

Andy Pavlo // 15-721 // Spring 2023

ADMINISTRIVIA

Project #2:

- \rightarrow Feedback Submission: Saturday April 1st
- \rightarrow Final Submission: Monday May 1st

Project #3

- \rightarrow Status Update Presentation: Wednesday April 5th
- → Final Presentations: Friday May 5th @ 5:30pm
- \rightarrow Please email me if you want to discuss your project!

QUERY OPTIMIZATION

For a given query, find a **<u>correct</u>** execution plan that has the lowest "cost".

This is the part of a DBMS that is the hardest to implement well (proven to be NP-Complete).

No optimizer truly produces the "optimal" plan \rightarrow Use estimation techniques to guess real plan cost.

 \rightarrow Use heuristics to limit the search space.



QUERY OPTIMIZATION STRATEGIES

Choice #1: Heuristics

 \rightarrow INGRES, Oracle (until mid 1990s)

Choice #2: Heuristics + Cost-based Join Search

 \rightarrow System R, early IBM DB2, most open-source DBMSs

Choice #3: Randomized Search

 \rightarrow Academics in the 1980s, current Postgres

Choice #4: Stratified Search

 \rightarrow IBM's STARBURST (late 1980s), now IBM DB2 + Oracle

Choice #5: Unified Search

 \rightarrow Volcano/Cascades in 1990s, now MSSQL + Greenplum

STRATIFIED SEARCH

First rewrite the logical query plan using transformation rules.

- \rightarrow The engine checks whether the transformation is allowed before it can be applied.
- \rightarrow Cost is never considered in this step.

Then perform a cost-based search to map the logical plan to a physical plan.

UNIFIED SEARCH

- Unify the notion of both logical \rightarrow logical and logical \rightarrow physical transformations. \rightarrow No need for separate stages because everything is
- → No need for separate stages because everything is transformations.

This approach generates many transformations, so it makes heavy use of memoization to reduce redundant work.



TOP-DOWN VS. BOTTOM-UP

Top-down Optimization

- → Start with the outcome that the query wants, and then work down the tree to find the optimal plan that gets you to that goal.
- \rightarrow **Examples**: Volcano, Cascades

Bottom-up Optimization

- \rightarrow Start with nothing and then build up the plan to get to the outcome that you want.
- \rightarrow **Examples**: System R, Starburst



TODAY'S AGENDA

Logical Query Optimization Cascades Real-World Implementations



LOGICAL QUERY OPTIMIZATION

- Transform a logical plan into an equivalent logical plan using pattern matching rules.
- The goal is to increase the likelihood of enumerating the optimal plan in the search.
- Cannot compare plans because there is no cost model but can "direct" a transformation to a preferred side.

LOGICAL QUERY OPTIMIZATION

Split Conjunctive Predicates Predicate Pushdown Replace Cartesian Products with Joins Projection Pushdown

Source: <u>Thomas Neumann</u> **CMU-DB** 15-721 (Spring 2023)

SPLIT CONJUNCTIVE PREDICATES

SELECT ARTIST.NAME
FROM ARTIST, APPEARS, ALBUM
WHERE ARTIST.ID=APPEARS.ARTIST_ID
AND APPEARS.ALBUM_ID=ALBUM.ID
AND ALBUM.NAME="Andy's OG Remix"

Decompose predicates into their simplest forms to make it easier for the optimizer to move them around.



PREDICATE PUSHDOWN

SELECT ARTIST.NAME
FROM ARTIST, APPEARS, ALBUM
WHERE ARTIST.ID=APPEARS.ARTIST_ID
AND APPEARS.ALBUM_ID=ALBUM.ID
AND ALBUM.NAME="Andy's OG Remix"

Move the predicate to the lowest point in the plan after Cartesian products.



REPLACE CARTESIAN PRODUCTS

SELECT ARTIST.NAME
FROM ARTIST, APPEARS, ALBUM
WHERE ARTIST.ID=APPEARS.ARTIST_ID
AND APPEARS.ALBUM_ID=ALBUM.ID
AND ALBUM.NAME="Andy's OG Remix"

Replace all Cartesian Products with inner joins using the join predicates.



PROJECTION PUSHDOWN

SELECT ARTIST.NAME
FROM ARTIST, APPEARS, ALBUM
WHERE ARTIST.ID=APPEARS.ARTIST_ID
AND APPEARS.ALBUM_ID=ALBUM.ID
AND ALBUM.NAME="Andy's OG Remix"

Eliminate redundant attributes before pipeline breakers to reduce materialization cost.



PHYSICAL QUERY OPTIMIZATION

Transform a query plan's logical operators into physical operators.

- \rightarrow Add more execution information
- \rightarrow Select indexes / access paths
- \rightarrow Choose operator implementations
- \rightarrow Choose when to materialize (i.e., temp tables).

This stage must support cost model estimates.

OBSERVATION

All the queries we have looked at so far have had the following properties:

- \rightarrow Equi/Inner Joins
- \rightarrow Simple join predicates that reference only two tables.
- \rightarrow No cross products

Real-world queries are more complex / nasty:

- \rightarrow Outer Joins
- \rightarrow Semi Joins
- \rightarrow Anti Joins
- \rightarrow Lateral Joins

REORDERING LIMITATIONS



No valid reordering is possible!

The A→B operator is not commutative with B→C. → The DBMS does not know the value of B. val until after computing the join with A.

Source: <u>Pit Fender</u> **ECMU-DB** 15-721 (Spring 2023)

PLAN ENUMERATION

Approach #1: Transformation

- \rightarrow Modify some part of an existing query plan to transform it into an alternative plan that is equivalent.
- \rightarrow Top-Down Search

Approach #2: Generative

- \rightarrow Iteratively assemble and add building blocks to generate a query plan.
- \rightarrow Bottom-Up Search





DYNAMIC PROGRAMMING OPTIMIZER

Model the query as a hypergraph and then incrementally expand to enumerate new plans.

Algorithm Overview:

- \rightarrow Iterate connected sub-graphs and incrementally add new edges to other nodes to complete query plan.
- \rightarrow Use rules to determine which nodes the traversal is allowed to visit and expand.



CASCADES OPTIMIZER

Object-oriented implementation of the Volcano query optimizer.

→ **Top-down approach** (backward chaining) using branchand-bound search.

Supports simplistic expression re-writing through a direct mapping function rather than an exhaustive search.



15-721 (Spring 2023



Graefe

CASCADES OPTIMIZER

Optimization tasks as data structures. Rules to place property enforcers. Ordering of moves by promise. Predicates as logical/physical operators.



CASCADES - EXPRESSIONS

An <u>expression</u> represents some operation in the query with zero or more input expressions.
 → Optimizer needs to be able to quickly determine whether two expressions are equivalent.

```
SELECT * FROM A
JOIN B ON A.id = B.id
JOIN C ON C.id = A.id;
```

Logical Expression: (A ⋈ B) ⋈ C Physical Expression: (A_{Seq} ⋈_{HJ} B_{Seq}) ⋈_{NL} C_{Idx}

CASCADES - GROUPS

A **group** is a set of logically equivalent logical and physical expressions that produce the same output. \rightarrow All logical forms of an expression.

 \rightarrow All physical expressions that can be derived from selecting the allowable physical operators for the corresponding logical forms.

Physical ExpsPhysical ExpsOutput: $1. (A \bowtie B) \bowtie C$ $1. (A_{Seq} \bowtie_{NL} B_{Seq}) \bowtie_{NL}$ $2. (B \bowtie C) \bowtie A$ $2. (B_{Seq} \bowtie_{NL} C_{Seq}) \bowtie_{NL}$ $3. (A \bowtie C) \bowtie B$ $3. (A_{Seq} \bowtie_{NL} C_{Seq}) \bowtie_{NL}$ $4. A \bowtie (B \bowtie C)$ $4. A_{Seq} \bowtie_{NL} (C_{Seq} \bowtie_{NL})$ $1. (A_{Seq} \bowtie_{NL} C_{Seq}) \bowtie_{NL}$ $1. (A_{Seq} \bowtie_{NL} C_{Seq}) \bowtie_{NL}$
--

15-721 (Spring 2023)

Equivalent Expressions

CASCADES - MULTI-EXPRESSION

Instead of explicitly instantiating all possible expressions in a group, the optimizer implicitly represents redundant expressions in a group as a **multi-expression**.

 \rightarrow This reduces the number of transformations, storage overhead, and repeated cost estimations.

Output: [ABC]	Logical Multi-Exps 1. [AB]⋈[C] 2. [BC]⋈[A] 3. [AC]⋈[B] 4. [A]⋈[BC]	Physical Multi-Exps 1. [AB]⋈ _{SM} [C] 2. [AB]⋈ _{HJ} [C] 3. [AB]⋈ _{NL} [C] 4. [BC]⋈ _{SM} [A]
	4. [A]⊠[BC]	4. [BC]⋈ _{SM} [A]
	•	

CASCADES - RULES

- A <u>**rule</u>** is a transformation of an expression to a logically equivalent expression.</u>
- → **Transformation Rule:** Logical to Logical
- → **Implementation Rule:** Logical to Physical

Each rule is represented as a pair of attributes:

- \rightarrow <u>**Pattern**</u>: Defines the structure of the logical expression that can be applied to the rule.
- \rightarrow **Substitute**: Defines the structure of the result after applying the rule.

CASCADES - RULES

Pattern

ECMU·DB 15-721 (Spring 2023)



26

CASCADES - MEMO TABLE

Stores all previously explored alternatives in a compact graph structure / hash table.

Equivalent operator trees and their corresponding plans are stored together in groups.

Provides memoization, duplicate detection, and property + cost management.

PRINCIPLE OF OPTIMALITY

Every sub-plan of an optimal plan is itself optimal.

This allows the optimizer to restrict the search space to a smaller set of expressions.

→ The optimizer never has to consider a plan containing subplan P1 that has a greater cost than equivalent plan P2 with the same physical properties.



15-721 (Spring 2023)

CASCADES - MEMO TABLE



CMU·DB 15-721 (Spring 2023)

SEARCH TERMINATION

Approach #1: Transformation Exhaustion

 \rightarrow Stop when there are no more ways to transform the target plan. Usually done per group.

Approach #2: Wall-clock Time

 \rightarrow Stop after the optimizer runs for some length of time.

Approach #3: Transformation Count

 \rightarrow Stop after a certain number of transformations have been considered.

Approach #4: Cost Threshold

 \rightarrow Stop when the optimizer finds a plan that has a lower cost than some threshold.

CASCADES IMPLEMENTATIONS

Standalone:

- \rightarrow <u>Wisconsin OPT++</u> (1990s)
- \rightarrow <u>Portland State Columbia</u> (1990s)
- \rightarrow <u>Greenplum Orca</u> (2010s)
- \rightarrow <u>Apache Calcite</u> (2010s)

Integrated:

- \rightarrow Microsoft SQL Server (1990s)
- \rightarrow <u>Tandem NonStop SQL</u> (1990s)
- $\rightarrow \underline{\text{Clustrix}}$ (2000s)
- \rightarrow CockroachDB (2010s)



REAL-WORLD IMPLEMENTATIONS

Microsoft SQL Server Apache Calcite Greenplum Orca CockroachDB SingleStore



MICROSOFT SQL SERVER

First Cascades implementation started in 1995.

- \rightarrow Derivatives are used in many MSFT database products.
- \rightarrow All transformations are written in C++. No DSL.
- \rightarrow Scalar / expression transformations are written in procedural code and not rules.

DBMS applies transformations in multiple stages with increasing scope and complexity.

→ The goal is to leverage domain knowledge to apply transformations that you always want to do first to reduce the search space.

MICROSOFT SQL SERVER



Source: Nico Bruno + Cesar Galindo-Legaria

15-721 (Spring 2023)

MICROSOFT SQL SERVER

Optimization #1: Timeouts are based on the number of transformations not wallclock time.
 → Ensures that overloaded systems do not generate different plans than under normal operations.

Optimization #2: Pre-populate the Memo Table with potentially useful join orderings.
→ Heuristics that consider relationships between tables.
→ Syntactic appearance in query.

APACHE CALCITE

Standalone extensible query optimization framework for data processing systems.

- → Support for pluggable query languages, cost models, and rules.
- → Does not distinguish between logical and physical operators. Physical properties are provided as annotations.

Originally part of <u>LucidDB</u>.

APACHE CALCITE: A FOUNDATIONAL FRAMEWORK FOR OPTIMIZED OUERY PROCESSING OVER HETEROGENEOUS DATA SOURCES SIGMOD 2018

15-721 (Spring 2023

GREENPLUM ORCA

Standalone Cascades implementation in C++.

- \rightarrow Originally written for <u>Greenplum</u>.
- \rightarrow Extended to support <u>HAWQ</u>.

A DBMS integrates Orca by implementing API to send catalog + stats + logical plans and then retrieve physical plans.

Supports multi-threaded search.



27 CMU·DB 15-721 (Spring 2023)

GREENPLUM ORCA - ENGINEERING

Issue #1: Remote Debugging

- → Automatically dump the state of the optimizer (with inputs) whenever an error occurs.
- \rightarrow The dump is enough to put the optimizer back in the exact same state later for further debugging.

Issue #2: Optimizer Accuracy

 \rightarrow Automatically check whether the ordering of the estimate cost of two plans matches their actual execution cost.

COCKROACHDB

Custom Cascades im All transformation r DSL (<u>OptGen</u>) and t → Can embed Go logic i analysis and modifica

Also considers scalar transformations toge

DSL: Optgen

// ConstructNot constructs an expression for the Not operator.
func (_f *Factory) ConstructNot(kinput opt.ScalarExpr) opt.ScalarExpr {

```
// [EliminateNot]
{
    _not, _ := input.(*memo.NotExpr)
    if _not != nil {
        input := _not.Input
        if _f.matchedRule == nil || _f.matchedRule(opt.EliminateNot) {
            _expr := input
            return _expr
        }
    }
    // ... other rules ...
    e := _f.mem.MemoizeNot(input)
    return _f.onConstructScalar(e)
}
Cockroach LABS
```

Source: <u>Rebecca Taft</u> SCMU-DB 15-721 (Spring 2023)

SINGLESTORE OPTIMIZER

Rewriter

 \rightarrow Logical-to-logical transformations with access to the cost-model.

Enumerator

- \rightarrow Logical-to-physical transformations.
- \rightarrow Mostly join ordering.

Planner

- \rightarrow Convert physical plans back to SQL.
- \rightarrow Contains SingleStore-specific commands for moving data.



EFCMU·DB 15-721 (Spring 2023)

SINGLESTORE OPTIMIZER



PARTING THOUGHTS

All of this relies on a good <u>cost model</u>. A good <u>cost model</u> needs good statistics.



NEXT CLASS

Cost Models

