

# COMP 3311: Database Management Systems

## Lecture 17 Exercises

### Query Processing: Expression Evaluation

**Exercise 1:** The Student relation consists of 10,000 tuples sorted on student id. Each student has 5 attributes, each 20 bytes, so the tuple size is 100 bytes. The page size is 1,000 bytes, so  $bf_{Student} = \lceil 1000/100 \rceil = 10$ . Therefore,  $B_{Student} = \lceil 10000/10 \rceil = 1000$  pages. Assume that the available main memory  $M$  is 100 pages and that there are 5,000 different student names. There is no index. We want to evaluate the query on the right.

select distinct name  
from Student;

a) Projection using external sorting

Pass 0:

Pass 1:

Total cost:

Number of final result pages written:

b) Projection using hashing (using 20 buckets/partitions)

Partitioning cost:

Duplicate elimination cost:

Total cost:

Number of final result pages written:

**Exercise 2:**

Sailor(sailorId, sName, rating, age)

Reserves(sailorId, boatId, rDate)

For the Sailor relation, each tuple is 50 bytes, a page can hold 80 tuples and there are 500 full pages. For the Reserves relation, each tuple is 40 bytes, a page can hold 100 tuples and there are 1,000 full pages. There are 10 different sailor ratings and 100 different boats. Assume that sailors are distributed uniformly over the 10 ratings and reservations are distributed uniformly over the 100 boats.

Our goal is to process the query:

```
select sName
from Sailor natural join Reserves
where boatId=30
and rating>5;
```

**Some useful statistics**

- On average, each sailor has 2.5 reservations.
- On average, each boat has 1,000 reservations.
- On average, for each rating there are 4,000 sailors.

Estimate the approximate cost of processing the query using a fully pipelined execution method (i.e., do not materialize anything except the final result).

Assume that the Sailor relation contains a clustering B<sup>+</sup>-tree index with 3 levels on sailorId and the Reserves relation contains a non-clustering hash index on boatId. Both the B<sup>+</sup>-tree and hash index can fit 400 index entries per page. For non-clustering indexes, each pointer is assumed to lead to a different page. Use Reserves as the outer relation in the join and take advantage of both indexes.

a) Cost to evaluate  $\sigma_{\text{boatId}=30}$

b) Cost to evaluate  $\sigma_{\text{rating}>5}$

c) Cost to evaluate Sailor  $\bowtie$  Reserves

d) Cost to evaluate  $\pi_{\text{sName}}$

Total cost:

Name: (1) \_\_\_\_\_ / \_\_\_\_\_ Student#: (1) \_\_\_\_\_ Date: \_\_\_\_\_  
Family/Given (PRINT) Given/First (PRINT)

Name: (2) \_\_\_\_\_ / \_\_\_\_\_ Student#: (2) \_\_\_\_\_  
Family/Given (PRINT) Given/First (PRINT)

**NOTE: You are highly encouraged to do this exercise with a partner.**

## COMP 3311: Database Management Systems

### Lecture 17 Exercises Query Processing: Expression Evaluation

**Exercise 3:** The Sailor relation consists of 40,000 tuples sorted on sailor id. Each sailor has 4 attributes, each 10 bytes, so the tuple size is 40 bytes. The page size is 800 bytes so  $bf_{\text{Sailor}} = \lceil 800/40 \rceil = 20$ . Therefore, the Sailor relation requires  $\lceil 40000/20 \rceil = 2000$  pages. Assume that the available main memory  $M$  is 100 pages and that 5% of sailor names are the same. There is no index. We want to evaluate the query on the right.

select distinct sName  
from Sailor;

a) Projection using external sorting

Pass 0:

Pass 1:

Total cost:

Number of final result pages written:

b) Projection using hashing (using 40 buckets/partitions); no optimization

Partitioning cost:

Duplicate elimination cost:

Total cost:

Number of final result pages written:

**Exercise 4:**                      Student(sld, name, deptId, address)                      EnrollsIn(courseId, sId, semester, grade)

The Student relation contains 10,000 tuples in 1,000 pages and the EnrollsIn relation contains 50,000 tuples in 5,000 pages. There are 25 different departments and 1,000 different courses. All attributes have the same length. Each available index is a B<sup>+</sup>-tree with 3 levels. For non-clustering indexes, each pointer is assumed to lead to a different page.

Our goal is to process the query:

```
select name
from Student natural join EnrollsIn
where courseId='COMP3311'
and deptId='COMP';
```

**Some useful statistics**

- On average, a student enrolls in 5 courses.
- On average, a department has 400 students.
- On average, each course has an enrollment of 50 students.

Estimate the approximate cost of processing the query using a fully pipelined execution method (i.e., do not materialize anything except the final result).

Assume that the Student relation contains a clustering index on deptId and the EnrollsIn relation contains a non-clustering index on sId. Use Student as the outer relation in the join and take advantage of both indexes.

a) Cost to evaluate  $\sigma_{\text{courseId}='COMP3311'}$

b) Cost to evaluate  $\sigma_{\text{deptId}='COMP'}$

c) Cost to evaluate Student  $\bowtie$  EnrollsIn

d) Cost to evaluate  $\pi_{\text{name}}$

Total cost: