

Outline of solutions to selected exercises in Mandatory assignment 1

Exercise 1 Some SQL operations

a)

```
SELECT S.Course_number, S.Semester, S.Year, COUNT(DISTINCT
Student_number)
FROM Section S, Grade_Report G
WHERE Instructor='King' AND S.Section_Identifier = G.Section_Identifier
GROUP BY S.Course_number, S.Semester, S.Year
```

b)

```
SELECT S.Name, S.Major
FROM Student S
WHERE NOT EXISTS (
    SELECT *
    FROM Grade_Report G
    WHERE G.StudentNumber = S.StudentNumber AND NOT
(G.Grade='A'))
```

c)

```
INSERT INTO Student VALUES ('Johnson', 25, 1, 'MATH')
```

d)

```
UPDATE Student SET CLASS = 2 WHERE Name = 'Smith'
```

e)

```
DELETE FROM Student WHERE Name = 'Smith' AND StudentNumber=17
```

Exercise 2 Introduction to DBMS Architecture

2.1 QTE schematics

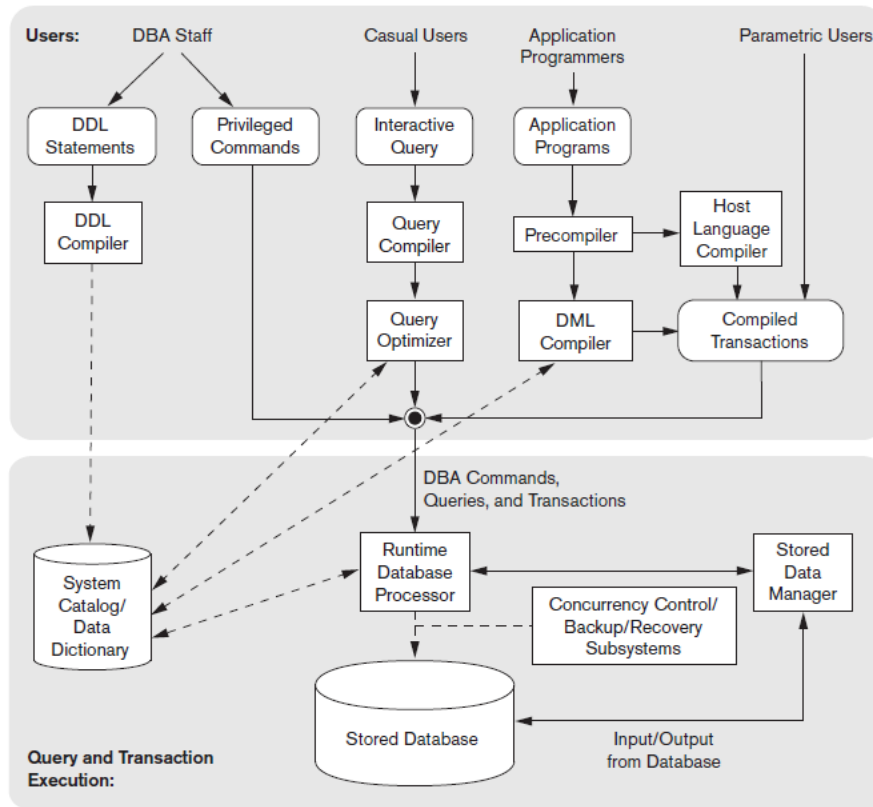


Figure 2.3
Component modules of a DBMS and their interactions.

Elmasri and Navathe. Fundamentals of Database Systems, Seventh edition.

[The task here was to describe the DBMS components in the lower box, following the books and/or the lectures while using your own words]

2.2 Central components and Higher-level components

[Please consult slides 7-11 in the link:

<https://www.uio.no/studier/emner/matnat/ifi/IN3020/v21/lectures/021-2---dbms-architecture-dependencies.pdf>]

Exercise 4 Relational Algebra

4.1

- a) Rewrite the query:

```
SELECT P.C
FROM Purchase P
WHERE Exists
  (SELECT S.A
   FROM Supply S
   WHERE P.A=S.A AND S.D>0)
```

RA expression for the nested subquery:

$$\pi_{P.A,P.B,P.C,S.A,S.D} [\sigma_{P.A=S.A \wedge S.D>0} \{ \rho_P \text{ Purchase} \times \rho_S \text{ Supply} \}]$$

[For Selection in the nested part, you keep all the attributes/columns, so you inadvertently do not remove any columns that might be needed for the final Selections as specified in the outer query]

The first **FROM** part is written as:

$$\rho_P \text{ Purchase}$$

Joining these two, making some optimization on projection, i.e. removing redundant columns, and applying the final selection corresponding to the first line:

$$\pi_{P.C}(\rho_P \text{ Purchase} \bowtie \pi_{P.A,P.B,P.C} [\sigma_{P.A=S.A \wedge S.D>0} \{ \rho_P \text{ Purchase} \times \rho_S \text{ Supply} \}])$$

- b) Optimized query:

$$\pi_{P.C}(\rho_P \text{ Purchase} \bowtie \pi_{P.A,P.B,P.C} [\rho_P \text{ Purchase} \bowtie_{P.A=S.A} \{ \sigma_{S.D>0} \rho_S(\text{Supply}) \}])$$

By pushing selection and by converting Selection and Cartesian-join into a Theta-join

- c) Selection has been applied on individual relation (S) rather than on a Join. Additionally, Cartesian-join is replaced by a Theta join.

4.2

- a) Sequence of steps for the following RA transformation:

$$\pi_M(\sigma_{\langle c \rangle \wedge \langle d \rangle} (P \bowtie Q \bowtie R)) \rightarrow \pi_M[\{ \pi_A(\sigma_{\langle c \rangle} P) \bowtie \pi_{A,H}(\sigma_{\langle d \rangle} Q) \} \bowtie R]$$

- Associativity of join: $\pi_M[\sigma_{\langle c \rangle \wedge \langle d \rangle} (P \bowtie Q \bowtie R)]$
- Pushing selection over join: $\pi_M[\{ \sigma_{\langle c \rangle \wedge \langle d \rangle} (P \bowtie Q) \} \bowtie R]$ $\langle c \rangle, \langle d \rangle$ not on R
- Pushing selection to respective relation: $\pi_M[\{ \sigma_{\langle c \rangle} (P) \bowtie \sigma_{\langle d \rangle} (Q) \} \bowtie R]$

Since $\langle c \rangle$ applies on P only and $\langle d \rangle$ on Q only.

- Distribution of projection over join:

$$\pi_M[\pi_H\{\sigma_{\langle c \rangle}(P) \bowtie \sigma_{\langle d \rangle}(Q)\} \bowtie \pi_{F,H,K,M}(R)]$$

- Since $\pi_{F,H,K,M} R = R$:

$$\pi_M[\pi_H\{\sigma_{\langle c \rangle}(P) \bowtie \sigma_{\langle d \rangle}(Q)\} \bowtie R]$$

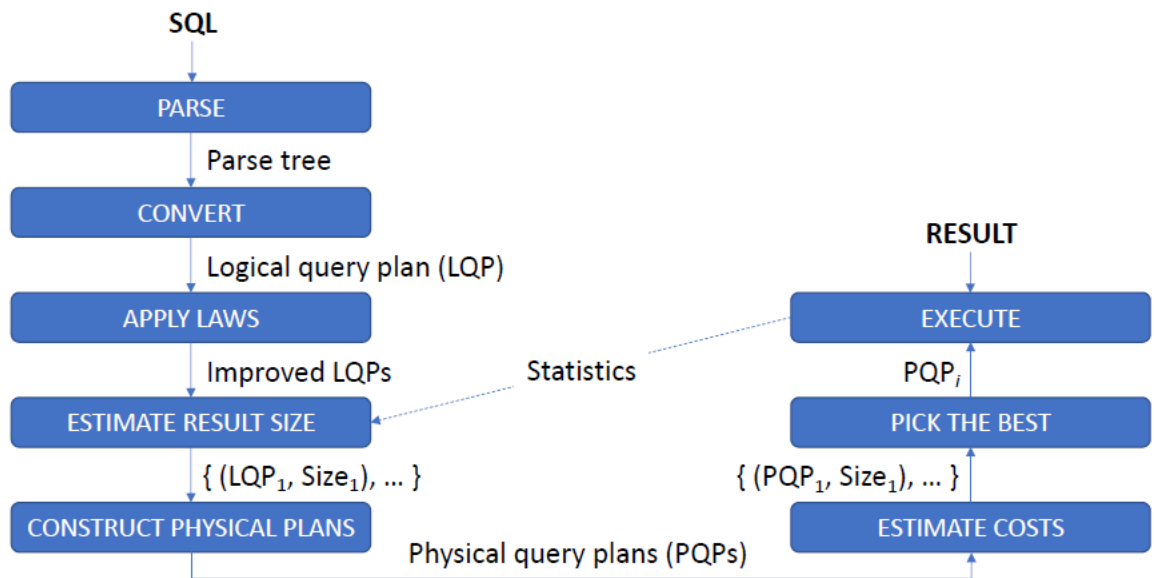
- Distribution of projection over join:

$$\pi_M[\{\pi_A(\sigma_{\langle c \rangle} P) \bowtie \pi_{A,H}(\sigma_{\langle d \rangle} Q)\} \bowtie R]$$

- b) The final expression appears better because Selection has been split and pushed down over the join, and the large join $P \bowtie Q \bowtie R$ at the very beginning has especially been avoided.

Exercise 5 – Query Compilation

5.1 Query processing



[The task was about describing this well-known process]

5.2 SQL, RA tree, Size computations

a) SQL query

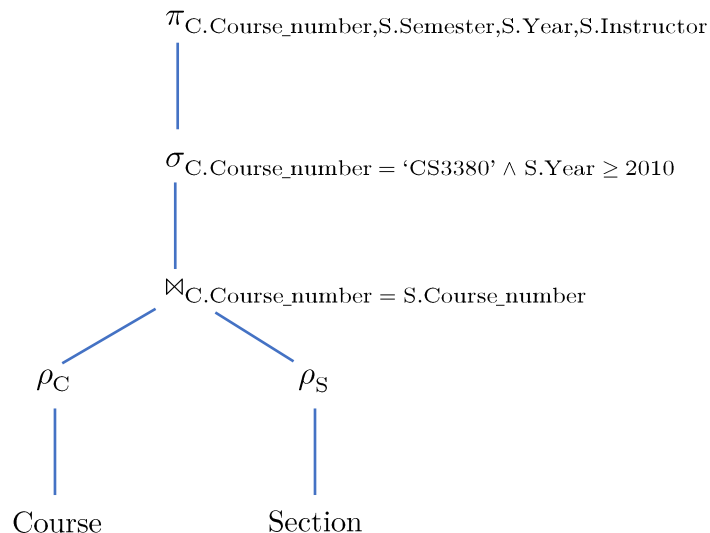
[Since this query involves joining two relations, it can be written in several possible ways, e.g., joining the relations first and filtering tuples/attributes later, or the other way around. About the join type, Cartesian product will be the simpler but costlier than other joins. Here, we start with the theta-join in the non-optimized version]

```
Select C.Course_name, S.Semester, S.Year, S.Instructor
From Course As C Join Section As S
On C.Course_number = S.Course_number
Where C.Course_number = 'CS3380' And S.Year >= 2010
```

b) RA expression

$$\pi_{C.Course_name, S.Semester, S.Year, S.Instructor} \left(\sigma_{C.Course_number = 'CS3380' \wedge S.Year \geq 2010} \left(\rho_C \text{ Course } \bowtie_{C.Course_number = S.Course_number} \rho_S \text{ Section} \right) \right)$$

c) RA tree



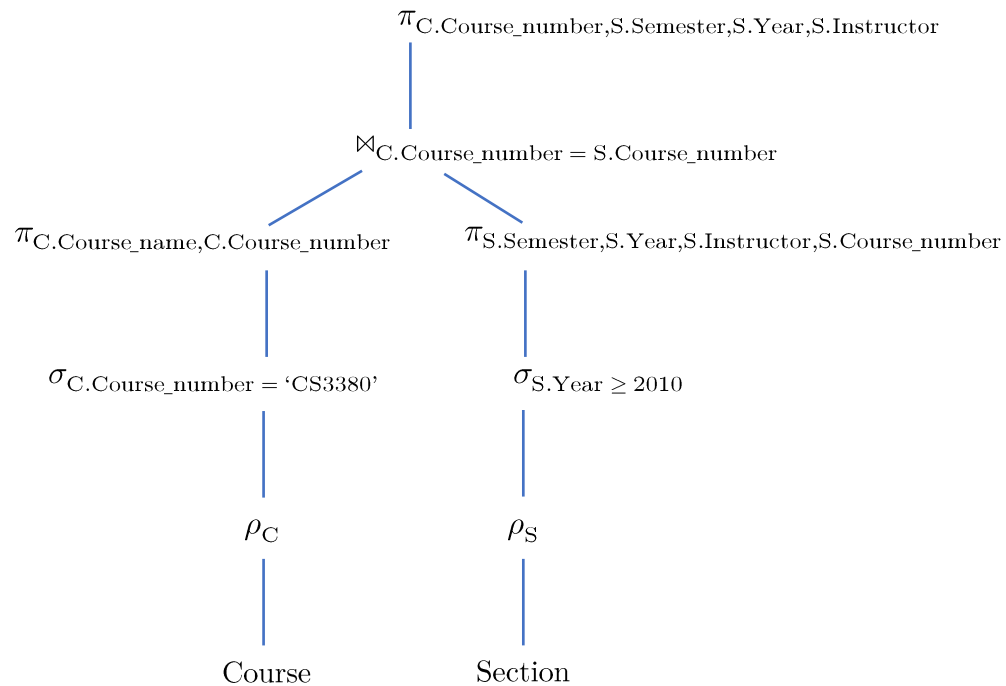
d) Improved RA expression and tree

[Apply selection and projection on the respective relations first, perform join afterward]

$$\pi_{C.Course_name,S.Semester,S.Year,S.Instructor} [\{ \pi_{C.C.Course_name,C.Course_number} (\sigma_{C.Course_number = 'CS3380'} \rho_C Course) \} \bowtie_{C.Course_number = S.Course_number} \{ \pi_{S.Semester,S.Year,S.Instructor,S.Course_number} (\sigma_{S.Year \geq 2010} \rho_S Section) \}]$$

[Projections in the 2nd and 4th lines improve the RA expression further. However, even if we did not perform these projections, this expression is still likely to be better than the expression in b)]

RA tree of the improved expression:



[It might be convenient to replace the attribute names in the expression by suitable symbols with prior declaration, e.g., Course_name = Cn, Course_number = Cm, Year = Yr, etc.]

e) Size calculations

[This part allowed the students to exercise freedom in choosing the tuple numbers, tuple sizes, selectivity factors etc. and demonstrate their understanding of the logical flow in cost/size calculations. We did not make a detailed official solution in this part. Martin Strøm Olsen from IN 4020 has made a thorough calculations here, capturing the key considerations in such computation, which is presented as inspiration for others. Thanks, Martin.]

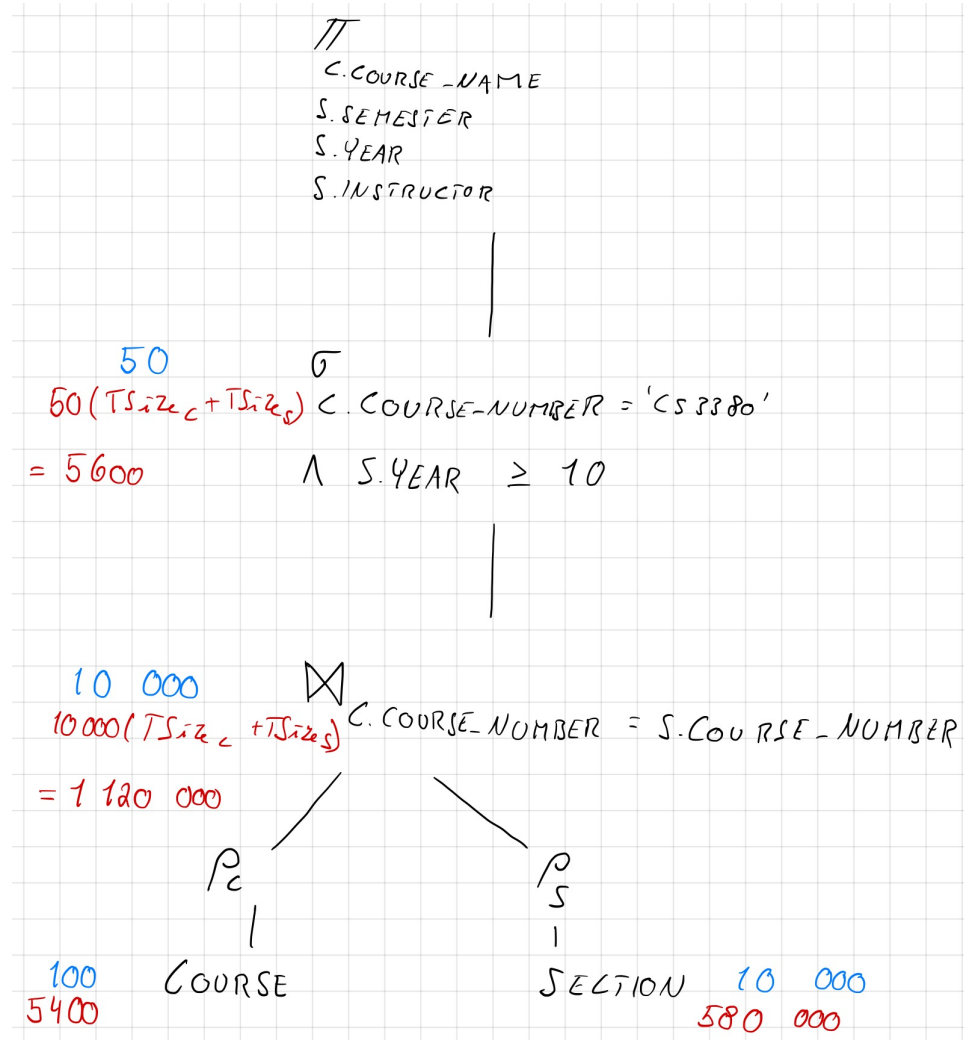
Course (C)	Section (S)
$\text{Tup}(C) = 100$	$\text{Tup}(S) = 10000$
$\text{Val}_C(\text{Course_Number}) = 100$	$\text{Val}_S(\text{Course_Number}) = 100$ $\text{Val}_S(\text{Year}) = 100$
$\text{TSize}_C = (20, 20, 4, 10)$ bytes, with the order of attributes as given in the relation Course.	$\text{TSize}_S = (4, 20, 10, 4, 20)$ bytes, with the order of attributes as given in the relation Section.

It is assumed that Course number is a PK of Course and FK in Section. The number of tuples (blue) and the estimated tuple size (red) for the original and the modified tree is shown below, together with the detailed computations.

Some notations defined:

$$\begin{aligned}
 c_1 &:= C. \text{Course_number} = \text{'CS3380'} \\
 c_2 &:= S. \text{Year} \geq 2010 \\
 c_3 &:= C. \text{Course_number} = S. \text{Course_number}
 \end{aligned}$$

Original tree



Estimates of number of tuples

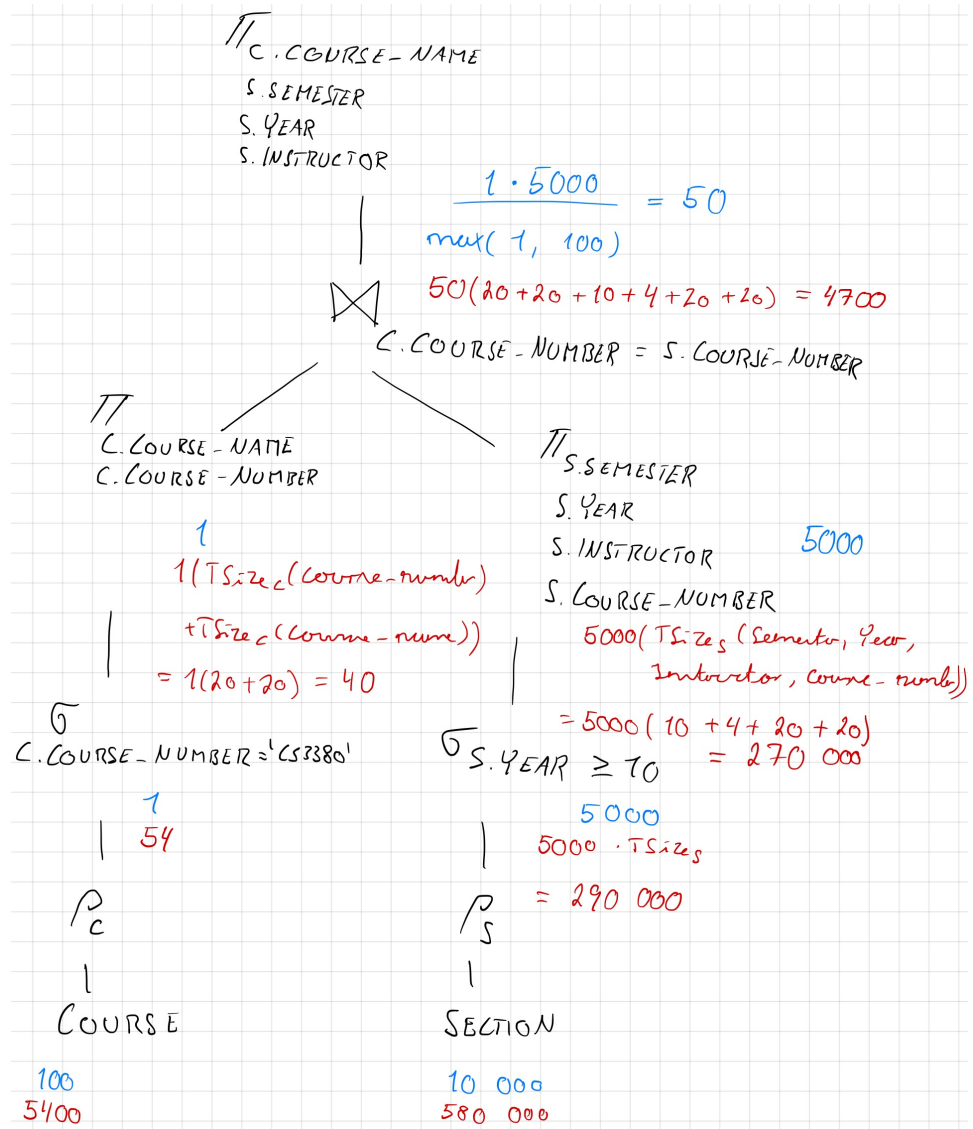
- $C \bowtie_{C.Course_number = S.Course_number} S$
 here Course_number is a PK in C and a FK in S which gives
 $Tup = Tup(S) = 10000$
- $\sigma_{c_1 \wedge c_2}$
 c_1 selectivity factor = $\frac{1}{Val_C(Course_number)} = \frac{1}{100}$
 c_2 selectivity factor: assume Year \sim Uniform(5, 14), so 14 - 5 + 1 = 10 different values and Year ≥ 10 gives (10, 11, 12, 13, 14) 5 possible values. Selectivity factor = $\frac{5}{10}$
 $Tup = Tup(\sigma_{c_1 \wedge c_2}(C \bowtie_{c_3} S)) = Tup(C \bowtie_{c_3} S) \frac{1}{100} \frac{5}{10} = 10000 \cdot \frac{1}{100} \cdot \frac{5}{10} = 50$

Estimating the size of the temporary relations

- $C \bowtie_{C.Course_number = S.Course_number} S$
 $\Sigma_1 = 10000(TSize_C + TSize_S) = 10000 \cdot (54 + 58) = 1120000$
- $\sigma_{c_1 \wedge c_2}$
 $\Sigma_2 = 50(TSize_C + TSize_S) = 50(54 + 58) = 5600$

then the estimated size (excluding root node and leaf nodes)
 $\Sigma = \Sigma_1 + \Sigma_2 = 1120000 + 5600 = 1125600$

Modified tree



Estimates of number of tuples

- σ_{c_1}
 $Tup(\sigma_{c_1}) = \frac{Tup_C}{Val_C(Course_number)} = \frac{100}{100} = 1$

- σ_{c_2}
 $\text{Tup}(\sigma_{c_2}) = \text{Tup}_S \frac{5}{10} = 10000 \cdot \frac{5}{10} = 5000$
- $\pi_{C.Course_name, C.Course_number}$
 $\text{Tup} = \text{Tup}(\sigma_{c_1}) = 1$
- $\pi_{S.Semester, S.Year, S.Instructor, S.Course_number}$
 $\text{Tup} = \text{Tup}(\sigma_{c_2}) = 5000$
- $C \bowtie_{C.Course_number = S.Course_number} S$

$$\text{Tup} = \frac{\text{Tup}_C \cdot \text{Tup}_S}{\max(\text{Val}_C(\text{Course_number}), \text{Val}_S(\text{Course_number}))} = \frac{1 \cdot 5000}{\max(1, 100)} = \frac{5000}{100} = 50$$

Estimating the size of the temporary relations

- σ_{c_1}
 $\Sigma_1 = 1 \cdot 54 = 54$
- σ_{c_2}
 $\Sigma_2 = 5000 \cdot 58 = 290000$
- $\pi_{C.Course_name, C.Course_number}$
 $\Sigma_3 = 1(20 + 20) = 40$
- $\pi_{S.Semester, S.Year, S.Instructor, S.Course_number}$
 $\Sigma_4 = 5000(20 + 10 + 4 + 20) = 270000$
- $C \bowtie_{C.Course_number = S.Course_number} S$
 $\Sigma_5 = 50(20 + 20 + 20 + 10 + 4 + 20) = 4700$

then the estimated size (excluding root node and leaf nodes)

$$\Sigma = \Sigma_1 + \Sigma_2 + \Sigma_3 + \Sigma_4 + \Sigma_5 = 54 + 290000 + 40 + 270000 + 4700 = 564794$$

so the estimated size of the temporary relations (that are not equal in the two trees) is likely smaller in the modified tree, implying that the modified tree is likely to be more efficient than the original.

References

Akkøk, M. N. (2021a). Database Systems Spring 2021, Week 2.1 (part 1) Summary of DBMS Architecture [Lecture slides].

<https://www.uio.no/studier/emner/matnat/ifi/IN3020/v21/lectures/021-1---dbms-architecture.pdf>

Akkøk, M. N. (2021). Database Systems Spring 2021, Week 2.1 (part 2) [Lecture slides].

<https://www.uio.no/studier/emner/matnat/ifi/IN3020/v21/lectures/021-2---dbms-architecture-dependencies.pdf>

Elmasri, R. & Navathe, S. (2016). Fundamentals of database systems (7th ed.). Boston , Mass: Pearson.

Kostylev, E. V. (2021). Database Systems Spring 2021, Week 2.2 Indexing [Lecture slides].

<https://www.uio.no/studier/emner/matnat/ifi/IN3020/v21/lectures/022---indexing.pdf>

Kostylev, E. V. (2021a). Database Systems Spring 2021, Week 4.2-5.2 Query Compilation - Parts 1-3 [Lecture slides].

<https://www.uio.no/studier/emner/matnat/ifi/IN3020/v21/lectures/042-052---qcomp.pdf>

Kostylev, E. V. (2021b). Database Systems Spring 2021, Weeks 6.1-6.2 Efficient Query Execution - Parts 1-2 [Lecture slides].

<https://www.uio.no/studier/emner/matnat/ifi/IN3020/v21/lectures/061-062---qeval.pdf>

Sam, S. (2018, 22. february). How can we convert subqueries to INNER JOIN? Retrieved from

<https://www.tutorialspoint.com/How-can-we-convert-subqueries-to-INNER-JOIN>