# DB 2

14 – Query Optimization

Summer 2018

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 **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*)

| **o_orderkey** | **o_custkey** | **o_totalprice** | **o_clerk** | … |
|:---:|:---:|:---:|:---:|:---:|
| *o* | *c* | | | |

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

| **l_orderkey** | **l_linenumber** | **l_partkey** | **l_quantity** | **l_extendedprice** | … |
|:---:|:---:|:---:|:---:|:---:|:---:|
| *o* | | | | | |

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

| **c_custkey** | **c_name** | **c_acctbal** | **c_nationkey** | … |
|:---:|:---:|:---:|:---:|:---:|
| *c* | | | *n* | |

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

| **n_nationkey** | **n_name** | **n_regionkey** | … |
|:---:|:---:|:---:|:---:|
| *n* | | *r* | |

nation (25 rows)

| **r_regionkey** | **r_name** | … |
|:---:|:---:|:---:|
| *r* | | |

region (5 rows)

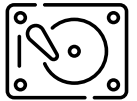## $Q_{14}$: Three-Way Join Against a TPC-H Instance

Price and quantity of parts orderd by customer #001:

```
SELECT l.l_partkey, l.l_quantity, l.l_extendedprice
FROM   lineitem AS l JOIN orders AS o        -- ⎱ l ⋈ o
         ON (l.l_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 ⋈ o) ⋈ c.
- Commutativity and associativity of ⋈ enable the RDBMS to **reorder** the joins—based on *estimated evaluation costs*.
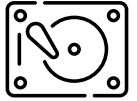    - ... unless we insist on the syntactic order. 🤖

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

```
    SELECT [DISTINCT] e, …, e
    FROM      ⏃, …, ⏃            -- ⏃ ≡ base table or (query)
[   WHERE     p AND ⋯ AND p ]      -- p ≡ predicate in DNF
[   GROUP BY  g, …, g              -- ⎤ e, p, g, o ≡
  [ HAVING    p AND ⋯ AND p ] ]    -- ⎬   atomic expression or
[   ORDER BY  o, …, o ]            -- ⎦   scalar (subquery)
[   OFFSET    n ]                  -- ⎱ n, m ≡ integer literal
[   LIMIT     m ]                  -- ⎰
```
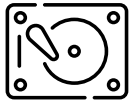
- Query clauses in [...] may be missing.

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

```
SELECT c.c_name                              ❶
FROM   customer AS c,
  ⚠ ⌠ (SELECT n.n_nationkey, n.n_name
    ⌡  FROM   nation AS n) AS t
WHERE c.c_nationkey = t.n_nationkey
  AND strpos(c.c_address, t.n_name) > 0
```

```
SELECT o.o_orderkey                  ❷
FROM   orders AS o
WHERE  o.o_custkey IN
   ⎰ (SELECT c.c_custkey
⚠ ⎱  FROM   customer AS c
   ⎱  WHERE  c.c_name = '…')
```

- 💡 If possible, **unnest** ⚠ queries and "inline" into parent query ⟹ ⚠ can participate in join reordering.

Perform **query unnesting** on the level of

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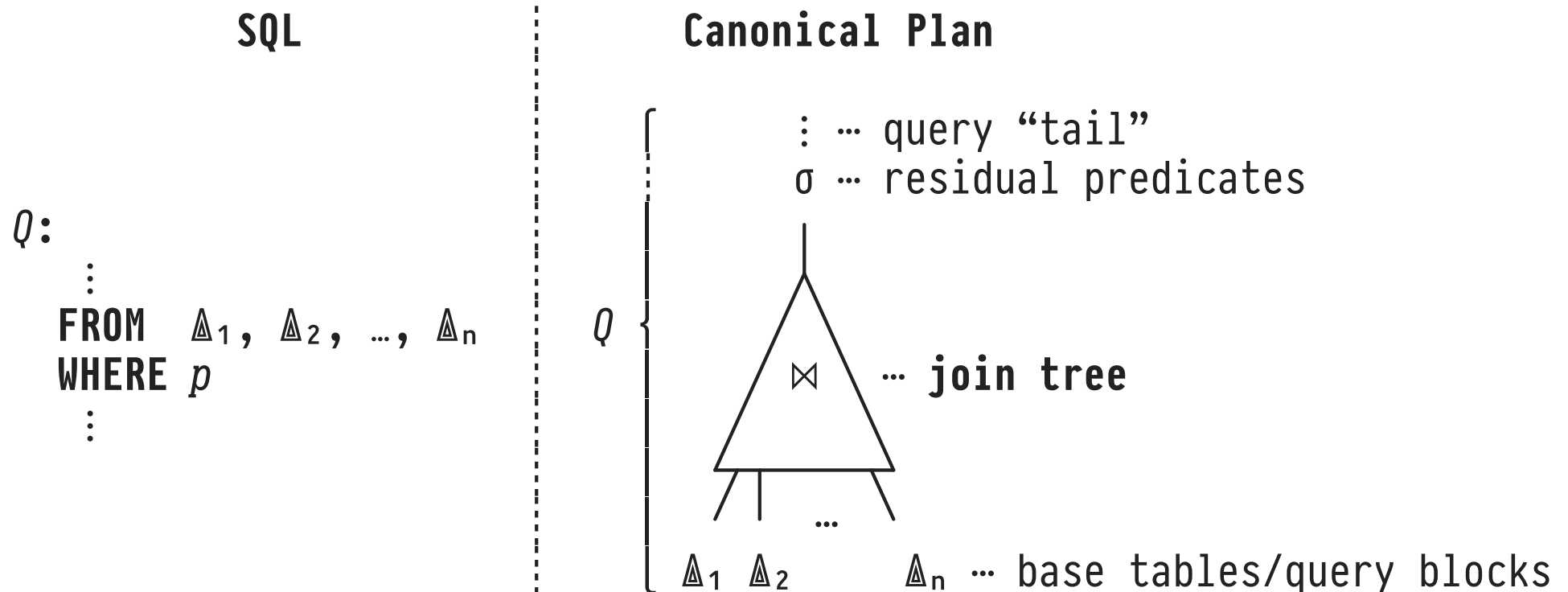- the internal AST representation of SQL. Re ❷:

| | |
|---|---|
| ```SELECT``` $e_1$ <br> ```FROM``` $q_1,\ldots,q_i$ <br> ```WHERE``` $p_1$ <br>   ```AND``` $e_2$ ```IN (SELECT``` $e_3$ <br>            ```FROM``` $q_{i+1},\ldots,q_n$ <br>            ```WHERE``` $p_3$```)``` | ⇒    ```SELECT DISTINCT``` $e_1$ <br> ```FROM``` $q_1,\ldots,q_i,q_{i+1},\ldots,q_n$ <br> ```WHERE``` $p_1$ <br>   ```AND``` $e_2 = e_3$ <br>   ```AND``` $p_3$ |

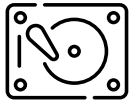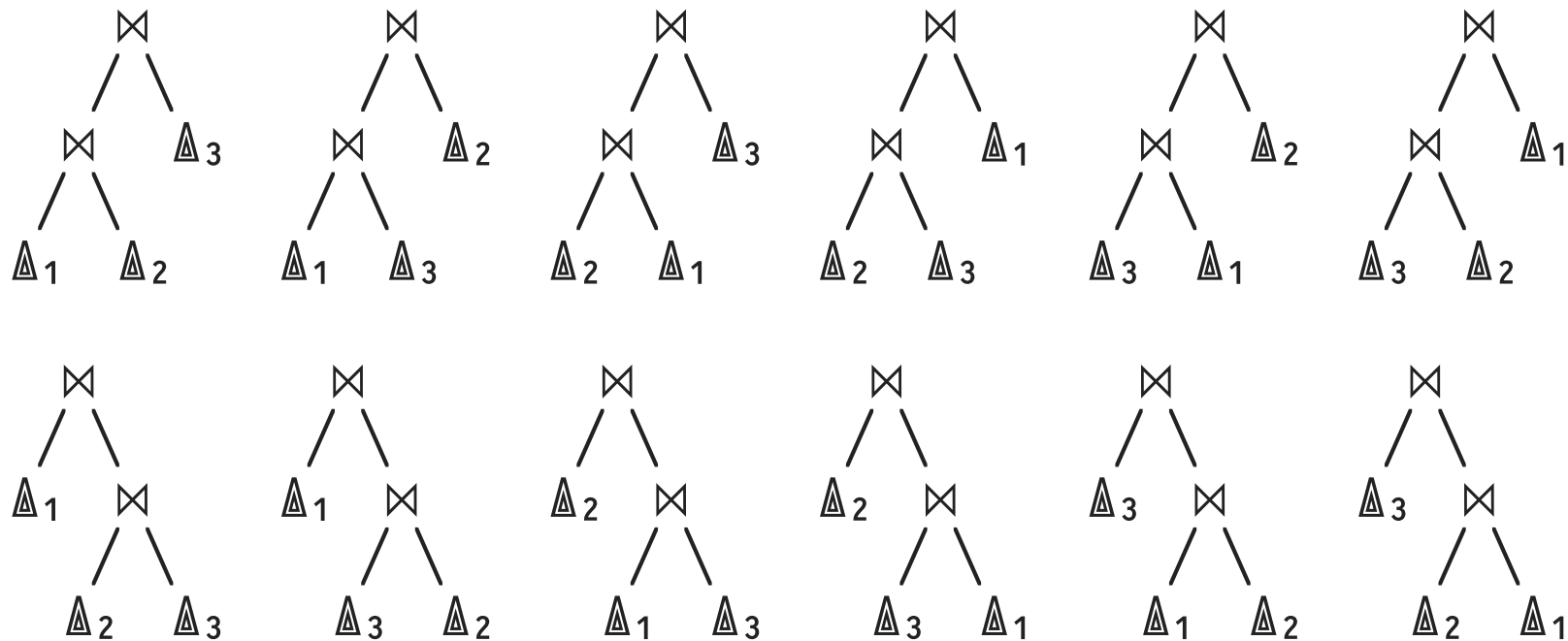[2] See *Unnesting Arbitrary Queries*, Thomas Neumann, Alfons Kemper. BTW 2015, Hamburg, Germany.

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

| SQL | Canonical Plan |
|---|---|

$Q$:
    ⋮
  **FROM**  △$_1$, △$_2$, …, △$_n$
  **WHERE** $p$
    ⋮

⋮ ⋯ query "tail"
σ ⋯ residual predicates

⋈ ⋯ **join tree**

△$_1$ △$_2$    △$_n$ ⋯ base tables/query blocks
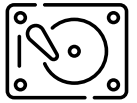
# Join Tree Optimization

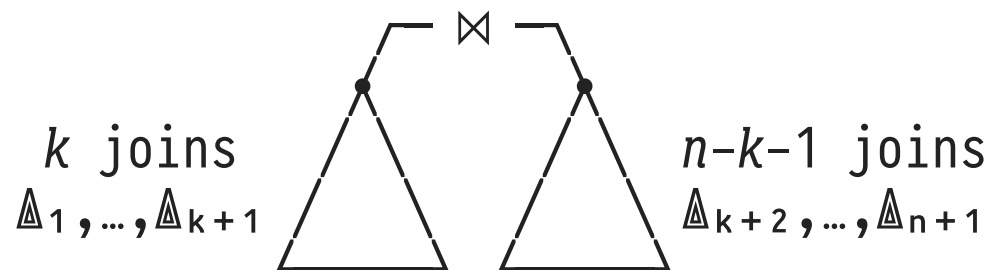Given $n$ join inputs, the number of possible **join tree shapes** is *huge*. Consider $n = 3$:



- Shapes based on associativity and commutativity of ⋈.

# How Many Possible Join Trees are There?

1. A join of $n+1$ inputs △ requires $n$ binary joins. The root ⋈ combines subtrees of $k$ and $n-k-1$ joins ($0 \leqslant k \leqslant n-1$):[3]

$k$ joins
$△_1, \ldots, △_{k+1}$

$n-k-1$ joins
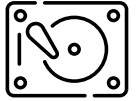$△_{k+2}, \ldots, △_{n+1}$

\# of join tree shapes:

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

2. Orderings of the △ 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. $C_0 = 1$.
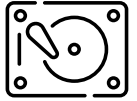
# How Many Possible Join Trees are There?

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

| # of ⚠ ($n$+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 $opt[\{⧌_1,\ldots,⧌_n\}]$ that joins $n$ inputs $⧌_1,\ldots,⧌_n$.

1. **Iteration 1:** For each $⧌_j$, find and memorize **best 1-input plan** $opt[\{⧌_j\}]$ that accesses $⧌_j$ only.

2. **Iteration $k > 1$:** Find and memorize **best $k$-input plans** that join $k \leqslant n$ inputs by combining (for $1 \leqslant i < k$)
   - the best $i$-input plans and ⎰ simple lookups in
   - the best $(k{-}i)$-input plans. ∫ $opt[\cdot]$ memo 👍

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

| $k$ | Possible $k$-input Access/Join Plans | if $A_i$ is complex |
|---|---|---|

**1**   $opt[\{A_1\}] \leftarrow prune(\{\text{Seq Scan } A_1, \text{ Index Scan } A_1, \text{ Bitmap Scan } A_1, A_1\})$
$opt[\{A_2\}] \leftarrow prune(\{\text{Seq Scan } A_2, \text{ Index Scan } A_2, \text{ Bitmap Scan } A_2, A_2\})$
$opt[\{A_3\}] \leftarrow prune(\{\text{Seq Scan } A_3, \text{ Index Scan } A_3, \text{ Bitmap Scan } A_3, A_3\})$

---

**2**   $opt[\{A_1,A_2\}] \leftarrow prune(opt[\{A_1\}] \otimes opt[\{A_2\}])$
$opt[\{A_1,A_3\}] \leftarrow prune(opt[\{A_1\}] \otimes opt[\{A_3\}])$
$opt[\{A_2,A_3\}] \leftarrow prune(opt[\{A_2\}] \otimes opt[\{A_3\}])$

---

**3**   $opt[\{A_1,A_2,A_3\}] \leftarrow prune(opt[\{A_1\}] \otimes opt[\{A_2,A_3\}] \cup$
$opt[\{A_2\}] \otimes opt[\{A_1,A_3\}] \cup$
$opt[\{A_3\}] \otimes opt[\{A_1,A_2\}] \ )$

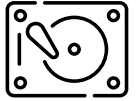$prune(P) \equiv$ best (= minimal cost + interestingly ordered) plans in set $P$
$l \otimes r \equiv \{l \bowtie^{nl} r, r \bowtie^{nl} l, \quad l \bowtie^{mj} r, r \bowtie^{mj} l, \quad l \bowtie^{hj} r, r \bowtie^{hj} l\}$

# Join Plan Generation (Notes)

- **Access plan choices** (*access(•)*):
  - Consider sequential/index scans if ▥ is a base table, otherwise simply consume ▥'s rows.

- **Join plan choices** (_ ⊛ _):
  - Considers all viable join algorithms (given $\theta$, 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 $\{▥_1, …, ▥_n\}$, memorize in *opt[•]*
    1. … its overall best plan and
    2. … its best plan satisfying each **interesting order**.

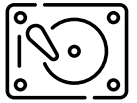# (Bushy) Join Plan Generation: Pseudo Code
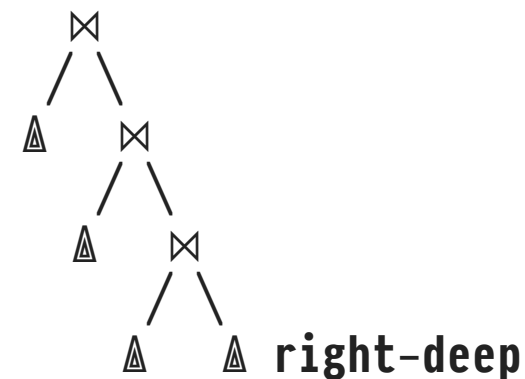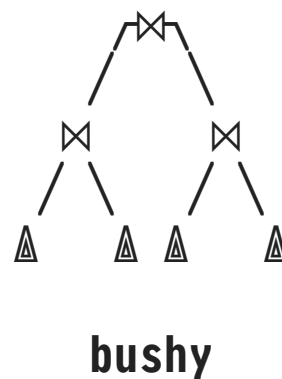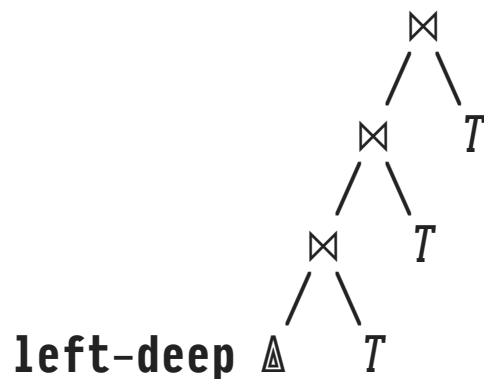
```
JoinPlan({𝔸₁,…,𝔸ₙ}):
 foreach p ∈ {𝔸₁,…,𝔸ₙ}                    } 1-input plans
 └ opt[{p}] ← prune(access(p));

 for k in 2,…,n                           } k-input plans
 │  foreach S ⊆ {𝔸₁,…,𝔸ₙ} with |S| = k } enumerate subsets
 │  │  opt[S] ← ∅;
 │  │  foreach T ⊂ S with T ≠ ∅        ⌐⋈ᵃ⌐
 │  │  └ opt[S] ← opt[S] ∪ { opt[T]    opt[S \ T] };
 │  └  opt[S] ← prune(opt[S]);

 return opt[{𝔸₁,…,𝔸ₙ}];
```

- *access*(•), *prune*(•) defined as above,
  ⌐⋈ᵃ⌐ builds all join algorithm choices ($a \in$ {nl,mj,hj}).
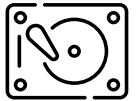
- Avoid generating costly **Cartesian products:** don't form joins between inputs w/o join predicate (_ $\theta$ _ = true).

- Generate **left-deep** join plans only: right join input (NL⋈: inner input) is a scan over base table $T$.
  - Admits use of Index Nested Loop Join.
  - Straightforward Volcano-style execution (reset inner).



left-deep          bushy          right-deep

The query optimizer explores the vast plan search space to find the **optimal ("best", "cheapest") plan.**

- Typically, RDBMSs measure **plan cost** in terms of *total execution time* (time until last result row delivered).
- These total plan costs are **estimated** *before* plan execution begins (EXPLAIN: … cost=$c_1$..$c_2$↩ …).
- A **cost model**—measured in abstract "space\$"—reflects the true costs (measured in *ms*, CPU time, # I/O ops, …) of plans $p_1$, $p_2$:

space\$$(p_1)$ < space\$$(p_2)$ $\Rightarrow$ true cost$(p_1)$ < true cost$(p_2)$

EXPLAIN shows estimated costs (unit: space\$) and cardinalities (# of rows):

```
                              QUERY PLAN
──────────────────────────────────────────────────────────────────
              startup cost        total cost

                    ▼                  ▼
  Hash Join  (cost=299.00..15443.31 rows=505183 width=50)
    ⋮                                     ▲
                                      cardinality
```

- **run cost** $\overset{\text{def}}{=}$ **total cost** – **startup cost**[4] (not shown).
- Optimizer decisions are based on estimated **total cost.**

[4] To implement set enable_<*op*> = off, PostgreSQL sets the operator's **startup cost** to $10^9$ ($\equiv \infty$).
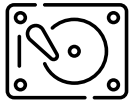
# Cost Model Configuration

| Model Configuration | Default | Description |
|---|---|---|
| seq_page_cost | 1.0 | I/O cost of one sequential page access |
| random_page_cost | 4.0 | I/O cost of one random page access |
| cpu_tuple_cost | 0.01 | CPU cost to process a heap file row |
| cpu_index_tuple_cost | 0.005 | CPU cost to process an index leaf entry |
| cpu_operator_cost | 0.0025 | CPU function/operator evaluation cost |
| parallel_tuple_cost | 0.1 | Cost of passing one row worker→leader |
| parallel_setup_cost | 1000.0 | Shared memory setup cost |

- Parameters are configurable:
  - Seek cost, thus random_page_cost » seq_page_cost. But...
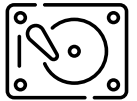  - ... if DB fits in RAM, random_page_cost = seq_page_cost may be more appropriate.

Given an occurrence of Seq Scan with arguments

- *in*: input table,
- *pred*: (optional) filter predicate on *in*,
- *expr*: SELECT clause expression(s),

how does PostgreSQL derive *startup_cost* and *total_cost*?

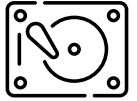|  | *in* | QUERY PLAN | *total_cost* |  |
|---|---|---|---|---|
| Seq Scan on public.indexed i | | (cost=0.00..22.75 rows=100 width=4) | | |
| Output: (a + 1) ⟵ *expr* | | | | |
| Filter: (i.a <= 100) ⟵ *pred* | | *startup_cost* | | #rows(*out*) |

# Cost of **Seq Scan** ❷

Cost calculation depends on the following parameters, mostly available in PostgreSQL's internal `pg_*` meta data tables:

| Parameter | Description | Available as… |
|---|---|---|
| #rows(*in*) | # rows (cardinality) of table *in* | pg_class.reltuples |
| #pages(*in*) | # pages in heap file of *in* | pg_class.relpages |
| sel(*pred*) | selectivity of filter *pred*[5] | see below |

- Meta data like #rows(*in*), #pages(*in*) and others are updated whenever the system performs an ANALYZE run on table *in*.
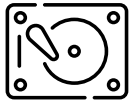- Predicate selectivity sel(*pred*) is estimated based on sampled table data and the syntactic structure of *pred*.

[5] sel(*pred*) ∈ {0,…,1} with sel(*pred*) = 0 ≡ no row satisfies filter *pred*.

$$\textbf{startup\_cost} \overset{\text{def}}{=} \overset{\lrcorner \quad \text{typically} = 0 \quad \lrcorner}{\text{startup\_cost}(\textit{pred}) + \text{startup\_cost}(\textit{expr})}$$

$$\text{cpu\_run\_cost} \overset{\text{def}}{=} \#\text{rows}(\textit{in}) \times (\overbrace{\text{cpu\_tuple\_cost}}^{\text{decode heap row}} + \overbrace{\text{run\_cost}(\textit{pred})}^{\text{evaluate filter}})$$

$$+ \underbrace{\#\text{rows}(\textit{in}) \times \text{sel}(\textit{pred})}_{= \#\text{rows}(\textit{out})} \times \underbrace{\text{run\_cost}(\textit{expr})}_{\text{evaluate SELECT clause}}$$

$$\text{disk\_run\_cost} \overset{\text{def}}{=} \underbrace{\#\text{pages}(\textit{in}) \times \text{seq\_page\_cost}}_{\text{sequentially read entire input heap file}}$$

$$\textbf{total\_cost} \overset{\text{def}}{=} \text{startup\_cost} + \underbrace{\text{cpu\_run\_cost} + \text{disk\_run\_cost}}_{= \textbf{run\_cost}}$$

# Cost of Index Scan ❶

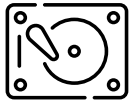Modeling the cost for an Index Scan has to reflect that *two* data structures (heap file & B⁺Tree) are involved:

```
                    idx          in      QUERY PLAN
                     ▼           ▼
  Index Scan using indexed_a on indexed i (cost=0.42..443.12 rows=10885 …
      Output: (c + '1'::numeric) ◀ expr                              ▲
      Index Cond: (i.a <= 10000) ◀ pred                     #rows(out)
```

The model separately accounts for

1.  the B⁺Tree descent (startup of the Index Scan),
2.  the index leaf level scan, and
3.  heap file access (clustered *vs*. non-clustered).

Cost model parameters:
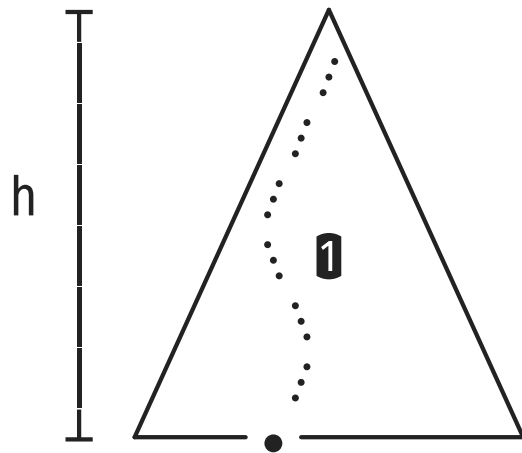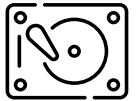
| Parameter | Description | Available as… |
|-----------|-------------|---------------|
| #rows($in$) | # rows (cardinality) of table $in$ | pg_class.reltuples |
| #pages($in$) | # pages in heap file of $in$ | pg_class.relpages |
| sel($pred$) | selectivity of filter $pred$ | see below |
| h($idx$) | height of B⁺Tree $idx$ | bt_metap($\cdot$) |
| #rows($idx$) | # leaf entries in index $idx$ | pg_class.reltuples |
| #pages($idx$) | # pages in leaf level of $idx$ | pg_class.relpages |
| corr($idx$) | $\approx$ clustering factor for index $idx$ | pg_stats.correlation |

- corr($idx$) $\in$ {−1.0,…,1.0} characterizes how much the physical orderings of index leaves and heap file deviate.
  - After CLUSTER $in$ ON $idx$, we have corr($idx$) = 1.0.

- B⁺Tree height $h = \log_{2 \times o}(\#rows(idx))$

⇒ # of key comparisons during B⁺Tree descent **❶**:

$$\underbrace{\lceil \log_2(2 \times o) \times h\rceil}_{} = \lceil \log_2(\#rows(idx))\rceil$$

binary search in inner B⁺Tree node with fan-out $F = 2 \times o$

---

**startup_cost** $\overset{\text{def}}{=}$
    startup_cost(*pred*) + startup_cost(*expr*)
 + ($\underbrace{\lceil \log_2(\#rows(idx))\rceil}_{\text{B⁺Tree descent}}$ + ($\underbrace{(h + 1)}_{\therefore + \bullet}$ × 50) × cpu_operator_cost
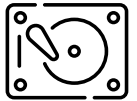
        B⁺Tree descent         ∴ + ●    index node processing

The index leaf level (sequence set) scan ❷ incurs CPU as well as I/O cost that contribute to the overall **run_cost:**

# rows in scanned range $|\cdots\blacktriangleright|$

$$\text{index\_cpu\_cost} \overset{\text{def}}{=} \underbrace{\text{sel}(\textit{pred}) \times \text{\#rows}(\textit{idx})}_{} \times (\underbrace{\text{cpu\_index\_tuple\_cost}}_{\text{decode index leaf entry}} + \underbrace{\text{run\_cost}(\textit{pred})}_{\text{evaluate} \leqslant \textit{hi}})$$

$|\cdots\blacktriangleright|$
$lo \quad ❷ \quad hi$

\# of pages • in scanned range       $B^+$Tree leaves not clustered

$$\text{index\_IO\_cost} \overset{\text{def}}{=} \underbrace{\lceil\text{sel}(\textit{pred}) \times \text{\#pages}(\textit{idx})\rceil}_{} \times \underbrace{\text{random\_page\_cost}}_{}$$

Unless we perform an Index **Only** Scan, we additionally pay
CPU and I/O costs for heap file accesses ❸:
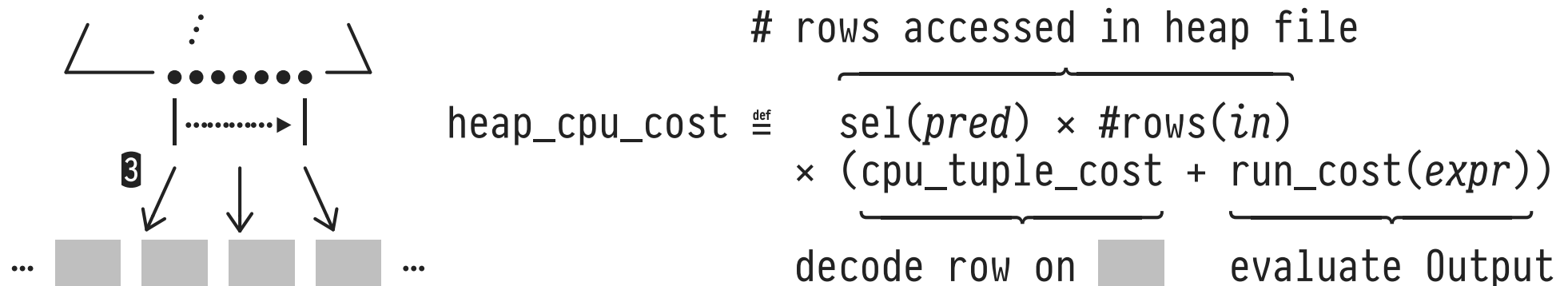


$$
\text{heap\_cpu\_cost} \stackrel{\text{def}}{=} \overbrace{\text{sel}(\textit{pred}) \times \text{\#rows}(\textit{in})}^{\text{\# rows accessed in heap file}}
$$

$$
\times \; (\underbrace{\text{cpu\_tuple\_cost}}_{\text{decode row on } \blacksquare} + \underbrace{\text{run\_cost}(\textit{expr})}_{\text{evaluate Output}})
$$

- The more **clustered** the index, the cheaper the heap I/O.  Linearly
  interpolate between the clustered and non-clustered scenarios:

$$
\text{heap\_IO\_cost} \stackrel{\text{def}}{=} \text{unclustered\_IO\_cost}
$$
$$
+ \underbrace{\text{corr}(\textit{idx})^2}_{\approx \text{ clustering factor} \in \{0,\dots,1\}} \times (\text{clustered\_IO\_cost} - \text{unclustered\_IO\_cost})
$$

first
access
is random ········

heap file

$$\text{clustered\_IO\_cost} \; \stackrel{\text{def}}{=} \; 1 \times \text{random\_page\_cost} \; (\blacksquare)$$
$$+ \; (\text{sel}(\textit{pred}) \times \#\text{pages}(\textit{in}) - 1) \times \text{seq\_page\_cost} \; (\blacksquare)$$

$$\text{unclustered\_IO\_cost} \; \stackrel{\text{def}}{=} \; \#\text{pages}(\textit{in}) \times \text{random\_page\_cost} \; (\blacksquare)$$

$$\textbf{total\_cost} \stackrel{\text{def}}{=} \textbf{startup\_cost} + \text{index\_cpu\_cost} + \text{index\_IO\_cost}$$
$$+ \; \text{heap\_cpu\_cost} \; + \text{heap\_IO\_cost}$$