# Evaluation of relational operations

Iztok Savnik, FAMNIT

# Slides & Textbook

- Textbook:
  - Raghu Ramakrishnan, Johannes Gehrke, *Database Management Systems, McGraw-Hill, 3ʳᵈ ed., 2007.*
- *Slides:*
  - *From „Cow Book":  R.Ramakrishnan, http://pages.cs.wisc.edu/~dbbook/*

# Overview of Query Evaluation

- *Plan*:  *Tree of R.A. ops, with choice of alg for each op.*
  - Each operator typically implemented using a `pull' interface: when an operator is `pulled' for the next output tuples, it `pulls' on its inputs and computes them.
- Two main issues in query optimization:
  - For a given query, what plans are considered?
    - Algorithm to search plan space for cheapest (estimated) plan.
  - How is the cost of a plan estimated?
- Ideally: Want to find best plan.  Practically: Avoid worst plans!
- We will study the System R approach.

# Some Common Techniques

- Algorithms for evaluating relational operators use some simple ideas extensively:
  - Indexing:  Can use WHERE conditions to retrieve small set of tuples (selections, joins)
  - Iteration:  Sometimes, faster to scan all tuples even if there is an index. (And sometimes, we can scan the data entries in an index instead of the table itself.)
  - Sorting: Many algorithms for the evaluation of relational operations evolved from the external merge sort algorithm: project, join, grouping, etc.
  - Partitioning: By hashing, we can partition the input tuples and replace an expensive operation by similar operations on smaller inputs.

  *Watch for these techniques as we discuss query evaluation!*

# Statistics and Catalogs

- Need information about the relations and indexes involved.  ***Catalogs*** typically contain at least:
    - # tuples (NTuples) and # pages (NPages) for each relation.
    - # distinct key values (NKeys) and NPages for each index.
    - Index height, low/high key values (Low/High) for each tree index.
- Catalogs updated periodically.
    - Updating whenever data changes is too expensive; lots of approximation anyway, so slight inconsistency ok.
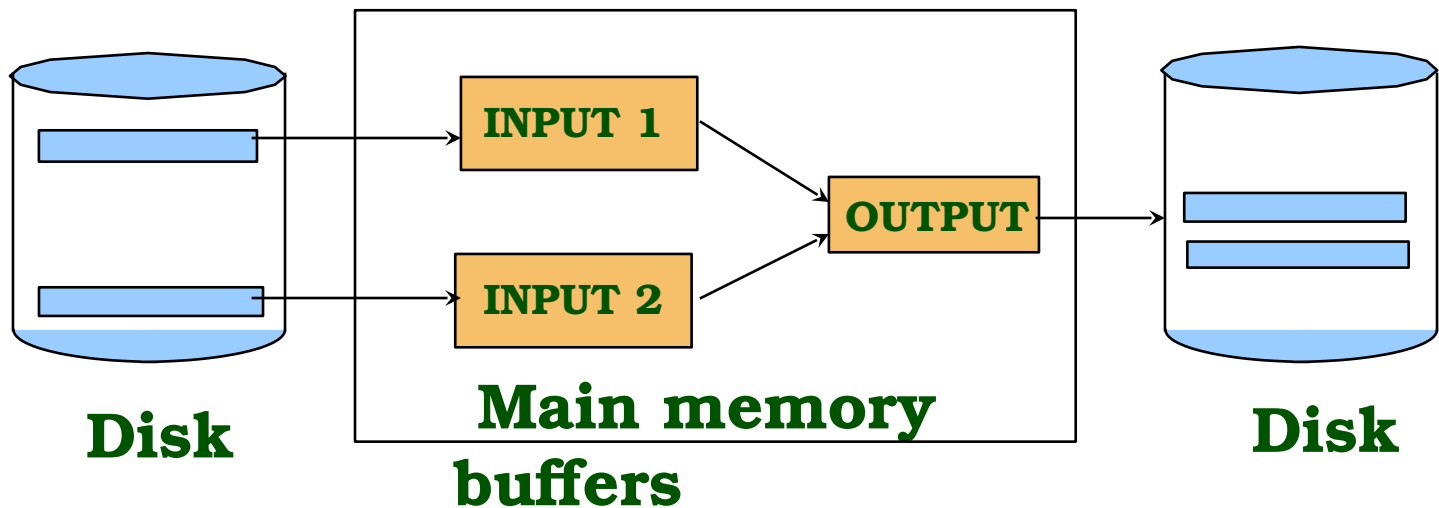- More detailed information (e.g., histograms of the values in some field) are sometimes stored.

IDB, Eval RA

# Access Paths

❖ An <u>access path</u> is a method of retrieving tuples:
  ▪ File scan, or index that matches a selection (in the query)

❖ A tree index _matches_ (a conjunction of) terms that involve only attributes in a _prefix_ of the search key.
  ▪ E.g., Tree index on _<a, b, c>_ matches the selection _a=5 AND b=3_, and _a=5 AND b>6_, but not _b=3_.

❖ A hash index _matches_ (a conjunction of) terms that has a term _attribute = value_ for every attribute in the search key of the index.
  ▪ E.g., Hash index on _<a, b, c>_ matches _a=5 AND b=3 AND c=5_; but it does not match _b=3, or a=5 AND b=3, or a>5 AND b=3 AND c=5_.

# Why Sort?

- A classic problem in computer science!
- Data requested in sorted order
  - e.g., find students in increasing *gpa* order
- Sorting is first step in *bulk loading* B+ tree index.
- Sorting useful for eliminating *duplicate copies* in a collection of records (Why?)
- *Sort-merge* join algorithm involves sorting.
- Problem: sort 1Gb of data with 1Mb of RAM.
  - why not virtual memory?

# 2-Way Sort: Requires 3 Buffers

- Pass 1: Read a page, sort it, write it.
  - only one buffer page is used
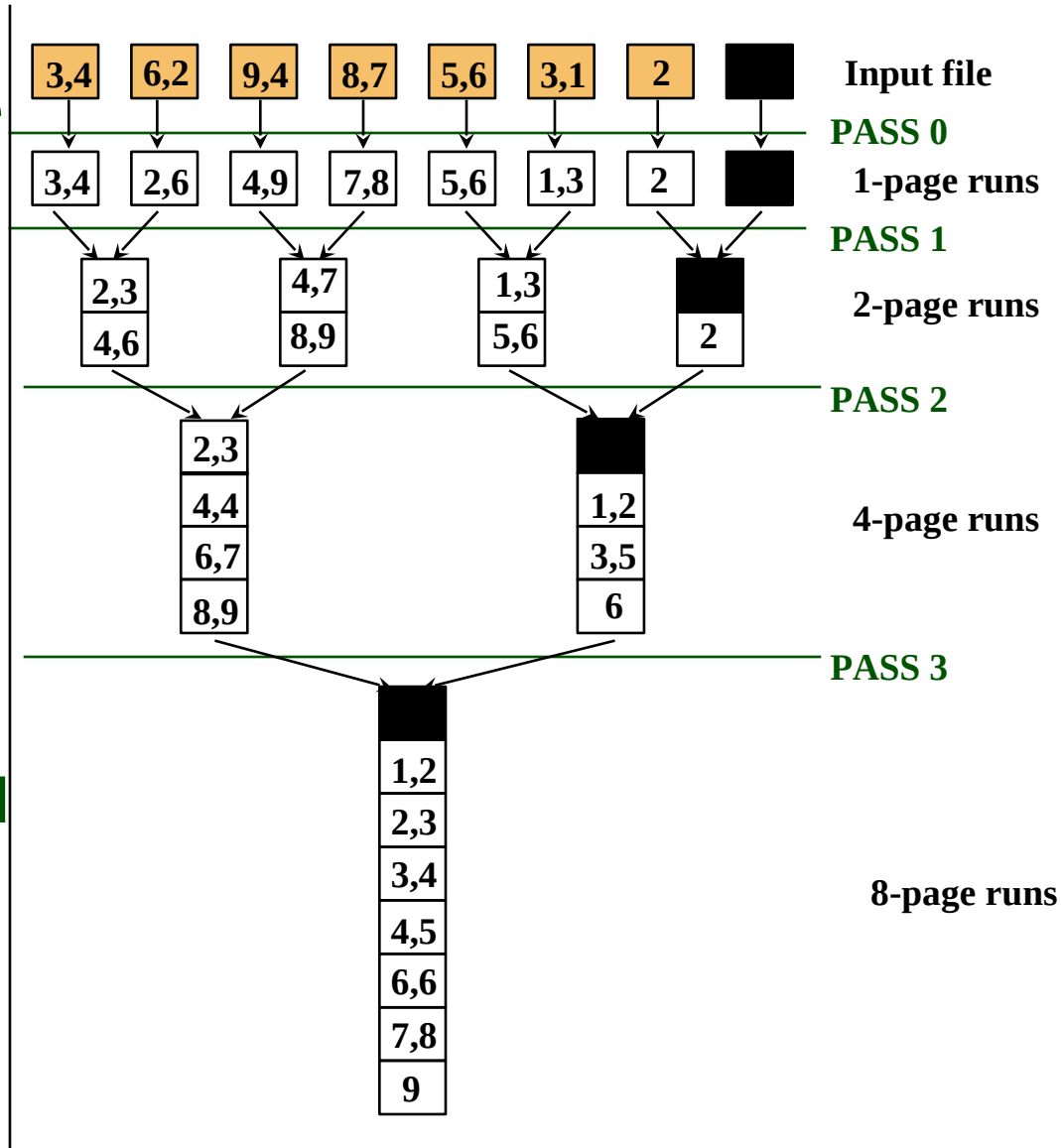- Pass 2, 3, …, etc.:
  - three buffer pages used.



**Disk**     **Main memory buffers**     **Disk**

INPUT 1

INPUT 2

OUTPUT

IDB, Eval RA

# Two-Way External Merge Sort

- Each pass we read + write each page in file.

- N pages in the file => the number of passes
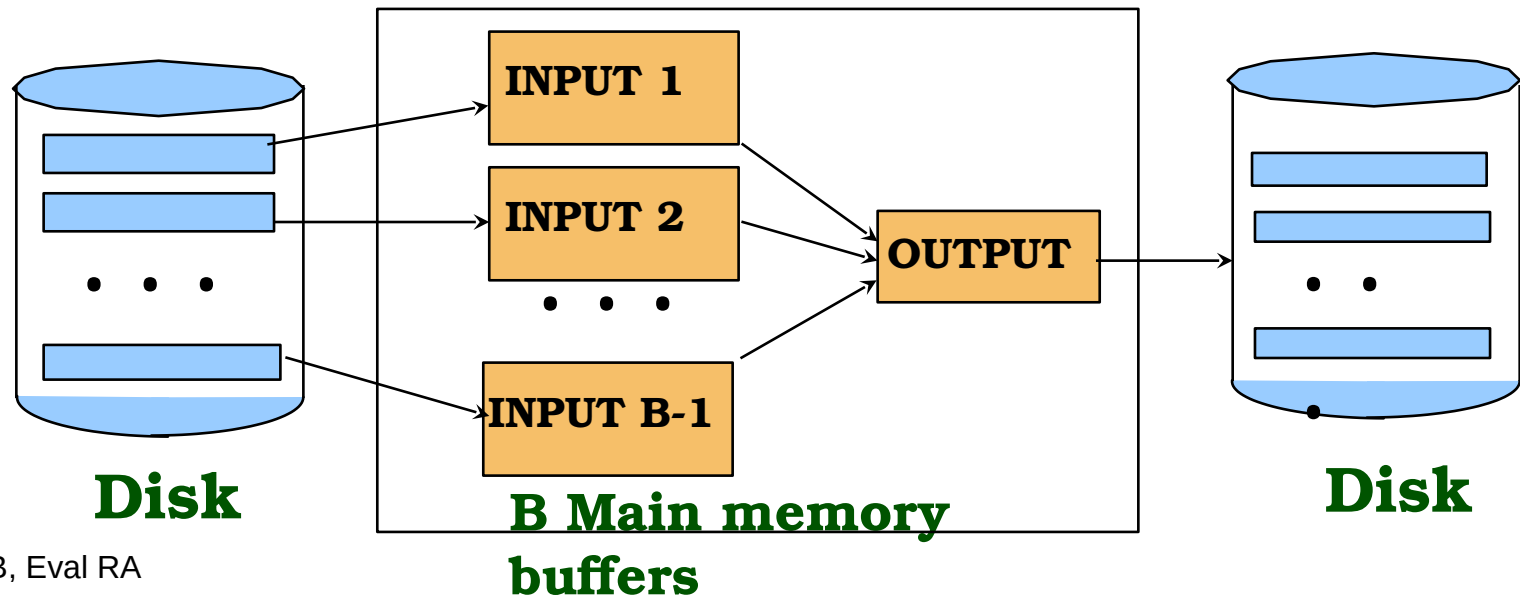
$$= \lceil \log_2 N \rceil + 1$$

- So total cost is:

$$2N \left( \lceil \log_2 N \rceil + 1 \right)$$

- *Idea:* **Divide and conquer:** sort subfiles and merge



| 3,4 | 6,2 | 9,4 | 8,7 | 5,6 | 3,1 | 2 | ■ | **Input file** |

**PASS 0**

| 3,4 | 2,6 | 4,9 | 7,8 | 5,6 | 1,3 | 2 | ■ | **1-page runs** |

**PASS 1**

2-page runs:
| 2,3 | | 4,7 | | 1,3 | | ■ |
| 4,6 | | 8,9 | | 5,6 | | 2 |

**2-page runs**

**PASS 2**

4-page runs:
| 2,3 | | ■ |
| 4,4 | | 1,2 |
| 6,7 | | 3,5 |
| 8,9 | | 6 |

**4-page runs**

**PASS 3**

8-page runs:
| ■ |
| 1,2 |
| 2,3 |
| 3,4 |
| 4,5 |
| 6,6 |
| 7,8 |
| 9 |

**8-page runs**

# General External Merge Sort

☛ *More than 3 buffer pages.  How can we utilize them?*

- To sort a file with *N* pages using *B* buffer pages:
  - Pass 0: use *B* buffer pages. Produce $\lceil N/B \rceil$ sorted runs of *B* pages each.
  - Pass 2, …,  etc.: merge *B-1* runs.



**Disk**          **B Main memory buffers**          **Disk**

# Cost of External Merge Sort

- Number of passes: $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$
- <span style="color:red">Cost = 2N * (# of passes)</span>
- E.g., with 5 buffer pages, to sort 108 page file:
  - Pass 0: $\lceil 108 / 5 \rceil$ = 22 sorted runs of 5 pages each (last run is only 3 pages)
  - Pass 1: $\lceil 22 / 4 \rceil$ = 6 sorted runs of 20 pages each (last run is only 8 pages)
  - Pass 2:  2 sorted runs, 80 pages and 28 pages
  - Pass 3:  Sorted file of 108 pages

# Number of Passes of External Sort

| N | B=3 | B=5 | B=9 | B=17 | B=129 | B=257 |
|---|---|---|---|---|---|---|
| 100 | 7 | 4 | 3 | 2 | 1 | 1 |
| 1,000 | 10 | 5 | 4 | 3 | 2 | 2 |
| 10,000 | 13 | 7 | 5 | 4 | 2 | 2 |
| 100,000 | 17 | 9 | 6 | 5 | 3 | 3 |
| 1,000,000 | 20 | 10 | 7 | 5 | 3 | 3 |
| 10,000,000 | 23 | 12 | 8 | 6 | 4 | 3 |
| 100,000,000 | 26 | 14 | 9 | 7 | 4 | 4 |
| 1,000,000,000 | 30 | 15 | 10 | 8 | 5 | 4 |

# Relational Operations

- We will consider how to implement:
  - *Selection* ( σ )   Selects a subset of rows from relation.
  - *Projection* ( π )   Deletes unwanted columns from relation.
  - *Join* ( ⋈ )  Allows us to combine two relations.
  - *Set-difference* ( / )  Tuples in reln. 1, but not in reln. 2.
  - *Union* ( ∪ )  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*!  After we cover the operations, we will discuss how to *optimize* queries formed by composing them.

# Schema for Examples

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

- Similar to old schema; *rname* added for variations.
- 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.

# A Note on Complex Selections

> *(day<8/9/94 AND rname='Paul') OR bid=5 OR sid=3*

- Selection conditions are first converted to *conjunctive normal form* (CNF):

  *(day<8/9/94 OR bid=5 OR sid=3 ) AND (rname='Paul' OR bid=5 OR sid=3)*

- We only discuss case with no ORs; see text if you are curious about the general case.

# Using an Index for Selections

- Cost depends on #qualifying tuples, and clustering.

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

  - 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). With a clustered index, cost is little more than 100 I/Os; if unclustered, upto 10000 I/Os!

- *Important refinement for unclustered indexes*:
  1. Find qualifying data entries.
  2. Sort the rid's of the data records to be retrieved.
  3. Fetch rids in order. This ensures that each data page is looked at just once (though # of such pages likely to be higher than with clustering).

IDB, Eval RA

# Two Approaches to General Selections

- <u>First approach:</u> 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.

# Intersection of Rids

- <u>Second approach</u> (if we have 2 or more matching indexes that use Alternatives (2) or (3) for data entries):
  - 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*, both using Alternative (2), 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.*

# Projection

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

- The expensive part is removing duplicates.
  - SQL systems don't remove duplicates unless the keyword DISTINCT is specified in a query.
- Sorting Approach: Sort on <sid, bid> and remove duplicates. (Can optimize this by dropping unwanted information while sorting.)
- Hashing Approach: Hash on <sid, bid> to create partitions. Load partitions into memory one at a time, build in-memory hash structure, and eliminate duplicates.
- If there is an index with both R.sid and R.bid in the search key, may be cheaper to sort data entries!

IDB, Eval RA

# The Projection Operation

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

- An approach based on sorting:
  - Modify Pass 0 of external sort to eliminate unwanted fields. Thus, runs of about B (2B with opt.) pages 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.  Using Reserves example, 1000 input pages reduced to 250 in Pass 0 if size ratio is 0.25

IDB, Eval RA

# Projection Based on Hashing

- *Partitioning phase*:  Read R using one input buffer. For each tuple, discard unwanted fields, apply hash function *h1* to choose one of B-1 output buffers.

  - Result is B-1 partitions (of tuples with no unwanted fields). 2 tuples from different partitions guaranteed to be distinct.

- *Duplicate elimination phase*:  For each partition, read it and build an in-memory hash table, using hash fn *h2* (*<> h1*) on all fields, while discarding duplicates.

  - If partition does not fit in memory, can apply hash-based projection algorithm recursively to this partition.

- Cost:  For partitioning, read R, write out each tuple, but with fewer fields.  This is read in next phase.

# Discussion of Projection

- Sort-based approach is the standard; better handling of skew and result is sorted.

- If an index on the relation contains all wanted attributes in its search key, can do *index-only* scan.
    - Apply projection techniques to data entries (much smaller!)

- If an ordered (i.e., tree) index contains all wanted attributes as *prefix* of search key, can do even better:
    - Retrieve data entries in order (index-only scan), discard unwanted fields, compare adjacent tuples to check for duplicates.

# Equality Joins With One Join Column

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

- In algebra: R ⋈ S.  Common!  Must be carefully optimized.  R x S is large; so, R x S followed by a selection is inefficient.

- Assume: M tuples in R, $p_R$ tuples per page, N tuples in S, $p_S$ tuples per page.
  - In our examples, R is Reserves and S is Sailors.

- We will consider more complex join conditions later.

- *Cost metric*:  # of I/Os.  We will ignore output costs.

# Simple Nested Loops Join

foreach tuple r in R do
      foreach tuple s in S do
           if $r_i == s_j$ 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.

- Page-oriented Nested Loops join: For each *page* of R, get each *page* of S, and write out matching pairs of tuples <r, s>, 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

IDB, Eval RA

# Index Nested Loops Join

foreach tuple r in R do
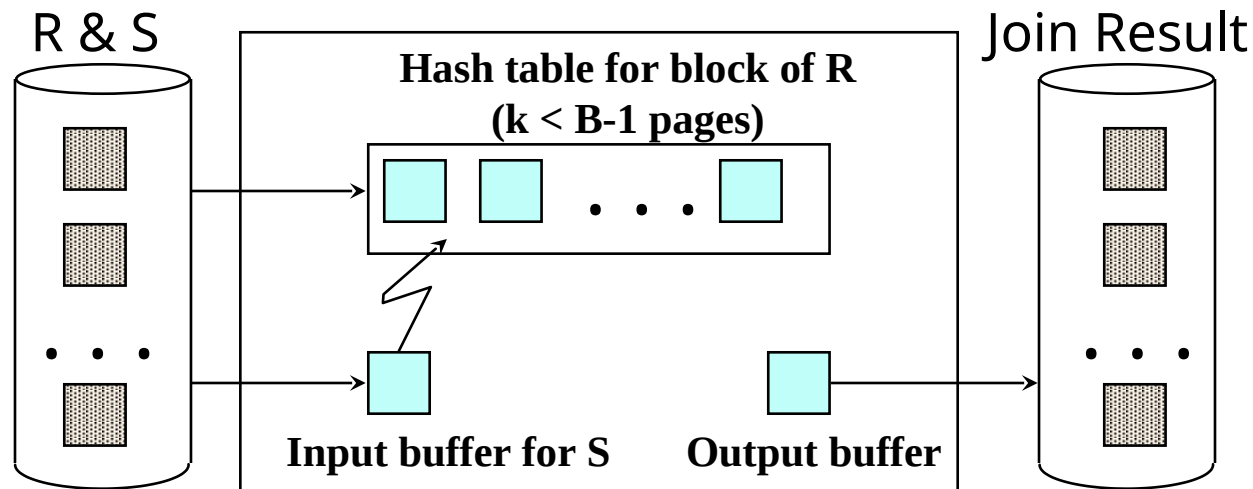  foreach tuple s in S where $r_i == s_j$ do
    add <r, s> to result

- 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*$p_R$) * 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 Alt. (2) or (3) for data entries) depends on clustering.
  - Clustered index:  1 I/O (typical), unclustered: upto 1 I/O per matching S tuple.

IDB, Eval RA

# Examples of Index Nested Loops

- Hash-index (Alt. 2) on *sid* of Sailors (as inner):
  - Scan Reserves:  1000 page I/Os, 100*1000 tuples.
  - For each Reserves tuple:  1.2 I/Os to get data entry in index, plus 1 I/O to get (the exactly one) matching Sailors tuple.  Total:  220,000 I/Os.
- Hash-index (Alt. 2) on *sid* of Reserves (as inner):
  - Scan Sailors:  500 page I/Os, 80*500 tuples.
  - For each Sailors tuple:  1.2 I/Os to find index page with data entries, plus cost of retrieving matching Reserves 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.

# 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 <r, s> to result.  Then read next R-block, scan S, etc.

R & S

Join Result

**Hash table for block of R**

**(k < B-1 pages)**

· · ·

**Input buffer for S**

**Output buffer**

· · ·

IDB, Eval RA

# Examples of Block Nested Loops

- Cost:  Scan of outer +  #outer blocks * scan of inner
  - #outer blocks = $\lceil \# \ of \ pages \ of \ outer \ / \ blocksize \rceil$
- 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.
- With *sequential reads* considered, analysis changes: may be best to divide buffers evenly between R and S.

# Sort-Merge Join   (R ⋈$_{i=j}$ 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 >= current S tuple, then advance scan of S until current S-tuple >= current R tuple; do this until current R tuple = current S tuple.

  - At this point, all R tuples with same value in Ri (*current R group*) and all S tuples with same value in Sj (*current S group*) *match*;  output <r, s> 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

| sid | 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 |

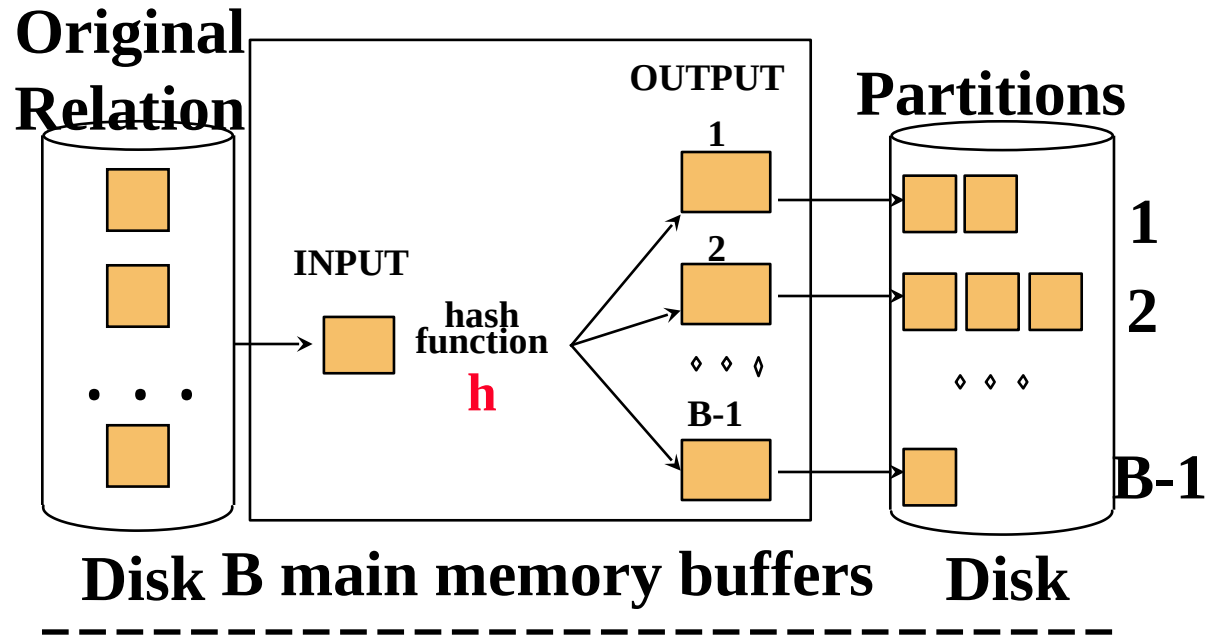| sid | bid | day | rname |
|-----|-----|---------|--------|
| 28 | 103 | 12/4/96 | guppy |
| 28 | 103 | 11/3/96 | yuppy |
| 31 | 101 | 10/10/96 | dustin |
| 31 | 102 | 10/12/96 | lubber |
| 31 | 101 | 10/11/96 | lubber |
| 58 | 103 | 11/12/96 | dustin |

- Cost:  M log M + N log N + (M+N)
  - The cost of scanning, M+N, could be M*N (very unlikely!)
- With 35, 100 or 300 buffer pages, both Reserves and Sailors can be sorted in 2 passes; total join cost: 7500.

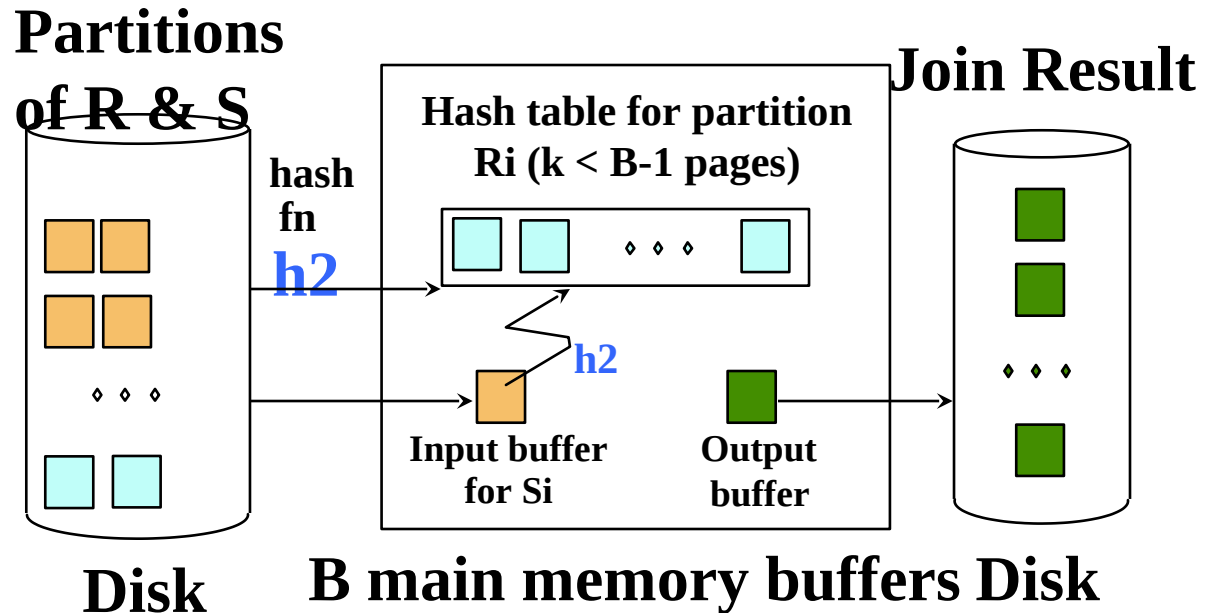*(BNL cost:  2500 to 15000 I/Os)*

# Refinement of Sort-Merge Join

- We can combine the merging phases in the *sorting* of R and S with the merging required for the join.
  - With B > √L, where *L* is the size of the larger relation, using the sorting refinement that produces runs of length B (2B with opt.) in Pass 0, #runs of each relation is < B/2.
  - #runs < √L for each relation. Now suppose B > 2*√L!  Allocate 1 page per run of each relation, and `merge' while checking the join condition.
  - Cost:  read+write each relation in Pass 0 + read each relation in (only) merging pass  (+ writing of result tuples).
  - In example, cost goes down from 7500 to 4500 I/Os.
- In practice, cost of sort-merge join, like the cost of external sorting, is *linear*.

# 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 (<> h!)**. Scan matching partition of S, search for matches.

**Original Relation**

**OUTPUT**

**Partitions**

1

**INPUT**

**hash function**

**h**

2

◊ ◊ ◊

B-1

1

2

◊ ◊ ◊

B-1

**Disk**  **B main memory buffers**  **Disk**

**Partitions of R & S**

**Join Result**

**Hash table for partition Ri (k < B-1 pages)**

**hash fn**

**h2**

**h2**

**Input buffer for Si**

**Output buffer**

**Disk**  **B main memory buffers Disk**

# Observations on Hash-Join

- # partitions k < B-1 (why?), and B-2 > 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.
- Sort-Merge Join vs. Hash Join:
  - Given a minimum amount of memory (*what is this, for each?*) both have a cost of 3(M+N) I/Os.  Hash Join superior on this count if relation sizes differ greatly.  Also, Hash Join shown to be highly parallelizable.
  - Sort-Merge less sensitive to data skew; result is sorted.

# General Join Conditions

- Equalities over several attributes (e.g.,  *R.sid=S.sid* AND *R.rname=S.sname*):
  - For Index NL, build index on *<sid, sname>* (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.

# Set Operations

- Intersection and cross-product special cases of join.
- Union (Distinct) and Except similar; we'll do union.
- Sorting based approach to union:
  - Sort both relations (on combination of all attributes).
  - Scan sorted relations and merge them.
  - *Alternative*:  Merge runs from Pass 0 for *both* relations.
- 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.

# 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.
- With grouping:
  - Sort on group-by attributes, then scan relation and compute aggregate for each group. (Can improve upon this by combining sorting and aggregate computation.)
  - 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.

# Impact of Buffering

- If several operations are executing concurrently, estimating the number of available buffer pages is guesswork.

- Repeated access patterns interact with buffer replacement policy.

  - e.g., Inner relation is scanned repeatedly in Simple Nested Loop Join. With enough buffer pages to hold inner, replacement policy does not matter. Otherwise, MRU is best, LRU is worst (*sequential flooding*).

  - Does replacement policy matter for Block Nested Loops?

  - What about Index Nested Loops? Sort-Merge Join?

IDB, Eval RA

# 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.