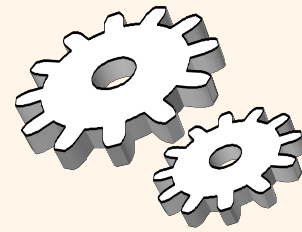


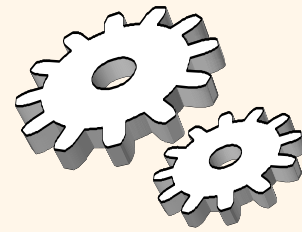
# *Evaluation of Relational Operations*

## Chapter 14, Part A (Joins)



# Relational Operations

- ❖ We will consider how to implement:
  - Selection (  $\sigma$  ) Selects a subset of rows from relation.
  - Projection (  $\pi$  ) ~~Deletes~~ unwanted columns from relation.
  - Join (  $\bowtie$  ) Allows us to combine two relations.
  - Set-difference (  $-$  ) Tuples in reln. 1, but not in reln. 2.
  - Union (  $\cup$  ) Tuples in reln. 1 and in reln. 2.
  - Aggregation (SUM, MIN, etc.) and GROUP BY



# Schema for Examples

Sailors (sid: integer, sname: string, rating: integer, age: real)

Reserves (sid: integer, bid: integer, day: dates, rname: string)

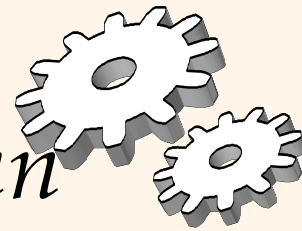
## ❖ Reserves:

- Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.

## ❖ Sailors:

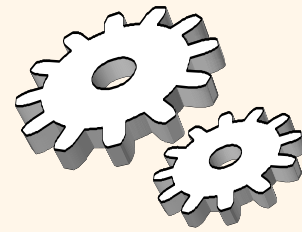
- Each tuple is 50 bytes long, 80 tuples per page, 500 pages.

# Equality Joins With One Join Column



```
SELECT *  
FROM   Reserves R1, Sailors S1  
WHERE  R1.sid=S1.sid
```

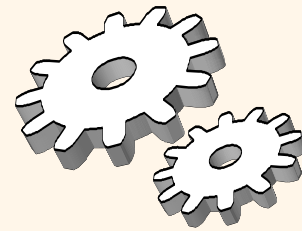
- ❖ In algebra:  $R \bowtie S$ . Common! Must be carefully optimized.  $R \times S$  is large; so,  $R \times S$  followed by a selection is inefficient.
- ❖ Assume:  $M$  pages for  $R$ ,  $p_R$  tuples per page,  $N$  pages for  $S$ ,  $p_S$  tuples per page.
  - In our examples,  $R$  is Reserves and  $S$  is Sailors.
- ❖ *Cost metric*: # of I/Os. We will ignore output costs.



# *Join Algorithms to Consider*

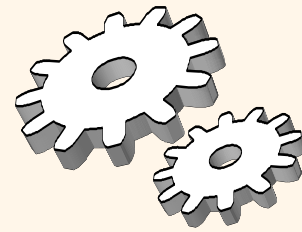
- ❖ Nested loop join
- ❖ Sort-merge join
- ❖ Hash join
- ❖ Index nested loop join

# Simple Nested Loops Join



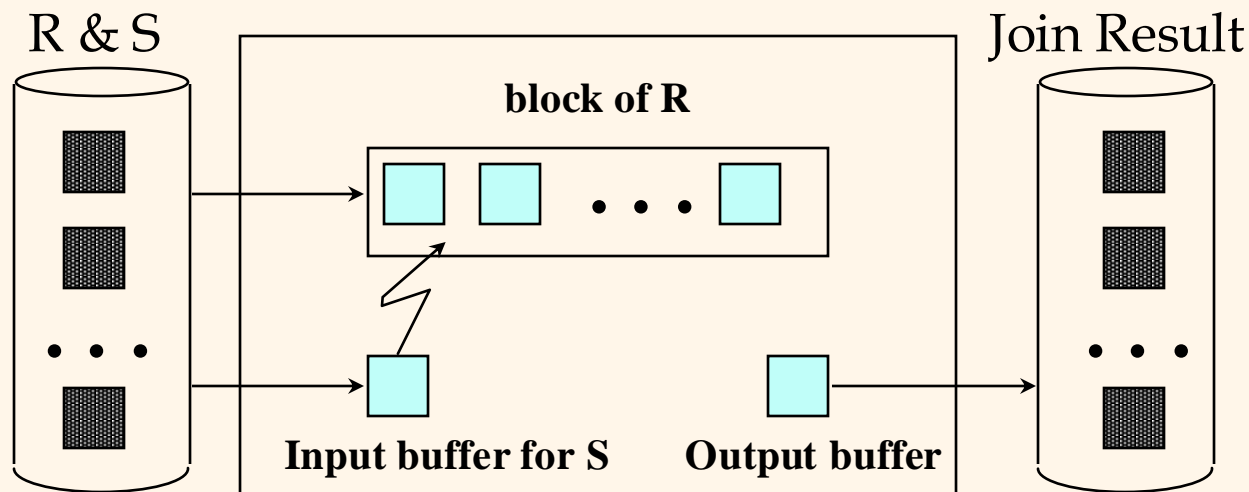
```
foreach tuple r in R do
    foreach tuple s in S do
        if  $r_i == s_j$  then add  $\langle r, s \rangle$  to result
```

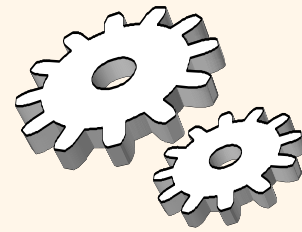
- ❖ Tuple-oriented NLJ: For each tuple in the *outer* relation R, we scan the entire *inner* relation S.
  - Cost:  $M + p_R * M * N = 1000 + 100 * 1000 * 500$  I/Os.
- ❖ Page-oriented NLJ: For each *page* of R, get each *page* of S, and write out matching pairs of tuples  $\langle r, s \rangle$ , where r is in R-page and S is in S-page.
  - Cost:  $M + M * N = 1000 + 1000 * 500$
  - If smaller relation (S) is outer, cost =  $500 + 500 * 1000$



# Block Nested Loops Join

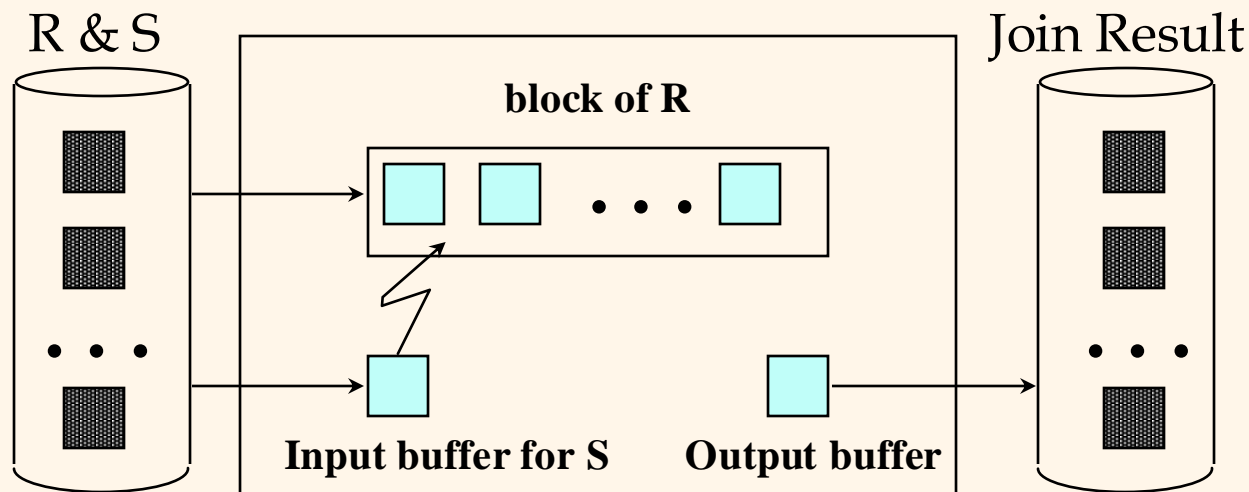
- ❖ Use one page as an input buffer for scanning the inner S, one page as the output buffer, and use all remaining pages to hold ``block'' of outer R.
- For each matching tuple  $r$  in R-block,  $s$  in S-page, add  $\langle r, s \rangle$  to result. Then read next R-block, scan S, etc.



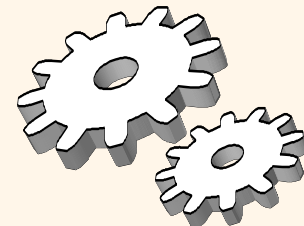


# *Block Nested Loops Join*

- ❖ R is scanned once, cost = M pages
- ❖ Each block has size B-2 (B is # of buffer pages), so need to read S a total of  $\text{ceiling}(M/[B-2])$  times
- ❖ Total cost:  $M + N * \text{ceiling}(M/[B-2])$







# Examples of Block Nested Loops

❖ Cost: Scan of outer + #outer blocks \* scan of inner

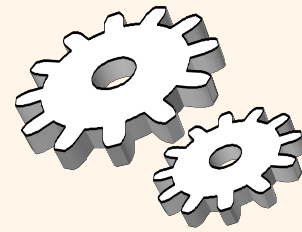
▪ #outer blocks =  ~~$\lceil \frac{\text{#pages of outer}}{\text{#pages of inner}} \rceil$~~

❖ With Reserves (R) as outer, and 100 pages of R:

- Cost of scanning R is 1000 I/Os; a total of 10 *blocks*.
- Per block of R, we scan Sailors (S); 10\*500 I/Os.
- If space for just 90 pages of R, we would scan S 12 times.

❖ With 100-page block of Sailors as outer:

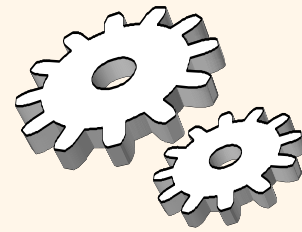
- Cost of scanning S is 500 I/Os; a total of 5 blocks.
- Per block of S, we scan Reserves; 5\*1000 I/Os.



## Sort-Merge Join $(R \bowtie_{i=j} S)$

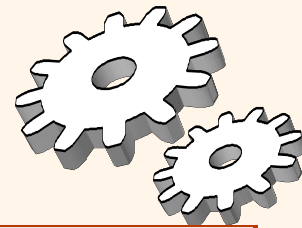
- ❖ Sort R and S on the join column, then scan them to do a “merge” (on join col.), and output result tuples.
  - Advance scan of R until current R-tuple  $\geq$  current S tuple, then advance scan of S until current S-tuple  $\geq$  current R tuple; do this until current R tuple = current S tuple.
  - At this point, all R tuples with same value in  $R_i$  (*current R group*) and all S tuples with same value in  $S_j$  (*current S group*) match; output  $\langle r, s \rangle$  for all pairs of such tuples.
  - Then resume scanning R and S.
- ❖ R is scanned once; each S group is scanned once per matching R tuple. (Multiple scans of an S group are likely to find needed pages in buffer.)

# Example of Sort-Merge Join



<u>sid</u>	<u>snam e</u>	<u>rating</u>	<u>age</u>
2 2	d u s t i n	7	4 5 .0
2 8	y u p p y	9	3 5 .0
3 1	l u b b e r	8	5 5 .5
4 4	g u p p y	5	3 5 .0
5 8	r u s t y	1 0	3 5 .0

<u>sid</u>	<u>snam e</u>	<u>rating</u>	<u>age</u>
2 2	d u s t i n	7	4 5 .0
2 8	y u p p y	9	3 5 .0
3 1	l u b b e r	8	5 5 .5
4 4	g u p p y	5	3 5 .0
5 8	r u s t y	1 0	3 5 .0

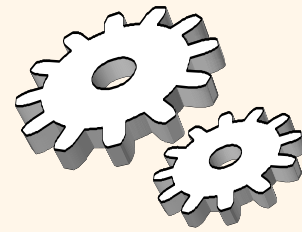


# Example of Sort-Merge Join

<u>s i d</u>	s n a m e	r a t i n g	a g e
2 2	d u s t i n	7	4 5 .0
2 8	y u p p y	9	3 5 .0
3 1	l u b b e r	8	5 5 .5
4 4	g u p p y	5	3 5 .0
5 8	r u s t y	1 0	3 5 .0

<u>s i d</u>	<u>s n a m e</u>	<u>r a t i n g</u>	<u>a g e</u>
2 2	d u s t i n	7	4 5 .0
2 8	y u p p y	9	3 5 .0
3 1	l u b b e r	8	5 5 .5
4 4	g u p p y	5	3 5 .0
5 8	r u s t y	1 0	3 5 .0

- ❖ Cost:  $M \log M + N \log N + (M+N)$  [this is incorrect!]
  - The cost of scanning,  $M+N$ , could be  $M*N$  (very unlikely!)

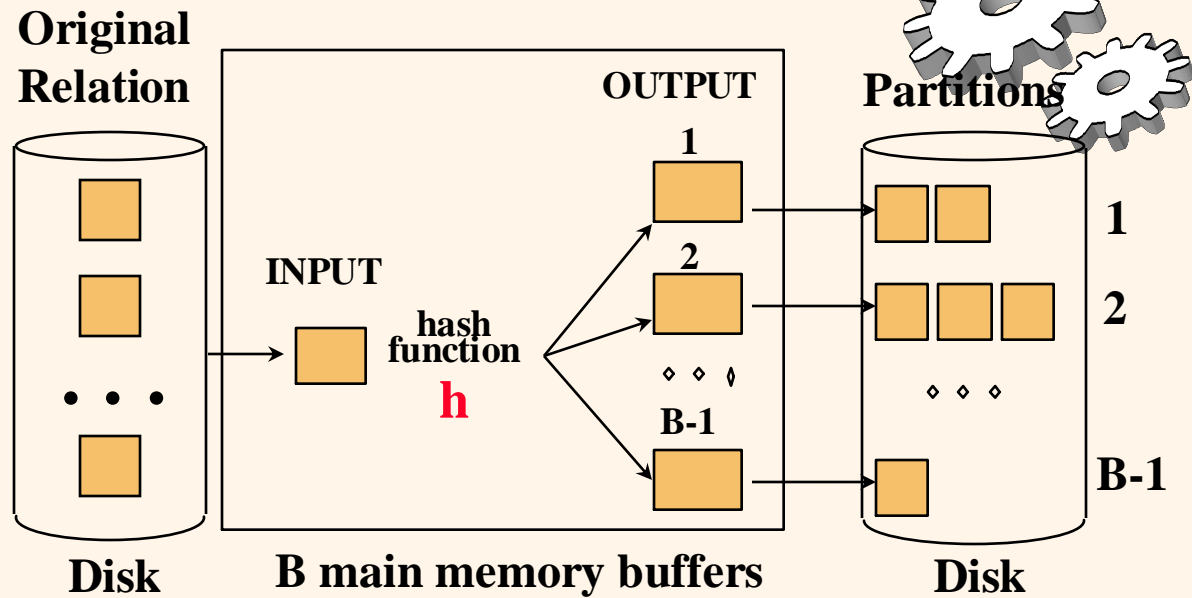


# *Cost of Sort-Merge Join*

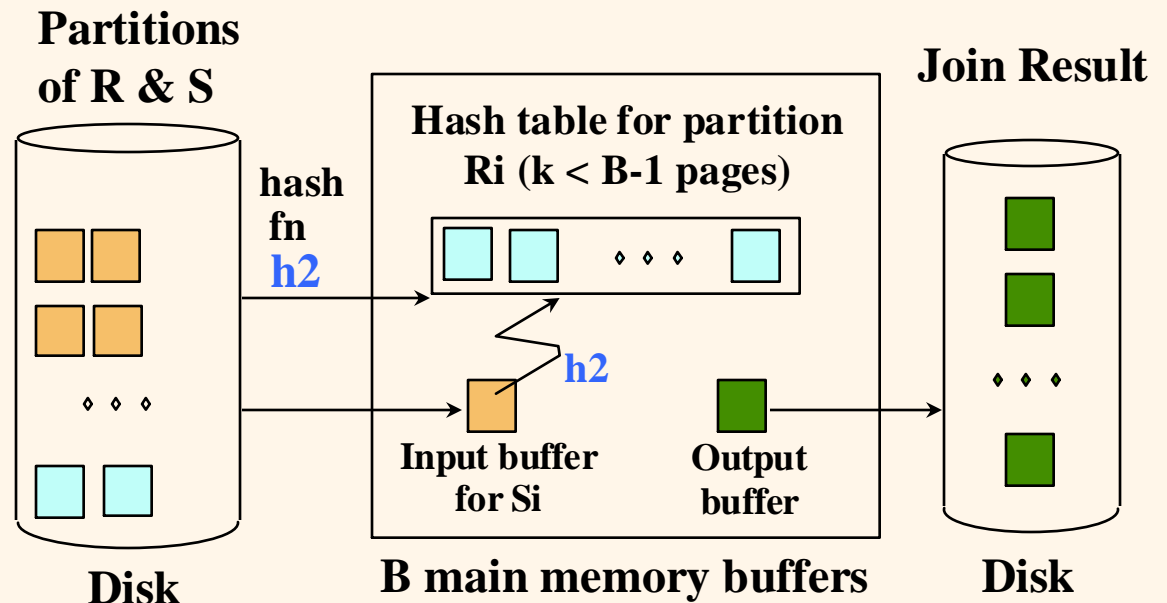
- ❖ Start by sorting both R and S on the join attribute:
  - Assumption:  $M \leq B^2$ ,  $N \leq B^2$
  - Cost:  $4M + 4N$
- ❖ Read both relations in sorted order, match tuples
  - Typical cost:  $M + N$  (but can be as high as  $M*N$ )
- ❖ Difficulty: many tuples in R may match many in S
  - If at least one set of tuples fits in memory, we are OK
  - Otherwise need nested loop, higher cost
- ❖ Total cost:  $5M + 5N$

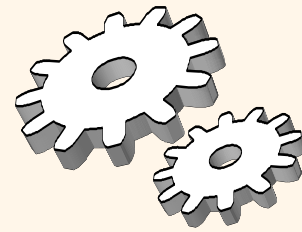
# Hash-Join

❖ Partition both relations using hash fn **h**:  $R$  tuples in partition  $i$  will only match  $S$  tuples in partition  $i$ .



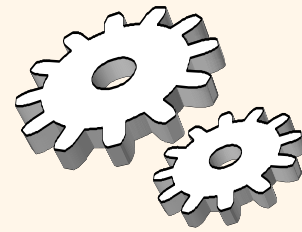
❖ Read in a partition of  $R$ , hash it using **h2** ( $\neq \mathbf{h}$ !). Scan matching partition of  $S$ , search for matches.





# *Observations on Hash-Join*

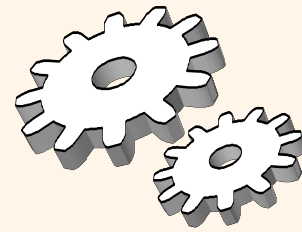
- ❖ #partitions  $k < B-1$  (why?), and  $B-2 > \text{size of largest partition}$  to be held in memory. Assuming uniformly sized partitions, and maximizing  $k$ , we get:
  - $k = B-1$ , and  $M/(B-1) < B-2$ , i.e.,  $B$  must be  $> \sqrt{M}$
- ❖ If we build an in-memory hash table to speed up the matching of tuples, a little more memory is needed.
- ❖ If the hash function does not partition uniformly, one or more  $R$  partitions may not fit in memory. Can apply hash-join technique recursively to do the join of this  $R$ -partition with corresponding  $S$ -partition.



# *Cost of Hash-Join*

- ❖ In partitioning phase, read+write both relns;  $2(M+N)$ .  
In matching phase, read both relns;  $M+N$  I/Os.
- ❖ In our running example, this is a total of 4500 I/Os.



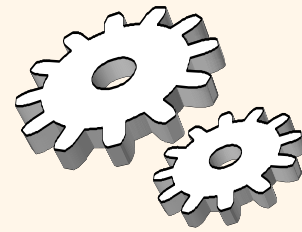


# *Index Nested Loops Join*

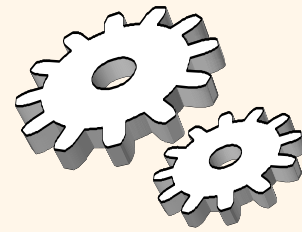
```
foreach tuple r in R do
    foreach tuple s in S where  $r_i == s_j$  do
        add  $\langle r, s \rangle$  to result
```

- ❖ If there is an index on the join column of one relation (say S), can make it the inner and exploit the index.
  - Cost:  $M + (M * p_R) * \text{cost of finding matching S tuples}$
- ❖ For each R tuple, cost of probing S index is about 1.2 for hash index, 2-4 for B+ tree.

# General Join Conditions [See Textbook]



- ❖ Equalities over several attributes (e.g.,  $R.sid = S.sid$  AND  $R.rname = S.sname$ ):
  - For Index NL, build index on  $\langle sid, sname \rangle$  (if S is inner); or use existing indexes on  $sid$  or  $sname$ .
  - For Sort-Merge and Hash Join, sort/partition on combination of the two join columns.
- ❖ Inequality conditions (e.g.,  $R.rname < S.sname$ ):
  - For Index NL, need (clustered!) B+ tree index.
    - Range probes on inner; # matches likely to be much higher than for equality joins.
  - Hash Join, Sort Merge Join not applicable.
  - Block NL quite likely to be the best join method here.



# *Key Things to Remember*

- ❖ Joins are very common, and very expensive
- ❖ Need to bring similar tuples together
- ❖ Can do this in a few ways
  - sort
  - hash
  - using an index
  - or just enumerate all pairs (nested loop joins)