

BD2 – April 9th, 2018 – solutions, V1.0

Please feel free to answer this test in English, Italian, or any mixture

1. Consider a schema $R(\underline{IdR}, A, B, \dots, IdS^*), S(\underline{IdS}, \dots), T(\underline{IdT}, C, D, \dots, IdS^*)$ and the following query

```

SELECT DISTINCT T.C. count(*)
FROM      R, T
WHERE     R.IdS = T.IdS And R.A ≤ 100 And T.D = 10
GROUP BY T.C, T.D
    
```

Assume that R is stored with a dense clustered index on A, while S and T are stored as heap files. Primary keys are R.IdR, S.IdS and T.IdT, while R.IdS and T.IdS are foreign keys that refers to S. Assume that unclustered RID-sorted index are defined on all the primary and foreign keys, and on T.D. Assume that the size of all indexes only depend on the number of RIDs, as indicated in the table below. Assume that every page is 4.000 bytes long, and that every attribute uses 4 bytes. Assume a buffer size of 100 pages. If you need Cardenas formula $\Phi(n,k)$, approximate it with $\min(n,k)$.

	NReg	NPag	NLeaf of Indexes	NKey	Min	Max
R	500.000	10.000	1.000			
S	100.000	2.000	200			
T	10.000.000	50.000	20.000			
Idx.R.A			See R	100	0	1.000
Idx.T.D			See T	100	0	1.000

- a) Is DISTINCT redundant?
 - b) Compute the selectivity factor of the three predicates in the condition
 - c) Compute the cost of accessing $\sigma_{A \leq 100}(R)$, $\sigma_{D \leq 10}(T)$ (useless for this query) and $\sigma_{D=10}(T)$, with and without index. Remember that the index on R.A is clustered.
 - d) Compute the cost of a MergeJoin plan that uses the cheapest plans for the selections (ignore the fact that IdS is not a key for neither table) (draw the plan first)
 - e) Compute the cost for an IndexNestedLoop plan where R is the outer (left-hand-side) leaf (draw the plan first)
2. a) Assume two relation $R(a,b)$ and $S(a,c)$. Define the left natural outer join(R,S) ($\text{leftjoin}(R,S)$) and give a very simple example.
- b) Write a two-lines (more or less) pseudo code specification for IndexNestedLoop. Briefly specify (using natural language or psudo-code) how one could generalize IndexNestedLoop to LeftIndexNestedLoop in order to compute left outer join
- c) Compare the cost of $\text{LeftIndexNestedLoop}(OE,OI,\text{condition})$ with the cost of $\text{IndexNestedLoop}(OE,OI,\text{condition})$
3. a) Explain (briefly) why heap organization is the one that is most commonly used
- b) In which situations the heap organization is the best one?
- c) In which situations is the hash (procedural) organization the best one?

BD2 – April 9th, 2018 – solutions – 0.1

Please feel free to answer this test in English, Italian, or any mixture

1. Consider a schema $R(\underline{IdR}, A, B, \dots, IdS^*), S(\underline{IdS}, \dots), T(\underline{IdT}, C, D, \dots, IdS^*)$ and the following query

```

SELECT DISTINCT T.C. count(*)
FROM R, T
WHERE R.IdS = T.IdS And R.A ≤ 100 And T.D = 10
GROUP BY T.C, T.D
    
```

Assume that R is stored with a dense clustered index on A, while S and T are stored as heap files. Primary keys are R.IdR, S.IdS and T.IdT, while R.IdS and T.IdS are foreign keys that refers to S.

Assume that unclustered RID-sorted index are defined on all the primary and foreign keys, and on T.D.

Assume that the size of all indexes only depend on the number of RIDs, as indicated in the table below.

Assume that every page is 4.000 bytes long, and that every attribute uses 4 bytes. Assume a buffer size of 100 pages. If you need Cardenas formula $\Phi(n,k)$, approximate it with $\min(n,k)$.

	NReg	NPag	NLeaf of Indexes	NKey	Min	Max
R	500.000	10.000	1.000			
S	100.000	2.000	200			
T	1.000.000	50.000	2.000			
Idx.R.A			See R	100	0	1.000
Idx.T.D			See T	100	0	1.000

a) Is DISTINCT redundant?

Yes, it is, since the closure of the projected attribute T.C is {T.C, T.D}, which includes the GROUP BY attributes

b) Compute the selectivity factor of the three predicates in the condition

$$s_f(R.IdS = T.IdS) = 1/NReg(S) = 1/100.000$$

$$s_f(R.A \leq 100) = 100/Max(R.A) = 100/1.000 = 1/10$$

$$s_f(T.D = 10) = 1/NKey(T.D) = 1/100$$

c) Compute the cost of accessing $\sigma_{A \leq 100}(R)$, $\sigma_{D \leq 10}(T)$ (useless for this query) and $\sigma_{D=10}(T)$, with and without index. Remember that the index on R.A is clustered.

$$C(Filter(TableScan(R), A \leq 100)) = NPag(R) = 10.000$$

$$C(IndexFilter(R, Idx.R.A, A \leq 100))$$

$$= CI + CD = s_f(R.A \leq 100) * NLeaf(Idx.R.A) + s_f(R.A \leq 100) * NPag(R) = 1/10 * 1.000 + 1/10 * 10.000$$

$$= 1.100$$

$$C(\text{Filter}(\text{TableScan}(T), D \leq 10)) = \text{NPag}(T) = 50.000$$

$$C(\text{IndexFilter}(T, \text{Idx.T.D}, D \leq 10))$$

$$= CI + CD$$

$$= \lceil s_f(T.D \leq 10) * N\text{Leaf}(\text{Idx.T.D}) \rceil + \lceil s_f(T.D \leq 10) * N\text{Key}(T.D) \rceil * \lceil \Phi(N\text{Reg}(T)/N\text{Key}(T.D), \text{NPag}(T)) \rceil$$

$$= 10/1.000 * 20.000 + 10/1.000 * 100 * \min(10.000.000/100, 50.000) = 200 + 50.000 = 50.200$$

$$C(\text{Filter}(\text{TableScan}(T), D=10)) = \text{NPag}(T) = 50.000$$

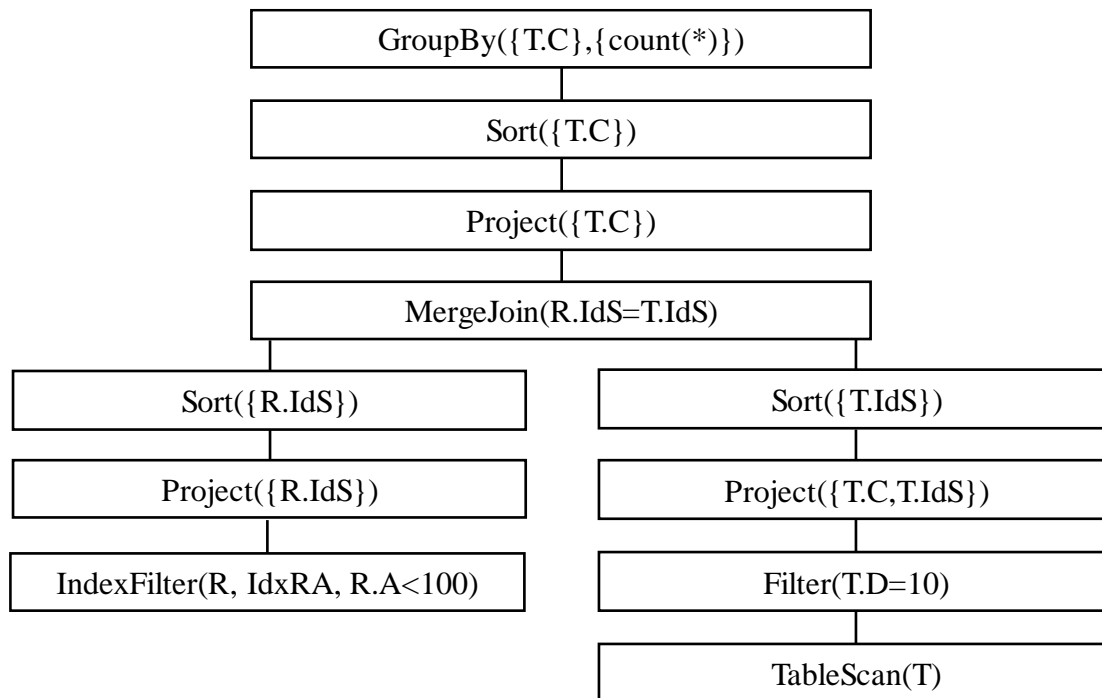
$$C(\text{IndexFilter}(T, \text{Idx.T.D}, D=10))$$

$$= CI + CD$$

$$= \lceil s_f(T.D=10) * N\text{Leaf}(\text{Idx.T.D}) \rceil + \lceil \Phi(N\text{Reg}(T)/N\text{Key}(T.D), \text{NPag}(T)) \rceil$$

$$= 1/100 * 20.000 + \min(10.000.000/100, 50.000) = 200 + 50.000 = 50.200$$

- d) Compute the cost of a MergeJoin plan that uses the cheapest plans for the selections (ignore the fact that IdS is not a key for neither table) (draw the plan first)



After the condition $T.D = 10$ we can forget the $T.D$ field – this is just a small optimization.

$$C(\text{IndexFilter}(R, \text{Idx.R.A}, A \leq 100)) = 1.100 \text{ (see the previous point)}$$

$$E\text{Reg}(\text{IndexFilter}(R, \text{Idx.R.A}, A \leq 100)) = s_f(R.A \leq 100) * N\text{Reg}(R) = 1/10 * 500.000 = 50.000$$

$$E\text{Pag}(\text{Project}(\text{IndexFilter}(\dots), \{R.IdS\})) = 50.000 * 4/4.000 = 50$$

$$C(\text{Sort}(\{R.IdS\})) = C(\text{IndexFilter}) = 1.100$$

$$C(\text{TableScan}(T)) = 50.000$$

$$E\text{Reg}(\text{Filter}(T.D=10)) = s_f(T.D=10) * N\text{Reg}(T) = 1/100 * 10.000.000 = 100.000$$

$$\text{NPag}(\text{Project}(\text{Filter}(\dots), \{T.\text{IdS}, T.C\})) = 100.000 * 8 / 4.000 = 200$$

$$C(\text{Sort}(T)) = 2 * \text{NPag}(\text{Project}(\text{Filter}(\dots), \{T.\text{IdS}, T.C\})) + C(\text{TableScan}(T)) = 400 + 50.000 = 50.400$$

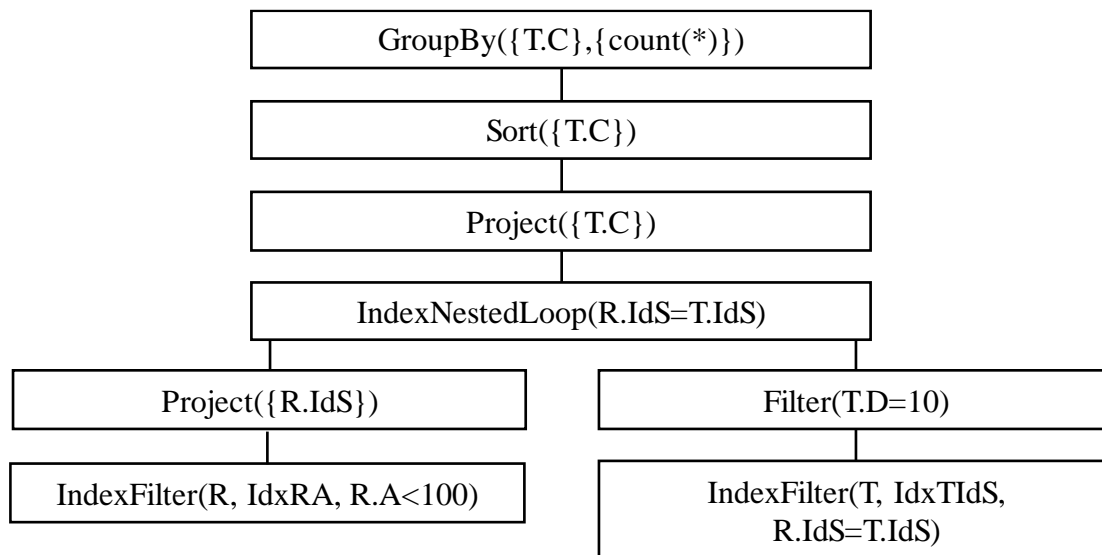
$$C(\text{MergeJoin}) = C(\text{Sort}(\dots)) + C(\text{Sort}(\dots)) = 1.100 + 50.400 = 51.500$$

$$\begin{aligned} \text{NReg}(\text{MergeJoin}) &= \text{NReg}(R) * \text{NReg}(T) * s_f(R.\text{IdS} = T.\text{IdS}) * s_f(R.A \leq 100) * s_f(T.D = 10) \\ &= 500.000 * 10.000.000 / (100.000 * 10 * 100) = 50.000 \end{aligned}$$

$$\text{NPag}(\text{Project}(\text{MergeJoin}, \{T.C\})) = 50.000 * 4 / 4.000 = 50$$

$$C(\text{Sort}(\text{Project}(\text{MergeJoin}, \{T.C\}))) = C(\text{MergeJoin}) = 51.500$$

e) Compute the cost for an IndexNestedLoop plan where R is the outer (left-hand-side) leaf.



$$C(\text{IndexFilter}(R, \text{Idx}.R.A, A \leq 100)) = 1.100$$

$$\text{EReg}(\text{IndexFilter}(R, \text{Idx}.R.A, A \leq 100)) = 50.000 \text{ (see above)}$$

$$C(\text{IndexFilter}(T, \dots, R.\text{IdS}=T.\text{IdS})) = C_I + C_D$$

$$= \lceil s_f(R.\text{IdS}=T.\text{IdS}) * \text{NLeaf}(\text{IdxTIdS}) \rceil + \lceil \Phi(\text{NReg}(T) / \text{NKey}(T.\text{IdS}), \text{NPag}(T)) \rceil$$

$$= \lceil 20.000 / 100.000 \rceil + \lceil \Phi(10.000.000 / 100.000, 20.000) \rceil$$

$$= 1 + \min(100, 20.000) = 101$$

$$C(\text{IndexNestedLoop})$$

$$= C(\text{OE}) + \text{EReg}(\text{OE}) * C(\text{OI}) = 1.100 + 50.000 * 101 = 1.100 + 5.050.000 = 5.051.100$$

2. a) Assume two relation R(a,b) and S(a,c). Define the left natural outer join(R,S) (leftjoin(R,S)) and give a very simple example.

leftjoin(R,S) contains all the tuple of join(R,S) plus all tuples t of R for which no tuple s exists in S with t.a=s.a, each extended with a null value in the c field.

Example:

R

a	b
1	X
1	Y
2	Z

S

a	c
1	A
1	B
3	C

LeftJoin(R,S)

a	b	c
1	X	A
1	X	B
1	Y	A
1	Y	B
2	Z	null

- b) Write a two-lines (more or less) pseudo code specification for IndexNestedLoop. Briefly specify (using natural language or pseudo-code) how one could generalize IndexNestedLoop to LeftIndexNestedLoop in order to compute left outer join

Index nested loop(R,S):

```
for r in R do for s in ( getIndex(S,IdxSa, r.a=s.a) ) do return(r+s);
```

LeftIndexNestedLoop(R,S):

```
for r in R do
  let Set = ( getIndex(S,IdxSa, r.a=s.a) )
  if Set = empty then return({a=r.a, b=r.b, c=null})
  else for s in Set do return(r+s);
```

- c) Compare the cost of LeftIndexNestedLoop(OE,OI,condition)) with the cost of IndexNestedLoop(OE,OI,condition)

LeftIndexNestedLoop(OE,OI,condition)) only adds some main-memory operations to IndexNestedLoop(OE,OI,condition), hence they have the same cost

3. a) Explain (briefly) why heap organization is the one that is most commonly used
 b) In which situations the heap organization is the best one?
 c) In which situations is the hash (procedural) organization the best one?