

Kamyar Karimifakhr 1068176

Assignment 3

Consider two relations A and B. A has 80,000 tuples, and B has 100,000 tuples. Both relations store 100 tuples per page. Consider the following SQL statement:

```
SELECT *  
FROM A INNER JOIN B  
ON A.a = B.a;
```

We wish to evaluate an equijoin between A and B, with an equality condition $A.a = B.a$. There are 102 buffer pages available for this operation. Both relations are stored as (unsorted) heap files. Neither relation has any indexes built on it.

Consider the alternative join strategies described below and calculate the cost of each alternative. Evaluate the algorithms using the number of disk I/O's (i.e., pages) as the cost. For each strategy, provide the formulae you use to calculate your cost estimates.

- a) Page-oriented Nested Loops Join. Consider A as the outer relation.

$$\begin{aligned}\text{Cost(PNLJ)} &= \\ &= \text{NPages(Outer)} + ((\text{NPages(Outer)}) * (\text{NPages(Inner)})) \\ &= 800 + (800 * 1000) = 800800 \text{ I/Os}\end{aligned}$$

- b) Block-oriented Nested Loops Join. Consider A as the outer relation.

$$\begin{aligned}\text{NBlocks(Outer)} &= (\text{NPages(Outer)}) / (\text{Block size} - 2) = 800 / 100 = 8 \text{ blocks} \\ \text{Cost(BNLJ)} &= \text{NPages(Outer)} + ((\text{NBlocks(Outer)}) * (\text{NPages(Inner)})) \\ &= 800 + (8 * 1000) = 8800 \text{ I/Os}\end{aligned}$$

- c) Sort-Merge Join. Assume that Sort-Merge Join can be done in 2 passes.

$$\begin{aligned}\text{Sort(A)} &= 2 * \text{Number Passes} * \text{NPages(A)} \\ &= 2 * 2 * 800 = 3200 \text{ I/Os}\end{aligned}$$

$$\begin{aligned}\text{Sort(B)} &= 2 * \text{Number Passes} * \text{NPages(B)} \\ &= 2 * 2 * 1000 = 4000 \text{ I/Os}\end{aligned}$$

$$\begin{aligned} \text{Cost}(A \text{ SMJ } B) \\ = \text{NPages}(A) + \text{NPages}(B) = 800 + 1000 = 1800 \end{aligned}$$

$$\begin{aligned} \text{Cost}(\text{SMJ}) &= \text{sort}(A) + \text{sort}(B) + \text{Cost}(A \text{ SMJ } B) = \\ 3200 + 4000 + 1800 &= 9000 \text{ I/Os} \end{aligned}$$

d) Hash Join

$$\begin{aligned} \text{Cost}(\text{HJ}) \\ = 2 * \text{NPages}(\text{Outer}) + 2 * \text{NPages}(\text{Inner}) + \text{NPages}(\text{Outer}) + \text{NPages}(\text{Inner}) = 3 * \\ \text{NPages}(\text{Inner}) + 3 * \text{NPages}(\text{Outer}) = \\ 3 * 800 + 3 * 100 = 5400 \text{ I/Os} \end{aligned}$$

e) What would be the lowest possible cost to perform this query, assuming that no indexes are built on any of the two relations, and assuming that sufficient buffer space is available? What would be the minimum buffer size required to achieve this cost? Explain briefly.

Total cost = Npages(outer) + Npages(inner) (Each relation only read once!)

$$= 800 + 1000 = 1800 \text{ I/Os}$$

$$\text{Minimum buffer needed} = \min(\text{Npages}(\text{outer}), \text{Npages}(\text{inner})) + 1 + 1 = 802$$

Question 2

SELECT * FROM JobSeekers

WHERE city = 'Melbourne' AND soughtsalary > 80,000;

- a) Compute the reduction factors and the estimated result size in number of tuples.

$$RF_{city} = 1/\text{key}(I) = 1/8$$

$$RF_{citysoughtsalary} = (\text{high} - \text{value})/(\text{high} - \text{low}) = (160000 - 80000)/100000 = 8/10$$

$$\text{Result size} = \text{Ntuples}(\text{Jobseeker}) * \text{product of RFs} = 1000000 * 1/8 * 8/10 = 100000$$

- b) Compute the estimated cost in number of disk I/O's of the best plan if a clustered B+ tree index on (city, soughtsalary) is the only index available. Suppose there are 2,000 index pages. Discuss and calculate alternative plans.

First calculating the reduction factors

$$RF_{city} = 1/\text{key}(I) = 1/8$$

$$RF_{soughtsalary} = 8/10$$

Cost of a clustered B+ tree index

$$\text{Cost} = \text{Product of RFs of matching condition} * (\text{NPages}(I) + \text{NPages}(\text{Jobseeker}))$$

$$= 1/8 * 8/10 * (2000 + 10000) = 1200 \text{ I/Os}$$

$$\text{Full table scan} = \text{Npages} = 10000 \text{ I/OS}$$

As a result, the Full table scan is a cheaper option.

- c) Compute the estimated cost in number of disk I/O's of the best plan if an unclustered B+ tree index on (soughtsalary) is the only index available. Suppose there are 2,000 index pages. Discuss and calculate alternative plans.

$$\text{Cost} = \text{Product of RFs of matching condition} * (\text{NPages}(I) + \text{Ntuples}(\text{Jobseeker})) = 8/10$$

$$*(1000000 + 2000) = 801600 \text{ I/Os}$$

Alternative plan: Full table scan = 10000 I/Os is cheapest option

- d) Compute the estimated cost in number of disk I/O's of the best plan if an unclustered Hash index on (city) is the only index available. Discuss and calculate alternative plans.

Cost of unclustered Hash index:

Cost = product of reduction factors * hash lookup cost * Ntuples(Jobseeker)

$$= 8/10 * 2.2 * (100 * 10000) = 1760000 \text{ I/Os}$$

Alternative plan : Full table scan = 10000 which is a the best plan here

- e) Compute the estimated cost in number of disk I/O's of the best plan if an unclustered Hash index on (soughtsalary) is the only index available. Discuss and calculate alternative plans.

Best plan if an unclustered Hash index on city is the only index available.

Cost = product of reduction factors * hash lookup cost * Ntuples(Jobseeker)

$$= 1/8 * 2.2 * (100 * 10000) = 275000 \text{ I/Os}$$

Alternative plan:

Full table scan = 10000 I/Os is the cheaper option

Question 3

```
SELECT D.dname, F.budget
FROM Emp E, Dept D, Finance F

WHERE E.did = D.did
AND D.did = F.did
AND E.sal > 100,000
AND E.hobby IN ('diving', 'soccer');
```

The system's statistics indicate that employee salaries range from 50,000 to 150,000, and employees enjoy 50 different hobbies. There is a total of 25,000 employees and 1,200 departments (each with corresponding financial record in the Finance relation) in the database. Each relation fits 100 tuples in a page. Suppose there exists a clustered B+ tree index on (Dept.did) and a clustered B+ tree index on (Emp.salary), both of size 50 pages.

a) Compute the reduction factors and the estimated result size in number of tuples.

Total tuple size = $N_{\text{tuple}}(\text{Emp}) * N_{\text{tuple}}(\text{Dep}) * N_{\text{tuple}}(\text{Fin})$

= $25000 * 1200 * 1200 = 36000000000$

First we calculate all the reduction factors

$RF(E.did = D.did) = 1/25000$

$RF(D.did = F.did) = 1/1200$

$RF(E.sal > 100,000) =$

$(\text{highest} - \text{value}) / (\text{highest} - \text{lowest}) = 150000 - 100000 / 150000 - 50000 = 50000 / 100000 = 0.5$

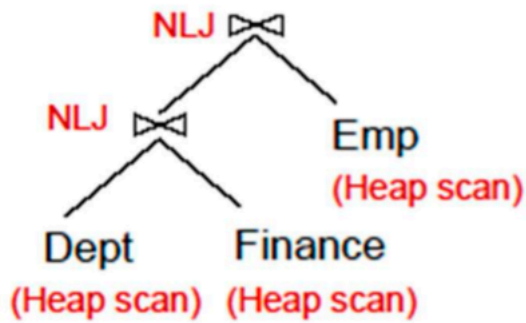
$RF(E.hobby \text{ IN } ('diving', 'soccer')) = 2/50 = 1/25$

ESTIMATION :

Tuples in total * Product of RFs =

$36000000000 * 1/25000 * 1/1200 * 0.5 * 1/25 = 23 \text{ Tuples}$

1)



Size of (D NJL F) = $N_{\text{tuples}}(D) * N_{\text{tuples}}(F) * 1/\text{key}(I)$
 $= 1200 * 1200 * 1/1200 = 1200 / 50 = 24$ pages

D scanning cost = $1200/50 = 24$ I/Os

Cost of (D NJL F) = $N_{\text{pages}}(D) + (N_{\text{pages}}(d) * N_{\text{pages}}(F)) = 24 + 24 * 12 = 312$ I/Os

We do not need to (D NJL F) as we already read that so we remove it from (NLJ E) calculation.

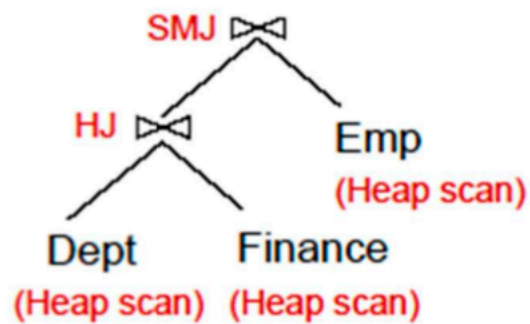
$N_{\text{pages}}(E) = 25000/100 = 250$ pages

Cost to join with employee = $N_{\text{pages}}(D \text{ NJL } F)$
 $N_{\text{pages}}(E) =$

Cost (NJL E) = $24 * 250 = 6000$ I/Os

Total cost = $312 + 6000 = 6312$ I/Os

2)



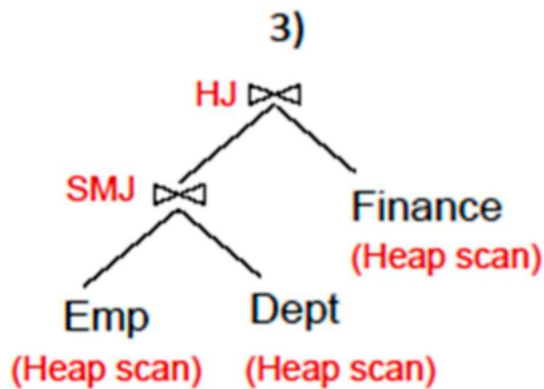
$N_{\text{tuples}}(D \text{ HJ } F) = N_{\text{tuples}}(D) * N_{\text{tuples}}(F) * 1/\text{key}(I) = 1200 * 1200 * 1/1200 = 1200/100 =$
 1200 Tuples so $1200/50 = 24$ pages

Cost of $(D \text{ HJ } F) = 3 * N_{\text{pages}}(\text{Outer}) + 3 * N_{\text{pages}}(\text{Inner}) = 3 * 50 + 3 * 12 = 186$ I/Os

Since the result of $(D \text{ HJ } F)$ has already been stored into memory once, it can be pipelined for the latter stage

Cost $(D \text{ HJ } F) \text{ SMJ } E$
 $= (2 * N_{\text{Passes}} * N_{\text{Pages}}(D \text{ HJ } F)) + N_{\text{Pages}}(E) = (2 * 2 * 24) + 250 = 346$ I/Os

Costs in total = $186 + 346 = 532$ I/Os



Number of resulting tuples for Emp SMJ Dept = $1 / NKeys(I) \times NTuples(Emp) \times NTuples(Dept)$
 $= 1 / 25000 \times 25000 \times 1200 = 1200$ Tuples

Number of pages of (Emp SMJ Dept) = $1200 / 50 = 24$ pages

Sort Emp = $2 \times NPasses \times NPages(Emp) = 2 \times 2 \times 250 = 1000$ I/O

Sort Dept = 0 (As it is already sorted since it's in a clustered Index)

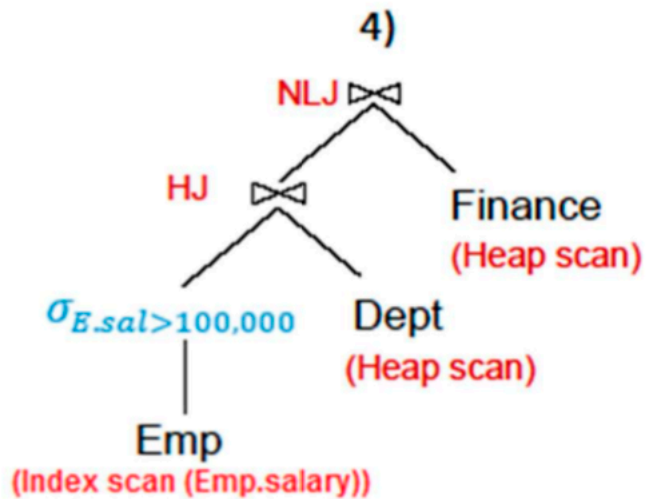
Cost of merge joining Emp and Dept = $NPages(Emp) + NPages(Dept)$
 $= 250 + 50 = 300$ I/O

Cost of SMJ = $1000 + 300 = 1300$ I/Os

We already stored the (E join D) once so we need to multiply it by 2.

Cost of (HJ F) = $2 * NPages(Emp JOIN Dept) + 3 * NPages(Finance) = 2 * 24 + 3 * 12 = 84$ I/Os

Total cost = $1300 + 84 = 1384$ I/Os



$$RF = (\text{High} - \text{Value}) / (\text{High} - \text{Value}) = 50000 / 100000 = 0.5$$

$$\text{Selection cost} = \text{Npages}(I) + \text{Npages}(E) * RF$$

$$= (50 + 250) * 0.5 = 150 \text{ I/Os}$$

$$\text{Number of pages after selection} = (\text{rf} * \text{Ntuples}) / 100$$

$$(0.5 * 25000) / 100 = 12500 / 100 = 125 \text{ pages}$$

$$\begin{aligned} \text{Size of } Emp_{E.sal > 100000} \text{ join Dept} &= \text{Ntuples}(Emp_{E.sal > 100000}) * \text{NTuples}(\text{Dept}) * 1 / \text{Nkeys}(I) \\ &= 12500 * 1200 / 12500 * 0.5 = 600 \text{ Tuples} \Rightarrow \end{aligned}$$

$$\text{Result size} = 600 / 50 = 12 \text{ Pages}$$

$$\text{Cost HJ}(Emp_{E.sal > 100000} \text{ join Dept}) =$$

$$3 * \text{Npages}(\text{Emp}) + 3 * \text{Npages}(\text{Dept}) = 3 * 125 + 3 * 50 = 525 \text{ I/Os}$$

$\text{Npages}(\text{Emp}_{E.\text{sal}>100000} \text{ join Dept})$ is stored in memory before so no need to calculate again.

$\text{Cost (NLJ Finance)} = \text{Npages}(\text{Emp}_{E.\text{sal}>100000} \text{ join Dept}) * \text{Npages}(\text{Fin}) = 12 * 12 = 144 \text{ I/Os}$

$\text{Total cost} = \text{cost of selection} + \text{cost}(\text{Emp}_{E.\text{sal}>100000} \text{ join Dept}) +$

$\text{cost}((\text{Emp}_{E.\text{sal}>100000} \text{ hash join Dept}) \text{ NLJ}(\text{Finance}))$

$\text{Total cost} = 150 + 144 + 525 = 819 \text{ I/Os}$