

# CMPT 354 –Database Systems I (Section D100)

---

## Assignment #4

Instructor: Richard Frank (rfrank@sfu.ca)

TA: Ankit Gupta (aga53@sfu.ca)

Total Marks: 50 (5% of the Individual Assignments)

Due Date: Oct 7, 14:30

For all questions, use the AdventureWorksLT database.

**Question 1** [12 marks] Find the names of customers who have ordered all products in category ‘Helmets’. First, write the relational algebra using the basic operations (selection, projection, Cartesian product, set-difference, union), then write the SQL code to perform the division. Hint, both are covered in the slides.

**Step 1** Project *Product* onto *ProductID*

$\pi_{ProductID}(Product)$

**Step 2** Perform Cartesian product with *ProductCategory*

$\pi_{ProductID}(Product) \times ProductCategory$

**Step 3** Subtract *Product*

$\pi_{ProductID}(Product) \times ProductCategory - Product$

**Step 4** Project onto the *ProductID* of *Product*

$\pi_{ProductID} (\pi_{ProductID}(Product) \times ProductCategory - Product)$

**Step 5** Subtract from *Product*, projected onto *ProductID*

$\pi_{ProductID} (Product) - \pi_{ProductID} (\pi_{ProductID}(Product) \times ProductCategory - Product)$

```
SELECT FirstName, LastName
FROM SalesLT.Customer C
WHERE NOT EXISTS
  (SELECT *
   FROM SalesLT.Product P, SalesLT.ProductCategory Cat
   WHERE P.ProductCategoryID = Cat.ProductCategoryID
   AND Cat.Name = 'Helmets'
   AND NOT EXISTS
     (SELECT *
      FROM SalesLT.SalesOrderHeader H, SalesLT.SalesOrderDetail D
      WHERE H.CustomerID = C.CustomerID
      AND H.SalesOrderID = D.SalesOrderID
      AND D.ProductID = P.ProductID));
```

**Question 2)** For each relational algebra expression below, answer the following questions:

- i) What is the meaning of the query in English? [1 mark each]
- ii) Can it be modified to be more efficient? If so, what is the most optimal solution? Rationalize each step in terms of the number of tuples at each step, number of cells, etc. Show work/logic! It is not necessary to go into the database and use exact numbers, just assume/estimate. [3 marks for algebra manipulation, 3 for rationalization]

For the below calculations/estimations, I'm using the following assumptions:

- All relations contain 10 tuples
- All joins of relation  $n$  and  $m$  result in the creation of  $nm/2$  tuples.
- All conditions/selections on relation  $n$  halve the number of tuples to  $n/2$ .

a)  $\pi_{Product.Name,Color}[\sigma_{weight < 10}(ProductCategory \bowtie Product)]$

i) Display name and colour of all products.

ii)  $\pi_{Name,Color}[\sigma_{weight < 10}(Product)]$

ProductCategory is not needed. All conditions and features exist in Product. Instead of creating 50 tuples, we need to analyze only 10.

b)  $\sigma_{Address1.City='Burnaby' \wedge SubTotal > 1000}(Customer \bowtie SalesOrderHeader$

$\bowtie_{SalesOrderHeader.ShipToAddressID=AddressID} (\rho_{Address1} Address)$

$\bowtie_{SalesOrderHeader.BillToAddressID=AddressID} Address)$

i) All details of Customers and their orders for those customers who had an order, with subtotal more than 1000 shipped to Burnaby.

ii)  $(Customer \bowtie \sigma_{SubTotal > 1000} SalesOrderHeader$

$\bowtie_{SalesOrderHeader.ShipToAddressID=AddressID} \sigma_{Address1.City='Burnaby'}(\rho_{Address1} Address)$

$\bowtie_{SalesOrderHeader.BillToAddressID=AddressID} Address)$

### Original Query:

Assume the join between Customer and SalesOrderHeader creates originally 50 tuples. Joining further to Address (renamed Address1) will result in  $50 * 10 / 5 = 250$  tuples. Finally, joining to Address will result in  $250 * 10 / 2 = 1250$  tuples. Selection is applied via  $City = 'Burnaby'$ , decreasing our dataset to 625 tuples. Then the selection  $SubTotal > 1000$  is applied, giving a final result of 312 tuples.

$\sigma_{Address1.City='Burnaby' \wedge SubTotal > 1000}(10 \bowtie 10$

$\bowtie_{SalesOrderHeader.ShipToAddressID=AddressID} (10) \bowtie_{SalesOrderHeader.BillToAddressID=AddressID} 10)$

$= \sigma_{Address1.City='Burnaby' \wedge SubTotal > 1000}(1250)$

$= \sigma_{SubTotal > 1000}(625)$

$= 312$  tuples

Note: This notation for calculating the complexity does not exist. I am just using it instead of writing things out. The number indicates the number of tuples.

**Optimal Query:**

$$\begin{aligned}
 & (10 \bowtie \sigma_{SubTotal > 1000}(10) \bowtie \sigma_{Address1.City='Burnaby'}(10) \bowtie 10) \\
 & = (10 \bowtie 5 \bowtie 5 \bowtie 10) \\
 & = (25 \bowtie 5 \bowtie 10) \\
 & = (62.5 \bowtie 10) \\
 & = (312)
 \end{aligned}$$

The maximum number of tuples is 312, whereas in the original is it 1250.

$$c) \pi_{FirstName}(Customer \bowtie CustomerAddress \bowtie \sigma_{City='Burnaby'}Address) \cup$$

$$\pi_{FirstName}(\sigma_{LastName='Smith'}Customer \bowtie CustomerAddress \bowtie \sigma_{City='Vancouver'}Address)$$

- i) The first-name of all customers who live in Burnaby and all customers with last name Smith who live in Vancouver.
- ii)  $\pi_{FirstName}\{(\pi_{FirstName, CustomerID}Customer \bowtie \pi_{CustomerID, AddressID}CustomerAddress \bowtie \sigma_{City='Burnaby'}Address) \cup (\pi_{FirstName, CustomerID}\sigma_{LastName='Smith'}Customer \bowtie \pi_{CustomerID, AddressID}CustomerAddress \bowtie \sigma_{City='Vancouver'}Address)\}$

<< 10,9 >> represents 10 tuples and 9 features, total of 90 cells in the relation.

**Original Query:**

$$\begin{aligned}
 & = \pi_{FirstName}(0(Customer \bowtie CustomerAddress \bowtie \sigma_{Address1.City='Burnaby'}Address) \cup \pi_{FirstName}(\sigma_{LastName='Smith'}Customer \bowtie CustomerAddress \bowtie \sigma_{City='Vancouver'}Address)) \\
 & = \pi_{FirstName}(\langle\langle 10,16 \rangle\rangle \bowtie \langle\langle 10,5 \rangle\rangle \bowtie \sigma_{Address1.City='Burnaby'} \langle\langle 10,9 \rangle\rangle) \cup \pi_{FirstName}(\sigma_{LastName='Smith'} \langle\langle 10,16 \rangle\rangle \bowtie \langle\langle 10,5 \rangle\rangle \bowtie \sigma_{City='Vancouver'} \langle\langle 10,9 \rangle\rangle) \\
 & = \pi_{FirstName}(\langle\langle 10,16 \rangle\rangle \bowtie \langle\langle 10,5 \rangle\rangle \bowtie \langle\langle 5,9 \rangle\rangle) \cup \pi_{FirstName}(\langle\langle 5,16 \rangle\rangle \bowtie \langle\langle 10,5 \rangle\rangle \bowtie \langle\langle 5,9 \rangle\rangle) \\
 & = \pi_{FirstName}(\langle\langle 125,28 \rangle\rangle) \cup \pi_{FirstName}(\langle\langle 62,28 \rangle\rangle) \\
 & = \pi_{FirstName}(\langle\langle 125,28 \rangle\rangle + \langle\langle 62,28 \rangle\rangle) \\
 & = \pi_{FirstName}(\langle\langle 187,28 \rangle\rangle) \\
 & = \langle\langle 187,1 \rangle\rangle \text{ that is, 187 tuples, and 1 feature}
 \end{aligned}$$

Until the final stage, the database is also carrying around 16 features from *Customer*, 3 from *CustomerAddress* (since both *CustomerID* and *AddressID* are part of the JOIN, they are projected out and do not count), and 9 from *Address*. So the largest intermediary relation contains 125\*28=3500 cells

**Optimal Query:**

$$\pi_{FirstName}\{(\pi_{FirstName, CustomerID}Customer \bowtie \pi_{CustomerID, AddressID}CustomerAddress \bowtie \pi_{AddressID}\sigma_{City='Burnaby'}Address) \cup (\pi_{FirstName, CustomerID}\sigma_{LastName='Smith'}Customer \bowtie \pi_{CustomerID, AddressID}CustomerAddress \bowtie \pi_{AddressID}\sigma_{City='Vancouver'}Address)\}$$

$$\begin{aligned}
 &= \pi_{\text{Firstname}} \{ (\pi_{\text{Firstname}, \text{CustomerID}} \ll 10, 16 \gg \triangleright \triangleleft \pi_{\text{CustomerID}, \text{AddressID}} \ll 10, 5 \gg \triangleright \triangleleft \\
 &\pi_{\text{AddressID}} \sigma_{\text{City}='Burnaby'} \ll 10, 9 \gg) \cup (\pi_{\text{Firstname}, \text{CustomerID}} \sigma_{\text{LastName}='Smith'} \ll 10, 16 \gg \triangleright \triangleleft \\
 &\pi_{\text{CustomerID}, \text{AddressID}} \ll 10, 5 \gg \triangleright \triangleleft \pi_{\text{AddressID}} \sigma_{\text{City}='Vancouver'} \ll 10, 9 \gg) \} \\
 &= \pi_{\text{Firstname}} \{ (\ll 10, 2 \gg \triangleright \triangleleft \ll 10, 2 \gg \triangleright \triangleleft \pi_{\text{AddressID}} \ll 5, 9 \gg) \cup (\pi_{\text{Firstname}, \text{CustomerID}} \ll 5, 16 \gg \triangleright \triangleleft \\
 &\pi_{\text{CustomerID}, \text{AddressID}} \ll 10, 5 \gg \triangleright \triangleleft \pi_{\text{AddressID}} \ll 5, 9 \gg) \} \\
 &= \pi_{\text{Firstname}} \{ (\ll 10, 2 \gg \triangleright \triangleleft \ll 10, 2 \gg \triangleright \triangleleft \ll 5, 1 \gg) \cup (\ll 5, 2 \gg \triangleright \triangleleft \ll 10, 2 \gg \triangleright \triangleleft \ll 5, 1 \gg) \} \\
 &= \pi_{\text{Firstname}} \{ (\ll 50, 3 \gg \triangleright \triangleleft \ll 5, 1 \gg) \cup (\ll 25, 3 \gg \triangleright \triangleleft \ll 5, 1 \gg) \} \\
 &= \pi_{\text{Firstname}} \{ (\ll 125, 3 \gg) \cup (\ll 62, 3 \gg) \} \\
 &= \pi_{\text{Firstname}} \{ (\ll 125, 3 \gg + \ll 62, 3 \gg) \} \\
 &= \pi_{\text{Firstname}} \{ (\ll 187, 3 \gg) \} \\
 &= \ll 187, 1 \gg
 \end{aligned}$$

The maximum amount of information we need to keep in memory is 187 tuples and 3 features.

d)  $\rho_{\text{Temp}}(\text{Customer} \triangleright \triangleleft \text{SalesOrderHeader} \triangleright \triangleleft \text{SalesOrderDetail} \triangleright \triangleleft \text{Product})$   
 $\pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}}(\text{Temp})$

- i) The first name, last name of each customer, along with the products and the price at which they ordered it.
- ii)  $\pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\pi_{\text{Firstname}, \text{LastName}, \text{CustomerID}} \text{Customer} \triangleright \triangleleft$   
 $\pi_{\text{SalesOrderID}, \text{CustomerID}} \text{SalesOrderHeader} \triangleright \triangleleft$   
 $\pi_{\text{SalesOrderID}, \text{ProductID}} \text{SalesOrderDetail} \triangleright \triangleleft$   
 $\pi_{\text{ProductID}, \text{Name}, \text{Price}} \text{Product})$

### Original Query

$$\begin{aligned}
 &\pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\text{Customer} \triangleright \triangleleft \text{SalesOrderHeader} \triangleright \triangleleft \text{SalesOrderDetail} \triangleright \triangleleft \\
 &\triangleleft \text{Product}) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 10, 16 \gg \triangleright \triangleleft \ll 10, 22 \gg \triangleright \triangleleft \ll 10, 9 \gg \triangleright \triangleleft \ll 10, 17 \gg) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 50, 37 \gg \triangleright \triangleleft \ll 10, 9 \gg \triangleright \triangleleft \ll 10, 17 \gg) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 250, 46 \gg \triangleright \triangleleft \ll 10, 17 \gg) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 1250, 62 \gg) \quad \leftarrow \text{that is 77500 cells!} \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 1250, 4 \gg)
 \end{aligned}$$

### Optimal Query

$$\begin{aligned}
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\pi_{\text{Firstname}, \text{LastName}, \text{CustomerID}} \ll 10, 16 \gg \triangleright \triangleleft \\
 &\pi_{\text{SalesOrderID}, \text{CustomerID}} \ll 10, 22 \gg \triangleright \triangleleft \\
 &\pi_{\text{SalesOrderID}, \text{ProductID}} \ll 10, 9 \gg \triangleright \triangleleft \\
 &\pi_{\text{ProductID}, \text{Name}, \text{Price}} \ll 10, 17 \gg) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 10, 3 \gg \triangleright \triangleleft \ll 10, 2 \gg \triangleright \triangleleft \ll 10, 2 \gg \triangleright \triangleleft \ll 10, 3 \gg) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 50, 4 \gg \triangleright \triangleleft \ll 10, 2 \gg \triangleright \triangleleft \ll 10, 3 \gg) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 250, 5 \gg \triangleright \triangleleft \ll 10, 3 \gg) \\
 &= \pi_{\text{Firstname}, \text{LastName}, \text{Product.Name}, \text{ListPrice}} (\ll 1250, 7 \gg) \\
 &= (\ll 1250, 4 \gg)
 \end{aligned}$$

**Question 3)** Write the *optimal* relational algebra statement for the following. [5 marks each]

a) Display the product descriptions for Mountain Tires (a product).

$$= \pi_{Description} (\pi_{ProductModelID} \sigma_{Name='Mountain Tires'} Product \triangleright \triangleleft$$

$$\pi_{ProductModelID, ProductDescriptionID} ProductModelProductDescription \triangleright \triangleleft$$

$$\pi_{ProductDescriptionID, Description} ProductDescription$$

b) Assume the ProductCategory has 3 levels. Write the query to display the product names at all levels within the hierarchy.

For example, assume ProductCategory looks like:

ProductCategoryID	ParentProductCategoryID	Name
1		Car
2	1	Sports
3	2	Mustang
4	2	Ferrari

Required result:

Car	Sports	Mustang
Car	Sports	Ferrari

$$= \pi_{P1.Name, P2.Name, P3.Name} (\rho_{P1} ProductCategory$$

$$\triangleright \triangleleft_{P1.ProductCategoryID = P2.ParentProductCategoryID} \rho_{P2} ProductCategory$$

$$\triangleright \triangleleft_{P2.ProductCategoryID = P3.ParentProductCategoryID} \rho_{P3} ProductCategory)$$