

Database Management Systems

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Lecture 9, 17 November 2023

Query processing

- Translate the query from SQL into relational algebra
- Evaluate the relational algebra expression

■ Challenges

- Many equivalent relational algebra expressions

$\sigma_{salary < 75000}(\pi_{salary}(instructor))$ vs $\pi_{salary}(\sigma_{salary < 75000}(instructor))$

- Many ways to evaluate a given expression

■ Query plan

- Annotate the expression with a detailed evaluation strategy key values
 - Use index on *salary* to find instructors with *salary* < 75000
 - Or, scan entire relation, discard rows with *salary* ≥ 75000

Query optimization

- Choose plan with lowest cost
- Maintain **database catalogue** — number of tuples in each relation, size of tuples, ...
- Assess cost in terms of disk access and transfer, CPU time, ...
- For simplicity, ignore in-memory costs (CPU time), restrict to disk access
- Disk accesses
 - Relation r occupies b_r blocks
 - **Disk seeks** — time t_S per seek
 - **Block transfers** — time t_T per transfer
- Other factors — buffer management etc

Selection

- $t_s + b_r b_r$ (A1) Linear search **b_r blocks** — I seek to start, b_r transfer
- (A2) Clustering index, equality on key — index height h_i
- (A3) Clustering index, equality on nonkey
- (A4) Secondary index (key, non-key)
- (A5) Clustering index, comparison — sorted on A
- (A6) Clustering index, comparison — not sorted on A
- (A7) Conjunctive selection using one index
- (A8) Conjunctive selection using composite index
- (A9) Conjunctive selection using intersection of pointers
- (A10) Disjunctive selection by union of pointers
- (Neg) Negation

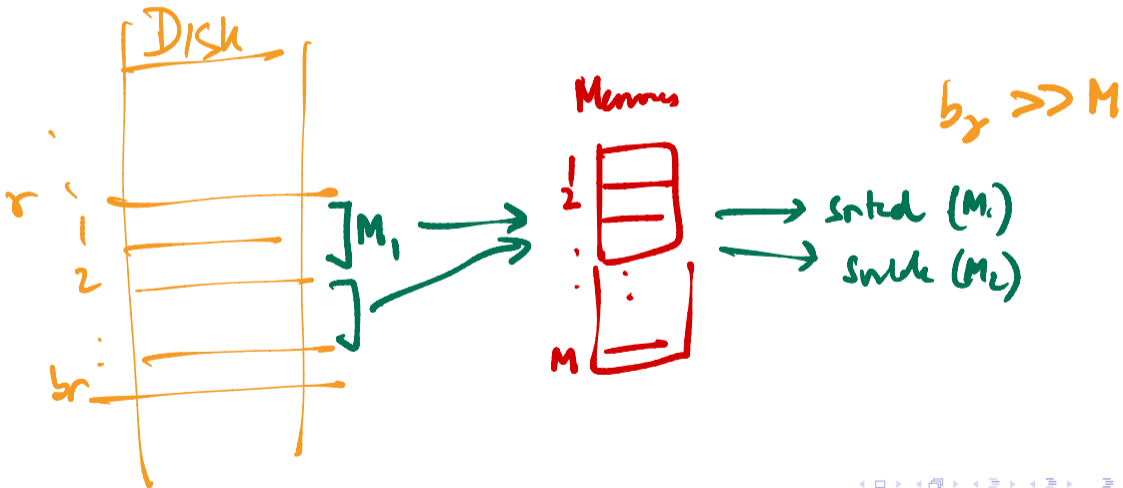
Ignore final
output
write

Sorting

- In-memory sorting vs sorting on disk
- Merging sorted lists — varieties
- Traditional merge sort

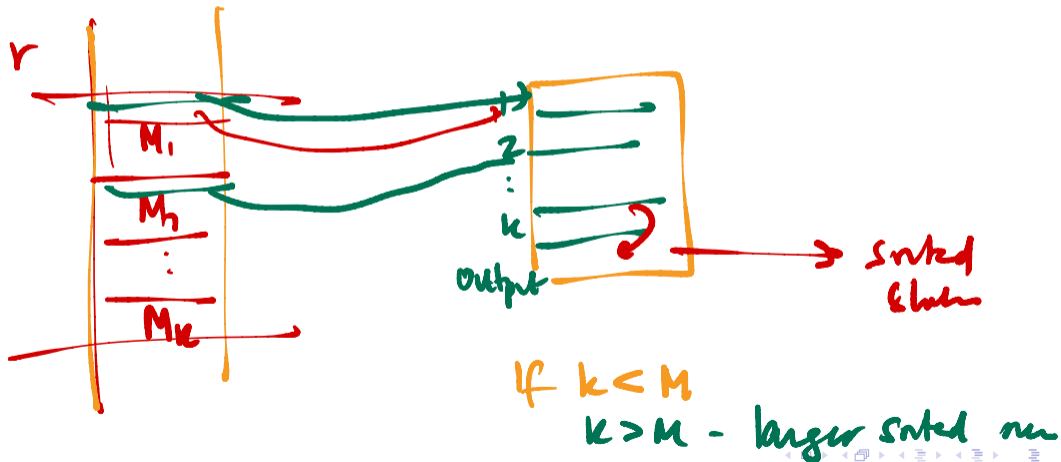
External merge sort

- N records, b_r blocks, M blocks in memory



External merge sort

- N records, b_r blocks, M blocks in memory

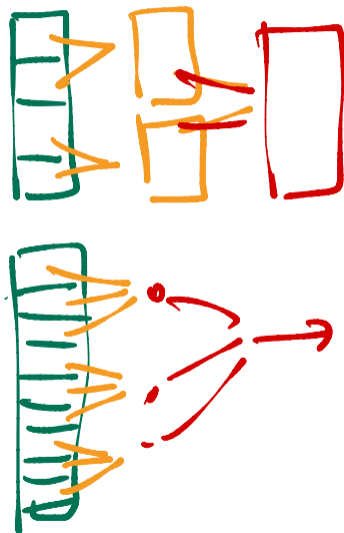


External merge sort

- N records, b_r blocks, M blocks in memory
- Compute sorted runs of size M
- Merge sorted runs, 1 block per run vs b_b blocks per run
- Complexity
 - b_r/M sorted runs, $\lceil \log_{\lfloor M/b_b \rfloor - 1}(b_r/M) \rceil$ merge passes

External merge sort

- N records, b_r blocks, M blocks in memory
- Compute sorted runs of size M
- Merge sorted runs, 1 block per run vs b_b blocks per run
- Complexity
 - b/M sorted runs, $\lceil \log_{M/b_b} \lceil b_r/M \rceil \rceil$ merge passes
 - Block transfers — $b_r (2 \lceil \log_{M/b_b} \lceil b_r/M \rceil \rceil + 1)$
 - Why not $b_r (2 \lceil \log_{M/b_b} \lceil b_r/M \rceil \rceil + 2)$?



External merge sort

- N records, b_r blocks, M blocks in memory
- Compute sorted runs of size M
- Merge sorted runs, 1 block per run vs b_b blocks per run
- Complexity
 - b_r/M sorted runs, $\lceil \log_{\lfloor M/b_b \rfloor - 1}(b_r/M) \rceil$ merge passes
 - Block transfers — $b_r (2 \lceil \log_{\lfloor M/b_b \rfloor - 1}(b_r/M) \rceil + 1)$
 - Why not $b_r (2 \lceil \log_{\lfloor M/b_b \rfloor - 1}(b_r/M) \rceil + 2)$?
 - Block seeks — $2 \lceil b_r/M \rceil + \lceil b_r/b_b \rceil (2(\lceil \log_{\lfloor M/b_b \rfloor - 1}(b_r/M) \rceil) - 1)$

Computing joins

- Running example

- *Student* ⋈ *Takes*

- *Student* — 5000 rows, 100 blocks

- *Takes* — 10000 rows, 400 blocks

50 rows/block
25 rows/block

Nested-loop join

- r (5000 rows, 100 blocks) *Student* ⋈ s *Takes* (10000 rows, 400 blocks)

n_r b_r
 n_s b_s

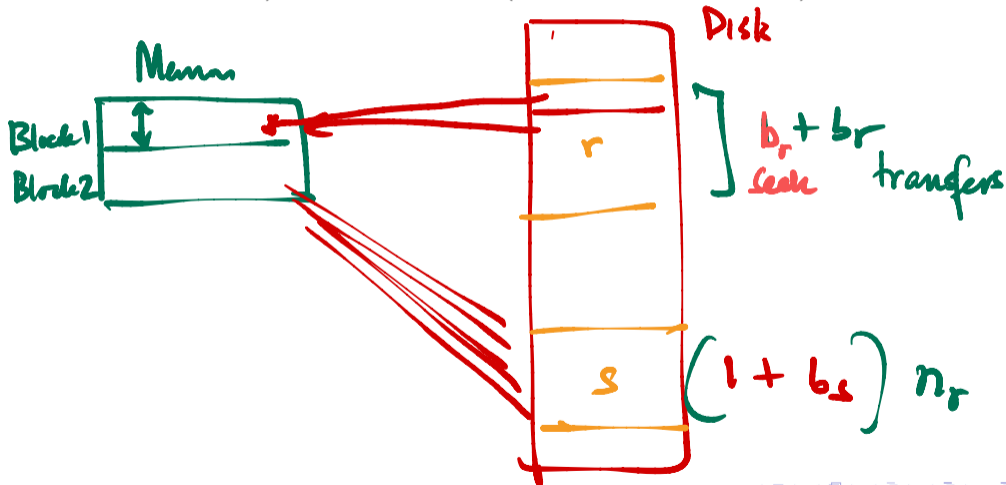
for each row r in *Student* r DAS = SDAR
 for each row s in *Takes* $\leftarrow \times n_r$

$\hookrightarrow 1 \text{ seek} + b_s$ 2×10^6
 $101 + 5000(401)$

$$\begin{aligned}
 & \times \boxed{1 \text{ seek}} + b_r \text{ transfers for } r \quad b_r \times 1 + b_r + n_r(1 + b_s) \\
 & + n_r (1 \text{ seek} + b_s \text{ transfers for } s) \quad b_s \times 1 + b_s + n_s(1 + b_r) \\
 & \qquad \qquad \qquad = 401 + 10000(101)
 \end{aligned}$$

Nested-loop join

- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)



Nested-loop join

- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)
- Complexity
 - $r \bowtie_{\theta} s$ — r is outer relation, s is inner relation
 - Block transfers: $b_r + n_r \cdot b_s$
 - Block seeks: $b_r + n_r$ — inner relation read sequentially
 - Special case: smaller relation fits in memory

As asked in class, you do re-seek τ for each fresh block

Block nested-loop join

- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)

for each block r in *Student*

$$(1 + b_r) + b_r(1 + b_s)$$

for each block s in *Takes*

$$(1 + b_s) + b_s(1 + b_r)$$

$$b_r(1 + b_s)$$

Compare all rows in r vs all rows in s

vs

$$n_r(1 + b_s)$$

Block nested-loop join

- (5000 rows, 100 blocks) *Student* \bowtie *Takes* (10000 rows, 400 blocks)
- Complexity
 - $r \bowtie_{\theta} s$ — r is outer relation, s is inner relation
 - Block transfers: $b_r + b_r \cdot b_s$
 - Block seeks: $b_r + b_r = 2b_r$

Indexed nested-loop join

- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)

for each r in *Student*

Use index to fetch matching rows in *Takes*

Indexed nested-loop join

- (5000 rows, 100 blocks) *Student* \bowtie *Takes* (10000 rows, 400 blocks)
- Complexity
 - $r \bowtie_{\theta} s$ — r is outer relation, s is inner relation
 - Total cost: $b_r(t_T + t_S) + n_r \cdot c$
 - c is cost of single selection on s

Merge join

- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)

Sort r
Sort s
Merge

Cost of sorting r & s
+ $O(b_r + b_s)$

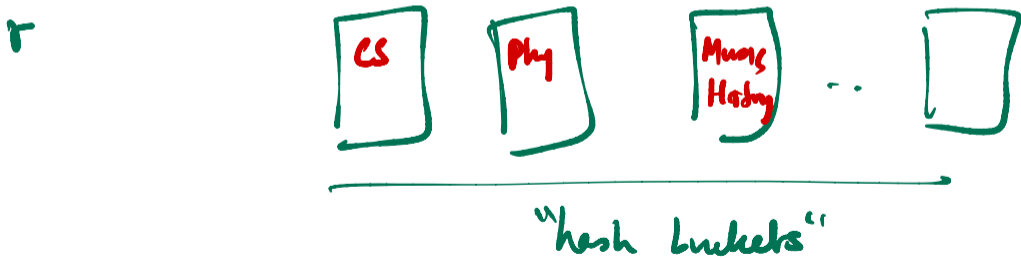
Merge join

- (5000 rows, 100 blocks) *Student* \bowtie *Takes* (10000 rows, 400 blocks)
- Complexity
 - $r \bowtie_{\theta} s$ — r is outer relation, s is inner relation
 - Block transfers: $b_r + b_s$
 - Block seeks: $\lceil b_r/b_b \rceil + \lceil b_s/b_b \rceil$

Hash join

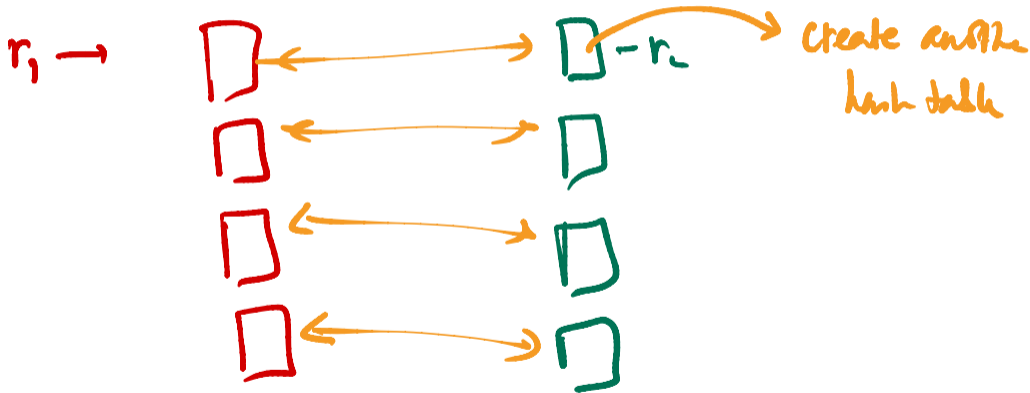
- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)

$k \mapsto h(k) \rightarrow i$ position in list



Hash join

- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)

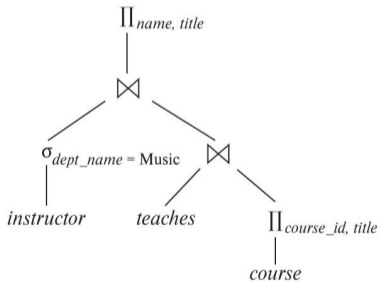
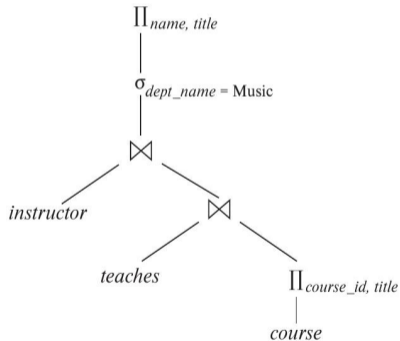


Other operations

- Duplicate removal
- Aggregate queries with grouping
 - Aggregate while sorting/hashing
- Set theoretic operations

Query optimization

- Choose plan with lowest cost
- *Find names and course titles of courses taught by instructors from Music Dept*



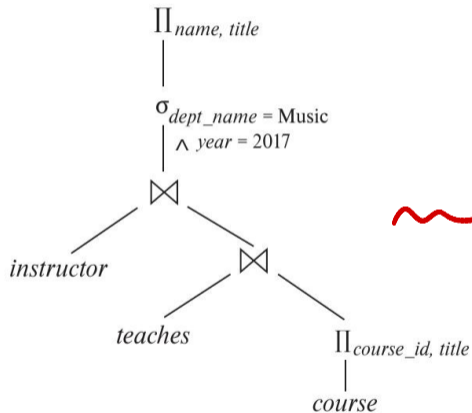
Transforming expressions

$$r \bowtie s = s \bowtie r$$

$$\sigma_{\theta}(r \bowtie s) \implies \sigma_{\theta}(r) \bowtie s \quad \text{if } \theta \text{ only refers to } r$$

Transformer

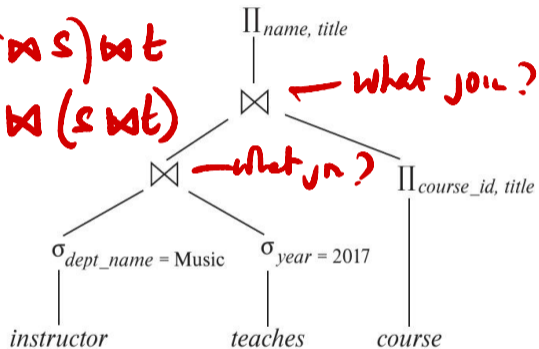
Transforming expressions



raswt

(ras)wt

rw (s wt)



Maintaining a database catalogue

- n_r — number of tuples in r
- b_r — number of blocks used by r
- l_r — size of a tuple in r
- f_r — blocking factor of r , how many tuples fit in a block
- $V(A, r)$ — number of distinct values of attribute A in r
 - Store distribution of values as histogram

Estimating output of an operation

- Selection
 - Simple, range, conjunction, disjunction
- Join
 - Keys and non-keys
- Projection
- Aggregation
- Set operations
- Outer joins

$$(M_1 \times M_2 \times \dots \times M_n)$$

$$r_1 \bowtie r_2 \bowtie \dots \bowtie r_n$$

Similar to
matrix multi-
ordering

Heuristics

- Perform selection early
- Perform projection early
- Perform most restrictive selection/join first