

Chapter 15: Query Processing

(and also chapter 22: Parallel Query Processing)

Sistemas de Bases de Dados 2019/20

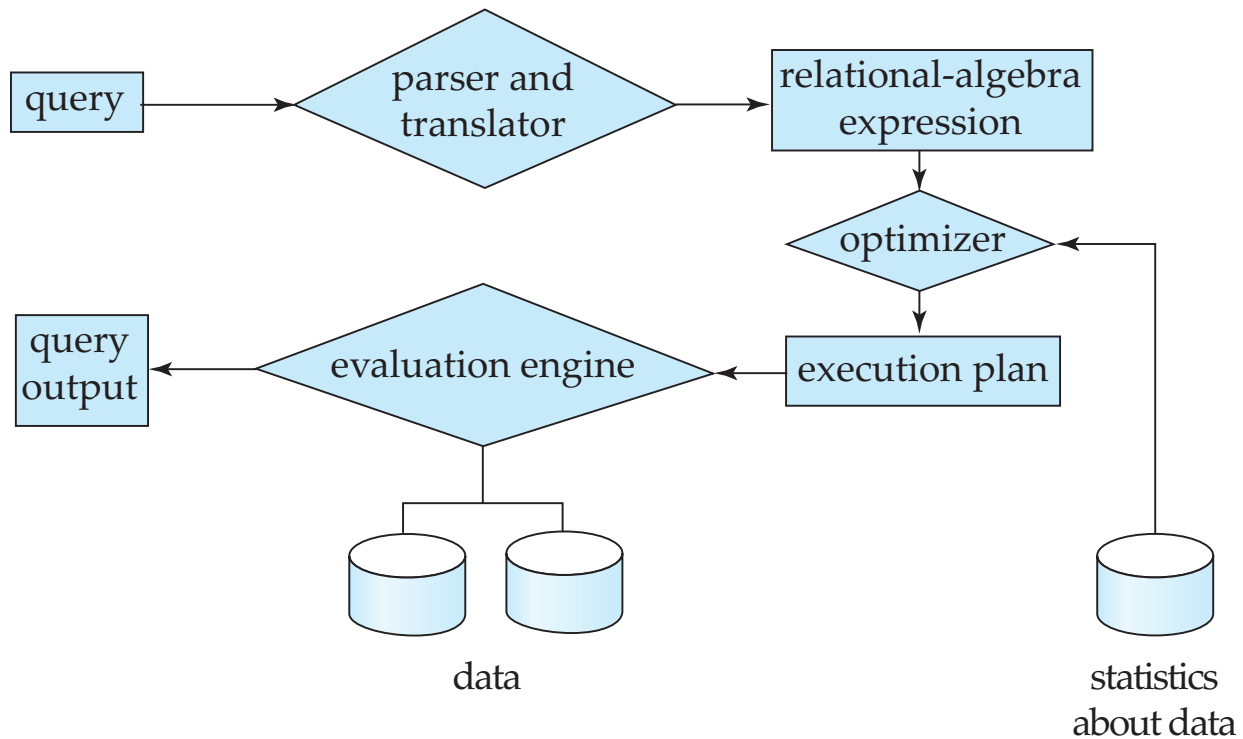
Capítulo refere-se a: Database System Concepts, 7th Ed

Chapter 15: Query Processing

- Overview
- Measures of Query Cost
- Selection Operation
- Sorting
- Join Operation
- Other Operations
- Evaluation of Expressions
- Parallel query processing

Basic Steps in Query Processing

1. Parsing and translation
2. Optimization
3. Evaluation



Join Operation

- Several different algorithms to implement joins
 - Nested-loop join
 - Block nested-loop join
 - Indexed nested-loop join
 - Merge-join
 - Hash-join
- Choice based on cost estimate
- Examples use the following information
 - Number of records of *student*: 5,000 *takes*: 10,000
 - Number of blocks of *student*: 100 *takes*: 400

Nested-Loop Join

- To compute the theta join $r \bowtie_{\theta} s$
 for each tuple t_r **in** r **do begin**
 for each tuple t_s **in** s **do begin**
 test pair (t_r, t_s) to see if they satisfy the join condition θ
 if they do, add $t_r \cdot t_s$ to the result.
 end
 end
- r is called the **outer relation** and s the **inner relation** of the join.
- Requires no indices and can be used with any kind of join condition.
- Expensive since it examines every pair of tuples in the two relations.

Nested-Loop Join Costs

- In the worst case, if there is enough memory only to hold one block of each relation, the estimated cost is
$$n_r * b_s + b_r \text{ block transfers, plus } n_r + b_r \text{ seeks}$$
- In general, it is much better to have the smaller relation as the outer relation
 - The number of block transfers is multiplied by the number of blocks of the inner relation
 - The number of seeks only depends on the outer relation
- However, if the smaller relation fits entirely in memory, one should use it as the inner relation!
 - Reduces cost to $b_r + b_s$ block transfers and 2 seeks
- The choice of the inner and outer relation strongly depends on the estimate of the size of each relation
 - Statics on the size of the relations, in run time, can be a great help!

Nested-Loop Join Costs

- For joining *student* and *takes*, assuming worst case memory availability, the cost estimate is
 - with *student* as outer relation:
 - $5000 * 400 + 100 = 2,000,100$ block transfers,
 - $5000 + 100 = 5100$ seeks
 - with *takes* as the outer relation
 - $10000 * 100 + 400 = 1,000,400$ block transfers and 10,400 seeks
- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers and 2 seeks
- Instead of iterating over records, one could iterate over blocks. This way, instead of $n_r * b_s + b_r$ we would have $b_r * b_s + b_r$ block transfers
- This is the basis of the block nested-loops algorithm (next slide).

Block Nested-Loop Join

- Variant of nested-loop join in which every block of inner relation is paired with every block of outer relation.

```
for each block  $B_r$  of  $r$  do begin  
  for each block  $B_s$  of  $s$  do begin  
    for each tuple  $t_r$  in  $B_r$  do begin  
      for each tuple  $t_s$  in  $B_s$  do begin  
        Check if  $(t_r, t_s)$  satisfy the join condition  
        if they do, add  $t_r \cdot t_s$  to the result.  
      end  
    end  
  end  
end
```


Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + $2 * b_r$ seeks
 - Each block in the inner relation s is read once for each *block* in the outer relation
- Best case (when smaller relation fits into memory): $b_r + b_s$ block transfers plus 2 seeks.
- In the running example the cost of *student* ⋈ *takes* is:
 - If *student* is outer: $100 * 400 + 100 = 40,100$ transfer + 200 seeks
 - If *takes* is outer: $400 * 100 + 400 = 40,400$ transfers + 400 seeks
- Improvements to nested loop and block nested loop algorithms:
 - If equijoin attribute forms a key or inner relation, stop inner loop on first match
 - Scan inner loop forward and backward alternately, to make use of the blocks remaining in buffer (with LRU replacement)
 - Use index on inner relation if available (next slide)

Indexed Nested-Loop Join

- Index lookups can replace file scans if
 - join is an equijoin or natural join and
 - an index is available on the inner relation's join attribute
 - In some cases, it pays to construct an index just to compute a join.
- For each tuple t_r in the outer relation r , use the index to look up tuples in s that satisfy the join condition with tuple t_r .
- Worst case: buffer has space for only one page of r , and, for each tuple in r , we perform an index lookup on s .
- Cost of the join: $b_r(t_T + t_S) + n_r * c$
 - Where c is the cost of traversing index and fetching all matching s tuples for one tuple of r
 - c can be estimated as cost of a single selection on s using the join condition (usually quite small compared to the join cost)
- If indices are available on join attributes of both r and s , use the relation with fewer tuples as the outer relation.

Example of Nested-Loop Join Costs

- Compute $student \bowtie takes$, with $student$ as the outer relation.
- Let $takes$ have a primary B⁺-tree index on the attribute ID , which contains 20 entries in each index node.
- Since $takes$ has 10,000 tuples, the height of the tree is 4, and one more access is needed to find the actual data
- $student$ has 5000 tuples
- As we've seen, the best cost of block nested loops join
 - $400 \times 100 + 100 = 40,100$ block transfers + $2 \times 100 = 200$ seeks
 - assuming worst case memory
 - may be significantly less with more memory
- Cost of indexed nested loops join
 - $100 + 5000 \times 5 = 25,100$ block transfers and seeks.
 - CPU cost likely to be less than that for block nested loops join
 - However in terms of time for transfers and seeks, in this case using the index doesn't pay (this is so because the relations are small)

Merge-Join

1. Sort both relations on their join attribute (if not already sorted on the join attributes).
2. Merge the sorted relations to join them
 1. Join step is like the merge stage of the sort-merge algorithm.
 2. Main difference is handling of duplicate values in join attribute — every pair with same value on join attribute must be matched
 3. Detailed algorithm in the book

	<i>a1</i>	<i>a2</i>
<i>pr</i>	a	3
	b	1
	d	8
	d	13
	f	7
	m	5
	q	6
	<i>r</i>	

	<i>a1</i>	<i>a3</i>
<i>ps</i>	a	A
	b	G
	c	L
	d	N
	m	B
	<i>s</i>	

Merge-Join (Cont.)

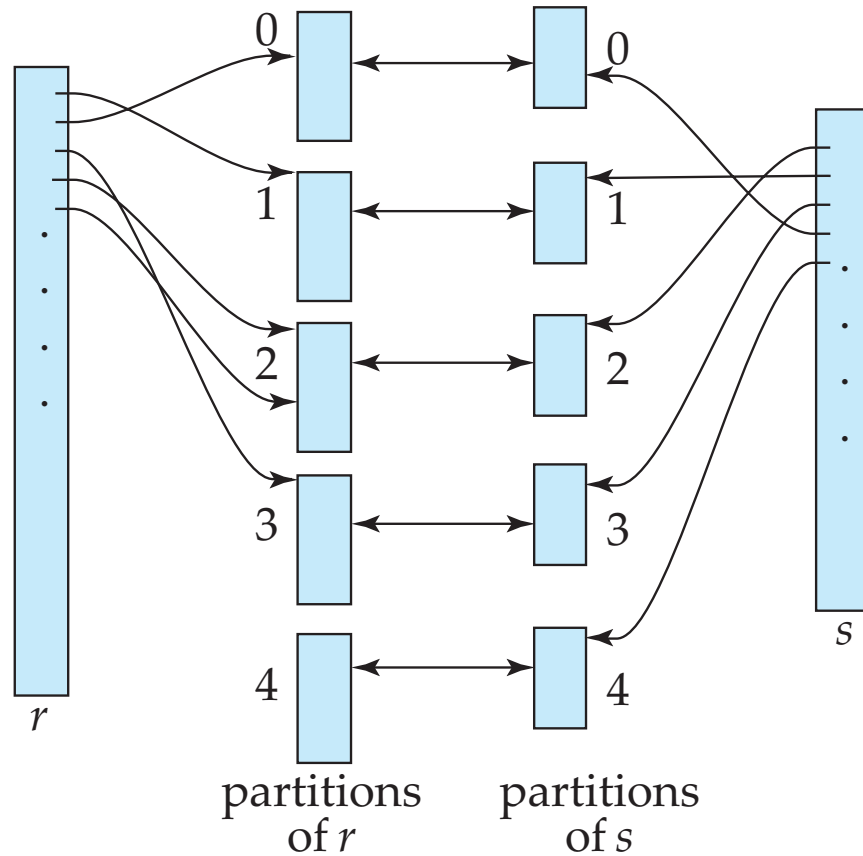
- Can be used only for equi-joins and natural joins
- Each block needs to be read only once (assuming all tuples for any given value of the join attributes fit in memory)
- Thus the cost of merge join is:
$$b_r + b_s \text{ block transfers} + \lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil \text{ seeks}$$

+ the cost of sorting if relations are unsorted.
- **hybrid merge-join:** If one relation is sorted, and the other has a secondary B⁺-tree index on the join attribute
 - Merge the sorted relation with the leaf entries of the B⁺-tree .
 - Sort the result on the addresses of the unsorted relation's tuples
 - Scan the unsorted relation in physical address order and merge with previous result, to replace addresses by the actual tuples
 - Sequential scan more efficient than random lookup

Hash-Join

- Applicable for equi-joins and natural joins.
- A hash function h is used to partition tuples of both relations
- h maps *JoinAttrs* values to $\{0, 1, \dots, n\}$, where *JoinAttrs* denotes the common attributes of r and s used in the natural join.
 - r_0, r_1, \dots, r_n denote partitions of r tuples
 - Each tuple $t_r \in r$ is put in partition r_i where $i = h(t_r[\text{JoinAttrs}])$.
 - s_0, s_1, \dots, s_n denotes partitions of s tuples
 - Each tuple $t_s \in s$ is put in partition s_i , where $i = h(t_s[\text{JoinAttrs}])$.
- General idea:
 - Partition the relations according to this
 - Then perform the join on each partition r_i and s_i
 - There is no need to compute the join between different partitions since an r tuple and an s tuple that satisfy the join condition will have the same value for the join attributes. If that value is hashed to some value i , the r tuple must be in r_i and the s tuple in s_i

Hash-Join (Cont.)



Hash-Join Algorithm

The hash-join of r and s is computed as follows.

1. Partition the relation s using hashing function h . When partitioning a relation, one block of memory is reserved as the output buffer for each partition.
2. Partition r similarly.
3. For each i :
 - (a) Load s_i into memory and build an in-memory hash index on it using the join attribute. *This hash index uses a different hash function than the earlier one h .*
 - (b) Read the tuples in r_i from the disk one by one. For each tuple t_r locate each matching tuple t_s in s_i using the in-memory hash index. Output the concatenation of their attributes.

Relation s is called the **build input** and r is called the **probe input**.

Hash-Join algorithm (Cont.)

- The value n and the hash function h is chosen such that each s_i should fit in memory.
 - Typically n is chosen as $\lceil b_s/M \rceil * f$ where f is a “**fudge factor**”, typically around 1.2
 - The probe relation partitions r_i need not fit in memory
- **Recursive partitioning** required if number of partitions n is greater than number of pages M of memory.
 - instead of partitioning n ways, use $M - 1$ partitions for s
 - Further partition the $M - 1$ partitions using a different hash function
 - Use same partitioning method on r
 - Rarely required: e.g., with block size of 4 KB, recursive partitioning not needed for relations of $< 1\text{TB}$ with memory size of 2GB
 - So we will not further consider it here (see the book for details on the associated costs)

Cost of Hash-Join

- If recursive partitioning is not required: cost of hash join is
$$3(b_r + b_s) + 4 * n \text{ block transfers} + 2(\lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil) \text{ seeks}$$
where b_b is the number of blocks allocated for the input and each output buffer
- If the entire build input can be kept in main memory no partitioning is required
 - Cost estimate goes down to $b_r + b_s$
- For the running example *student* ⋈ *takes*
 - Assume that memory size is 20 blocks
 - $b_{student} = 100$ and $b_{takes} = 400$.
 - *student* is to be used as build input. Partition it into five partitions, each of size 20 blocks. This partitioning can be done in one pass.
 - Similarly, partition *takes* into five partitions, each of size 80. This is also done in one pass.
 - Therefore total cost, ignoring cost of writing partially filled blocks (and assuming 3 blocks for input and each partition buffer – so that they fit in memory):
 - $3(100 + 400) = 1500 \text{ block transfers} + 2(\lceil 100/3 \rceil + \lceil 400/3 \rceil) = 336 \text{ seeks}$
 - The best we had was 40,100 block transfer+200 seek (block nested loop) or 25,100 block transfers and seeks (index nested loop).

Complex Joins

- Join with a conjunctive condition:

$$r \bowtie_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n} s$$

- Either use nested loops/block nested loops, or
- Compute the result of one of the simpler joins $r \bowtie_{\theta_i} s$
 - the final result comprises those tuples in the intermediate result that satisfy the remaining conditions

$$\theta_1 \wedge \dots \wedge \theta_{i-1} \wedge \theta_{i+1} \wedge \dots \wedge \theta_n$$

- Join with a disjunctive condition

$$r \bowtie_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n} s$$

- Either use nested loops/block nested loops, or
- Compute as the union of the records in individual joins $r \bowtie_{\theta_i} s$:

$$(r \bowtie_{\theta_1} s) \cup (r \bowtie_{\theta_2} s) \cup \dots \cup (r \bowtie_{\theta_n} s)$$

Joins over Spatial Data

- No simple sort order for spatial joins
- Indexed nested loops join with spatial indices
 - R-trees, quad-trees, k-d-B-trees

Other Operations

- **Duplicate elimination** can be implemented via hashing or sorting.
 - On sorting duplicates will come adjacent to each other, and all but one set of duplicates can be deleted.
 - *Optimization:* duplicates can be deleted during run generation as well as at intermediate merge steps in external sort-merge.
 - Hashing is similar – duplicates will come into the same bucket.
- **Projection:**
 - perform projection on each tuple
 - followed by duplicate elimination.

Other Operations : Aggregation

- **Aggregation** can be implemented in a manner like duplicate elimination.
 - **Sorting** or **hashing** can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.
 - Optimization: **partial aggregation**
 - combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values
 - For count, min, max, sum: keep aggregate values on tuples found so far in the group.
 - When combining partial aggregate for count, add up the partial aggregates
 - For avg, keep sum and count, and divide sum by count at the end

Other Operations : Set Operations

- **Set operations** (\cup , \cap and $-$): can either use variant of merge-join after sorting, or variant of hash-join.
- E.g., Set operations using hashing:
 1. Partition both relations using the same hash function
 2. Process each partition i as follows.
 1. Using a different hashing function, build an in-memory hash index on r_i .
 2. Process s_i as follows
 - $r \cup s$:
 1. Add tuples in s_i to the hash index if they are not already in it.
 2. At end of s_i add the tuples in the hash index to the result.

Other Operations : Set Operations

- E.g., Set operations using hashing:
 1. as before partition r and s ,
 2. as before, process each partition i as follows
 1. build a hash index on r_i
 2. Process s_i as follows
 - $r \cap s$:
 1. output tuples in s_i to the result if they are already there in the hash index
 - $r - s$:
 1. for each tuple in s_i , if it is there in the hash index, delete it from the index.
 2. At end of s_i add remaining tuples in the hash index to the result.

Other Operations : Outer Join

- **Outer join** can be computed either as
 - A join followed by addition of null-padded non-participating tuples.
 - by modifying the join algorithms.
- Modifying merge join to compute $r \bowtie s$
 - In $r \bowtie s$, nonparticipating tuples are those in $r - \Pi_R(r \bowtie s)$
 - Modify merge-join to compute $r \bowtie s$:
 - During merging, for every tuple t_r from r that do not match any tuple in s , output t_r padded with nulls.
 - Right outer-join and full outer-join can be computed similarly.
- Modifying hash join to compute $r \bowtie s$
 - If r is probe relation, output non-matching r tuples padded with nulls
 - If r is build relation, when probing keep track of which r tuples matched s tuples. At end of s_i output non-matched r tuples padded with nulls

Evaluation of Expressions

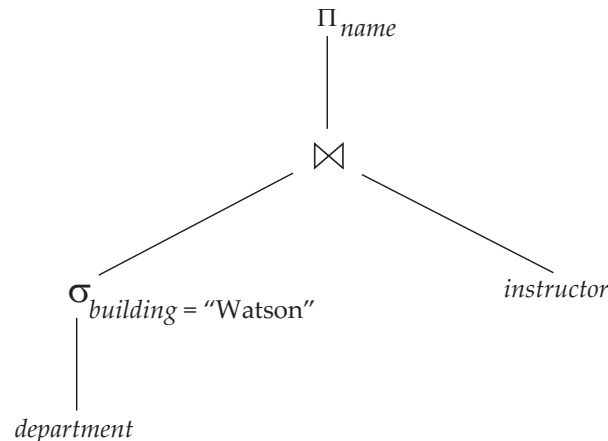
- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree
 - **Materialization**: generate results of an expression whose inputs are relations or are already computed, **materialize** (store) it on disk. Repeat.
 - **Pipelining**: pass on tuples to parent operations even as an operation is being executed
- We study above alternatives in more detail

Materialization

- **Materialized evaluation:** evaluate one operation at a time, starting at the lowest-level. Use intermediate results materialized into temporary relations to evaluate next-level operations.
- E.g., in figure below, compute and store

$$\sigma_{building="Watson"}(department)$$

then compute the store its join with *instructor*, and finally compute the projection on *name*.



Materialization (Cont.)

- Materialized evaluation is always applicable
- Cost of writing results to disk and reading them back can be quite high
 - Our cost formulas for operations ignore cost of writing results to disk, so
 - Overall cost = Sum of costs of individual operations + cost of writing intermediate results to disk
- **Double buffering**: use two output buffers for each operation, when one is full write it to disk while the other is getting filled
 - Allows overlap of disk writes with computation and reduces execution time

Pipelining

- **Pipelined evaluation:** evaluate several operations simultaneously, passing the results of one operation on to the next.

- E.g., in previous expression tree, don't store result of

$$\sigma_{building="Watson"}(department)$$

- instead, pass tuples directly to the join.. Similarly, don't store result of join, pass tuples directly to projection.
- Much cheaper than materialization: no need to store a temporary relation to disk.
- Pipelining may not always be possible – e.g., sort, hash-join.
- For pipelining to be effective, use evaluation algorithms that generate output tuples even as tuples are received for inputs to the operation.
- Pipelines can be executed in two ways: **demand driven** and **producer driven**

Pipelining (Cont.)

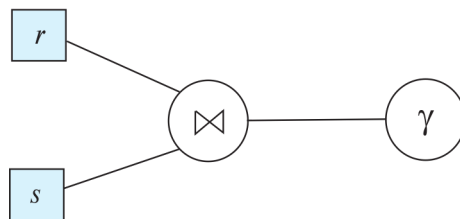
- In **demand driven** (or **lazy** or **pull**) evaluation
 - system repeatedly requests next tuple from top level operation
 - Each operation requests next tuple from children operations as required, in order to output its next tuple
 - In between calls, operation has to maintain “**state**” so it knows what to return next
- In **producer-driven** (or **eager** or **push**) pipelining
 - Operators produce tuples eagerly and pass them up to their parents
 - Buffer maintained between operators, child puts tuples in buffer, parent removes tuples from buffer
 - if buffer is full, child waits till there is space in the buffer, and then generates more tuples
 - System schedules operations that have space in output buffer and can process more input tuples

Pipelining (Cont.)

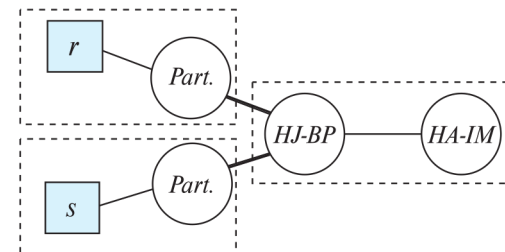
- Implementation of demand-driven pipelining
 - Each operation is implemented as an **iterator** implementing the following operations
 - **open()**
 - E.g., file scan: initialize file scan
 - state: pointer to beginning of file
 - E.g., merge join: sort relations;
 - state: pointers to beginning of sorted relations
 - **next()**
 - E.g., for file scan: Output next tuple, and advance and store file pointer
 - E.g., for merge join: continue with merge from earlier state till next output tuple is found. Save pointers as iterator state.
 - **close()**

Blocking Operations

- **Blocking operations:** cannot generate any output until all input is consumed
 - E.g., sorting, aggregation, ...
- But can often consume inputs from a pipeline, or produce outputs to a pipeline
- Key idea: blocking operations often have two suboperations
 - E.g., for sort: run generation and merge
 - For hash join: partitioning and build-probe
- Treat them as separate operations



(a) Logical Query

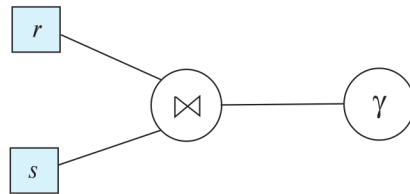


(b) Pipelined Plan

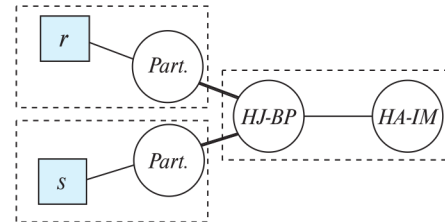
Pipeline Stages

■ Pipeline stages:

- All operations in a stage run concurrently
- A stage can start only after preceding stages have completed execution



(a) Logical Query



(b) Pipelined Plan

Evaluation Algorithms for Pipelining

- Some algorithms are not able to output results even as they get input tuples
 - E.g., merge join, or hash join
 - intermediate results written to disk and then read back
- Algorithm variants to generate (at least some) results on the fly, as input tuples are read in
 - E.g., hybrid hash join generates output tuples even as probe relation tuples in the in-memory partition (partition 0) are read in
- It is clear that pipelining could greatly benefit from parallel processing, especially if there are sufficiently independent sub-expressions
 - And this is not the only chance for parallelism in query processing!