c.email
FROM orders o JOIN customers c
  ON c.id = o.customer_id
WHERE o.flagged IS TRUE

UNION

SELECT o.id, o.total, c.email
FROM orders o JOIN customers c
  ON c.id = o.customer_id
WHERE c.is_blocked IS TRUE
```

Rewrite OR to UNION: Parse Tree Comparison



Rewrite OR to UNION to enable index scans across JOINS



```
# 1. Walk the tree to find the SelectStmt to rewrite, and remember where it lives
target_select = nil
target_loc = nil
parsed_query.walk! do |_, _, node, loc| # walk! visits every node in the tree
  next unless is_target_select?(node)
  target_select = node
  target_loc = loc
end
# 2. Walk the subtree to find the OR node
or_node = nil
parsed_query.walk_subtree!(target_select) do |node| # visits every node in the subtree
  or_node = node if node.is_a?(PgQuery::BoolExpr) && node.boolop == :OR_EXPR
end
```

Rewrite OR to UNION to enable index scans across JOINS



```
# 3. Turn each OR branch into its own SelectStmt
left_select = Google::Protobuf.deep_copy(target_select)
left_select.where_clause = or_node.args[0]
right_select = Google::Protobuf.deep_copy(target_select)
right_select.where_clause = or_node.args[1]

# 4. Combine the branches with UNION and write it back to the recorded location
union = PgQuery::SelectStmt.new(larg: left_select, rarg: right_select, op: :SETOP_UNION)
parsed_query.find_tree_location(parsed_query.tree, target_loc) do |parent_node, _, _|
  parent_node.select_stmt = union
end
```

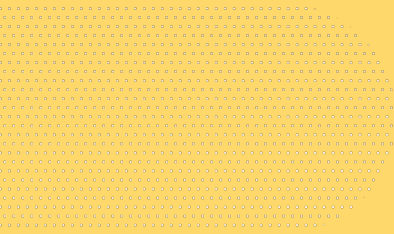
```
# 5. Deparse
```

```
parsed_query.deparse
```

```
#⇒ "SELECT o.id, o.total, c.email FROM orders o JOIN customers c ON c.id = o.customer_id  
WHERE o.flagged IS TRUE UNION SELECT o.id, o.total, c.email FROM orders o JOIN customers  
c ON c.id = o.customer_id WHERE c.is_blocked IS TRUE"
```



Formatting any valid Postgres Query



- **Problem:** Rewrite diff is messy and unclear what changed
- **Solution:** Teach deparser to format before/after queries

```
1 SELECT databases.id, s.reset_value::bigint FROM databases JOIN
  postgres_settings s USING (server_id) WHERE databases.id =
  ANY($1::bigint[]) AND s.name = $2 AND s.removed_at IS NULL
```

```
1 WITH __postgres_settings AS MATERIALIZED (SELECT s.reset_value,
  s.server_id FROM postgres_settings s WHERE s.name = $2) SELECT
  databases.id, s.reset_value::bigint FROM databases JOIN
  __postgres_settings s USING (server_id) WHERE databases.id =
  ANY($1::bigint[]) AND s.removed_at IS NULL
```

```
1 SELECT databases.id, s.reset_value::bigint
2 FROM
3   databases
4 JOIN postgres_settings s USING (server_id)
5 WHERE
6   databases.id = ANY($1::bigint[])
7   AND s.name = $2
8   AND s.removed_at IS NULL
9
```

```
1 WITH __postgres_settings AS MATERIALIZED (
2   SELECT s.reset_value, s.server_id
3   FROM postgres_settings s
4   WHERE s.name = $2
5 )
6 SELECT databases.id, s.reset_value::bigint
7 FROM
8   databases
9 JOIN __postgres_settings s USING (server_id)
10 WHERE
11   databases.id = ANY($1::bigint[])
12   AND s.removed_at IS NULL
13
```

- Pass options to deparse method

```
opts = PgQuery::DeparseOpts.new(  
  pretty_print: true,           # default: false  
  indent_size: 2,              # default: 4  
  max_line_length: 80,         # default: 80  
  trailing_newline: true,      # default: false  
  commas_start_of_line: false # default: false  
)  
parsed_query.deparse(opts: opts)
```

- Comments can be preserved

```
query_text = "/* my comment */\nSELECT /* my inline comment */1"
parsed_query = PgQuery.parse(query_text)
opts = PgQuery::DeparseOpts.new(...)
parsed_query.deparse(opts: opts)
#⇒ "SELECT 1\n"
opts.comments = PgQuery.deparse_comments_for_query(query_text)
parsed_query.deparse(opts: opts)
#⇒ "/* my comment */\nSELECT /* my inline comment */1\n"
```

- Available on <https://pganalyze.com/format>

pganalyze Query Formatter

```
1 SELECT i.inventory_id, s.store_id
2 FROM
3   store s
4   JOIN inventory i ON i.store_id = s.store_id
5 WHERE
6   st_dwithin(s.geog, (
7     SELECT geog
8     FROM customer
9     WHERE customer_id = $1
10  ), $2)
11 AND i.status_id = $3
12 AND NOT EXISTS (
13   SELECT $4
14   FROM rental r
15   WHERE
16     r.inventory_id = i.inventory_id
17     AND upper(r.rental_period) IS NULL
18  )
19 )
```

Formatted!

FORMAT

About the Query Formatter

The pganalyze Query Formatter utilizes [libpg_query](#) - a library developed by pganalyze to parse, deparse and format SQL queries.

We developed the formatter capabilities for [pganalyze Workbooks](#) to aid in query benchmarking and diffing, and are making it publicly available for anyone to use.

All query formatting happens locally in your browser (using WASM) and queries are not shared with pganalyze or stored.

Based on PostgreSQL 17.7



Thank you!

Learn more at github.com/pganalyze/libpq_query

Contact us: lukas@pganalyze.com, keiko@pganalyze.com