

CMPT606 – Advanced Database
Fall 2019 – Homework 2

Due Thursday 21/11/2019 midnight – Submit your softcopy to blackboard.

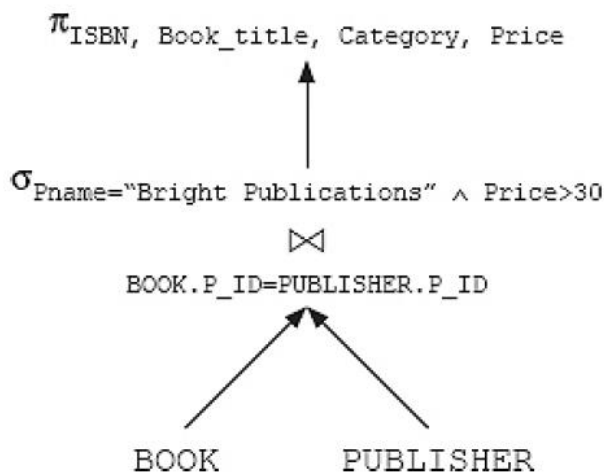
1. [20 pts] Consider the following SQL queries on Online Book Store Application.
 - a. [3 pts] Transform these SQL queries into relational algebra expressions.
 - b. [3 pts] Draw the initial query trees for each of these expressions.
 - c. [4 pts] Derive their optimized query trees after applying heuristics on them. You can just show the last optimal query tree but briefly explain the heuristics applied to arrive to the optimal solution.

Solution

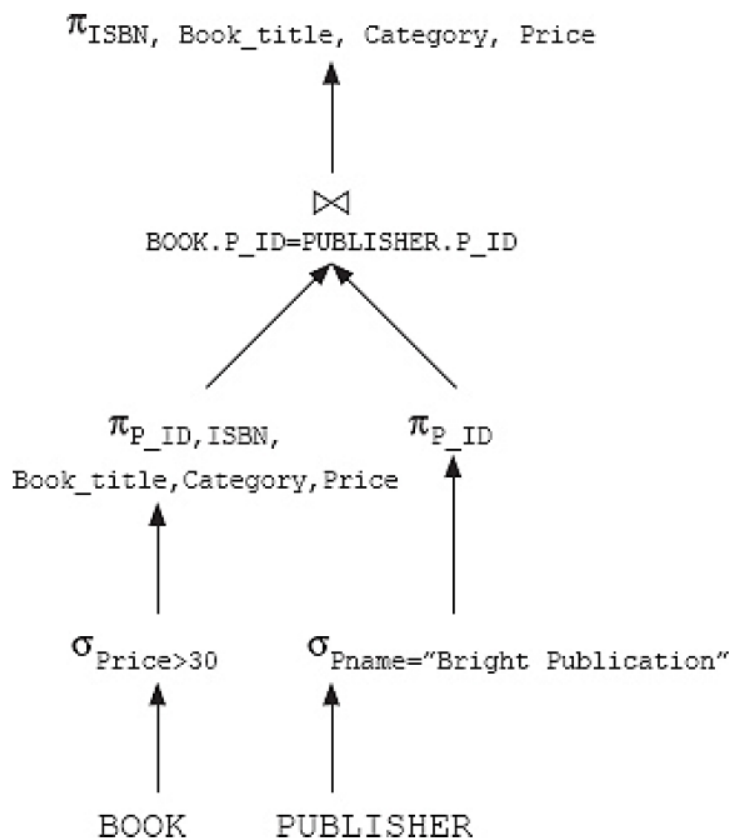
1.1) SELECT ISBN, Book_title, Category, Price
FROM BOOK B, PUBLISHER P
WHERE P.P_ID = B.P_ID AND Pname='Bright Publications'
AND Price>30;

Ans: $\pi_{\text{ISBN, Book_title, Category, Price}} (\sigma_{\text{Pname="Bright Publications"} \wedge \text{Price}>30} (\text{BOOK} \bowtie_{\text{B.P_ID=P.P_ID}} \text{PUBLISHER}))$

The initial query tree of this expression is given below:



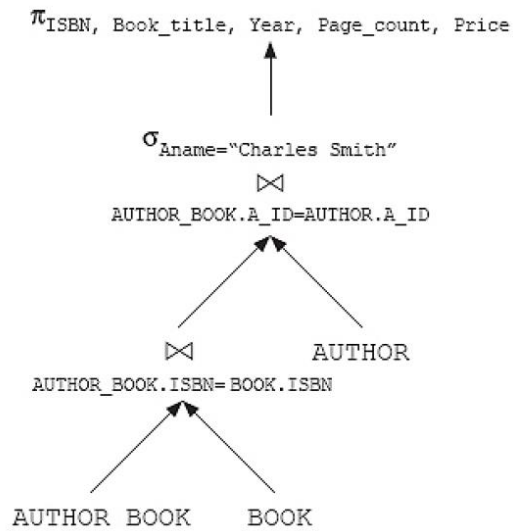
The optimized query tree is given below:



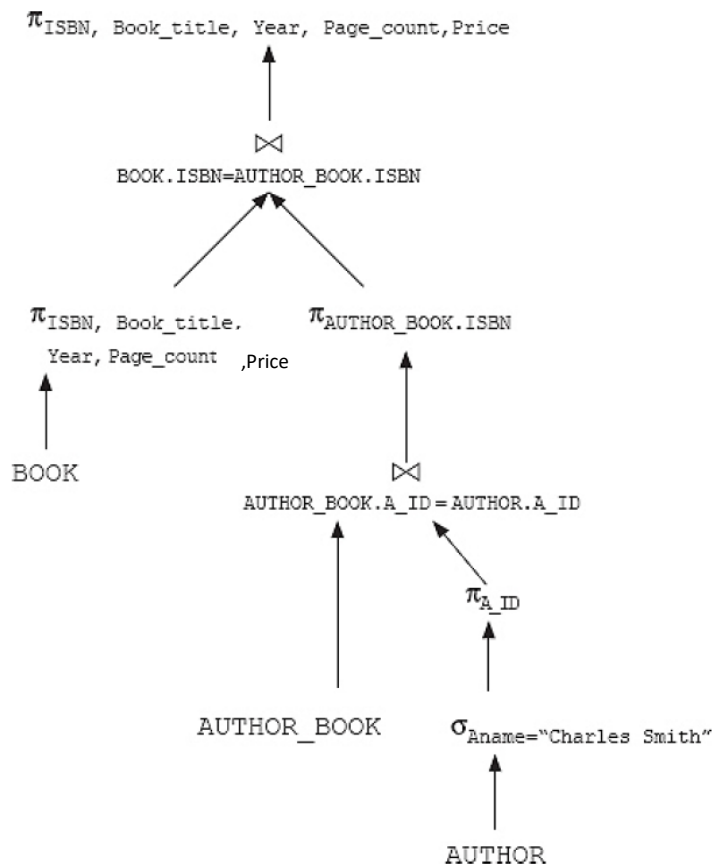
1.2) **SELECT** ISBN, Book_title, Year, Page_count, Price
FROM BOOK B, AUTHOR A, AUTHOR_BOOK AB
WHERE B.ISBN = AB.ISBN **AND** AB.A_ID = A.A_ID
AND Aname = 'Charles Smith';

Ans: $\pi_{ISBN, Book_title, Year, Page_count, Price} (\sigma_{Aname = 'Charles Smith'} ((\text{BOOK} \bowtie_{B.ISBN = AB.ISBN} \text{AUTHOR_BOOK}) \bowtie_{AB.A_ID = A.A_ID} \text{AUTHOR}))$

The initial query tree of this expression is given below:



The optimized query tree is given below:



2. [16 pts] Sort the following values using external merge-sort assuming that each data block stores 2 values and that 3 memory blocks are available:
 44, 6, 55, 4, 3, 7, 11, 2, 16, 8, 21, 9, 12, 5, 7, 22, 33, 6

Sort phase:

$$\text{Number of initial runs} = n_r = \left\lceil \frac{B_r}{M} \right\rceil = \left\lceil \frac{9}{3} \right\rceil = 3$$

Initial sorted runs:

3, 4	6, 7	44, 55
------	------	--------

2, 8	9, 11	16, 21
------	-------	--------

5, 6	7, 12	22, 33
------	-------	--------

Merge phase: Number of passes = $\lceil \log_{M-1}(n_r) \rceil = \lceil \log_2 3 \rceil = 2$

First pass

3, 4	6, 7	44, 55
------	------	--------

2, 8	9, 11	16, 21
------	-------	--------

5, 6	7, 12	22, 33
------	-------	--------

2,3	4, 6	7, 8	9, 11	16, 21	44, 55
-----	------	------	-------	--------	--------

5, 6	7, 12	22, 33
------	-------	--------

Second pass

2,3	4, 6	7, 8	9, 11	16, 21	44, 55
-----	------	------	-------	--------	--------

5, 6	7, 12	22, 33
------	-------	--------

2,3	4, 5	6, 6	7, 7	8, 9	11, 12	16, 21	22, 33	44, 55
-----	------	------	------	------	--------	--------	--------	--------

3. [12 pts] Consider a file with 10,000 blocks and 3 available buffer blocks. Assume that the external sort-merge algorithm is used to sort the file.
 - a. [4 pts] Calculate the initial number of runs produced in the first pass.
 - b. [4 pts] How many passes will it take to sort the file completely?
 - c. [4 pts] How many buffer blocks are required to sort the file completely in two passes?

Solution

a) In the first pass (Pass 0), the number of initial runs (n_r) can be computed as $n_r = \lceil B/M \rceil$, that is, $10000/3 = 3334$ initial sorted runs each containing 3 pages except the last that will have only one page.

b) The number of passes required to sort the file completely, including the initial sorting pass, can be calculated as $\lceil \log_{M-1}(n_r) \rceil + 1$. Here, $M-1$ is 2 and n_r is 3334. Therefore, $\lceil \log_2(3334) \rceil + 1 = 12+1=13$ passes are required.

c) In Pass 0, $\lceil B/M \rceil$ runs are produced. In Pass 1, we must be able to merge these runs, that is, $M-1 \geq \lceil B/M \rceil$. This implies that M must at least be large enough to satisfy $M * (M-1) \geq B$. Thus, with $B = 10000$, $M = 101$.

Thus, we need 101 buffer pages.

4. [15 pts] Query Cost Estimation

Consider two relations, $R(A,B,C)$ and $S(B,D,E)$. R has 1140 tuples. Each block can hold 15 tuples per block. Assume that S is the smaller relation and R is the larger relation.

- [5 pts] a. Assume we have a memory buffer that can hold 25 blocks ($M=25$), and the cost of joining R and S using a block nested-loop join is 228. How many blocks are used to store the tuples in S ?
- [5 pts] b. What is the cost of joining R and S using a simple sort-merge join, using the $B(S)$ you found in question 1?
- [5 pts] c. What is the cost of joining R and S using a hash-based join, using the $B(S)$ you found in question 1?

Solution

a) $B(R) = 1140 / 15 = 76$

The formula is: $B(S) + (B(S).B(R))/(M-2) = C$ where S is the smaller relation and R is the larger relation. Any fractional values should be rounded **UP**.

$$\begin{aligned} B(S) + (B(S).B(R))/(M-2) &= C \\ B(S) + B(S)*76 / 23 &= 228 \\ B(S)(1 + 76/23) &= 228 \\ B(S) = 52.9 &\Rightarrow \text{rounded up} = \mathbf{53} \end{aligned}$$

b)

$$\underbrace{2b_R(\lceil \log_{M-1}(b_R/M) \rceil + 1) + 2b_S(\lceil \log_{M-1}(b_S/M) \rceil + 1)}_{\text{Sort Cost}} + \underbrace{(b_R + b_S)}_{\text{Merge Cost}} = \mathbf{645}$$

c)

$$3 * (B(R) + B(S)) = 3 * (76 + 53) = \mathbf{387}$$

5. [15 pts] Query Cost Estimation

Consider the following two relations:

Customer(cid, name, phone, etc.)

Order(oid, orderDate, cid, etc.)

Customer (C) has 10,000 tuples, with 25 tuples fitting on a block.

Order (O) has 5,000 tuples, with 50 tuples fitting on a block.

There is no index on any attribute of the relations. You may assume the table blocks are contiguous. In your answer always provide the general formula.

(a) Suppose the memory buffer has 101 blocks, compute the cost of using a block-nested loop join to join the above two relations.

(b) Suppose we wanted to join the two relations using a block-nested loop join and limit the cost to 900. What is the smallest value M can be?

(c) What is the cost of joining Customer and Order using a hash-based join?

Solution

(a)

$$B(C) = 10,000 / 25 = 400$$

$$B(O) = 5,000 / 50 = 100$$

$$\text{Cost} = \lceil B(O) + B(C) \times (B(O) / (M-2)) \rceil = \lceil 100 + 400 \times (100 / (101-2)) \rceil = 505$$

(b)

$$B(C) = 10,000 / 25 = 400$$

$$B(O) = 5,000 / 50 = 100$$

$$B(O) + B(C) \times (B(O) / (M-2)) \leq 900$$

$$100 + 400 \times (100 / (M-2)) \leq 900$$

$$100/(M-2) \leq 2$$

$$50 \leq M - 2$$

$$52 \leq M$$

So M must be at least 52 blocks.

$$(c) 3B(C) + 3B(O) = 1,500$$

6. [22 pts] Join Algorithms

Consider two relations R and S. The tuples in each relation are listed in the following table.

R	S
7	8
2	4
9	2
8	1
3	3
9	2
1	7
3	3
6	

We want to do the natural join of R and S based on different join algorithms. For each algorithm listed as following, give the join results in the order that they would be output by the corresponding join algorithm.

(a) [7 pts] The one tuple at a time nested-loop join algorithm. Suppose R is used for the outer loop and S is used for the inner loop.

(b) [7 pts] The merge-sort join algorithm.

(c) [8 pts] The hash join algorithm. We assume only two hash buckets exist, numbered 0 and 1, respectively. The hash function hashes even values to bucket 0 and odd values to bucket 1. Moreover, we assume that in the join phase of the hash join algorithm, **R** is used as the “load” relation and **S** is used as the “stream” relation (i.e., First load R bucket 0, and for each entry scan through S bucket 0 to find matches. Then do the same for the bucket 1). Furthermore, we assume the content of a bucket are read in the same order as they were written.

Solution

(a) For every tuple in R (in order), find matches in S (in order).

Thus we get 7; 2; 2; 8; 3; 3; 1; 3; 3.

(b) First we sort R and S. The result of the join come out in sorted order on the join attributes.

Thus we get 1; 2; 2; 3; 3; 3; 3; 7; 8.

(c) After first phase:

R bucket 0: 2, 8, 6

R bucket 1: 7, 9, 3, 9, 1, 3

S bucket 0: 8, 4, 2, 2

S bucket 1: 1, 3, 7, 3

First load R bucket 0 and then scan through S bucket 0 to find matches: 2, 2, 8

Then load R bucket 1 and then scan through S bucket 1 to find matches: 7, 3, 3, 1, 3, 3

Final result:

2, 2, 8, 7, 3, 3, 1, 3, 3