## RDBMS and SQL

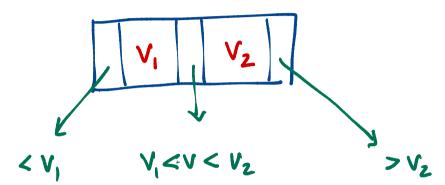
Madhavan Mukund

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

Lecture 10, 6 November 2025

Recep - Store tables on disk - Transfer data in blocks - seek, transfer - Maintain index k lards -> 2°+2'+ ..+2k-1 = 2k-1-size N L'/R' k = log N

# Generalize to 3-ary trees



#### Search trees

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height
- Block-based access

Attnhute

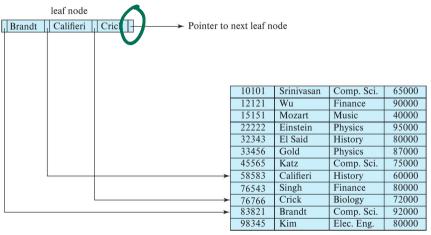
- Binary tree node has one search key value, two pointers
- Block can hold much more
- Generalize to multiple key values, multiple pointers



1 TB duste - 
$$10^{12}$$
 $\Rightarrow$  40916 Byte Stocks  $\Rightarrow$   $10^9$  blocks  $10^9 \approx 30$ 

#### B+ trees

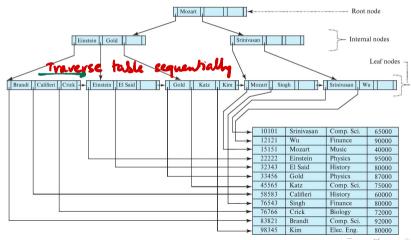
■ Leaf nodes form a dense index — linked list of leaves, each one block



RDBMS-SQL. Lecture 10, 6 Nov 2025

#### B+ trees

- Leaf nodes form a dense index linked list of leaves
- Non-Leaf nodes form a sparse index



V Name = "Brandt"

Name starts with D-L Raye of value Search DAA-A

#### B+ trees

- Leaf nodes form a dense index linked list of leaves
- Non-leaf nodes form a sparse index
- $\blacksquare$  Constraints assume n keys and pointers can fit in a block
  - Each leaf has at least  $\lceil (n-1)/2 \rceil$  key values
  - Each non-leaf has at least  $\lceil n/2 \rceil$  pointers
  - Height of the tree is proportional to  $log_{n/2}(n)$

## Query processing

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

## 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 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

4 / 20

#### Query optimization

- Choose plan with lowest cost
- Maintain database catalogue number of tuples in each relationn, 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

#### Query optimization

- Choose plan with lowest cost
- Maintain database catalogue number of tuples in each relationn, 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<sub>S</sub>* per seek
  - Block transfers time  $t_T$  per transfer

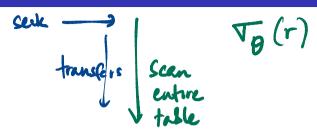
5/20

#### Query optimization

- Choose plan with lowest cost
- Maintain database catalogue number of tuples in each relationn, 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<sub>S</sub>* per seek
  - Block transfers time  $t_T$  per transfer
- Other factors buffer management etc

5 / 20

■ Linear search



Linear search

Clustering index — index height h

■ Equality on key

■ Equality on nonkey

- Linear search
- Clustering index index height h<sub>i</sub> table is redeal by this
  - Equality on key
  - Equality on nonkey
- Secondary index (key, non-key)

- Linear search
- Clustering index index height h<sub>i</sub>
  - Equality on key
  - Equality on nonkey
- Secondary index (key, non-key)
- Clustering index, comparison on A
  - Sorted on A
  - Not sorted on A

$$V_1 \le h \le V_2 = Size n$$
blocks
$$h(t_5 + t_7) + t_5 + n \cdot t_7$$

- Linear search
- Clustering index index height h<sub>i</sub>
  - Equality on key
  - Equality on nonkey
- Secondary index (key, non-key)
- Clustering index, comparison on A
  - Sorted on A
  - Not sorted on A
- Boolean combinations
  - Conjunctive selection using one index
  - Conjunctive selection using composite index
  - Disjunction, negation, . . .





, smalle of two combine

Join = Va (r, xr2)

## Sorting

■ In-memory sorting vs sorting on disk

# Sorting

- In-memory sorting vs sorting on disk
- Merging sorted lists varieties

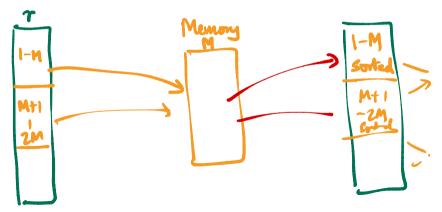
Merging - Combine 2 sorted lists into 1 Duplicates - Eliminate - Union - Keep mly - Intersech duphieter - Set difference

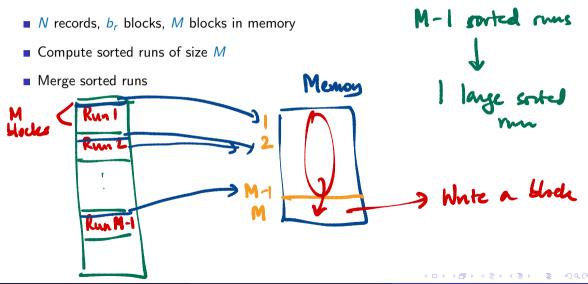
## Sorting

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

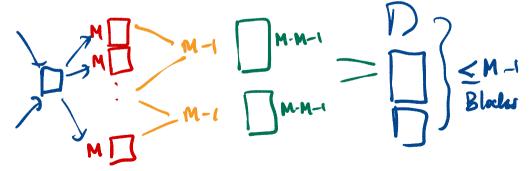
 $\blacksquare$  N records,  $b_r$  blocks, M blocks in memory

- $\blacksquare$  N records,  $b_r$  blocks, M blocks in memory
- Compute sorted runs of size M

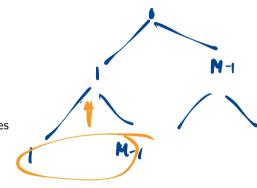




- $\blacksquare$  N records,  $b_r$  blocks, M blocks in memory
- Compute sorted runs of size M
- Merge sorted runs

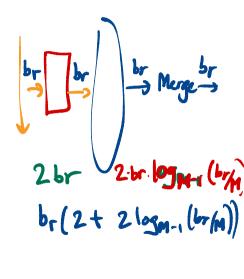


- $\blacksquare$  N records,  $b_r$  blocks, M blocks in memory
- Compute sorted runs of size M
- Merge sorted runs
- Complexity
  - $b_r/M$  sorted runs,  $\lceil \log_{M-1}(b_r/M) \rceil$  merge passes



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

last nound of writing is ignored - "pipeline"



- $\blacksquare$  N records,  $b_r$  blocks, M blocks in memory
- Compute sorted runs of size M
- Merge sorted runs
- Complexity
  - $b_r/M$  sorted runs,  $\lceil \log_{M-1}(b_r/M) \rceil$  merge passes
  - Block transfers  $b_r (2\lceil \log_{M-1}(b_r/M) \rceil + 1)$ 
    - Why not  $b_r (2\lceil \log_{M-1}(b_r/M) \rceil + 2)$ ?
  - Block seeks  $2\lceil b_r/M \rceil + b_r \left(2(\lceil \log_{M-1}(b_r/M) \rceil 1\right)$

# Computing joins

- Running example
  - Student ⋈ Takes

# Computing joins

- Running example
  - Student ⋈ Takes
  - Student 5000 rows, 100 blocks
  - *Takes* 10000 rows, 400 blocks

50 rows / block 25 rows / bock

# 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
  - $ightharpoonup r\bowtie_{\theta} s r$  is outer relation, s is inner relation

for each row in or for each vow in s

10 / 20

#### 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:  $b_r + n_r \cdot b_s$



### 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

secks of s



### 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:  $b_r + \chi \cdot b_s$
  - Block seeks:  $b_r + n_r$  inner relation read sequentially
  - Special case: smaller relation fits in memory



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

- (5000 rows, 100 blocks) Student ⋈ Takes (10000 rows, 400 blocks)
- Complexity
  - $r\bowtie_{\theta} s-r$  is outer relation, s is inner relation

- (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:  $b_r b_s$

- (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:  $b_r + b_r \cdot b_s$
  - Block seeks:  $b_r + b_r = 2b_r$



### Indexed nested-loop join

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

### Indexed nested-loop join

- (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)
- Complexity
  - $ightharpoonup r\bowtie_{\theta} s r$  is outer relation, s is inner relation

for each row in r seach for A in S

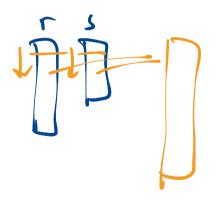
### Indexed nested-loop join

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

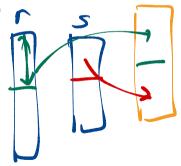
■ (5000 rows, 100 blocks) *Student* ⋈ *Takes* (10000 rows, 400 blocks)

- (5000 rows, 100 blocks) Student ⋈ Takes (10000 rows, 400 blocks)
- Complexity
  - $r\bowtie_{\theta} s-r$  is outer relation, s is inner relation
    - Assume relations are sorted (add cost of sorting)

- (5000 rows, 100 blocks) Student ⋈ Takes (10000 rows, 400 blocks)
- Complexity
  - $r \bowtie_{\theta} s r$  is outer relation, s is inner relation
    - Assume relations are sorted (add cost of sorting)
  - Block transfers:  $b_r + b_s$

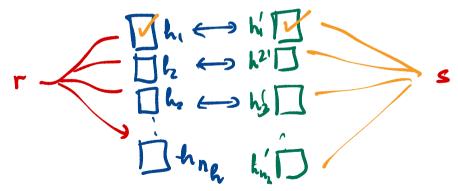


- (5000 rows, 100 blocks) Student ⋈ Takes (10000 rows, 400 blocks)
- Complexity
  - $r\bowtie_{\theta} s-r$  is outer relation, s is inner relation
    - Assume relations are sorted (add cost of sorting)
  - Block transfers:  $b_r + b_s$
  - Block seeks:  $\lceil b_r/b_b \rceil + \lceil b_s/b_b \rceil$ 
    - Read a chunk of blocks  $b_b$  at a time



### Hash join

- (5000 rows, 100 blocks) Student ⋈ Takes (10000 rows, 400 blocks)
- Hash function on join attribute A,  $n_h$  output values



### Hash join

- (5000 rows, 100 blocks) Student ⋈ Takes (10000 rows, 400 blocks)
- Hash function on join attribute A,  $n_h$  output values
- Join each pair of hash buckets build index and probe

