Philipp Fent
Technical University of Munich
fent@in.tum.de

# Umbra

- TUM's first DBMS acquired by Salesforce

- Rewrite from scratch

- Cutting-edge database research

- Disk-based with in-memory performance



## CSRankings: Computer Science Rankings

CSRankings is a metrics-based ranking of top computer science institutions around the world. **Click on a triangle (▶)** to expand areas or institutions. **Click on a name** to go to a faculty member's home page. **Click on a chart icon** (the 📊 after a name or institution) to see the distribution of their publication areas as a bar chart ▾. **Click on a Google Scholar icon** (🎓) to see publications, and **click on the DBLP logo** (▶) to go to a DBLP entry. *Applying to grad school? Read this first.* **Do you find CSRankings useful? Sponsor CSrankings on GitHub.**

Rank institutions in [ the world ▾ ] by publications from [ 2017 ▾ ] to [ 2023 ▾ ]

**All Areas** [off | on]

**AI** [off | on]

- ▶ Artificial intelligence ☐
- ▶ Computer vision ☐
- ▶ Machine learning ☐
- ▶ Natural language processing ☐
- ▶ The Web & information retrieval ☐

**Systems** [off | on]

- ▶ Computer architecture ☐
- ▶ Computer networks ☐
- ▶ Computer security ☐
- ▶ Databases ☑
- ▶ Design automation ☐
- ▶ Embedded & real-time systems ☐

| #  | Institution |
|----|-------------|
| 1  | ▶ TU Munich 🇩🇪 📊 |
| 2  | ▶ HKUST 🇭🇰 📊 |
| 3  | ▶ Tsinghua University 🇨🇳 📊 |
| 4  | ▶ University of Waterloo 🇨🇦 📊 |
| 5  | ▶ National University of Singapore 🇸🇬 📊 |
| 6  | ▶ Duke University 🇺🇸 📊 |
| 7  | ▶ Chinese University of Hong Kong 🇭🇰 📊 |
| 8  | ▶ Nanyang Technological University 🇸🇬 📊 |
| 9  | ▶ Univ. of California - San Diego 🇺🇸 📊 |
| 10 | ▶ Univ. of California - Berkeley 🇺🇸 📊 |
| 11 | ▶ Peking University 🇨🇳 📊 |

# Performance



TPC-H SF10

# Performance



TPC-H SF10

# What makes Umbra fast?

# What makes Umbra fast?

- Pipelined execution
  - Keeps values in registers
  - Minimizes materialization

$$\bowtie_{a=b}$$

$$\sigma_{x=7}$$

$$R_1$$

$$\bowtie_{z=c}$$

$$\Gamma_{z;count(*)}$$

$$\sigma_{y=3}$$

$$R_2$$

$$R_3$$

# What makes Umbra fast?

- Pipelined execution

- Data-centric code generation

  - Efficient code for complex expressions

```
%1 = zext i64 %int1;                         Zero extend to 64 bit
%2 = zext i64 %int2;
%3 = rotr i64 %2, 32;                              Rotate right
%v = or i64 %1, %3;                       Combine int1 and int2
%5 = crc32 i64 6763793487589347598, %v;           First crc32
%6 = crc32 i64 4593845798347983834, %v;          Second crc32
%7 = rotr i64 %6, 32;                        Shift second part
%8 = xor i64 %5, %7;                       Combine hash parts
%hash = mul i64 %8, 11400714819323198485;          Mix parts
```

# What makes Umbra fast?

- Pipelined execution

- Data-centric code generation

- Fully parallel algorithms

  - Allows scaling

  - Benefits from new hardware

# What makes Umbra fast?

- Pipelined execution

- Data-centric code generation

- Fully parallel algorithms

- **State-of-the-art query optimizer**

# What makes Umbra fast?

- Pipelined execution

- Data-centric code generation

- Fully parallel algorithms

- **State-of-the-art query optimizer**

Research system with all custom advanced parts

We're commercializing soon!

# Query Optimization

- PostgreSQL grammar

- Parsed into relational algebra

  - Example: TPC-H Q17

  - https://umbra-db.com/interface/

# Query Optimization

- PostgreSQL grammar

- Parsed into relational algebra

- Optimizer passes over algebra

| |
|---|
| 1: Unoptimized Plan |
| 2: Expression Simplification |
| 3: Unnesting |
| 4: Predicate Pushdown |
| 5: Initial Join Tree |
| 6: Sideway Information Passing |
| 7: Operator Reordering |
| 8: Early Probing |
| 9: Common Subtree Elimination |
| 10: Physical Operator Mapping |

# Query Optimization

- PostgreSQL grammar

- Parsed into relational algebra

- Optimizer passes over algebra

1: Unoptimized Plan

2: Expression Simplification

3: Unnesting

4: Predicate Pushdown

5: Initial Join Tree

Rule-based Canonicalization

6: Sideway Information Passing

7: Operator Reordering

8: Early Probing

9: Common Subtree Elimination

10: Physical Operator Mapping

Cost-based Optimization

# Expression Simplification

- Fold constants

- Canonicalize expressions

```
      o_orderdate >= date '1994-01-01'
  and o_orderdate <  date '1994-01-01' + interval '1' year

                       ==

  o_orderdate between date '1994-01-01' and date '1994-12-31'
```

- Execute in evaluation engine

# Query Unnesting & Decorrelation

- Unnesting Arbitrary Queries



## Unnesting Arbitrary Queries

Thomas Neumann and Alfons Kemper
Technische Universität München
Munich, Germany
neumann@in.tum.de, kemper@in.tum.de

**Abstract:** SQL-99 allows for nested subqueries at nearly all places within a query. From a user's point of view, nested queries can greatly simplify the formulation of complex queries. However, nested queries that are correlated with the outer queries frequently lead to dependent joins with nested loops evaluations and thus poor performance.

Existing systems therefore use a number of heuristics to *unnest* these queries, i.e., de-correlate them. These unnesting techniques can greatly speed up query processing, but are usually limited to certain classes of queries. To the best of our knowledge no existing system can de-correlate queries in the general case. We present a generic approach for unnesting arbitrary queries. As a result, the de-correlated queries allow for much simpler and much more efficient query evaluation.

### 1 Introduction

Subqueries are frequently used in SQL queries to simplify query formulation. Consider for our running examples the following schema:

- students: {[id, name, major, year, ...]}
- exams: {[sid, course, curriculum, date, ...]}

Then the following is a nested query to find for each student the best exams (according to the German grading system where lower numbers are better):

```
Q1: select s.name,e.course
    from   students s,exams e
    where  s.id=e.sid and
           e.grade=(select min(e2.grade)
                    from exams e2
                    where s.id=e2.sid)
```

Conceptually, for each student, exam pair $(s, e)$ it determines, in the subquery, whether or not this particular exam $e$ has the best grade of all exams of this particular student $s$.

From a performance point of view the query is not so nice, as the subquery has to be re-evaluated for every student, exam pair. From a technical perspective the query contains a

383



● DuckDB                    Documentation ˅    Blog

## Blog

2023-05-26    Mark Raasveldt

## Correlated Subqueries in SQL

Subqueries in SQL are a powerful abstraction that allow simple queries to be used as composable building blocks. They allow you to break down complex problems into smaller parts, and subsequently make it easier to write, understand and maintain large and complex queries.

DuckDB uses a state-of-the-art subquery decorrelation optimizer that allows subqueries to be executed very efficiently. As a result, users can freely use subqueries to create expressive queries without having to worry about manually rewriting subqueries into joins. For more information, skip to the Performance section.

### Types of Subqueries

SQL subqueries exist in two main forms: subqueries as *expressions* and subqueries as *tables*. Subqueries that are used as expressions can be used in the `SELECT` or `WHERE` clauses. Subqueries that are used as tables can be used in the `FROM` clause. In this blog post we will focus on subqueries used as *expressions*. A future blog post will discuss subqueries as *tables*.

Subqueries as expressions exist in three forms.

- Scalar subqueries
- `EXISTS`
- `IN / ANY / ALL`

All of the subqueries can be either *correlated* or *uncorrelated*. An uncorrelated subquery is a query that is independent from the outer query. A correlated subquery is a subquery that contains expressions from the outer query. Correlated subqueries can be seen as *parameterized subqueries*.

# Query Unnesting

- Unnesting Arbitrary Queries
  - $O(n^2)$

# Query Unnesting

- Unnesting Arbitrary Queries
  - O(n²)

# Query Unnesting

- Unnesting Arbitrary Queries
  - $O(n^2) \rightarrow O(n)$
  - Huge improvement

# Predicate Pushdown

- Place predicates at scan

- Propagate & fold constants

# Predicate Pushdown

- Place predicates at scan

- Propagate & fold constants

# Predicate Pushdown

- Place predicates at scan

- Propagate & fold constants



```
where p_partkey = 42
```

# Predicate Pushdown

- Place predicates at scan

- Propagate & fold constants

# Initial Join Tree

- Push joins through aggregates

- Expand transitive join conditions

```
        c_nationkey = s_nationkey
    and s_nationkey = n_nationkey


                ==


        c_nationkey = s_nationkey
    and s_nationkey = n_nationkey
    and c_nationkey = n_nationkey
```

# Initial Join Tree

- Push joins through aggregates

- Expand transitive join conditions

- Drop unnecessary joins

```
select sum(o_totalprice)
  from customer, orders
 where c_custkey = o_custkey


              ==


select sum(o_totalprice)
  from orders
```

# Cost-Based Optimization

- Heuristics vs. statistics

# Cost-Based Optimization

- Heuristics vs. statistics

- Statistics in Umbra:

  - Samples

  - Distinct counts

  - Numerical statistics (mean, variance) for aggregates

  - Functional dependencies

⇒ Estimate execution cost

# Sample Evaluation

- Maintain uniform reservoir sample

- Evaluate scan predicates σ on sample

- Execute in evaluation engine

- Surprisingly accurate
    - 1024 tuples ~ 0.1% error

```
select count(*)
  from lineitem
 where l_commitdate < l_receiptdate
   and l_shipdate < l_commitdate
```

# Sample Evaluation

```
for l in lineitem:
  if not l_shipdate < l_commitdate:
    continue  -- 51% taken
  if not l_commitdate < l_receiptdate:
    continue  -- 75% taken

  counter++
```
**Variant Ⓐ : Separate branches**

```
for l in lineitem:
  if not l_commitdate < l_receiptdate:
    continue  -- 37% taken
  if not l_shipdate < l_commitdate:
    continue  -- 81% taken

  counter++
```
**Variant Ⓑ : Separate branches**

```
for l in lineitem:
  if not (l_shipdate < l_commitdate
    and l_commitdate < l_receiptdate):
    continue  -- 88% taken

  counter++
```
**Variant Ⓒ : Combined branch**

# Sample Evaluation

```
for l in lineitem:
  if not l_shipdate < l_commitdate:
    continue  -- 51% taken
  if not l_commitdate < l_receiptdate:
    continue  -- 75% taken

  counter++
```
**Variant Ⓐ : Separate branches**

```
for l in lineitem:
  if not l_commitdate < l_receiptdate:
    continue  -- 37% taken
  if not l_shipdate < l_commitdate:
    continue  -- 81% taken

  counter++
```
**Variant Ⓑ : Separate branches**

```
for l in lineitem:
  if not (l_shipdate < l_commitdate
    and l_commitdate < l_receiptdate):
    continue  -- 88% taken

  counter++
```
**Variant Ⓒ : Combined branch**

| Variant | branch-misses | instructions | loads | exec. time |
|---------|---------------|--------------|-------|------------|
| Ⓐ | 0.63 / tpl | 7.62 / tpl | 2.85 / tpl | 18.4 ms |
| Ⓑ | 0.58 / tpl | 7.91 / tpl | 3.00 / tpl | 17.7 ms |
| Ⓒ | 0.13 / tpl | 11.67 / tpl | 3.37 / tpl | **12.7 ms** |

# Sample Evaluation

- Estimate (correlated) predicates with confidence

- Any combination of predicates

- Tricky when 0 / 1024 tuples qualify

- Can do better for conjunctions

# Sample Evaluation

- Calculate matches-bitsets
- Combine them to optimize ordering
  - TPC-H Q12:

```sql
where l_shipmode in ('MAIL', 'SHIP')
  and l_commitdate < l_receiptdate
  and l_shipdate < l_commitdate
  and l_receiptdate between date '1994-01-01'
                        and date '1994-12-31'
```

```
  0100'0011'1010'0100'1110'1011'1011'1100'1010'1010'1011'0000'1011'0011'1100'0000
& 0000'1111'0000'1111'0000'1111'0000'1111'0000'1111'0000'1111'0000'1111'0000'1111
& 1111'0000'1111'0000'1111'0000'1111'0000'1111'0000'1111'0000'1111'0000'1111'0000
& 1010'0110'1110'1110'1000'0011'0111'0101'0110'1111'1001'1101'1110'0011'1000'0001
```

# Early Execution

- Size of sample > table size

- Allows a third round of constant propagation
  - Especially for small fact tables

```
select r_regionkey
  from region
 where r_name = 'Europe'


     ==


select 3
```

# Join Ordering

- Hash Joins rule
    - Indexes don't allow bushy plans -> less useful

# Join Ordering

- ● Hash Joins rule
  - ○ Indexes don't allow bushy plans -> less useful

# Join Ordering

- Hash Joins rule
  - Indexes don't allow bushy plans -> less useful
- Distinct count estimates with Pat Sellinger's equations
- HyperLogLog intersections
- Mean & stddev approximations for `l_quantity` `<` `0.2` `*` `avg(l_quantity)`

# Early Probing

- Semijoin reduction

- Reuses existing hash tables

- Can use bloom filters if beneficial

# Physical Optimization

- Indexes

- Worst-case optimal join

# Physical Optimization

- Indexes

- Worst-case optimal join

- Groupjoin

# Physical Optimization

- Indexes

- Worst-case optimal join

- Groupjoin

- Range join

- Join micro-optimizations

  - Multiset semantics

  - Allocation sizes

# Recap

- Query compilation & optimization
  - Optimizer passes
  - Rule-based canonicalization
  - Cost-based optimization
- Cutting-edge research
  - Join ordering
  - Cardinality estimation
  - Integrated in a running system

1: Unoptimized Plan

2: Expression Simplification

3: Unnesting

4: Predicate Pushdown

5: Initial Join Tree

6: Sideway Information Passing

7: Operator Reordering

8: Early Probing

9: Common Subtree Elimination

10: Physical Operator Mapping

# Conclusion

- Low latency analytical queries

- Also works excellent for transactional and graph workloads

We are commercializing

Reach out:

fent@in.tum.de

TUM Open Source Project

LingoDB