



# Architektur und Implementierung von Datenbanksystemen

Task 6 – recursive query execution in PostgreSQL

Team 3 - Gründlinger Diana, Huber Marcel, Klotz Thomas, Targa Aaron, Thalmann Matthias

# Content

- Syntax of a recursive query
- Execution plan
- Operators used and execution control flow

## Syntax of a recursive query

```
WITH RECURSIVE RecRel (attr1, attr2, ..., attrn)  
  AS (SFW-Statement1           -- Rekursionsinitialisierung  
    UNION ALL  
    SFW-Statement2           -- Rekursionsschritt  
  )  
  
SELECT [ DISTINCT ] attributliste  
  FROM RecRel  
  [ WHERE ... ] [ GROUP BY ... ] [ HAVING ... ] [ ORDER BY ... ]
```

(2)

# Syntax of a recursive query

- WITH → common table expression (CTE), basically defining a temporary table that just exists for one query
- Optional RECURSIVE modifier → WITH query can refer to its own output

General Form of a recursive WITH query:

- A non-recursive term
- UNION or UNION ALL
- A recursive term

```
WITH RECURSIVE RecRel (attr1, attr2, ..., attrn)  
  AS (SFW-Statement1           -- Rekursionsinitialisierung  
      UNION ALL  
      SFW-Statement2           -- Rekursionsschritt  
  )  
SELECT [ DISTINCT ] attributliste  
FROM RecRel  
[ WHERE ... ] [ GROUP BY ... ] [ HAVING ... ] [ ORDER BY ... ]
```

Only the recursive term can reference the query's own output.

(1,2)

# Execution plan - example

```
Query:
WITH RECURSIVE recursive_query(id, child_id) AS (
  SELECT id, child_id
  FROM my_tree
  WHERE id = 0
  UNION
  SELECT t.id, t.child_id
  FROM my_tree t, recursive_query
  WHERE recursive_query.child_id = t.id
)
SELECT * FROM recursive_query
```

Explain output:

```
CTE Scan on recursive_query (cost=113.31..115.33 rows=101 width=8) (actual time=0.686..1257.405 rows=500 loops=1)
  CTE recursive_query
    -> Recursive Union (cost=0.27..113.31 rows=101 width=8) (actual time=0.680..1254.706 rows=500 loops=1)
      -> Index Scan using idx on my_tree (cost=0.27..8.29 rows=1 width=8) (actual time=0.664..0.670 rows=1 loops=1)
        Index Cond: (id = 0)
      -> Hash Join (cost=0.33..10.30 rows=10 width=8) (actual time=1.301..2.487 rows=1 loops=500)
        Hash Cond: (t.id = recursive_query_1.child_id)
        -> Seq Scan on my_tree t (cost=0.00..8.00 rows=500 width=8) (actual time=0.009..1.225 rows=500 loops=500)
        -> Hash (cost=0.20..0.20 rows=10 width=4) (actual time=0.013..0.013 rows=1 loops=500)
          Buckets: 1024 Batches: 1 Memory Usage: 9kB
          -> WorkTable Scan on recursive_query recursive_query_1 (cost=0.00..0.20 rows=10 width=4) (actual
time=0.003..0.005 rows=1 loops=500)
Planning time: 0.688 ms
Execution time: 1259.438 ms
```

(3)

# Execution plan - example

```
Query:
WITH RECURSIVE recursive_query(id, child_id) AS (
  SELECT id, child_id
  FROM my_tree
  WHERE id = 0
  UNION
  SELECT t.id, t.child_id
  FROM my_tree t, recursive_query
  WHERE recursive_query.child_id = t.id
)
SELECT * FROM recursive_query
```

Explain output:

```
CTE Scan on recursive_query (cost=113.31..115.33 rows=101 width=8) (actual time=0.686..1257.405 rows=500 loops=1)
  CTE recursive_query
    -> Recursive Union (cost=0.27..113.31 rows=101 width=8) (actual time=0.680..1254.706 rows=500 loops=1)
      -> Index Scan using idx on my_tree (cost=0.27..8.29 rows=1 width=8) (actual time=0.664..0.670 rows=1 loops=1)
        Index Cond: (id = 0)
      -> Hash Join (cost=0.33..10.30 rows=10 width=8) (actual time=1.301..2.487 rows=1 loops=500)
        Hash Cond: (t.id = recursive_query_1.child_id)
        -> Seq Scan on my_tree t (cost=0.00..8.00 rows=500 width=8) (actual time=0.009..1.225 rows=500 loops=500)
        -> Hash (cost=0.20..0.20 rows=10 width=4) (actual time=0.013..0.013 rows=1 loops=500)
          Buckets: 1024 Batches: 1 Memory Usage: 9kB
          -> WorkTable Scan on recursive_query recursive_query_1 (cost=0.00..0.20 rows=10 width=4) (actual
time=0.003..0.005 rows=1 loops=500)
Planning time: 0.688 ms
Execution time: 1259.438 ms
```

(3)

# Execution plan - example

```
Query:
WITH RECURSIVE recursive_query(id, child_id) AS (
  SELECT id, child_id
  FROM my_tree
  WHERE id = 0
  UNION
  SELECT t.id, t.child_id
  FROM my_tree t, recursive_query
  WHERE recursive_query.child_id = t.id
)
SELECT * FROM recursive_query
```

Explain output:

```
CTE Scan on recursive_query (cost=113.31..115.33 rows=101 width=8) (actual time=0.686..1257.405 rows=500 loops=1)
  CTE recursive_query
    -> Recursive Union (cost=0.27..113.31 rows=101 width=8) (actual time=0.680..1254.706 rows=500 loops=1)
      -> Index Scan using idx on my_tree (cost=0.27..8.29 rows=1 width=8) (actual time=0.664..0.670 rows=1 loops=1)
        Index Cond: (id = 0)
      -> Hash Join (cost=0.33..10.30 rows=10 width=8) (actual time=1.301..2.487 rows=1 loops=500)
        Hash Cond: (t.id = recursive_query_1.child_id)
        -> Seq Scan on my_tree t (cost=0.00..8.00 rows=500 width=8) (actual time=0.009..1.225 rows=500 loops=500)
        -> Hash (cost=0.20..0.20 rows=10 width=4) (actual time=0.013..0.013 rows=1 loops=500)
          Buckets: 1024 Batches: 1 Memory Usage: 9kB
          -> WorkTable Scan on recursive_query recursive_query_1 (cost=0.00..0.20 rows=10 width=4) (actual
time=0.003..0.005 rows=1 loops=500)
Planning time: 0.688 ms
Execution time: 1259.438 ms
```

(3)

# Execution plan - example

```
Query:
WITH RECURSIVE recursive_query(id, child_id) AS (
  SELECT id, child_id
  FROM my_tree
  WHERE id = 0
  UNION
  SELECT t.id, t.child_id
  FROM my_tree t, recursive_query
  WHERE recursive_query.child_id = t.id
)
SELECT * FROM recursive_query
```

Explain output:

```
CTE Scan on recursive_query (cost=113.31..115.33 rows=101 width=8) (actual time=0.686..1257.405 rows=500 loops=1)
  CTE recursive_query
    -> Recursive Union (cost=0.27..113.31 rows=101 width=8) (actual time=0.680..1254.706 rows=500 loops=1)
      -> Index Scan using idx on my_tree (cost=0.27..8.29 rows=1 width=8) (actual time=0.664..0.670 rows=1 loops=1)
        Index Cond: (id = 0)
      -> Hash Join (cost=0.33..10.30 rows=10 width=8) (actual time=1.301..2.487 rows=1 loops=500)
        Hash Cond: (t.id = recursive_query_1.child_id)
        -> Seq Scan on my_tree t (cost=0.00..8.00 rows=500 width=8) (actual time=0.009..1.225 rows=500 loops=500)
        -> Hash (cost=0.20..0.20 rows=10 width=4) (actual time=0.013..0.013 rows=1 loops=500)
          Buckets: 1024 Batches: 1 Memory Usage: 9kB
          -> WorkTable Scan on recursive_query recursive_query_1 (cost=0.00..0.20 rows=10 width=4) (actual
time=0.003..0.005 rows=1 loops=500)
Planning time: 0.688 ms
Execution time: 1259.438 ms
```

(3)



# Operators used

- CTE Scan
- Recursive Union
- WorkTable Scan

# Operators used and execution control flow

- **CTE Scan** – Scans results of a CTE as temporary table
- **WorkTable Scan** - Scans the work table used in evaluating a recursive CTE
- **Recursive Union** - Returns the union of the recursive and non recursive subplan

(4)

# Operators used and execution control flow

1. evaluate non-recursive term  
for UNION discard duplicate rows  
include remaining rows in result of recursive query, and also place them in a temporary working table
2. while working table not empty, repeat:
  - a) Evaluate recursive term, substituting current contents of working table for recursive self-reference.  
For UNION (but not UNION ALL), discard duplicate rows and rows that duplicate any previous result row.  
Include all remaining rows in result of the recursive query, and also place them in temporary intermediate table.
  - b) Replace contents of working table with contents of intermediate table, then empty intermediate table.

Execution terminates when no further tuples are added in a recursive step.

Note: It's called RECURSIVE but the process is really iterative.

(5)

```

56  /* -----
57  *      ExecRecursiveUnion(node)
58  *
59  *      Scans the recursive query sequentially and returns the next
60  *      qualifying tuple.
61  *
62  * 1. evaluate non recursive term and assign the result to RT
63  *
64  * 2. execute recursive terms
65  *
66  * 2.1 WT := RT
67  * 2.2 while WT is not empty repeat 2.3 to 2.6. if WT is empty returns RT
68  * 2.3 replace the name of recursive term with WT
69  * 2.4 evaluate the recursive term and store into WT
70  * 2.5 append WT to RT
71  * 2.6 go back to 2.2
72  * -----
73  */
74 static TupleTableSlot *
75 ExecRecursiveUnion(PlanState *pstate)
76 {
77     RecursiveUnionState *node = castNode(RecursiveUnionState, pstate);
78     PlanState *outerPlan = outerPlanState(node);
79     PlanState *innerPlan = innerPlanState(node);
80     RecursiveUnion *plan = (RecursiveUnion *) node->ps.plan;
81     TupleTableSlot *slot;
82     bool        isNew;
83
84     CHECK_FOR_INTERRUPTS();
85
86     /* 1. Evaluate non-recursive term */
87     if (!node->recurring)
88     {
89         for (;;)
90         {
91             slot = ExecProcNode(outerPlan);
92             if (TupIsNull(slot))
93                 break;
94
95             if (plan->numCols > 0)
96             {
97                 /* Find or build hashtable entry for this tuple's group */
98                 LookupTupleHashEntry(node->hashtable, slot, &isNew, NULL);
99                 /* Must reset temp context after each hashtable lookup */
100                MemoryContextReset(node->tempContext);
101                /* Ignore tuple if already seen */
102                if (!isNew)
103                    continue;
104            }
105            /* Each non-duplicate tuple goes to the working table ... */
106            tuplestore_puttupleslot(node->working_table, slot);
107            /* ... and to the caller */
108            return slot;
109        }
110        node->recurring = true;
111    }
112
113    /* 2. Execute recursive term */
114    for (;;)
115    {
116        slot = ExecProcNode(innerPlan);
117        if (TupIsNull(slot))
118        {
119            /* Done if there's nothing in the intermediate table */
120            if (node->intermediate_empty)
121                break;
122
123            /* done with old working table ... */
124            tuplestore_end(node->working_table);
125
126            /* intermediate table becomes working table */
127            node->working_table = node->intermediate_table;
128
129            /* create new empty intermediate table */
130            node->intermediate_table = tuplestore_begin_heap(false, false,
131                                                            work_mem);
132            node->intermediate_empty = true;
133
134            /* reset the recursive term */
135            innerPlan->chgParam = bms_add_member(innerPlan->chgParam,
136                                                plan->wtParam);
137
138            /* and continue fetching from recursive term */
139            continue;
140        }
141
142        if (plan->numCols > 0)
143        {
144            /* Find or build hashtable entry for this tuple's group */
145            LookupTupleHashEntry(node->hashtable, slot, &isNew, NULL);
146            /* Must reset temp context after each hashtable lookup */
147            MemoryContextReset(node->tempContext);
148            /* Ignore tuple if already seen */
149            if (!isNew)
150                continue;
151        }
152
153        /* Else, tuple is good; stash it in intermediate table ... */
154        node->intermediate_empty = false;
155        tuplestore_puttupleslot(node->intermediate_table, slot);
156        /* ... and return it */
157        return slot;
158    }
159
160    return NULL;

```

(5)

# References

- (1) <https://www.postgresql.org/docs/current/queries-with.html>
- (2) Lecture slides from course “Database Systems”
- (3) <https://medium.com/swlh/postgres-recursive-query-cte-or-recursive-function-3ea1ea22c57c>
- (4) <https://pganalyze.com/docs/explain#postgres-plan-nodes>
- (5) <https://github.com/postgres/postgres/tree/master/src/backend/executor>
- (6) <https://gitlab.com/postgres/postgres/blob/master/src/include/nodes/plannodes.h>

