# CS315: DATABASE SYSTEMS
# QUERY PROCESSING

## Arnab Bhattacharya
arnabb@cse.iitk.ac.in

Computer Science and Engineering,
Indian Institute of Technology, Kanpur
http://web.cse.iitk.ac.in/~cs315/

2nd semester, 2018-19
Mon 12:00-13:15, Tue 9:00-10:15

- Query in a high-level language is first parsed and validated

  **select** cpi **from** students **where** cpi $> 9$;

# Query

- Query in a high-level language is first parsed and validated

  **select** cpi **from** students **where** cpi > 9;

- Query planner next outputs an equivalent expression in relational algebra
  - $\Pi_{cpi}(\sigma_{cpi>9}(students))$
  - A query order tree is created

# Query

- Query in a high-level language is first parsed and validated

  **select** cpi **from** students **where** cpi > 9;

- Query planner next outputs an equivalent expression in relational algebra
  - $\Pi_{cpi}(\sigma_{cpi>9}(students))$
  - A query order tree is created
- All equivalent evaluation plans are then generated
  - $\Pi_{cpi}(\sigma_{cpi>9}(students))$
  - $\sigma_{cpi>9}(\Pi_{cpi}(students))$

# Query

- Query in a high-level language is first parsed and validated

  **select** cpi **from** students **where** cpi > 9;

- Query planner next outputs an equivalent expression in relational algebra
  - $\Pi_{cpi}(\sigma_{cpi>9}(students))$
  - A query order tree is created
- All equivalent evaluation plans are then generated
  - $\Pi_{cpi}(\sigma_{cpi>9}(students))$
  - $\sigma_{cpi>9}(\Pi_{cpi}(students))$
- Query optimizer next decides on the best evaluation plan among all equivalent plans
  - Uses statistics of data
  - Actual algorithms also influence the choice

# Query

- Query in a high-level language is first parsed and validated

  **select** cpi **from** students **where** cpi > 9;

- Query planner next outputs an equivalent expression in relational algebra
  - $\Pi_{cpi}(\sigma_{cpi>9}(students))$
  - A query order tree is created
- All equivalent evaluation plans are then generated
  - $\Pi_{cpi}(\sigma_{cpi>9}(students))$
  - $\sigma_{cpi>9}(\Pi_{cpi}(students))$
- Query optimizer next decides on the best evaluation plan among all equivalent plans
  - Uses statistics of data
  - Actual algorithms also influence the choice
- Query code is finally generated and processed

# Query Cost

- Factors that affect the runtime of the query
  - Disk accesses
  - CPU time
  - Network communication
- In a non-dstributed setting, disk access is the most dominant factor

# Query Cost

- Factors that affect the runtime of the query
  - Disk accesses
  - CPU time
  - Network communication
- In a non-dstributed setting, disk access is the most dominant factor
- It can be estimated as a total of
  - Number of seeks times average seek cost
  - Number of blocks read times average block read cost
  - Number of blocks written times average block write cost
    - Write cost can be more since data is verified

# Query Cost

- Factors that affect the runtime of the query
    - Disk accesses
    - CPU time
    - Network communication
- In a non-dstributed setting, disk access is the most dominant factor
- It can be estimated as a total of
    - Number of seeks times average seek cost
    - Number of blocks read times average block read cost
    - Number of blocks written times average block write cost
        - Write cost can be more since data is verified
- For $s$ seeks and $b$ block transfers, simply estimated as $s \times t_s + b \times t_b$
- Ignores CPU time and buffer management issues

# Selection

- LINEAR SEARCH
- *Always* applicable

# Selection

- LINEAR SEARCH
- *Always* applicable
- Scan all file blocks and test each record

# Selection

- LINEAR SEARCH
- *Always* applicable
- Scan all file blocks and test each record
- Cost for a relation containing $b$ blocks
    - 1 seek
    - $b$ transfers
- If equality on key, then $b/2$ transfers on average

# Selection (contd.)

- BINARY SEARCH
- Applicable for comparison on the ordering attribute
- *Equality*

# Selection (contd.)

- BINARY SEARCH
- Applicable for comparison on the ordering attribute
- *Equality*
- Cost for locating the first record is
  - $\lceil \log_2 b \rceil$ seeks
  - $\lceil \log_2 b \rceil$ transfers
- If one block does not contain all records, add required number of additional transfers to the cost
- *Greater than*

# Selection (contd.)

- BINARY SEARCH
- Applicable for comparison on the ordering attribute
- *Equality*
- Cost for locating the first record is
  - $\lceil \log_2 b \rceil$ seeks
  - $\lceil \log_2 b \rceil$ transfers
- If one block does not contain all records, add required number of additional transfers to the cost
- *Greater than*
  - Add required number of additional transfers
- *Lesser than*

# Selection (contd.)

- BINARY SEARCH
- Applicable for comparison on the ordering attribute
- *Equality*
- Cost for locating the first record is
  - $\lceil \log_2 b \rceil$ seeks
  - $\lceil \log_2 b \rceil$ transfers
- If one block does not contain all records, add required number of additional transfers to the cost
- *Greater than*
  - Add required number of additional transfers
- *Lesser than*
  - Scan from beginning till matching record

# Selection (contd.)

- INDEX SEARCH
- Applicable for comparison on the indexed attribute
- *Equality* using primary index

# Selection (contd.)

- INDEX SEARCH
- Applicable for comparison on the indexed attribute
- *Equality* using primary index
    - $h + 1$ seeks, where $h$ is the height of B+-tree
    - $h + m$ transfers

  where $m$ is the total number of blocks containing matching records
- For key, $m = 1$
- *Equality* using secondary index

# Selection (contd.)

- INDEX SEARCH
- Applicable for comparison on the indexed attribute
- *Equality* using primary index
    - $h + 1$ seeks, where $h$ is the height of B+-tree
    - $h + m$ transfers

    where $m$ is the total number of blocks containing matching records
- For key, $m = 1$
- *Equality* using secondary index
    - $h + n$ seeks, where $h$ is the height of the index tree
    - $h + n$ transfers

    where $n$ is the total number of matching records, each in a separate block
- For key, $n = 1$

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index

# Selection (contd.)

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index
    - $h + 1$ seeks and $h + 1$ transfers to locate $v$
- Scan for the rest of the relation resulting in $n$ more transfers

# Selection (contd.)

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index
  - $h + 1$ seeks and $h + 1$ transfers to locate $v$
- Scan for the rest of the relation resulting in $n$ more transfers
- *Lesser than* ($A \leq v$) using primary index

# Selection (contd.)

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index
  - $h + 1$ seeks and $h + 1$ transfers to locate $v$
- Scan for the rest of the relation resulting in $n$ more transfers
- *Lesser than* ($A \leq v$) using primary index
  - 1 seek to the first block
  - $n$ transfers till $v$ is located

# Selection (contd.)

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index
  - $h + 1$ seeks and $h + 1$ transfers to locate $v$
- Scan for the rest of the relation resulting in $n$ more transfers
- *Lesser than* ($A \leq v$) using primary index
  - 1 seek to the first block
  - $n$ transfers till $v$ is located
- *Lesser than* ($A \leq v$) using secondary index

# Selection (contd.)

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index
  - $h + 1$ seeks and $h + 1$ transfers to locate *v*
- Scan for the rest of the relation resulting in *n* more transfers
- *Lesser than* ($A \leq v$) using primary index
  - 1 seek to the first block
  - *n* transfers till *v* is located
- *Lesser than* ($A \leq v$) using secondary index
  - $h + 1$ seeks and $h + 1$ transfers to locate *v*
- Use sibling leaf pointers to locate all other matching records resulting in *n* more seeks and *n* more transfers

# Selection (contd.)

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index
  - $h + 1$ seeks and $h + 1$ transfers to locate *v*
- Scan for the rest of the relation resulting in *n* more transfers
- *Lesser than* ($A \leq v$) using primary index
  - 1 seek to the first block
  - *n* transfers till *v* is located
- *Lesser than* ($A \leq v$) using secondary index
  - $h + 1$ seeks and $h + 1$ transfers to locate *v*
- Use sibling leaf pointers to locate all other matching records resulting in *n* more seeks and *n* more transfers
- *Greater than* ($A \geq v$) using secondary index

# Selection (contd.)

- INDEX SEARCH
- *Greater than* ($A \geq v$) using primary index
  - $h + 1$ seeks and $h + 1$ transfers to locate *v*
- Scan for the rest of the relation resulting in *n* more transfers
- *Lesser than* ($A \leq v$) using primary index
  - 1 seek to the first block
  - *n* transfers till *v* is located
- *Lesser than* ($A \leq v$) using secondary index
  - $h + 1$ seeks and $h + 1$ transfers to locate *v*
- Use sibling leaf pointers to locate all other matching records resulting in *n* more seeks and *n* more transfers
- *Greater than* ($A \geq v$) using secondary index
  - $h + 1$ seeks and $h + 1$ transfers to locate *v*
- Use sibling leaf pointers to locate all other matching records resulting in *n* more seeks and *n* more transfers

- *Conjunction*: AND

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
  - For each record satisfying it, test the other attributes

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
  - For each record satisfying it, test the other attributes
- Multiple attribute index

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
    - For each record satisfying it, test the other attributes
- Multiple attribute index
- Intersection
    - If some attribute does not have an index, test explicitly

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
  - For each record satisfying it, test the other attributes
- Multiple attribute index
- Intersection
  - If some attribute does not have an index, test explicitly
- *Disjunction*: OR

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
  - For each record satisfying it, test the other attributes
- Multiple attribute index
- Intersection
  - If some attribute does not have an index, test explicitly
- *Disjunction*: OR
- Union
  - If some attribute does not have an index, then linear scan

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
  - For each record satisfying it, test the other attributes
- Multiple attribute index
- Intersection
  - If some attribute does not have an index, test explicitly
- *Disjunction*: OR
- Union
  - If some attribute does not have an index, then linear scan
- Negation of conjunction
- *Negation* of equality

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
  - For each record satisfying it, test the other attributes
- Multiple attribute index
- Intersection
  - If some attribute does not have an index, test explicitly
- *Disjunction*: OR
- Union
  - If some attribute does not have an index, then linear scan
- Negation of conjunction
- *Negation* of equality
  - Linear scan
  - Index selects leaf pointers for which corresponding records will not be retrieved
- Negation of comparison

# Selection on Multiple Attributes

- *Conjunction*: AND
- Single index
  - For each record satisfying it, test the other attributes
- Multiple attribute index
- Intersection
  - If some attribute does not have an index, test explicitly
- *Disjunction*: OR
- Union
  - If some attribute does not have an index, then linear scan
- Negation of conjunction
- *Negation* of equality
  - Linear scan
  - Index selects leaf pointers for which corresponding records will not be retrieved
- Negation of comparison is just another comparison

# Sorting

- Display purposes
- Certain operations such as join can be implemented efficiently on sorted relations

# Sorting

- Display purposes
- Certain operations such as join can be implemented efficiently on sorted relations
- Index provides a *logical* sorted view
- Tuples need to be *physically* sorted

# Sorting

- Display purposes
- Certain operations such as join can be implemented efficiently on sorted relations
- Index provides a *logical* sorted view
- Tuples need to be *physically* sorted
- When the relation fits in memory, QUICKSORT can be used
- When it does not, external sorting algorithms are used

# Sorting

- Display purposes
- Certain operations such as join can be implemented efficiently on sorted relations
- Index provides a *logical* sorted view
- Tuples need to be *physically* sorted
- When the relation fits in memory, QUICKSORT can be used
- When it does not, external sorting algorithms are used
- EXTERNAL MERGESORT or EXTERNAL SORT-MERGE is the most used

# External Mergesort

- Assume only $m$ blocks can be put into memory
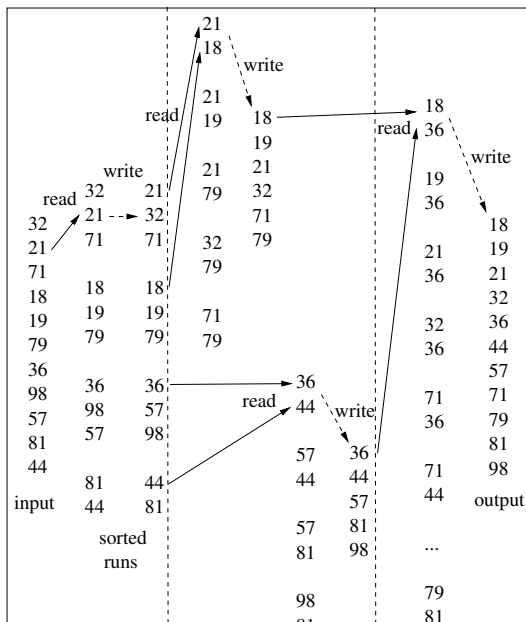- Size of relation is more than $m$ blocks

# External Mergesort

- Assume only $m$ blocks can be put into memory
- Size of relation is more than $m$ blocks
- Create sorted <span style="color:red">runs</span>
  - Read $m$ blocks at a time
  - Sort them in-memory using any algorithm such as quicksort
  - Write them back to disk

# External Mergesort

- Assume only $m$ blocks can be put into memory
- Size of relation is more than $m$ blocks
- Create sorted runs
  - Read $m$ blocks at a time
  - Sort them in-memory using any algorithm such as quicksort
  - Write them back to disk
- Merge $m - 1$ runs (($m - 1$)-way merge)
  - Read in first block of $m - 1$ runs
  - Output the first record to *buffer* block ($m$-th block in memory)
  - Continue till buffer block is full
  - Write buffer block to disk
  - When a block of a particular run is exhausted, read in the next block of the run

# External Mergesort

- Assume only $m$ blocks can be put into memory
- Size of relation is more than $m$ blocks
- Create sorted runs
  - Read $m$ blocks at a time
  - Sort them in-memory using any algorithm such as quicksort
  - Write them back to disk
- Merge $m - 1$ runs (($m - 1$)-way merge)
  - Read in first block of $m - 1$ runs
  - Output the first record to *buffer* block ($m$-th block in memory)
  - Continue till buffer block is full
  - Write buffer block to disk
  - When a block of a particular run is exhausted, read in the next block of the run
- Continue with ($m - 1$)-way merge till the number of sorted runs is *less than m*
- The last ($m - 1$)-way merge sorts the relation

# Example

# Cost of External Mergesort

- Total number of blocks is $b$
- Initial number of sorted runs is

# Cost of External Mergesort

- Total number of blocks is *b*
- Initial number of sorted runs is $n = \lceil b/m \rceil$
- In each merge pass, $m - 1$ runs are sorted
- Therefore, total number of merge passes required is

# Cost of External Mergesort

- Total number of blocks is $b$
- Initial number of sorted runs is $n = \lceil b/m \rceil$
- In each merge pass, $m - 1$ runs are sorted
- Therefore, total number of merge passes required is $r = \lceil \log_{m-1} n \rceil$
- During each of these passes and the first pass, all blocks are read and written
- There is an initial pass of creating runs
- Hence, total number of block transfers is

# Cost of External Mergesort

- Total number of blocks is *b*
- Initial number of sorted runs is $n = \lceil b/m \rceil$
- In each merge pass, $m - 1$ runs are sorted
- Therefore, total number of merge passes required is $r = \lceil \log_{m-1} n \rceil$
- During each of these passes and the first pass, all blocks are read and written
- There is an initial pass of creating runs
- Hence, total number of block transfers is $2b(r + 1)$

# Cost of External Mergesort (contd.)

- Initial pass to created sorted runs reads *m* blocks at a time
- Therefore, number of seeks is 2*n* for reading and writing

# Cost of External Mergesort (contd.)

- Initial pass to created sorted runs reads $m$ blocks at a time
- Therefore, number of seeks is $2n$ for reading and writing
- During the merge passes, blocks from different runs may not be read consecutively
- Consequently, each read and write for another run may move the disk head away, thereby requiring a seek every time
- Hence, number of seeks for these passes is $2br$
- Therefore, total number of seeks is

# Cost of External Mergesort (contd.)

- Initial pass to created sorted runs reads $m$ blocks at a time
- Therefore, number of seeks is $2n$ for reading and writing
- During the merge passes, blocks from different runs may not be read consecutively
- Consequently, each read and write for another run may move the disk head away, thereby requiring a seek every time
- Hence, number of seeks for these passes is $2br$
- Therefore, total number of seeks is $2n + 2br$

# Join

- Different join algorithms
  - NESTED-LOOP JOIN
  - BLOCK NESTED-LOOP JOIN
  - INDEXED NESTED-LOOP JOIN
  - MERGE JOIN
  - HASH JOIN
- Choice depends on cost estimates

# Nested-Loop Join

- Applicable for any kind of join
- For each record $t_r \in r$ and for each record $t_s \in s$, if $t_r \bowtie t_s$ satisfies the join condition, add it to result
- Outer relation $r$: outer loop; inner relation $s$: inner loop

# Block Nested-Loop Join

- Applicable for any kind of join
- Block version of the nested-loop algorithm
- For each block $l_r \in r$ and for each block $l_s \in s$, test if every record $t_r \in l_r$ and $t_s \in l_s$ satisfies the join condition; if so, add to the result

- Worst-case assumption that only two blocks fit in memory

# Cost of Block Nested-Loop Join

- Worst-case assumption that only two blocks fit in memory
- Block transfers
  - $b_s$ transfers every time a block $l_r$ is read
  - $b_r$ transfers for blocks in $r$
  - Therefore, total is $b_r + b_r \times b_s$ transfers

# Cost of Block Nested-Loop Join

- Worst-case assumption that only two blocks fit in memory
- Block transfers
    - $b_s$ transfers every time a block $l_r$ is read
    - $b_r$ transfers for blocks in $r$
    - Therefore, total is $b_r + b_r \times b_s$ transfers
- Seeks
    - 1 seek for records in $s$ every time a block $l_r$ is read
    - $b_r$ seeks for records in $r$
    - Therefore, total is $b_r + b_r = 2b_r$ seeks
- Smaller relation should be

# Cost of Block Nested-Loop Join

- Worst-case assumption that only two blocks fit in memory
- Block transfers
  - $b_s$ transfers every time a block $l_r$ is read
  - $b_r$ transfers for blocks in $r$
  - Therefore, total is $b_r + b_r \times b_s$ transfers
- Seeks
  - 1 seek for records in $s$ every time a block $l_r$ is read
  - $b_r$ seeks for records in $r$
  - Therefore, total is $b_r + b_r = 2b_r$ seeks
- Smaller relation should be outer

- Assume *m* blocks of memory to be available

# Cost of Block Nested-Loop Join (contd.)

- Assume $m$ blocks of memory to be available
- 1 left for inner relation $s$
- $m - 1$ blocks read at a time from outer relation $r$

# Cost of Block Nested-Loop Join (contd.)

- Assume $m$ blocks of memory to be available
- 1 left for inner relation $s$
- $m - 1$ blocks read at a time from outer relation $r$
- Number of such runs is $n = \lceil b_r / (m - 1) \rceil$

# Cost of Block Nested-Loop Join (contd.)

- Assume $m$ blocks of memory to be available
- 1 left for inner relation $s$
- $m - 1$ blocks read at a time from outer relation $r$
- Number of such runs is $n = \lceil b_r / (m - 1) \rceil$
- Block transfers
    - For $r$: $b_r$
    - For $s$: $n.b_s$

# Cost of Block Nested-Loop Join (contd.)

- Assume $m$ blocks of memory to be available
- 1 left for inner relation $s$
- $m - 1$ blocks read at a time from outer relation $r$
- Number of such runs is $n = \lceil b_r / (m - 1) \rceil$
- Block transfers
    - For $r$: $b_r$
    - For $s$: $n.b_s$
- Seeks
    - For $r$: $n$
    - For $s$: $n$
- Smaller relation should be

# Cost of Block Nested-Loop Join (contd.)

- Assume $m$ blocks of memory to be available
- 1 left for inner relation $s$
- $m - 1$ blocks read at a time from outer relation $r$
- Number of such runs is $n = \lceil b_r / (m-1) \rceil$
- Block transfers
    - For $r$: $b_r$
    - For $s$: $n.b_s$
- Seeks
    - For $r$: $n$
    - For $s$: $n$
- Smaller relation should be outer

# Example

- $r$: 40000 tuples at 200 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$

# Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be *r*

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be *r*
- Number of runs is

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be *r*
- Number of runs is $n = \lceil 200/25 \rceil = 8$

## Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be $r$
- Number of runs is $n = \lceil 200/25 \rceil = 8$
- In each run, $s$ requires 1 seek and 250 transfers
- Total is 8 seeks and 2000 transfers

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be *r*
- Number of runs is $n = \lceil 200/25 \rceil = 8$
- In each run, *s* requires 1 seek and 250 transfers
- Total is 8 seeks and 2000 transfers
- *r* requires 8 seeks and 200 transfers

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be *r*
- Number of runs is $n = \lceil 200/25 \rceil = 8$
- In each run, *s* requires 1 seek and 250 transfers
- Total is 8 seeks and 2000 transfers
- *r* requires 8 seeks and 200 transfers
- Total is 16 seeks and 2200 transfers

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be *r*
- Number of runs is $n = \lceil 200/25 \rceil = 8$
- In each run, *s* requires 1 seek and 250 transfers
- Total is 8 seeks and 2000 transfers
- *r* requires 8 seeks and 200 transfers
- Total is 16 seeks and 2200 transfers
- If *s* made outer,

# Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 26$
- Outer should be $r$
- Number of runs is $n = \lceil 200/25 \rceil = 8$
- In each run, $s$ requires 1 seek and 250 transfers
- Total is 8 seeks and 2000 transfers
- $r$ requires 8 seeks and 200 transfers
- Total is 16 seeks and 2200 transfers
- If $s$ made outer, 20 seeks and 2250 transfers

# Indexed Nested-Loop Join

- Indexed version of the block nested-loop algorithm
- Applicable when inner relation has an index on the joining attribute
- For each block $l_r \in r$ and for each record $t_r \in l_r$, use index on $s$ to locate records $t_s \in s$ that satisfies the join condition

# Indexed Nested-Loop Join

- Indexed version of the block nested-loop algorithm
- Applicable when inner relation has an index on the joining attribute
- For each block $l_r \in r$ and for each record $t_r \in l_r$, use index on $s$ to locate records $t_s \in s$ that satisfies the join condition
- Most effective when the join condition is equality

# Cost of Indexed Nested-Loop Join

- Assumption is that $m$ blocks are available in memory

# Cost of Indexed Nested-Loop Join

- Assumption is that $m$ blocks are available in memory
- $m - 1$ blocks from outer relation $r$ can be read at a time

# Cost of Indexed Nested-Loop Join

- Assumption is that $m$ blocks are available in memory
- $m - 1$ blocks from outer relation $r$ can be read at a time
- Assume cost of searching index on $s$ is $c_s$ seeks and $c_t$ transfers
  - Typically, $c_s$ and $c_t$ are height of B+-tree

# Cost of Indexed Nested-Loop Join

- Assumption is that $m$ blocks are available in memory
- $m - 1$ blocks from outer relation $r$ can be read at a time
- Assume cost of searching index on $s$ is $c_s$ seeks and $c_t$ transfers
  - Typically, $c_s$ and $c_t$ are height of B+-tree
- Cost for $r$ is $n = \lceil b_r/(m - 1) \rceil$ seeks and $b_r$ transfers
- Therefore, total is $n + n_r \times c_s$ seeks and $b_r + n_r \times c_t$ transfers

# Cost of Indexed Nested-Loop Join

- Assumption is that $m$ blocks are available in memory
- $m - 1$ blocks from outer relation $r$ can be read at a time
- Assume cost of searching index on $s$ is $c_s$ seeks and $c_t$ transfers
    - Typically, $c_s$ and $c_t$ are height of B+-tree
- Cost for $r$ is $n = \lceil b_r/(m - 1) \rceil$ seeks and $b_r$ transfers
- Therefore, total is $n + n_r \times c_s$ seeks and $b_r + n_r \times c_t$ transfers
- If index is available on both relations, like block nested-loop, smaller relation should be

# Cost of Indexed Nested-Loop Join

- Assumption is that $m$ blocks are available in memory
- $m - 1$ blocks from outer relation $r$ can be read at a time
- Assume cost of searching index on $s$ is $c_s$ seeks and $c_t$ transfers
    - Typically, $c_s$ and $c_t$ are height of B+-tree
- Cost for $r$ is $n = \lceil b_r/(m - 1) \rceil$ seeks and $b_r$ transfers
- Therefore, total is $n + n_r \times c_s$ seeks and $b_r + n_r \times c_t$ transfers
- If index is available on both relations, like block nested-loop, smaller relation should be outer

# Cost of Indexed Nested-Loop Join

- Assumption is that $m$ blocks are available in memory
- $m - 1$ blocks from outer relation $r$ can be read at a time
- Assume cost of searching index on $s$ is $c_s$ seeks and $c_t$ transfers
  - Typically, $c_s$ and $c_t$ are height of B+-tree
- Cost for $r$ is $n = \lceil b_r/(m-1) \rceil$ seeks and $b_r$ transfers
- Therefore, total is $n + n_r \times c_s$ seeks and $b_r + n_r \times c_t$ transfers
- If index is available on both relations, like block nested-loop, smaller relation should be outer
- Generally, all levels of B+-tree are held in memory except the last
  - Then, cost of index search falls to $c_s = c_t = 1$

# Example

- $r$: 40000 tuples at 200 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner
- Outer should be

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner
- Outer should be *s*

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner
- Outer should be *s*
- Number of runs is

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner
- Outer should be *s*
- Number of runs is $n = \lceil 250/25 \rceil = 10$

## Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner
- Outer should be $s$
- Number of runs is $n = \lceil 250/25 \rceil = 10$
- In each run, $r$ requires 1 seek and 1 transfer
- Total is 25000 seeks and 25000 transfers

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner
- Outer should be *s*
- Number of runs is $n = \lceil 250/25 \rceil = 10$
- In each run, *r* requires 1 seek and 1 transfer
- Total is 25000 seeks and 25000 transfers
- *s* requires 10 seeks and 250 transfers

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m = 25$ for outer and enough for inner
- Outer should be *s*
- Number of runs is $n = \lceil 250/25 \rceil = 10$
- In each run, *r* requires 1 seek and 1 transfer
- Total is 25000 seeks and 25000 transfers
- *s* requires 10 seeks and 250 transfers
- Total is 25010 seeks and 25250 transfers

# Merge Join or Sort-Merge Join

- Applicable only when the join condition is *equality*
- Relations need to be sorted according to the joining attribute

# Merge Join or Sort-Merge Join

- Applicable only when the join condition is *equality*
- Relations need to be sorted according to the joining attribute
- Proceed in sorted order on two relations
- If records match, output; otherwise, advance to next record
- Join step is similar to merge step in mergesort

# Merge Join or Sort-Merge Join

- Applicable only when the join condition is *equality*
- Relations need to be sorted according to the joining attribute
- Proceed in sorted order on two relations
- If records match, output; otherwise, advance to next record
- Join step is similar to merge step in mergesort
- If relations are not sorted, secondary index on attributes can be used
- HYBRID MERGE JOIN algorithm merges sorted records in one relation with B+-tree leaves of other relation

# Cost of Merge Join

- Each record is read only once
- Consequently, each block is read only once

# Cost of Merge Join

- Each record is read only once
- Consequently, each block is read only once
- However, blocks may be read in an interleaved manner

# Cost of Merge Join

- Each record is read only once
- Consequently, each block is read only once
- However, blocks may be read in an interleaved manner
- If $m' = m/2$ blocks of memory are available for each relation, cost is $\lceil b_r/m' \rceil + \lceil b_s/m' \rceil$ seeks and $b_r + b_s$ transfers

# Example

- $r$: 40000 tuples at 200 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks

# Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m' = 13$ for each

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m' = 13$ for each
- Number of seeks for *r* is

## Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m' = 13$ for each
- Number of seeks for $r$ is $\lceil 250/13 \rceil = 20$
- Number of seeks for $s$ is

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m' = 13$ for each
- Number of seeks for *r* is $\lceil 250/13 \rceil = 20$
- Number of seeks for *s* is $\lceil 200/13 \rceil = 16$

# Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m' = 13$ for each
- Number of seeks for $r$ is $\lceil 250/13 \rceil = 20$
- Number of seeks for $s$ is $\lceil 200/13 \rceil = 16$
- Total is

# Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is $m' = 13$ for each
- Number of seeks for $r$ is $\lceil 250/13 \rceil = 20$
- Number of seeks for $s$ is $\lceil 200/13 \rceil = 16$
- Total is 36 seeks and 450 transfers

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into *n* partitions: $h : t \in (r \cup s) \rightarrow \{0, \ldots, n-1\}$

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into *n* partitions: $h : t \in (r \cup s) \to \{0, \ldots, n - 1\}$
- *Partitioning* (or *building*) phase partitions the relations and writes them to disk

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into *n* partitions: $h : t \in (r \cup s) \rightarrow \{0, \ldots, n-1\}$
- *Partitioning* (or *building*) phase partitions the relations and writes them to disk
- *Matching* (or *probing*) phase
  - Partition $P_i(s)$ is read into memory
  - A finer-grained in-memory hash may be constructed on $P_i(s)$

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into *n* partitions: $h : t \in (r \cup s) \rightarrow \{0, \ldots, n-1\}$
- *Partitioning* (or *building*) phase partitions the relations and writes them to disk
- *Matching* (or *probing*) phase
  - Partition $P_i(s)$ is read into memory
  - A finer-grained in-memory hash may be constructed on $P_i(s)$
  - Partition $P_i(r)$ is read block by block
  - For each tuple in $P_i(r)$, in-memory hash of $P_i(s)$ is searched

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into *n* partitions: $h : t \in (r \cup s) \rightarrow \{0, \ldots, n-1\}$
- *Partitioning* (or *building*) phase partitions the relations and writes them to disk
- *Matching* (or *probing*) phase
  - Partition $P_i(s)$ is read into memory
  - A finer-grained in-memory hash may be constructed on $P_i(s)$
  - Partition $P_i(r)$ is read block by block
  - For each tuple in $P_i(r)$, in-memory hash of $P_i(s)$ is searched
  - After partition $i$ is done, new partitions are loaded

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into *n* partitions: $h : t \in (r \cup s) \rightarrow \{0, \dots, n-1\}$
- *Partitioning* (or *building*) phase partitions the relations and writes them to disk
- *Matching* (or *probing*) phase
    - Partition $P_i(s)$ is read into memory
    - A finer-grained in-memory hash may be constructed on $P_i(s)$
    - Partition $P_i(r)$ is read block by block
    - For each tuple in $P_i(r)$, in-memory hash of $P_i(s)$ is searched
    - After partition *i* is done, new partitions are loaded
- *s* is called build input
- *r* is called probe input

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into $n$ partitions: $h : t \in (r \cup s) \rightarrow \{0, \ldots, n-1\}$
- *Partitioning* (or *building*) phase partitions the relations and writes them to disk
- *Matching* (or *probing*) phase
    - Partition $P_i(s)$ is read into memory
    - A finer-grained in-memory hash may be constructed on $P_i(s)$
    - Partition $P_i(r)$ is read block by block
    - For each tuple in $P_i(r)$, in-memory hash of $P_i(s)$ is searched
    - After partition $i$ is done, new partitions are loaded
- *s* is called build input
- *r* is called probe input
- Each partition of build relation must be stored in memory
- Build relation should be

# Hash Join

- Applicable only when the join condition is *equality*
- If record $t_r$ and $t_s$ match, they must hash to same value, and thus, only partitions with the same hash value need to be compared
- PARTITION HASH JOIN Hash function $h$ to partition records of *both* relations into *n* partitions: $h : t \in (r \cup s) \rightarrow \{0, \ldots, n-1\}$
- *Partitioning* (or *building*) phase partitions the relations and writes them to disk
- *Matching* (or *probing*) phase
  - Partition $P_i(s)$ is read into memory
  - A finer-grained in-memory hash may be constructed on $P_i(s)$
  - Partition $P_i(r)$ is read block by block
  - For each tuple in $P_i(r)$, in-memory hash of $P_i(s)$ is searched
  - After partition $i$ is done, new partitions are loaded
- *s* is called build input
- *r* is called probe input
- Each partition of build relation must be stored in memory
- Build relation should be smaller

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$

## Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$
- During partitioning, at least 1 block of each partition should fit in memory
- Thus, $m > n$

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$
- During partitioning, at least 1 block of each partition should fit in memory
- Thus, $m > n$
- Combining,

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$
- During partitioning, at least 1 block of each partition should fit in memory
- Thus, $m > n$
- Combining, $m > n > b_s/m$ or, $m > \sqrt{b_s}$
- If not so, then partitioning *cannot* be done in one pass

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$
- During partitioning, at least 1 block of each partition should fit in memory
- Thus, $m > n$
- Combining, $m > n > b_s/m$ or, $m > \sqrt{b_s}$
- If not so, then partitioning *cannot* be done in one pass
- Recursive partitioning is employed

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$
- During partitioning, at least 1 block of each partition should fit in memory
- Thus, $m > n$
- Combining, $m > n > b_s/m$ or, $m > \sqrt{b_s}$
- If not so, then partitioning *cannot* be done in one pass
- Recursive partitioning is employed
- Initially, the number of partitions $n_s$ is chosen to be at most $m$
- However, each partition now does not fit into memory
- Thus, each partition is read and re-partitioned till each smaller partition fits into memory

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$
- During partitioning, at least 1 block of each partition should fit in memory
- Thus, $m > n$
- Combining, $m > n > b_s/m$ or, $m > \sqrt{b_s}$
- If not so, then partitioning *cannot* be done in one pass
- Recursive partitioning is employed
- Initially, the number of partitions $n_s$ is chosen to be at most $m$
- However, each partition now does not fit into memory
- Thus, each partition is read and re-partitioned till each smaller partition fits into memory
- Number of passes is

# Recursive Partitioning

- Each of the $n$ partitions should fit into memory
- If build relation has $b_s$ blocks, each partition has roughly $b_s/n$ blocks
- Assume an available memory of $m$ blocks for build
- If each partition fits, then $m > b_s/n$ or $n > b_s/m$
- During partitioning, at least 1 block of each partition should fit in memory
- Thus, $m > n$
- Combining, $m > n > b_s/m$ or, $m > \sqrt{b_s}$
- If not so, then partitioning *cannot* be done in one pass
- Recursive partitioning is employed
- Initially, the number of partitions $n_s$ is chosen to be at most $m$
- However, each partition now does not fit into memory
- Thus, each partition is read and re-partitioned till each smaller partition fits into memory
- Number of passes is $p = \lfloor \log_m b_s \rfloor$

- Suppose, build input has $b_s = 120$ blocks

# Example

- Suppose, build input has $b_s = 120$ blocks
- Available memory is $m = 12$ blocks

# Example

- Suppose, build input has $b_s = 120$ blocks
- Available memory is $m = 12$ blocks
  - Partitioning produces $120/12 = 10$ blocks that fit
  - Thus, only 1 pass is required

# Example

- Suppose, build input has $b_s = 120$ blocks
- Available memory is $m = 12$ blocks
  - Partitioning produces $120/12 = 10$ blocks that fit
  - Thus, only 1 pass is required
- Available memory is $m = 10$ blocks

# Example

- Suppose, build input has $b_s = 120$ blocks
- Available memory is $m = 12$ blocks
  - Partitioning produces $120/12 = 10$ blocks that fit
  - Thus, only 1 pass is required
- Available memory is $m = 10$ blocks
  - Partitioning produces $120/10 = 12$ blocks that do not fit
  - Thus, 12 blocks are partitioned again to $12/10 = 2$ blocks each
  - This now fits
  - Thus, 2 passes are required

# Example

- Suppose, build input has $b_s = 120$ blocks
- Available memory is $m = 12$ blocks
    - Partitioning produces $120/12 = 10$ blocks that fit
    - Thus, only 1 pass is required
- Available memory is $m = 10$ blocks
    - Partitioning produces $120/10 = 12$ blocks that do not fit
    - Thus, 12 blocks are partitioned again to $12/10 = 2$ blocks each
    - This now fits
    - Thus, 2 passes are required
- Available memory is $m = 4$ blocks

# Example

- Suppose, build input has $b_s = 120$ blocks
- Available memory is $m = 12$ blocks
    - Partitioning produces $120/12 = 10$ blocks that fit
    - Thus, only 1 pass is required
- Available memory is $m = 10$ blocks
    - Partitioning produces $120/10 = 12$ blocks that do not fit
    - Thus, 12 blocks are partitioned again to $12/10 = 2$ blocks each
    - This now fits
    - Thus, 2 passes are required
- Available memory is $m = 4$ blocks
    - Partitioning produces $120/4 = 30$ and then $30/4 = 8$ and then $8/4 = 2$ blocks
    - Thus, 3 passes are required

# Cost of Hash Join

- Initial reading of a relation $i$ requires $b_i$ transfers
- Total number of blocks after partitioning is at most

# Cost of Hash Join

- Initial reading of a relation $i$ requires $b_i$ transfers
- Total number of blocks after partitioning is at most $b_i + n$ as each of $n$ partitions may be partially full
- Writing them back requires

# Cost of Hash Join

- Initial reading of a relation $i$ requires $b_i$ transfers
- Total number of blocks after partitioning is at most $b_i + n$ as each of $n$ partitions may be partially full
- Writing them back requires $(b_i + n)$ transfers
- Reading them again during matching requires

# Cost of Hash Join

- Initial reading of a relation $i$ requires $b_i$ transfers
- Total number of blocks after partitioning is at most $b_i + n$ as each of $n$ partitions may be partially full
- Writing them back requires $(b_i + n)$ transfers
- Reading them again during matching requires $(b_i + n)$ transfers
- Therefore, total number of transfers is $3(b_r + b_s) + 4n$

# Cost of Hash Join

- Initial reading of a relation $i$ requires $b_i$ transfers
- Total number of blocks after partitioning is at most $b_i + n$ as each of $n$ partitions may be partially full
- Writing them back requires $(b_i + n)$ transfers
- Reading them again during matching requires $(b_i + n)$ transfers
- Therefore, total number of transfers is $3(b_r + b_s) + 4n$
- Assume memory buffer of $m$ blocks
- Partitioning requires

# Cost of Hash Join

- Initial reading of a relation $i$ requires $b_i$ transfers
- Total number of blocks after partitioning is at most $b_i + n$ as each of $n$ partitions may be partially full
- Writing them back requires $(b_i + n)$ transfers
- Reading them again during matching requires $(b_i + n)$ transfers
- Therefore, total number of transfers is $3(b_r + b_s) + 4n$
- Assume memory buffer of $m$ blocks
- Partitioning requires $\lceil b_i/m \rceil$ seeks for reading and $\lceil (b_i + n)/m \rceil$ seeks for writing
- Reading $n$ partitions during matching requires $n$ seeks per relation
- Therefore, total number of seeks is
  $\left\lceil \frac{b_r}{m} \right\rceil + \left\lceil \frac{b_s}{m} \right\rceil + \left\lceil \frac{b_r+n}{m} \right\rceil + \left\lceil \frac{b_s+n}{m} \right\rceil + 2n$

- Number of passes in recursive partitioning is

# Cost of Recursive Partitioning

- Number of passes in recursive partitioning is $p = \lfloor \log_m b_s \rfloor$

# Cost of Recursive Partitioning

- Number of passes in recursive partitioning is $p = \lfloor \log_m b_s \rfloor$
- During partitioning, total number of block transfers is

# Cost of Recursive Partitioning

- Number of passes in recursive partitioning is $p = \lfloor \log_m b_s \rfloor$
- During partitioning, total number of block transfers is
  $2p(b_r + n + b_s + n)$

# Cost of Recursive Partitioning

- Number of passes in recursive partitioning is $p = \lfloor \log_m b_s \rfloor$
- During partitioning, total number of block transfers is $2p(b_r + n + b_s + n)$
- Number of seeks during partitioning is

# Cost of Recursive Partitioning

- Number of passes in recursive partitioning is $p = \lfloor \log_m b_s \rfloor$
- During partitioning, total number of block transfers is $2p(b_r + n + b_s + n)$
- Number of seeks during partitioning is $2p(\lceil (b_r + n)/m \rceil + \lceil (b_s + n)/m \rceil)$

# Cost of Recursive Partitioning

- Number of passes in recursive partitioning is $p = \lfloor \log_m b_s \rfloor$
- During partitioning, total number of block transfers is
  $2p(b_r + n + b_s + n)$
- Number of seeks during partitioning is
  $2p(\lceil (b_r + n)/m \rceil + \lceil (b_s + n)/m \rceil)$
- Matching requires

# Cost of Recursive Partitioning

- Number of passes in recursive partitioning is $p = \lfloor \log_m b_s \rfloor$
- During partitioning, total number of block transfers is $2p(b_r + n + b_s + n)$
- Number of seeks during partitioning is $2p(\lceil (b_r + n)/m \rceil + \lceil (b_s + n)/m \rceil)$
- Matching requires $(b_r + n + b_s + n)$ transfers and $n + n$ seeks

# Example

- *r*: 40000 tuples at 200 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block:

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering

# Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$
- Build input should be

# Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of $s$ is $\lceil 200/15 \rceil = 14$
- Size of each partition of $r$ is $\lceil 250/15 \rceil = 17$
- Build input should be $s$

## Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of $s$ is $\lceil 200/15 \rceil = 14$
- Size of each partition of $r$ is $\lceil 250/15 \rceil = 17$
- Build input should be $s$
- Partitioning $s$ requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for $s$ is $(200/10) + (215/10) = 42$

## Example

- $r$: 40000 tuples at 200 per block: 200 blocks
- $s$: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of $s$ is $\lceil 200/15 \rceil = 14$
- Size of each partition of $r$ is $\lceil 250/15 \rceil = 17$
- Build input should be $s$
- Partitioning $s$ requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for $s$ is $(200/10) + (215/10) = 42$
- Partitioning $r$ requires

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$
- Build input should be *s*
- Partitioning *s* requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for *s* is $(200/10) + (215/10) = 42$
- Partitioning *r* requires 1 pass since it is probe input
- Number of transfers for *r* is $250 + (250 + 15) = 515$
- Number of seeks for *r* is $(250/10) + (265/10) = 52$

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$
- Build input should be *s*
- Partitioning *s* requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for *s* is $(200/10) + (215/10) = 42$
- Partitioning *r* requires 1 pass since it is probe input
- Number of transfers for *r* is $250 + (250 + 15) = 515$
- Number of seeks for *r* is $(250/10) + (265/10) = 52$
- Matching phase requires reading all blocks of *s* and *r*:

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$
- Build input should be *s*
- Partitioning *s* requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for *s* is $(200/10) + (215/10) = 42$
- Partitioning *r* requires 1 pass since it is probe input
- Number of transfers for *r* is $250 + (250 + 15) = 515$
- Number of seeks for *r* is $(250/10) + (265/10) = 52$
- Matching phase requires reading all blocks of *s* and *r*: $215 + 265 = 480$
- Number of seeks in matching phase is

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$
- Build input should be *s*
- Partitioning *s* requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for *s* is $(200/10) + (215/10) = 42$
- Partitioning *r* requires 1 pass since it is probe input
- Number of transfers for *r* is $250 + (250 + 15) = 515$
- Number of seeks for *r* is $(250/10) + (265/10) = 52$
- Matching phase requires reading all blocks of *s* and *r*: $215 + 265 = 480$
- Number of seeks in matching phase is $15 + 15 = 30$

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$
- Build input should be *s*
- Partitioning *s* requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for *s* is $(200/10) + (215/10) = 42$
- Partitioning *r* requires 1 pass since it is probe input
- Number of transfers for *r* is $250 + (250 + 15) = 515$
- Number of seeks for *r* is $(250/10) + (265/10) = 52$
- Matching phase requires reading all blocks of *s* and *r*: $215 + 265 = 480$
- Number of seeks in matching phase is $15 + 15 = 30$
- Total is

## Example

- *r*: 40000 tuples at 200 per block: 200 blocks
- *s*: 25000 tuples at 100 per block: 250 blocks
- Available memory is 15 for partitions and $m = 10$ for buffering
- Number of partitions is $n = 15$
- Size of each partition of *s* is $\lceil 200/15 \rceil = 14$
- Size of each partition of *r* is $\lceil 250/15 \rceil = 17$
- Build input should be *s*
- Partitioning *s* requires $200 + (200 + 15) = 415$ transfers
- Number of seeks for *s* is $(200/10) + (215/10) = 42$
- Partitioning *r* requires 1 pass since it is probe input
- Number of transfers for *r* is $250 + (250 + 15) = 515$
- Number of seeks for *r* is $(250/10) + (265/10) = 52$
- Matching phase requires reading all blocks of *s* and *r*:
  $215 + 265 = 480$
- Number of seeks in matching phase is $15 + 15 = 30$
- Total is 124 seeks and 1410 transfers

# Other Operations

- Complex join

# Other Operations

- Complex join
  - Use block nested-loop for conjunctive and/or disjunctive selection
  - If only equality, union/intersection of results of merge/hash join may be used

# Other Operations

- Complex join
  - Use block nested-loop for conjunctive and/or disjunctive selection
  - If only equality, union/intersection of results of merge/hash join may be used
- Outer join

# Other Operations

- Complex join
  - Use block nested-loop for conjunctive and/or disjunctive selection
  - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
  - Nested-loop algorithms require almost no modification
  - Select all records while scanning in merge join
  - If $r \supset\!\!\!\bowtie s$, $r$ should be the

# Other Operations

- Complex join
    - Use block nested-loop for conjunctive and/or disjunctive selection
    - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
    - Nested-loop algorithms require almost no modification
    - Select all records while scanning in merge join
    - If $r \sqsupset\!\!\bowtie s$, $r$ should be the probe relation
    - If $r$ is the build relation, keep track of which records in hash index have been used; output all non-used records
    - If $r \sqsupset\!\!\bowtie\!\!\sqsubset s$, use both techniques

# Other Operations

- Complex join
    - Use block nested-loop for conjunctive and/or disjunctive selection
    - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
    - Nested-loop algorithms require almost no modification
    - Select all records while scanning in merge join
    - If $r \; ⟕ \; s$, $r$ should be the probe relation
    - If $r$ is the build relation, keep track of which records in hash index have been used; output all non-used records
    - If $r \; ⟗ \; s$, use both techniques
- Set operations

# Other Operations

- Complex join
    - Use block nested-loop for conjunctive and/or disjunctive selection
    - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
    - Nested-loop algorithms require almost no modification
    - Select all records while scanning in merge join
    - If $r \sqsupset\!\!\bowtie s$, $r$ should be the probe relation
    - If $r$ is the build relation, keep track of which records in hash index have been used; output all non-used records
    - If $r \sqsupset\!\!\bowtie\!\!\sqsubset s$, use both techniques
- Set operations
    - If relations are sorted, scan in order
    - Build hash index on one relation; test records from other relation

# Other Operations

- Complex join
  - Use block nested-loop for conjunctive and/or disjunctive selection
  - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
  - Nested-loop algorithms require almost no modification
  - Select all records while scanning in merge join
  - If $r \sqsupset\!\!\bowtie s$, $r$ should be the probe relation
  - If $r$ is the build relation, keep track of which records in hash index have been used; output all non-used records
  - If $r \sqsupset\!\!\bowtie\!\sqsubset s$, use both techniques
- Set operations
  - If relations are sorted, scan in order
  - Build hash index on one relation; test records from other relation
- Aggregation with grouping

# Other Operations

- Complex join
    - Use block nested-loop for conjunctive and/or disjunctive selection
    - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
    - Nested-loop algorithms require almost no modification
    - Select all records while scanning in merge join
    - If $r \sqsupset\!\bowtie s$, $r$ should be the probe relation
    - If $r$ is the build relation, keep track of which records in hash index have been used; output all non-used records
    - If $r \sqsupset\!\bowtie\!\sqsubset s$, use both techniques
- Set operations
    - If relations are sorted, scan in order
    - Build hash index on one relation; test records from other relation
- Aggregation with grouping
    - Use hashing to organize into groups; then apply aggregation

# Other Operations

- Complex join
  - Use block nested-loop for conjunctive and/or disjunctive selection
  - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
  - Nested-loop algorithms require almost no modification
  - Select all records while scanning in merge join
  - If $r \sqsupset\!\!\bowtie s$, $r$ should be the probe relation
  - If $r$ is the build relation, keep track of which records in hash index have been used; output all non-used records
  - If $r \sqsupset\!\!\bowtie\!\!\sqsubset s$, use both techniques
- Set operations
  - If relations are sorted, scan in order
  - Build hash index on one relation; test records from other relation
- Aggregation with grouping
  - Use hashing to organize into groups; then apply aggregation
- Duplicate detection and elimination

# Other Operations

- Complex join
  - Use block nested-loop for conjunctive and/or disjunctive selection
  - If only equality, union/intersection of results of merge/hash join may be used
- Outer join
  - Nested-loop algorithms require almost no modification
  - Select all records while scanning in merge join
  - If $r \rtimes s$, $r$ should be the probe relation
  - If $r$ is the build relation, keep track of which records in hash index have been used; output all non-used records
  - If $r \bowtie s$, use both techniques
- Set operations
  - If relations are sorted, scan in order
  - Build hash index on one relation; test records from other relation
- Aggregation with grouping
  - Use hashing to organize into groups; then apply aggregation
- Duplicate detection and elimination
  - Use hashing or sorting