

Database Management Systems

Lecture 3

Winter 2004

CMPUT 391: Query Processing: The Basics

Dr. Osmar R. Zaiane



University of Alberta

Chapter 13
of Textbook

Based on slides by Lewis, Bernstein and Kifer.

External Sorting

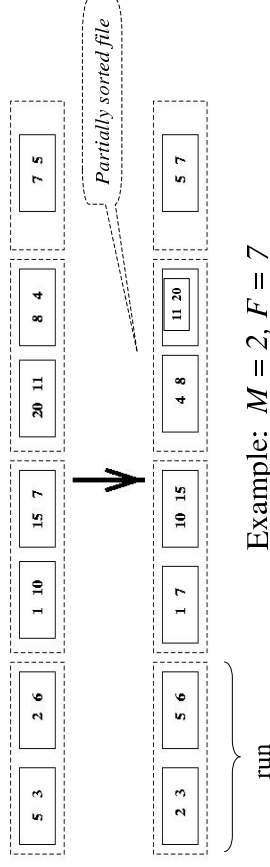
- Sorting is used in implementing many relational operations
- Problem:
 - Relations are typically large, do not fit in main memory
 - So cannot use traditional in-memory sorting algorithms
- Approach used:
 - Combine in-memory sorting with clever techniques aimed at minimizing I/O
 - I/O costs dominate \rightarrow cost of sorting algorithm is measured in the number of page transfers

External Sorting (cont'd)

- External sorting has two main components:
 - Computation involved in sorting records in buffers in main memory
 - I/O necessary to move records between mass store and main memory

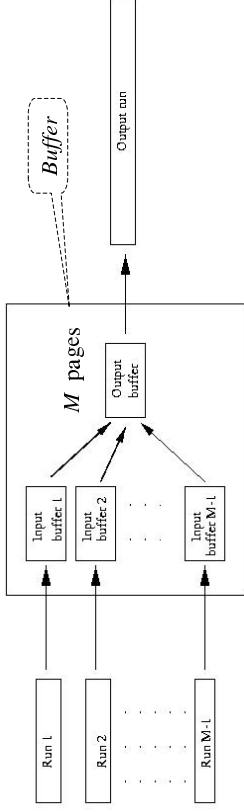
Simple Sort Algorithm

- M = number of main memory page buffers
- F = number of pages in file to be sorted
- Typical algorithm has two phases:
 - **Partial sort phase:** sort M pages at a time; create F/M sorted runs on mass store, cost = $2F$
(read the file and write the file $\rightarrow 2F$ I/Os)



K-Way Merge Algorithm

- **Merge Phase:** merge all runs into a single run using $M-1$ buffers for input and 1 output buffer
- Merge step: divide runs into groups of size $M-1$ and merge each group into a run; cost = $2F$
each step reduces number of runs by a factor of $M-1$



K-Way Merge Algorithm

- Cost of merge phase:
 - $(F/M)/(M-1)^k$ runs after k merge steps
 - $\lceil \text{Log}_{M-1}(F/M) \rceil$ merge steps needed to merge an initial set of F/M sorted runs
- Cost = # merge steps * $2F$ (R & W every page)
= $2F \lceil \text{Log}_{M-1}(F/M) \rceil$
- Total cost = cost of partial sort phase + cost of merge phase
 $2F + 2F \lceil \text{Log}_{M-1}(F/M) \rceil = 2F(1 + \lceil \text{Log}_{M-1}(F/M) \rceil)$
- *Simplification in text book (not to use in 391 tests)*
 $\text{cost} = \lceil 2F \text{Log}_{M-1}(F/M) \rceil \approx 2F(\text{Log}_{M-1} F - 1) \approx 2F \text{Log}_{M-1} F$

Duplicate Elimination

- A major step in computing *projection*, *union*, and *difference* relational operators
- Algorithm:
 - Sort
 - At the last stage of the merge step eliminate duplicates on the fly
 - No additional cost (with respect to sorting) in terms of I/O

Sort-Based Projection

- Algorithm:
 - Sort rows of relation at cost of $2F(1 + \lceil \text{Log}_{M-1}(F/M) \rceil)$ eliminate unwanted columns in partial sort phase (no additional cost)
 - Eliminate duplicates on completion of last merge step (no additional cost)
- Cost: the cost of sorting

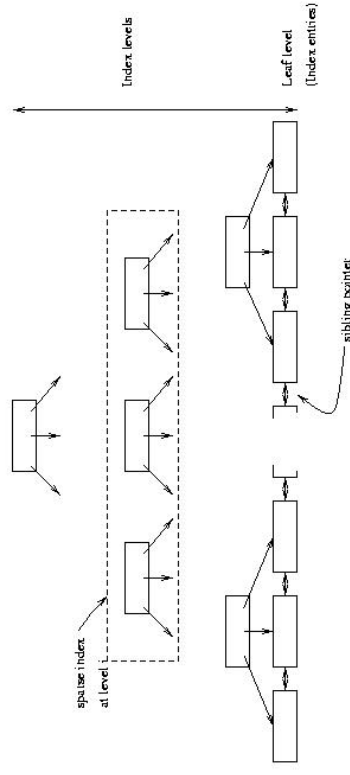
Quick review ...

- At this point many decisions (regarding query processing/optimizations) will be depend on what type of index(es), if any, are available.
- Hence, a quick review is in order ... (for details refer to Chapter 11 in the textbook).

B⁺ Tree

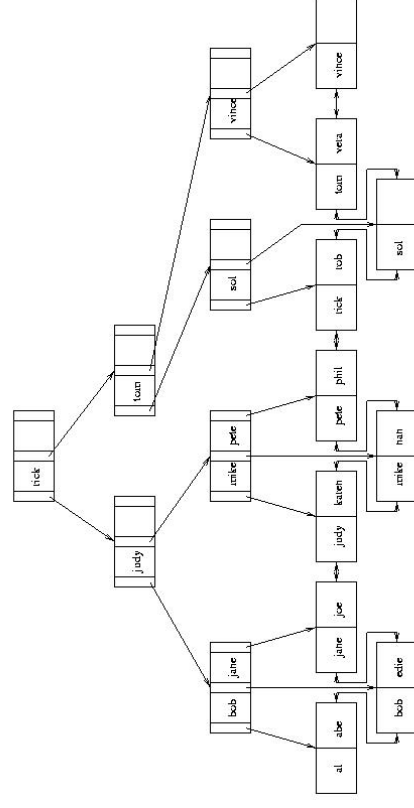
- Supports equality and range searches, multiple attribute keys and partial key searches
 - Either a secondary index (in a separate file) or the basis for an integrated storage structure
- *Responds to dynamic changes in the table*

B⁺ Tree Structure



- Leaf level is a (sorted) linked list of index entries
- Sibling pointers support range searches in spite of allocation and deallocation of leaf pages (but leaf pages might not be physically contiguous on disk)

Example B⁺ Tree



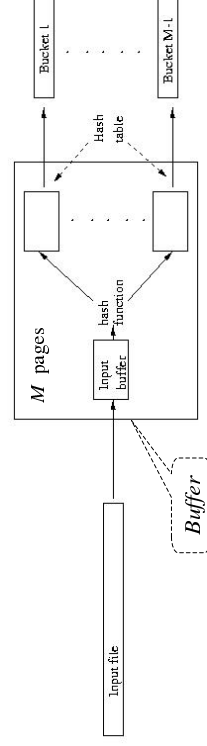
Hash Index

- Index entries partitioned into *buckets* in accordance with a *hash function*, $h(v)$, where v ranges over search key values
 - Each bucket is identified by an address, a
 - Bucket at address a contains all index entries with search key v such that $h(v) = a$
- Each bucket is stored in a page (with possible overflow chain)
- If index entries contain rows, set of buckets forms an integrated storage structure; else set of buckets forms an (unclustered) secondary index

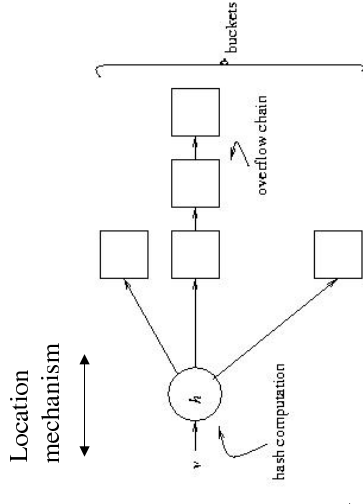


Hash-Based Projection

- *Phase 1:*
 - Input rows
 - Project out columns
 - Hash remaining columns using a hash function with range $1 \dots M-1$ creating $M-1$ buckets on disk
 - **Cost** = $2F$
- *Phase 2:*
 - Sort each bucket to eliminate duplicates
 - **Cost** (assuming a bucket fits in $M-1$ buffer pages) = $2F$
- **Total cost** = $4F$



Equality Search with Hash Index



Given v :

1. Compute $h(v)$
2. Fetch bucket at $h(v)$
3. Search bucket

There are techniques to avoid/minimize overflow (see textbook)



Computing Selection $\sigma_{(attr \text{ op } value)}$

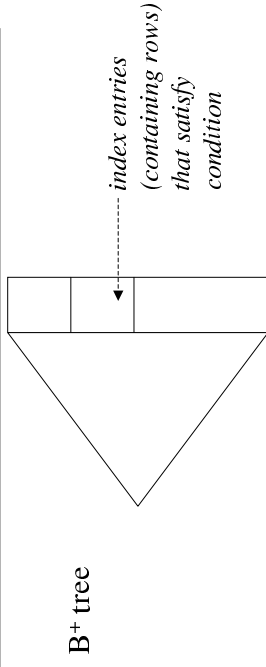
- **No index on $attr$:**
 - If rows are not sorted on $attr$:
 - Scan all data pages to find rows satisfying selection condition
 - **Cost** = F
 - If rows are sorted on $attr$ and op is $=, >, <$ then:
 - Use *binary search* (at $\log_2 F$) to locate first data page containing row in which ($attr = value$)
 - Scan further to get all rows satisfying ($attr \text{ op } value$)
 - **Cost** = $\log_2 F + (\text{cost of scan})$



Computing Selection $\sigma_{(attr \text{ op } value)}$

- Clustered B⁺ tree index on *attr* (for “=” or range search):
 - Locate first index entry corresponding to a row in which (*attr = value*). **Cost = depth of tree**
 - Rows satisfying condition packed in sequence in successive data pages; scan those pages.

Cost: number of pages occupied by qualifying rows



Computing Selection $\sigma_{(attr \text{ op } value)}$

- Unclustered B⁺ tree index on *attr* (for “=” or range search):
 - Locate first index entry corresponding to a row in which (*attr = value*).

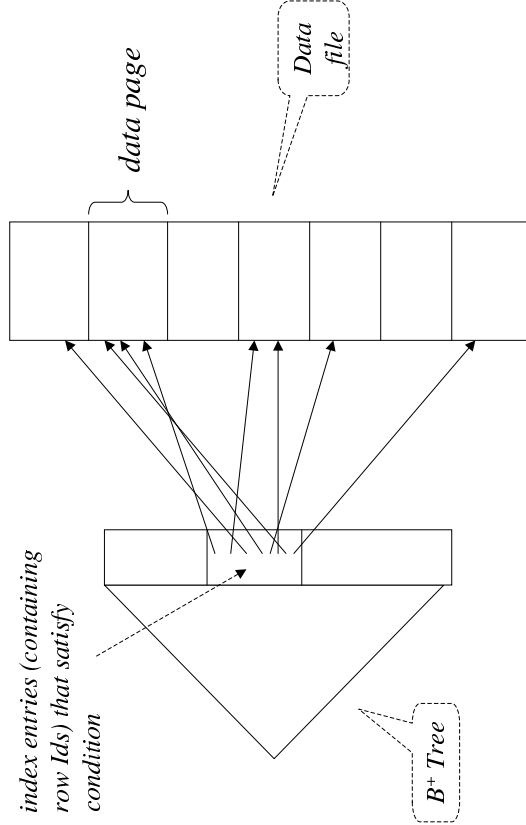
Cost = depth of tree

- Index entries with pointers to rows satisfying condition are packed in sequence in successive index pages
- Scan entries and sort record Ids to identify table data pages with qualifying rows
- Any page that has at least one such row must be fetched once.

• Cost: number of rows that satisfy selection condition



Unclustered B⁺ Tree Index



Computing Selection $\sigma_{(attr = value)}$

- Hash index on *attr* (for “=” search only):

– Hash on *value*. **Cost ≈ 1.2**

- 1.2 – typical average cost of hashing (> 1 due to possible overflow chains)
- Finds the (unique) bucket containing all index entries satisfying selection condition
- Clustered index – all qualifying rows packed in the bucket (a few pages)

Cost: number of pages occupies by the bucket

- Unclustered index – sort row Ids in the index entries to identify data pages with qualifying rows

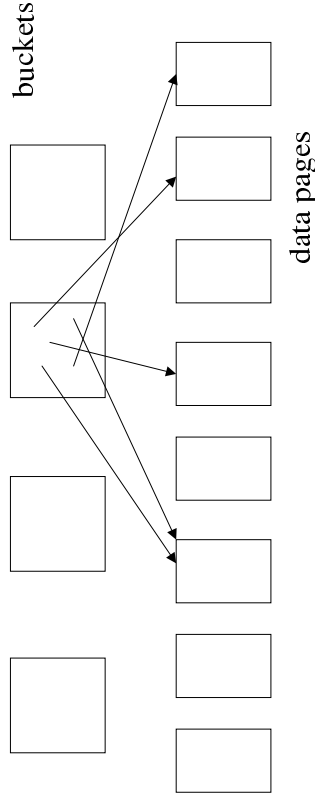
Each page containing at least one such row must be fetched once

Cost: number of rows in bucket



Computing Selection $\sigma_{(attr = value)}$

- Unclustered hash index on *attr* (for equality search)



Access Paths Supported by B⁺ tree

- Example:** Given a B⁺ tree whose search key is the sequence of attributes *a2*, *a1*, *a3*, *a4*
 - Access path for search $\sigma_{a1>5 \wedge a2=3 \wedge a3='x'}$ (*R*): find first entry having *a2*=3 \wedge *a1*>5 \wedge *a3*='x' and scan leaves from there until entry having *a2*>3 or *a3* \neq 'x'. Select satisfying entries
 - Access path for search $\sigma_{a2=3 \wedge a3 > 'x'}$ (*R*): locate first entry having *a2*=3 and scan leaves until entry having *a2*>3. Select satisfying entries
 - Access path for search $\sigma_{a1>5 \wedge a3 = 'x'}$ (*R*): Scan of *R*



Access Path

- Access path* is the notion that denotes *algorithm* + *data structure* used to locate rows satisfying some condition
- Examples:**
 - File scan:* can be used for any condition
 - Hash:* equality search; *all* search key attributes of hash index are specified in condition
 - B⁺ tree:* equality *or* range search; a *prefix* of the search key attributes are specified in condition
 - B⁺ tree supports a variety of access paths
 - Binary search:* Relation sorted on a sequence of attributes and some *prefix* of that sequence is specified in condition



Choosing an Access Path

- Selectivity* of an access path = number of pages retrieved using that path
- If several access paths support a query, DBMS chooses the one with *lowest* selectivity
- Size of domain of attribute is an indicator of the selectivity of search conditions that involve that attribute
- Example:** $\sigma_{CrsCode='CS305' \wedge Grade='B'}$ (Transcript)
 - a B⁺ tree with search key *CrsCode* has lower selectivity than a B⁺ tree with search key *Grade*



Computing Joins

- The cost of joining two relations makes the choice of a join algorithm crucial
- *Block-nested loops* join algorithm for computing $\mathbf{r} \bowtie_{A=B} \mathbf{s}$

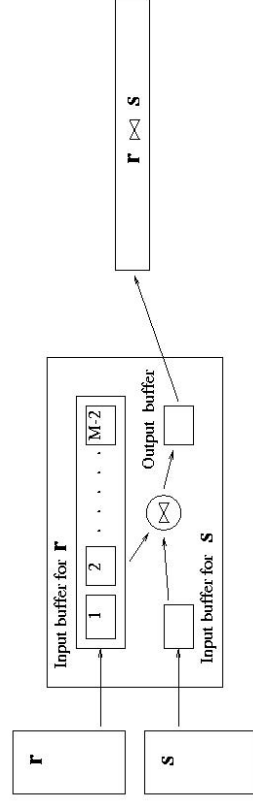
```

foreach page  $p_r$  in  $\mathbf{r}$  do
  foreach tuple  $t_r$  in  $p_r$  do
    foreach page  $p_s$  in  $\mathbf{s}$  do
      foreach tuple  $t_s$  in  $p_s$  do
        if  $t_r.A = t_s.B$  then output  $(t_r, t_s)$ 
  
```



Block-Nested Loops Join

- Cost can be reduced to $\beta_r + (\beta_r / (M-2)) * \beta_s + \text{cost of outputting final result}$ by using M buffer pages instead of 1.



Block-Nested Loops Join

- If β_r and β_s are the number of pages in \mathbf{r} and \mathbf{s} , the cost of algorithm is $\beta_r + \beta_r * \beta_s + \text{cost of outputting final result}$
 - Number of scans of relation \mathbf{s}*
- If \mathbf{r} and \mathbf{s} have 10^3 pages each, cost is $10^3 + 10^3 * 10^3$
- Choose *smaller relation for the outer loop*:
 - If $\beta_r < \beta_s$ then $\beta_r + \beta_r * \beta_s < \beta_s + \beta_r * \beta_s$



Index-Nested Loop Join $\mathbf{r} \bowtie_{A=B} \mathbf{s}$

- Use an index on \mathbf{s} with search key B (instead of scanning \mathbf{s}) to find rows of \mathbf{s} that match t_r
 - Cost = $\beta_r + \tau_r * \omega + \text{cost of outputting final result}$
 - Number of rows in \mathbf{r}*
 - avg cost of retrieving all rows in \mathbf{s} that match t_r*
- Effective if number of rows of \mathbf{s} that match tuples in \mathbf{r} is small (i.e., $\omega \ll 1$) and index is clustered

```

foreach tuple  $t_r$  in  $\mathbf{r}$  do {
  use index to find all tuples  $t_s$  in  $\mathbf{s}$  satisfying  $t_r.A = t_s.B$ ;
  output  $(t_r, t_s)$ 
}
  
```

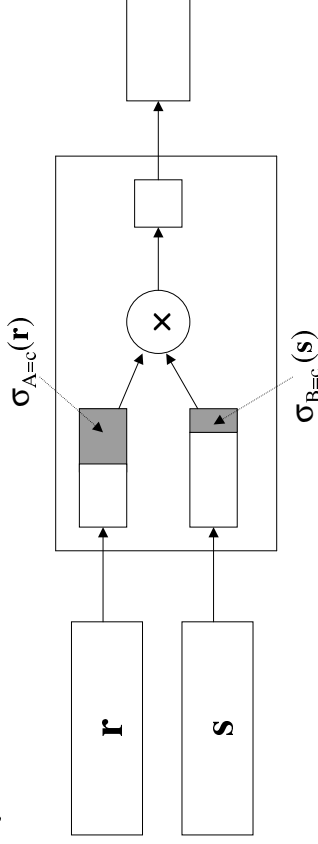


Sort-Merge Join $\mathbf{r} \bowtie_{A=B} \mathbf{s}$

```

sort r on A;
sort s on B;
while !eof(r) and !eof(s) do {
    scan r and s concurrently until  $t_r.A=t_s.B$ ;
    if ( $t_r.A=t_s.B=c$ ) { output  $\sigma_{A=c}(\mathbf{r}) \times \sigma_{B=c}(\mathbf{s})$  }
}

```



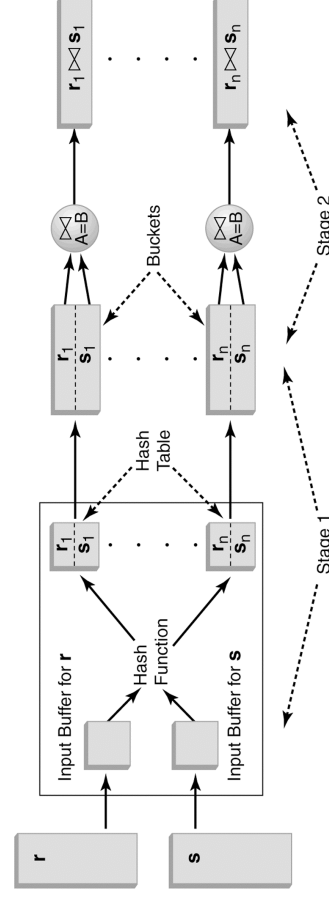
Sort-Merge Join

- Cost of *sorting* assuming M buffers: $2\beta_r \log_{M-1} \beta_r + 2\beta_s \log_{M-1} \beta_s + \text{cost of final result}$
 - Cost of *merging* depends on whether $\sigma_{A=c}(\mathbf{r})$ and $\sigma_{B=c}(\mathbf{s})$ can be fit in buffers
 - If yes, merge step takes $\beta_r + \beta_s$
 - In no, then each $\sigma_{A=c}(\mathbf{r}) \times \sigma_{B=c}(\mathbf{s})$ has to be computed using *block-nested join*.
- (Think why indexed methods or sort-merge are inapplicable.)

Hash-Join $\mathbf{r} \bowtie_{A=B} \mathbf{s}$

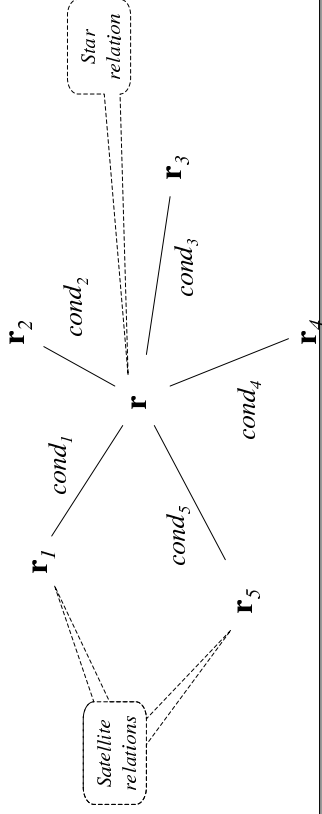
- *Step 1*: Hash \mathbf{r} on A and \mathbf{s} on B into the same set of buckets
- *Step 2*: Since matching tuples must be in same bucket, read each bucket in turn and output the result of the join
- *Cost*: $3(\beta_r + \beta_s) + \text{cost of output of final result}$
 - assuming each bucket fits in memory

Hash Join

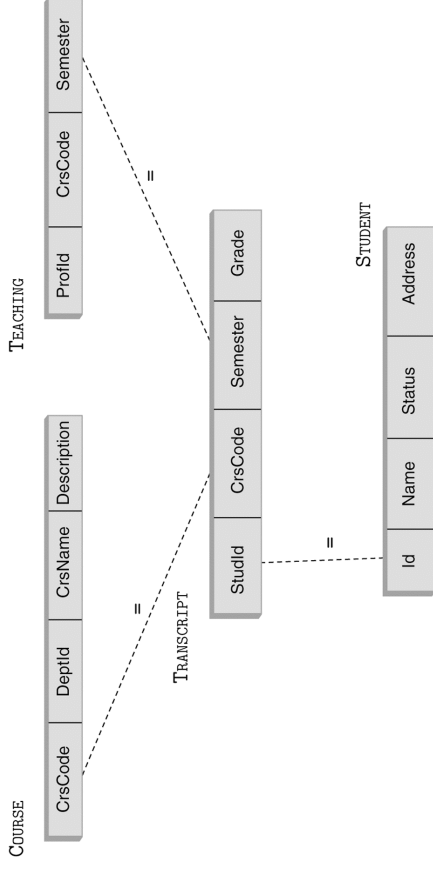


Star Joins

- $\mathbf{r} \bowtie_{cond_1} \mathbf{r}_1 \bowtie_{cond_2} \dots \bowtie_{cond_n} \mathbf{r}_n$
 – Each $cond_i$ involves only the attributes of \mathbf{r}_i and \mathbf{r}



Star Join



Computing Star Joins

- Use *join index* (Chapter 11)
 - Scan \mathbf{r} and the join index $\{ \langle r, r_1, \dots, r_n \rangle \}$ (which is a set of tuples of rids) in one scan
 - Retrieve matching tuples in $\mathbf{r}_1, \dots, \mathbf{r}_n$
 - Output result

Computing Star Joins

- Use *bitmap indices* (Chapter 11)
 - Use one bitmapped join index, J_j , per each partial join $\mathbf{r} \bowtie_{cond_j} \mathbf{r}_j$
 - Recall: J_j is a set of $\langle v, bitmap \rangle$, where v is an rid of a tuple in \mathbf{r}_j and $bitmap$ has 1 in k -th position iff k -th tuple of \mathbf{r} joins with the tuple pointed to by v

 1. Scan J_j and logically OR all bitmaps. We get all rids in \mathbf{r} that join with \mathbf{r}_j
 2. Now logically AND the resulting bitmaps for J_1, \dots, J_n .
 3. Result: a subset of \mathbf{r} , which contains all tuples that can possibly be in the star join
 - *Rationale:* only a few such tuples survive, so can use indexed loops

Choosing Indices

- DBMSs may allow user to specify
 - Type (hash, B⁺ tree) and search key of index
 - Whether or not it should be clustered
- Using information about the frequency and type of queries and size of tables, designer can use cost estimates to choose appropriate indices
- Several commercial systems have tools that suggest indices
 - Simplifies job, but index suggestions must be *verified*

Choosing Indices – Example

- If a frequently executed query that involves selection or a join and has a large result set, use a clustered B⁺ tree index
 - *Example:* Retrieve all rows of Transcript for *StudId*
- If a frequently executed query is an equality search and has a small result set, an unclustered hash index is best
 - Since only one clustered index on a table is possible, choosing unclustered allows a different index to be clustered
 - *Example:* Retrieve all rows of Transcript for (*StudId*, *CrsCode*)