



THE UNIVERSITY OF  
MELBOURNE

# INFO20003 Database Systems

Xiuge Chen

Tutorial 8  
2020.05.11



- 1. Notice - 2min**
- 2. Estimate cost of single-relation plans -  
20min**
- 3. Estimate cost of multi-relation plans -  
35min**
- 4. Lab - 1 hour**

1. Assignment 2 has released - LMS Assessments
2. due date: **6:00pm Friday 15 May**
3. Tips:
  - Follow the submission instruction and format
  - Try SQL practice first - LMS Practice on your own / Lab
  - Might involve some SQL functions not taught - Google
  - Complex queries - break down into sub tasks - nest
  - Always check solutions manually

1. Consider a relation with this schema:

Employees (*eid*: integer, *ename*: string, *sal*: integer, *title*: string, *age*: integer)

Suppose that the following indexes exist:

- a. An unclustered hash index on *eid*
- b. An unclustered B+ tree index on *sal*
- c. An unclustered hash index on *age*
- d. A clustered B+ tree index on (*age*, *sal*)

## 1. Consider a relation with this schema:

Employees relation contains 10,000 pages and each page contains 20 tuples. Suppose there are 500 index pages for B+ tree indexes and 500 index pages for hash indexes. There are 40 distinct values of *age*, ranging from 20 to 60, in the relation. Similarly, *sal* ranges from 0 to 50,000 and there are up to 50,000 distinct values. *eid* is a candidate key; its value ranges from 1 to 200,000 and there are 200,000 distinct values.

For each of the following selection conditions, compute the **Reduction Factor (selectivity)** and the **cost of the cheapest access path** for retrieving all tuples from Employees that satisfy the condition:

1. Consider a relation with this schema:

Employee: 10,000 pages, 20 tuples / page

*eid*: integer, 200,000 distinct values, from 1 to 200,000

*ename*: string,

*sal*: integer, 50,000 distinct values, from 0 to 50,000

*title*: string,

*age*: integer, 40 distinct values, from 20 to 60

- a. An unclustered hash index (500 pages) on *eid*
- b. An unclustered B+ tree index (500 pages) on *sal*
- c. An unclustered hash index (500 pages) on *age*
- d. A clustered B+ tree index (500 pages) on (*age*, *sal*)

What is RF (Reduction Factor, selectivity)?

**Reduction factor (RF)** associated with each predicate reflects the impact of the predicate in reducing the result size

1. Col = value

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

2. Col > value

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

3. Col < value

$$RF = (val - Low(Col)) / (High(Col) - Low(Col))$$

4. Col\_A = Col\_B (for joins)

$$RF = 1 / (Max (NKeys(Col\_A), NKeys(Col\_B)))$$

5. In no information about Nkeys or interval, use a “magic number”

$$RF = 1/10$$



1. Consider a relation with this schema:

Employee: 10,000 pages, 20 tuples / page

*eid*: integer, 200,000 distinct values, from 1 to 200,000

*ename*: string,

*sal*: integer, 50,000 distinct values, from 0 to 50,000

*title*: string,

*age*: integer, 40 distinct values, from 20 to 60

- a. An unclustered hash index (500 pages) on *eid*
- b. An unclustered B+ tree index (500 pages) on *sal*
- c. An unclustered hash index (500 pages) on *age*
- d. A clustered B+ tree index (500 pages) on (*age*, *sal*)



a.  $sal > 20,000 \rightarrow \sigma_{sal > 20000}(R)$

$$\begin{aligned} RF &= (\text{High}(\text{Col}) - \text{value}) / (\text{High}(\text{Col}) - \text{Low}(\text{Col})) \\ &= (50,000 - 20,000) / (50,000 - 0) \\ &= 0.6 \end{aligned}$$

There are **2** possible access paths for this query:

1. The unclustered B+ tree index on *sal*:

$$\begin{aligned} \text{Cost} &:= \text{product of RFs of matching selects} \times (\text{NTuples}(R) + \text{NPages}(I)) \\ &= 0.6 \times ((20 \times 10,000) + 500) \\ &= 120,300 \text{ I/Os} \end{aligned}$$

2. Full table scan

$$\text{Cost: Number of pages} = 10,000 \text{ I/Os}$$

b. Age = 25  $\rightarrow \sigma_{\text{age} = 25} (R)$

$$\text{RF:} = 1/\text{NKeys}(\text{Col}) = 1 / 40$$

There are 3 possible access paths for this query:

1. The clustered B+ tree index on (*age*, *sal*):

$$\begin{aligned}\text{Cost:} &= \text{product of RFs of matching} \times (\text{NPages}(R) + \text{NPages}(I)) \\ &= 1 / 40 * (500 + 10,000) \\ &= 263 \text{ I/Os}\end{aligned}$$

2. The unclustered hash index on *age*:

$$\begin{aligned}\text{Cost:} &= \text{product of RFs of matching} \times \text{hash lookup cost} \times \text{NTuples}(R) \\ &= 1 / 40 * 2.2 * (20 * 10,000) \\ &= 11,000\end{aligned}$$

3. Full table scan

$$\text{Cost: Number of pages} = 10,000 \text{ I/Os}$$

c. Age > 30     $\rightarrow$      $\sigma_{\text{age} > 30} (R)$

$$\begin{aligned}\text{RF} &= (\text{High}(\text{Col}) - \text{value}) / (\text{High}(\text{Col}) - \text{Low}(\text{Col})) \\ &= (60 - 30) / (60 - 20) \\ &= 0.75\end{aligned}$$

There are **2** possible access paths for this query:

1. The clustered B+ tree index on (*age*, *sal*):

$$\begin{aligned}\text{Cost} &= \text{product of RFs of matching} \times (\text{NPages}(R) + \text{NPages}(I)) \\ &= 0.75 * (500 + 10,000) \\ &= 7875 \text{ I/Os}\end{aligned}$$

2. Full table scan

$$\text{Cost: Number of pages} = 10,000 \text{ I/Os}$$



$$d. \text{eid} = 1000 \quad \rightarrow \quad \sigma_{\text{eid} = 1000} (R)$$

$$RF: = 1 / NKeys(Col) = 1 / 200,000$$

There are **2** possible access paths for this query:

1. The unclustered hash index on *eid*:

$$\begin{aligned} \text{Cost:} &= \text{product of RFs of matching} \times \text{hash lookup cost} \times NTuples(R) \\ &= 1 / 200,000 * 2.2 * (20 * 10,000) \\ &= 2.2 \text{ I/Os} \end{aligned}$$

2. Full table scan

$$\text{Cost: Number of pages} = 10,000 \text{ I/Os}$$

**When query a single instance (baed on primary key):**

$$\text{Cost} = \text{hash lookup cost} + 1 \text{ data page access} = 1.2 + 1 = 2.2$$

e.  $sal > 20,000 \wedge age > 30 \rightarrow \sigma_{sal > 20,000 \wedge age > 30}(R)$

$$\begin{aligned} RF &:= RF_{age > 30} * RF_{sal > 20,000} \\ &= 0.75 * 0.6 \end{aligned}$$

There are **3** possible access paths for this query:

1. The clustered B+ tree index on  $(age, sal)$ :

$$\begin{aligned} \text{Cost} &:= \text{product of RFs of matching} \times (\text{NPages}(R) + \text{NPages}(I)) \\ &= 0.75 * 0.6 * (500 + 10,000) \\ &= 4725 \text{ I/Os} \end{aligned}$$

2. Full table scan

$$\text{Cost: Number of pages} = 10,000 \text{ I/Os}$$

3. The unclustered B+ tree index on  $sal$ :

$$\text{Cost: } 0.6 * ((20 \times 10,000) + 500) = 120,300 \text{ I/Os}$$



**Any questions?**

## 2. Estimate cost of multi-relation plans

Consider the following schema:

Emp (eid, sal, age, **did**)

Dept (did, **projid**, budget, status)

Proj (projid, code, report)

The number of tuples in **Emp** is 20,000 and each page can hold 20 records. The **Dept** relation has 5000 tuples and each page contains 40 records. There are 500 distinct *did*s in **Dept**. One page can fit 100 resulting tuples of **Dept JOIN Emp**. Similarly, **Proj** has 1000 tuples and each page can contain 10 tuples. Assuming that *projid* is the candidate key of Proj, there can be 1000 unique values for *projid*. The number of available buffer pages is 50, and Sort-Merge Join can be done in 2 passes. Let's assume that, if we **join Proj with Dept**, 50 resulting tuples will fit on a page.

## 2. Estimate cost of multi-relation plans

Consider the following query:

**SELECT** E.eid, D.did, P.projid

**FROM** Emp **AS** E, Dept **AS** D, Proj **AS** P **WHERE** E.did = D.did

**AND** D.projid = P.projid;

For this query, estimate the cost of the following plans, focusing on the join order and join types:



## 2. Estimate cost of multi-relation plans

Emp: 20,000 tuples, 20 records / page

Dept: 5000 tuples, 40 records / page

did: 500 distinct

Proj: 1000 tuples, 10 records / page

projid: 1000 distinct

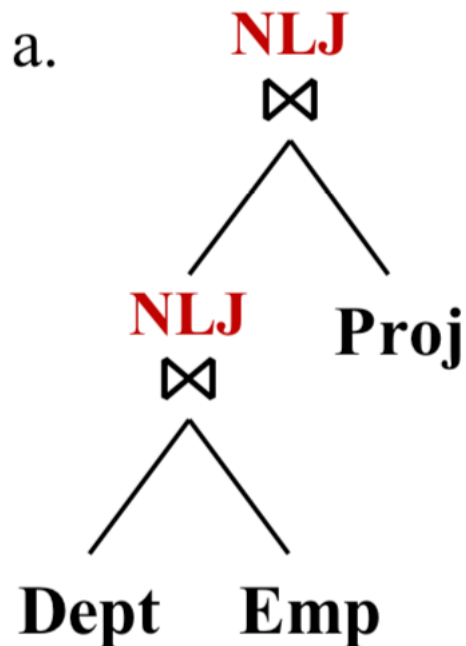
**Dept JOIN Emp:** 100 tuples / page

**Proj JOIN Dept:** 50 tuples / page

Buffer: 50 pages

Sort-Merge Join can be done in 2 passes.

## 2. Estimate cost of multi-relation plans



Total cost of NLJ

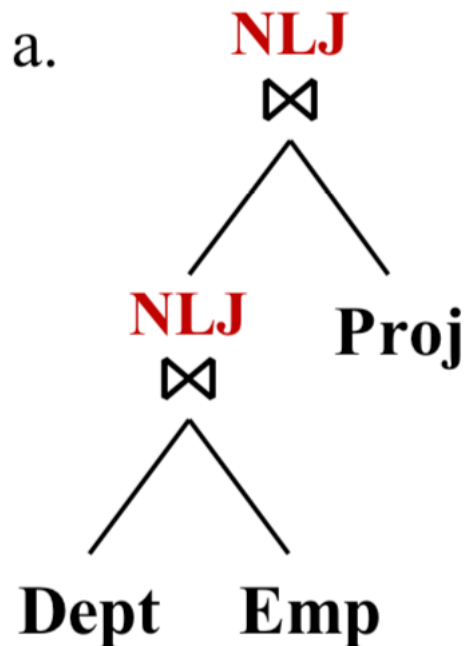
$$= (\# \text{ of pages in outer}) + (\# \text{ of pages in outer} \times \# \text{ of pages in inner})$$

cost of scan whole outer

For each outer page

Cost of scan whole inner

## 2. Estimate cost of multi-relation plans



1. Cost of NLJ Dept and Emp

= Cost of scanning Smaller + Cost of join with larger

= Cost of scanning Dept + Cost to join with Emp

=  $125 + 125 * 1,000$

= 125,125 I/Os

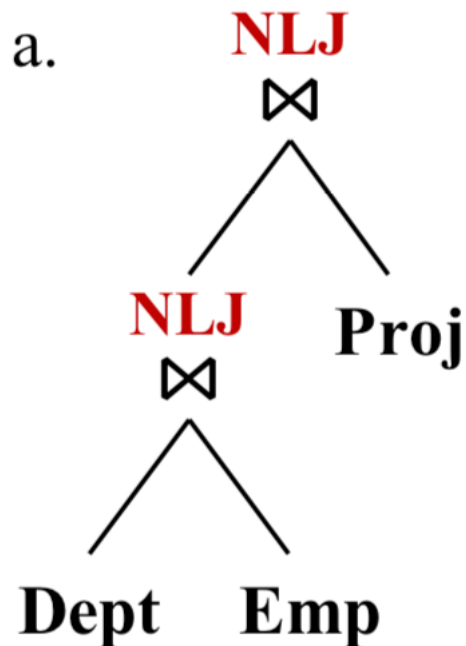
2. Cost of NLJ Proj and Result of (Dept JOIN Emp)

= Cost of join two relation (no need to scan)

= Cost to join Result (D Join E) with Proj

= Number of resulting pages \* Npages(Proj)

## 2. Estimate cost of multi-relation plans



2. Cost of NLJ Proj and Result of (Dept JOIN Emp)

= Cost of join two relation (no need to scan)

= Cost to join Result (D Join E) with Proj

= Number of resulting pages \* Npages(Proj)

= 2,000 \* 100

= 200,000 I/Os

Number of resulting tuples

=  $1 / \text{NKeys(Bigger)} * \text{NTuples(Dept)} * \text{NTuples(Emp)}$

=  $1 / 500 * 5,000 * 20,000$

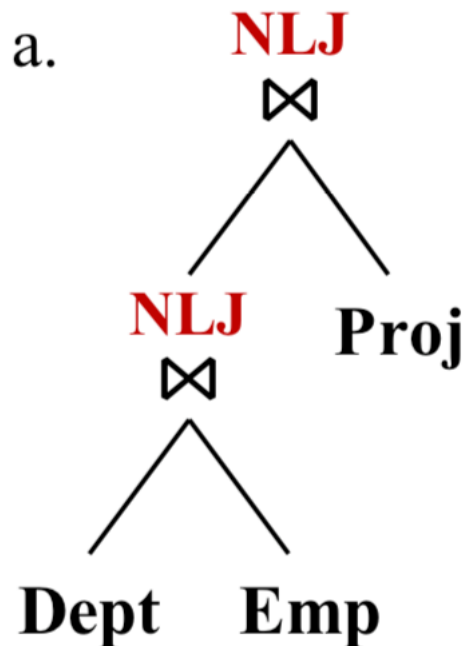
= 200,000 tuples

Number of resulting pages

=  $\text{NTuples}(r) / \text{NTuplesAPage}(r)$

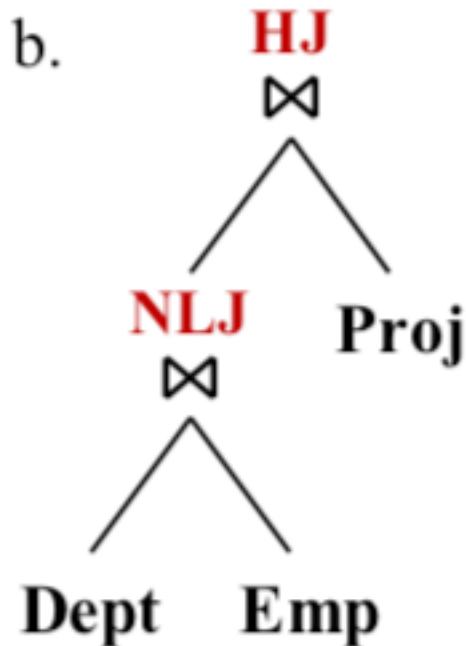
=  $200,000 / 100 = 2000$  pages

## 2. Estimate cost of multi-relation plans



$$\begin{aligned}\text{Total Cost} &= \text{Cost of first NLJ} + \text{Cost of second NLJ} \\ &= 125,125 + 200,000 \\ &= 325,125 \text{ I/Os}\end{aligned}$$

## 2. Estimate cost of multi-relation plans

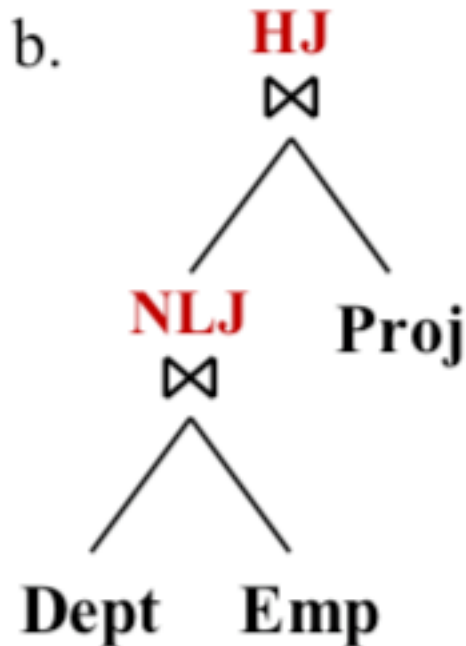


Total cost of HJ

$$= 3(\text{NPages}(I) + \text{NPages}(O))$$

- 1: Cost of scan whole table
- 2: Write hash results into another table
- 3: Scan new hash table for comparing

## 2. Estimate cost of multi-relation plans



1. Cost of NLJ Dept and Emp

= Cost of scanning Smaller+ Cost of join with larger

= Cost of scanning Dept + Cost to join with Emp

=  $125 + 125 * 1,000$

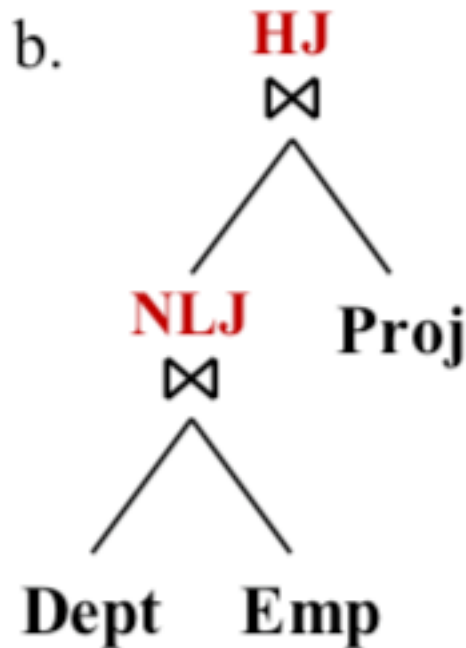
= 125,125 I/Os

2. Cost of HJ Proj and Result of (Dept JOIN Emp)

=  $2 \times \text{NPages}(\text{Dept JOIN Emp}) + 3 \times \text{NPages}(\text{Proj})$

=  $2 * \text{Number of resulting pages} + 3 * \text{Npages}(\text{Proj})$

## 2. Estimate cost of multi-relation plans



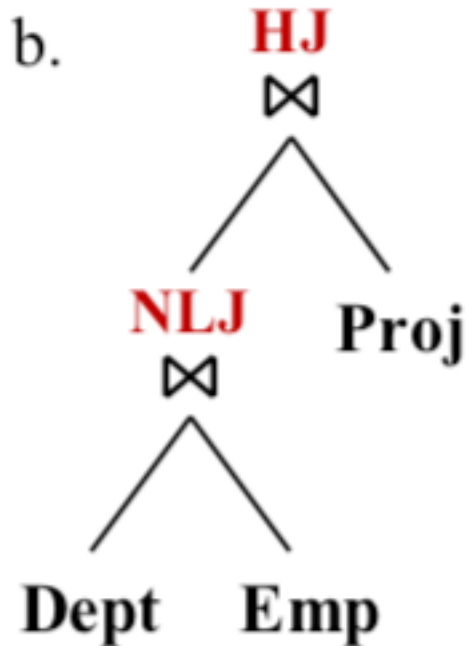
$$\begin{aligned}
 &2. \text{ Cost of NLJ Proj and Result of (Dept JOIN Emp)} \\
 &= 2 \times \text{NPages}(\text{Dept JOIN Emp}) + 3 \times \text{NPages}(\text{Proj}) \\
 &= 2 * \text{Number of resulting pages} + 3 * \text{Npages}(\text{Proj}) \\
 &= 2 * 2,000 + 3 * 100 \\
 &= 4,300 \text{ I/Os}
 \end{aligned}$$

$$\begin{aligned}
 &\text{Number of resulting tuples} \\
 &= 1 / \text{NKeys}(\text{Bigger}) * \text{NTuples}(\text{Dept}) * \text{NTuples}(\text{Emp}) \\
 &= 1 / 500 * 5,000 * 20,000 \\
 &= 200,000 \text{ tuples}
 \end{aligned}$$

$$\begin{aligned}
 &\text{Number of resulting pages} \\
 &= \text{NTuples}(r) / \text{NTuplesAPage}(r) \\
 &= 200,000 / 100 = 2000 \text{ pages}
 \end{aligned}$$

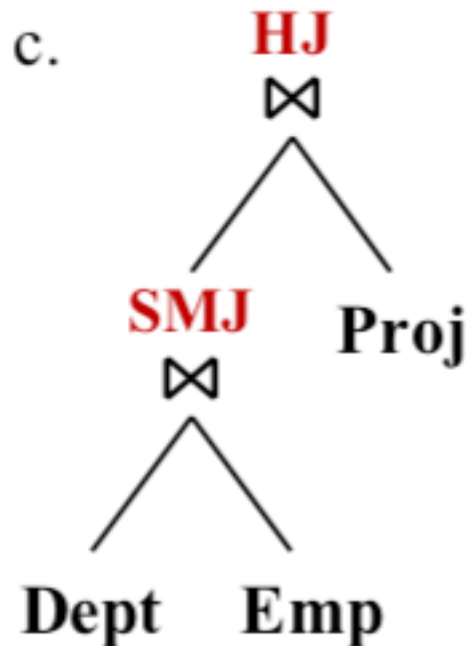


## 2. Estimate cost of multi-relation plans



$$\begin{aligned}\text{Total Cost} &= \text{Cost of first NLJ} + \text{Cost of second HJ} \\ &= 125,125 + 4,300 \\ &= 129,425 \text{ I/Os}\end{aligned}$$

## 2. Estimate cost of multi-relation plans



Total cost of SMJ

= Cost of sorting R + Cost of sorting S +  
Cost of JOIN sorted R and S

Cost of sorting R

= 2 × # of passes × # of pages of R



1: Cost of scan whole table

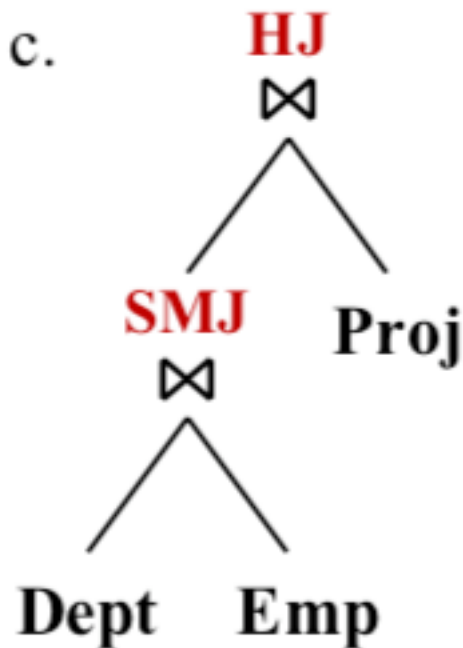
2: Write sorted results into another table

Cost of JOIN sorted R and S =  $NPages(R) + NPages(S)$

Scan new sorted table



## 2. Estimate cost of multi-relation plans



1. Cost of SMJ Dept and Emp

= Cost of sorting Dept + Cost of sorting Emp + Cost of join sorted Dept and Emp

=  $2 \times \text{NPASSES} \times \text{NPAGES}(\text{Dept}) + 2 \times \text{NPASSES} \times \text{NPAGES}(\text{Emp}) + \text{NPAGES}(\text{Dept}) + \text{NPAGES}(\text{Emp})$

=  $2 * 2 * 125 + 2 * 2 * 1,000 + 125 + 1,000$

=  $500 + 4,000 + 1,125$

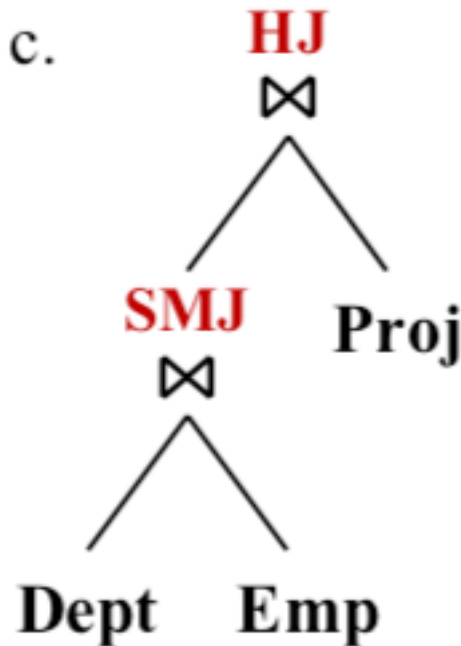
= 5,625 I/Os

2. Cost of HJ Proj and Result of (Dept JOIN Emp)

=  $2 \times \text{NPAGES}(\text{Dept JOIN Emp}) + 3 \times \text{NPAGES}(\text{Proj})$

=  $2 * \text{Number of resulting pages} + 3 * \text{Npages}(\text{Proj})$

## 2. Estimate cost of multi-relation plans



2. Cost of SMJ Proj and Result of (Dept JOIN Emp)

$$\begin{aligned} &= 2 \times \text{NPages}(\text{Dept JOIN Emp}) + 3 \times \text{NPages}(\text{Proj}) \\ &= 2 * \text{Number of resulting pages} + 3 * \text{Npages}(\text{Proj}) \\ &= 2 * 2,000 + 3 * 100 \\ &= 4,300 \text{ I/Os} \end{aligned}$$

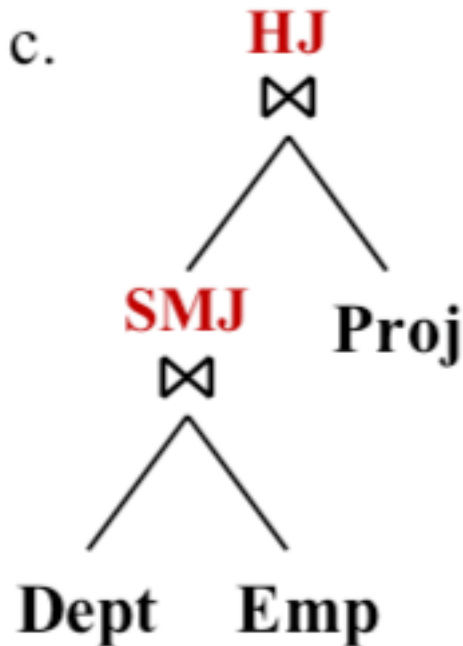
Number of resulting tuples

$$\begin{aligned} &= 1 / \text{NKeys}(\text{Bigger}) * \text{NTuples}(\text{Dept}) * \text{NTuples}(\text{Emp}) \\ &= 1 / 500 * 5,000 * 20,000 \\ &= 200,000 \text{ tuples} \end{aligned}$$

Number of resulting pages

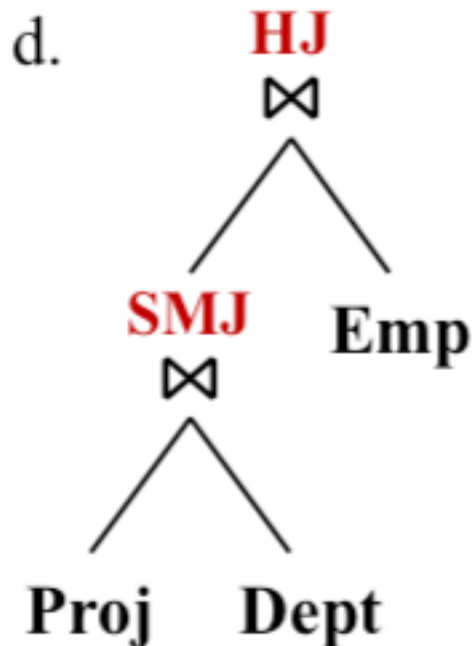
$$\begin{aligned} &= \text{NTuples}(r) / \text{NTuplesAPage}(r) \\ &= 200,000 / 100 = 2000 \text{ pages} \end{aligned}$$

## 2. Estimate cost of multi-relation plans



$$\begin{aligned}\text{Total Cost} &= \text{Cost of first SMJ} + \text{Cost of second HJ} \\ &= 5,625 + 4,300 \\ &= 9,925 \text{ I/Os}\end{aligned}$$

## 2. Estimate cost of multi-relation plans



1. Cost of SMJ Dept and Proj

= Cost of sorting Dept + Cost of sorting Proj + Cost of join sorted Dept and Proj

=  $2 \times \text{NPASSES} \times \text{NPAGES}(\text{Dept}) + 2 \times \text{NPASSES} \times \text{NPAGES}(\text{Proj}) + \text{NPAGES}(\text{Dept}) + \text{NPAGES}(\text{Proj})$

=  $2 * 2 * 125 + 2 * 2 * 100 + 125 + 100$

=  $500 + 400 + 225$

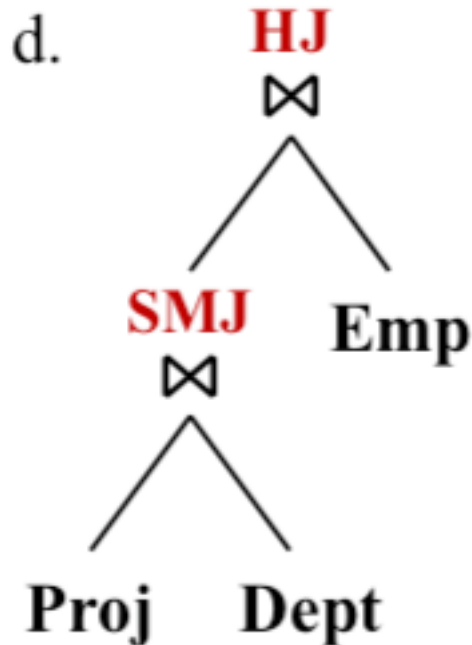
= 1,125 I/Os

2. Cost of HJ Emp and Result of (Dept JOIN Proj)

=  $2 \times \text{NPAGES}(\text{Dept JOIN Proj}) + 3 \times \text{NPAGES}(\text{Emp})$

=  $2 * \text{Number of resulting pages} + 3 * \text{Npages}(\text{Emp})$

## 2. Estimate cost of multi-relation plans



2. Cost of HJ Emp and Result of (Dept JOIN Proj)

$$\begin{aligned}
 &= 2 \times \text{NPages}(\text{Dept JOIN Proj}) + 3 \times \text{NPages}(\text{Emp}) \\
 &= 2 * \text{Number of resulting pages} + 3 * \text{Npages}(\text{Emp}) \\
 &= 2 * 100 + 3 * 1,000 \\
 &= 3,200 \text{ I/Os}
 \end{aligned}$$

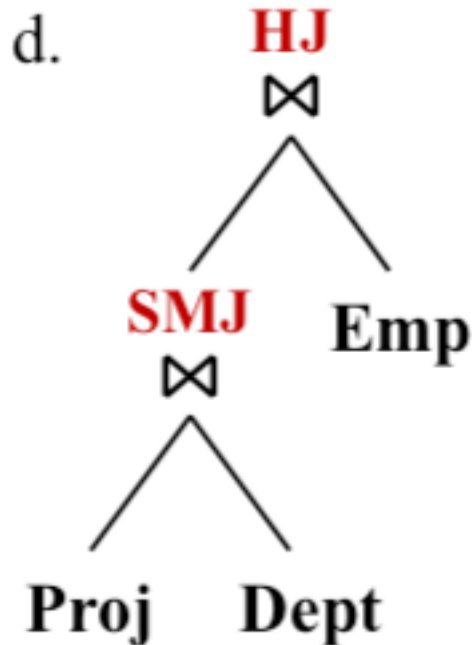
Number of resulting tuples

$$\begin{aligned}
 &= 1 / \text{NKeys}(\text{Bigger}) * \text{NTuples}(\text{Dept}) * \text{NTuples}(\text{Proj}) \\
 &= 1 / 1,000 * 5,000 * 1,000 \\
 &= 5,000 \text{ tuples}
 \end{aligned}$$

Number of resulting pages

$$\begin{aligned}
 &= \text{NTuples}(r) / \text{NTuplesAPage}(r) \\
 &= 5,000 / 50 = 100 \text{ pages}
 \end{aligned}$$

## 2. Estimate cost of multi-relation plans



$$\begin{aligned}\text{Total Cost} &= \text{Cost of first SMJ} + \text{Cost of second HJ} \\ &= 1,125 + 3,200 \\ &= 4,325 \text{ I/Os}\end{aligned}$$





**Any questions?**



## Please refer to Lab 8 on LMS

Let me know if you encounter with  
any problem

**Query Optimisation in MySQL**