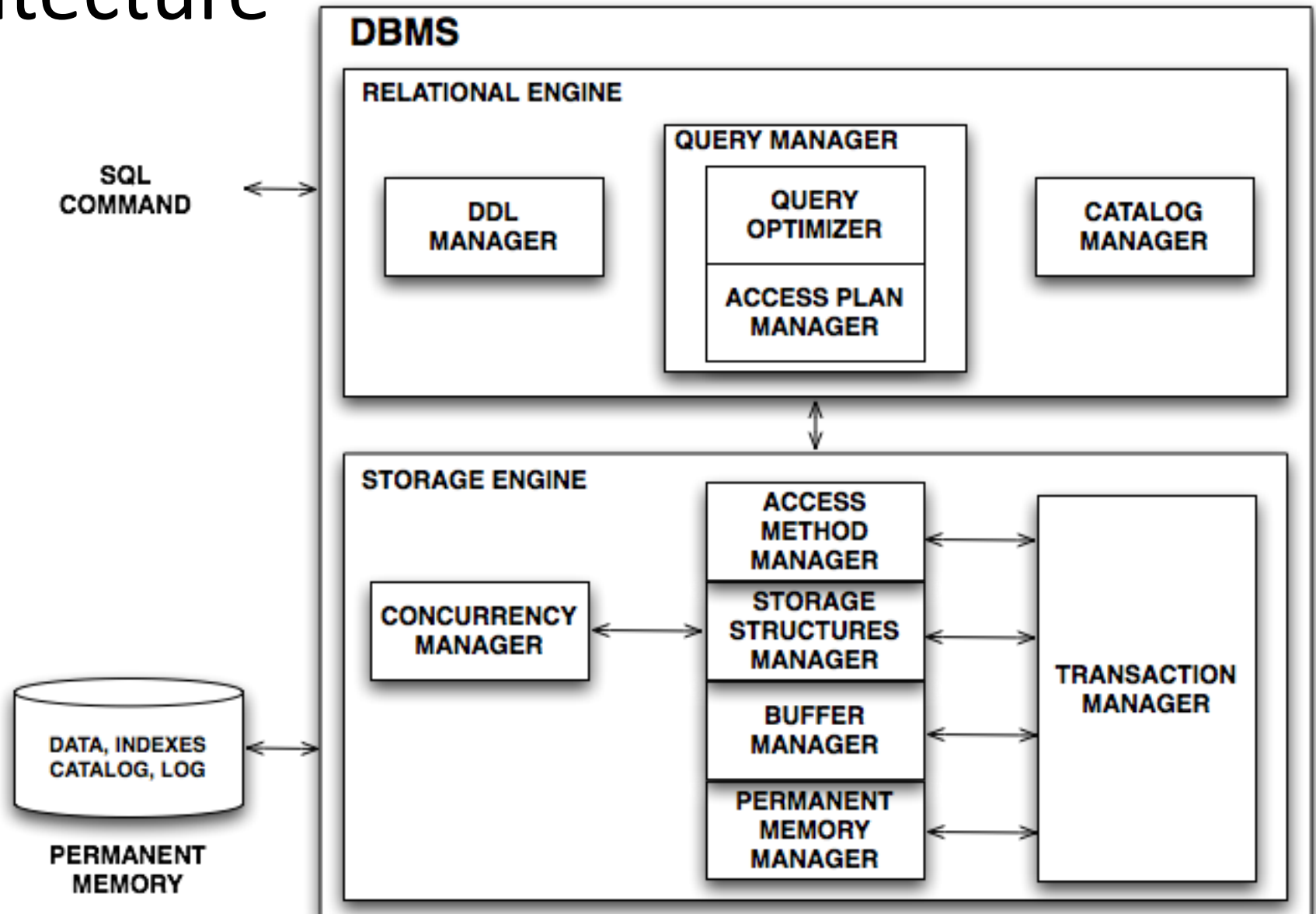


Architecture



Relational operations

- We will consider how to implement and define physical operators for:
 - Projection
 - Selection
 - Grouping
 - Set operators
 - Join
- Then we will discuss how the optimizer use them to generate physical query plans

Java Relational System (JRS):

- <http://www.di.unipi.it/~albano/JRS/toStart.html>

Selectivity of selection $\sigma_{\psi}(E)$

- The DBMS Catalog stores information about database tables and indexes.

$$s_f(A = v) = \frac{1}{N_{\text{key}}(A)} \quad (1/10)$$

$$s_f(A > v) = \frac{\max(A) - v}{\max(A) - \min(A)} \quad (1/3)$$

$$s_f(A < v) = \frac{v - \min(A)}{\max(A) - \min(A)} \quad (1/3)$$

$$s_f(v_1 < A < v_2) = \frac{v_2 - v_1}{\max(A) - \min(A)} \quad (1/4)$$

Selectivity of selection $\sigma_{\psi}(E)$

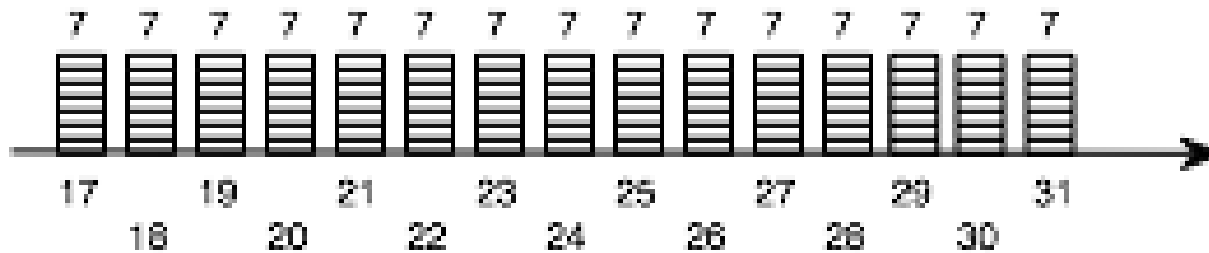
$$s_f(A = B) = \frac{1}{\max(N_{\text{key}}(A), N_{\text{key}}(B))} \quad (1/10)$$

$$s_f(\psi_1 \wedge \psi_2) = s_f(\psi_1) \times s_f(\psi_2)$$

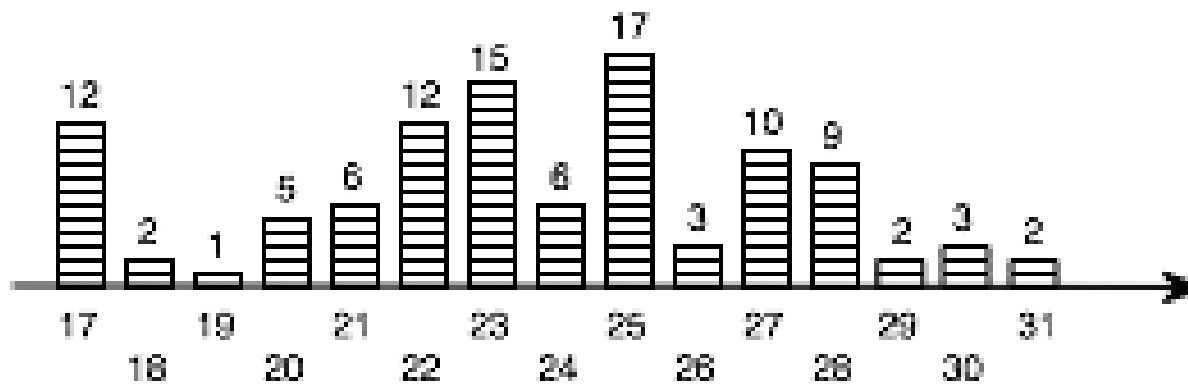
$$s_f(\psi_1 \vee \psi_2) = s_f(\psi_1) + s_f(\psi_2) - s_f(\psi_1) \times s_f(\psi_2)$$

Hystograms

105 records and A with 15 possible values in the range 17 to 31



Uniform distribution

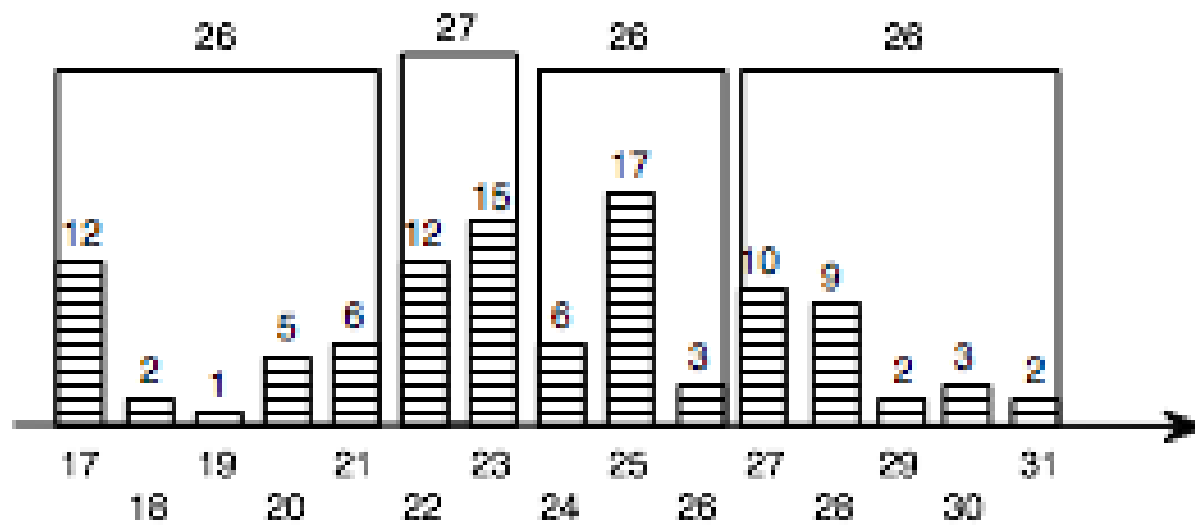


NonUniform distribution

Selectivity of selection with histograms

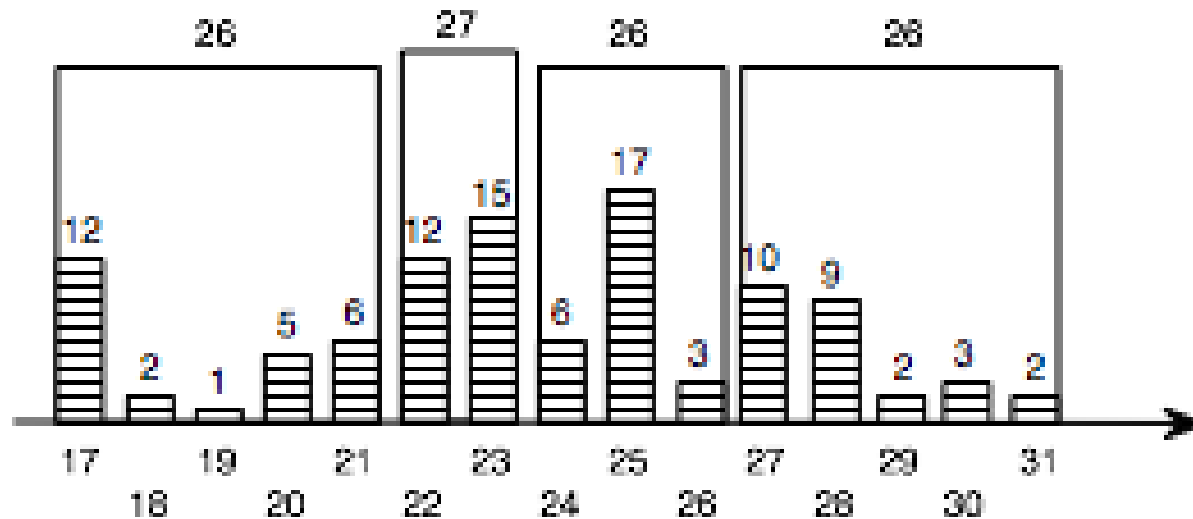
- An Equi-Height histogram is used as an approximation of the actual distribution
- The active domain of A is divided into k intervals, containing a number of records of about N_{rec}/k
- For each interval h_i are known:
 - $\min(h_i)$, $\max(h_i)$, $N_{\text{key}}(h_i)$, $N_{\text{rec}}(h_i)$

Equi-Height histograms



Equi-Height histograms

- $s_f(V = 24) = N_{\text{rec}}(V=24)/N_{\text{rec}}$
- $s_f \text{ real} = 6/105 = 0,057$
- $s_f \text{ uniform} = 1/15 = 7/105 = 0,066$
- $s_f \text{ histEqH} = (26/3)/105 = 0,082$



Physical operators for tables and sort

- Operators for R :
 - TableScan (R)
 - SortScan (R, {A_i})
 - IndexScan (R, Idx)
 - IndexSequentialScan (R, Idx)
- Operator to sort (τ {A_i}):
 - Sort (O, {A_i})

Physical operators for $\pi_b\{A_i\}$

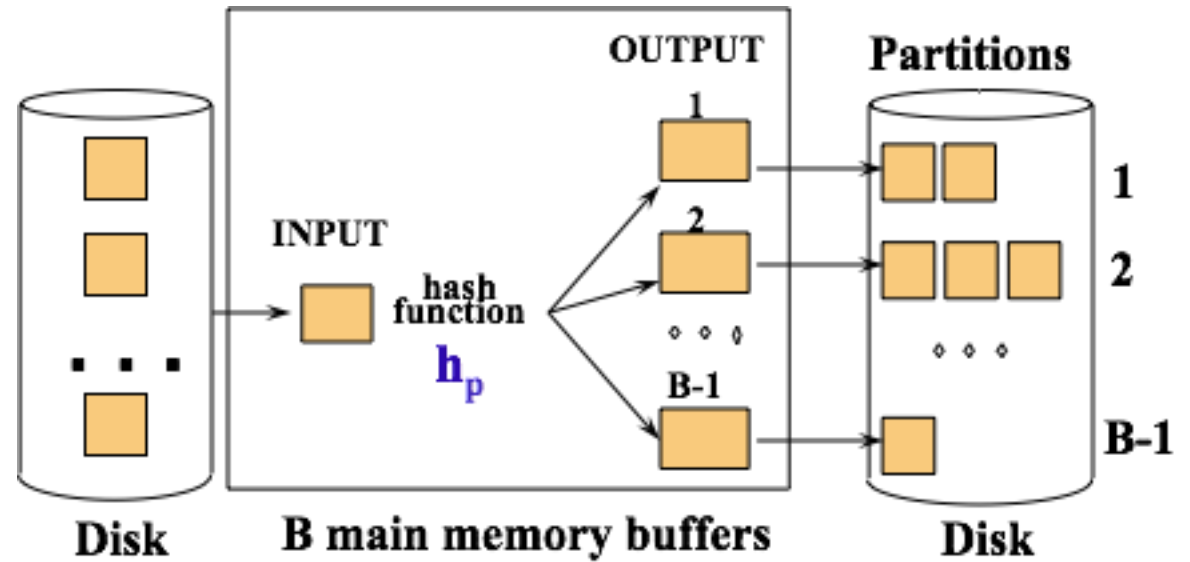
- Project ($O, \{A_i\}$): to project the records of O without duplicates elimination
 - $C = C(O)$
- IndexOnlyScan($R, \text{Idx}, \{A_i\}$):
 - Idx an index on the attributes to project (or that contains them as prefix)
 - $C = N_{\text{leaf}}(\text{Idx})$
- Is the result equivalent to a Project or to a Project+Distinct ?

Physical ops for duplicate elimination

- Distinct (O): to eliminate duplicates from *sorted* records of O
 - $C = C(O)$
- Actually, we only need them to be grouped:
 - $r_i = r_j$ and $i < j \Rightarrow r_i = r_i = r_j$
- HashDistinct(O): to eliminate duplicates from records of O;
 - $C = ?$

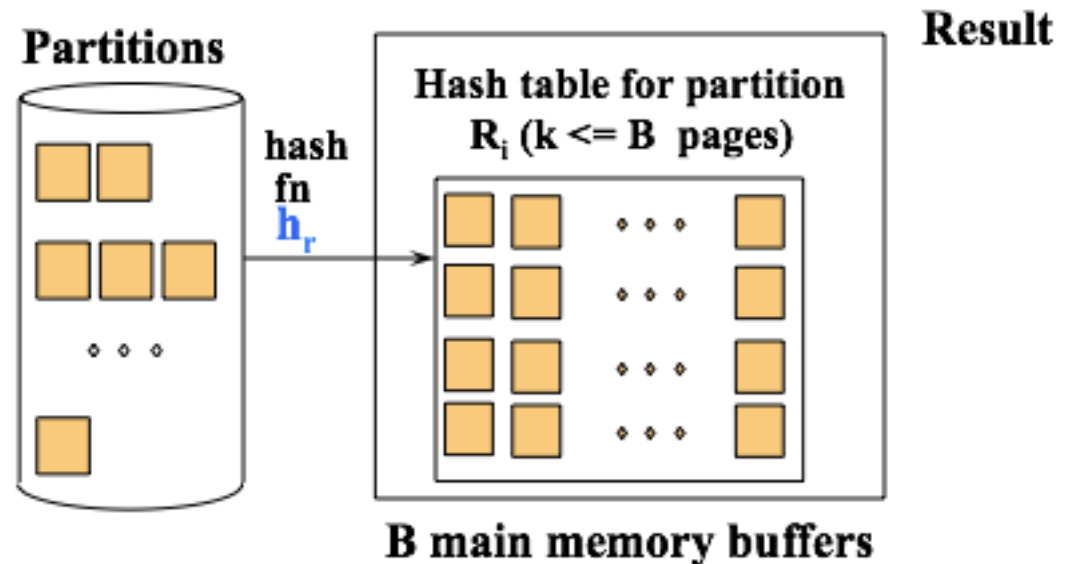
Hashdistinct

- **partitioning phase:**

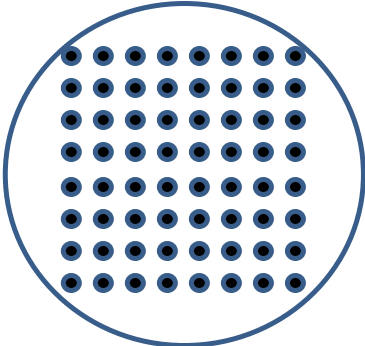
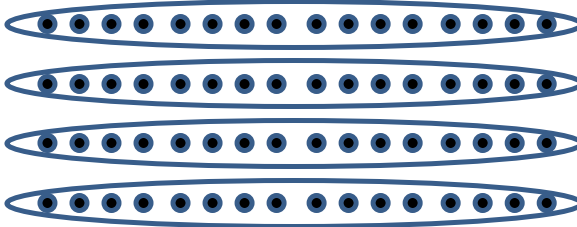
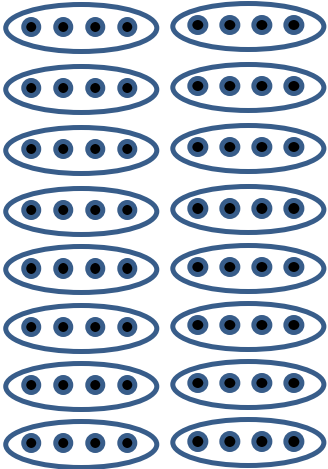
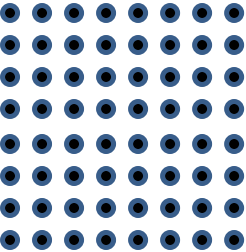


- **duplicate elimination phase ($h_r \neq h_p$):**

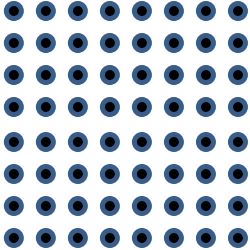
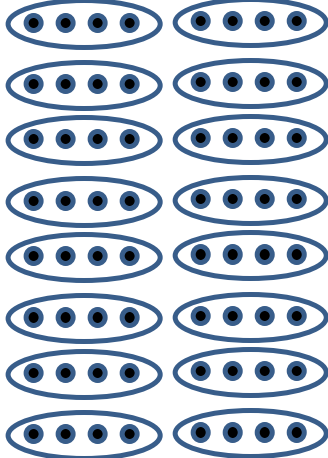
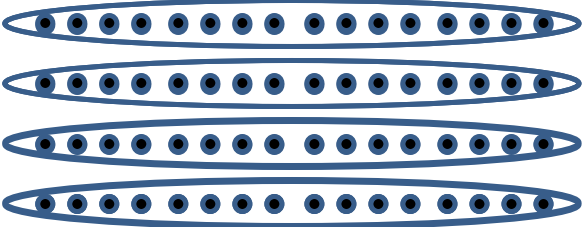
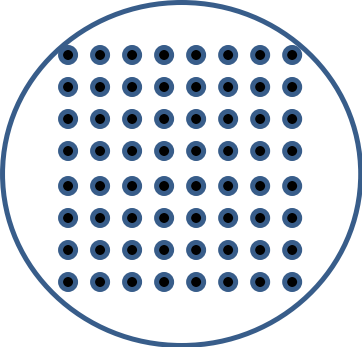
$$C = C(O) + 2 N_{\text{pag}}(O)$$



Sort



Hash



Result size: Distinct(O)

- If A is the only O attribute
 - $|Q| = N_{\text{key}}(A)$
- If $\{A_1, A_2, \dots, A_n\}$ are the attributes
 - $|Q| = \min(|O|/2, \prod_i N_{\text{key}}(A_i))$
- Why not just $\prod_i N_{\text{key}}(A_i)$?

Conclusion

- The sort method often preferred because produces a sorted result and because sort is heavily optimized
- DBMS:
 - Informix uses Hash
 - DB2, Oracle, Sybase ASE use Sort
 - SQL Server e Sybase ASIQ use both methods

Simple selection

```
SELECT *  
FROM Students  
WHERE City = 'PI'
```

- With no index, and data unsorted:
 - Relation scan with cost $N_{\text{pag}}(R)$.
- With an index on selection attribute:
 - Use index to retrieve RIDs, then retrieve data records: $C_I + C_D$.
- Cost depends on qualifying tuples (size of $R \times s_f$), and clustering
- General selection conditions...

Physical operators for selection

- Filter (O, ψ)
 - $C = C(O)$
 - $E_{\text{rec}} = s_f(\psi) \times N_{\text{rec}}(O)$
- IndexFilter(R, Idx, ψ)
 - = RidIndexFilter(Idx, ψ) + TableAccess(O, R)
 - $C = C_I + C_D$
 - $E_{\text{rec}} = s_f(\psi) \times N_{\text{rec}}(R)$

Physical operators for selection

- $\text{IndexSequentialFilter}(R, \text{Idx}, \psi)$
 - $C = s_f(\psi) \times N_{\text{leaf}}(\text{Idx})$
- $\text{IndexOnlyFilter}(R, \text{Idx}, \{A_i\}, \psi)$
 - $C = s_f(\psi) \times N_{\text{leaf}}(\text{Idx})$
- $\text{AndIndexFilter}(R, \{\text{Idx}_i, \psi_i\})$
 - $C_I = \sum_i C_I(\text{Idx}_i)$
 - $C_D = \Phi(E_{\text{rec}}, N_{\text{pag}}(R))$
 - $E_{\text{rec}} = s_f(\psi) \times N_{\text{rec}}(R)$
- $\text{OrIndexFilter}(R, \{\text{Idx}_i, \psi_i\})$

Conjunctive selection in DBMSs

- Informix, DB2 use intersection of RID sets
- Oracle, Sybase ASIQ use bitmaps
- Oracle use also HashJoin with index (later) on the attribute RID
- Sybase ASE uses one index only
- SQL Server use index join (later)

Exercises

SELECT A, B Idx an index on A
FROM R
WHERE (A **BETWEEN** 50 **AND** 100) **AND** B > 20;

SELECT A, B One index on A, B
FROM R
WHERE (A **BETWEEN** 50 **AND** 100) **AND** B > 20
ORDER BY A;

SELECT **DISTINCT** A, B Two indexes on A and on B
FROM R
WHERE (A **BETWEEN** 50 **AND** 100) **AND** B > 20
ORDER BY B;

Physical operators for grouping

- As for duplicate elimination:
 - Sorting
 - Hashing
 - Using an index on the grouping attributes
- (Duplicate elimination is grouping on all attributes with no aggregation)

Physical operators for ($\{A_i\} \vee \{f_j\}$)

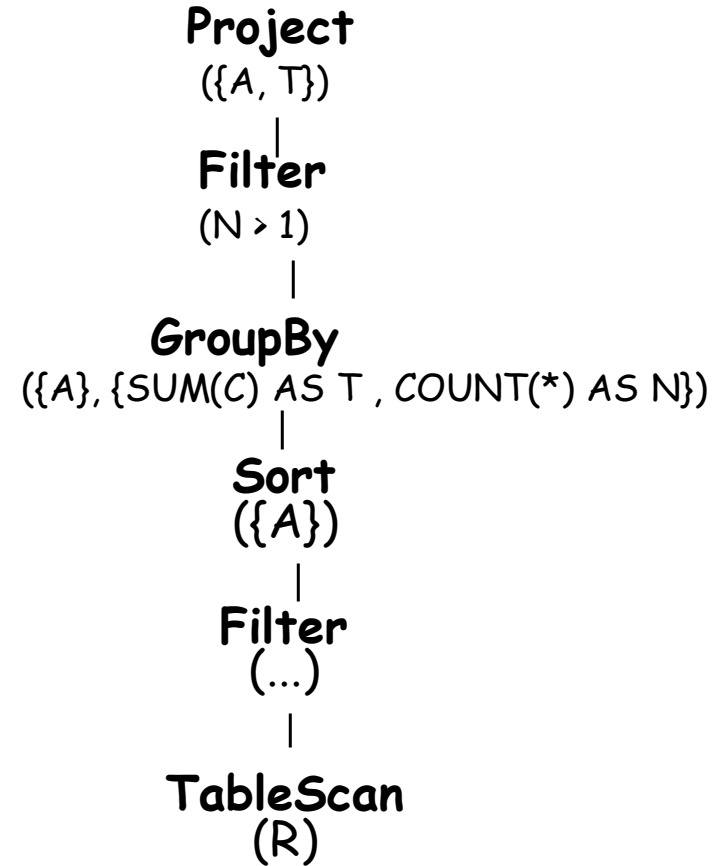
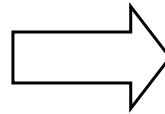
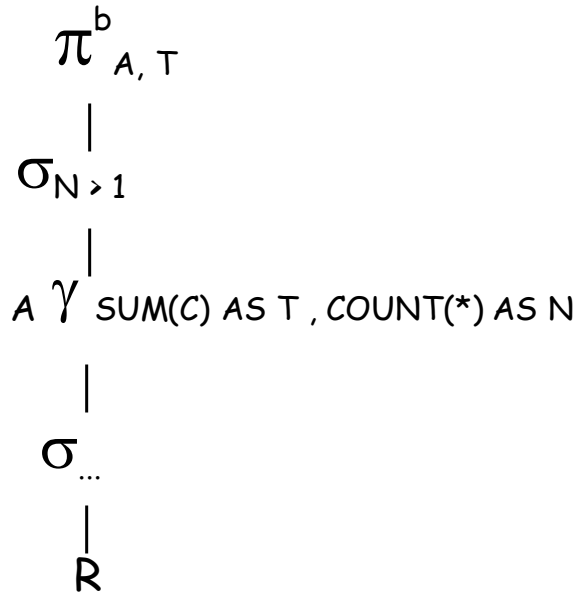
- GroupBy (O, $\{A_i\}$, $\{f_j\}$): to group the records of O **sorted** on the $\{A_i\}$, computing functions in $\{f_j\}$.
 - $\{f_j\}$ contains the aggregation functions present in the SELECT and HAVING clauses.
 - The operator returns records with attributes $\{A_i\} \cup \{f_j\}$
 - The records of O must be sorted on the $\{A_i\}$
 - Any permutation of $\{A_i\}$ is ok
 - (Actually, they only need to be grouped on the $\{A_i\}$)
- HashGroupBy (O, $\{A_i\}$, $\{f_j\}$)
- C? Cardinality?

Computing aggregations

- Each aggregation function is an object `o` with methods `start()`, `next(v)`, `end()`
- With the first record of a group: `o.start()`
- For each record of the group `o.next(v)` is invoked
- `o.end()` computes the final aggregate value.

Example

```
SELECT  A, SUM(C) AS T
FROM    R
WHERE   A BETWEEN 50 AND 100
GROUP BY A
HAVING  COUNT(*) > 1;
```



Logical Tree

Physical Tree

... = A BETWEEN 50 AND 100

Physical ops for join: nested loops

- foreach r in O_E do
 - foreach s in O_I do
 - if $r.r1 = s.s1$ then add $\langle r, s \rangle$ to result
- $C = C(O_E) + E_{rec}(O_E) \times C(O_I)$
- $E_{rec} = s_f(C_j) \times E_{rec}(O_E) \times E_{rec}(O_I)$
- Inequality join conditions ?

Cost of nested loops

```
foreach r in R do
  foreach s in S do
    if r.r1 = s.s1 then add <r, s> to result
```

With R **external**:

$$C = \text{Npag}(R) + \text{Nrec}(R) \times \text{Npag}(S) \approx \text{Npag}(R) \times \frac{\text{Nrec}(R)}{\text{Npag}(R)} \times \text{Npag}(S)$$

With S **external**:

$$C = \text{Npag}(S) + \text{Nrec}(S) \times \text{Npag}(R) \approx \text{Npag}(S) \times \frac{\text{Nrec}(S)}{\text{Npag}(S)} \times \text{Npag}(R)$$

Hence...

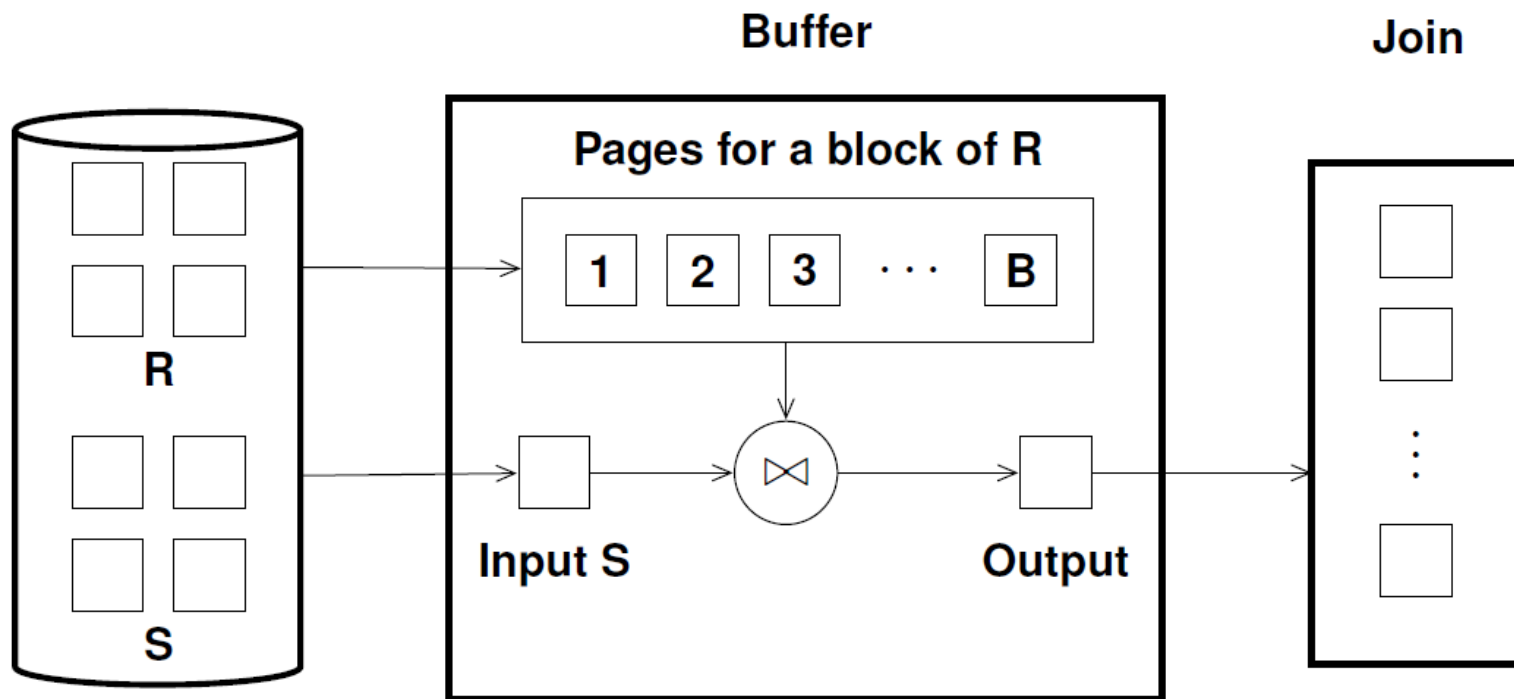
Page nested loops

R	<table border="1"><thead><tr><th>Pk</th><th>B</th></tr></thead><tbody><tr><td>10</td><td>a</td></tr><tr><td>20</td><td>b</td></tr><tr><td>30</td><td>c</td></tr><tr><td>40</td><td>d</td></tr></tbody></table>	Pk	B	10	a	20	b	30	c	40	d	Join	S	<table border="1"><thead><tr><th>Fk</th><th>C</th></tr></thead><tbody><tr><td>10</td><td>f</td></tr><tr><td>20</td><td>g</td></tr><tr><td>30</td><td>g</td></tr><tr><td>10</td><td>k</td></tr><tr><td>40</td><td>f</td></tr><tr><td>20</td><td>k</td></tr></tbody></table>	Fk	C	10	f	20	g	30	g	10	k	40	f	20	k	=	<table border="1"><thead><tr><th>Pk</th><th>B</th><th>Fk</th><th>C</th></tr></thead><tbody><tr><td>10</td><td>a</td><td>10</td><td>f</td></tr><tr><td>20</td><td>b</td><td>20</td><td>g</td></tr><tr><td>10</td><td>a</td><td>10</td><td>k</td></tr><tr><td>20</td><td>b</td><td>20</td><td>k</td></tr><tr><td>30</td><td>c</td><td>30</td><td>g</td></tr><tr><td>40</td><td>d</td><td>40</td><td>f</td></tr></tbody></table>	Pk	B	Fk	C	10	a	10	f	20	b	20	g	10	a	10	k	20	b	20	k	30	c	30	g	40	d	40	f
Pk	B																																																									
10	a																																																									
20	b																																																									
30	c																																																									
40	d																																																									
Fk	C																																																									
10	f																																																									
20	g																																																									
30	g																																																									
10	k																																																									
40	f																																																									
20	k																																																									
Pk	B	Fk	C																																																							
10	a	10	f																																																							
20	b	20	g																																																							
10	a	10	k																																																							
20	b	20	k																																																							
30	c	30	g																																																							
40	d	40	f																																																							

$$C = \text{Npag}(R) + \text{Npag}(R) * \text{Npag}(S)$$

The smallest relation is used as **external**

Block nested loops join



$$C_{BNL} = N_{pag}(R) + \lceil N_{pag}(R)/B \rceil \times N_{pag}(S)$$

Index nested loop

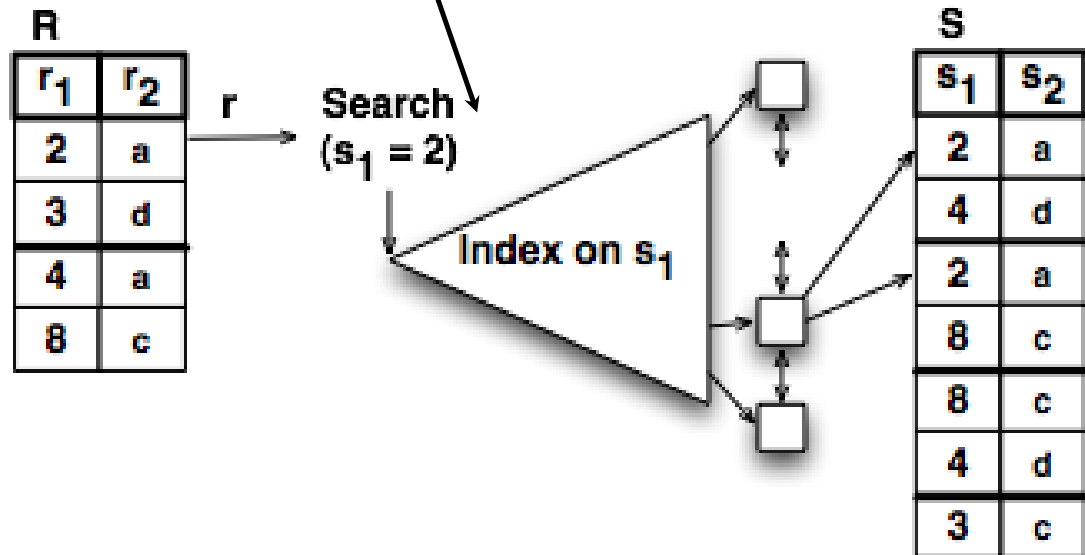
- Hyp: There is an index on the join column of the internal relation (S)

foreach r in R do

 foreach s in *IndexFilter*(S, I, s.s1=r.r1) do

 add <r, s> to result

$$R(r_1, r_2) \bowtie_{r_1=s_1} S(s_1, s_2)$$



Index nested loop

foreach r in R do

 foreach s in *IndexFilter*(S,l,s1 =r.r1) do

 add <r, s> to result

- Cost for R join S:

$$- C = N_{\text{pag}}(R) + N_{\text{rec}}(R) \times \text{CaWithIdx}(S)$$

- General Case:

$$- C = C(O_E) + E_{\text{rec}}(O_E) \times (C_I + C_D)$$

Merge join

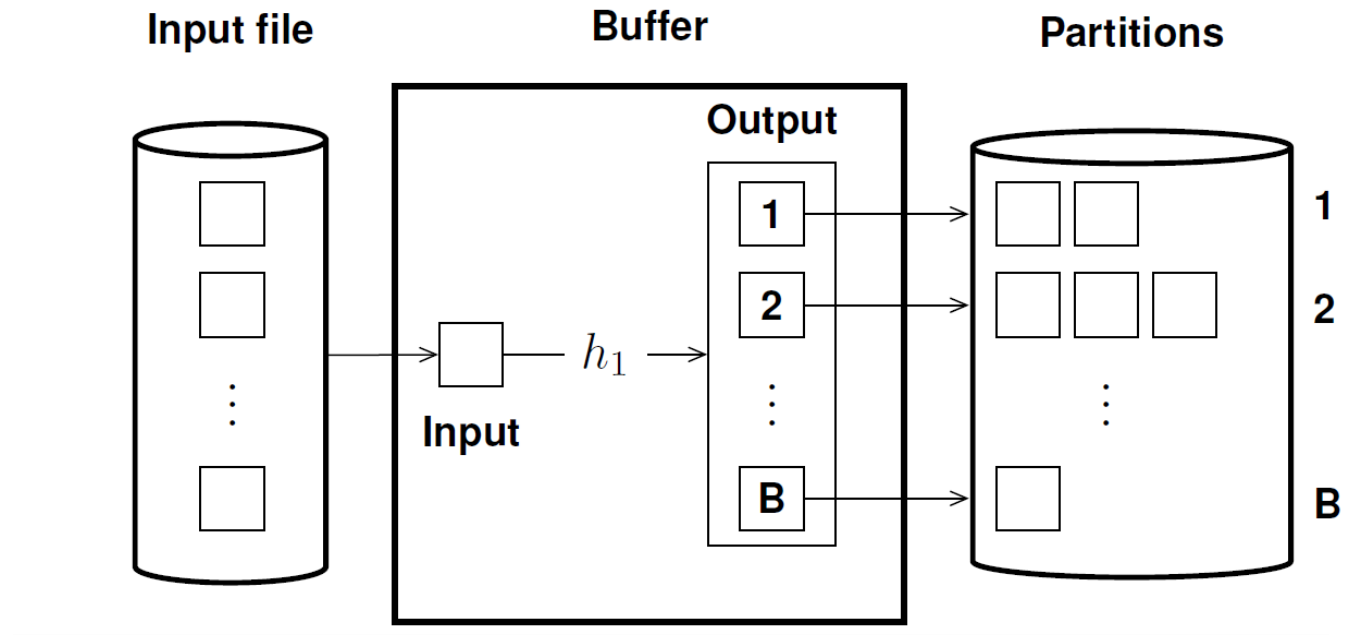
- Hyp: R and S are sorted on the join attribute, a key of the external relation
- $C = C(O_E) + C(O_I)$
- Inequality join conditions ?

...	ri
...	1
...	2
...	3
...	4

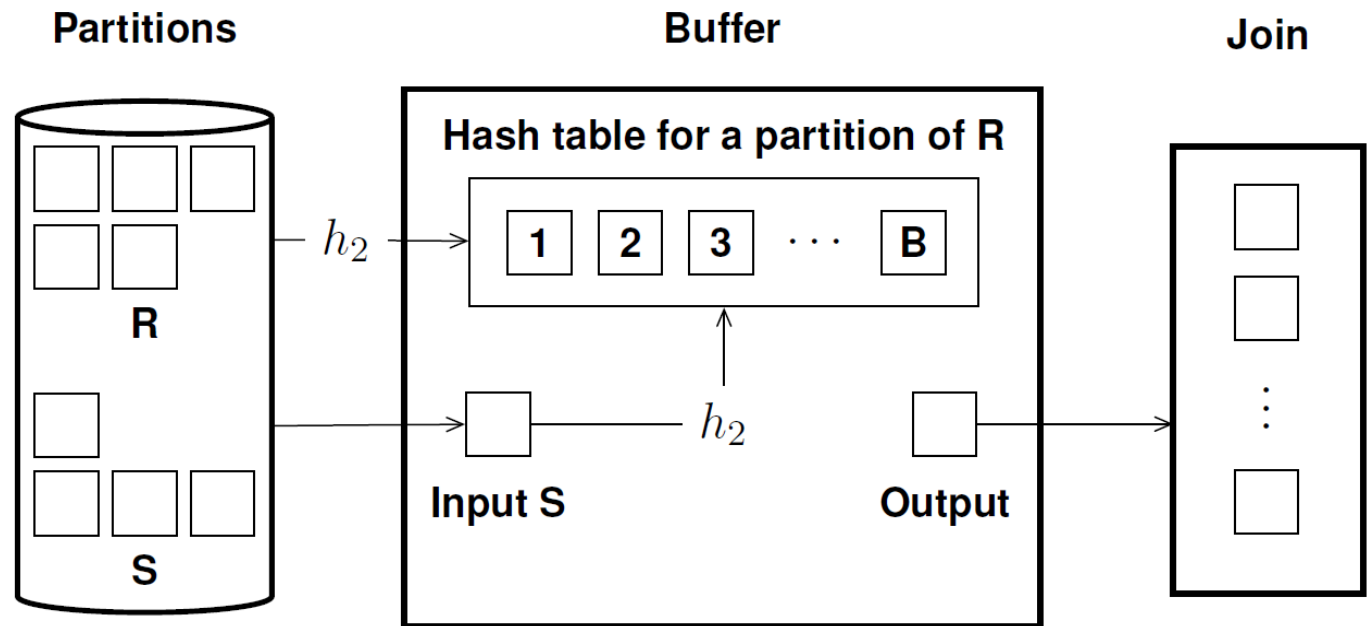
sj	...
1	...
1	...
3	...
3	...
3	...
4	...

Hash join

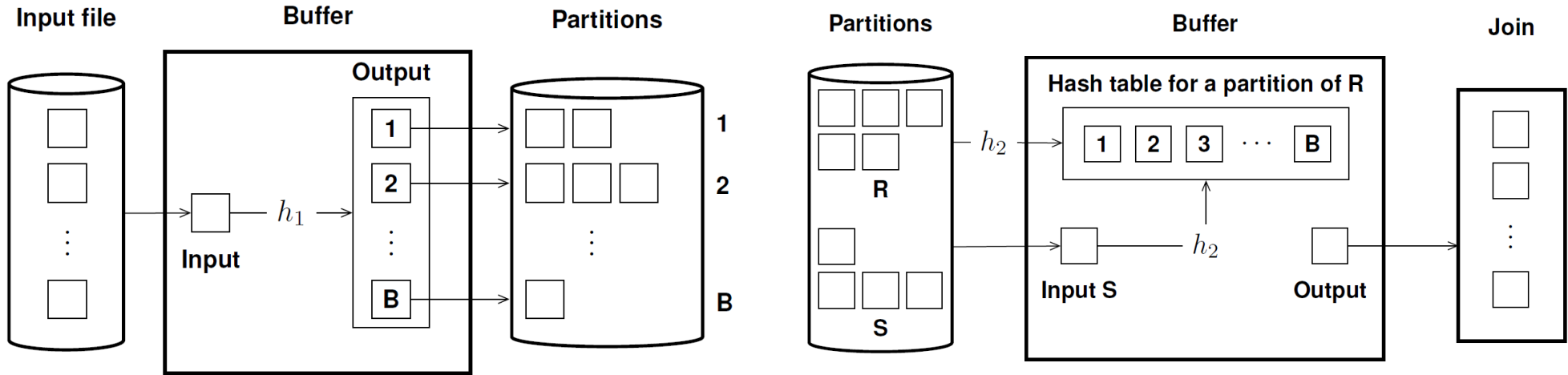
- Partitioning (R and S)



- Matching



Hash join: cost



- Assume $N_{\text{pag}}(R)/B < B$ and uniformity
- $C = C(O_E) + C(O_I) + 2(N_{\text{pag}}(O_E) + N_{\text{pag}}(O_I))$
- $C = (\log_B(N_{\text{pag}}(O_E)) \times 2-2)^*(N_{\text{pag}}(O_E) + N_{\text{pag}}(O_I))$
- $C(R \text{ join } S) = 3 \times (N_{\text{pag}}(R) + N_{\text{pag}}(S))$
- What if $N_{\text{pag}}(R) < B$?

Complex joins

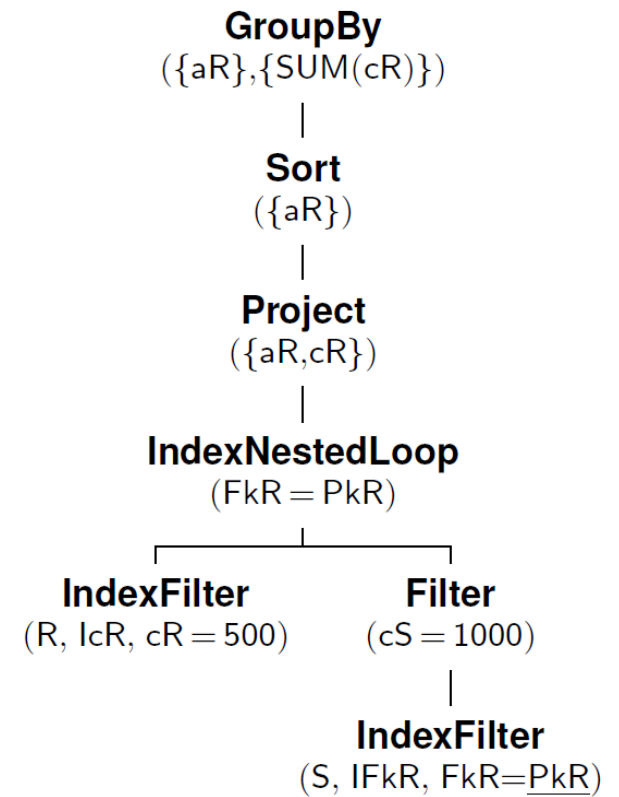
- Join with a conjunctive condition
- Join with a disjunctive condition

Physical operators for join

- NestedLoop (O_E, O_I, ψ_J)
- PageNestedLoop (O_E, O_I, ψ_J)
- IndexNestedLoop (O_E, O_I, ψ_J)
- MergeJoin (O_E, O_I, ψ_J)
- HashJoin (O_E, O_I, ψ_J)

Example

SELECT aR, Sum(cR)
FROM R, S
WHERE R.PkR=S.FkR **AND** cS=1000
GROUP BY aR
ORDER BY aR



$$C_{GroupBy} = C_{Sort} = C_{IndexNestedLoop} + 2 \times N_{pag}(Project)$$

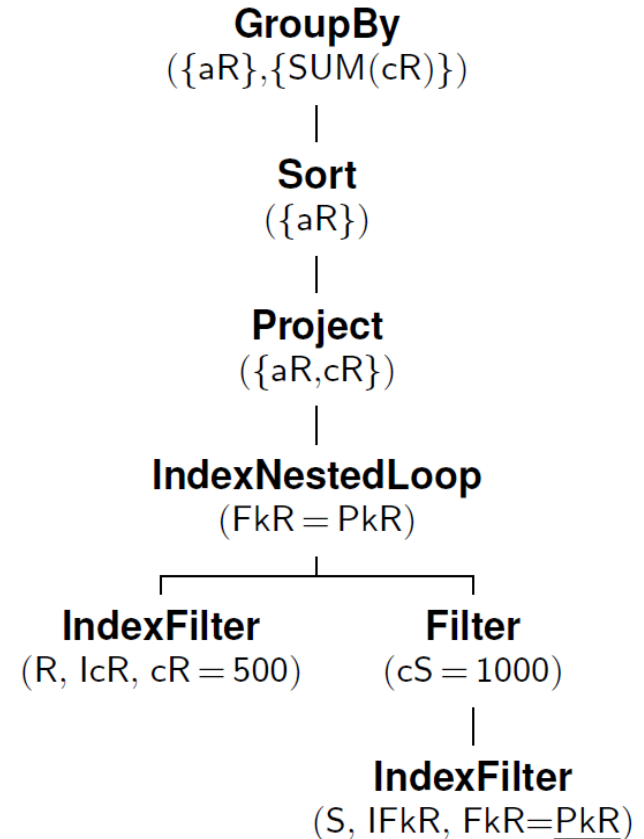
$$N_{Pag}(Project) = (E_{Rec}(IndNestedLoop) \times (L(aR) + L(cR))) / D_{pag}$$

$$C_{INL} = C(IndexFilterOnR) + E_{rec}(IndexFilterOnR) \times C(Filter)$$

$$E_{rec}(INL) = s_f(PkR = FkR) \times E_{rec}(IndexFilterOnR) \times E_{rec}(Filter)$$

Example

- $C(\text{IndFilOnR}) = C_I + C_D$
 - $C_I = \lceil N_{\text{leaf}}(\text{IcR}) / N_{\text{key}}(\text{IcR}) \rceil$
 - $C_D = \Phi(N_{\text{rec}}(R) / N_{\text{key}}(\text{IcR}), N_{\text{pag}}(R))$
- $E_{\text{rec}}(\text{IndFilOnR}) = N_{\text{rec}}(R) / N_{\text{key}}(\text{IcR})$
- $C(\text{IndFilOnS}) = C_I + C_D$
 - $C_I = \lceil N_{\text{leaf}}(\text{IFkR}) / N_{\text{key}}(\text{IFkR}) \rceil$
 - $C_D = \Phi(N_{\text{rec}}(S) / N_{\text{key}}(\text{IFkR}), N_{\text{pag}}(S))$
- $E_{\text{rec}}(\text{Filter}) = s_f(cS=1000) \times N_{\text{rec}}(S)$

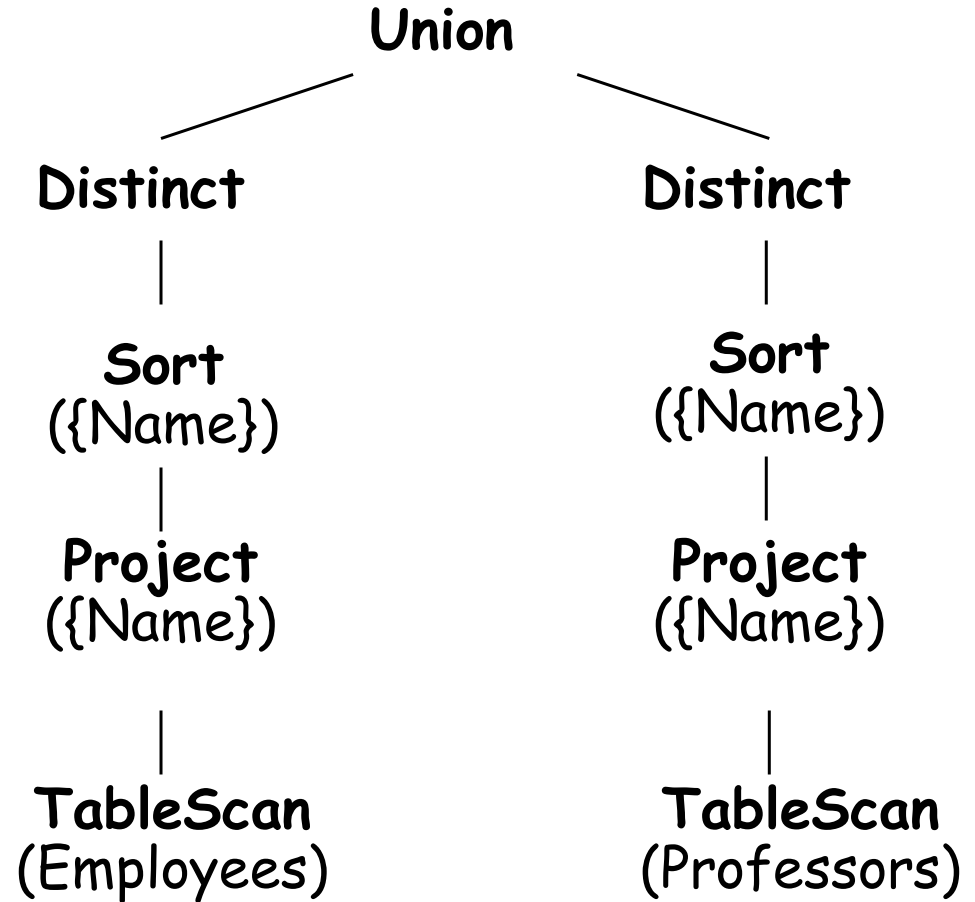


Set operators

- $\text{Union}(O_E, O_I)$, $\text{Except}(O_E, O_I)$, $\text{Intersect}(O_E, O_I)$
 - Operand sorted and without duplicates.
- $\text{HashUnion}(O_E, O_I)$, $\text{HashExcept}(O_E, O_I)$,
 $\text{HasIntersect}(O_E, O_I)$
 - Using hash.
- $\text{UnionAll}(O_E, O_I)$
 - trivial
- $\text{ExceptAll}(O_E, O_I)$, $\text{IntersectAll}(O_E, O_I)$
 - Not always supported

Example

```
SELECT Name
FROM Employees
UNION
SELECT Name
FROM Professors ;
```



Summary

- Very few basic operators, hence the implementation of these operators can be carefully tuned.
- No universally superior technique for operators with many implementations
- We must consider available alternatives and select the best one: “Query optimization”
- ... the next topic