

DB 2

14 – Query Optimization

Summer 2020

Torsten Grust
Universität Tübingen, Germany

1 | One Query — Millions of Plans

Q: Given a SQL query Q , what is *the optimal* (a reasonable)¹ plan to evaluate it? — **A:** It depends:

- Can we **simplify** (flatten, unnest) Q ?
- How can we **access the tables** referenced in Q ?
- How do **CPU and (sequential, random) I/O cost** compare?
- What is the **selectivity of the predicates** used in Q ?
- Which plan **operator implementations** are applicable?
- Can we **regroup/reorder the joins** in Q ?

¹ Here: focus on reducing the overall query evaluation time. The optimum is, generally, not reached.

Excerpt of the TPC-H Benchmark (at Scale Factor SF)

<u>o_orderkey</u>	<u>o_custkey</u>	<u>o_totalprice</u>	<u>o_clerk</u>	...
o	c			

orders ($\approx SF \times 1.5 \times 10^6$ rows)

<u>l_orderkey</u>	<u>l_linenum</u>	<u>l_partkey</u>	<u>l_quantity</u>	<u>l_extendedprice</u>	...
o					

lineitem ($\approx SF \times 6 \times 10^6$ rows)

<u>c_custkey</u>	<u>c_name</u>	<u>c_acctbal</u>	<u>c_nationkey</u>	...
c			n	

customer ($\approx SF \times 150000$ rows)

<u>n_nationkey</u>	<u>n_name</u>	<u>n_regionkey</u>	...
n		r	

nation (25 rows)

<u>r_regionkey</u>	<u>r_name</u>	...
r		

region (5 rows)



Q₁₄: Three-Way Join Against a TPC-H Instance

Price and quantity of parts ordered by customer #001:

```

SELECT 1.1_partkey, 1.1_quantity, 1.1_extendedprice
FROM   lineitem AS l JOIN orders AS o      -- } l ⋈ o
      ON (l.1_orderkey = o.o_orderkey)    -- }
      JOIN customer AS c                  -- } ⋈ c
      ON (o.o_custkey = c.c_custkey)      -- }
WHERE  c.c_name = 'Customer#001';

```

- Above SQL syntax suggests the **join order** $(l \bowtie o) \bowtie c$.
- Commutativity and associativity of \bowtie enable the RDBMS to **reorder** the joins—based on *estimated evaluation costs*.
 - ... unless we insist on the syntactic order. 🧐



2 | Pre-Processing: Query Normalization

Transform the input SQL query such that it features **SELECT-FROM-WHERE** (SFW) blocks of the following shape:

```

SELECT [ DISTINCT ] e, ..., e
FROM       $\Delta$ , ...,  $\Delta$            --  $\Delta \equiv$  base table or (query)
[ WHERE    p AND ... AND p ]      -- p  $\equiv$  predicate in DNF
[ GROUP BY g, ..., g                -- { e, p, g, o  $\equiv$ 
[ HAVING   p AND ... AND p ] ]    --   atomic expression or
[ ORDER BY o, ..., o ]              --   scalar (subquery)
[ OFFSET   n ]                       -- { n, m  $\equiv$  integer literal
[ LIMIT    m ]                       -- }
```

- Query clauses in [...] may be missing.

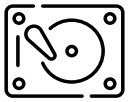


3 | Pre-Processing: Query Unnesting

Nested SQL queries suggest a (naïve, inefficient) nested-loop-style evaluation strategy. Consider:

<div style="display: flex; justify-content: space-between;"> SELECT c.c_name 1 </div> <div style="display: flex; justify-content: space-between;"> FROM customer AS c, </div> <div style="display: flex; justify-content: space-between;"> \triangle { (SELECT n.n_nationkey, n.n_name </div> <div style="display: flex; justify-content: space-between;"> { FROM nation AS n) AS t </div> <div style="display: flex; justify-content: space-between;"> WHERE c.c_nationkey = t.n_nationkey </div> <div style="display: flex; justify-content: space-between;"> AND strpos(c.c_address, t.n_name) > 0 </div>	<div style="display: flex; justify-content: space-between;"> SELECT o.o_orderkey 2 </div> <div style="display: flex; justify-content: space-between;"> FROM orders AS o </div> <div style="display: flex; justify-content: space-between;"> WHERE o.o_custkey IN </div> <div style="display: flex; justify-content: space-between;"> \triangle { (SELECT c.c_custkey </div> <div style="display: flex; justify-content: space-between;"> { FROM customer AS c </div> <div style="display: flex; justify-content: space-between;"> WHERE c.c_name = '...') </div>
---	--

- 💡 If possible, **unnest** \triangle queries and “inline” into parent query \Rightarrow \triangle can participate in join reordering.



Pre-Processing: Query Unnesting

Perform **query unnesting** on the level of

- the operator-based plan representation of the query,² or
- the internal AST representation of SQL. Re **2**:

<pre> SELECT e₁ FROM q₁, ..., q_i WHERE p₁ AND e₂ IN (SELECT e₃ FROM q_{i+1}, ..., q_n WHERE p₃) </pre>	\equiv^*	<pre> SELECT DISTINCT e₁ FROM q₁, ..., q_i, q_{i+1}, ..., q_n WHERE p₁ AND e₂ = e₃ AND p₃ </pre>
--	------------	---

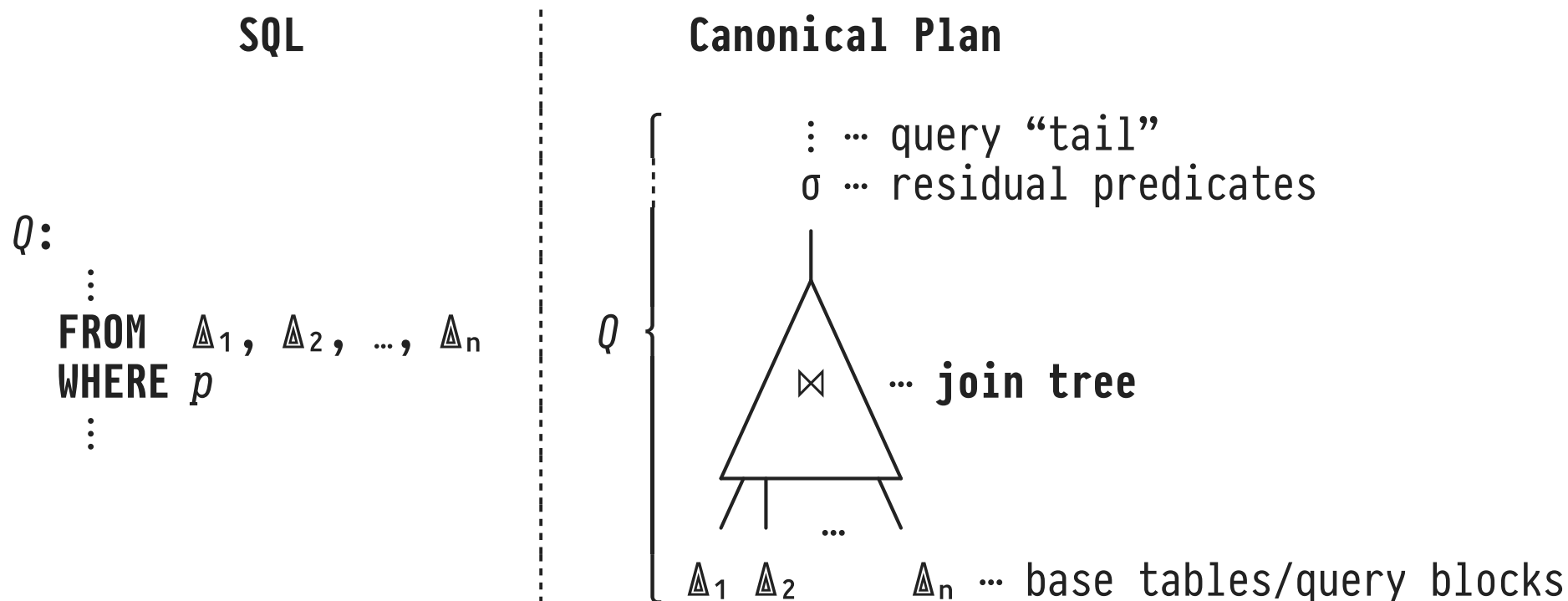
* Precondition: e_1 is key in the left-hand side query

² See *Unnesting Arbitrary Queries*, Thomas Neumann, Alfons Kemper. BTW 2015, Hamburg, Germany.



4 | Join Tree Optimization

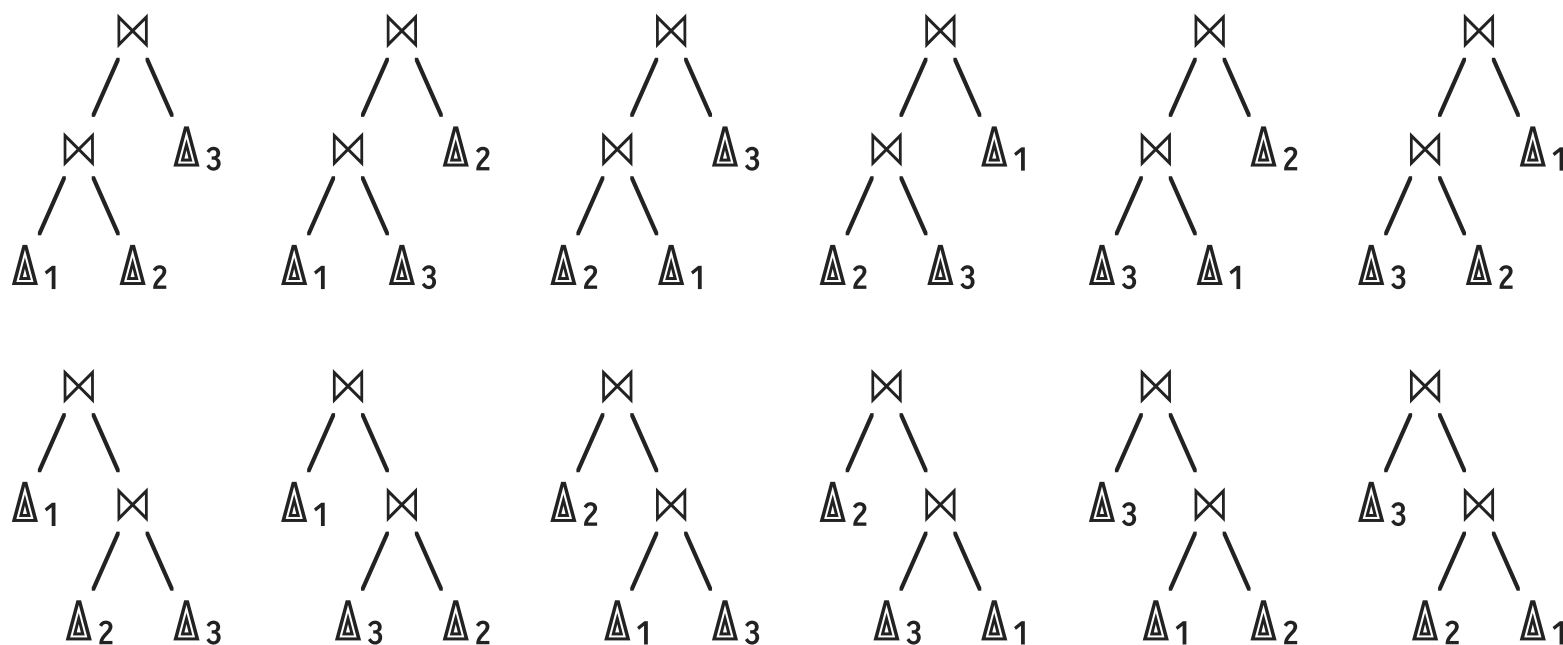
Processing a SQL query Q starts out with its **FROM** and **WHERE** clauses which describe a **join tree** over Q 's inputs:





Join Tree Optimization

Given n join inputs, the number of possible **join tree shapes** is *huge*. Consider $n = 3$:

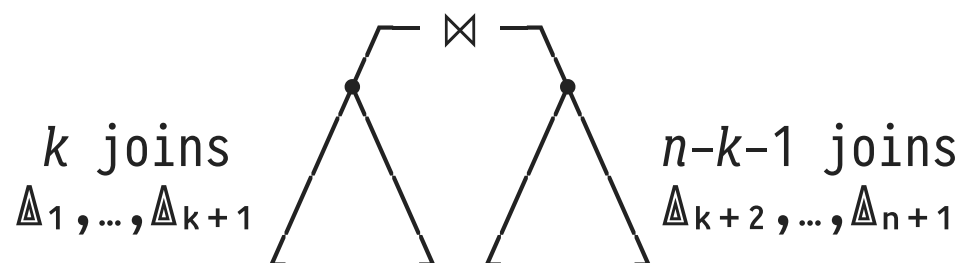


- Shapes based on associativity and commutativity of \Join .



How Many Possible Join Trees Are There?

1. A join of $n+1$ inputs Δ requires n binary joins. The root \bowtie combines subtrees of k and $n-k-1$ joins ($0 \leq k \leq n-1$):³



of join tree shapes:

$$C_n = \sum_{k=0}^{n-1} C_k \times C_{n-k-1}$$

2. Orderings of the Δ at the join tree leaf level: $(n+1)!$.
3. Join algorithm choices (a available algorithms): a^n .

³ C_n are the *Catalan numbers*, the number of ordered binary trees with $n+1$ leaves. $C_0 = 1$.



How Many Possible Join Trees Are There?

Number of possible join trees given n binary joins with $a = 3$ implementation choices:

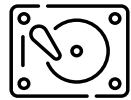
# of Δ ($n+1$)	C_n	# of join trees
2	1	6
3	2	108
4	5	3240
5	14	136080
6	42	7384320
7	132	484989120
8	429	37829151360
9	1430	3404623622400
10	4862	347271609484800

- A search space of this size is impossible to fully explore for any query optimizer.



Join Plan Generation Through Dynamic Programming

- **Problem:** Find optimal query plan $opt[\{\Delta_1, \dots, \Delta_n\}]$ that joins n inputs $\Delta_1, \dots, \Delta_n$.
 1. **Iteration 1:** For each Δ_j , find and memorize **best 1-input plan** $opt[\{\Delta_j\}]$ that accesses Δ_j only.
 2. **Iteration $k > 1$:** Find and memorize **best k -input plans** that join $k \leq n$ inputs by combining (for $1 \leq i < k$)
 - the best i -input plans and \setminus simple lookups in
 - the best $(k-i)$ -input plans. $\int opt[\cdot]$ memo 👍



Bottom-Up Dynamic Programming ($n = 3$)

k Possible k -input Access/Join Plans if Δ_i is complex

1 $opt[\{\Delta_1\}] \leftarrow prune(\{Seq\ Scan\ \Delta_1, Index\ Scan\ \Delta_1, Bitmap\ Scan\ \Delta_1, \overbrace{\Delta_1}^{complex}\})$
 $opt[\{\Delta_2\}] \leftarrow prune(\{Seq\ Scan\ \Delta_2, Index\ Scan\ \Delta_2, Bitmap\ Scan\ \Delta_2, \Delta_2\})$
 $opt[\{\Delta_3\}] \leftarrow prune(\{Seq\ Scan\ \Delta_3, Index\ Scan\ \Delta_3, Bitmap\ Scan\ \Delta_3, \Delta_3\})$

2 $opt[\{\Delta_1, \Delta_2\}] \leftarrow prune(opt[\{\Delta_1\}] \otimes opt[\{\Delta_2\}])$
 $opt[\{\Delta_1, \Delta_3\}] \leftarrow prune(opt[\{\Delta_1\}] \otimes opt[\{\Delta_3\}])$
 $opt[\{\Delta_2, \Delta_3\}] \leftarrow prune(opt[\{\Delta_2\}] \otimes opt[\{\Delta_3\}])$

3 $opt[\{\Delta_1, \Delta_2, \Delta_3\}] \leftarrow prune($
 $\quad opt[\{\Delta_1\}] \otimes opt[\{\Delta_2, \Delta_3\}] \cup$
 $\quad opt[\{\Delta_2\}] \otimes opt[\{\Delta_1, \Delta_3\}] \cup$
 $\quad opt[\{\Delta_3\}] \otimes opt[\{\Delta_1, \Delta_2\}] \quad)$

$prune(P) \equiv$ best (= minimal cost + interestingly ordered) plans in set P

$l \otimes r \equiv \{l \bowtie^{n1} r, r \bowtie^{n1} l, l \bowtie^{mj} r, r \bowtie^{mj} l, l \bowtie^{hj} r, r \bowtie^{hj} l\}$



Join Plan Generation (Notes)

- **Access plan choices** (*access(·)*):
 - Consider sequential/index scans if Δ is a base table, otherwise simply consume Δ 's rows.
- **Join plan choices** (*_ \Join _*):
 - Considers all viable join algorithms (given θ , available indexes, ...) and left/right input orders.
- **Principle of Optimality** (*prune(·)*): A globally optimal plan is built from optimal subplans. Thus:
 - 💡 For each subset of $\{\Delta_1, \dots, \Delta_n\}$, memorize in *opt[·]*
 1. ... its overall best plan and
 2. ... its best plan satisfying each **interesting order**.



(Bushy) Join Plan Generation: Pseudo Code

```

JoinPlan( $\{\Delta_1, \dots, \Delta_n\}$ ):
  foreach  $p \in \{\Delta_1, \dots, \Delta_n\}$  } 1-input plans
  |  $opt[\{p\}] \leftarrow prune(access(p))$ ;

  for  $k$  in  $2, \dots, n$  }  $k$ -input plans
  |   foreach  $S \subseteq \{\Delta_1, \dots, \Delta_n\}$  with  $|S| = k$  } enumerate subsets
  |   |    $opt[S] \leftarrow \phi$ ;
  |   |   foreach  $T \subset S$  with  $T \neq \phi$   $\Join^a$ 
  |   |   |    $opt[S] \leftarrow opt[S] \cup \{ opt[T] \Join^a opt[S \setminus T] \}$ ;
  |   |   |    $opt[S] \leftarrow prune(opt[S])$ ;

return  $opt[\{\Delta_1, \dots, \Delta_n\}]$ ;

```

- $access(\cdot)$, $prune(\cdot)$ defined as above,
 \Join^a builds all join algorithm choices ($a \in \{nl, mj, hj\}$).



Reducing the Search Space

- Avoid generating costly **Cartesian products**: don't form joins between inputs w/o join predicate (`_ θ _ = true`).
- Generate **left-deep** join plans only: right join input (NL \bowtie : inner input) is a scan over base table *T*.
 - Admits use of Index Nested Loop Join.
 - Straightforward Volcano-style execution (reset inner).

