

## Review (1)

- We would like to sort the tuples of a relation R on a given key. The following is known about the relation: R contains 100,000 tuples. The size of a page on disk is 4000 bytes. The size of each R tuple is 400 bytes. R is clustered, i.e., each disk page holding R tuples is full of R tuples. The size of the sort key is 32 bytes. A record pointer is 8 bytes. Answer the following questions:
  - If we use a two pass sorting algorithm, what is the minimum amount of main memory (in terms of number of pages) required?
  - What is the cost of the two pass sorting algorithm in terms of number of disk I/Os? Include the cost of writing the sorted file to disk.
  - Consider the following variant of the sorting algorithm. Instead of sorting the entire tuple, we just sort the (key, recordPointer) for each tuple. As in the conventional two pass sorting algorithm, we sort chunks of (key, recordPointer) in main memory and write the chunks to the tuple (from the original copy of R) and write the sorted relation to disk. What is the minimum amount of main memory required for this operation? What is the cost in terms of number of disk I/Os?
  - Keeping all other parameters constant, for what values of tuple size is the variant discussed above better (in the number of I/Os)?

## Review (2)

- $\sqrt{|R|} + 1 = 101$ , where  $|R|$  denotes the size of R in pages
- $2 \times 2 \times |R| = 40000$
- Memory required = 34 (an additional page is needed for the random access step in the second phase)
- This is an optimized version. The I/Os of the sorting scheme is 122000. This includes 10000 for initially reading R and constructing (key, recordPointer) pairs; 1000 I/Os for writing the sorted runs of (key, recordPointer) pairs to disk; 1000 for reading the same from disk to merge the runs; 100000 I/Os for random access to retrieve the tuples pointed by the record pointer; and finally 10000 I/Os to write the sorted relation R to disk
- Assume that records are unspanned, then tuplesize > 2001

# Relational Operators

First comes thought; then organization of that thought, into ideas and plans; then transformation of those plans into reality. The beginning, as you will observe, is in your imagination.

Napolean Hill

CS5208

3

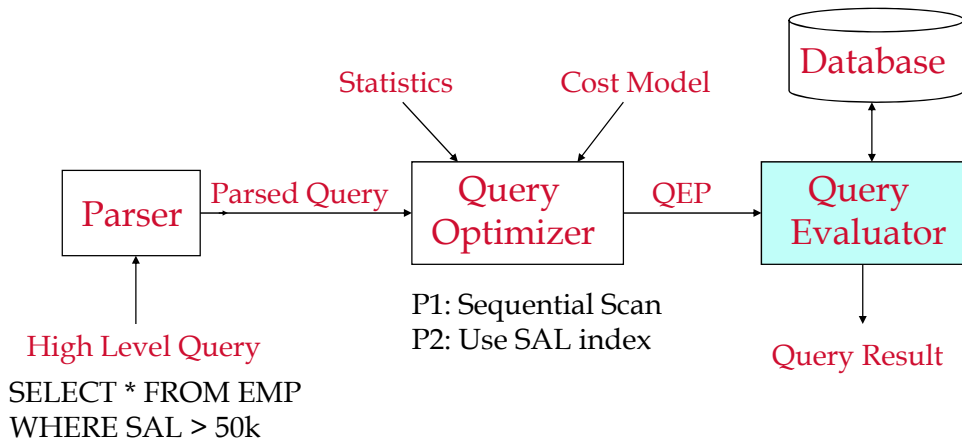
## Introduction

- We've covered the basic underlying storage, buffering, and indexing technology.
  - Now we can move on to query processing.
- Some database operations are EXPENSIVE
- Can greatly improve performance by being “smart”
  - e.g., can speed up 1,000,000x over naïve approach
- Main weapons are:
  - clever implementation techniques for operators
  - exploiting “equivalences” of relational operators
  - using statistics and cost models to choose among these.

CS5208

4

## Steps of processing a high-level query



CS5208

5

## 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

Since each op returns a relation, ops can be *composed*!

Queries that require multiple ops to be composed may be composed in different ways - thus *optimization* is necessary for good performance

CS5208

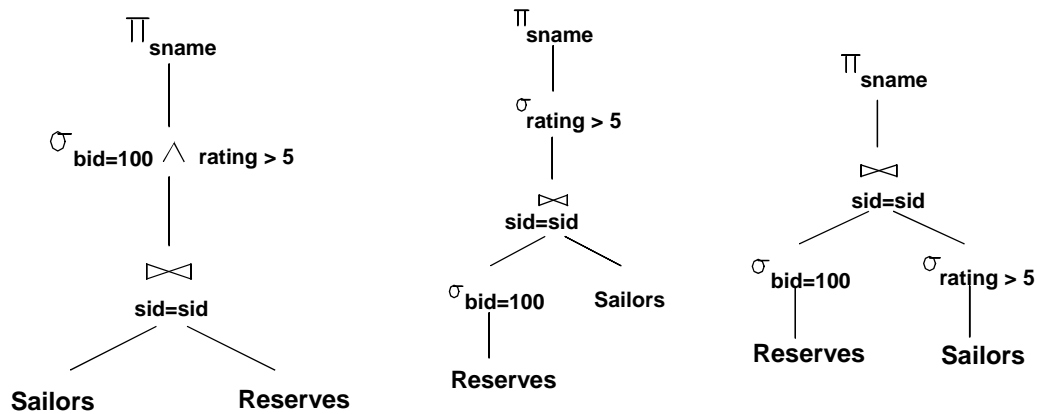
6

```

SELECT S.sname
FROM Reserves R, Sailors S
WHERE R.sid=S.sid AND
      R.bid=100 AND S.rating>5

```

## Example



CS5208

## Paradigm

- Cross product
- Index
  - B+-tree, Hash
    - assume index entries to be (rid,pointer) pair
  - Clustered, Unclustered
- Sort
- Hash

CS5208

8

## Schema for Examples

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

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

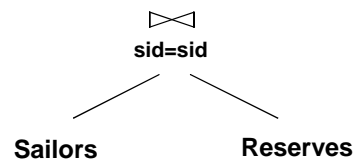
- Reserves (R):
  - $p_R$  tuples per page, M pages.  $p_R = 100$ .  $M = 1000$ .
- Sailors (S):
  - $p_S$  tuples per page, N pages.  $p_S = 80$ .  $N = 500$ .
- **Cost metric:** # of I/Os (pages)
  - We will ignore output costs in the following discussion.

CS5208

9

## Equality Joins With One Join Column

```
SELECT *  
FROM   Reserves R, Sailors S  
WHERE  R.sid=S.sid
```



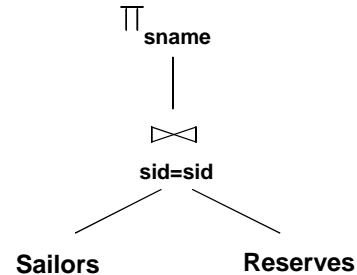
- In algebra:  $R \bowtie S$ .
- Most frequently used operation; very costly operation.
- $\text{join\_selectivity} = \text{join\_size} / (\#R \text{ tuples} \times \#S \text{ tuples})$

CS5208

10

## Equality Joins With One Join Column

```
SELECT sname
FROM   Reserves R, Sailors S
WHERE  R.sid=S.sid
```



- In algebra:  $R \bowtie S$ .
- Most frequently used operation; very costly operation.
- $join\_selectivity = join\_size / (\#R \text{ tuples} \times \#S \text{ tuples})$

CS5208

11

## Join Example

Sailor

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

Reserve

<u>sid</u>	<u>bid</u>	<u>day</u>	rname
31	101	10/11/96	lubber
58	103	11/12/96	dustin

CS5208

12

## Join Example

Sailor

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

Reserve

<u>sid</u>	<u>bid</u>	<u>day</u>	rname
31	101	10/11/96	lubber
58	103	11/12/96	dustin

Query (join) output

<u>sid</u>	<u>sname</u>	<u>rating</u>	<u>age</u>	<u>bid</u>	<u>day</u>	rname
31	lubber	8	55.5	101	10/11/96	lubber
58	rusty	10	35.0	103	11/12/96	dustin

CS5208

13

## Simple Nested Loops Join

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

- For each tuple in the outer relation R, we scan the entire inner relation S.
  - I/O Cost?
  - Memory?

CS5208

14

## *Simple Nested Loops Join*

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

- 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.

## *Simple Nested Loops Join*

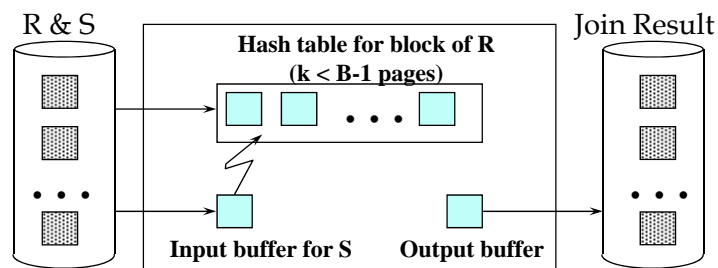
```
foreach tuple r in R do
  foreach tuple s in S do
    if r.sid == s.sid then add <r, s> to result
```

- 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.
  - Memory: 3 pages!



## 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.



CS5208

17

## Examples of Block Nested Loops

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

CS5208

18

## *Examples of Block Nested Loops*

- Cost: Scan of outer + #outer blocks \* scan of inner
  - #outer blocks =  $\lceil \text{no. of pages in outer relation} / \text{block size} \rceil$

## *Examples of Block Nested Loops*

- Cost: Scan of outer + #outer blocks \* scan of inner
  - #outer blocks =  $\lceil \text{no. of pages in outer relation} / \text{block size} \rceil$
- With R as outer, block size of 100 pages:
  - Cost of scanning R is 1000 I/Os; a total of 10 blocks.
  - Per block of R, we scan S;  $10 \times 500$  I/Os.
  - If block size for just 90 pages of R, scan S 12 times.

## *Examples of Block Nested Loops*

- Cost: Scan of outer + #outer blocks \* scan of inner
  - #outer blocks =  $\lceil \text{no. of pages in outer relation} / \text{block size} \rceil$
- With R as outer, block size of 100 pages:
  - Cost of scanning R is 1000 I/Os; a total of 10 blocks.
  - Per block of R, we scan S;  $10 \times 500$  I/Os.
  - If block size for just 90 pages of R, scan S 12 times.
- With 100-page block of S as outer?

CS5208

21

## *Examples of Block Nested Loops*

- Cost: Scan of outer + #outer blocks \* scan of inner
  - #outer blocks =  $\lceil \text{no. of pages in outer relation} / \text{block size} \rceil$
- With R as outer, block size of 100 pages:
  - Cost of scanning R is 1000 I/Os; a total of 10 blocks.
  - Per block of R, we scan S;  $10 \times 500$  I/Os.
  - If block size for just 90 pages of R, scan S 12 times.
- With 100-page block of S as outer:
  - Cost of scanning S is 500 I/Os; a total of 5 blocks.
  - Per block of S, we scan R;  $5 \times 1000$  I/Os.

CS5208

22

## Sort-Merge Join

- 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.)

CS5208

23

## Example of Sort-Merge Join

<u>sid</u>	sname	rating	age	<u>sid</u>	<u>bid</u>	<u>day</u>	rname
22	dustin	7	45.0	28	103	12/4/96	guppy
28	yuppy	9	35.0	28	103	11/3/96	yuppy
31	lubber	8	55.5	31	101	10/10/96	dustin
44	guppy	5	35.0	31	102	10/12/96	lubber
58	rusty	10	35.0	31	101	10/11/96	lubber
				58	103	11/12/96	dustin

- Cost?

CS5208

24

## Example of Sort-Merge Join

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

<u>sid</u>	<u>bid</u>	<u>day</u>	rname
28	103	12/4/96	guppy
28	103	11/3/96	yuppy
31	101	10/10/96	dustin
31	102	10/12/96	lubber
31	101	10/11/96	lubber
58	103	11/12/96	dustin

- Cost:  $2M \cdot K_1 + 2N \cdot K_2 + (M+N)$ 
  - $K_1$  and  $K_2$  are the number of passes to sort R and S respectively
  - The cost of scanning,  $M+N$ , could be  $M \cdot N$  (very unlikely!)

CS5208

25

## Example of Sort-Merge Join

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

<u>sid</u>	<u>bid</u>	<u>day</u>	rname
28	103	12/4/96	guppy
28	103	11/3/96	yuppy
31	101	10/10/96	dustin
31	102	10/12/96	lubber
31	101	10/11/96	lubber
58	103	11/12/96	dustin

- Cost:  $2M \cdot K_1 + 2N \cdot K_2 + (M+N)$ 
  - $K_1$  and  $K_2$  are the number of passes to sort R and S respectively
  - The cost of scanning,  $M+N$ , could be  $M \cdot N$  (very unlikely!)
- With 35, 100 or 300 buffer pages, both R and S can be sorted in 2 passes; total join cost: 7500.  
(BNL cost: 2500 to 15000 I/Os)

CS5208

26

## GRACE Hash-Join

bucketID =  $X \bmod 4$

		S			
		0	1	2	3
R	0	XXX XXX XXX			
	1		XXX XXX XXX		
	2			XXX XXX XXX	
	3				XXX XXX XXX

CS5208

27

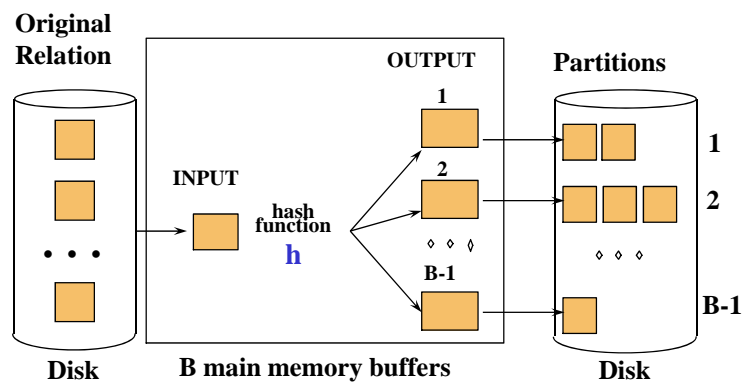
## GRACE Hash-Join

- Operates in two phases:
  - Partition phase
    - Partition relation R using hash fn  $h$ .
    - Partition relation S using hash fn  $h$ .
    - R tuples in partition i will only match S tuples in partition i.
  - Join phase
    - Read in a partition of R
    - Hash it using  $h_2 (<> h!)$
    - Scan matching partition of S, search for matches.

CS5208

28

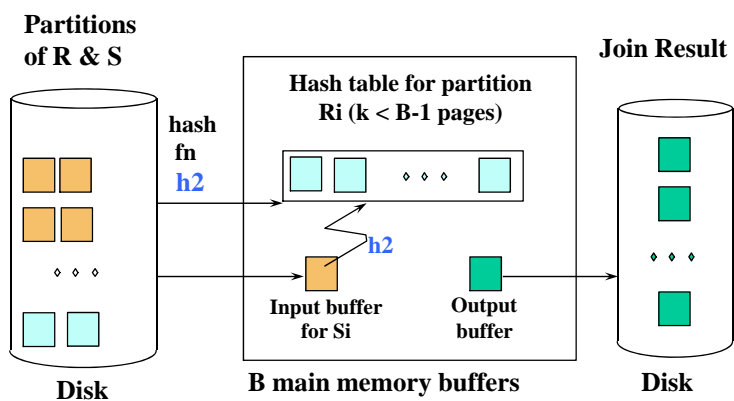
## Partitioning Phase



CS5208

29

## Joining Phase



CS5208

30

## *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.

CS5208

31

## *Observations on Hash-Join*

- #partitions  $k \leq B-1$  (why?), and  $B-2 \geq \text{size of largest partition}$  to be held in memory. Assuming uniformly sized partitions, and maximizing  $k$ , we get:
  - $k = B-1$ , and  $M/(B-1) \leq B-2$ , i.e.,  $B$  must be  $\geq \sqrt{M}$
- 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.
- What if  $B < \sqrt{M}$ ?

CS5208

32



## *Index Nested Loops Join*

```
foreach tuple r in R do
    search index of S on sid using  $S_{\text{search-key}} = r.\text{sid}$ 
    for each matching key
        retrieve s; 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. Cost of then finding S tuples (assuming leaf data entries are pointers) depends on clustering.
  - Clustered index: 1 I/O (typical), unclustered: upto 1 I/O per matching S tuple.

CS5208

33

## *Schema for Examples*

Sailors (sid: integer, sname: string, rating: integer, age: real)  
Reserves (sid: integer, bid: integer, day: dates, rname: string)

- Reserves (R):
  - $p_R$  tuples per page, M pages.  $p_R = 100$ .  $M = 1000$ .
- Sailors (S):
  - $p_S$  tuples per page, N pages.  $p_S = 80$ .  $N = 500$ .
- Cost metric: # of I/Os (pages)
  - We will ignore output costs in the following discussion.

CS5208

34

## *Examples of Index Nested Loops*

- Hash-index on sid of S (as inner):
  - Scan R: 1000 page I/Os, 100\*1000 tuples.
  - For each R tuple: 1.2 I/Os to get data entry in index, plus 1 I/O to get (the exactly one) matching S tuple. Total: 220,000 I/Os.
- Hash-index on sid of R (as inner)?

CS5208

35

## *Examples of Index Nested Loops*

- Hash-index on sid of S (as inner):
  - Scan R: 1000 page I/Os, 100\*1000 tuples.
  - For each R tuple: 1.2 I/Os to get data entry in index, plus 1 I/O to get (the exactly one) matching S tuple. Total: 220,000 I/Os.
- Hash-index on sid of R (as inner):
  - Scan S: 500 page I/Os, 80\*500 tuples.
  - For each S tuple: 1.2 I/Os to find index page with data entries, plus cost of retrieving matching R tuples.
  - Assuming uniform distribution, 2.5 reservations per sailor (100,000 / 40,000). Cost of retrieving them is 1 or 2.5 I/Os depending on whether the index is clustered.

CS5208

36

## General Join Conditions

- Equalities over several attributes (e.g.,  $R.sid = S.sid$  AND  $R.name = S.name$ ):
  - Join on one predicate, and treat the rest as selections;
  - For Index NL, build index on  $\langle sid, name \rangle$  (if S is inner); use existing indexes on  $sid$  or  $name$ .
  - For Sort-Merge and Hash Join, sort/partition on combination of the two join columns
- Inequality join ( $R.sid < S.sid$ )?

CS5208

37

## 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

Since each op returns a relation, ops can be *composed*!

Queries that require multiple ops to be composed may be composed in different ways - thus *optimization* is necessary for good performance

CS5208

38

## Simple Selections

- Of the form:  $\sigma_{R.attr \text{ op } value}(R)$
- **selectivity** = Size of result / **R**
- **With no index, unsorted**: Must essentially scan the whole relation; **cost is M** (#pages in R).
- Sorted?
- **With an index on selection attribute**: Use index to find qualifying data entries, then retrieve corresponding data records. (Hash index useful only for equality selections.)

```
SELECT *  
FROM   Reserves R  
WHERE  R.rname < 'C%'
```

CS5208

39

## Using an Index for Selections

- Cost depends on #qualifying tuples, and clustering.
  - Cost of finding qualifying data entries (typically small) plus cost of retrieving records (could be large w/o clustering).
  - In example, assuming uniform distribution of names, about 10% of tuples qualify (100 pages, 10000 tuples).
    - Clustered index?
    - Unclustered index?

CS5208

40

## *Using an Index for Selections*

- Cost depends on #qualifying tuples, and clustering.
  - Cost of finding qualifying data entries (typically small) plus cost of retrieving records (could be large w/o clustering).
  - In example, assuming uniform distribution of names, about 10% of tuples qualify (100 pages, 10000 tuples).
    - Clustered index: ~ 100 I/Os
    - Unclustered: upto 10000 I/Os!

CS5208

41

## *Two Approaches to General Selections*

- First approach: Find the **most selective access path**, retrieve tuples using it, and apply any remaining terms that don't **match** the index:
  - **Most selective access path**: An index or file scan that we estimate will require the fewest page I/Os.
  - Terms that match this index reduce the number of tuples retrieved; other terms are used to discard some retrieved tuples, but do not affect number of tuples/pages fetched.
  - Consider **day<8/9/94 AND bid=5 AND sid=3**. A B+ tree index on day can be used; then, bid=5 and sid=3 must be checked for each retrieved tuple. Similarly, a hash index on <bid, sid> could be used; day<8/9/94 must then be checked.

CS5208

42

## *Intersection of Rids*

- Second approach (if we have 2 or more matching indexes (assuming leaf data entries are pointers):
  - Get sets of rids of data records using each matching index.
  - Then **intersect** these **sets of rids** (we'll discuss intersection soon!)
  - Retrieve the records and apply any remaining terms.
  - Consider **day<8/9/94 AND bid=5 AND sid=3**. If we have a B+ tree index on day and an index on sid, we can retrieve rids of records satisfying day<8/9/94 using the first, rids of recs satisfying sid=3 using the second, intersect, retrieve records and check bid=5.

CS5208

43

```
SELECT  DISTINCT  
        R.sid, R.bid  
FROM    Reserves R
```

## *The Projection Operation (Duplicate Elimination)*

- An approach based on sorting:
  - **Modify Pass 0 of external sort to eliminate unwanted fields.** Thus, runs are produced, but tuples in runs are smaller than input tuples. (Size ratio depends on # and size of fields that are dropped.)
  - **Modify merging passes to eliminate duplicates.** Thus, number of result tuples smaller than input. (Difference depends on # of duplicates.)
  - **Cost:** In Pass 0, read original relation (size M), write out same number of smaller tuples. In merging passes, fewer tuples written out in each pass.
- Hash-based scheme?

CS5208

44

## *Set Operations*

- Intersection and cross-product special cases of join.
- Union (Distinct) and Difference similar.
- Sorting based approach to union:
  - Sort both relations (on combination of all attributes).
  - Scan sorted relations and merge them.
- Hash based approach to union?

CS5208

45

## *Set Operations*

- Intersection and cross-product special cases of join.
- Union (Distinct) and Difference similar.
- Sorting based approach to union:
  - Sort both relations (on combination of all attributes).
  - Scan sorted relations and merge them.
- Hash based approach to union:
  - Partition R and S using hash function h.
  - For each S-partition, build in-memory hash table (using h2), scan corr. R-partition and add tuples to table while discarding duplicates.

CS5208

46

## Aggregate Operations (AVG, MIN, etc.)

- Without grouping:

- In general, requires scanning the relation.
- Given index whose search key includes all attributes in the SELECT or WHERE clauses, can do index-only scan.

```
SELECT AVG(SALARY)
FROM EMPLOYEE
```

- With grouping:

- Sort on group-by attributes, then scan relation and compute aggregate for each group.
- Similar approach based on hashing on group-by attributes.
- Given tree index whose search key includes all attributes in SELECT, WHERE and GROUP BY clauses, can do index-only scan; if group-by attributes form prefix of search key, can retrieve data entries/tuples in group-by order.

```
SELECT DEPT, AVG(SALARY)
FROM EMPLOYEE
GROUP BY DEPT
```

CS5208

47

## Iterators for Implementation of Operators

- Most operators can be implemented as an *iterator*
- An iterator allows a **consumer of the result** of the operator to get the result **one tuple at a time**
  - Open – starts the process of getting tuples, but does not get a tuple. It initializes any data structures needed.
  - GetNext – **returns the next tuple in the result** and adjusts the data structures as necessary to allow subsequent tuples to be obtained. It may call GetNext one or more times on its arguments. It also signals whether a tuple was produced or there were no more tuples to be produced.
  - Close – ends the iteration after all tuples have been obtained.

CS5208

48



## *Iterators*

```
Open();  
While condition is true do {  
    GetNext();  
    perform other operations  
}  
Close();
```

CS5208

49

## *More on Iterators*

- Why iterators?
  - Do not need to *materialize* (i.e., store on disk) intermediate results
  - Many operators are active at once, and tuples flow from one operator to the next, thus reducing the need to store intermediate results
- In some cases (e.g., sort), almost all the work would need to be done by the Open function, which is tantamount to materialization
- We shall regard Open, GetNext, Close as overloaded names of methods.
  - Assume that for each physical operator, there is a class whose objects are the relations that can be produced by this operator. If R is a member of such a class, then we use R.Open(), R.GetNext, and R.Close() to apply the functions of the iterator for R.

CS5208

50

## *An iterator for table-scan operator*

```

Open(R) {
    b := first block of R;
    t := first tuple of block b;
    Found := TRUE;
}

Close(R) {
}

GetNext(R) {
    If (t is past the last tuple on b)
        b := next block
        If (there is no next block)
            Found := FALSE;
            RETURN;
        Else
            t := first tuple in b;
    oldt := t;
    t := next tuple of b
    RETURN oldt;
}

```

CS5208

51

## *An iterator for tuple-based nested-loops join operator (assumes R and S are non-empty)*

```

Open(R,S) {
    R.Open();
    S.Open();
    s := S.GetNext();
}

Close(R,S) {
    R.Close();
    S.Close();
}

GetNext(R,S) {
    REPEAT
        r := R.GetNext();
        If (NOT Found) {
            R.Close();
            s := S.GetNext();
            IF (NOT Found)
                Return;
            R.Open();
            r := R.GetNext();
        }
    UNTIL (r and s join);
    Return the join of r and s;
}

```

CS5208

52

## *Summary*

- A virtue of relational DBMSs: queries are composed of a few basic operators; the implementation of these operators can be carefully tuned (and it is important to do this!).
- Many alternative implementation techniques for each operator; no universally superior technique for most operators.
- Must consider available alternatives for each operation in a query and choose best one based on system statistics, etc. This is part of the broader task of optimizing a query composed of several ops.