# COL 362 & COL 632

Pipelining, Optimization

28 Mar 2023

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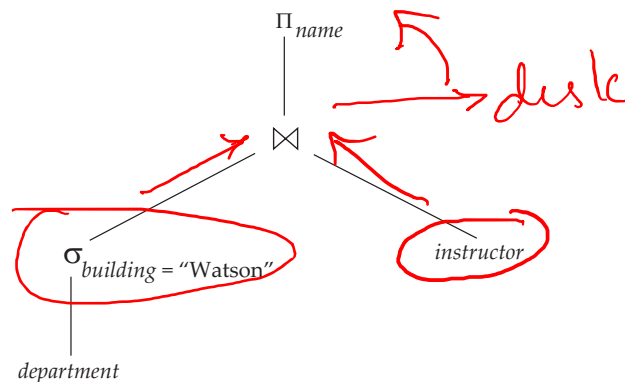
- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree
  - **Materialization**:  generate results of an expression whose inputs are relations or are already computed, **materialize** (store) it on disk. Repeat.
  - **Pipelining**:  pass on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

# 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
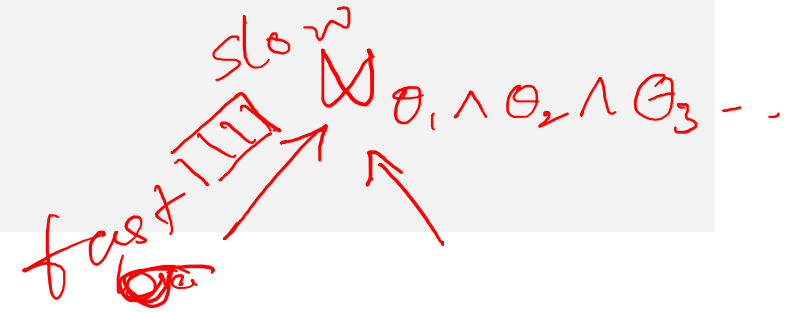
# Pipelining

- **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**
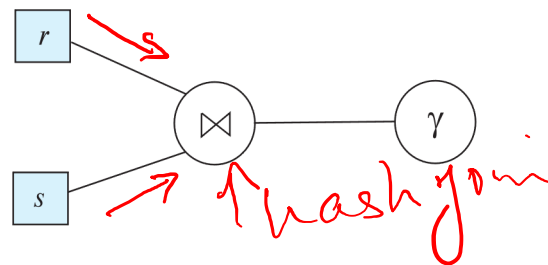
# Pipelining (Cont.)

- In **demand driven** or **lazy** 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** 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

- Alternative name: **pull** and **push** models of pipelining
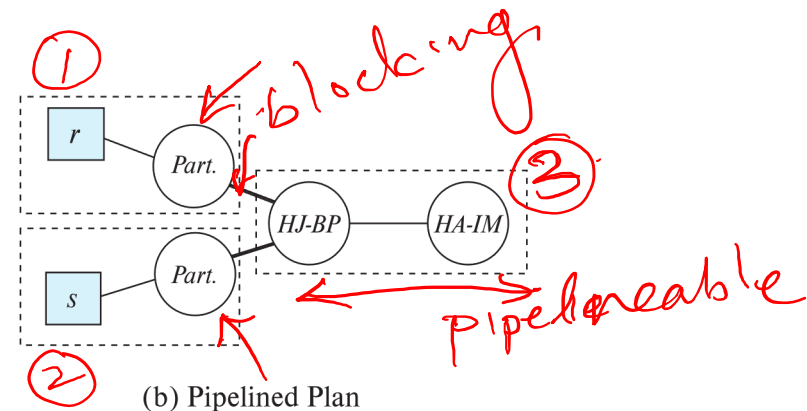
# Pipelining (Cont.)

- 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
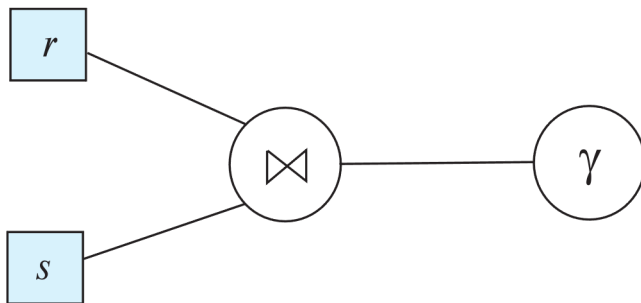- Treat them as separate operations



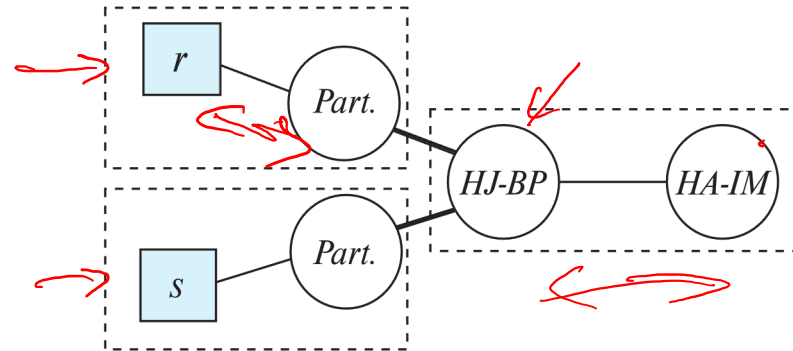(a) Logical Query

(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



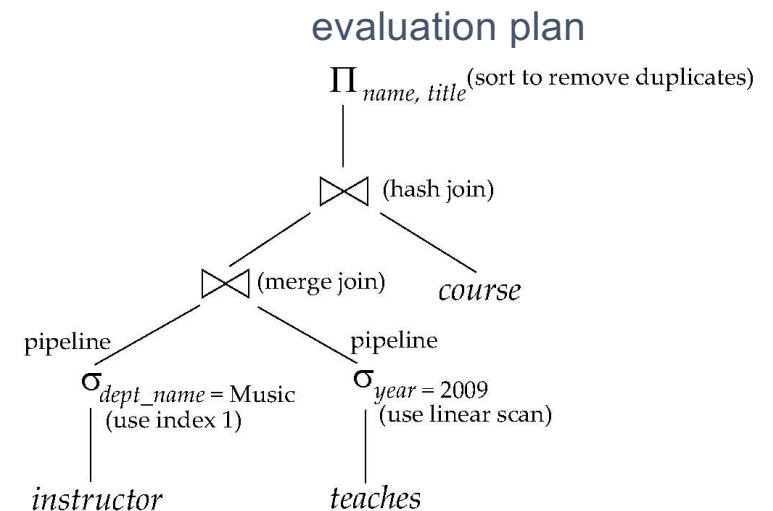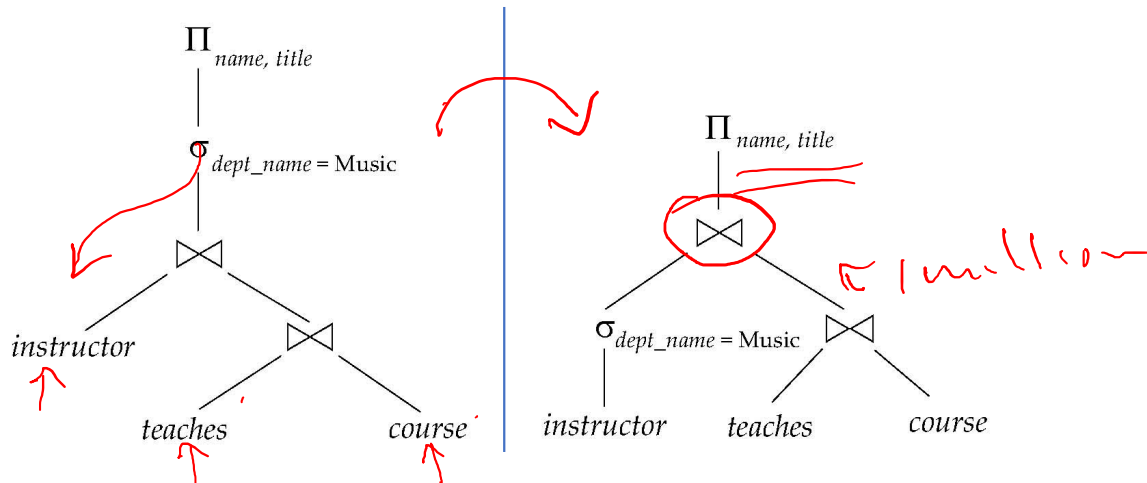(a) Logical Query

(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
  - **Double-pipelined join technique**: Hybrid hash join, modified to buffer partition 0 tuples of both relations in-memory, reading them as they become available, and output results of any matches between partition 0 tuples
    - When a new $r_0$ tuple is found, match it with existing $s_0$ tuples, output matches, and save it in $r_0$
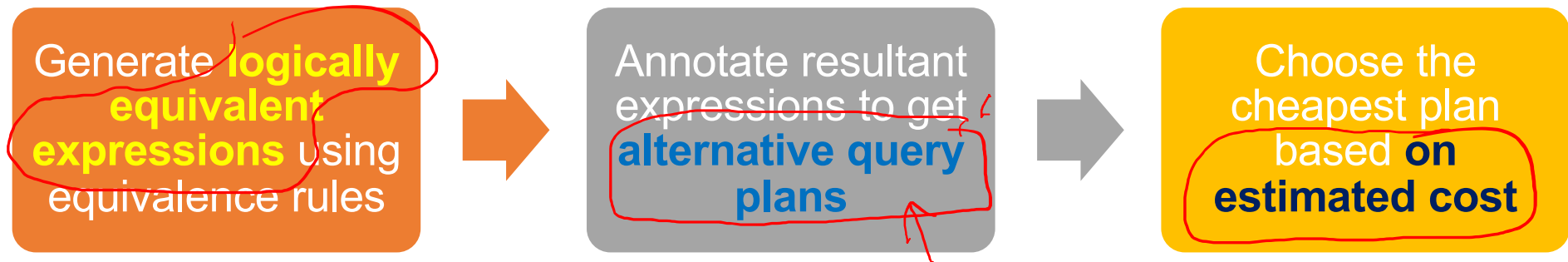    - Symmetrically for $s_0$ tuples

# Query Optimization

# Introduction

- Alternative ways of evaluating a given query
  - Equivalent expressions
  - Different algorithms for each operation

evaluation plan

$\Pi_{name,\ title}$

$\sigma_{dept\_name}$ = Music

$\bowtie$

$instructor$

$teaches$ $course$

$\Pi_{name,\ title}$

$\bowtie$

$\sigma_{dept\_name}$ = Music $\bowtie$

$instructor$ $teaches$ $course$

$\sigma$ 1 million

$\Pi_{name,\ title}$ (sort to remove duplicates)

$\bowtie$ (hash join)

$\bowtie$ (merge join) $course$

pipeline pipeline

$\sigma_{dept\_name}$ = Music (use index 1)

$\sigma_{year}$ = 2009 (use linear scan)

$instructor$ $teaches$

Cost difference between evaluation plans for a query can be enormous
E.g. seconds vs. days in some cases

# Cost Based Query Optimization

Generate **logically equivalent expressions** using equivalence rules

→

Annotate resultant expressions to get **alternative query plans**

→

Choose the cheapest plan based **on estimated cost**

- Estimation of plan cost based on
  - Statistical information about relations
  - Statistics estimation for intermediate results
  - Cost formulae for algorithms, computed using statistics