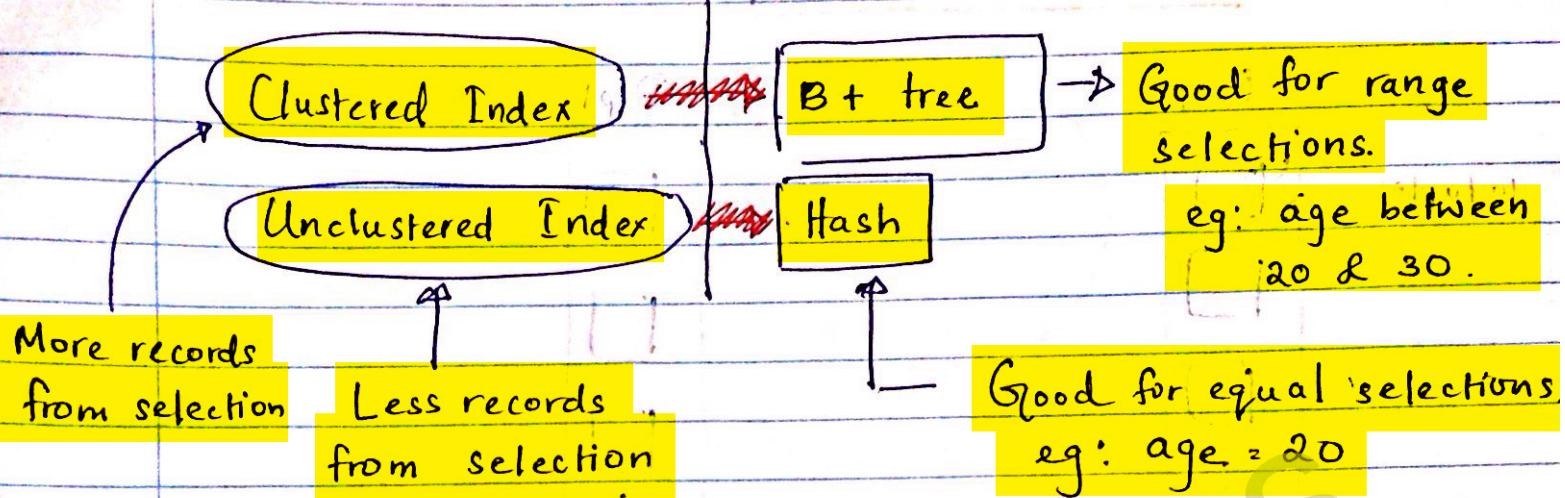


Query Optimization.

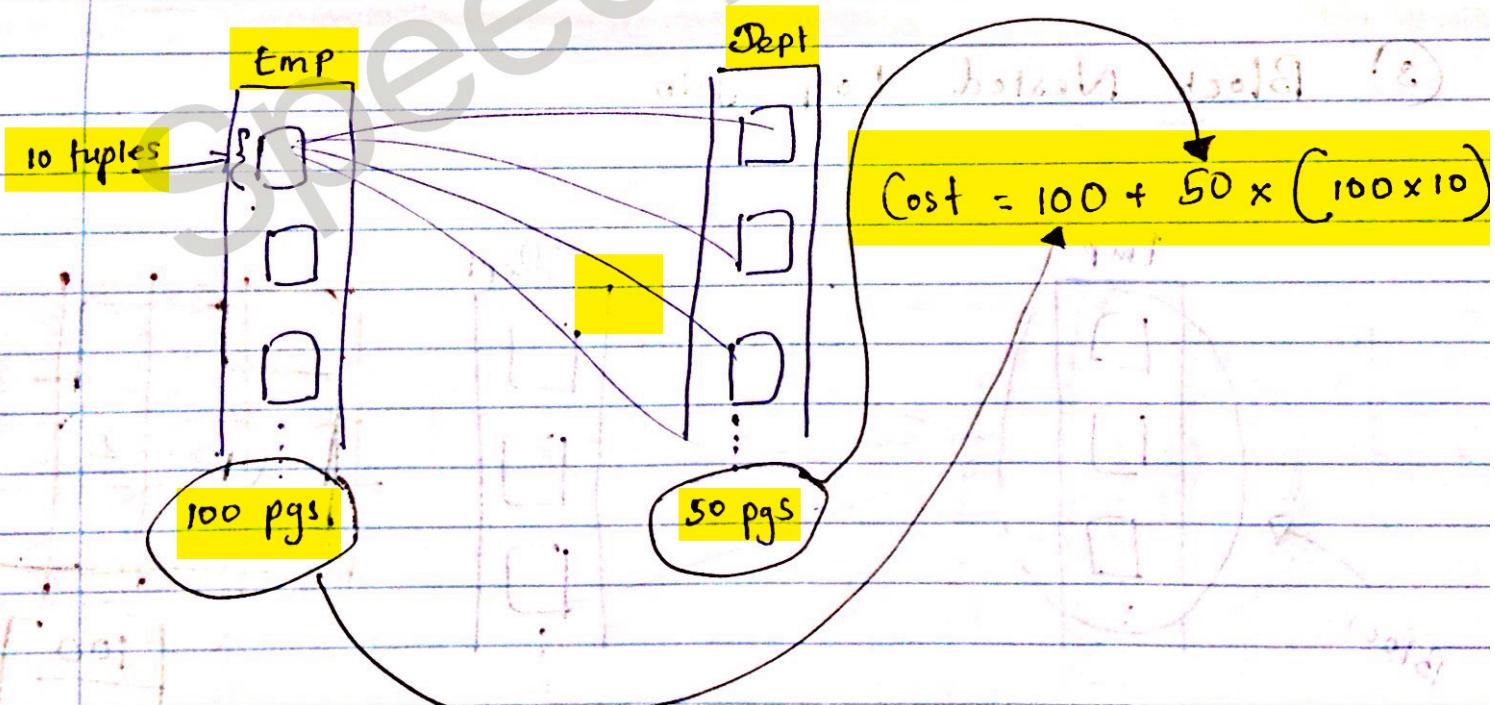
1



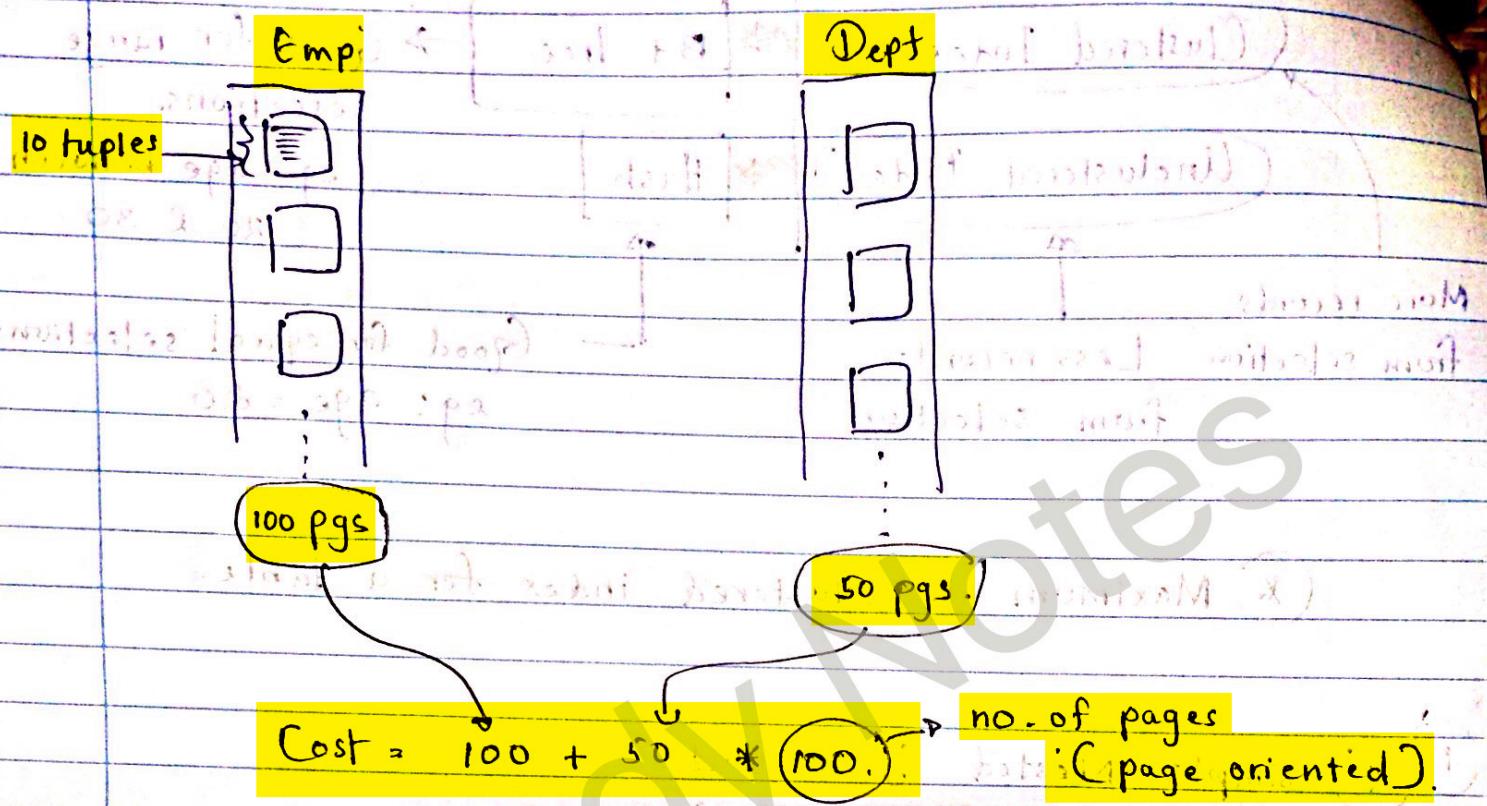
(*) Maximum clustered index for a table.

1 (*)

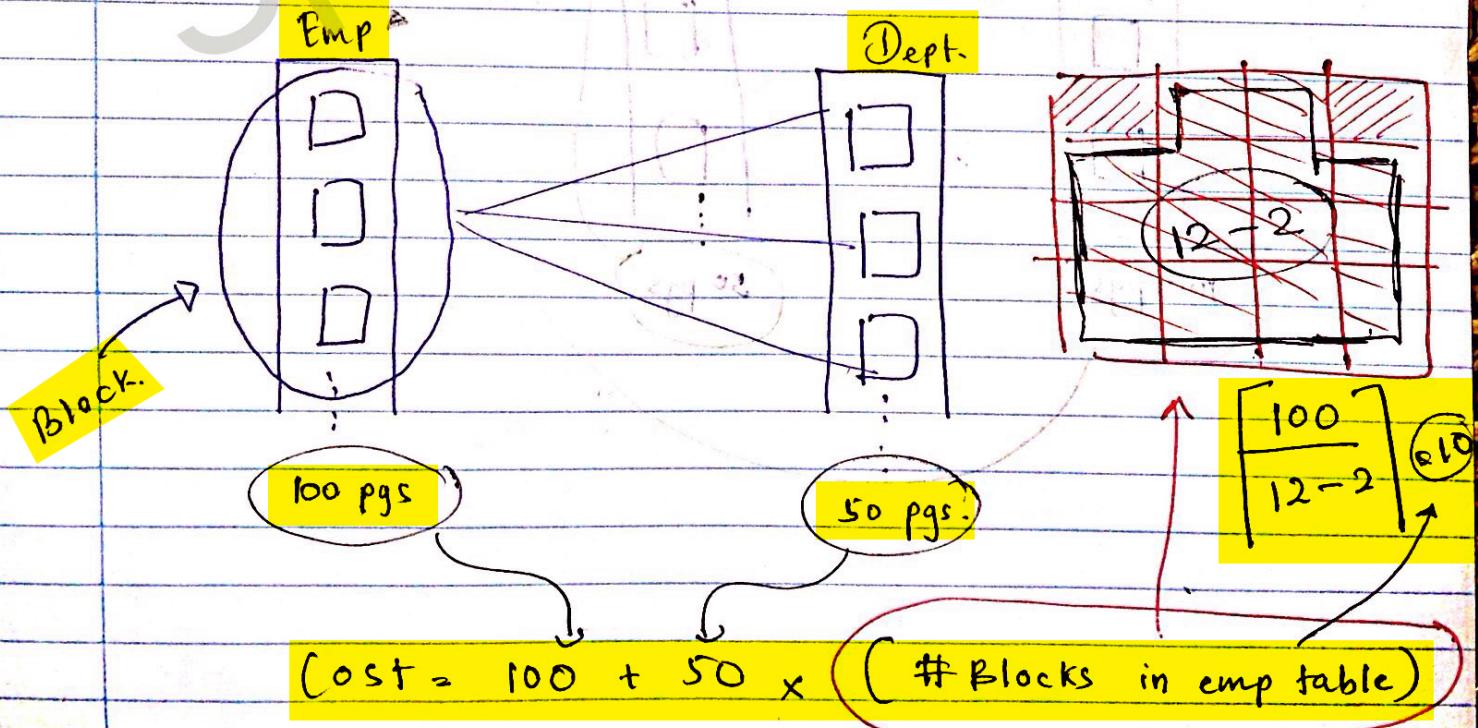
Simple Nested Join Cost



② Page Oriented Nested Loop Join.

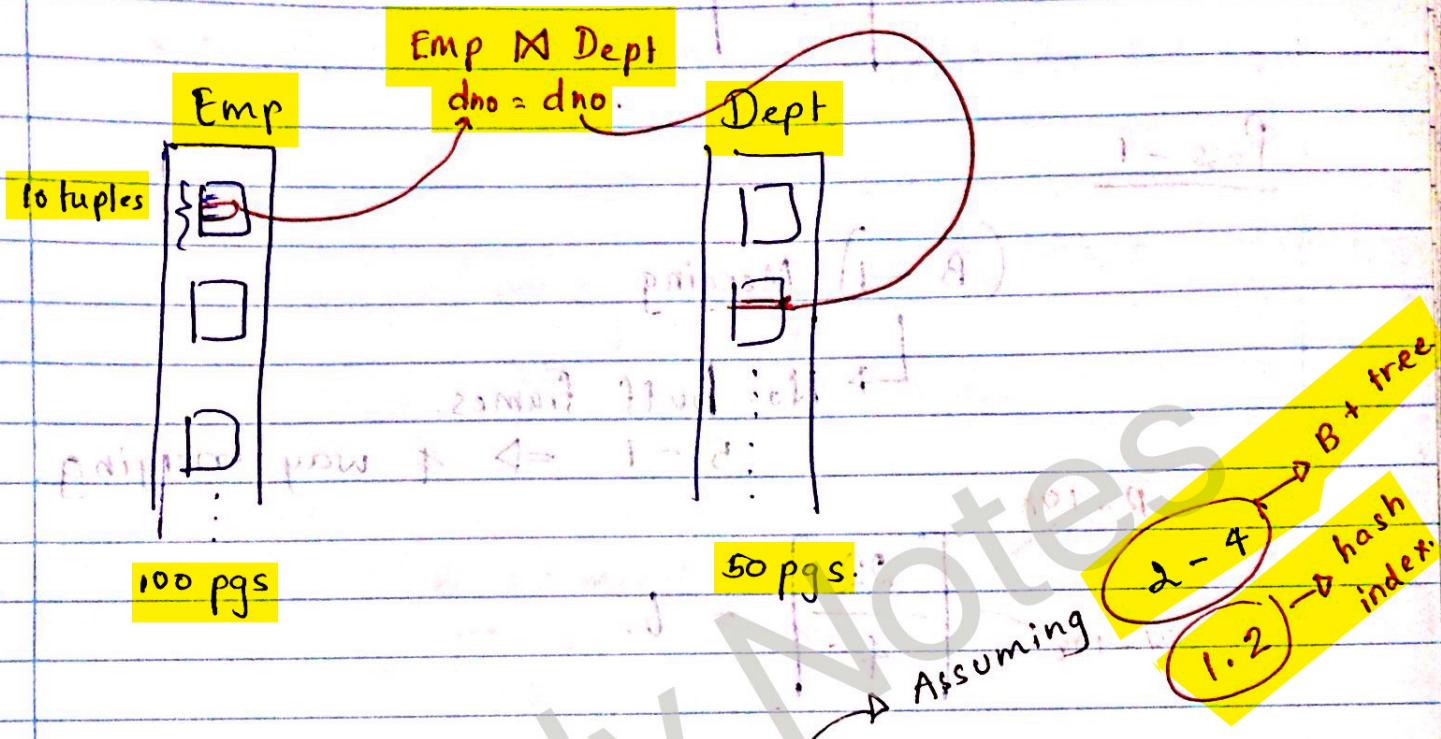


③ Block Nested Loop Join



4*

Index Nested Loop Join.

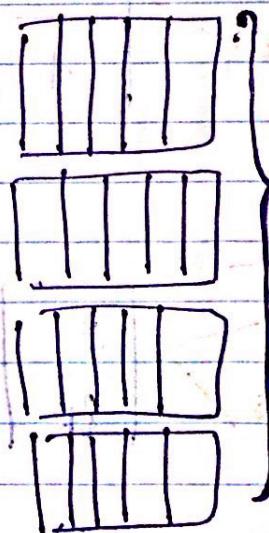
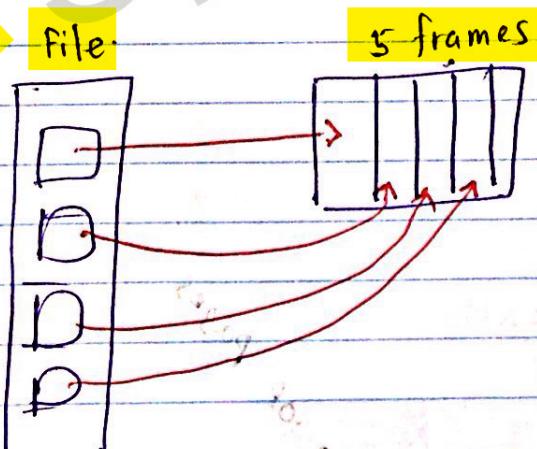


$$\text{Cost} = 100 + (\text{Index Cost}) * (100 * 10)$$

*)

External Merge Sort.

108 pages



(B-1) way merging

individually sorted.

Pass - 0

$$\left[\frac{108}{5} \right] = 22$$

R - 108
W - 108.

Pass - 1

(B - 1) Merging

→ No. buff frames.

5 - 1 \Rightarrow 4 way merging.

R - 108

W - 108

$$\left[\frac{22}{4} \right] = 6$$

Pass - 2

(B - 1) Merging

R - 108

W - 108

$$\left[\frac{6}{4} \right] = 2$$

Pass - 3

R - 108

W - 108

$$\left[\frac{2}{4} \right] = 1$$

No. of pages

External Merge Sort Cost = $2N * (\# \text{ passes})$.

$$= 2 \times 108 \times 4$$

(3).

External Merge Sort Questions.

1000 pages / 50 frames what is the sorting cost?

Pass - 0

$$\left[\frac{1000}{50} \right] = 20 //$$

Pass - 1

$$\left[\frac{20}{49} \right]$$

$$\text{Cost} = 2N \times (\# \text{ of passes})$$

$$= 2 \times 1000 \times 2$$

$$= \underline{\underline{4000}}$$

Sort - Merge Join

RMS

$$SMJ = (\text{Sorting Cost}) + (\text{Merging Cost})$$

$$= (2M \times \# \text{ passes}) + (2N \times \# \text{ passes}) + (M+N)$$

$$\downarrow \quad \downarrow \quad \downarrow$$

sort R + sort S + Merge R & S

$$(1000 \times 2) + 1000 = 11000$$

$$(1000 \times 2) \times 2 = 4000$$



Database Management Systems III

Semester 2, 2016

1. Consider two relations R and S with the following information about them:
Relation R consists of 2,000 pages with 10 tuples per page, and relation S consists of 500 pages with 50 tuples per page. If S has a B+ tree index on the join attribute, what is the I/O cost of performing an index nested loops join? Explain your answer.
2. Consider the join of two relations R and S on attributes R.a and S.b. R contains 10,000 tuples and has 10 tuples per page. S contains 2,000 tuples also with 10 tuples per page. There are 52 buffer pages available. Estimate the minimum number of page I/Os for performing a Block Nested Loop join of R and S. Ignore the I/O cost of writing out the result. Explain your steps clearly.
3. Consider the join of two relations R and S on attributes R.a and S.b. R contains 20,000 tuples and has 10 tuples per page. S contains 1,000 tuples also with 10 tuples per page. Assume that there are clustered B+ tree indexes on R.a and S.b. Estimate the minimum number of page I/Os for performing a Sort-Merge join of R and S. Ignore the I/O cost of writing out the result. Explain your steps clearly.
4. The relation Executives has attributes ename, title, dname, and address; all are string fields of the same length. The ename attribute is a candidate key. The relation contains 10,000 pages. There are 10 buffer pages. Consider the query:

```
SELECT E.Ename  
FROM Executives E  
WHERE E.title = 'CFO' and E.dname = 'Toy';
```

Assume that only 10% of Executives tuples meet the selection conditions. If a clustered B+ tree index on <dname, title> is (the only index) available, what is the cost of the best plan for this query? Clearly describe how you arrived at the answer.

5. The relation Executives has attributes ename, title, dname, and address; all are string fields of the same length. The ename attribute is a candidate key. The relation contains 20,000 pages. There are 20 buffer pages. Consider the following query:

```
SELECT E.title, COUNT(*)  
FROM Executives E  
WHERE E.dname < 'C'  
GROUP BY E.title;
```

Assume that 5% of the tuples satisfy the selection condition. If the only index available is an unclustered B+ tree on <dname, title>, what is the best plan for executing this query? Indicate the alternative plans you have considered, the cost of each alternative, and how you estimated the cost.

6. The relation Executives has attributes ename, title, dname, and address; all are string fields of the same length. The ename attribute is a candidate key. The relation contains 10,000 pages. There are 10 buffer pages. Consider the query:

```
SELECT E.Ename  
FROM Executives E  
WHERE E.title = 'CFO' and E.dname = 'Toy';
```

Assume that only 10% of Executives tuples meet the selection conditions. If a clustered B+ tree index on <fname, title, ename> is (the only index) available, what is the cost of the best plan for this query? Clearly describe how you arrived at the answer.

7. The relation Emp(empno: string, name: string, age: integer, salary: real) contains 4,000 pages. There are 10 buffer pages. The attribute empno takes up 10 bytes, name 20 bytes, age 2 bytes and salary 8 bytes. A clustered B+ tree on <age,salary> is the only index available. In estimating the cost of query plans, ignore the cost of writing the final result. Consider the following query:

```
Select age, avg(salary)  
From emp  
Where salary > 40000  
Group by age;
```

Assume that 10% of the tuples satisfy the selection condition.

- (a) Describe the best plan and show an estimation of its cost.
(b) If this is the most important of all queries on the Emp table, what additional index if any would you implement to speed up this query? Justify your answer.

8. The relation Emp(empno: char(10), name: char(20), age: number(2), salary: real) contains 3,000 pages. There are 10 buffer pages. The attribute empno takes up 5 bytes, name 15 bytes, age 2 bytes and salary 8 bytes. A clustered B+ tree on <age> is the only index available. In estimating the cost of query plans, ignore the cost of writing the final result. For convenience, consider only the space for storing data ignoring control information and pointers.

- (a) Consider the following query:

```
Select empno, name, salary  
From emp  
Where age = 30 and salary > 40000;
```

Assume that 20% of tuples satisfy each of the selection conditions and 10% of tuples satisfy both conditions. What is the cost of the best plan? Describe the plan you have chosen.

- (b) Suppose the query is as follows:

```
Select age, avg(salary)  
From emp  
Group by age;
```

What is the best plan and its cost?

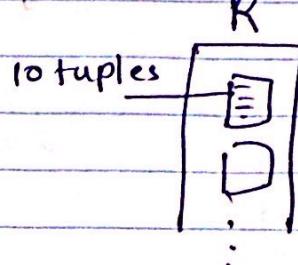
- (c) Suppose the following is an important query on Emp table:

```
Select empno, name, age, salary  
From emp  
Where age > 16 and salary >= 20000  
Order by salary;
```

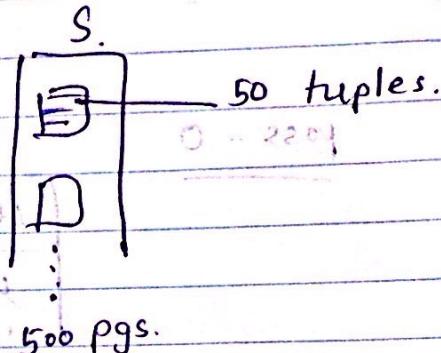
If 60% of tuples satisfy the selection condition, is it useful for this query to have a clustered index on <salary, age>? Explain your answer in terms of page I/Os.

Query Optimization Questions

1. R \bowtie S.



2000 pgs.



500 pgs.

Index Nested Loop Join. If S has B+ tree index.

$$\text{Cost} = \text{Scan R} + (\text{Index Cost}) * (2000 \times 10).$$

$$= 2000 + 3 * 20000$$

$$= 62000 \text{ I/Os} + 10 \text{ B+ tree}$$

R \bowtie S.

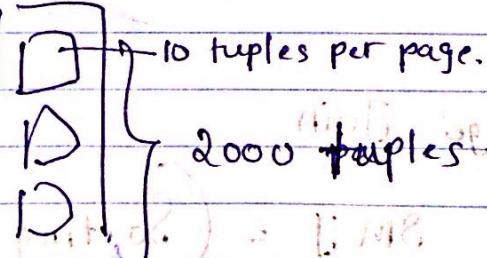
$$R.a = S.b.$$

2. 10 tuples per page.



1000 pgs.

8000 pgs.



200 pages

$$\text{S.NLJ} = (1000 \text{ I/Os} + (200 \times 10000)) + (200 \text{ I/Os} + (1000 \times 200)) =$$

2000 pgs. + 200000 pgs. = 202000 pgs. = 202 buffer pages.

$$\text{S.NLJ} = 1000 + (200 \times 10000)$$

OR

$$= 1000 + 200 \times (-1000 \times 10)$$

$$= 2001000 \text{ I/Os}$$

4

Contd.

PONLJ

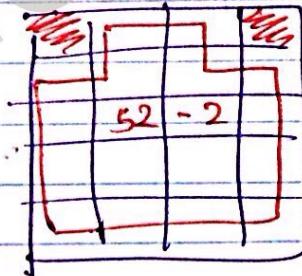
$$\text{Cost} = 1000 + (200 \times 1000)$$

$$= \underline{\underline{201000}} + 1/0$$

BNLJ

$$\text{Cost} = 1000 + 200 \times (\# \text{Blocks in } R)$$

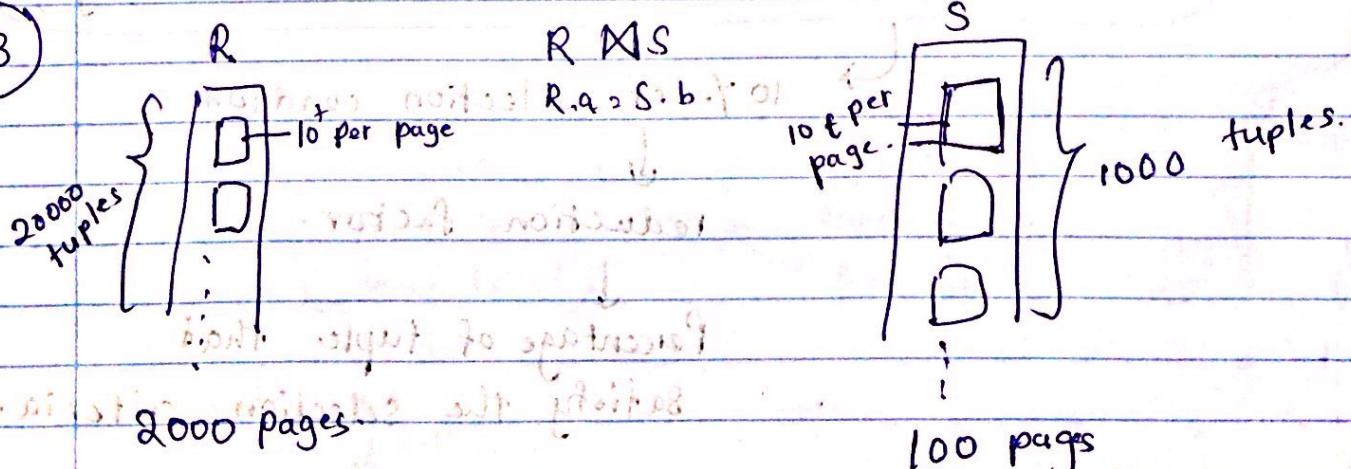
$$\left\lceil \frac{1000}{52-2} \right\rceil = 20$$



$$= 1000 + 200 \times 20$$

$$= \underline{\underline{5000}} + 1/0$$

3



(split, insert) sort n & insert

Clustered B + tree \rightarrow Already sorted.

$$SMJ = (\text{Sorting Cost}) + (\text{Merging Cost}).$$

$$= (\text{Sorting Cost of } R) + (\text{Merging Cost})$$

↓

Sorting Cost of S

$$= O + O + (2000 + 100)$$

\downarrow Clustered
Already Sorted.

\downarrow No. of pages in R

\downarrow No. of pages in S .

$$\approx 2100 \text{ I./O}$$

④ Executives (ename, title, dname, address)

Select e.ename
from Executives e

Where e.title = 'CFO' and e.dname = 'Toy'

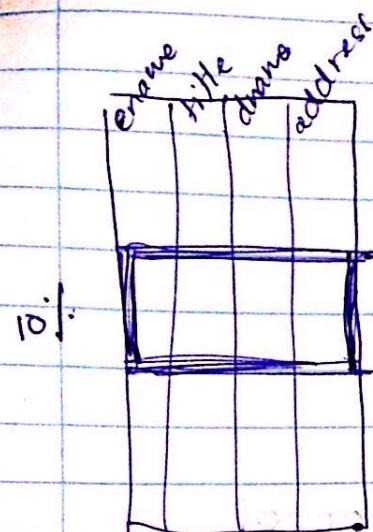
→ 10% of selection condition

reduction factor.

Percentage of tuple that
satisfy the selection criteria.

Clustered B + tree (dname, title)

Time complexity $\rightarrow O(n \log n)$



Executive: 10%.

10,000 pgs.

Cdname, title).



Index.

10 buffer pages.

T.eENAME.

O.e.title = 'CFO'

AND

T.e.dNAME = 'Toy'.

} Index Scan.

Executive \rightarrow

$$\text{Cost} = \left(\begin{array}{l} \text{To travels} \\ \text{root to} \\ \text{leaf level} \end{array} \right) + \left(\begin{array}{l} \text{To travels} \\ \text{through B+} \\ \text{tree leaf level.} \end{array} \right) + \left(\begin{array}{l} \text{table} \\ \times 10 \\ \text{size} \end{array} \right)$$

$(2-4) \downarrow$ \downarrow
 \rightarrow height of
the B+ tree.

index \times 10%
Size.

$$= 3 + \left(10,000 \times \frac{2}{4} \right) \frac{10}{100} \left(10,000 \times \frac{10}{100} \right)$$

$$= 3 + 500 + 1000$$

$$= 1503 \text{ I/O}$$

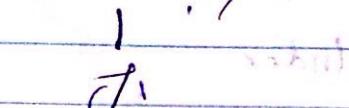
5) select e.title, count(*)

from executives E

where e.dname < 'c'

group by e.title.

Re.title, count



Pages = 20000

buffer pgs = 20

$$\text{cost} = (\text{height of the BT tree}) + \text{index size} \times 5\% + \text{Sorting Cost}$$

$$\text{Index Size} = \text{pages} \times \frac{2}{4} = 20,000 \times \frac{2}{4}$$

$$= 20000 \times \frac{2}{4}$$

$$= 10,000$$

Sorting cost = $2N * \# \text{ of passes}$.

$$\text{Dataset for sorting} = 10,000 \times \frac{5}{100} = 500$$

No. of passes:

$$\text{parse}(0) \rightarrow \left[\frac{500}{25} \right] = 25,$$

$$\text{parse}(1) \rightarrow 25$$

Q.

Suppose, avg page size is 1000 bytes

$\pi_{e.ename}$ 10,000 pgs
| $\sigma_{e.title = 'CFO'}$ AND $e.ename = 'Tay'$ $\sigma_{e.ename, e.title, e.ename}$

Executive

No sorting cost bcz clustered.

cost = Height of B+ tree + index size $\times \frac{10}{100}$ + table size

B+ tree

$$= 2 + \left(10,000 \times \frac{3}{4} \right) \times \frac{10}{100} + 10,000 \times \frac{10}{100}$$

$$= 2 + 7500 \times \frac{10}{100} + 1000$$

$$= 2 + 750 + 1000$$

~~1752 I/O~~

Ans

7. Emp(Empno, name, age, salary)

Empno → 10 bytes.

Name → 20 bytes.

Age → 2 bytes.

Salary → 8 bytes

40

} field length.

clustered B + tree < age, salary >

4000 pgs

10 buffer pgs.

Select age, avg(salary).

from emp

where salary > 40000

group by age; \rightarrow 10%.

9)

π age, avg(salary).

size of age + π avg(salary) \rightarrow 10 buffer pgs

○ salary > 4000

$\frac{1}{10} \times 600$ Emp. \rightarrow 60 pgs $1000 + 600 = 1600$ $\rightarrow 1600/10 = 160$ pgs

Cost = (Scan index) + (sorting cost).

$$\text{Index size} = \left(\frac{1000}{4000} \times \frac{10}{10} \right) +$$

$$\text{data set for sorting} = \left(1000 \times \frac{10}{100} \right) = 100 \text{ pgs}$$

$$\text{Sorting cost} = \left[\frac{100}{9} \right] = 11 \quad 2N \text{ passes}$$

$$\left[\frac{11}{9} \right] = R \quad 2 \times 100 \times 3$$

(+) .

(8) Emp (empno, name, age, salary) \rightarrow 3000 pgs.
 $\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
5 bytes 15 bytes 2 bytes 8 bytes. \rightarrow 10 buff pgs.

Clustered <index=age>
B+ tree index

(9).
Here we cannot use
index only plan. Because
select condition we have some more
like salary. \rightarrow Empno, name, age, salary.

\exists age = 30 AND salary > 40000

Emp

Cost : B+ tree height + Index File Size + Pages Of table.

$$2 \times 3 + \left[3000 \times \frac{2}{30} \times \frac{20}{100} \right] + 3000 \times \frac{20}{100}$$

$$2 \times 3 + 40 + 600$$

$$\underline{2 \times 643 \text{ I/O}}$$

b).

\exists age, avg(sal), - No sorting cost because
here the index is sorted
according to age wise.
So no sorting cost.

- Use index only plan.

$$\text{Cost} = \text{Index File} + \text{Salary search cost}$$

$$2 \times \left(3000 \times \frac{2}{30} \right) + 3000$$

$$2 \times 3200 \text{ I/O}$$

Atlas

c). ~~ppn 2008-4-6 ppn 2008-4-6~~ (original) and
by name, age, salary.

J Satam

Age > 16 AND Salary > 20000

Emp.

Up to 60%.

$$\text{Cost}_2 = \left(\text{memory backed} \right) + \text{Sorting cost.}$$

$$2 \left(3000 \times \frac{10}{30} \times \frac{60}{100} \right) + \text{Sorting}$$

$$2.600 + 0 + 3000 \times \frac{60}{100} =$$

22400