

## Tutorial Week 8

### ① Single-relation plans

$$NTuples(employees) = 10,000 \text{ pages} \times 20 \text{ tuples/page} \\ = 200,000 \text{ tuples}$$

a)  $sal > 20,000$

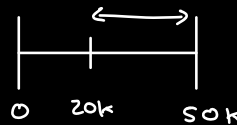
sal ranges from 0 to 50,000

- Reduction factor (RF) :

$$RF = \frac{High(sal) - value}{High(sal) - low(sal)}$$

$$= \frac{50,000 - 20,000}{50,000 - 0}$$

$$= 0.6$$



- 2 possible access paths :

- Full table Scan
- unclustered B<sup>+</sup> tree index on sal

- Full table scan / sequential scan / heap scan

$$cost = NPages(employees)$$

$$= 10,000 \text{ I/O}$$

careful : find the product of the reduction factors of matching predicates

- Unclustered B<sup>+</sup> tree index on sal

$$Cost = (NPages(Index) + NTuples(employees)) \times TTRF,$$

$$= (500 + 200,000) \times 0.6$$

$$= 120,300 \text{ I/O}$$

conditions that involve attributes in a prefix of the searchkey

cheapest access path is full-table scan

b) age = 25

$$RF = \frac{1}{N_{\text{keys}}(\text{age})}$$

$$= \frac{1}{40}$$

40 distinct values of age

• 3 possible access paths:

- full table scan
- unclustered hash index on age
- clustered B<sup>+</sup> tree index on (age, sal) → age is within prefix of search key

• Full table scan

$$\text{cost} = (10,000 \text{ I/O})$$

Same as before

• unclustered hash index on age

RF of matching conditions

$$\text{Cost} = N_{\text{Tuples}}(\text{employees}) \times \prod RF_i \times 2.2$$

$$= 200,000 \times \frac{1}{40} \times 2.2$$

$$= 11,000 \text{ I/O}$$

↳ 1.2 for the bucket check

1 to fetch page from disk

• clustered B<sup>+</sup> tree index on (age, sal)

$$\text{cost} = (N_{\text{Pages}}(\text{index}) + N_{\text{Pages}}(\text{employees})) \times \prod RF_i$$

$$= (500 + 10,000) \times \frac{1}{40}$$

$$= 262.5$$

$$\approx 263 \text{ I/O}$$

↓  
RF of matching conditions

cheapest access path is clustered B<sup>+</sup> tree index

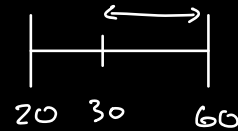
c) age > 30

age ranges from 20 to 60

$$RF = \frac{\text{high}(\text{age}) - \text{value}}{\text{high}(\text{age}) - \text{low}(\text{age})}$$

$$= \frac{60 - 30}{60 - 20}$$

$$= 0.75$$



2 possible access paths:

- full table scan
- clustered B+ tree index on (age, sal) → age is within prefix of search key

Can't use unclustered hash index on age  
since we can't use hash index over a range

Full table scan

$$\text{cost} = 10,000 \text{ I/O} \quad \text{Same as before}$$

clustered B+ tree index on (age, sal)

$$\begin{aligned} \text{cost} &= (N_{\text{Pages}}(\text{index}) + N_{\text{Pages}}(\text{employees})) \times \text{TRF}_i \\ &= (500 + 10,000) \times 0.75 \end{aligned}$$

$$= 7,875 \text{ I/O}$$

↓  
RF of  
matching  
conditions

cheapest access path is clustered B+ tree index

d)  $eid = 1000$

• 2 possible access paths :

- full table scan
- unclustered hash index on  $eid$

- $eid$  is a candidate key, so here, we're looking for just a single tuple
- unclustered hash index on  $eid$  is a primary index i.e. the search key of the index is the primary key of the employees relation
- selecting single tuple over primary key so :

unclustered hash index on  $eid$

$$\begin{aligned} \text{cost} &= 1.2 + 1 \quad \text{data page access} \\ &= 2.2 \text{ I/O} \quad \text{hash lookup cost} \end{aligned}$$

• Full table scan

$$\text{cost} = 10,000 \text{ I/O} \quad \text{Same as before}$$

cheapest access path is unclustered hash index on  $eid$

e)  $Sal > 20,000 \wedge age > 30$

$$RF_{Sal} = \frac{High(Sal) - value}{High(Sal) - low(Sal)}$$
$$= \frac{50,000 - 20,000}{50,000 - 0}$$

$$= 0.6$$

as calculated before

$$RF_{age} = \frac{high(age) - value}{high(age) - low(age)}$$
$$= \frac{60 - 30}{60 - 20}$$

$$= 0.75$$

3 possible access paths:

- full table scan
- unclustered B<sup>+</sup> tree index on Sal
- clustered B<sup>+</sup> tree index on (age, Sal)

filter age "on the fly" afterwards



only  $RF_{Sal}$  will be used  
since age condition  
doesn't match index

→ we will use product of  $RF$ s  
of both conditions, since  
attributes are in prefix  
of search key

can't use unclustered hash index on age  
since we can't use hash index over a range

- Full table scan

$$\text{cost} = 10,000 \text{ I/O}$$

Same as before

- Unclustered B+ tree index on sal

$$\text{cost} = (\text{NPages}(\text{index}) + \text{NTuples}(\text{employees})) \times \text{TRF}_i$$

$$= (500 + 200,000) \times 0.6$$

sal > 20,000 is the only matching condition

$$= 120,300 \text{ I/O}$$

Same as question 1a

matching conditions

- clustered B+ tree index on (age, sal)

$$\text{cost} = (\text{NPages}(\text{index}) + \text{NPages}(\text{employees})) \times \text{TRF}_i$$

$$= (500 + 10,000) \times 0.75 \times 0.6$$

matching conditions

$$= 4,725 \text{ I/O}$$

sal > 20,000  $\wedge$

age > 30  $\wedge$

a matching predicate

cheapest access path is clustered B+ tree index

② Multi-relation plans

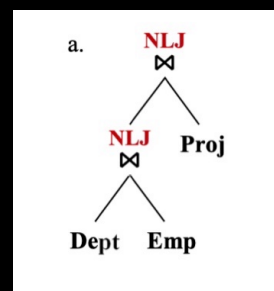
⑥ find number of pages for all relations

$$\begin{aligned} \text{NPages (Emp)} &= \text{NTuples (Emp)} / \text{NTuples/page} \\ &= 20,000 / 20 \\ &= 1,000 \text{ Pages} \end{aligned}$$

$$\begin{aligned} \text{NPages (Dept)} &= \text{NTuples (Dept)} / \text{NTuples/page} \\ &= 5,000 / 40 \\ &= 125 \text{ Pages} \end{aligned}$$

$$\begin{aligned} \text{NPages (Proj)} &= \text{NTuples (Proj)} / \text{NTuples/page} \\ &= 1,000 / 10 \\ &= 100 \text{ Pages} \end{aligned}$$

a) this left-deep plan is joining Dept with Emp using Nested Loop Join and then joining the result with Proj using Nested Loop Join



cost analysis :

① find cost of child join (join b/w first 2 relations)

$$\begin{aligned}
 \text{Cost (D JOIN E)} &: \text{NPages (R)} + \text{NPages (R)} \times \text{NPages (S)} \\
 &\quad \text{PNLJ} \qquad \qquad \text{let D be the outer relation (R)} \\
 &= \text{NPages (D)} + \text{NPages (D)} \times \text{NPages (E)} \\
 &= 125 + 125 \times 1000 \\
 &= 125,125 \text{ I/O}
 \end{aligned}$$

② find result size of child join in PAGES!!!

$$\begin{aligned}
 \text{result size (D JOIN E)} &= \prod \text{NTuples (R}_i\text{)} \times \prod \text{RF}_i \\
 &= \text{NTuples (D)} \times \text{NTuples (E)} \times \frac{1}{\text{NKeys (J)}} \rightarrow \begin{array}{l} \text{\# distinct keys /} \\ \text{values of the} \\ \text{joining attribute} \end{array} \\
 &= 5000 \times 20,000 \times \frac{1}{500} \rightarrow \begin{array}{l} \text{\# distinct val since that's} \\ \text{what you're joining E and D on} \end{array} \\
 &= 200,000 \text{ tuples}
 \end{aligned}$$

$$\begin{aligned}
 \text{NPages (D JOIN E)} &= 200,000 \text{ tuples} / 100 \text{ tuples/page} \\
 &= 2,000 \text{ pages}
 \end{aligned}$$

③ find cost of parent join

$$\text{Cost (JOIN P)} = \cancel{\text{NPages (D JOIN E)}} + \text{NPages (D JOIN E)} \times \text{NPages (P)}$$

outer relation is now the result of child join

$$= 2000 \times 100$$

$$= 200,000 \text{ I/O}$$

→ apply pipelining when performing parent join, so we can discard cost of 1st read of left input

→ pipelining is the direct "streaming" in memory of the output of one operation → (in this case, a join) as the input of another operation (eg another join) without writing the output to disk

④ add cost of child join and parent join

$$\text{Total cost} = 125,125 + 200,000$$

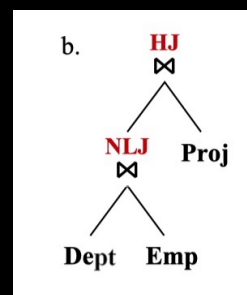
$$= 325,125 \text{ I/O}$$



b) this left-deep plan is joining Dept with Emp using Nested Loop Join and then joining the result with Proj using Hash Join

cost analysis:

→ shorthand for "Join"



①  $\text{cost}(D \times E) \underset{\text{NLJ}}{=} \boxed{125,125 \text{ I/O}}$  same as before

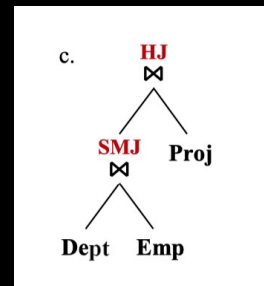
② result size  $(D \times E)$  in pages = 2,000 pages same as before

③  $\underset{\text{HJ}}{\text{Cost}(\times P)} = 3 \times (\text{NPages}(D \times E) + \text{NPages}(P)) - \text{NPages}(D \times E)$   
 $= 3 \times (2,000 + 100) - 2,000$   
 $\boxed{= 4,300 \text{ I/O}}$

↳ discard cost of first read of left input due to pipelining

④ Total cost =  $125,125 + 4,300$   
 $\boxed{= 129,425 \text{ I/O}}$

c) this left-deep plan is joining Dept with Emp using Sort-Merge join and joining results with Proj using Hash Join



$$\begin{aligned}
 \textcircled{1} \quad \text{cost}(D \bowtie E) &= \text{NPages}(D) + \text{NPages}(E) \\
 &\quad \text{SMJ} \quad + 2 \times \text{NPages}(D) \times \text{num\_passes}(D) \\
 &\quad + 2 \times \text{NPages}(E) \times \text{num\_passes}(E) \\
 &= 125 + 1000 + 2 \times 125 \times 2 + 2 \times 1000 \times 2 \\
 &= \boxed{5,625 \text{ I/O}}
 \end{aligned}$$

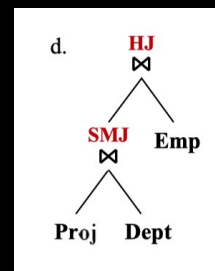
$\textcircled{2}$  result size  $(D \bowtie E)$  in pages = 2,000 pages same as before

$$\begin{aligned}
 \textcircled{3} \quad \text{cost}(R \bowtie P) &= \boxed{4,300 \text{ I/O}} \quad \text{Same as before} \\
 &\quad \text{HJ}
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{4} \quad \text{Total cost} &= 5,625 + 4,300 \\
 &= \boxed{9,925 \text{ I/O}}
 \end{aligned}$$

Note: if data file is already sorted, cost of sorting can be ignored in SMJ

d) this left-deep plan is joining Proj with Dept using Sort-merge join and joining results with Emp using Hash Join



$$\textcircled{1} \text{ cost } (P \times D) = \text{NPages}(P) + \text{NPages}(D) + 2 \times \text{NPages}(P) \times \text{num. passes}(P) + 2 \times \text{NPages}(D) \times \text{num. passes}(D)$$

$$= 100 + 125 + 2 \times 100 \times 2 + 2 \times 125 \times 2$$

$$= 1,125 \text{ I/O}$$

$$\textcircled{2} \text{ result size } (P \times D) = \prod \text{NTuples}(R_i) \times \prod R F_i$$

$$= \text{NTuples}(P) \times \text{NTuples}(D) \times \frac{1}{\text{NKeys}(J)} \rightarrow \# \text{ unique values of joining attribute}$$

$$= 1,000 \times 5,000 \times \frac{1}{1,000} \rightarrow \# \text{ distinct projid since that's what you're joining P and D on}$$

$$= 5,000 \text{ Tuples}$$

$$\text{NPages}(P \times D) = 5,000 \text{ tuples} / 50 \text{ tuples/page}$$

$$= 100 \text{ Pages}$$

$$\textcircled{3} \text{ cost } (\times E) = 3 \times (\text{NPages}(P \times D) + \text{NPages}(E)) - \text{NPages}(P \times D)$$

$$= 3 \times (100 + 1,000) - 100$$

$$= 3,200 \text{ I/O}$$

↳ discard cost of first read of left input due to pipelining

$$\textcircled{4} \text{ Total cost} = 1,125 + 3,200$$

$$= 4,325 \text{ I/O}$$