

Data Management for Big Data 2019-2020 (spring semester)

Dario Della Monica

# **Chapter 15: Query Processing**

These slides are a modified version of the slides provided with the book: (however, chapter numeration refers to 7<sup>th</sup> Ed.)

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan See <u>www.db-book.com</u> for conditions on re-use

The original version of the slides is available at: https://www.db-book.com/



# **Chapter 15: Query Processing**

#### Overview

- How to measure query costs
  - Establishing a cost model
- Algorithms for evaluating relational algebra operations (cost estimates)
  - Selection
  - Sorting
  - Join
- Evaluation of Expressions (How to combine algorithms for individual operations in order to evaluate a complex expression)
  - Materialization
  - Pipelining



Silberschatz, Korth, Sudarshan, Database System Concepts, 7° edition, 2011



# **Basic Steps in Query Processing**





# **Basic Steps in Query Processing (cont.)**

- Parser and translator
  - Translate the (SQL) query into relational algebra
  - Parser checks syntax (e.g., correct relation and operator names)
- Evaluation engine
  - The query-execution engine takes a query-evaluation plan, executes that plan, and returns the answers to the query
- Optimizer (in a nutshell more details in the next slides)
  - Chooses the most efficient implementation to execute the query
    - Produces equivalent relational algebra expressions
    - Annotates them with instructions (algorithms): query execution plan (QEP)
    - Estimates the cost of each equivalent QEP, according to a given cost model
    - Choose the "best" QEP



# **Basic Steps: Optimization**

- Ist level of optimization: an SQL query has many equivalent relational algebra expressions
  - $\sigma_{salary < 75000}(\prod_{salary}(instructor))$  and  $\prod_{salary}(\sigma_{salary < 75000}(instructor))$  are equivalent
  - They both correspond to
     SELECT salary
     FROM instructor
     WHERE salary < 75000</li>
- 2nd level of optimization: a relational algebra operation can be evaluated using one of several different algorithms
  - e.g., block nested-loop join VS. merge-join; file scan VS. index scan
- Input of optimization: a query in the form of an algebra expression
- Output of optimization: the "best" annotated relational algebra expression specifying detailed evaluation strategy (query evaluation plan or query execution plan – QEP) answering the input query



# **Basic Steps: Optimization (Cont.)**

- Different query evaluation plans have different costs
  - User is not expected to specify least-cost plans
- Query Optimization: amongst all equivalent QEP choose the one with lowest cost
  - Cost is estimated using statistical information from the database catalog
     # of tuples in relations, tuple sizes, # of distinct values for a given attribute, etc.
- We study... (Chapter  $15^*$  evaluation of QEP)
  - How to measure query costs (establish a cost model)
  - Algorithms for evaluating relational algebra operations and their cost
  - How to combine algorithms for individual operations in order to evaluate a complex expression (QEP)
- and (Chapter 16<sup>\*</sup> choosing the best QEP)
  - How to optimize queries, that is, how to find a QEP with lowest estimated cost

<sup>&</sup>lt;sup>\*</sup> Silberschatz, Korth, and Sudarshan, *Database System Concepts*, 7° ed.



# How to measure query costs (cost model)

These slides are a modified version of the slides provided with the book: (however, chapter numeration refers to 7<sup>th</sup> Ed.)

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan See <u>www.db-book.com</u> for conditions on re-use

The original version of the slides is available at: https://www.db-book.com/



# **Measures of Query Cost**

Response time (wall-clock time needed to execute a plan) depends on several factors

- system configuration
  - amount of dedicated buffer in RAM (aka, memory, main memory)
  - whether or not indices are (partially) stored permanently in the buffer
- runtime conditions
  - amount of free buffer at the time the plan is executed
  - content of the buffer at the time the plan is executed
  - parameters, embedded in queries, which are resolved at runtime only

SELECT salary FROM instructor WHERE salary < \$a

where \$a is a variable provided by the application (user)

#### Thus

- 1. cost models (like ours) focus on resource consumption rather than response time (optimizers minimize resource consumption rather than response time)
- 2. different optimizers may make different assumptions (parameters): every theoretical analysis must be recast with the actual parameters used by the concrete system (optimizer) to which the analysis is going to be applied



# Measures of Query Cost (Cont.)

Query cost (total elapsed time for answering a query) is measured in terms of different resources

- disk access (I/O operation on disk)
- CPU usage
- (*network communication* for distributed DBMS later in this course)
- Typically disk access is the predominant cost, and is also relatively easy to estimate. Measured by taking into account
  - Number of seeks (number of random I/O accesses)
  - Number of blocks read
  - Number of blocks written
    - It is generally assumed cost for writing to be twice as the cost for reading (data is read back after being written to ensure the write was successful)

#### **VERY IMPORTANT!!!**

- "disk" refers to permanent drive for file storage, hard-disk, secondary memory, permanent memory

- "memory" refers to volatile drive for data storage, RAM, main memory, buffer
- These are all used as synonims

This is a **so far** accepted choice for measuring query costs (**cost model**). New technologies: faster hard-disks (solid-state drives – SSD) and cheaper (thus bigger) RAM might direct towards different cost models (e.g., based also on CPU usage or RAM I/O operations)



# Measures of Query Cost (Cont.)

- We ignore difference between writing and reading: we just consider
  - *t*<sub>s</sub> time for one **seek**
  - $t_T$  time to **transfer** one block
  - Example: cost for **b** block transfers plus **S** seeks

#### $b * t_T + S * t_S$

- Values of  $t_{\tau}$  and  $t_{s}$  must be calibrated for the specific disk system
- Typical values (2018):  $t_s = 4 \text{ ms}, t_T = 0.1 \text{ ms}$
- Some DBMS performs, during installation, seeks and block transfers to estimate average values
- We ignore CPU costs for simplicity
  - Real systems usually do take CPU cost into account
- We do not include cost to writing output to disk in our cost formulae



### Algorithms for evaluating relational algebra operations

These slides are a modified version of the slides provided with the book: (however, chapter numeration refers to 7<sup>th</sup> Ed.)

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan See <u>www.db-book.com</u> for conditions on re-use

The original version of the slides is available at: https://www.db-book.com/



# **Physical organization of records**

- At the physical level, records are stored (on permanent disks) in files (managed and organized by the filesystem)
- We assume files are organized according to sequential file organization
  - i.e., a file is stored in contiguous blocks, with records ordered according to some attribute(s) – not necessarily ordered by primary key
- Other file organization techniques exist (e.g., B<sup>+</sup>-tree file organization), leading to different formulas for cost estimate



### **Selection Operation**

- **File scan** (relation scan without indices)
  - PROs: can be applied to any file, regardless of its ordering, availability of indices, nature of selection operation, etc.

CONs: it is slow

- Algorithm A1 (linear search). Retrieve and scan each file block and test all records to see whether they satisfy the selection condition
  - b<sub>r</sub> denotes number of blocks containing records from relation r
  - *Cost estimate???* (selection on a generic, non-key attribute)
    - cost =  $b_r$  block transfers + 1 seek =  $t_s + b_r * t_T$

We assume blocks are stored contiguously so 1 seek operation is enough (disk head does not need to move to seek next block)

- Selection on a key attribute. *Cost estimate???* 
  - stop on finding record
  - cost =  $(b_r/2)$  block transfers + 1 seek =  $t_s + (b_r/2) t_T$



# **Selections Using Indices**

- Index scan (relation scan using an index)
  - selection condition must be on search-key of index
  - *h<sub>i</sub>*: height of the B<sup>+</sup>-tree (# of accesses to traverse the index before accessing the data)
- A2 (primary index, equality on key). Retrieve a single record that satisfies the corresponding equality condition. Cost?

• cost =  $(h_i + 1) * (t_T + t_S)$ 

- A3 (primary index, equality on nonkey). Retrieve multiple records. Cost?
  - Let *b* = number of blocks containing matching records
  - Records will be on consecutive blocks

• cost =  $h_i * (t_T + t_S) + t_S + t_T * b_{\sim}$ 

There is a mistake in the 6<sup>th</sup> ed. of the book\* (Fig. 12.3): the " $t_s$ " summand is omitted

Silberschatz, Korth, and Sudarshan, *Database System Concepts*, 6° ed.



# **Selections Using Indices**

- A4 (secondary index, equality on key). Cost?
  - Equal to A2
    - $to solve time state (h_i + 1) * (t_T + t_S)$

#### A4 (secondary index, equality on nonkey)

- Retrieve multiple records. *Cost?* 
  - each of *n* matching records may be on a different block
  - Cost =  $(h_i + n) * (t_T + t_S)$ 
    - Can be very expensive! Can be worse than file scan



# **Selections Involving Comparisons**

Can implement selections of the form  $\sigma_{A \le V}(r)$  or  $\sigma_{A \ge V}(r)$  by using

- a linear file scan,
- or by using indices in the following ways:

#### **A5** (primary index, comparison).

- $\sigma_{A \ge V}(r)$ 
  - $\blacktriangleright$  use index to find first tuple  $\geq v~$  and scan relation sequentially from there
  - ▶ RECALL: *b* is the number of blocks containing matching records
  - Equal to A3:

```
Cost = h_i * (t_T + t_S) + t_S + t_T * b
```

- $\sigma_{A \leq V}(r)$ 
  - just scan relation sequentially till first tuple > v; do not use the index
  - Similar to A1 (file scan, equality on key):  $Cost = t_s + b^* t_T$



# **Selections Involving Comparisons**

#### • A6 (secondary index, comparison). Cost?

- For  $\sigma_{A \ge V}(r)$  use index to find first index entry  $\ge v$  and scan index sequentially from there, to find pointers to records.
- For  $\sigma_{A \le V}(r)$  just scan leaf pages of index finding pointers to records, till first entry > *v*
- In either case, retrieve records that are pointed to
  - requires an I/O for each record
  - Equal to A4, equality on nonkey: c

$$cost = (h_i + n) * (t_T + t_S)$$

Linear file scan may be cheaper



### **Summary of costs for selections**

	Algorithm	Cost	Reason			
A1	Linear Search	$t_S + b_r * t_T$	One initial seek plus $b_r$ block transfers, where $b_r$ denotes the number of blocks in the file.			
A1	Linear Search, Equality on Key	Average case $t_S + (b_r/2) * t_T$	Since at most one record satisfies the con- dition, scan can be terminated as soon as the required record is found. In the worst case, $b_r$ block transfers are still required.			
A2	Clustering B <sup>+</sup> -tree Index, Equality on Key	$(h_i + 1) * (t_T + t_S)$	(Where $h_i$ denotes the height of the in- dex.) Index lookup traverses the height of the tree plus one I/O to fetch the record; each of these I/O operations requires a seek and a block transfer.			
A3	Clustering B <sup>+</sup> -tree Index, Equality on Non-key	$h_i * (t_T + t_S) + t_S + b * t_T$	One seek for each level of the tree, one seek for the first block. Here <i>b</i> is the num- ber of blocks containing records with the specified search key, all of which are read. These blocks are leaf blocks assumed to be stored sequentially (since it is a clustering index) and don't require additional seeks.			
A4	Secondary B <sup>+</sup> -tree Index, Equality on Key	$\begin{array}{c} (h_i + 1) * \\ (t_T + t_S) \end{array}$	This case is similar to clustering index.			
A4	Secondary B <sup>+</sup> -tree Index, Equality on Non-key	$(h_i + n) * (t_T + t_S)$	(Where $n$ is the number of records fetched.) Here, cost of index traversal is the same as for A3, but each record may be on a different block, requiring a seek per record. Cost is potentially very high if $n$ is large.			
A5	Clustering B <sup>+</sup> -tree Index, Comparison	$ \begin{aligned} h_i * (t_T + t_S) + \\ t_S + b * t_T \end{aligned} $	Identical to the case of A3, equality on non-key.			
A6	Secondary B <sup>+</sup> -tree Index, Comparison	$(h_i + n) * (t_T + t_S)$	Identical to the case of A4, equality on non-key.			

Figure 15.3	Cost	estimates	for	selection	algorithms.
-------------	------	-----------	-----	-----------	-------------



### **Complex Selections**

conjunctions, disjunctions, and negation of simple conditions

#### A7 (conjunctive selection using 1 index)

- $\theta_1 AND \theta_2 AND \dots AND \theta_n$
- If there is at least 1 index useful for 1 simple condition  $\theta_i$ , then
  - use the right algorithm among A2-A6 to retrieve tuple satisfying  $\theta_i$
  - > and in the meantime check for the other simple conditions on records selected .
  - Cost?
- Cost is given by the cost of chosen algorithm for the chosen condition
  - cost depend on the choice of the condition (and the choice of the algorithm)
  - example: σ<sub>id=x AND dept=y</sub> (teacher)
  - primary index over id and secondary index over dept
  - *id* is primary key
  - it is convenient to choose id=x (with algorithm A2)

#### A8 (conjunctive selection using composite index)

- use a composite index over attributes involved in all or some of the simple conditions, if any
  - ▶ example: o<sub>name=x AND dept=y</sub> (teacher)
  - use composite index over pair (name, dept) with algorithm A4 (secondary index, equality on non-key)



# **Complex Selections (cont'd)**

A9 (conjunctive selection by intersection of identifiers)

- If there are indices with *pointers to records* (rather than *actual records*) this is our assumption so far anyway
- Scan indices but do not access records, just collect sets of pointers (one per index)
- Compute the intersection, and then access records. Cost?
- Cost: cost of scanning all indices plus cost of accessing records
- Optimization: order records in the intersection and then access them in sorted order. Advantages:
  - no block is accessed twice (2 records in the same block are retrieved together)
  - some seek time is saved as blocks are transferred in sorted order (disk-arm is minimized)
- A10 (disjunctive selection by union of identifiers)
  - If ALL conditions can be checked through some index, then similar to A9
    - Scan indices but do not access records, just collect sets of pointers (one per index)
    - Compute the union, and then access records (in sorted order)
    - Cost: cost of scanning all indices plus cost of accessing records
  - If even only 1 condition has no associate index, then A1 (linear scan)



# 2 more things on selections

- 1. Negation of a simple condition
  - NOT (Attr < v) is equivalent to Attr >= v
     NOT (Attr > v) is equivalent to Attr <= v</li>
  - NOT (Attr <= v) is equivalent to Attr > v

```
NOT (Attr >= v) is equivalent to Attr < v
```

- NOT (Attr = v) is equivalent (Attr < v) OR (Attr > v)
- 2. A4, A6, A9, A10 are very inefficient (possibly worse than linear scan) due to some blocks possibly accessed more than once
  - few records to be retrieved: better to use index scan, a lot of records to be retrieved: better to use liner scan
  - Solution: collect and sort pointers before accessing records (see optimization for A9)
  - Improved solution, based on bitmap structure: bitmap is a vector of bits (as many as number of blocks used by the relation)
    - visit index without accessing records: in the bitmap set to 1 the bits corresponding to blocks to be accessed
    - Inear scan guided by bitmap (sorted order access thanks to bitmap without actually performing a sorting)
    - hybrid solution (mix between linear scan and index access)
    - few records to be retrieved: slightly worse than index access, a lot of records to be retrieved: slightly worse than liner scan
    - thus, the cost is slightly worse than the optimal plan (linear or index scan)





- Reasons for sorting
  - Explicitly requested by SQL query
    - SELECT ...
       FROM ...
       SORT BY ...
  - Needed to efficient executions of join operations
- We may build an index on the relation, and then use the index to read the relation in sorted order. May lead to one disk block access for each tuple
- For relations that fit in memory, standard sorting techniques like quick-sort can be used. For relations that don't fit in memory, external sort-merge algorithm is a good choice



# **External Sort-Merge**

Let *M* denote number of blocks that can fit in memory.

1. Create sorted **runs** (files containing sorted pieces of relation)

Let *i* be 0 initially.

Repeatedly do the following till the end of the relation:

- (a) Read *M* blocks of relation into memory
- (b) Sort the in-memory blocks
- (c) Write sorted data to run  $R_i$
- (d) Increment i

Let the final value of *i* be *N* (number of runs)

1. Merge the runs (next slide).....



# **External Sort-Merge (Cont.)**

- 2. Merge the runs (N-way merge). We assume (for now) that N < M.
  - 1. Use *N* blocks of memory to buffer input runs, and 1 block to buffer output. Read the first block of each run into its buffer page
  - 2. repeat
    - 1. Select the first record (in sort order) among the N blocks for the runs
    - 2. Write the record to the output buffer. If the output buffer is full write it to disk.
    - Delete the record from its input buffer block.
       If the buffer block becomes empty then transfer the next block (if any) of the run into the buffer.
  - **3. until** all blocks for the runs are empty:



# External Sort-Merge (Cont.)

- If  $N \ge M$ , several merge *passes* are required.
  - In each pass, contiguous groups of *M* 1 runs are merged.
  - A pass reduces the number of runs by a factor of M -1 (and creates runs longer by the same factor)
    - E.g. If M=11, and there are 90 runs, one pass merge together groups of 10 runs into 9 new runs
       Thus, one pass reduces the number of runs to 9, each 10 times the size of the initial runs
  - Repeated passes are performed till all runs have been merged into one.



### **Example: External Sorting Using Sort-Merge**

M = 3





# **External Sort-Merge: Cost Analysis**

- Cost of block transfers:
  - Number of block transfers (read and write) for initial run creation: 2b<sub>r</sub>
  - Total number of merge passes required:  $\log_{M-1}(b_r/M)$
  - Number of block transfers (read and write) in each pass: 2b<sub>r</sub>
  - For final pass, we don't count write cost:  $b_r$  (i.e, subtract  $b_r$ )
    - we ignore final write cost for all operations since the output of an operation may be sent to the parent operation without being written to disk or just shown to video
  - Thus, total number of block transfers for external sorting:

$$2b_{r} + 2 b_{r} \lceil \log_{M-1} (b_{r}/M) \rceil - b_{r} = b_{r} (2 \lceil \log_{M-1} (b_{r}/M) \rceil + 1)$$

• Seeks: next slide



# External Sort-Merge: Cost Analysis (cont.)

- Cost of seeks
  - During run generation: one seek to read each run and one seek to write each run

▶  $2 \lceil b_r / M \rceil$ 

- Total number of merge passes required:  $\log_{M-1}(b_r/M)$
- During each pass (merge phases)
  - I seek for reading each block and 1 seek for writing each block
    - $-2 b_r$  seeks for each merge pass
  - except the final one which does not require a write

 $- - b_r$  (i.e, subtract  $b_r$ )

• Total number of seeks:

 $2\lceil b_r/M\rceil + 2b_r\lceil \log_{M-1}(b_r/M)\rceil - b_r =$ 

=  $2 \lceil b_r / M \rceil$  +  $b_r (2 \lceil \log_{M-1}(b_r / M) \rceil - 1)$ 



#### An improved version of the algorithm \*

- Number of seeks can be reduced by using b<sub>b</sub> many blocks (instead of 1) for each run during the run merge phase
  - Using 1 block per run leads to too many seeks
  - Instead, using  $b_b$  buffer blocks per run  $\rightarrow$  read/write  $b_b$  blocks with only 1 seek
    - Number of runs merged together:  $\lfloor M/b_b \rfloor 1$  runs (instead of M 1)
      - Scaling factor is  $\lfloor M / b_b \rfloor 1$  instead of M 1
    - Number of passes required:  $\lceil \log_{\lfloor M/b_h \rfloor 1}(b_r/M) \rceil$  instead of  $\lceil \log_{M-1}(b_r/M) \rceil$
    - During the merge phase:  $2 \left[ \frac{b_r}{b_p} \right]$  seeks for each pass (instead of 2  $b_r$ )
      - Except the final one (we assume final result is not written to disk)
- Thus total number of block transfers for external sorting:  $b_r(2 \lceil \log_{|M/b_r|-1} (b_r/M) \rceil + 1)$

Total number of seeks:

### $2\lceil b_r/M\rceil + \lceil b_r/b_b\rceil (2\lceil \log_{\lfloor M/b_b\rfloor - 1}(b_r/M)\rceil - 1)$

<sup>\*</sup> In Silberschatz, Korth, and Sudarshan, Database System Concepts, 6° ed., the non-improved version of the algorithm is given only, but the cost analysis mixes elements from both versions of the algorithm



# **Join Operation**

- Several different algorithms to implement joins
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join
- Choice based on cost estimate
- Running example : students × takes where
  - Number of *records* of *student*: 5,000
  - Number of *blocks* of *student*: 100
  - Number of *records* of *takes*: 10,000
  - Number of *blocks* of *takes*:

400



# **Nested-Loop Join**

To compute the theta join  $r \Join_{\theta} s$ 

```
for each tuple t_r in r do
for each tuple t_s in s do
test pair (t_r, t_s) to see if they satisfy the join condition \theta
if it does, add t_r \cdot t_s to the result
end
end
```

- *r* is called the **outer relation** and *s* the **inner relation** of the join
- Requires no indices and can be used with any kind of join condition
- Expensive since it examines every pair of tuples in the two relations



# **Nested-Loop Join (Cont.)**

- If the smaller relation fits entirely in memory, use that as the inner relation
  - $b_r + b_s$  block transfers and 2 seeks

(same cost in the best case scenario, when both relations fit in memory)

- Worst case (there is enough memory forn only one block for each relation)
  - # of block transfer:

 $n_r * b_s + b_r$  (*b<sub>r</sub>* transfers to read relation *r* +  $n_r * b_s$  transfers to read *s* for each tuple in *r*)

# of seeks:

- $n_r + b_r$  (*b<sub>r</sub>* seeks to read relation  $r + n_r$  seeks to read *s* for each tuple in *r*)
- Running example (join between *students* and *takes* assuming worst case memory availability)
  - with *student* as outer relation:
    - ▶ 5,000 \* 400 + 100 = 2,000,100 block transfers
    - ▶ 5,000 + 100 = 5,100 seeks
  - with *takes* as the outer relation
    - ▶ 10,000 \* 100 + 400 = 1,000,400 block transfers
    - ▶ 10,000 + 400 = 10,400 seeks
  - if smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers and 2 seeks
- Block nested-loops algorithm (next slide) is preferable



# **Block Nested-Loop Join**

Variant of nested-loop join in which every block of inner relation is paired with every block of outer relation.

```
for each block B_r of r do
for each block B_s of s do
for each tuple t_r in B_r do
for each tuple t_s in B_s do
check if (t_r, t_s) satisfy the join condition
if it does, add t_r \cdot t_s to the result
end
end
end
```



# **Block Nested-Loop Join (Cont.)**

- Worst case estimate (memory holds one block for each relation):
  - Each block in the inner relation is read once for each block in the outer relation
  - # of block transfers:

$$b_r * b_s + b_r$$

• # of seeks:

▶ 2 \* b<sub>r</sub>

- (block) nested-loop improvements
  - join attributes form a key for the inner relation:
    - inner loop terminates when first match is found
  - If there is more space in memory
    - read as many blocks as possible for the outer relation (block nested-loop):
    - # of block transfer:  $[b_r / (M-2)] * b_s + b_r$
    - # of seeks:  $2 * \left[ \frac{b_r}{M-2} \right]$
  - alternate forward and backward scan for inner relation:
    - Use blocks already in buffer: save some block transfer

**Database System Concepts - 7th Edition** 



# Indexed Nested-Loop Join

- If an index is available for one of the relations on the attribute of the join condition
  - then use such a relation as inner relation in a nested-loop join
- Instead of doing a linear scan as inner loop, do an index scan

for each tuple  $t_r$  in r do index scan over s to find tuples  $t_s$  satisfying the join condition with tuple  $t_r$ end

(basically, for every tuple in *r*, do a selection on *s* using the index)

- It might be convenient to create an ad-hoc index for the join if it does not exist
- Cost (worst case: space in memory for only 1 block for each relation)
  - $b_r$  seeks and block transfers to read r.  $b_r * (t_T + t_S)$
  - for each record in r, index scan on s: nr \* c (where c is the cost of index scan on s)
  - thus, total cost =  $b_r * (t_T + t_S) + n_r * c$
  - NOTICE: if there are index for both relations, then better to use relation with less records as outer relation



# Indexed vs. Block Nested-Loop Join

- Block nested-loop join
  - $b_r * b_s + b_r$  block transfer and  $2 * b_r$  seeks, that is,  $(b_r * b_s + b_r) * t_T + 2 * b_r * t_s$
- Indexed nested-loop join
  - $b_r * (t_T + t_S) + n_r * c$
- Running example : *students ≥ takes* 
  - $n_{students} = 5,000$  $b_{students} = 100$  $n_{takes} = 10,000$  $b_{takes} = 400$

students is used as outer relation as it has less records

- $B^+$ -index on *takes.id* with *fan* = 20, and thus height h = 4
- cost of each index scan  $c = (h + n) * (t_T + t_S)$  [A4: secondary index, eq. on non-key]
- $n = n_{students} / n_{takes} = 2$  [average]
- cost indexed nested-loop join:  $100 * (t_T + t_S) + 5,000 * ((4+2) * (t_T + t_S)) = 30,100* (t_T + t_S)$
- cost block nested-loop join: (100 \* 400 + 100) \*  $t_T$  + 2 \* 100 \*  $t_S$  = 40,100 \*  $t_T$  + 200 \*  $t_S$



### How to combine algorithms for individual operations in order to evaluate a complex expression

These slides are a modified version of the slides provided with the book: (however, chapter numeration refers to 7<sup>th</sup> Ed.)

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan See <u>www.db-book.com</u> for conditions on re-use

The original version of the slides is available at: https://www.db-book.com/



# **Evaluation of Expressions**

- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree
  - Materialization: store (materialize) on disk results of evaluation of sub-expressions into temporary relations for subsequent use
    - Disadvantage: several disk writing and reading to store temporary relations
    - Always possible
  - Pipelining: pass on tuples to parent operations as they are generated by inner operations being executed
    - Advantage: less disk writing
    - Not always possible



### **Materialization**

Materialized evaluation: evaluate one operation at a time, starting at the lowest-level. Use intermediate results materialized into temporary relations to evaluate next-level operations

E.g., 
$$\prod_{name} (\sigma_{building="Watson"}(department) \bowtie \text{ instructor})$$





# **Pipelining**

- Pipelined evaluation: evaluate several operations simultaneously, passing (partial) results of one operation on to the next as they are generated (es., single records), without writing them on disk
- E.g., in previous expression tree, don't store result of

 $\sigma_{building="Watson"}(department)$ 

Instead, pass tuples directly to the join as they are found

Similarly, don't store result of join, pass tuples directly to projection as they are generated

- Cheaper than materialization: no need to store a temporary relation to disk
- (partial) Results are output earlier, before waiting for complete query execution
- Parallelization of operations
- Pipelining may not always be possible
  - some sorting algorithms cannot output tuples early, only after all input tuples have been examined
  - indexed nested-loop join cannot have its input inner relation pipelined as the whole relation with associated index must be available



### **End of Chapter**

These slides are a modified version of the slides provided with the book: (however, chapter numeration refers to 7<sup>th</sup> Ed.)

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan See <u>www.db-book.com</u> for conditions on re-use

The original version of the slides is available at: https://www.db-book.com/