# DB 2

---

## 14 – Query Optimization

**Summer 2024**

**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*)

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

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

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

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

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

customer (≈ $SF \times 150000$ rows)

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

nation (25 rows)

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

region (5 rows)

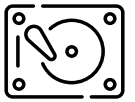# $Q_{14}$: Three-Way Join Against a TPC-H Instance

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. 🤖
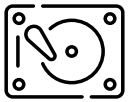
## 2 ┊ Pre-Processing: Query Normalization

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

```
   SELECT [DISTINCT] e, …, e
   FROM       ◬, …, ◬          -- ◬ ≡ 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 ]                 -- ⎰
```
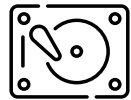
- Query clauses in […] may be missing.

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

```
SELECT c.c_name                                  ❶    SELECT o.o_orderkey                    ❷
FROM   customer AS c,                                 FROM    orders AS o
  ⚠ ⎧ (SELECT n.n_nationkey, n.n_name                 WHERE  o.o_custkey IN
     ⎩  FROM    nation AS n) AS t                        ⎧ (SELECT c.c_custkey
WHERE c.c_nationkey = t.n_nationkey              ⚠  ⎨   FROM    customer AS c
  AND strpos(c.c_address, t.n_name) > 0             ⎩    WHERE  c.c_name = '…')
```

- 💡 If possible, **unnest** ⚠ queries and "inline" into parent query ⇒ ⚠ 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 ❷:

```
SELECT  e₁                              SELECT DISTINCT e₁
FROM    q₁,…,qᵢ                         FROM    q₁,…,qᵢ,qᵢ₊₁,…,qₙ
WHERE   p₁                      ≅       WHERE   p₁
  AND   e₂ IN (SELECT e₃                  AND   e₂ = e₃
              FROM   qᵢ₊₁,…,qₙ            AND   p₃
              WHERE  p₃)
```

$$\text{SELECT } e_1$$
$$\text{FROM } q_1,\ldots,q_i$$
$$\text{WHERE } p_1$$
$$\text{AND } e_2 \text{ IN (SELECT } e_3$$
$$\text{FROM } q_{i+1},\ldots,q_n$$
$$\text{WHERE } p_3)$$

$$\overset{*}{\cong}$$

$$\text{SELECT DISTINCT } e_1$$
$$\text{FROM } q_1,\ldots,q_i,q_{i+1},\ldots,q_n$$
$$\text{WHERE } p_1$$
$$\text{AND } e_2 = e_3$$
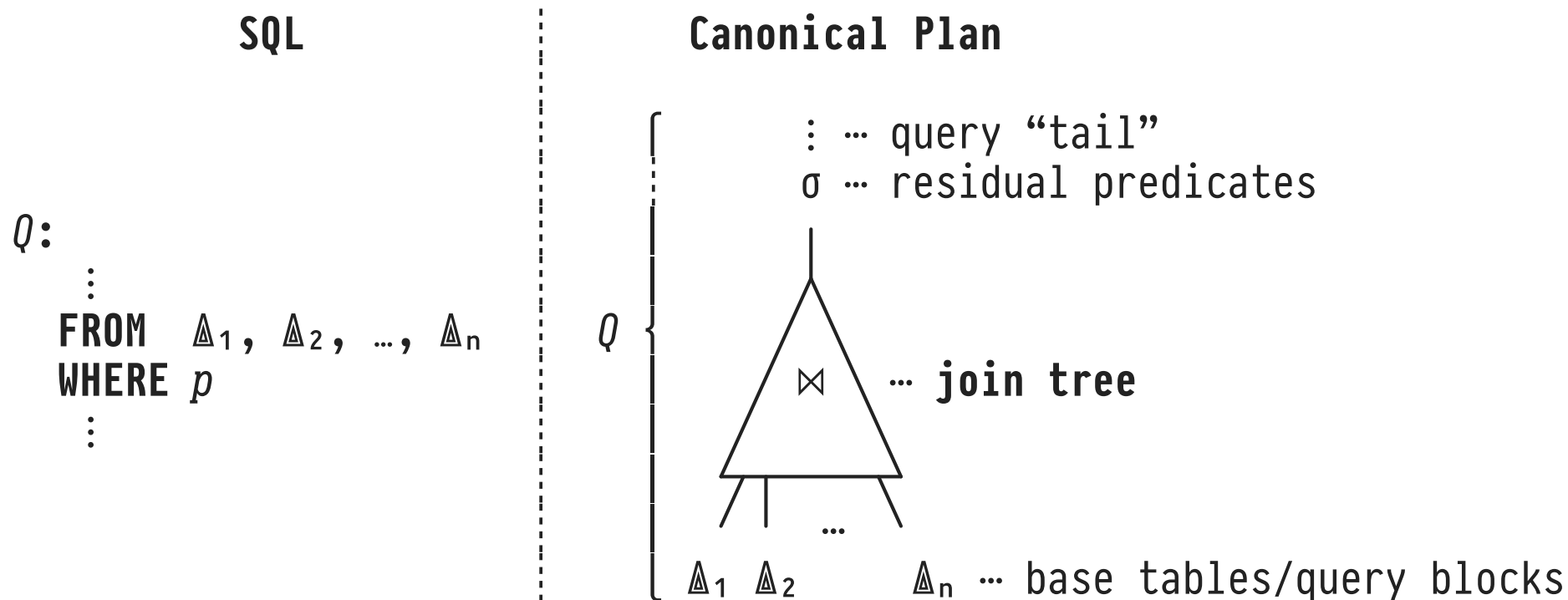$$\text{AND } p_3$$

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

[2] See *Unnesting Arbitrary Queries*, Thomas Neumann, Alfons Kemper. BTW 2015, Hamburg, Germany.
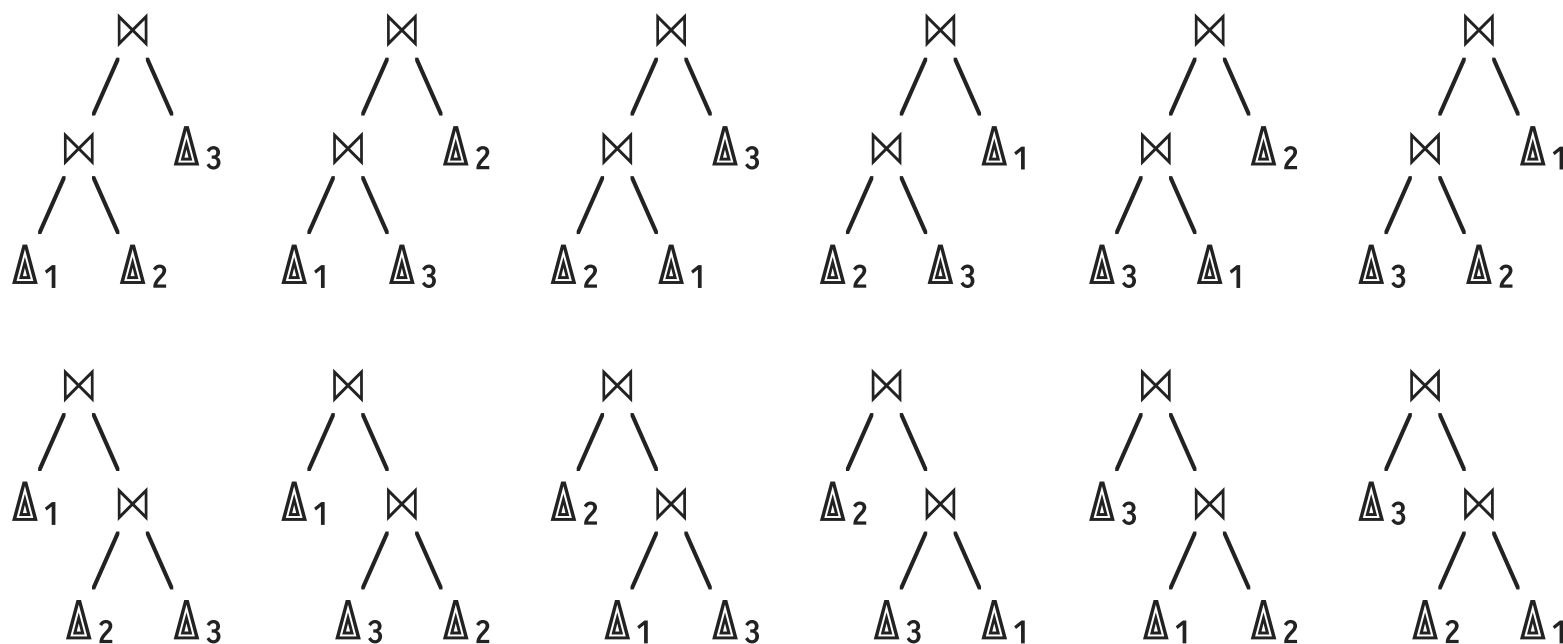
# 4 ⋮ Join Tree Optimization

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

**SQL**

**Canonical Plan**

$Q$:
⋮
**FROM** $\triangle_1$, $\triangle_2$, …, $\triangle_n$
**WHERE** $p$
⋮

⋮ ⋯ query "tail"
σ ⋯ residual predicates

$Q$

⋈ ⋯ **join tree**

…

$\triangle_1$ $\triangle_2$     $\triangle_n$ ⋯ base tables/query blocks

# Join Tree Optimization

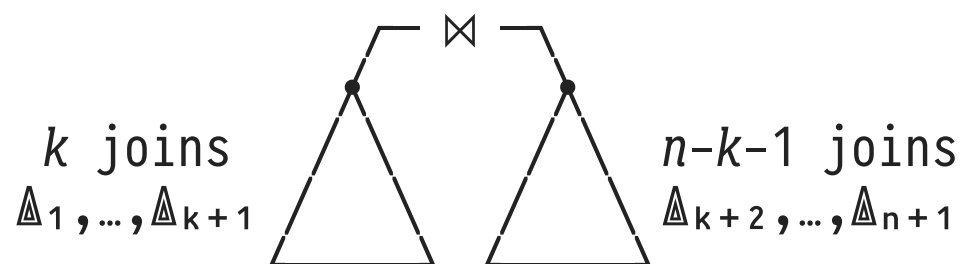Given $n$ join inputs, the number of possible **join tree shapes** is *huge*. Consider $n = 3$:



- Shapes based on associativity and commutativity of ⋈.

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
△$_1$,…,△$_{k+1}$

$n-k-1$ joins
△$_{k+2}$,…,△$_{n+1}$

\# 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.

- **Problem:** Find optimal query plan $opt[\{▲_1,…,▲_n\}]$ that joins $n$ inputs $▲_1,…,▲_n$.

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. ⎱ $opt[·]$ memo 👍

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

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

**1**  $opt[\{\triangle_1\}] \leftarrow prune(\{$Seq Scan $\triangle_1$, Index Scan $\triangle_1$, Bitmap Scan $\triangle_1$, $\triangle_1\})$
$opt[\{\triangle_2\}] \leftarrow prune(\{$Seq Scan $\triangle_2$, Index Scan $\triangle_2$, Bitmap Scan $\triangle_2$, $\triangle_2\})$
$opt[\{\triangle_3\}] \leftarrow prune(\{$Seq Scan $\triangle_3$, Index Scan $\triangle_3$, Bitmap Scan $\triangle_3$, $\triangle_3\})$

**2**  $opt[\{\triangle_1,\triangle_2\}] \leftarrow prune(opt[\{\triangle_1\}] \circledast opt[\{\triangle_2\}])$
$opt[\{\triangle_1,\triangle_3\}] \leftarrow prune(opt[\{\triangle_1\}] \circledast opt[\{\triangle_3\}])$
$opt[\{\triangle_2,\triangle_3\}] \leftarrow prune(opt[\{\triangle_2\}] \circledast opt[\{\triangle_3\}])$

**3**  $opt[\{\triangle_1,\triangle_2,\triangle_3\}] \leftarrow prune(opt[\{\triangle_1\}] \circledast opt[\{\triangle_2,\triangle_3\}] \cup$
$opt[\{\triangle_2\}] \circledast opt[\{\triangle_1,\triangle_3\}] \cup$
$opt[\{\triangle_3\}] \circledast opt[\{\triangle_1,\triangle_2\}] \quad )$

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

- **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**.

```
JoinPlan({▲₁,…,▲ₙ}):
 foreach p ∈ {▲₁,…,▲ₙ}                        } 1-input plans
  ⌊ opt[{p}] ← prune(access(p));

 for k in 2,…,n                               } k-input plans
  │ foreach S ⊆ {▲₁,…,▲ₙ} 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[{▲₁,…,▲ₙ}];
```

- *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 $p_1$, $p_2$:

$$\text{space}\$(p_1) < \text{space}\$(p_2) \Rightarrow \text{true cost}(p_1) < \text{true cost}(p_2)$$

# PostgreSQL: Plan Cost

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 $10^{10}$ ($\equiv \infty$).

# 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 | Cost of spawning a parallel worker |

- 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.

# Cost of **Seq Scan** 1

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** ②

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*.

**startup_cost** $\stackrel{\text{def}}{=}$ startup_cost(*pred*) + startup_cost(*expr*)

↳   typically = 0   ↴

cpu_run_cost $\stackrel{\text{def}}{=}$ #rows(*in*) × (cpu_tuple_cost + run_cost(*pred*))
+ #rows(*in*) × sel(*pred*) × run_cost(*expr*)

decode heap row    evaluate filter

= #rows(*out*)    evaluate SELECT clause

disk_run_cost $\stackrel{\text{def}}{=}$ #pages(*in*) × seq_page_cost

sequentially read entire input heap file

**total_cost** $\stackrel{\text{def}}{=}$ startup_cost + cpu_run_cost + disk_run_cost

= **run_cost**

# Cost of Index Scan ▫

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
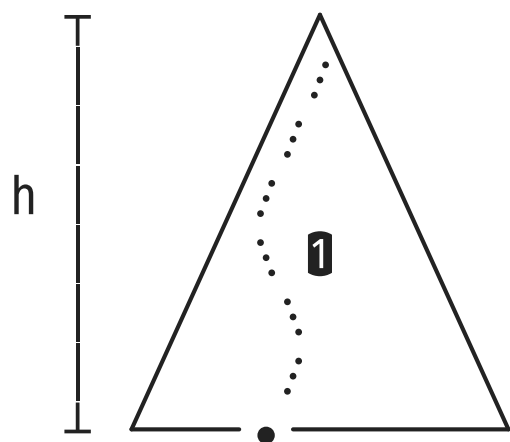3.  heap file access (clustered *vs*. non-clustered).

# Cost of Index Scan 2

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(·) |
| #rows(*idx*) | # leaf entries in index *idx* | pg_class.reltuples |
| #pages(*idx*) | # pages in leaf level of *idx* | pg_class.relpages |
| corr(*idx*) | ≈ clustering factor for index *idx* | pg_stats.correlation |

- corr(*idx*) ∈ {−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:

$$\lceil \log_2(2 \times o) \times h \rceil = \lceil \log_2(\#rows(idx)) \rceil$$

$\underbrace{\log_2(2 \times o)}$

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*)
 + ($\underbrace{\lceil \log_2(\#rows(idx)) \rceil}_{\text{B⁺Tree descent}}$ + $\underbrace{(h + 1)}_{\vdots\ +\ \bullet}$ × $\underbrace{50}$) × cpu_operator_cost
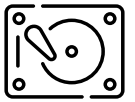
          B⁺Tree descent         $\vdots$ + ●     index node processing

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

$$
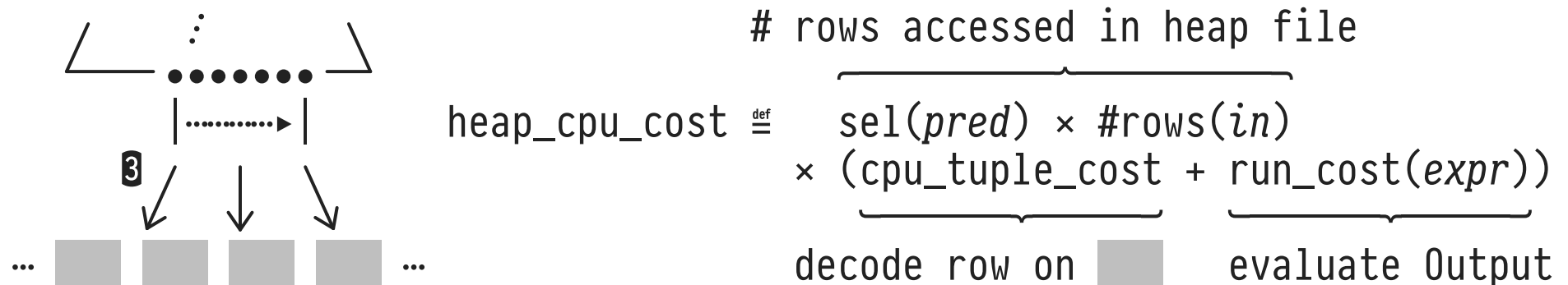\text{index\_cpu\_cost} \stackrel{\text{def}}{=} \overbrace{\underbrace{sel(pred) \times \#rows(idx)}_{\text{\# rows in scanned range } |\cdots\blacktriangleright|}} \times (\underbrace{\text{cpu\_index\_tuple\_cost}}_{\text{decode index leaf entry}} + \underbrace{run\_cost(pred)}_{\text{evaluate } \leqslant hi})
$$

$$
\text{index\_IO\_cost} \stackrel{\text{def}}{=} \lceil \underbrace{sel(pred) \times \#pages(idx)}_{\text{\# of pages} \bullet \text{ in scanned range}} \rceil \times \underbrace{\text{random\_page\_cost}}_{\text{B}^{+}\text{Tree leaves not clustered}}
$$

*lo* ❷ *hi*

# Cost of Index Scan 5 (Heap File Access)

Heap file accesses ❸ incur additional CPU and I/O costs (no I/O cost if we perform an Index **Only** Scan):

$$\text{heap\_cpu\_cost} \stackrel{\text{def}}{=} \overbrace{\text{sel}(\textit{pred}) \times \#\text{rows}(\textit{in})}^{\text{\# rows accessed in heap file}} \times (\underbrace{\text{cpu\_tuple\_cost}}_{\text{decode row on }\blacksquare} + \underbrace{\text{run\_cost}(\textit{expr})}_{\text{evaluate Output}})$$
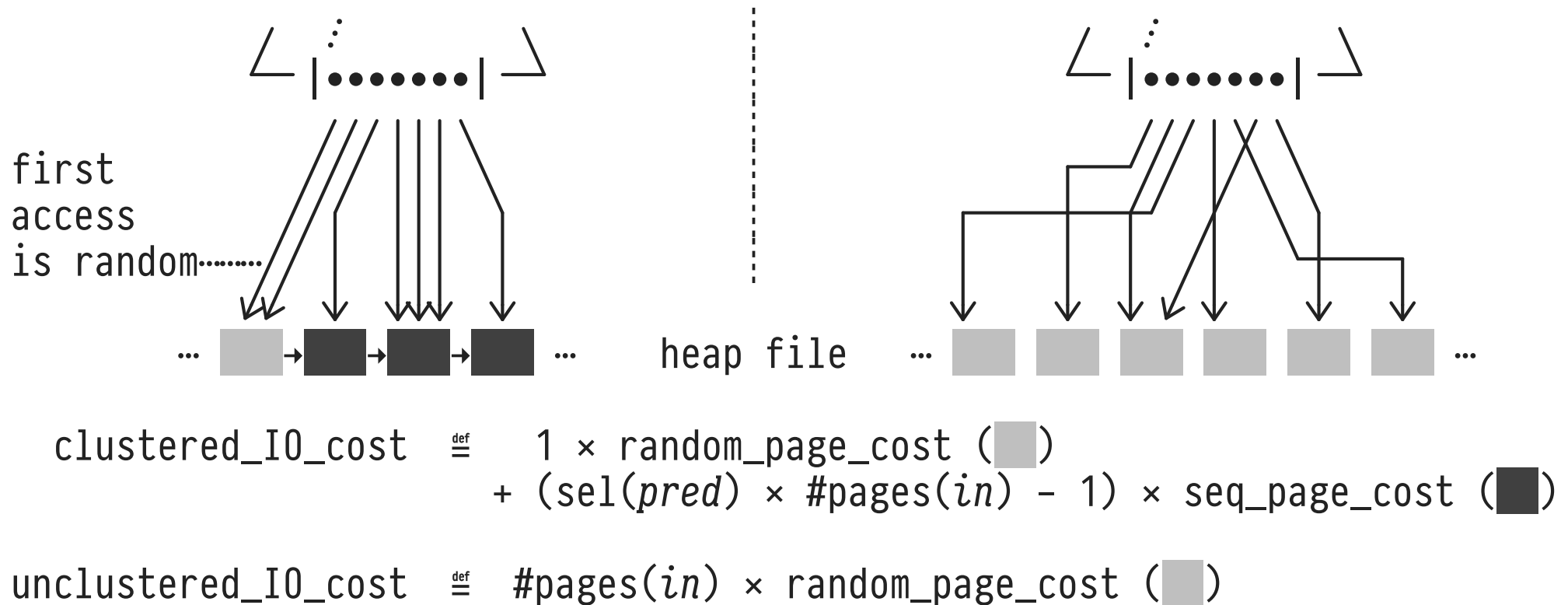
- The more **clustered** the index, the cheaper the heap I/O. Linearly interpolate between the clustered and non-clustered scenarios:

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

# Cost of Index Scan 6 ([Non-]Clustered Heap File Access)



first
access
is random ·········

heap file

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

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

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

Given ordered index *idx* over column A with values $a_1 \leqslant a_2 \leqslant \cdots \leqslant a_n$, where $\text{pos}(a_i) \in \{1,\ldots,n\}$ gives the position of $a_i$ in the heap file for A.[6]

- **Index Correlation** $\text{corr}(idx) \in \{-1,\ldots,1\}$ measures how far $[\text{pos}(a_1),\ldots,\text{pos}(a_n)]$ deviates from $[1,\ldots,n]$, *i.e.*, *idx*'s clustering degree:

$$\text{corr}(idx) = \frac{n \times \left(\sum_{i=1\ldots n} i \times \text{pos}(a_i)\right) - \left(\sum_{i=1\ldots n} i\right)^2}{n \times \left(\sum_{i=1\ldots n} i \times i \right) - \left(\sum_{i=1\ldots n} i\right)^2}$$

[6] After `CLUSTER <table> USING` *idx*, we have $\text{pos}(a_i) = i$ and thus $\text{corr}(idx) = 1$.