

# QUERY EXECUTION

We want to translate the SQL query into an algebraic query to define the order of operations.

Then we do a logical transformation, typically push the selection operation towards the leaves, so that the join is computed later.

Finally, you generate a feasible access plan, choosing the correct algorithm, meaning that we annotate each node with the algo that we are going to use.

The order of the join is the most important decision in terms of speed of the query.



## Physical operators

Each physical operator is a specific implementation of a logical operator.

- $R \rightarrow \text{TableScan}(R)$
- $\pi_x(E) \rightarrow \text{Project}(O_E, X) \rightarrow \text{argument } O_E \text{ and parameter } X$
- $E_1 \bowtie_{\text{cond}} E_2 \rightarrow \text{Nested Loop}(O_1, O_2, \text{cond})$ 
  - $\rightarrow \text{Index Nested Loop}(O_1, O_2, \text{cond})$
  - $\rightarrow \text{MergeJoin}(O_1, O_2, \text{cond})$

How do you implement a tree of operators?

One way is "one operator at a time"

or "materialization of the result of each operator."

Very simple to implement, but requires a lot of read/write.

Not how a DBMS works.

What is used is the iterator open/next/close,  
called "iterator style" or "cursor interface".

For every node we use the iterator, for example

$\text{filter}(\text{Province} = P1) \rightarrow$  it just calls next on the TableScan,  
which remembers the last cursor position

TableScan(student)

The filter doesn't need to keep everything in memory.

The TableScan just stores the position of the cursor.

We write the access plan like a tree,

where the argument of an operator is its child.

## Children of operators

Observe that each operator has  
a fixed number of children.

For example

- TableScan is always a leaf.

- Index Filter ( $R, Ix, \psi$ )  $\leadsto$  it gets out of the index  $Ix$   
all the values that satisfy  
the condition  $\psi$

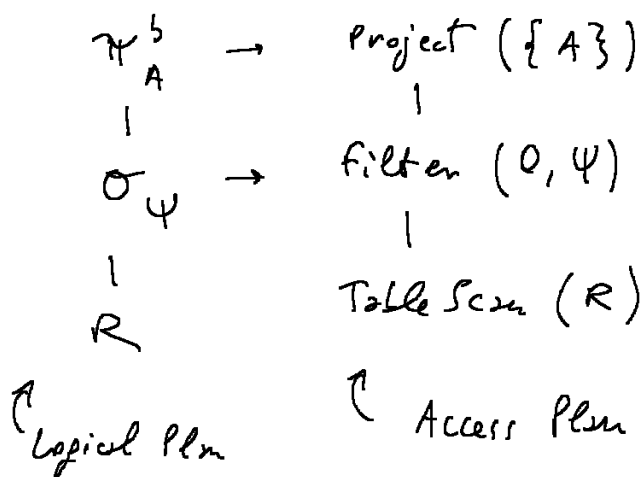
$\leadsto$  Index Filter (Students, IdAgeStud, age > 25)

open the index IdAgeStud,  
retrieve RIDs with age > 25,  
reads corresponding students from  
the table Students

$\leadsto$  this is ALWAYS a leaf,  
because it always reads from disk,  
cannot be combined a priori

## Project and Filter

They are the simplest operations.

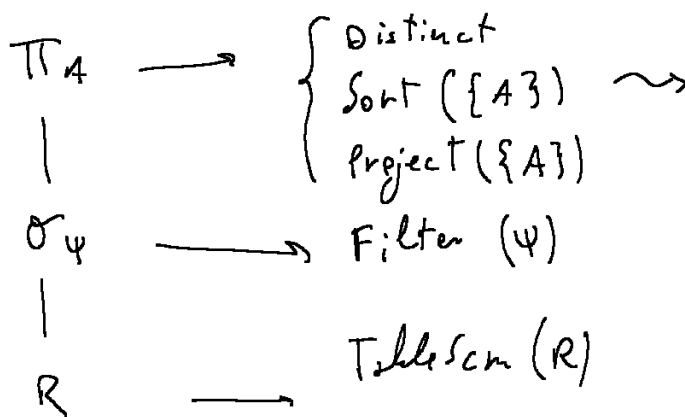


$\pi_A^b$  because SQL produces mult. sets  
(no duplicate elimination)

~ this is a basic  
SELECT  
FROM  
WHERE

## Distinct elimination

If we want to do duplicate elimination, it is more complicated



This is needed because whenever you ask 'next' it remembers only the last element.  
So distinct can compare only with the last element

What is the cost of distinct?  $\emptyset$  because it happens in the main memory.

And Sort? It is  $2 \cdot N \log N$ , because it needs the disk,

so we try to reduce the use of sort,

and avoid it if data is already sorted

And Project?  $\emptyset$

And Filter?  $\emptyset$

And TableScan?  $N \log N$

## Selection and Index Filter

SELECT \*  
FROM R  
WHERE A BETWEEN 50 AND 100

→ with idx as index on A

$\pi_x^b \rightarrow \text{Project} (*)$   
 $\left. \begin{array}{c} | \\ \sigma_\psi \\ | \\ R \end{array} \right\} \rightarrow \text{IndexFilter}(R, \text{idx}, \psi)$

'IndexOnlyFilter' combines access to index with projection.

It just gives the values that it finds in the index,  
without going to the data. Much more efficient.

Go to the index and project to the same attributes of the index.  
It does not read the table.

What is the difference between  $\text{Filter}(O, \psi)$  and  $\text{IndexFilter}(R, \text{idx}, \psi)$ ?

The first filters the records on the given argument  $O$ ,  
while the second access the records of  $R$   
through the index, with no argument.

all depends on  
the selectivity factor

About the cost,

$\text{Filter}(\text{TableScan}(R), \psi) \rightsquigarrow N_{\text{page}}(R)$  ✓

$\text{IndexFilter}(R, \text{idx}, \psi) \rightsquigarrow C_I + C_D = C_I + sf * N_{\text{Rec}}(R)$

In addition,  $\text{IndexFilter}(R, \text{idx}, \psi)$  gives data already sorted by A,  
then it might be not need to sort again.

## Group By

The canonical way is by sorting and then break when the attribute changes.

```
SELECT A, COUNT(*)  
FROM R  
WHERE  
GROUP BY
```

$\sigma_{\text{COUNT}(*)>1} \rightarrow \text{Filter}(\text{COUNT}(*), 1)$   
|

$\{A\} \bowtie \{\text{COUNT}(*)\} \rightarrow \left\{ \begin{array}{l} \text{GroupBy}(\{A\}, \text{COUNT}(*)) \\ \text{Sort} \end{array} \right.$   
|

$\sigma_{\psi} \rightarrow \text{Filter}(\psi)$   
|

$R \rightarrow \text{Table Scan}(R)$

GROUPED means that when there is a break I will never get that value again.

SORTING implies GROUPED, it is stronger.

If you are not sure that data is grouped, then you need to sort.

But since grouped data is enough, it might not be required to sort.

## Sort

The sort operator is  $\tau_{\{A_i\}}$   
and its cost is  $2 * N \log$

## Nested Loop (JOIN)

This is the simplest and most expensive algo.

Algo:

for each record  $r$  in  $R$  do  $\leftarrow$  outer relation  
    for each record in  $S$  do  $\leftarrow$  inner relation  
        if  $r_i = s_j$  then add  $\langle r, s \rangle$  to the result

The cost is  $N_{\text{page}}(R) + N_{\text{Rec}}(R) * N_{\text{Page}}(S)$

$\uparrow$  rest left  
side

$\uparrow$  for every record  
on the left side  
(outer rel.)

$\uparrow$  we scan the entire  
inner relation

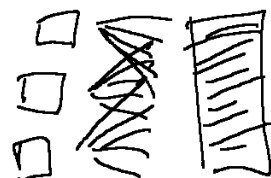
$\therefore$  the cost is quadratic, and you will never use it.

## Page Nested Loop (JOIN)

Instead of scanning the entire relation once for every record of the outer relation, you read in entire page of the outer relation and join the entire page.

More difficult to implement and still quadratic,  
but less expensive since  $N_{\text{page}} < N_{\text{records}}$ .

The cost is  $N_{\text{page}}(R) + N_{\text{Page}}(R) * N_{\text{Page}}(S)$



## Index Nested Loop

for every tuple of  $R$ , read the tuple  $R.A$ ,  
then use the index of  $S.B$ .

Index Nested Loop requires two arguments,  
on the left, something that allows "next", like TableScan,  
on the right,  $\text{IndexFilter}(S, \text{Idx}, \underline{S.B = R.A})$  ←

This is the only operator where  
data flows also from left to right,  
not not only bottom-up.  
It is the 'open' operator that carries  
this information.

observe that here  $S.B$   
opens the index, is a constant,  
while  $R.A$  is the attribute

→ remember that for every  
operator we do  
open | next | close

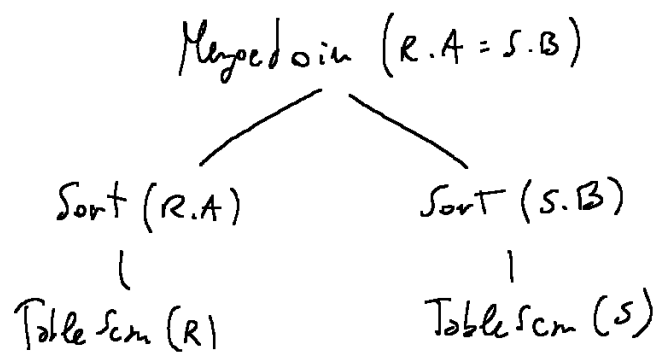
→ the condition is

77' and before



## Merge Join (Join 4) (also Soft Merge)

If both tables have many records.  
The cost is just that of sort,  
that is  $2 * N_{pg}$ , linear in the  
number of pages.



! Before MJ, be sure that  
the input is sorted  
exactly on the join attribute.

Should we use Index Nested Loop or Merge join?

The cost is respectively  $N_{pg}(R) + k * N_{rec}(R)$  and  $2 * N_{pg}$ .

So if S has  $N_{rec} = 100.000$  and  $N_{pg} = 1.000$   
and R has  $N_{rec} = 3$  and  $N_{pg} = 1$ .

Then INL is very convenient, because the cost is  $1 + k * 3$ .  
While MJ's cost is  $2 * 1.000$ .

But if R is as big as S, like

R with  $N_{rec} = 100.000$  and  $N_{pg} = 1.000$ ,

Then INL will be around  $200.000$

and MJ is  $2 * 1.000 = 2.000$

In the transactional application, using few records, use INL,  
while analytical applications use MJ.

Of course, to use INL, we need an INL.  
But typically every DBMS put an index on every key and foreign key.

## Operations of JOIN (review)

1. Nested Loop ( $O_E, O_I, \Psi$ )
  2. PsyeNested Loop ( $O_E, O_I, \Psi$ )
  3. Index Nested Loop ( $O_E, O_I, \Psi$ )
  4. Sort Merge ( $O_E, O_I, \Psi$ )
- } NOT USED
- } USED

## Query Optimization

Typically means moving the restriction  $\sigma$  before the join, that is expensive. This is the logical optimization.

Then we have the physical plan generation of the plan with the optimal cost.

Actually we just avoid terrible plans.