

## Question 1

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

$A = 80000 \text{ tuples} = 800 \text{ pages}$

$B = 100000 \text{ tuples} = 1000 \text{ pages}$

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

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

102 buffer pages available

$$\text{NBlocks(Outer)} = \text{NPages(Outer)} / (B - 2) = 800 / (102 - 2) = 8 \text{ blocks}$$

$$\begin{aligned}\text{Cost (BNJL)} &= \text{NPages(Outer)} + \text{NBlocks(Outer)} * \text{NPages(Inner)} \\ &= 800 + 8 * 1000 \\ &= 8800 \text{ I/O}\end{aligned}$$

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

$$\text{NumPasses} = 2$$

$$\text{Sort(R)} = \text{External Sort Cost} = 2 * \text{NumPasses} * \text{NPages(R)}$$

$$\begin{aligned}\text{Cost (SMJ)} &= \text{Sort(Outer)} + \text{Sort(Inner)} + \text{NPages(Outer)} + \text{NPages(Inner)} \\ &= 5 * 800 + 5 * 1000 \\ &= 9000 \text{ I/O}\end{aligned}$$

d) Hash Join

$$\begin{aligned}\text{Cost (HJ)} &= 2 * \text{NPages(Outer)} + 2 * \text{NPages(Inner)} + \text{NPages(Outer)} + \text{NPages(Inner)} \\ &= 3 * 800 + 3 * 1000 \\ &= 5400 \text{ I/O}\end{aligned}$$

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.

Out of the join strategies calculated, a Hash Join initially seems to have the lowest cost.

However, given a sufficient number of buffer pages the Block-oriented Nested Loops Join becomes optimal by solving for the number of blocks required to match the Hash Join:

$$800 + 800 / (x - 2) * 1000 = 5400 \text{ I/O}$$

$$800 / (x - 2) = 4.6$$

$$x = 800 / 4.6 + 2 \approx 175.913 = 176 \text{ buffer pages (need whole number of pages)}$$

## Question 2

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

8 cities

Salary ranges from 60000 to 100000

$NTuples(\text{JobSeekers}) = 10000 \text{ pages} * 100 \text{ tuples/page} = 1000000 \text{ tuples}$

Col = value:  $RF = 1/NKeys(\text{Col})$

$RF(\text{city}) = 1 / 8$

$= 0.125$

$Col > \text{value } RF = (\text{High}(\text{Col}) - \text{value}) / (\text{High}(\text{Col}) - \text{Low}(\text{Col}))$

$RF(\text{soughtsalary}) = (160000 - 80000) / (160000 - 60000)$

$= 80000 / 100000$

$= 0.8$

$\text{ResultSize} = NTuples(R) \prod_{i=1..n} RF_i$

$= 1000000 * 0.125 * 0.8$

$= 100,000 \text{ tuples}$

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.

$\text{Cost (B+Tree)} = (NPages(I) + NPages(R)) * \prod_{i=1..n} RF_i$

$= (2000 + 10000) * 0.125 * 0.8 = 1200 \text{ I/O}$

$\text{Cost (Heap)} = NPages(R)$

$= 10000 \text{ I/O}$

The cost of using the clustered B+ tree is far less than the heap scan, so using the index is the optimal plan in this case.

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.

$\text{Cost (B+Tree)} = (NPages(I) + NTuples(R)) * \prod_{i=1..n} RF_i$

$= (2000 + 1000000) * 0.8$

$= 801600 \text{ I/O}$

$\text{Cost (Heap)} = NPages(R)$

$= 10000 \text{ I/O}$

With the provided unclustered index, the resulting cost is far greater than if a simple Heap scan was used, making a Heap scan the optimal plan in this case.

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.

$$\begin{aligned}\text{Cost (HashIndex)} &= \text{NTuples(R)} * \prod_{i=1..n} \text{RF}_i * 2.2 \\ &= 1000000 * 0.125 * 2.2 \\ &= 275000 \text{ I/O}\end{aligned}$$

$$\begin{aligned}\text{Cost (Heap)} &= \text{NPages(R)} \\ &= 10000 \text{ I/O}\end{aligned}$$

The cost of a Heap scan is again much less than the provided unclustered Hash index, making a Heap scan the optimal plan in this case.

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.

$$\begin{aligned}\text{Cost (HashIndex)} &= \text{NTuples(R)} * \prod_{i=1..n} \text{RF}_i * 2.2 \\ &= 1000000 * 0.8 * 2.2 \\ &= 1760000 \text{ I/O}\end{aligned}$$

$$\begin{aligned}\text{Cost (Heap)} &= \text{NPages(R)} \\ &= 10000 \text{ I/O}\end{aligned}$$

The cost of a Heap scan is again far less than the provided unclustered Hash index, making a Heap scan the optimal plan in this case.

### Question 3

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

Sal ranges from 50000 to 150000

50 different hobbies

Col > value:  $RF = (High(Col) - value) / (High(Col) - Low(Col))$

$$\begin{aligned} RF(sal) &= (150000 - 100000) / (150000 - 50000) \\ &= 50000 / 100000 \\ &= 0.5 \end{aligned}$$

Col = value:  $RF = 1/NKeys(Col)$

$$\begin{aligned} RF(hobby) &= 2 / 50 \text{ ('diving' and 'soccer')} \\ &= 0.04 \end{aligned}$$

Col\_A = Col\_B (for joins):  $RF = 1 / (Max(NKeys(Col_A), NKeys(Col_B)))$

$$\begin{aligned} RF(did\ Emp \ \& \ Dept) &= 1 / (Max(25000, 1200)) \\ &= 1 / 25000 \end{aligned}$$

Col\_A = Col\_B (for joins):  $RF = 1 / (Max(NKeys(Col_A), NKeys(Col_B)))$

$$\begin{aligned} RF(did\ Dept \ \& \ Finance) &= 1 / (Max(1200, 1200)) \\ &= 1 / 1200 \end{aligned}$$

$$\begin{aligned} ResultSize &= \prod_{j=1..k} NTuples(R_j) \prod_{i=1..n} RF_i \\ &= (25000 * 1200 * 1200) * 0.5 * 0.04 * (1 / 1200) * (1 / 25000) \\ &= 24 \text{ tuples} \end{aligned}$$

b) Compute the cost in number of disk I/O's of the plans shown below. Assume that sorting of any relation (if required) can be done in 2 passes. NLJ is a Page-oriented Nested Loops Join. Assume that did is the candidate key, and that 50 tuples of a resulting join between Emp and Dept fit in a page. Similarly, 50 tuples of a resulting join between Finance and Dept fit in a page. Any selections/projections not indicated on the plan are performed “on the fly” after all joins have been completed.

1)

Calculating Cost:

Dept x Finance

$$1200 \text{ Dept tuples} / 100 \text{ tuples/page} = 12 \text{ pages}$$

$$1200 \text{ Finance tuples} / 100 \text{ tuples/page} = 12 \text{ pages}$$

$$25000 \text{ Emp tuples} / 100 \text{ tuples/page} = 250 \text{ pages}$$

$$\text{Cost (PNJL)} = NPages(Outer) + NPages(Outer) * NPages(Inner)$$

$$\begin{aligned} \text{Cost (Dept x Finance)} &= 12 + 12 * 12 \\ &= 156 \end{aligned}$$

(Dept x Finance) x Emp

$$\begin{aligned} ResultSize \text{ (Dept x Finance)} &= \prod_{j=1..k} NTuples(R_j) \prod_{i=1..n} RF_i \\ &= 1200 * 1200 * 1 / 1200 = 1200 \text{ tuples} / 50 \text{ tuples/page} = 24 \text{ pages} \end{aligned}$$

$$\text{Cost (PNJL)} = NPages(Outer) + NPages(Outer) * NPages(Inner)$$

$$\begin{aligned} \text{Cost (x Emp)} &= 24 + 24 * 250 \\ &= 6000 \end{aligned}$$

$$\begin{aligned} \text{Total Cost} &= 12 + 12 * 12 + 24 * 250 \\ &= 6156 \text{ I/O} \end{aligned}$$

2)

Calculating Cost:

Dept x Finance

1200 Dept tuples / 100 tuples/page = 12 pages

1200 Finance tuples / 100 tuples/page = 12 pages

25000 Emp tuples / 100 tuples/page = 250 pages

2 sorting passes

$$\begin{aligned}\text{Cost (HJ)} &= 2 * \text{NPages(Outer)} + 2 * \text{NPages(Inner)} + \text{NPages(Outer)} + \text{NPages(Inner)} \\ &= 3 * 12 + 3 * 12 \\ &= 72\end{aligned}$$

(Dept x Finance) x Emp

$$\begin{aligned}\text{ResultSize (Dept x Finance)} &= \prod_{j=1..k} \text{NTuples(R}_j) \prod_{i=1..n} \text{RF}_i \\ &= 1200 * 1200 * 1 / 1200 = 1200 \text{ tuples} / 50 \text{ tuples/page} = 24 \text{ pages}\end{aligned}$$

Cost (SMJ) = Sort(Outer) + Sort(Inner) + NPages(Outer) + NPages(Inner) Sort(R) = External

Sort Cost = 2\*NumPasses\*NPages(R)

$$\begin{aligned}\text{Cost (x Emp)} &= 2 * 2 * 24 - 24 + 2 * 2 * 250 + 24 + 250 - 250 \\ &= 72 + 1000 + 24 \\ &= 1096\end{aligned}$$

$$\begin{aligned}\text{Total Cost} &= 3 * 12 + 3 * 12 + 2 * 2 * 24 - 24 + 2 * 2 * 250 + 24 + 250 - 250 \\ &= 1168 \text{ I/O}\end{aligned}$$

3)

Calculating Cost:

Emp x Dept

1200 Dept tuples / 100 tuples/page = 12 pages

1200 Finance tuples / 100 tuples/page = 12 pages

25000 Emp tuples / 100 tuples/page = 250 pages

2 sorting passes

Cost (SMJ) = Sort(Outer) + Sort(Inner) + NPages(Outer) + NPages(Inner) Sort(R) = External

Sort Cost = 2\*NumPasses\*NPages(R)

$$\begin{aligned}\text{Cost (Emp x Dept)} &= 5 * 250 + 5 * 12 \\ &= 1310\end{aligned}$$

(Emp x Dept) x Finance

$$\begin{aligned}\text{ResultSize (Emp x Dept)} &= \prod_{j=1..k} \text{NTuples(R}_j) \prod_{i=1..n} \text{RF}_i \\ &= 25000 * 1200 * 1 / 25000 = 1200 \text{ tuples} / 50 \text{ tuples/page} = 24 \text{ pages}\end{aligned}$$

Cost (HJ) = 2 \* NPages(Outer) + 2 \* NPages(Inner) + NPages(Outer) + NPages(Inner)

$$\begin{aligned}\text{Cost (x Finance)} &= 3 * 2 * 24 + 3 * 12 \\ &= 84\end{aligned}$$

$$\begin{aligned}\text{Total Cost} &= 5 * 250 + 5 * 12 + 2 * 24 + 3 * 12 = 84 \\ &= 1394 \text{ I/O}\end{aligned}$$

4)

Calculating Cost:

Emp

1200 Dept tuples / 100 tuples/page = 12 pages

1200 Finance tuples / 100 tuples/page = 12 pages

25000 Emp tuples / 100 tuples/page = 250 pages

Cost (B+Tree) = (NPages(I) + NPages(R)) \*  $\prod_{i=1..n} RF_i$

Cost ( $\sigma_{Emp.sal > 100000}$ ) = (50 + 250) \* 0.5  
= 150

$\sigma_{Emp.sal > 100000}$  x Dept

ResultSize ( $\sigma_{Emp.sal > 100000}$ ) =  $NTuples(R_j) \prod_{i=1..n} RF_i$   
= 25000 \* 0.5  
= 12500 tuples / 100 tuples/page = 125 pages

Cost (HJ) = 2 \* NPages(Outer) + 2 \* NPages(Inner) + NPages(Outer) + NPages(Inner)

Cost ( $\sigma_{Emp.sal > 100000}$  x Dept) = 3 \* 2 \* 125 + 3 \* 12  
= 286

( $\sigma_{Emp.sal > 100000}$  x Dept) x Finance

ResultSize ( $\sigma_{Emp.sal > 100000}$  x Dept) =  $\prod_{j=1..k} NTuples(R_j) \prod_{i=1..n} RF_i$   
= 12500 \* 1200 \* 1 / 25000  
= 600 tuples / 50 tuples/page = 12 pages

Cost (PNJL) = NPages(Outer) + NPages(Outer) \* NPages(Inner)

Cost (x Finance) = 12 + 12 \* 12  
= 144

Total Cost = (50 + 250) \* 0.5 + 2 \* 125 + 3 \* 12 + 12 \* 12  
= 580 I/O