

COMP9120

Week 11: Query Processing and Evaluation

Semester 1, 2022

Professor Athman Bouguettaya
School of Computer Science



THE UNIVERSITY OF
SYDNEY



Acknowledgement of Country

I would like to acknowledge the Traditional Owners of Australia and recognise their continuing connection to land, water and culture. I am currently on the land of the Darug people and pay my respects to their Elders, past, present and emerging.

I further acknowledge the Traditional Owners of the country on which you are on and pay respects to their Elders, past, present and future.

› **Quiz** next week (week 12)

- Released on Thursday, 19 May at 21:00
- Due date Friday, 20 May at 23:59
- Quiz duration → 1.5 hours
 - Only one attempt is allowed
- 10 MCQ questions (covering week 6, 8, 9, 10, 11 contents)
- 3 essay questions
 - Integrity Constraint
 - Transaction
 - Normalization
 - Storage and Indexing
 - Query Processing



COMMONWEALTH OF AUSTRALIA

Copyright Regulations 1969

WARNING

This material has been reproduced and communicated to you by or on behalf of the University of Sydney pursuant to Part VB of the Copyright Act 1968 (**the Act**).

The material in this communication may be subject to copyright under the Act. Any further copying or communication of this material by you may be the subject of copyright protection under the Act.

Do not remove this notice.



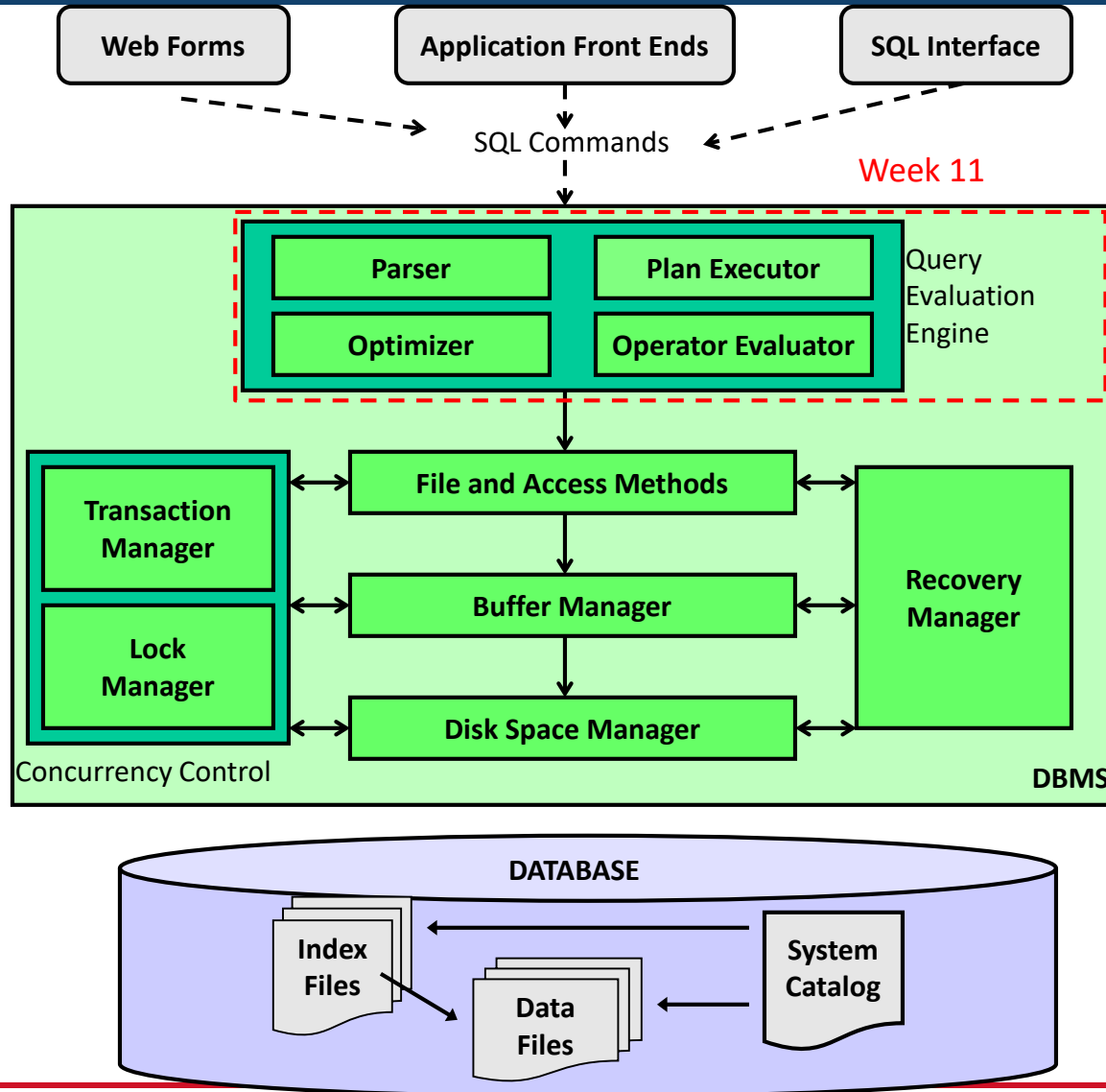
- › **Basic Steps in Query Processing**
- › Heuristic-based optimization
- › Data physical organization: Sorting
- › Cost estimate optimization

› Problems of processing queries:

- How is a query transformed in a form understood by the DBMS? (processing phase)
 - What is the best strategy to execute a query? (optimization phase)
 - Criteria used to select a strategy: I/Os and to a lesser extent, cpu processing (evaluation phase)
-



Internal Structure of a DBMS



Basic Steps in Query Processing

› Step 1: Parsing and Translation

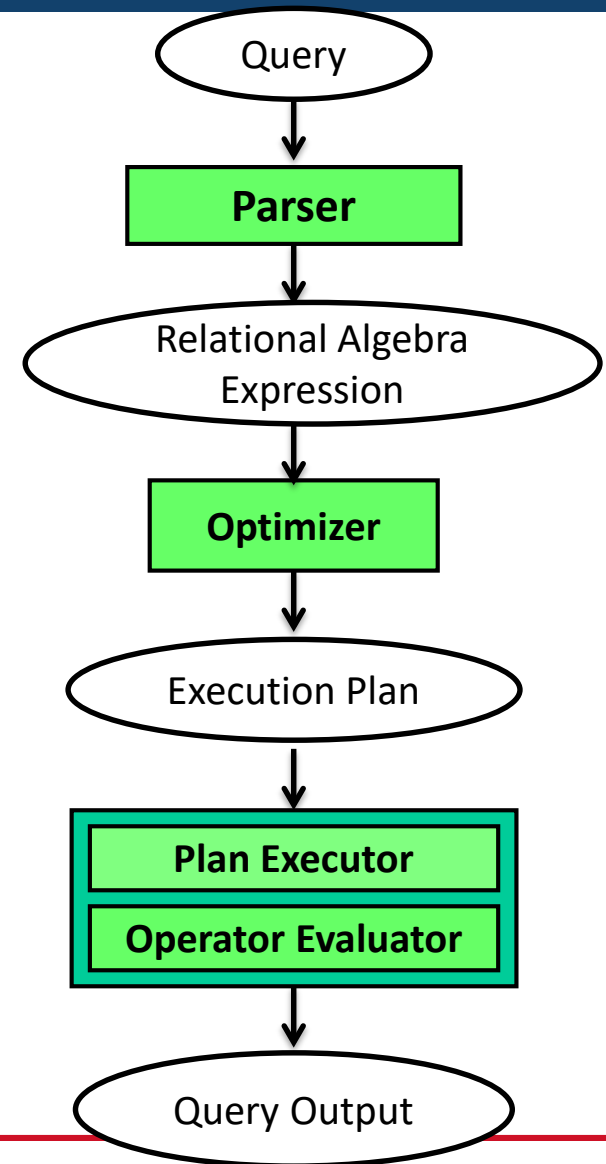
- Check for syntactic and semantic errors
- Translate the SQL query into relational algebra.
- Query Rewriter (e.g. views) is replaced with actual sub-query on relations

› Step 2: Query Optimization

- Amongst all equivalent query-evaluation plans choose the one with the lowest cost.
- Use heuristics to optimize at the relational algebra level
- Select a query execution strategy based on cost estimate

› Step 3: Query Execution

- Depends on available storage structure and access methods

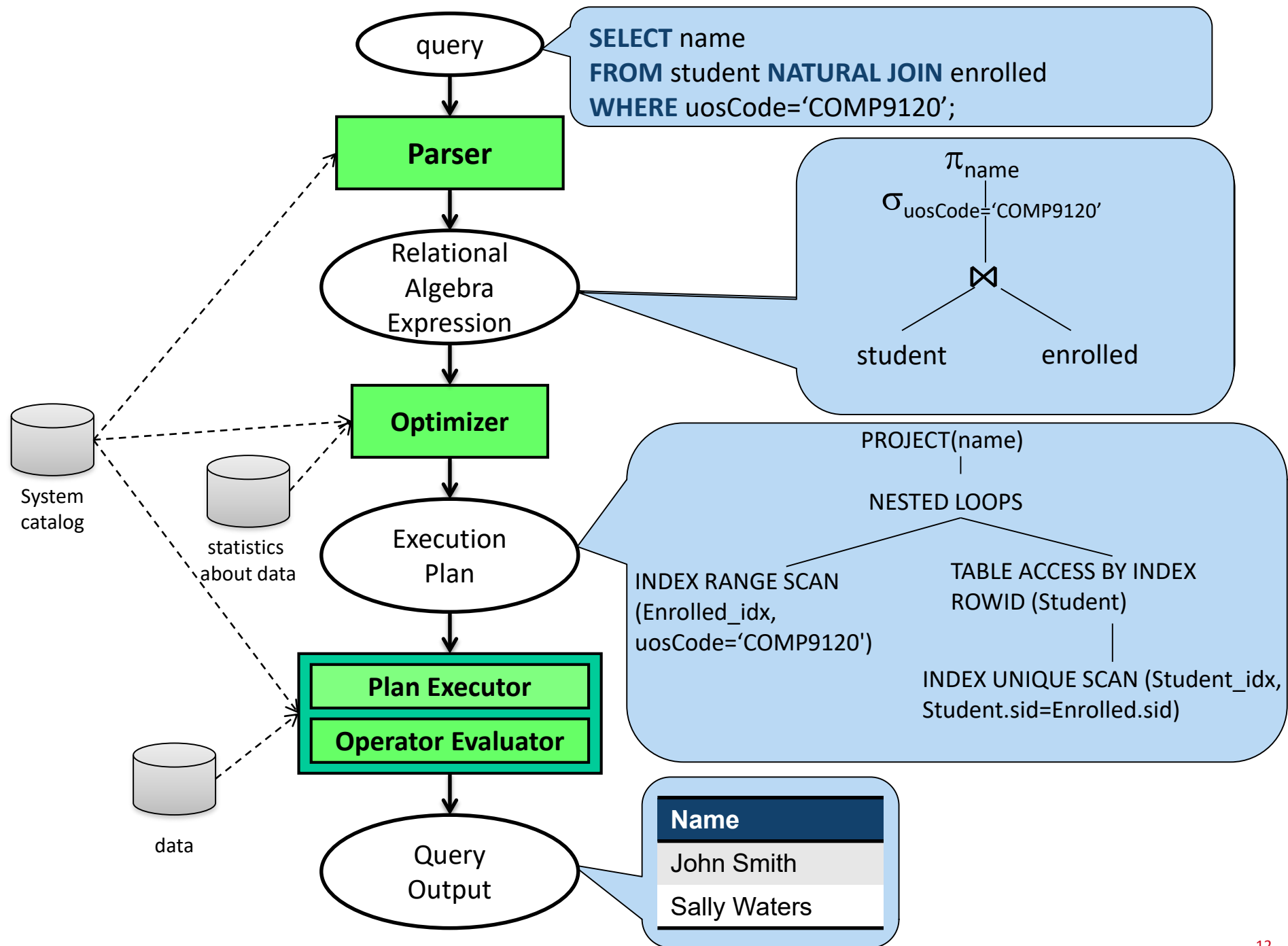


Relational model:

- Unlike early database models (e.g., hierarchical/network) where optimization was mostly left to the application programmer, relational models are based on the declarative (SQL queries - user interface) and equivalent algebraic (execution) models. These lend themselves to computer based optimization.
 - Provides more abstraction (through the declarative model) to make databases more usable by the general users.
 - High level abstraction also makes it hard for non-expert users to optimize queries.
-

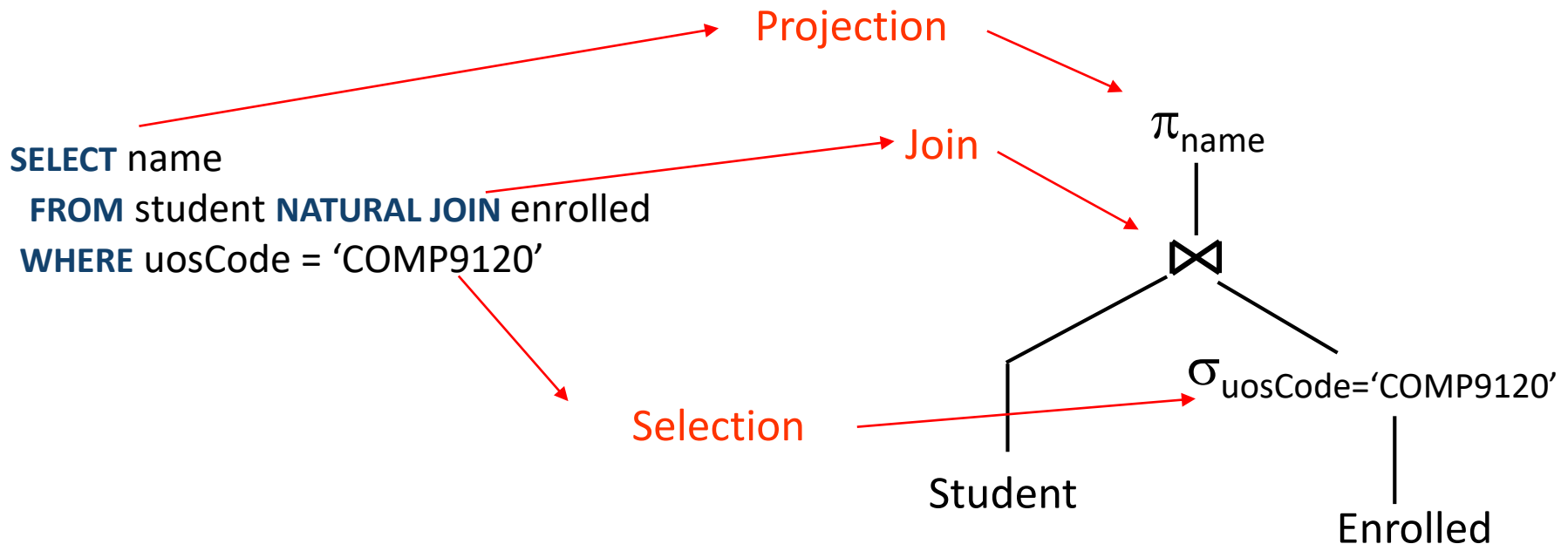
Output of each phase:

- Parser: A parse tree consisting of a SFW (Select-From-Where) expression
 - Processing: converts the parse tree into a **logical query plan**
 - Optimization:
 - heuristics: **efficient** logical query plans
 - cost estimate: **physical query plan**
-



Step 1: Parsing and Translation

- › SQL query gets translated into **Relational Algebra (RA)**, which is represented by a **logical query plan** (also called expression tree).
- Operators have one or more input sets and return one output set.
- Leafs correspond to (base) tables.
- › Example:



Query optimization:

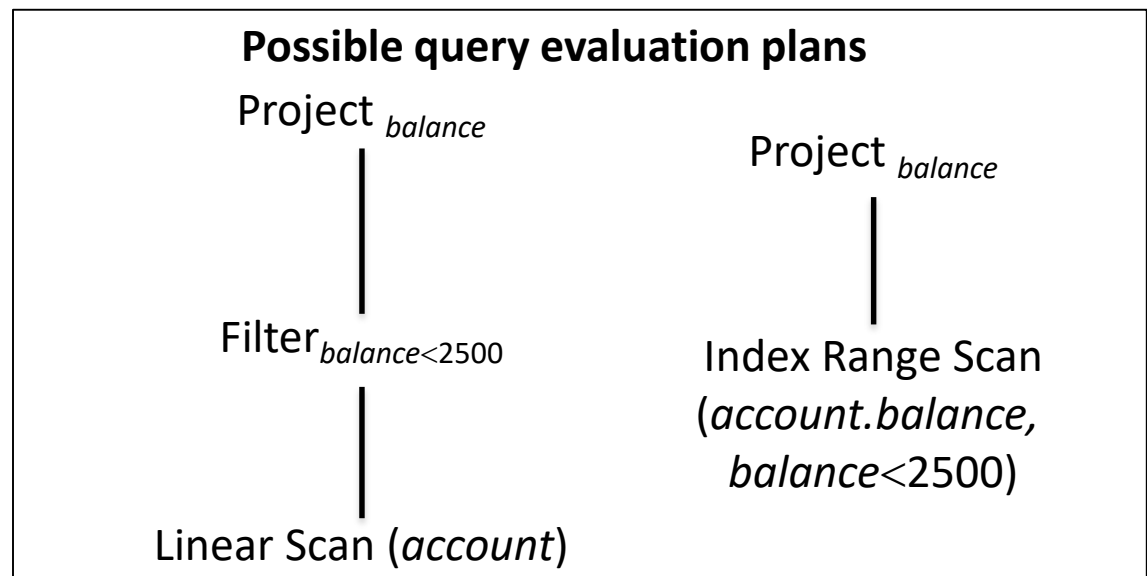
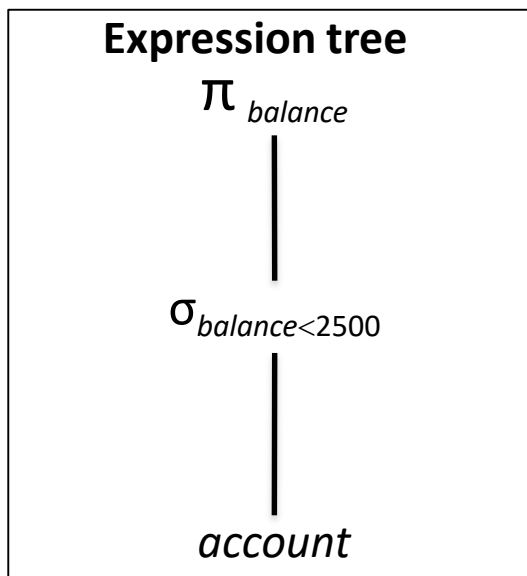
Two (2) main thrusts:

- *Heuristic* rules to rearrange operations in a query tree: output is an efficient ***logical query plan***, followed by
- *Cost estimate* of different execution strategies to pick the one with minimal cost: output is a ***physical query plan***

Heuristics: *minimize* the size of intermediate results (relations) using the manipulation of *equivalent* algebraic expressions.

Cost estimate: *minimize* the number of disk I/Os.

- › An annotated expression tree which specifies a *detailed evaluation* strategy with *physical operators* is called an **evaluation plan** or **query plan**
 - **RA operators** are logical operators
 - **Physical operators** show how query is evaluated
- › Multiple possible query evaluation plans for the expression tree are considered by the query optimizer



› Data organization and access

- **Organization:**

- Heap file (unsorted)
- Sorted file
- Mixed

- **Access**

- Heap files can be accessed
 - with a linear/file/table scan
 - Indexes
- Sorted files
 - Specialized access algorithms
 - Indexes

› An **access path** is a method of retrieving tuples

Step 3: Evaluation of Expressions

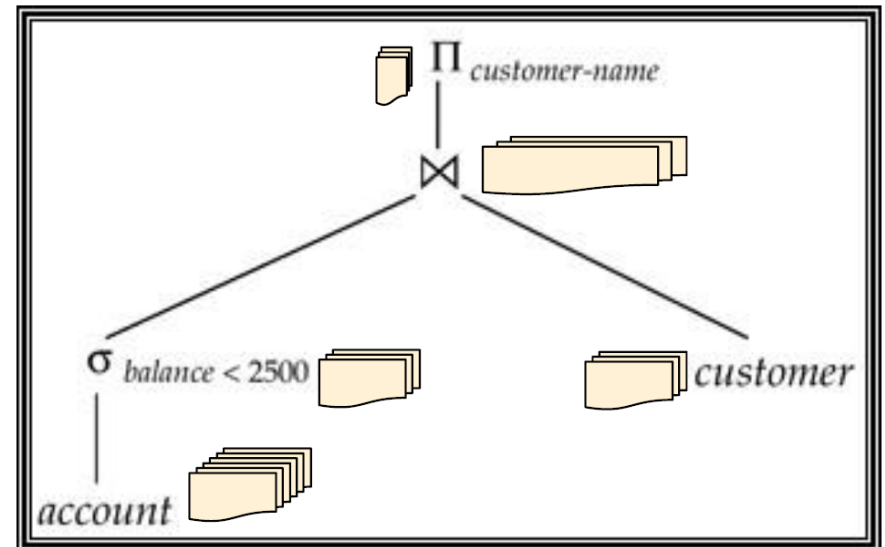
- › **Materialization** (also: **set-at-a-time**):
simply evaluate one operation at a time. The result of each evaluation is materialized (stored) in a temporary relation for subsequent use.
- › **Pipelining** (also: **tuple-at-a-time**):
evaluate several operations simultaneously in a pipeline

› **Materialized evaluation:** evaluate one operation at a time, starting at the lowest-level. Use intermediate materialized (stored) results into temporary tables to evaluate next-level operations.

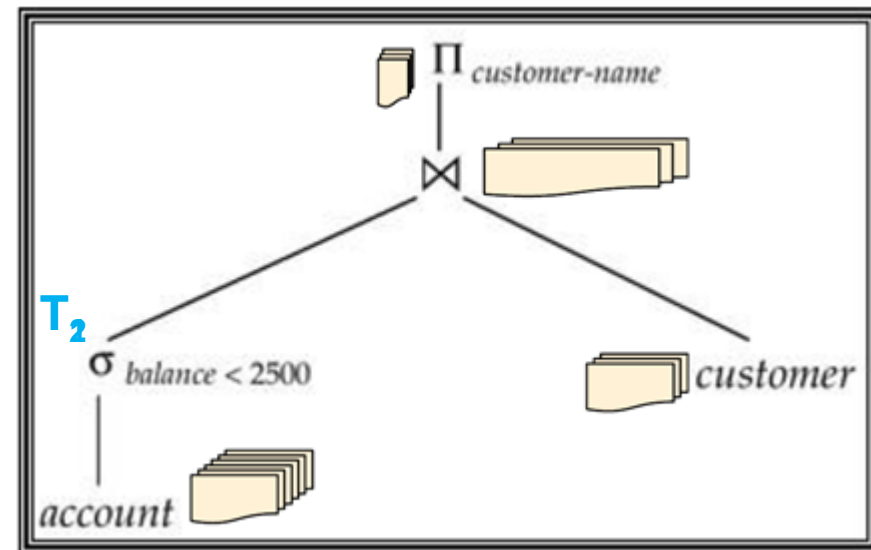
› E.g., in figure below:

1. compute and store new table for $\sigma_{balance < 2500}(account)$
2. Compute and store result of materialized result joined with **customer**
3. Read back new materialized result and compute the projections on **customer-name**.

- Materialized evaluation is always applicable
- Costs can be quite high



- › **Pipelined evaluation:** evaluate several operations simultaneously, passing the results of one operation on to the next.
- › E.g., in the expression tree
 - 1. Find tuple matching $\sigma_{balance < 2500}(account)$
 - a) Join matching tuple with tuples of customer until a new tuple is generated
 - b) Project customer name from joined tuple
 - c) Repeat for all join results
 - 2. Repeat for next selection result
- › Much cheaper than materialization:
 - no need to store a temporary table to disk.





- › Basic Steps in Query Processing
- › **Heuristic-based optimization**
- › Data physical organization: Sorting
- › Cost estimate optimization

Equivalence of expressions

Note:

- We can transform any tuple calculus expression (i.e., SQL query) to an equivalent algebraic expression
- We then make sure the equivalent algebraic expression is executed efficiently
- Heuristic optimization is mostly concerned with unary operations (e.g., selection, projection)
- Strategy is to always pick a sequence of operations that would minimize size of intermediate results

Consider the following query:

“find the assets and names of all banks which have depositors living in Sydney”.
assume we have three relations: *customer*, *deposit*, and *branch*.

Schema:

Table 1: DEPOSIT

<i>Branchname</i>	<i>Account#</i>	<i>Customername</i>	<i>Balance</i>
-------------------	-----------------	---------------------	----------------

Table 2: CUSTOMER

<i>Customername</i>	<i>Street</i>	<i>Customercity</i>
---------------------	---------------	---------------------

Table 3: BRANCH

<i>Branchname</i>	<i>Assets</i>	<i>Branchcity</i>
-------------------	---------------	-------------------

This query is equivalent to the following (*unoptimized*) algebraic expression:

$$\Pi_{\text{branchname, assets}} (\sigma_{\text{customercity}=\text{Sydney}} (\text{Customer} \bowtie \text{Deposit} \bowtie \text{Branch}))$$

The join of the three relations may yield a large relation that may not fit in memory.

Note that:

We only need a handful of tuples to begin with (those with customercity = Sydney). Further, we are only interested in two attributes (branchname and assets).

Question: Could we make the evaluation more "intelligent"?

Answer: YES!

but how?

Using some rearrangements of the operations (algebraic manipulation)

How do we insure the rearrangement is equivalent to the original arrangement?

Using knowledge about the rules governing algebraic operations.

For instance: the previous expression is equivalent to:

$$\Pi_{\text{branchname, assets}}((\sigma_{\text{customercity}=\text{Sydney}}(\text{Customer})) \bowtie \text{Deposit} \bowtie \text{Branch})$$

Example:

Query: “find the *assets* and *names* of all banks which have depositors living in Sydney and have a *balance* of more than \$500”.

This query is equivalent to the following algebraic expression:

$$\Pi_{\text{branchname, assets}} (\sigma_{\text{customercity}=\text{Sydney} \wedge \text{balance} > 500} (\text{Customer} \bowtie \text{Deposit} \bowtie \text{Branch}))$$

Problem: cannot do the selection on *Customer* only, because *balance* is an attribute of *Deposit*.

What is the solution then?

Do the selection *after* doing a join on *Customer* and *Deposit*. the resulting expression is then:

$$\Pi_{\text{branchname, assets}} (\sigma_{\text{customercity}=\text{Sydney} \wedge \text{balance} > 500} (\text{Customer} \bowtie \text{Deposit}) \bowtie \text{Branch})$$

The intermediate result has now been reduced.

Can we do even better?

The answer is YES!

but how?

Break up the selection condition into two selections and we get:

$$\Pi_{\text{branchname, assets}} (\sigma_{\text{customercity}=\text{Sydney}} (\sigma_{\text{balance} > 500} (\text{Customer} \bowtie \text{Deposit})) \bowtie \text{Branch})$$

And then move the second selection past the first join:

$$\Pi_{\text{branchname, assets}} (\sigma_{\text{customercity}=\text{Sydney}} (\text{Customer}) \bowtie \sigma_{\text{balance} > 500} (\text{Deposit})) \bowtie \text{Branch})$$

Projection optimization:

Whenever possible, do a projection as soon as possible. Consider the query in the first example:

$$\Pi_{\text{branchname, assets}} ((\sigma_{\text{customercity}=\text{Sydney}} (\text{Customer}) \bowtie \text{Deposit}) \bowtie \text{Branch})$$

When we compute the subexpression:

$$(\sigma_{\text{customercity}=\text{Sydney}} (\text{Customer}) \bowtie \text{Deposit})$$

We would like to eliminate those attributes that will not play any role in the remaining operations.

In the example above, the attribute *branchname* is the only attribute we need in the remaining operations. Do a projection on the subresult.

The transformation is as follows:

$$\Pi_{\text{branchname, assets}} (\Pi_{\text{branchname}} (\sigma_{\text{customercity}=\text{Sydney}} (\text{Customer}) \bowtie \text{Deposit}) \bowtie \text{Branch})$$

Rules for equivalent transformations

The internal form of a query is usually implemented using a *logical query plan*.

Heuristics based optimization consists of applying rules that yield equivalent transformations to obtain more efficient logical query plans.

In the previous optimization steps, we intuitively came up with transformation rules. What are the formal algebraic transformation rules that enabled us to do the previous transformations?

Example of rules for algebraic transformations:

1. Commutative laws for joins:

$$(E1 \bowtie E2) = (E2 \bowtie E1)$$

2. Associative laws for joins

$$(E1 \bowtie E2) \bowtie E3 = E1 \bowtie (E2 \bowtie E3)$$

3. Cascade of projections: if attributes $A1, \dots, An$ are a subset of $B1, \dots, Bn$ then

$$\Pi_{A1, \dots, An} (\Pi_{B1, \dots, Bn} (R)) = \Pi_{A1, \dots, An} (R)$$

4. Cascade of selections

$$\sigma_{\Theta_1} (\sigma_{\Theta_2} (R)) = \sigma_{\Theta_2} (\sigma_{\Theta_1} (R)) = \sigma_{\Theta_1 \wedge \Theta_2} (R)$$

5. Commuting selections and projections: If formula Θ involves only attributes $B1, \dots, Bn$ that are in the set $A1, \dots, An$ then

$$\Pi_{A1, \dots, An} (\sigma_{\Theta} (R)) = \sigma_{\Theta} (\Pi_{A1, \dots, An} (R))$$

If formula Θ involves attributes $B1, \dots, Bn$ that are not in the set $A1, \dots, An$ then

$$\Pi_{A1, \dots, An} (\sigma_{\Theta} (R)) = \Pi_{A1, \dots, An} (\sigma_{\Theta} (\Pi_{A1, \dots, An, B1, \dots, Bn} (R)))$$

Short 5 mn break:

please stand up, stretch, and move around



THE UNIVERSITY OF
SYDNEY



- › Basic Steps in Query Processing
- › Heuristic-based optimization
- › **Data physical organization: Sorting**
- › Cost estimate optimization

- › Importance of sorting for DBMS:
 - SQL queries can specify that the output be sorted (**ORDER BY**)
 - Some SQL operators (e.g., selection, join) can be implemented efficiently if the input are sorted
- › For small tables that fit in memory, techniques like *QuickSort* can be used.
- › Challenge: sort 10Gb of data with 4Gb of RAM...
=> **external merge-sort**

Let B denote memory size (in pages).

1. Create sorted **runs**. (A run is a sorted subset of records)

Let i be 0 initially. **Repeatedly** do the following till the end of the file:

- (a) Read B pages of records from disk into memory
- (b) Sort the in-memory pages
- (c) Write sorted data to run R_i ; increment i by 1.

Let the final value of i be $m (= \lceil N / B \rceil)$; there are m sorted runs

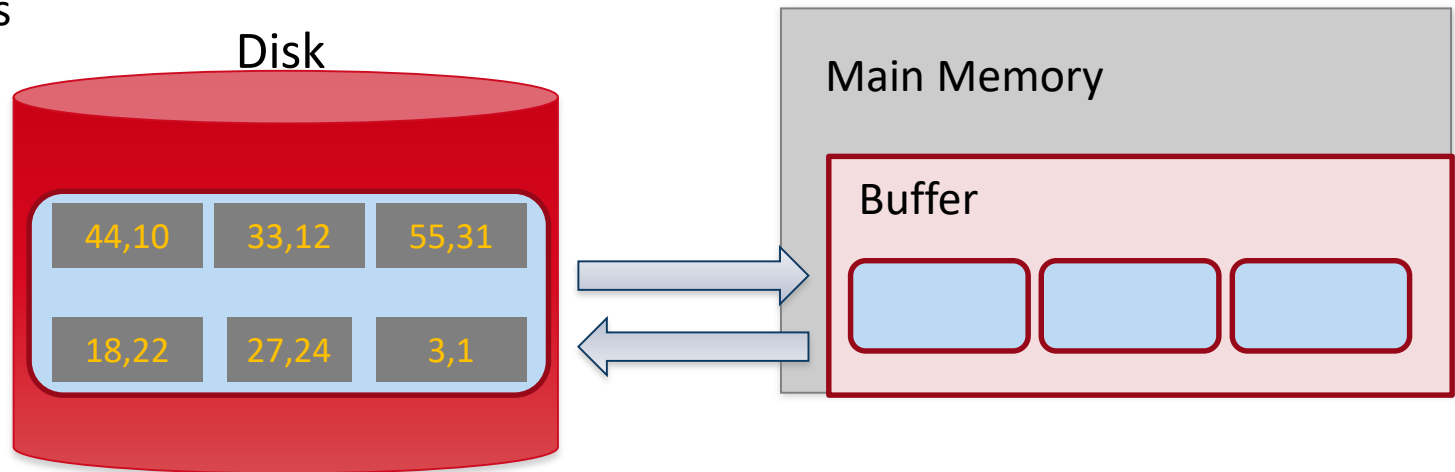
2. Merge each contiguous group of $B-1$ runs into 1 run: $(B-1)$ -way merge.

3. After each merge pass, the number of runs is reduced by a factor of $B-1$. If $m \geq B$, several merge passes are required.

Example Step 1: Create Sorted Runs

Example:

- A file consists of 6 pages
- 3 buffer pages

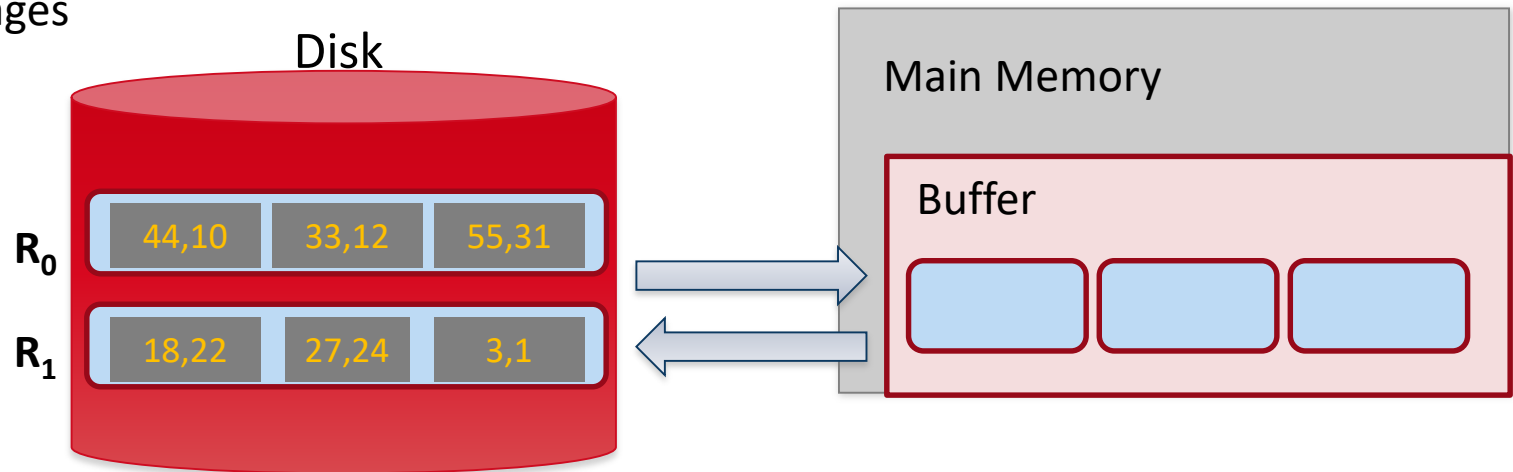


1. Split into runs small enough to **sort in memory**

Example Step 1: Create Sorted Runs

Example:

- A file consists of 6 pages
- 3 buffer pages

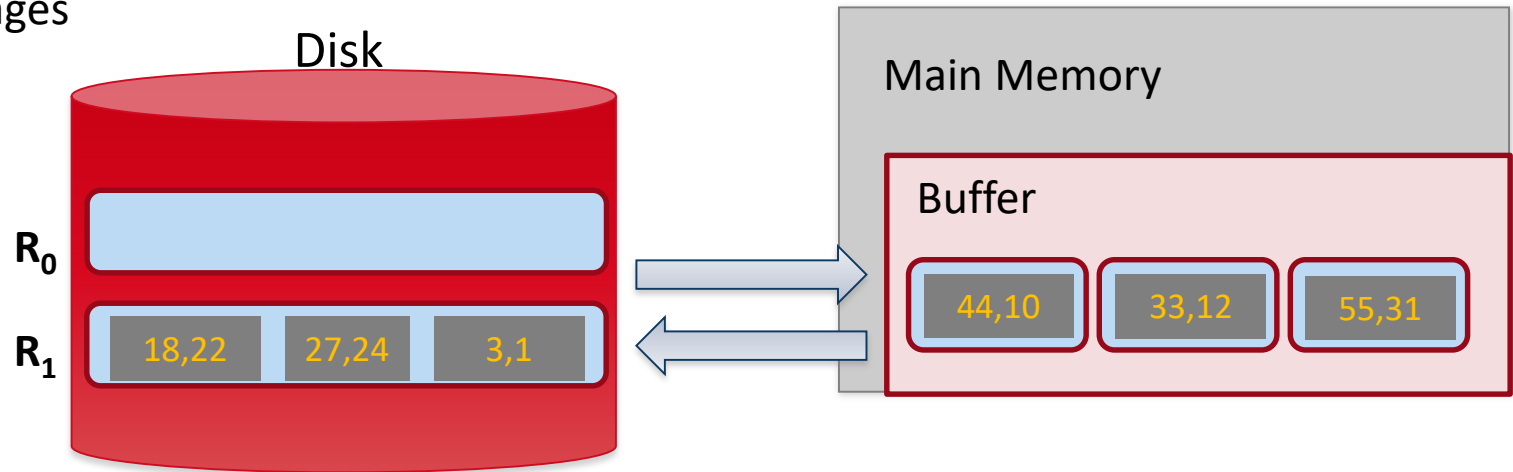


1. Split into runs small enough to **sort in memory**

Example Step 1: Create Sorted Runs

Example:

- A file consists of 6 pages
- 3 buffer pages

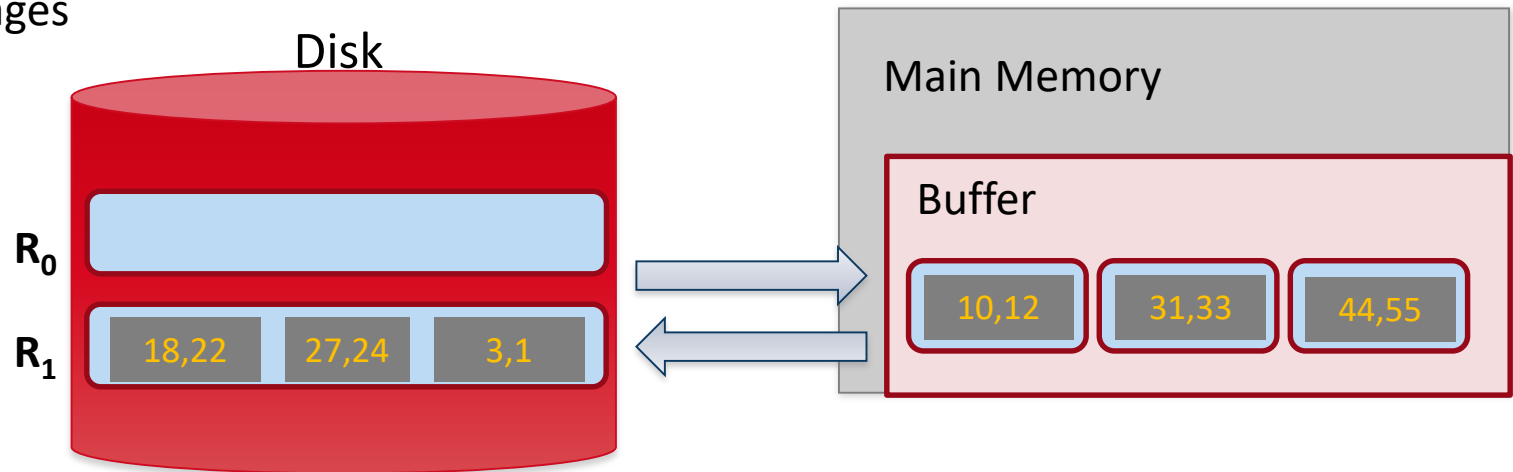


2. Load run R_0 into main memory

Example Step 1: Create Sorted Runs

Example:

- A file consists of 6 pages
- 3 buffer pages

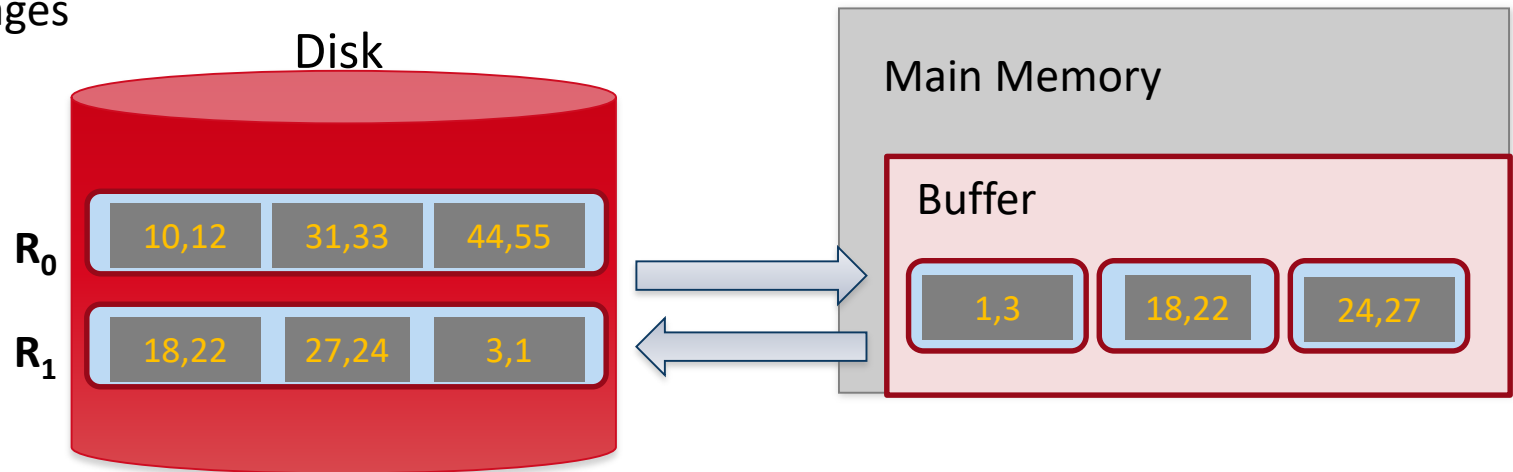


3. Sort run R_0 in main memory, and write back R_0 to disk

Example Step 1: Create Sorted Runs

Example:

- A file consists of 6 pages
- 3 buffer pages

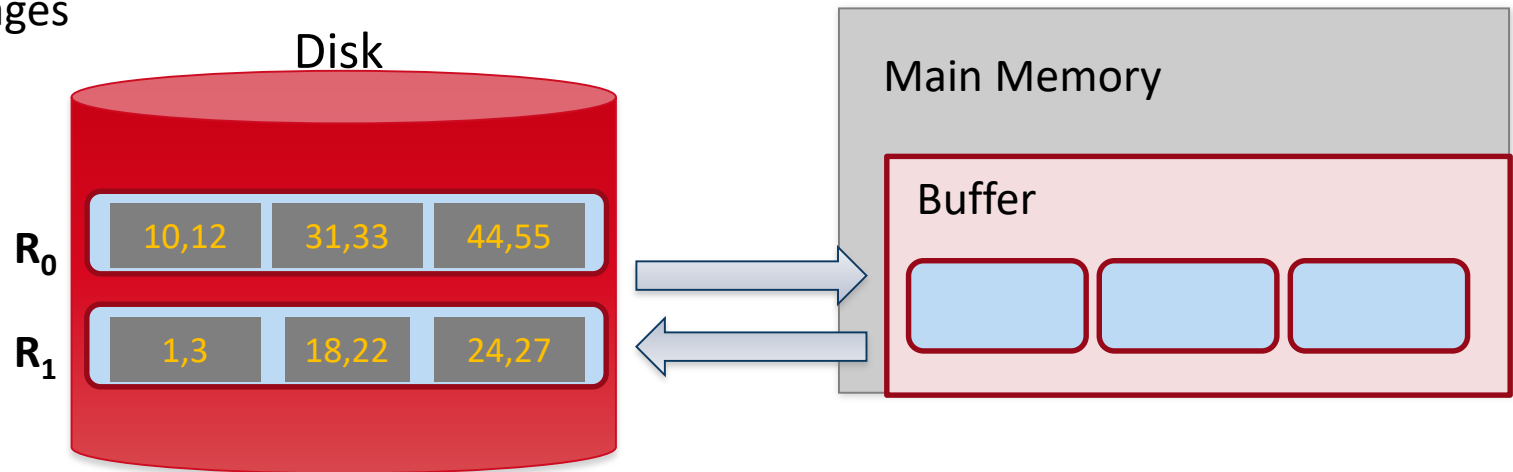


4. Similarly, load R_1 into main memory, sort it, and write it back to disk

Example Step 1: Create Sorted Runs

Example:

- A file consists of 6 pages
- 3 buffer pages



5. Now, we have sorted runs, and we next run the second step of the external merge sort algorithm

Let B denote memory size (in pages).

1. Create sorted **runs**. (A run is a sorted subset of records)

Let i be 0 initially. **Repeatedly** do the following till the end of the file:

- (a) Read B pages of records from disk into memory
- (b) Sort the in-memory blocks
- (c) Write sorted data to run i ; increment i by 1.

Let the final value of i be $m (= \lceil N / B \rceil)$; there are m sorted runs

2. Merge each contiguous group of $B-1$ runs into 1 run: $(B-1)$ -way merge.

- i. Use $B-1$ pages of memory to buffer input runs, and 1 page to buffer output.
Read the first page of each run into its buffer page

ii. repeat

1. Select the first record (in sort order) among all input buffer pages
2. Write the record to the output buffer. If output is full, write it to disk.
3. **If** this is the last record of the input buffer page
then read the next page (if any) of the run into the buffer.

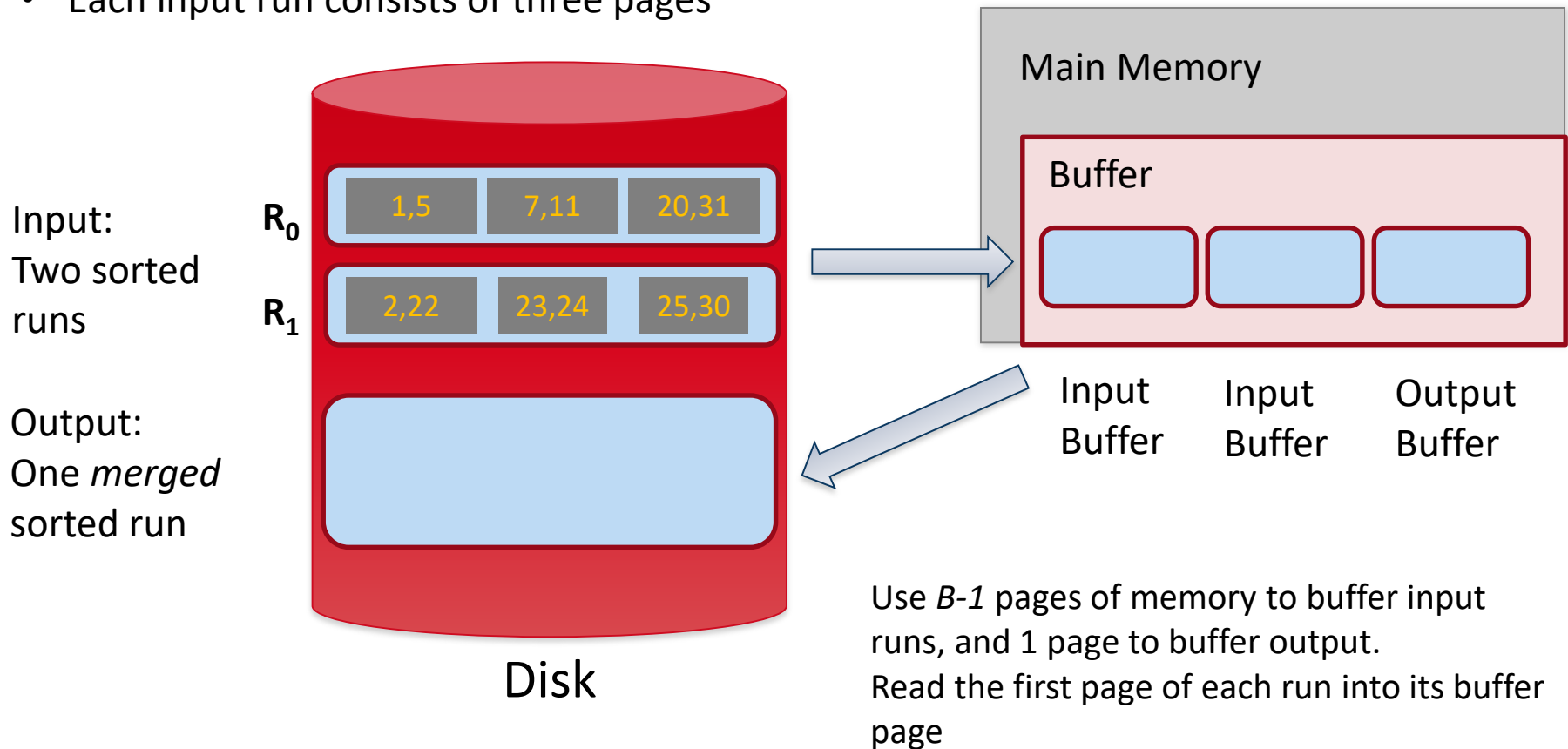
until all input buffer pages are empty:

3. After each merge pass, the number of runs is reduced by a factor of $B-1$. If $m \geq B$, several merge passes are required.

Example Step-2: Merge Sorted Runs

Example:

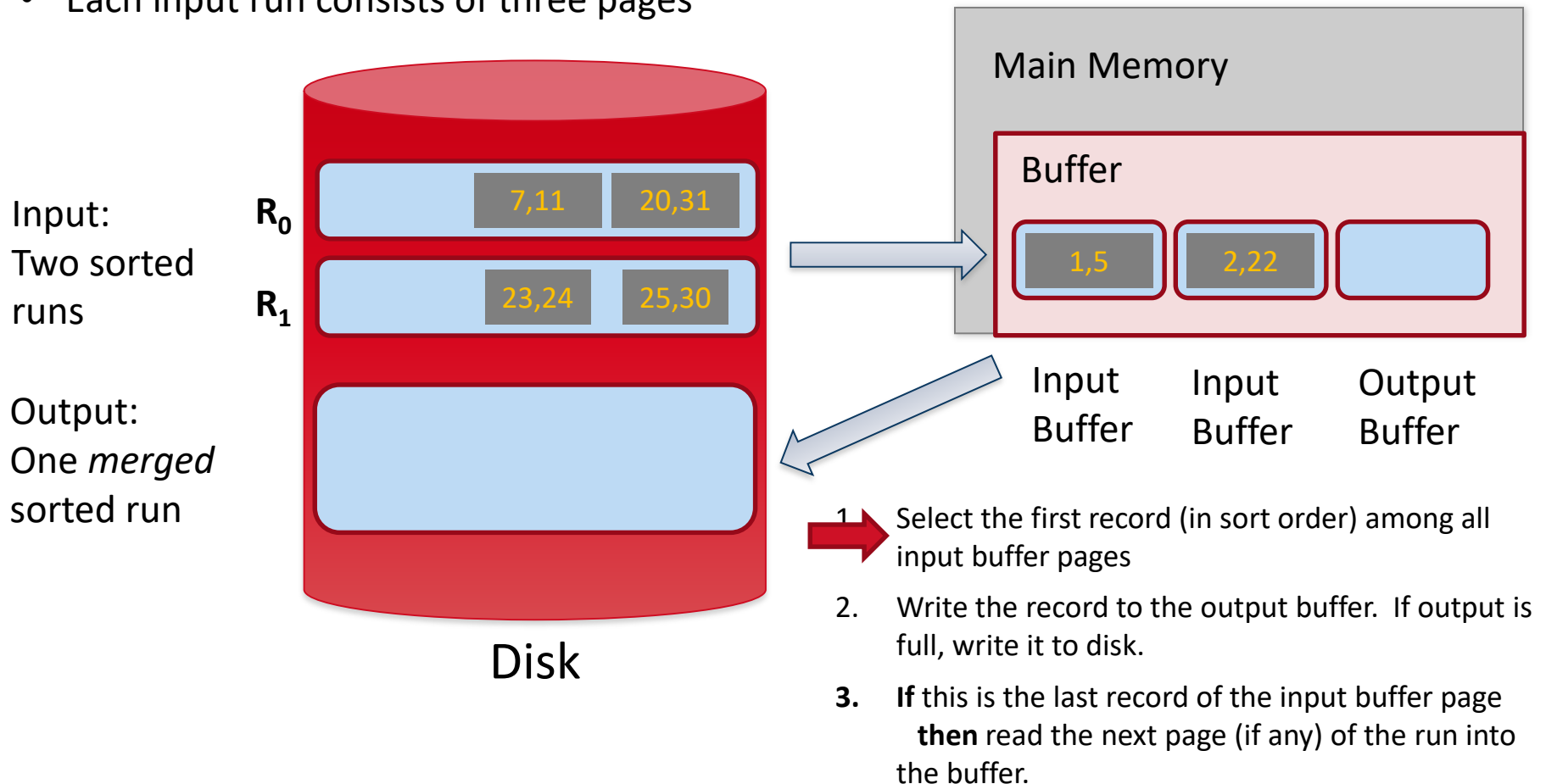
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

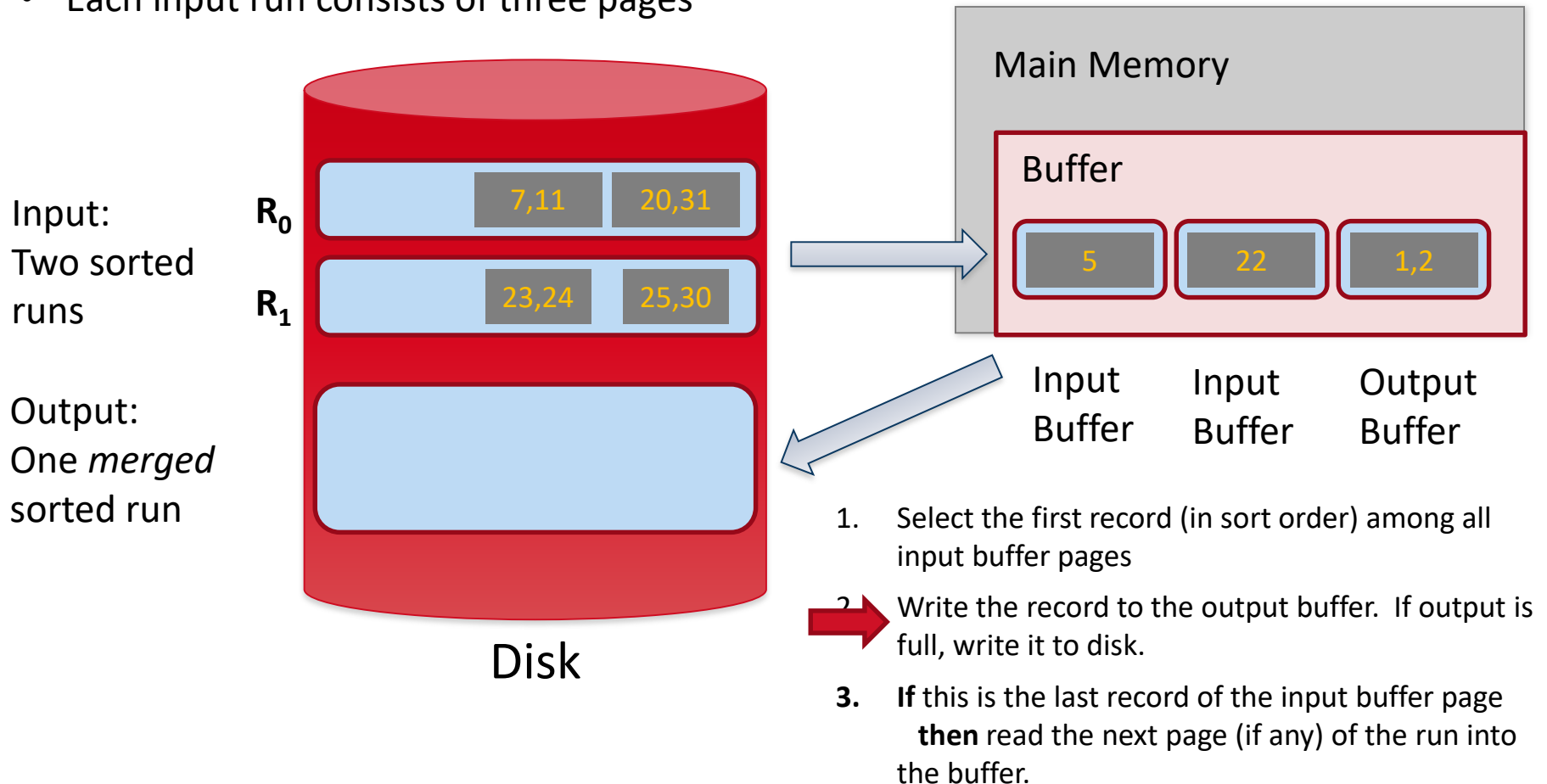
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

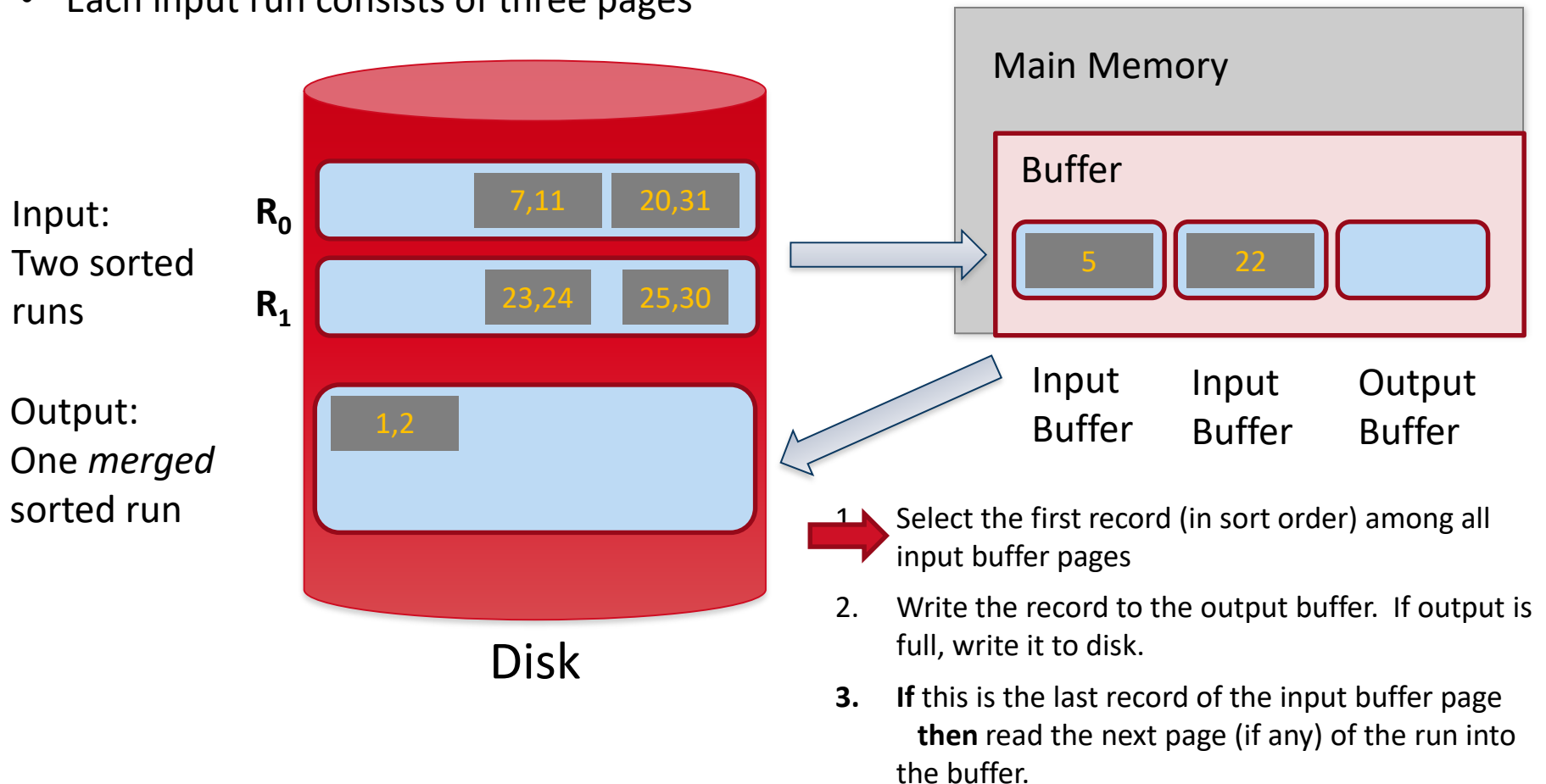
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

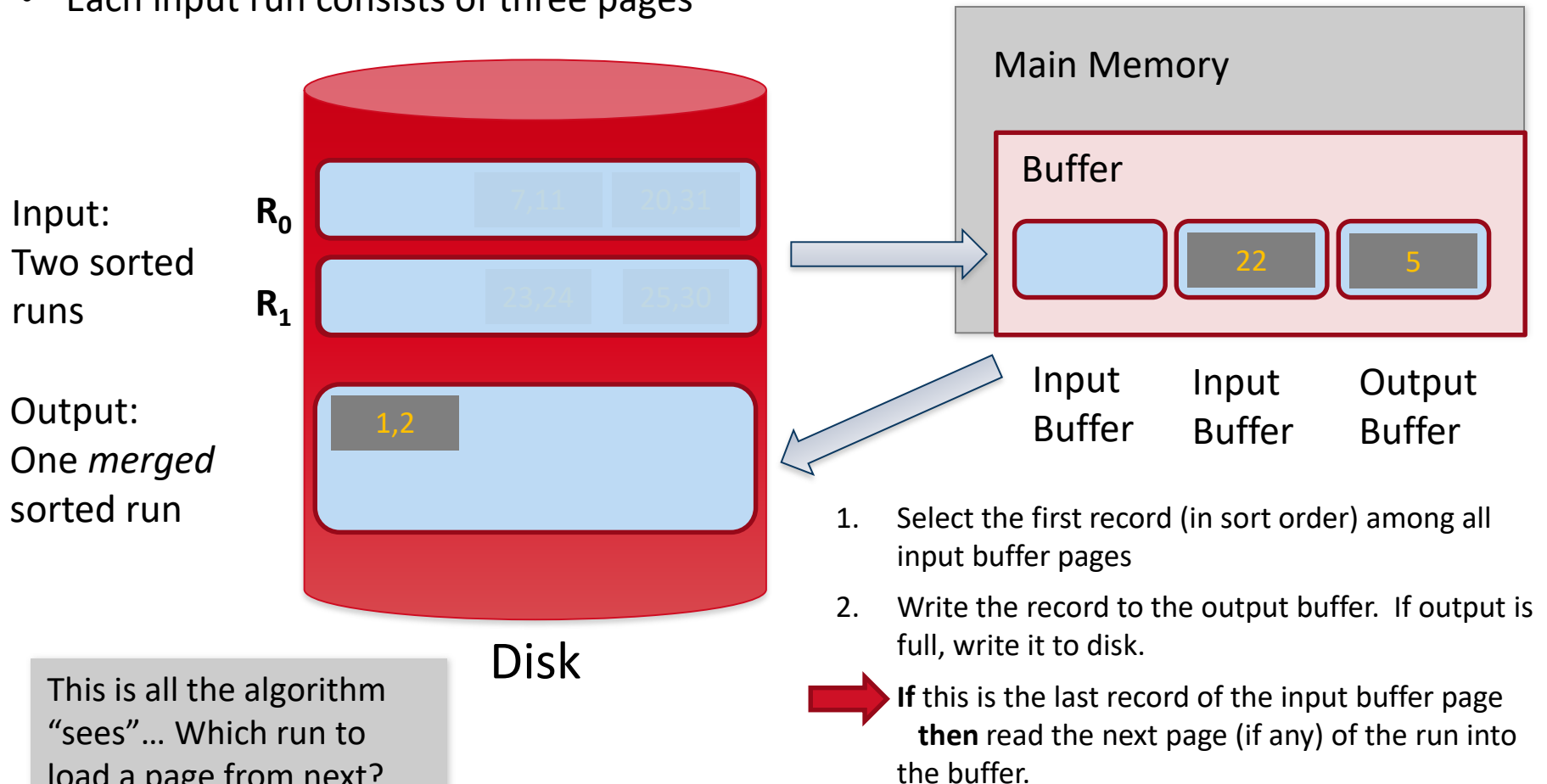
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

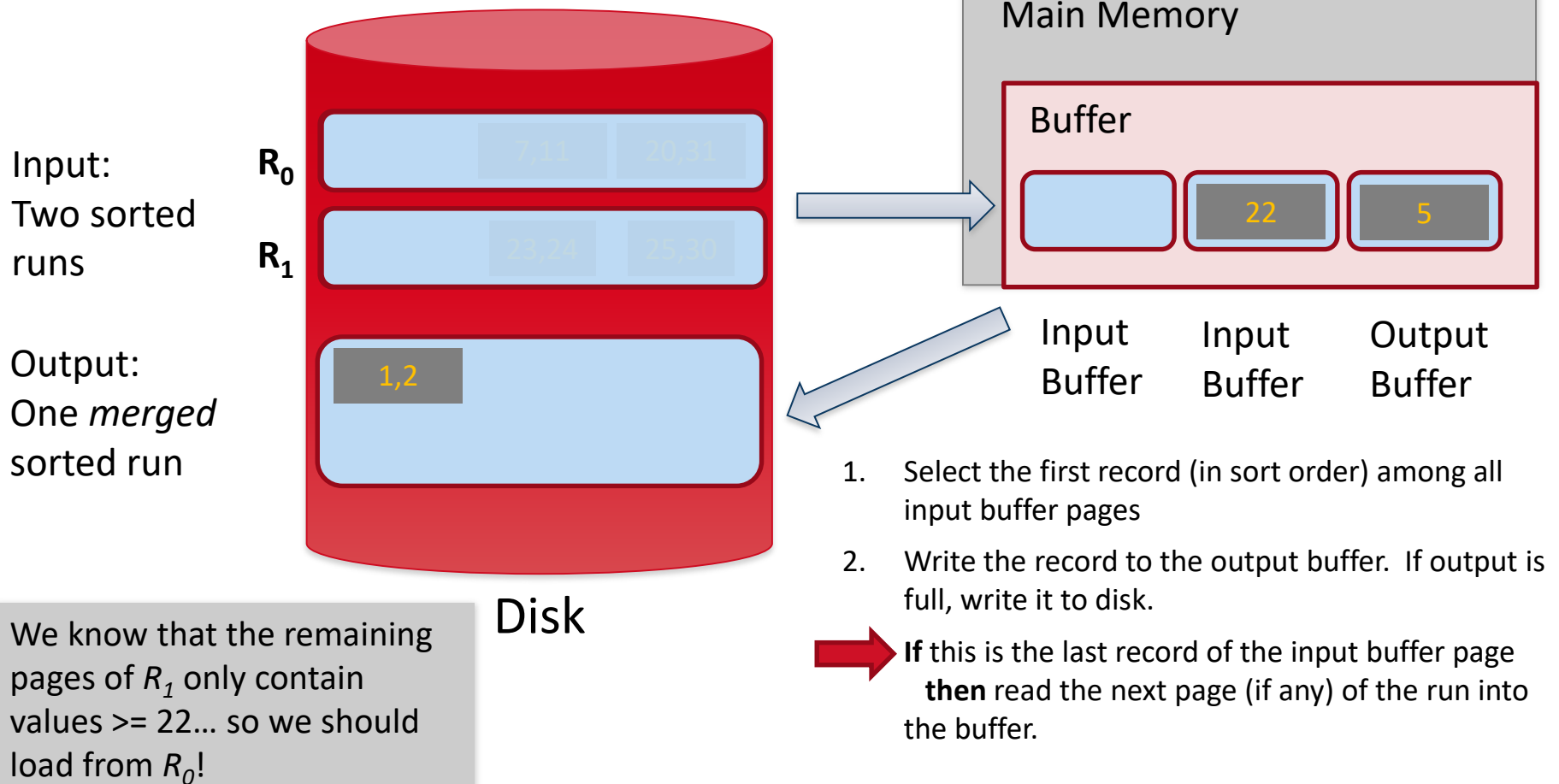
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

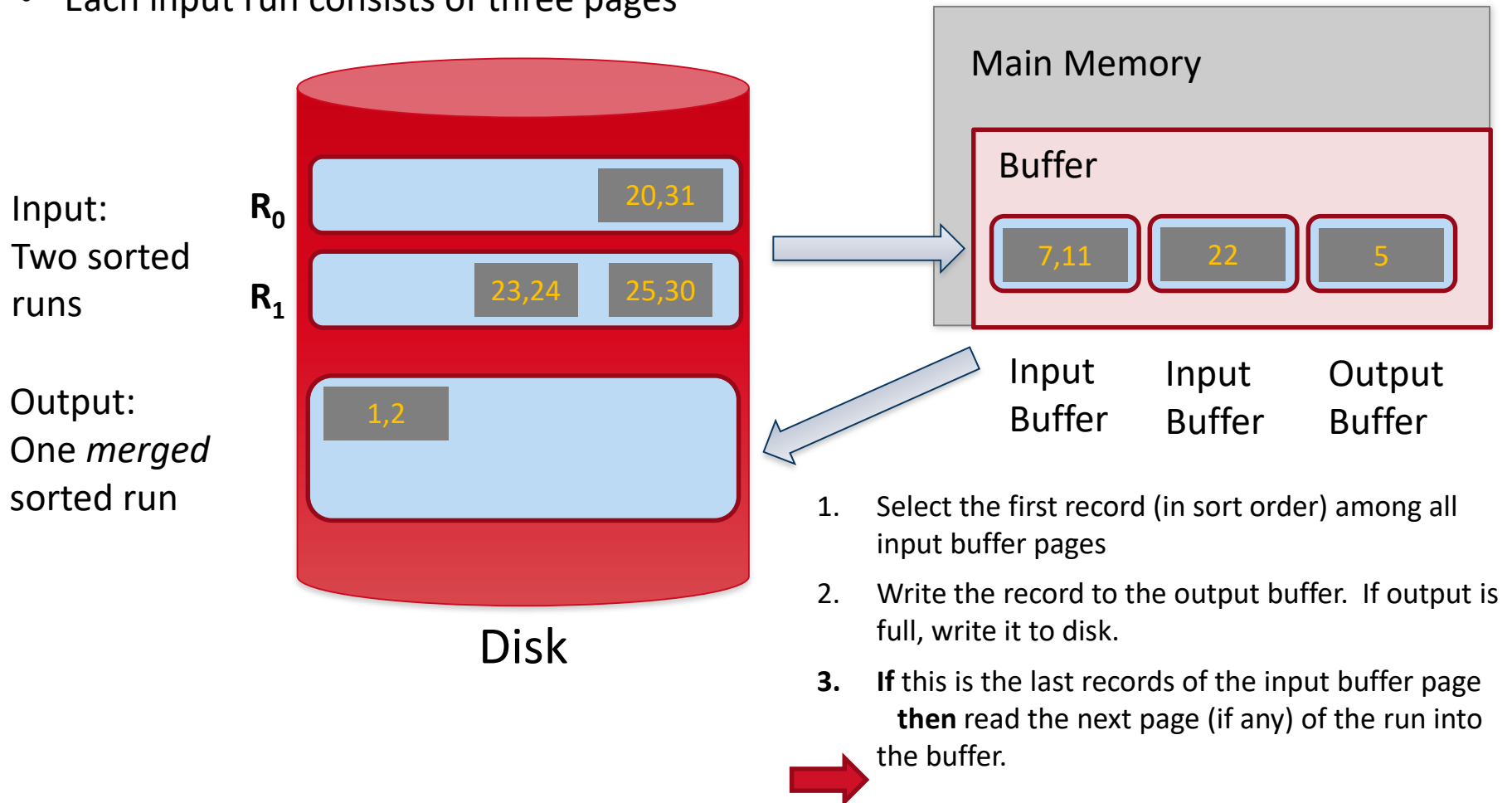
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

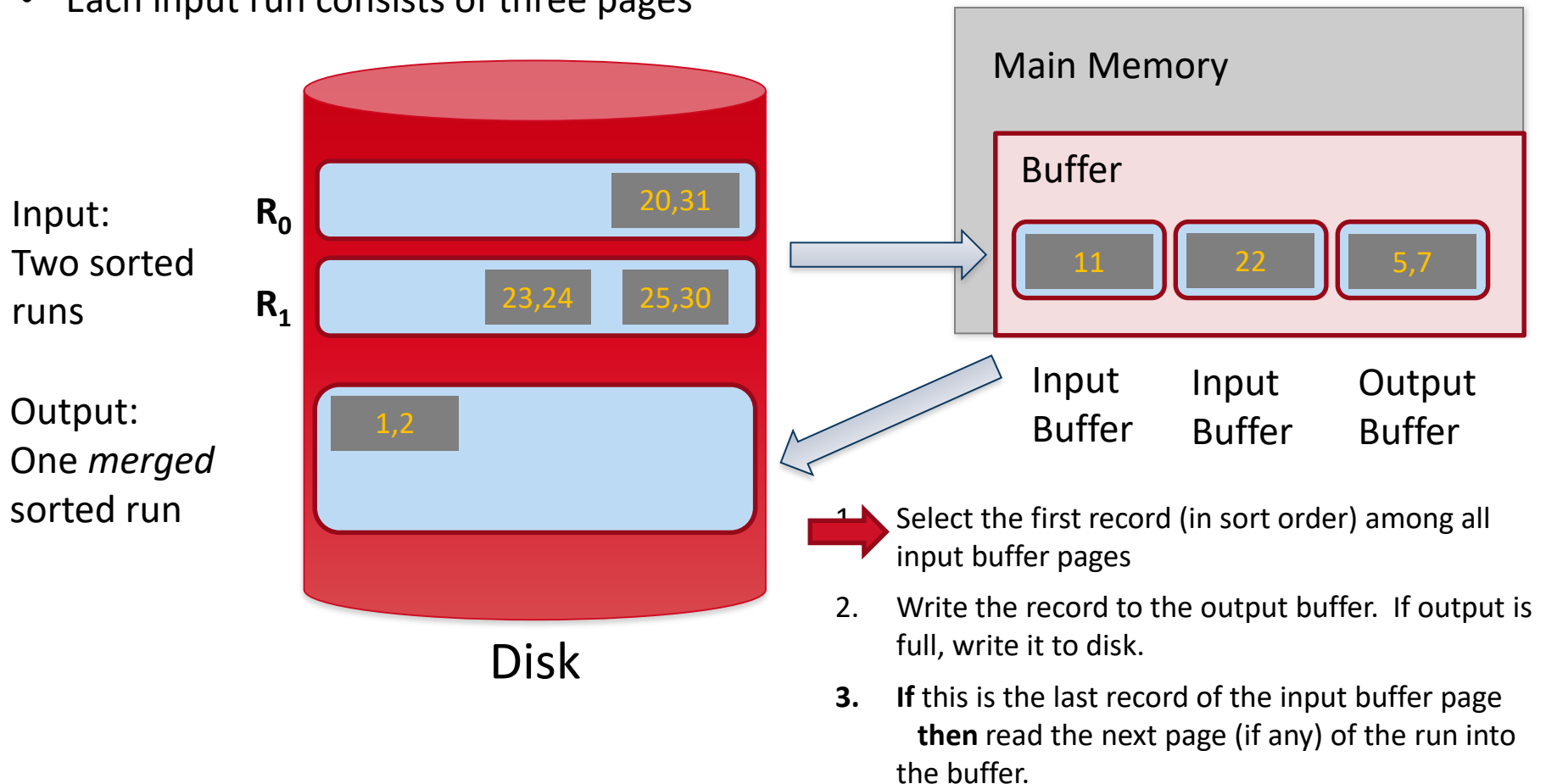
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

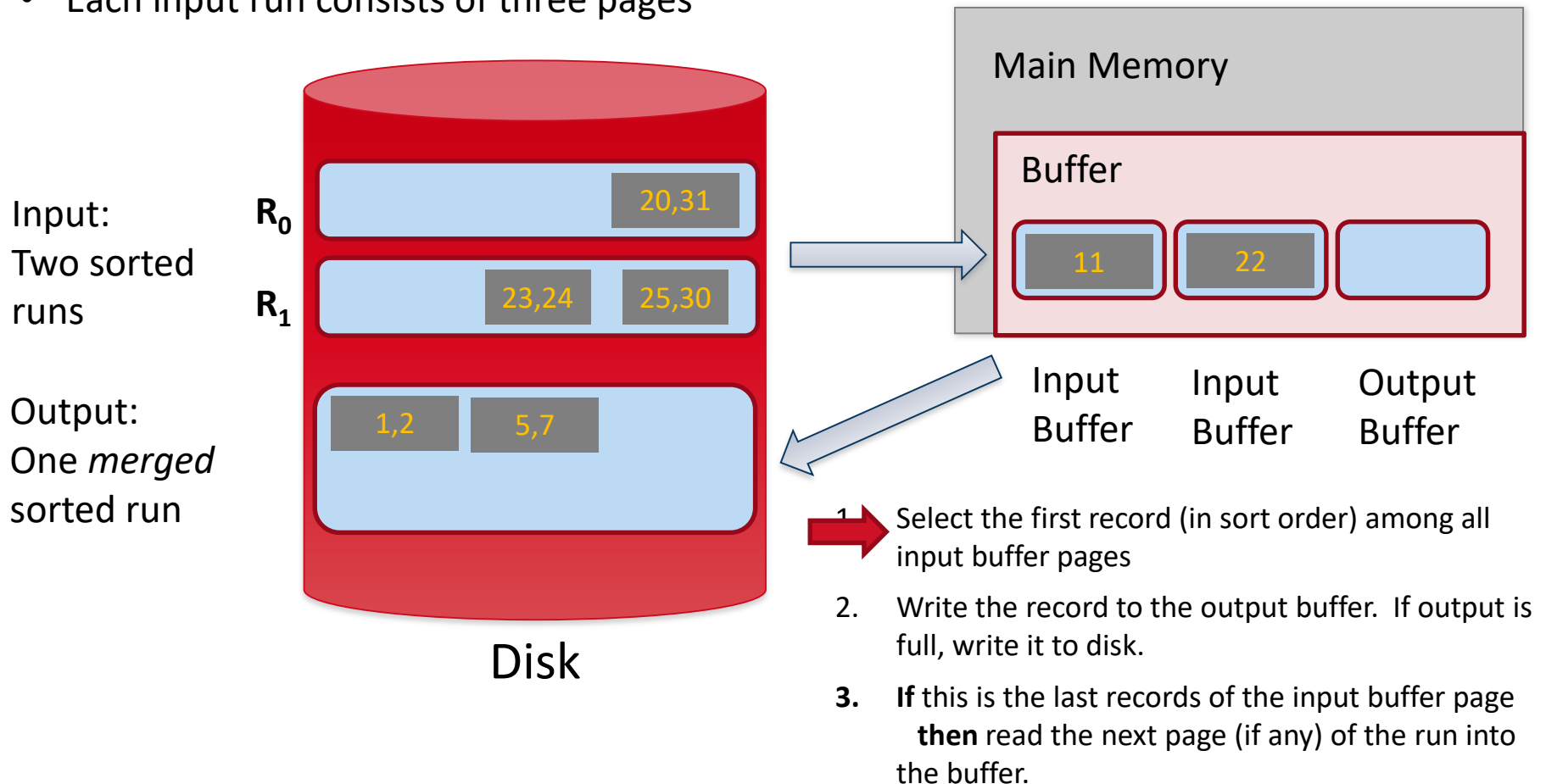
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

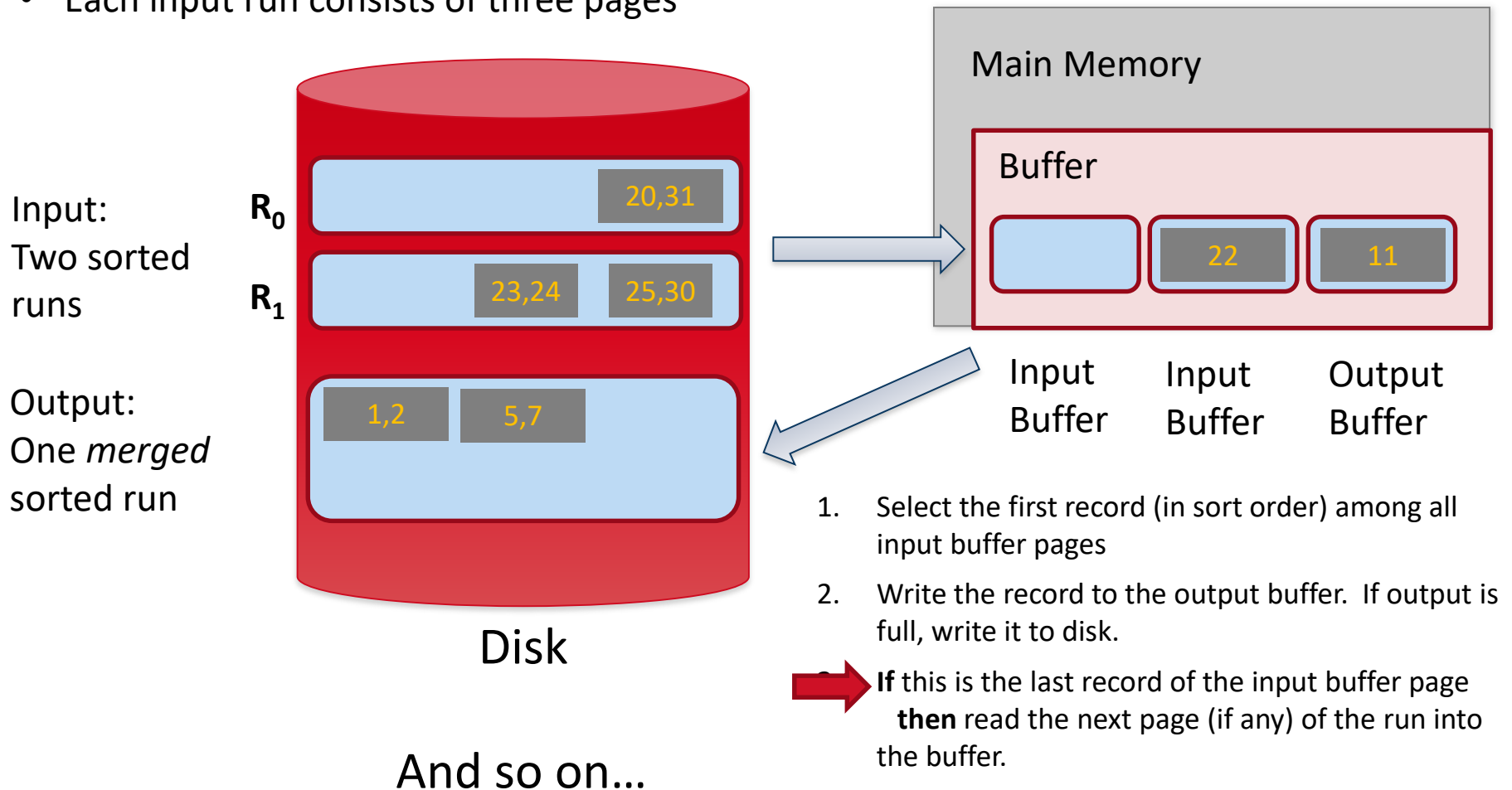
- Each input run consists of three pages



Example Step-2: Merge Sorted Runs

Example:

- Each input run consists of three pages





- › Basic Steps in Query Processing
- › Heuristic-based optimization
- › Data physical organization: Sorting
- › **Cost estimate optimization**

Cost of query processing:

- Cost of every algebraic operation in terms of I/Os

Use

- Access methods available
- Data physical organization: collected facts (e.g., blocking factors)
- Using statistics (e.g., selection cardinality)

Output of the cost estimate optimization: Efficient *physical query plan*

All joins combine tuples from two tables into a new table

- › **Theta join:** Tuples from R and S combine if some condition θ is true
 - e.g. $\theta = R.age > S.minimum\ age$
- › **Equi-join:** Tuples from R and S combine if specified attributes match in value
 - e.g: $R.x = S.y \text{ AND } R.m=S.n$
 - (special case of θ)
- › **Natural join:** Tuples from R and S combine if all attributes of same name and compatible types match in value
 - e.g $R.x=S.x \text{ AND } R.m=S.m$
 - (special case of equi-join)



Example: Joins

Assessment

sid	uosCode	sem	year	mark
316424328	INFO2120	S1	2012	72
305678453	INFO2120	S1	2012	86
316424328	INFO3005	S1	2010	63
305678453	COMP5138	S1	2012	94

UoSlecturer

uosCode	sem	year	lecturer
INFO2120	S1	2012	Uwe Roehm
INFO3005	S1	2010	Irena Koprinska
COMP5138	S2	2012	Bryn Jeffries

sid	mark	uosCode	sem	year	lecturer
316424328	72	INFO2120	S1	2012	Uwe Roehm
305678453	86	INFO2120	S1	2012	Uwe Roehm
316424328	63	INFO3005	S1	2010	Irena Koprinska

Example SQL:

```
SELECT * FROM Assessment NATURAL JOIN UoSlecturer;
```

RA:

Assessment \bowtie UosLecturer

› Join is very common!

- In SQL:

SELECT * **FROM** Students R, Enrolled S **WHERE** R.sid=S.sid

- In algebra: $R \bowtie S$

› Join must be carefully optimized.

- Semantically, $R \bowtie S$ is the same as $R \times S$ (cartesian product) followed by a selection
- However, the result of $R \times S$ usually is *significantly* larger than $R \bowtie S$; so, $R \times S$ followed by a selection is inefficient.

- › Several different algorithms to implement joins
 - *Nested loops* join
 - *Block-nested* loops join
 - *Indexed-nested* loops join

- › Choice based on cost estimate (i.e., choose the strategy with the smallest cost)
 - *Cost metric*: # of I/Os. By convention ignore output costs since they may be pipelined (i.e., not need to write back to disk).

Example Table Sizes for Cost Estimates

› The following examples refer to:

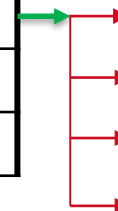
- $|R|$ tuples in R , stored in b_R pages
- $|S|$ tuples in S , stored in b_S pages
- In our examples, R refers to the relation *Students* and S refers to the relation *Enrolled*.

› Specific values:

- Number of tuples of
 - **students ($|R|$):** 1,000
 - **enrolled ($|S|$):** 10,000
- Number of pages of
 - **students (b_R):** 100
 - **enrolled (b_S):** 400

Student \bowtie Enrolled

Student			
<u>sid</u>	name	gender	country
1001	Ian	M	AUS
1002	Ha Tschì	F	ROK
1003	Grant	M	AUS



Enrolled		
<u>sid</u>	<u>uos_code</u>	semester
1001	COMP5138	2020-S2
1002	COMP5702	2020-S2
1003	COMP5138	2020-S2
1006	COMP5318	2020-S2

- › To compute the theta join $R \bowtie_{\theta} S$

for each tuple r in B_R do

for each tuple s in B_S do

if $\theta(r,s)=\text{true}$ then add $\langle r,s \rangle$ to the result

Requires
I/O

R is called the *outer table*,
 S the *inner table* of the join

- › For each tuple in the *outer* table R , we scan the entire *inner* table S .
- › **Pro:** Requires no indexes and can be used with any kind of join condition.
- › **Con:** Expensive since it examines every pair of tuples in the two tables.
- › The number of I/Os of table R is b_R
 - each page of R is read only once
- › The number of I/Os of table S is $|R| * b_S$
 - each page of S is read once for every tuple of R

- › The estimated cost of nested loops join is

$$b_R + |R| * b_S$$

- › Example:

- students (R) as outer table: $100 + 1000 * 400 = \mathbf{400,100}$ disk I/Os
- enrolled (S) as outer table: $400 + 10000 * 100 = \mathbf{1,000,400}$ disk I/Os $(b_S + |S| * b_R)$

Number of tuples of **students** ($|R|$): 1,000

enrolled ($|S|$): 10,000

Number of pages of **students** (b_R): 100

enrolled (b_S): 400

- Variant of nested loops join in which every page of inner table is paired with every page (or multiple pages) of outer table.
- For each *page* of R, get each *page* of S, and write out matching pairs of tuples $\langle r, s \rangle$, where r is in R-page and S is in S-page.

```
for each page  $B_R$  of R do
  for each page  $B_S$  of S do
    for each tuple  $r$  in  $B_R$  do
      for each tuple  $s$  in  $B_S$  do
        if  $\theta(r,s)=\text{true}$  then output  $\langle r,s \rangle$ 
```

Requires
I/O

Nested
Loops
Join

```
for each page  $B_R$  of R do
  for each tuple  $r$  in  $B_R$  do
    for each page  $B_S$  of S do
      for each tuple  $s$  in  $B_S$  do
        if  $\theta(r,s)=\text{true}$  then add  $\langle r,s \rangle$  to the result
```

Cost Analysis: Block-Nested Loops Join

- › The number of I/Os of table R is b_R (each page of R is read only once)
- › The number of I/Os of table S is $b_R * b_S$ (each page of S is read once for every page of R)

- › Cost of block-nested loops join is

$$b_R + b_R * b_S$$

- › Example:

- students (R) as outer table: $100 + 100 * 400 = \mathbf{40,100 \text{ disk I/Os}}$
- enrolled (S) as outer table: $400 + 400 * 100 = \mathbf{40,400 \text{ disk I/Os}}$

Given an index *idx* built on the join attribute of S

for each page B_R of R do

for each tuple r in B_R do

for each tuple s in *idx*(r) do

add $\langle r, s \rangle$ to result

Requires
I/O

- › To use index-nested loops join, the following conditions must satisfy:
 - join is an **equi-join or natural join**, and
 - an **index is available** on the inner table's join attribute
- › For each tuple r in the outer table R , use the index to look up tuples in S that satisfy the join condition with tuple r .

Cost Analysis: Index-Nested Loops Join

- › For each tuple in R , we perform an index lookup on S .
 - Cost: $b_R + (|R| * c)$
 - Where c is the cost of traversing index and fetching all matching S tuples for one tuple of R
 - c can be estimated as the cost of a single selection on S using the join condition.

- › If indexes are available on join attributes of both R and S , use the table with fewer tuples as the outer table.

- › Understanding of Role and Structure of Query Processing
 - From SQL to physical data access
 - 3 Steps: Query Parsing, Optimization, Execution
 - Expression Tree vs. Evaluation Plan
 - Query Execution Algorithms
- › Operator Algorithms
 - External Merge Sort
 - Joins
 - Simple Nested
 - Block Nested
 - Index Nested



- › Discussion regarding Final Exam
 - Instructions
 - Question types
- › Content Review

See you next week!



THE UNIVERSITY OF
SYDNEY