

Optimizing Queries on SQL-Level

Thomas Glas

Guided Research

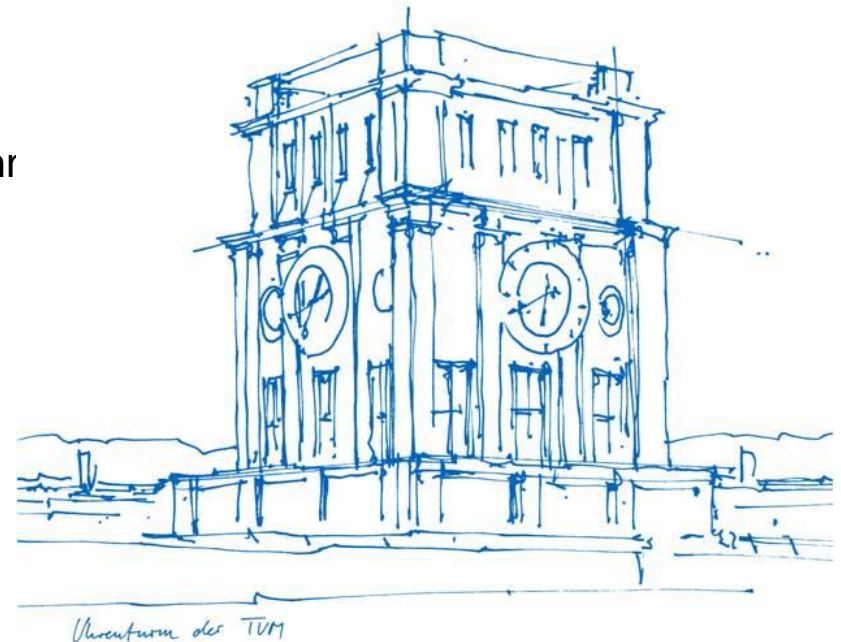
Supervised by Alexander Beischl

Technische Universität München

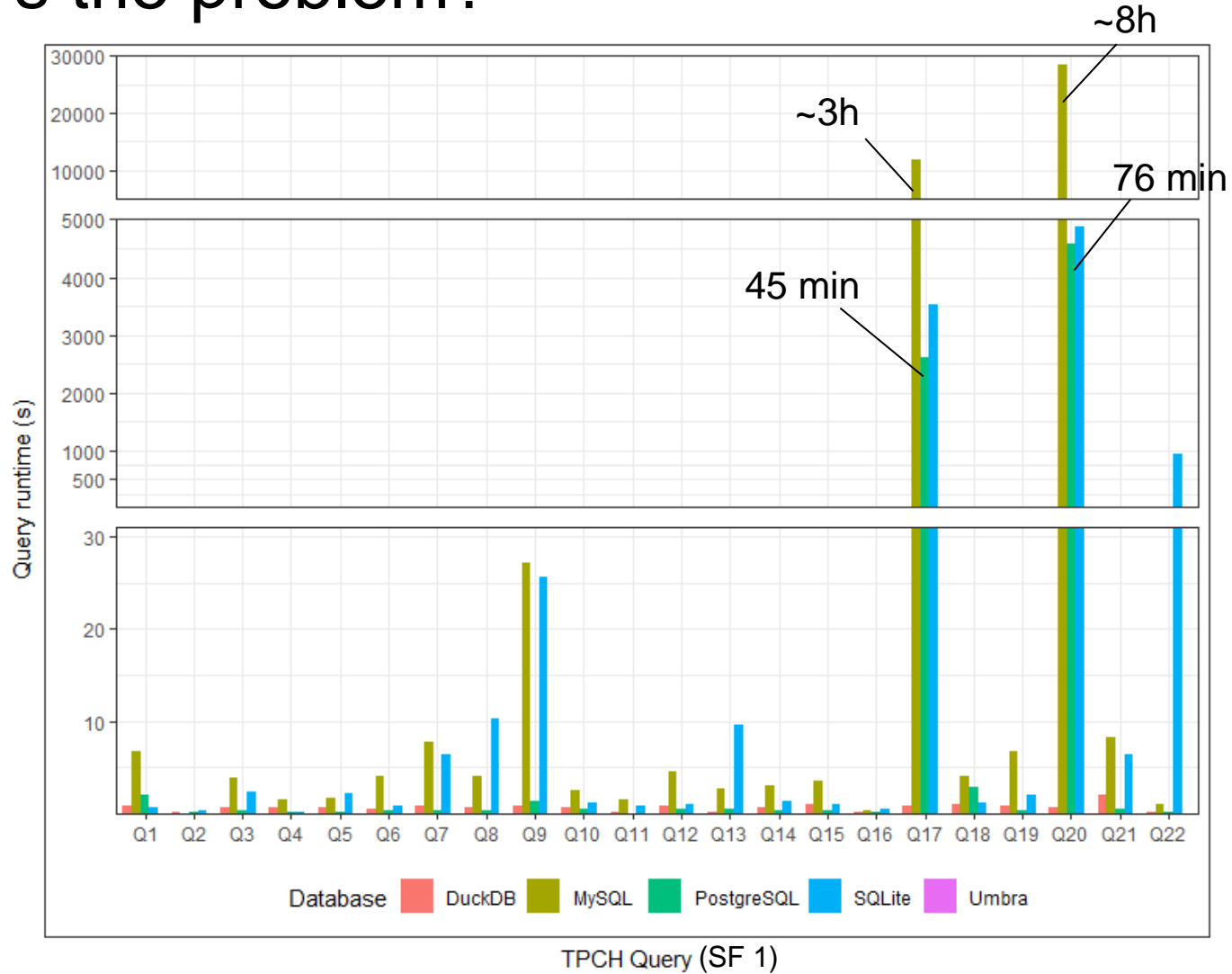
TUM School of Computation, Information and Techn

Lehrstuhl für Datenbanksysteme

München, 04. April 2023



What's the problem?



Correlated Subqueries...

```
1  -- TPC-H Query 17
2  select
3      sum(l_extendedprice) / 7.0 as avg_yearly
4  from
5      lineitem,
6      part
7  where
8      p_partkey = l_partkey
9      and p_brand = 'Brand#23'
10     and p_container = 'MED BOX'
11     and l_quantity < (
12         select
13             0.2 * avg(l_quantity)
14         from
15             lineitem
16         where
17             l_partkey = p_partkey
18     )
```

We can decorrelate this query in SQL

```

1  -- TPC-H Query 17
2  select
3      sum(l_extendedprice) / 7.0 as avg_yearly
4  from
5      lineitem,
6      part
7  where
8      p_partkey = l_partkey
9      and p_brand = 'Brand#23'
10     and p_container = 'MED BOX'
11     and l_quantity < (
12         select
13             0.2 * avg(l_quantity)
14         from
15             lineitem
16         where
17             l_partkey = p_partkey
18     )

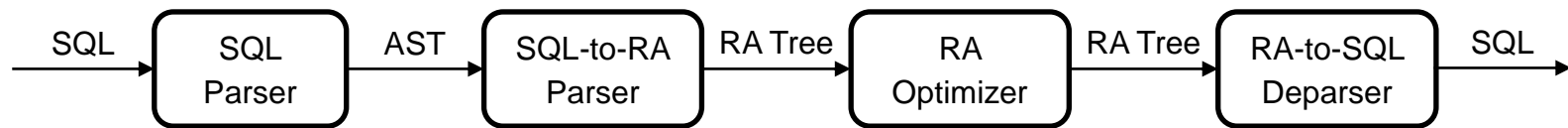
```

```

1  select
2      sum(l_extendedprice)/7.0 as avg_yearly
3  from
4      part,
5      lineitem,
6      (
7          select 0.2*avg(l_quantity),
8              l_partkey
9          from
10             lineitem
11          group by
12              l_partkey
13          ) as t1(m,t1_p_partkey)
14  where
15      l_quantity<t1.m
16      and p_partkey=t1.t1_p_partkey
17      and p_partkey=l_partkey
18      and p_brand='Brand#23'
19      and p_container='MED BOX';

```

Architecture



SQL Parser

pganalyze / libpg_query Public

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

BSD-3-Clause license

834 stars 135 forks

Star

Watch

Code

Issues 26

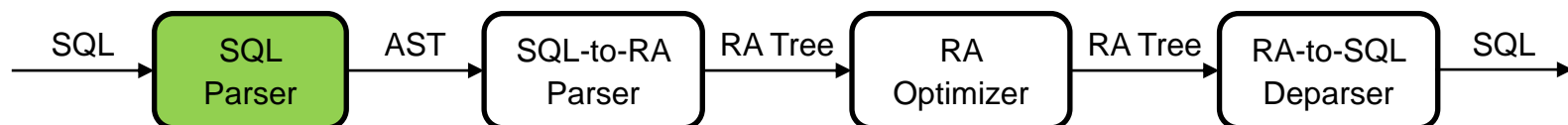
Pull requests 2

Actions

Projects

Wiki

...

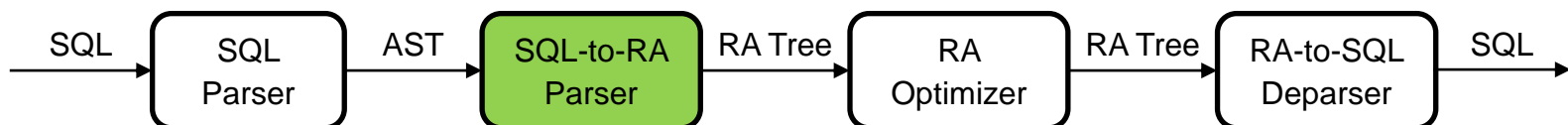
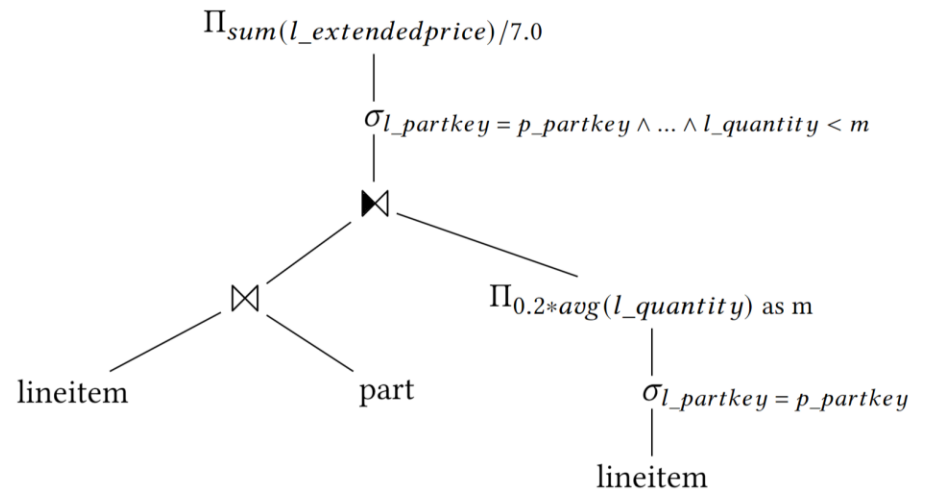


SQL-to-RA Parser

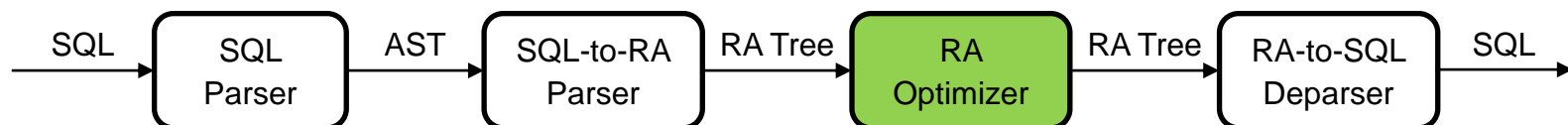
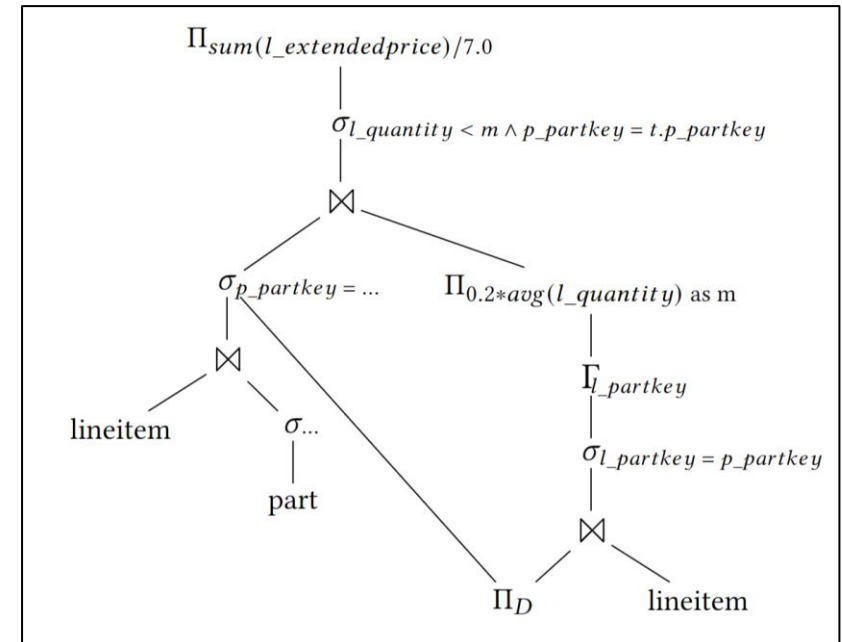
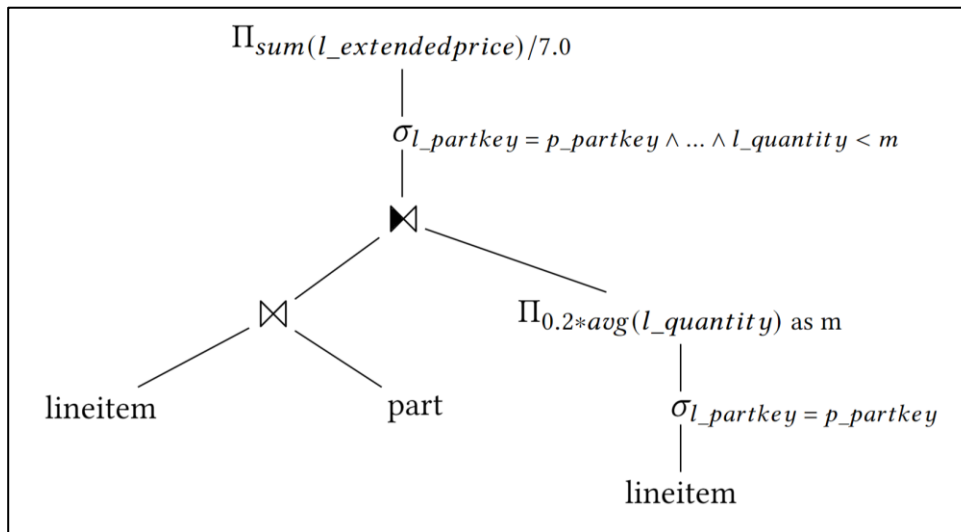
```

1  -- TPC-H Query 17
2  select
3      sum(l_extendedprice) / 7.0 as avg_yearly
4  from
5      lineitem,
6      part
7  where
8      p_partkey = l_partkey
9      and p_brand = 'Brand#23'
10     and p_container = 'MED BOX'
11     and l_quantity < (
12         select
13             0.2 * avg(l_quantity)
14         from
15             lineitem
16         where
17             l_partkey = p_partkey
18     )

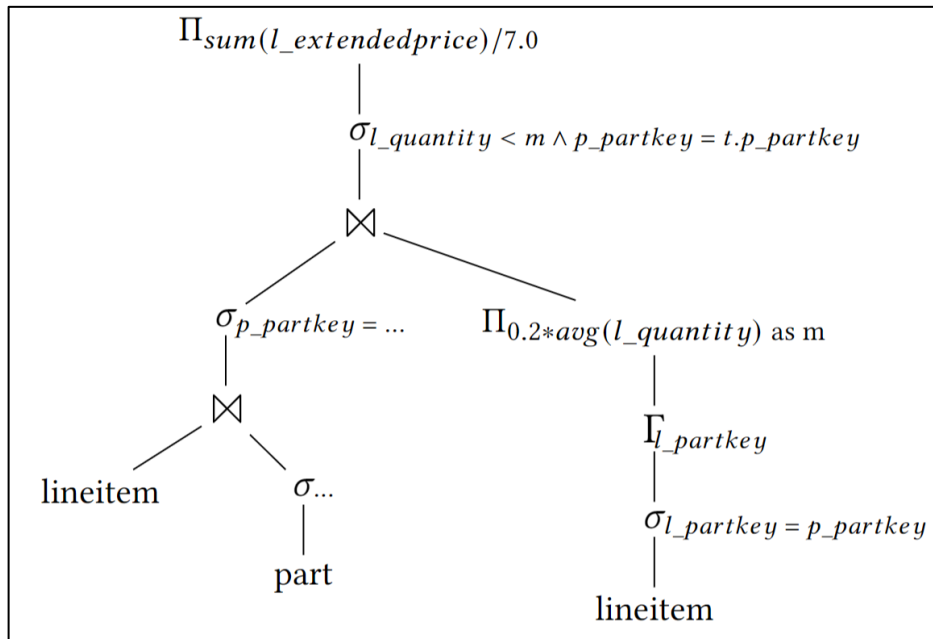
```



RA Optimizer



RA-to-SQL Deparser



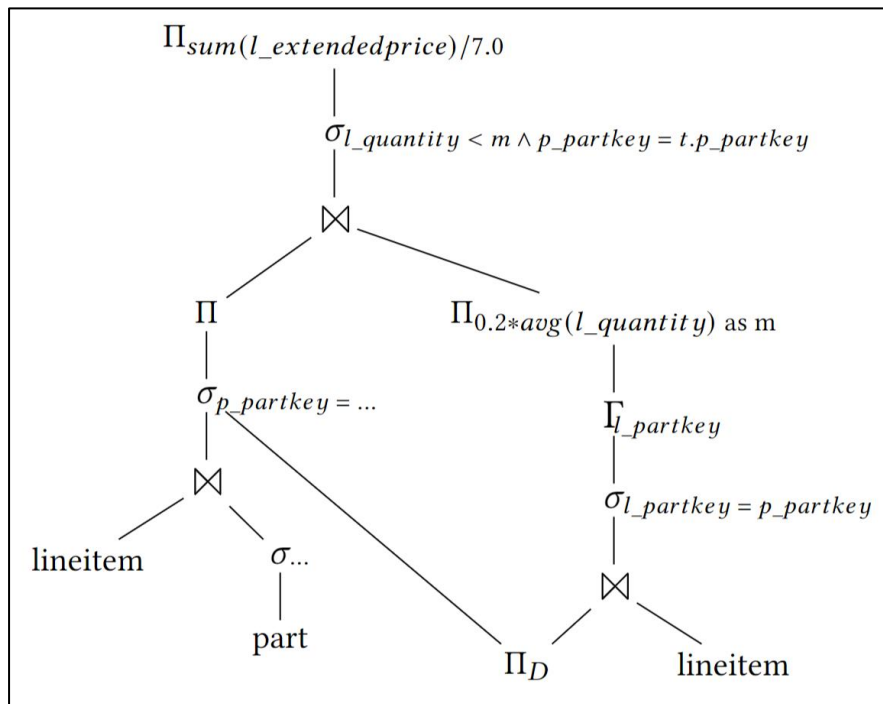
```

1  select
2      sum(l_extendedprice)/7.0 as avg_yearly
3  from
4      part,
5      lineitem,
6      (
7          select 0.2*avg(l_quantity),
8                  l_partkey
9          from
10             lineitem
11          group by
12              l_partkey
13          ) as t1(m,t1_p_partkey)
14  where
15      l_quantity < t1.m
16      and p_partkey = t1.t1_p_partkey
17      and p_partkey = l_partkey
18      and p_brand = 'Brand#23'
19      and p_container = 'MED BOX';

```



To decouple, or not to decouple



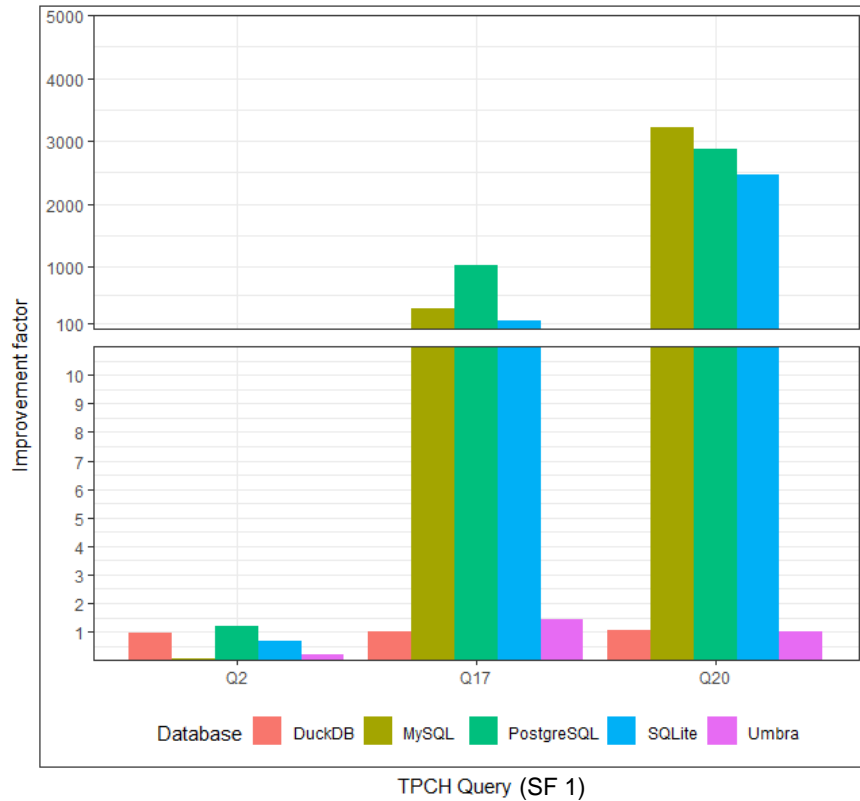
```

1  -- TPC-H Query 17
2  with cte_2 as (
3      select
4          l_extendedprice,
5          l_quantity,
6          p_partkey,
7          l_partkey
8      from
9          part,
10         lineitem
11     where
12         p_partkey=l_partkey
13         and (p_brand='Brand#23' and p_container='MED BOX')
14 )
15 select
16     (sum(l_extendedprice)/7.0) as avg_yearly
17 from
18     cte_2,
19     (select
20         (0.2*avg(l_quantity)),
21         d.p_partkey
22     from
23         (select
24             p_partkey
25         from cte_2
26         ) as d(p_partkey), lineitem
27     where l_partkey=d.p_partkey
28     group by d.p_partkey
29     ) as t1(m,t1_p_partkey)
30 where
31     l_quantity<t1.m
32     and p_partkey=t1.t1_p_partkey;

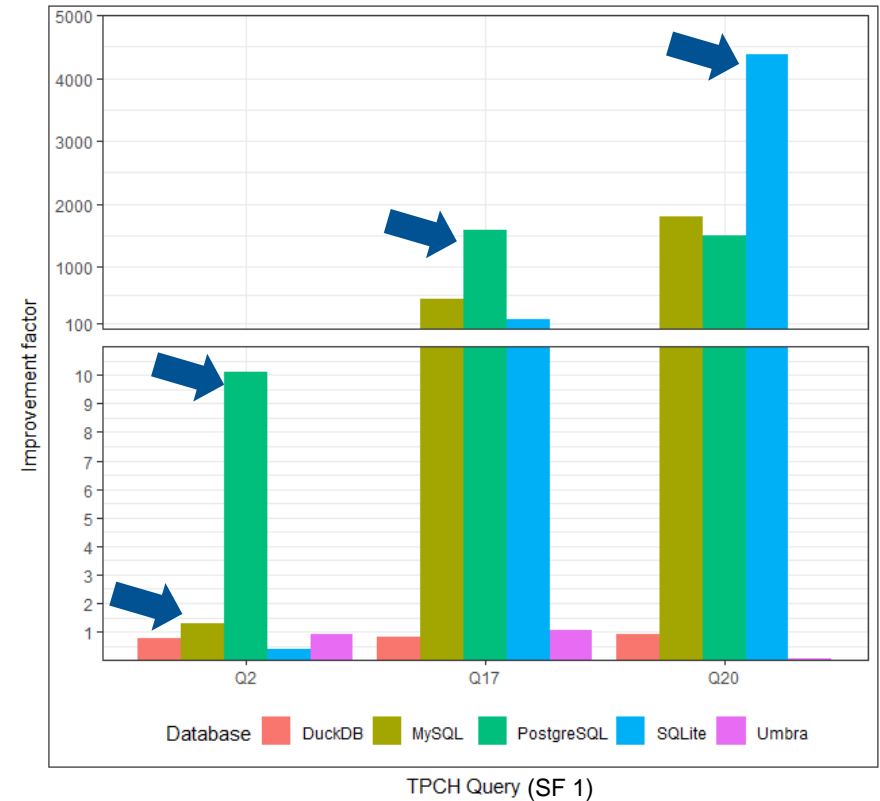
```

To decouple, or not to decouple

Decoupled

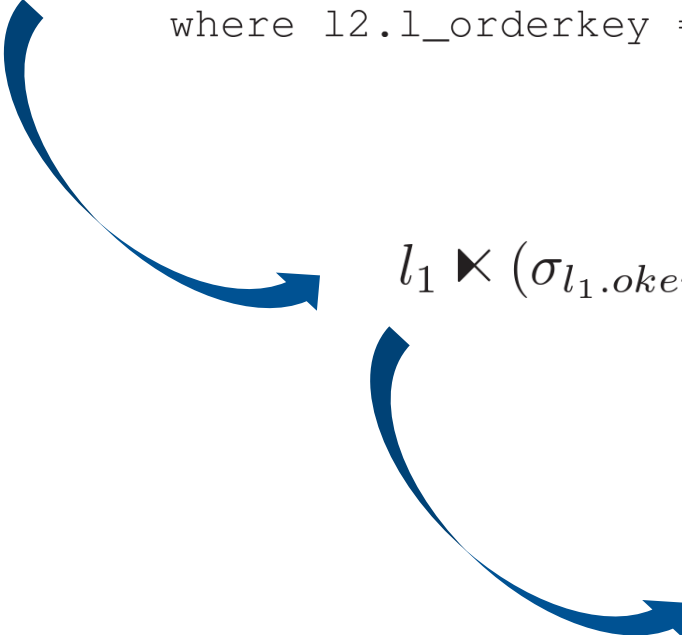


Not decoupled



Correlated *exists* and *in* subqueries

```
select ...  
from   lineitem l1 ...  
where  exists (select *  
               from lineitem l2  
               where l2.l_orderkey = l1.l_orderkey)
```



The diagram illustrates the transformation of a correlated `exists` subquery into a join operation. A blue curved arrow originates from the `exists` clause in the SQL query and points to the first relational algebra expression. A second blue curved arrow originates from the inner subquery of the first expression and points to the second, more simplified relational algebra expression.

$$l_1 \bowtie (\sigma_{l_1.okey=l_2.okey}(l_2))$$

$$l_1 \bowtie_{l_1.okey=l_2.okey} (l_2)$$

Correlated *exists* and *in* subqueries

```

1  -- TPC-H Query 4
2  select
3      o_orderpriority,
4      count(*) as order_count
5  from
6      orders
7  where
8      o_orderdate >= date '1993-07-01'
9      and o_orderdate < date '1993-07-01' + interval '3' month
10     and exists (
11         select
12             *
13         from
14             lineitem
15         where
16             l_orderkey = o_orderkey
17             and l_commitdate < l_receiptdate
18     )
19 group by
20     o_orderpriority
21 order by
22     o_orderpriority

```

```

1  with cte_1(l_orderkey) as (
2      select l_orderkey
3      from lineitem
4      where l_commitdate < l_receiptdate
5  )
6  select
7      o_orderpriority,
8      count(*) as order_count
9  from orders
10 where exists (
11     select *
12     from cte_1
13     where cte_1.l_orderkey = o_orderkey
14 )
15 and o_orderdate >= date '1993-07-01'
16 and o_orderdate < date '1993-07-01' + interval '3' month
17 group by o_orderpriority
18 order by o_orderpriority;

```

Results

