

University of Michigan

17

# Sorting & Aggregations



Database Management Systems  
EECS 484  
Fall 2024

**LM**

Lin Ma  
Computer Science and  
Engineering Division

# LOGISTICS

---

## Alternative final exam

- <https://forms.gle/FRub1ex753QhpeHd8>
- December 9<sup>th</sup>, 2024 @ 1pm – 3pm
- For exceptional situations

## Purposes of course resources

- Piazza: Discussions of high-level questions (NOT for remote debugging)
- Office hour: Individual questions (including helping with homework and project debugging)
- Staff email: Logistics questions (including account/autograder access)



# COURSE STATUS

---

We are now going to talk about how to execute queries using the DBMS components we have discussed so far.

Next few lectures:

- Operator Algorithms
- Query Processing Models
- Runtime Architectures

Log Manager

Transaction Manager

Query Planning

Operator Execution

Access Methods

Disk Manager

# COURSE STATUS

---

We are now going to talk about how to execute queries using the DBMS components we have discussed so far.

Next few lectures:

- Operator Algorithms
- Query Processing Models
- Runtime Architectures

Log Manager

Transaction Manager

Query Planning

Operator Execution

Access Methods

Disk Manager

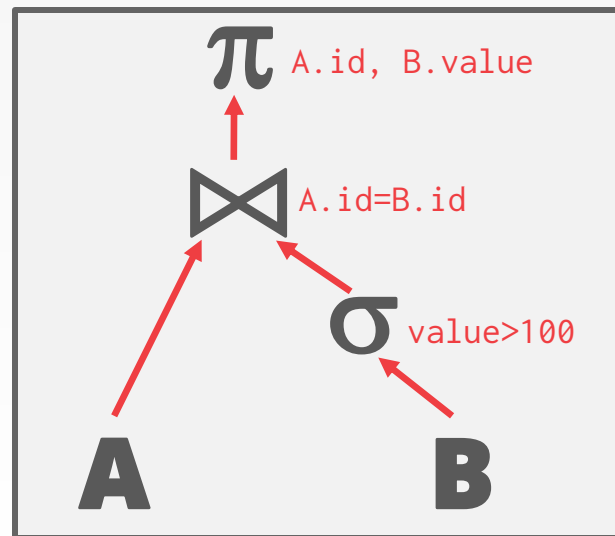
# QUERY PLAN

The operators are arranged in a tree.

Data flows from the leaves of the tree up towards the root.

The output of the root node is the result of the query.

```
SELECT A.id, B.value  
FROM A, B  
WHERE A.id = B.id  
AND B.value > 100
```



## DISK-ORIENTED DBMS

---

Just like it cannot assume that a table fits entirely in memory, a disk-oriented DBMS cannot assume that query results fit in memory.

We are going to rely on the buffer pool to implement algorithms that need to spill to disk.

We are also going to prefer algorithms that maximize the amount of sequential I/O.



# TODAY'S AGENDA

---

External Merge Sort  
Aggregations



# WHY DO WE NEED TO SORT?

---

Relational model/SQL is unsorted.





# WHY DO WE NEED TO SORT?

---

Relational model/SQL is unsorted.

Queries may request that tuples are sorted in a specific way (**ORDER BY**).



# WHY DO WE NEED TO SORT?

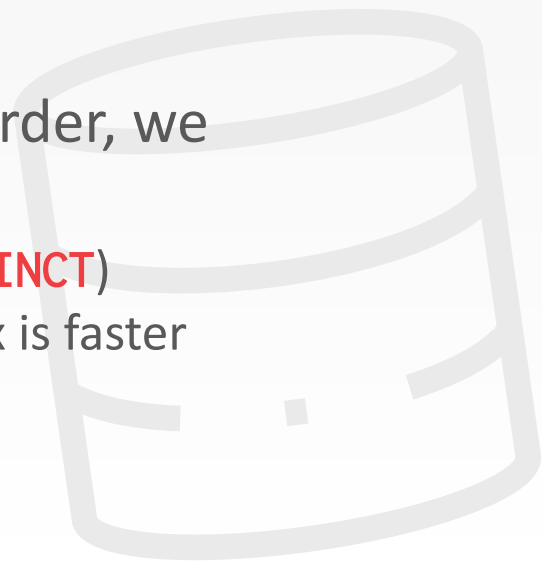
---

Relational model/SQL is unsorted.

Queries may request that tuples are sorted in a specific way (**ORDER BY**).

But even if a query does not specify an order, we may still want to sort to do other things:

- Trivial to support duplicate elimination (**DISTINCT**)
- Bulk loading sorted tuples into a B+Tree index is faster
- Aggregations (**GROUP BY**)
- ...



# SORTING ALGORITHMS

---

If data fits in memory, then we can use a standard sorting algorithm like quicksort.



# SORTING ALGORITHMS

---

If data fits in memory, then we can use a standard sorting algorithm like quicksort.

If data does not fit in memory, then we need to use a technique that is aware of the cost of reading and writing disk pages...



# EXTERNAL MERGE SORT

---

Divide-and-conquer algorithm that splits data into separate runs, sorts them individually, and then combines them into longer sorted runs.

## Phase #1 – Sorting

→ Sort chunks of data that fit in memory and then write back the sorted chunks to a file on disk.

## Phase #2 – Merging

→ Combine sorted runs into larger chunks.



## SORTED RUN

---

A run is an ordered list of key/value pairs.

**Key:** The attribute(s) to compare to compute the sort order.

**Value:** Two choices

- Tuple (*early materialization*).
- Record ID (*late materialization*).



# SORTED RUN

---

A run is an ordered list of key/value pairs.

**Key:** The attribute(s) to compare to compute the sort order.

**Value:** Two choices

- Tuple (*early materialization*).
- Record ID (*late materialization*).

## *Early Materialization*

K1	<Tuple Data>
K2	<Tuple Data>

⋮



# SORTED RUN

A run is an ordered list of key/value pairs.

**Key:** The attribute(s) to compare to compute the sort order.

**Value:** Two choices

- Tuple (*early materialization*).
- Record ID (*late materialization*).

## Early Materialization

K1	<Tuple Data>
K2	<Tuple Data>

⋮

## Late Materialization

K1	❏	K2	❏	...	Kn	❏
----	---	----	---	-----	----	---



# SORTED RUN

A run is an ordered list of key/value pairs.

**Key:** The attribute(s) to compare to compute the sort order.

**Value:** Two choices

- Tuple (*early materialization*).
- Record ID (*late materialization*).

## Early Materialization

K1	<Tuple Data>
K2	<Tuple Data>

⋮

## Late Materialization

K1	☒	K2	☒	...	Kn	☒
----	---	----	---	-----	----	---

→ *Record ID*

## 2-WAY EXTERNAL MERGE SORT

---

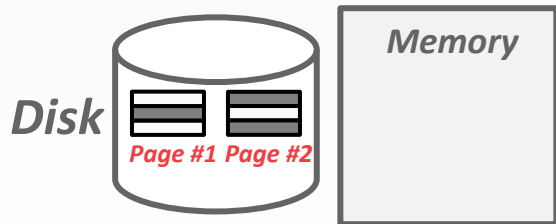
We will start with a simple example of a 2-way external merge sort.

→ “2” is the number of runs that we are going to merge into a new run for each pass.



# 2-WAY EXTERNAL MERGE SORT

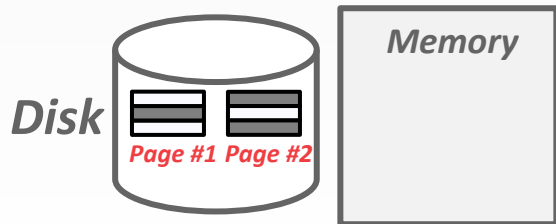
---



# 2-WAY EXTERNAL MERGE SORT

## Pass #0

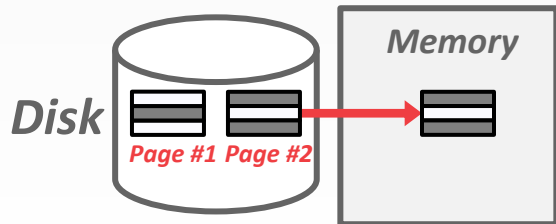
- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk



## 2-WAY EXTERNAL MERGE SORT

### Pass #0

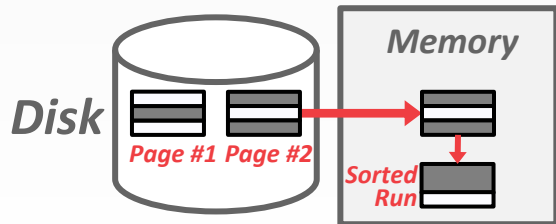
- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk



## 2-WAY EXTERNAL MERGE SORT

### Pass #0

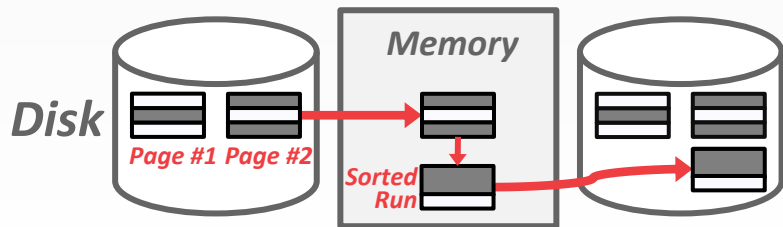
- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk



## 2-WAY EXTERNAL MERGE SORT

### Pass #0

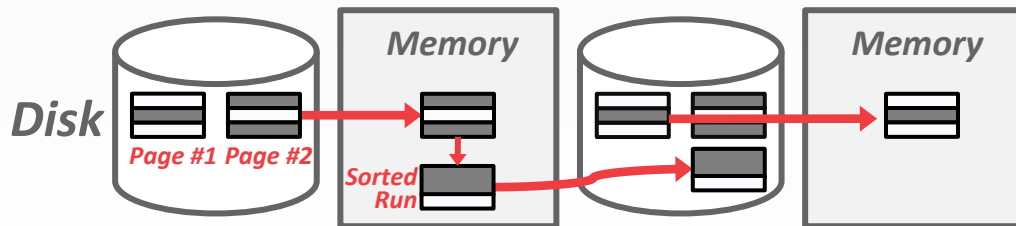
- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk



## 2-WAY EXTERNAL MERGE SORT

### Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk

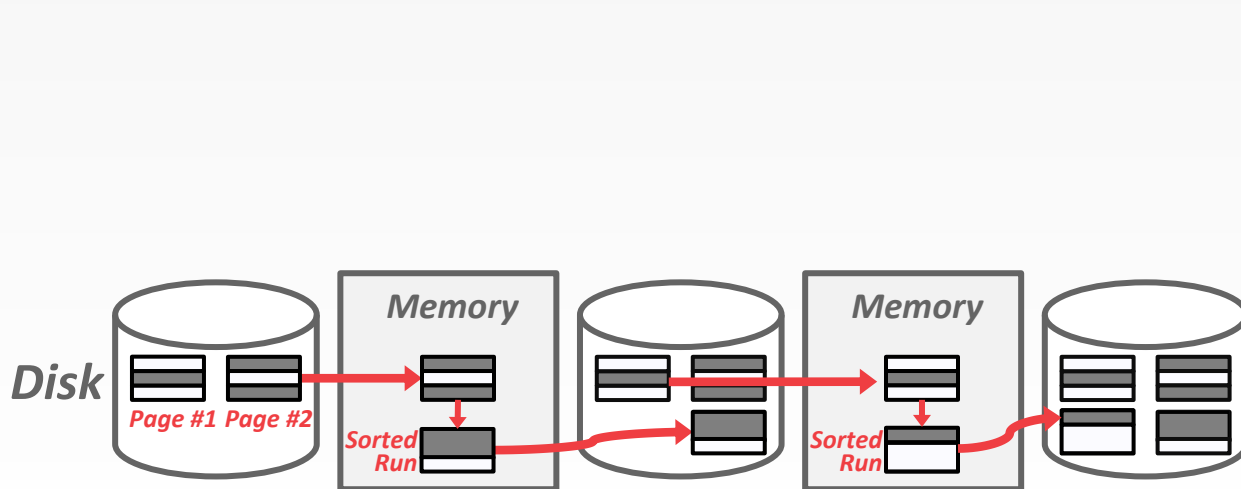




## 2-WAY EXTERNAL MERGE SORT

### Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk



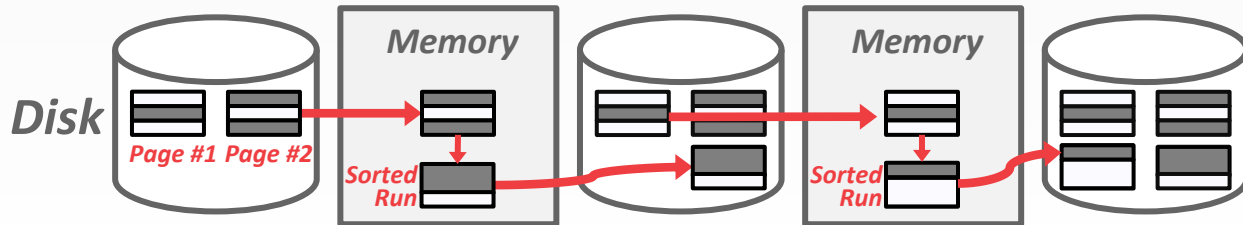
# 2-WAY EXTERNAL MERGE SORT

## Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk

## Pass #1,2,3,...

- Recursively merge pairs of runs into runs twice as long



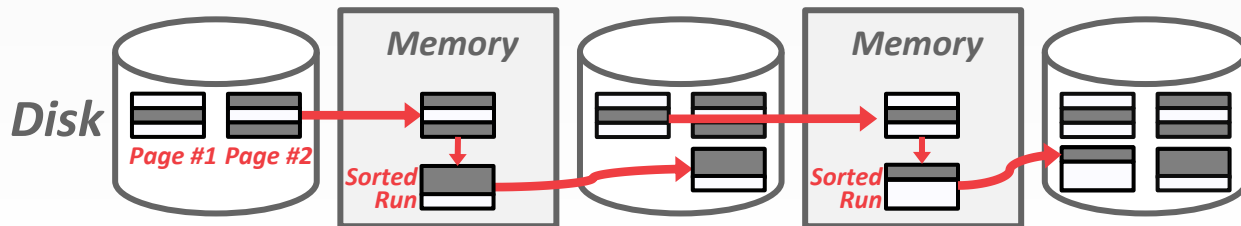
# 2-WAY EXTERNAL MERGE SORT

## Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk

## Pass #1,2,3,...

- Recursively merge pairs of runs into runs twice as long
- Uses three buffer pages (2 for input pages, 1 for output)



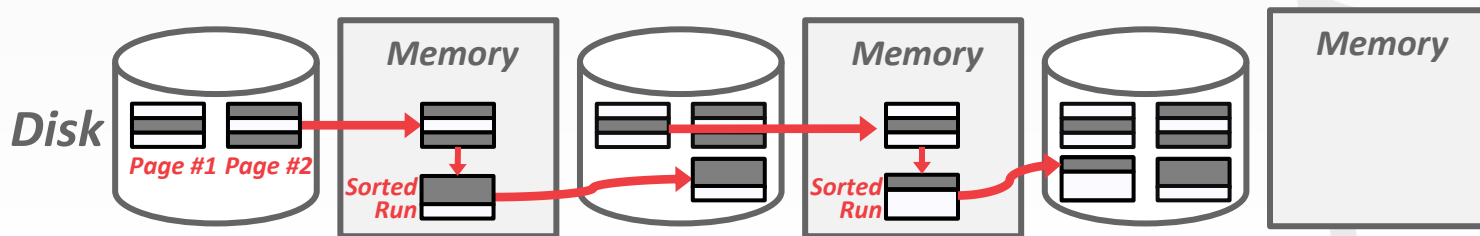
## 2-WAY EXTERNAL MERGE SORT

### Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk

### Pass #1,2,3,...

- Recursively merge pairs of runs into runs twice as long
- Uses three buffer pages (2 for input pages, 1 for output)



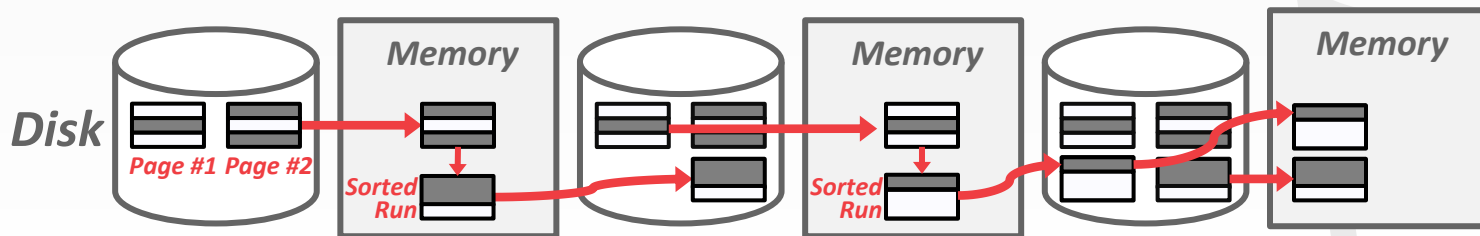
## 2-WAY EXTERNAL MERGE SORT

### Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk

### Pass #1,2,3,...

- Recursively merge pairs of runs into runs twice as long
- Uses three buffer pages (2 for input pages, 1 for output)



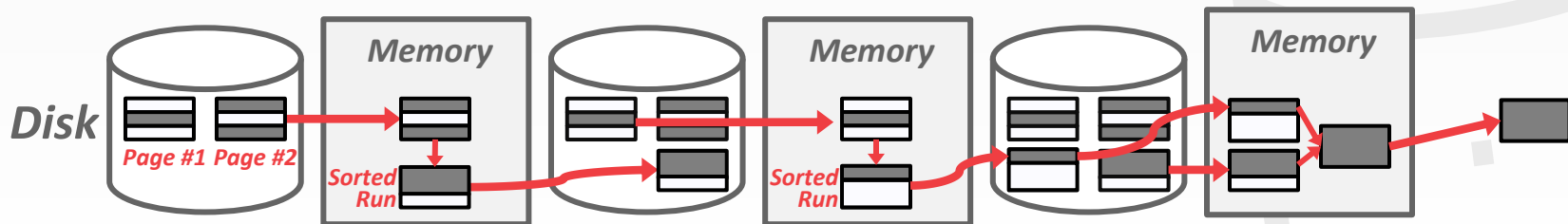
## 2-WAY EXTERNAL MERGE SORT

### Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk

### Pass #1,2,3,...

- Recursively merge pairs of runs into runs twice as long
- Uses three buffer pages (2 for input pages, 1 for output)



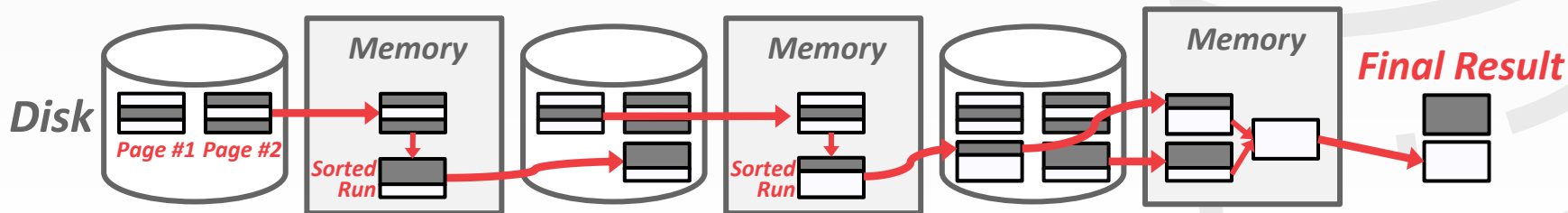
## 2-WAY EXTERNAL MERGE SORT

### Pass #0

- Read pages of the table into memory, one at a time
- Sort the page into a run and write it back to disk

### Pass #1,2,3,...

- Recursively merge pairs of runs into runs twice as long
- Uses three buffer pages (2 for input pages, 1 for output)



## 2-WAY EXTERNAL MERGE SORT

---

In each pass, we read and write every page in the file.

For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\text{\# of passes})$$





## 2-WAY EXTERNAL MERGE SORT

3,4	6,2	9,4	8,7	5,6	3,1	2	∅
-----	-----	-----	-----	-----	-----	---	---

In each pass, we read and write every page in the file.

For  $N$  pages of data:

Number of passes

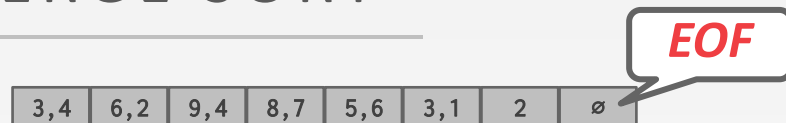
$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\# \text{ of passes})$$



## 2-WAY EXTERNAL MERGE SORT



In each pass, we read and write every page in the file.

For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\text{\# of passes})$$



## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

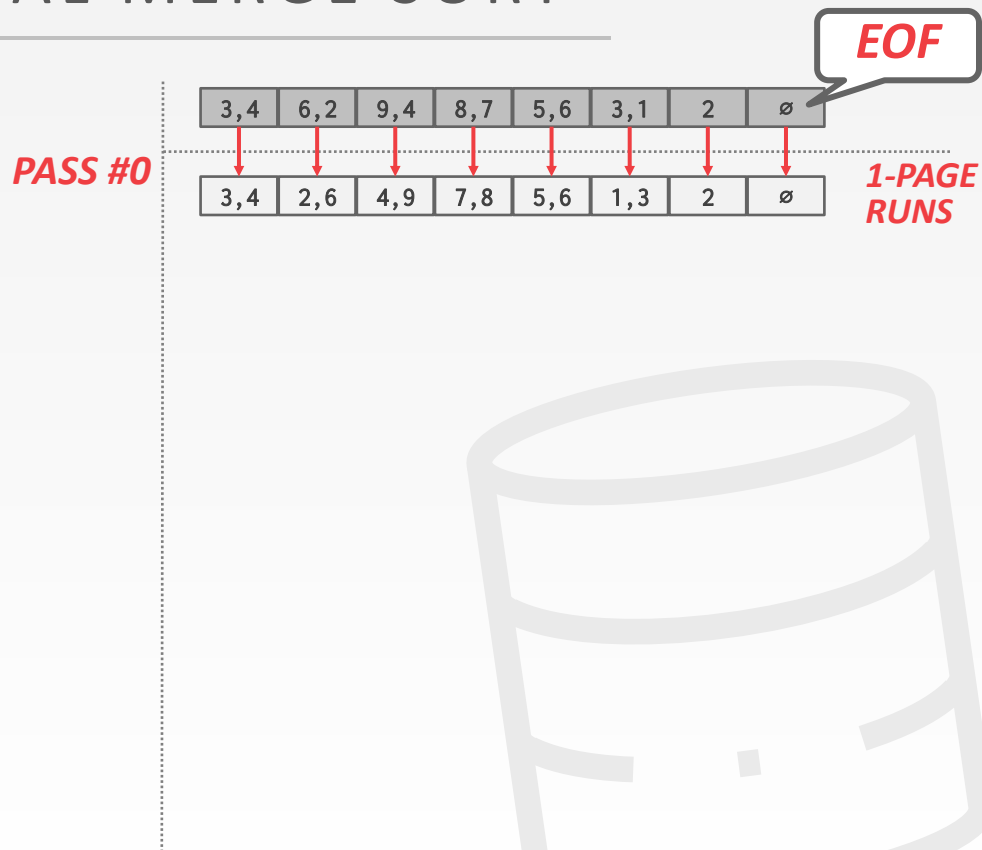
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\text{\# of passes})$$



## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

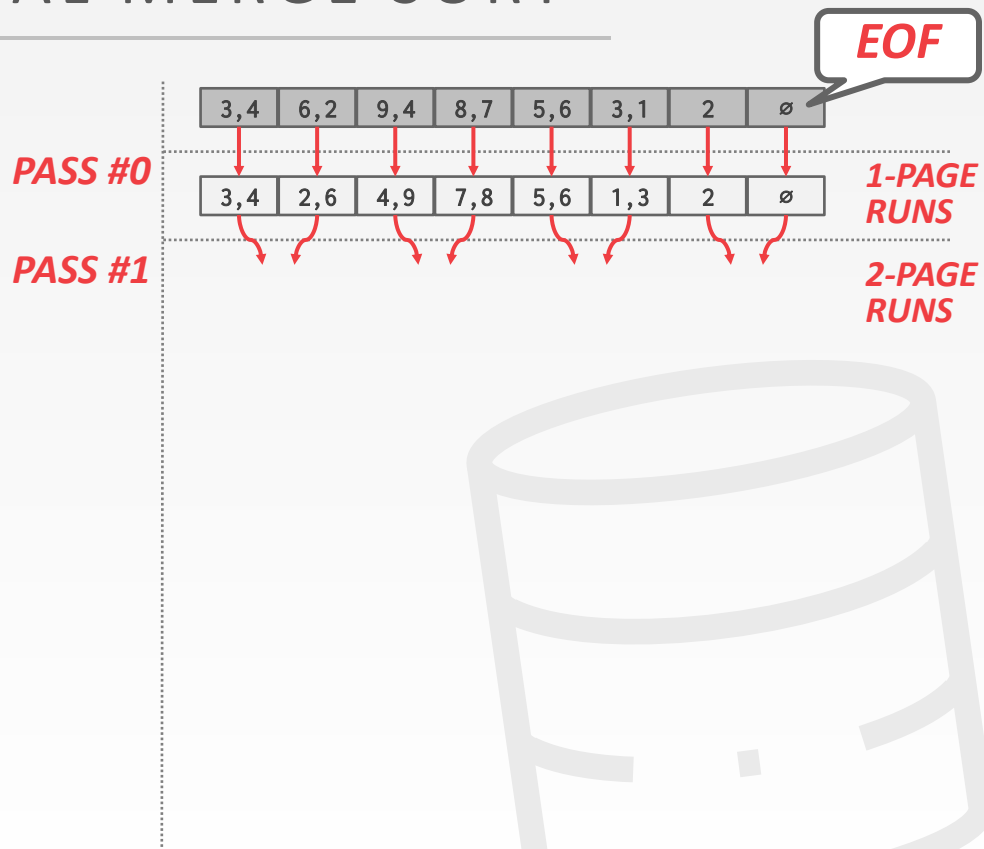
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\# \text{ of passes})$$



## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

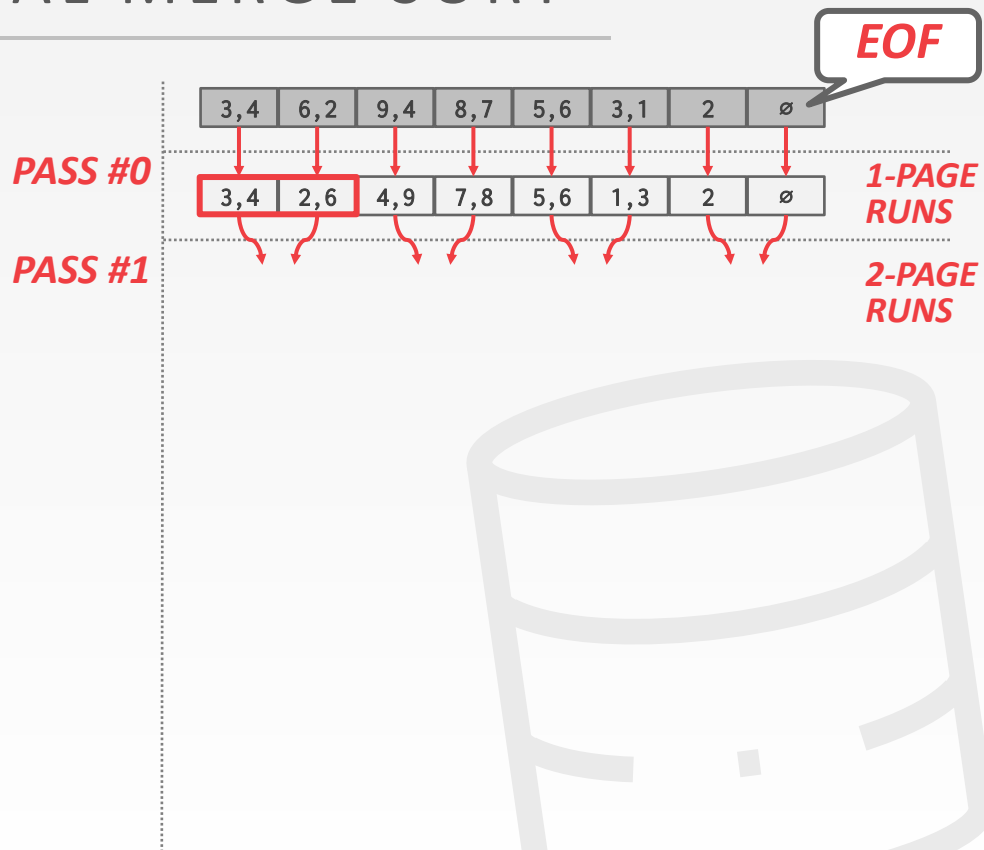
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\# \text{ of passes})$$



## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

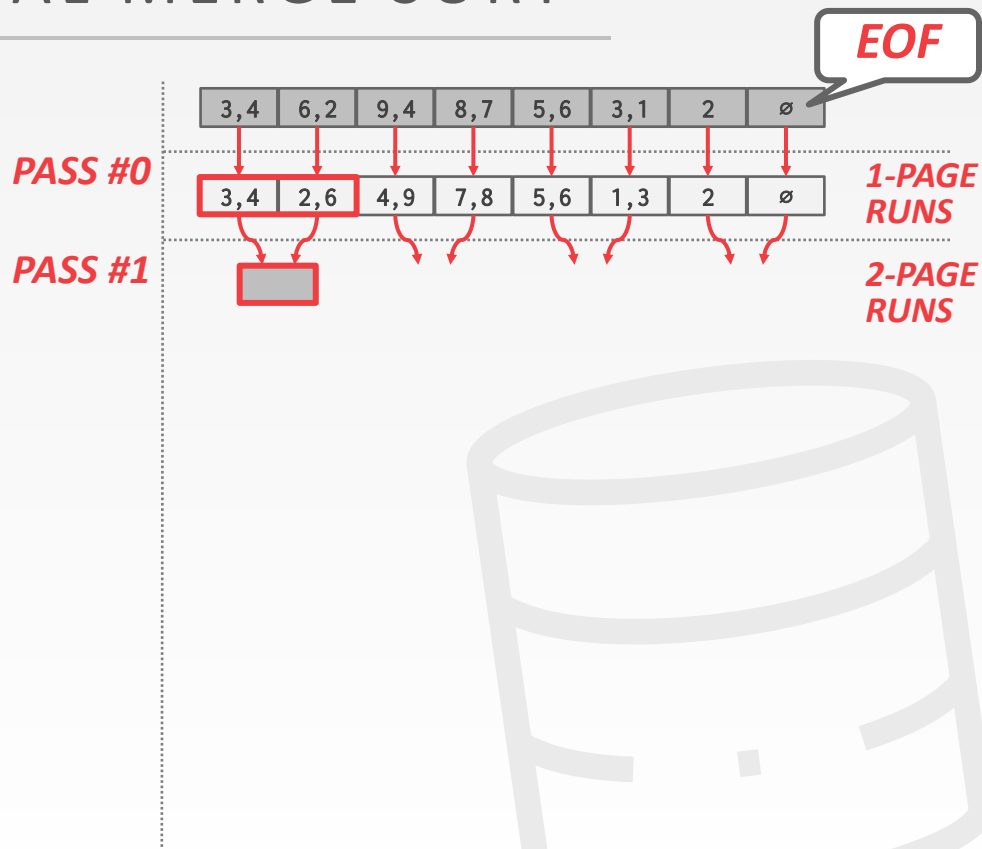
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\# \text{ of passes})$$



## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

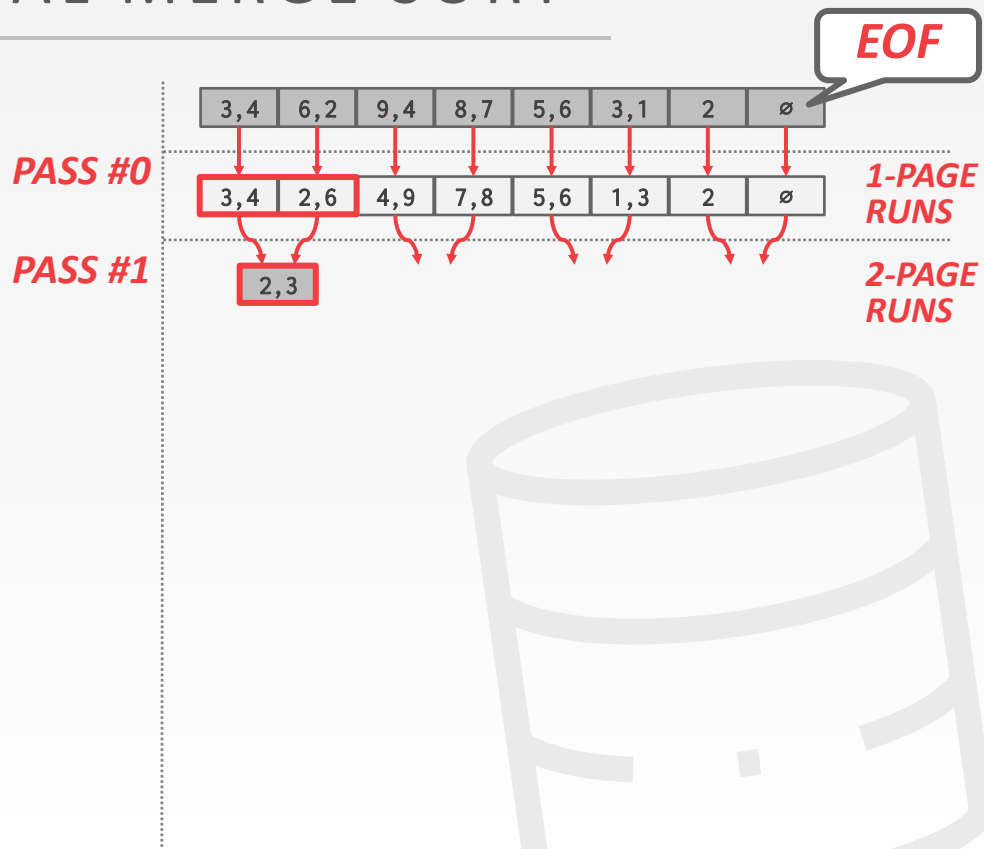
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\# \text{ of passes})$$



## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

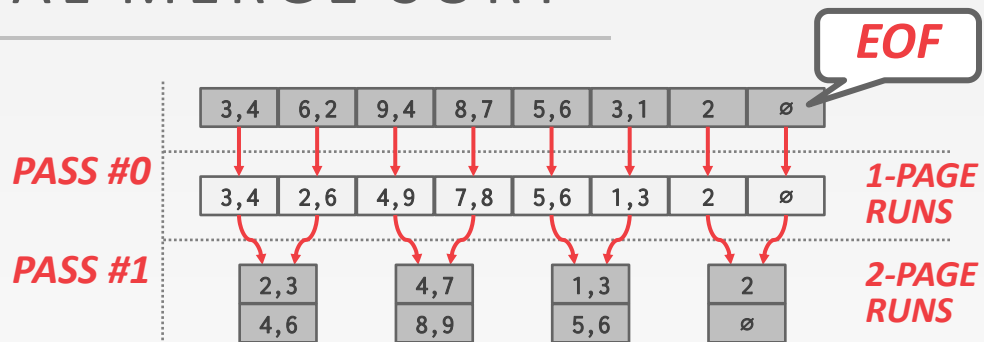
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\# \text{ of passes})$$





## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

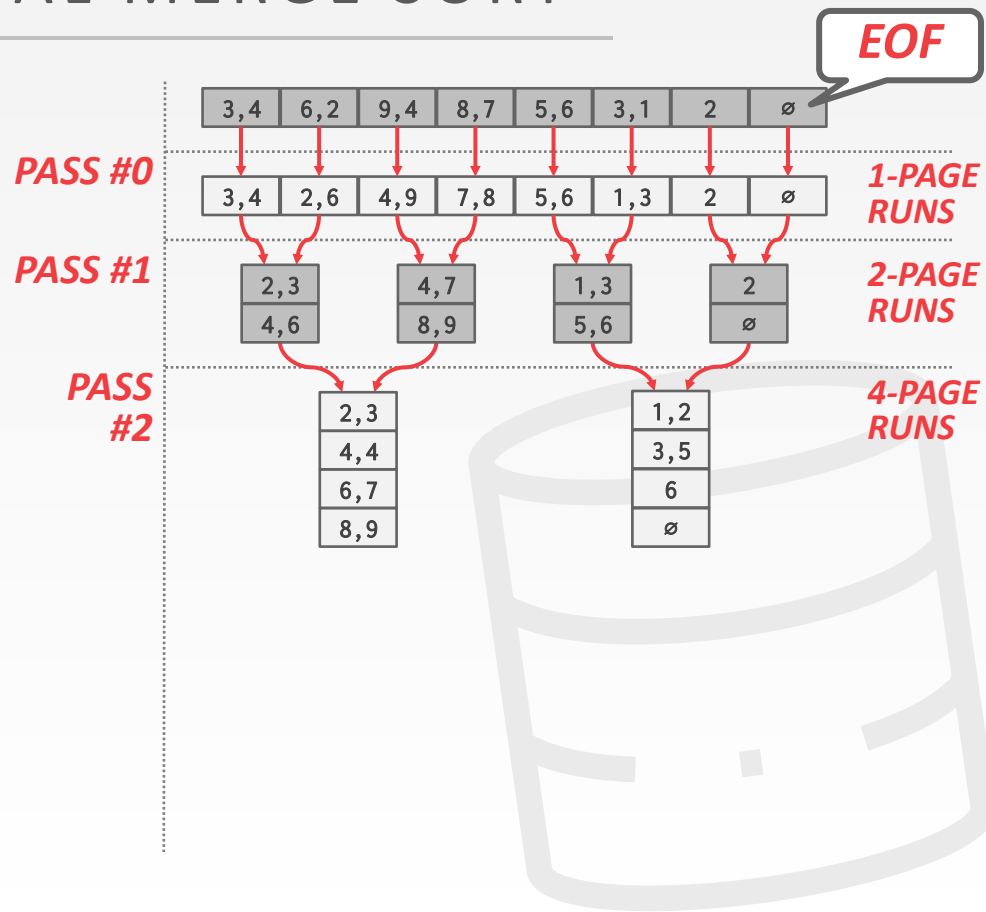
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\# \text{ of passes})$$



## 2-WAY EXTERNAL MERGE SORT

In each pass, we read and write every page in the file.

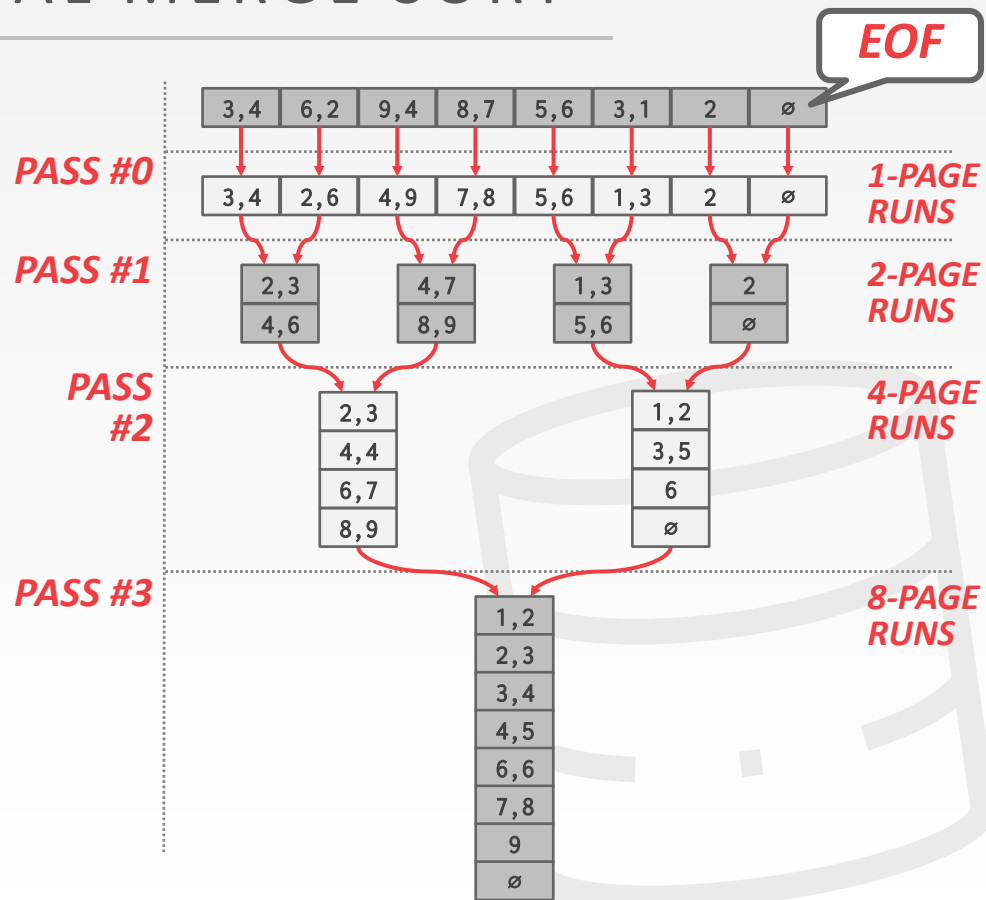
For  $N$  pages of data:

Number of passes

$$= 1 + \lceil \log_2 N \rceil$$

Total I/O cost

$$= 2N \cdot (\text{\# of passes})$$



## 2-WAY EXTERNAL MERGE SORT

---

This algorithm only requires three buffer pool pages to perform the sorting

→ Two input pages, one output page

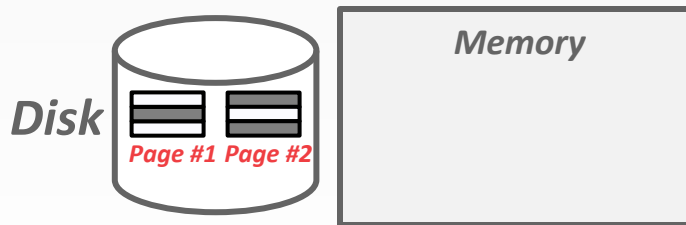
But even if we have more buffer space available, it does not effectively utilize them if the worker must block on disk I/O...



# DOUBLE BUFFERING OPTIMIZATION

Prefetch the next run in the background and store it in a second buffer while the system is processing the current run.

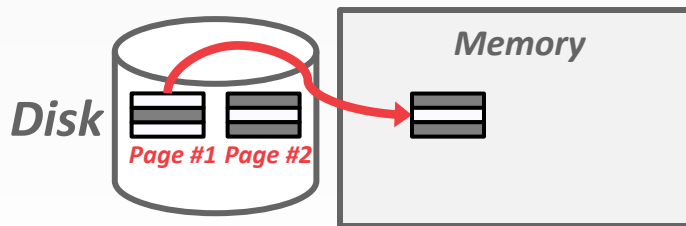
→ Reduces the wait time for I/O requests at each step by continuously utilizing the disk.



# DOUBLE BUFFERING OPTIMIZATION

Prefetch the next run in the background and store it in a second buffer while the system is processing the current run.

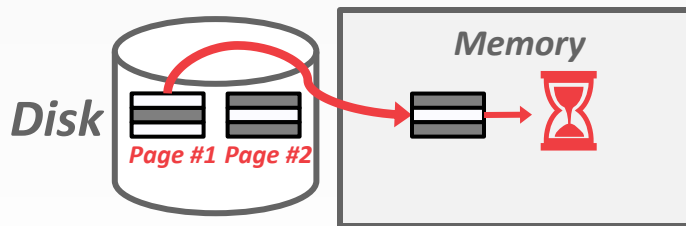
→ Reduces the wait time for I/O requests at each step by continuously utilizing the disk.



# DOUBLE BUFFERING OPTIMIZATION

Prefetch the next run in the background and store it in a second buffer while the system is processing the current run.

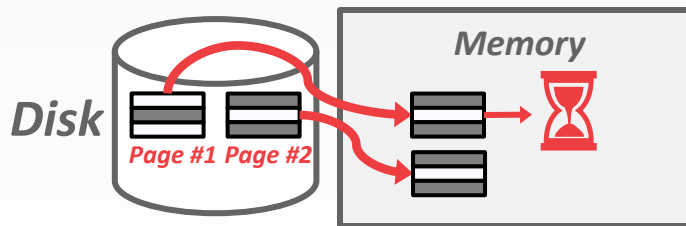
→ Reduces the wait time for I/O requests at each step by continuously utilizing the disk.



# DOUBLE BUFFERING OPTIMIZATION

Prefetch the next run in the background and store it in a second buffer while the system is processing the current run.

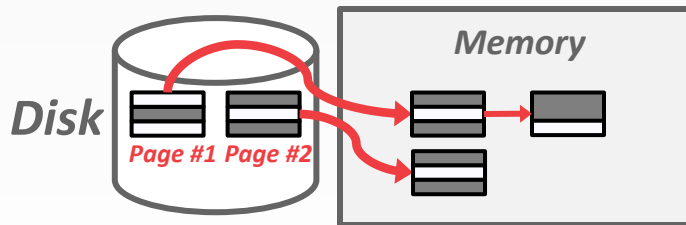
→ Reduces the wait time for I/O requests at each step by continuously utilizing the disk.



# DOUBLE BUFFERING OPTIMIZATION

Prefetch the next run in the background and store it in a second buffer while the system is processing the current run.

→ Reduces the wait time for I/O requests at each step by continuously utilizing the disk.

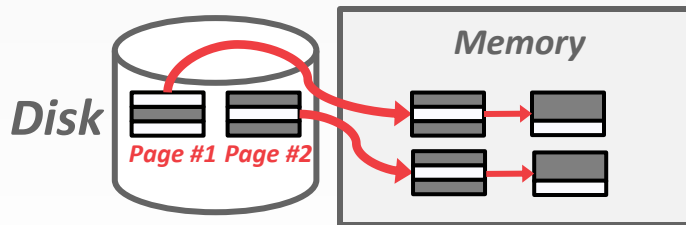




# DOUBLE BUFFERING OPTIMIZATION

Prefetch the next run in the background and store it in a second buffer while the system is processing the current run.

→ Reduces the wait time for I/O requests at each step by continuously utilizing the disk.



# GENERAL EXTERNAL MERGE SORT

## Pass #0

- Assume the DBMS has a finite number of  $B$  buffer pool pages to hold input and output data.
- Read  $B$  pages of the table into memory at a time. Sort them and write back to disk.
- Produce  $\lceil N / B \rceil$  sorted runs of size  $B$

## Pass #1,2,3,...

- Merge  $B-1$  runs (i.e., K-way merge)

Number of passes =  $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$

Total I/O Cost =  $2N \cdot (\# \text{ of passes})$



# GENERAL EXTERNAL MERGE SORT

## Pass #0

- Assume the DBMS has a finite number of  $B$  buffer pool pages to hold input and output data.
- Read  $B$  pages of the table into memory at a time. Sort them and write back to disk.
- Produce  $\lceil N / B \rceil$  sorted runs of size  $B$

## Pass #1,2,3,...

- Merge  $B-1$  runs (i.e., K-way merge)

Number of passes =  $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$

Total I/O Cost =  $2N \cdot (\# \text{ of passes})$



# GENERAL EXTERNAL MERGE SORT

## Pass #0

- Assume the DBMS has a finite number of  $B$  buffer pool pages to hold input and output data.
- Read  $B$  pages of the table into memory at a time. Sort them and write back to disk.
- Produce  $\lceil N / B \rceil$  sorted runs of size  $B$

## Pass #1,2,3,...

- Merge  $B-1$  runs (i.e., K-way merge)

Number of passes =  $1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil$

Total I/O Cost =  $2N \cdot (\# \text{ of passes})$



## EXAMPLE

---

Determine how many passes it takes to sort 108 pages with 5 buffer pool pages:  **$N=108$ ,  $B=5$**



## EXAMPLE

---

Determine how many passes it takes to sort 108 pages with 5 buffer pool pages:  **$N=108$ ,  $B=5$**

→ **Pass #0:**  **$\lceil N / B \rceil$**  =  $\lceil 108 / 5 \rceil$  = 22 sorted runs of 5 pages each (last run is only 3 pages).



## EXAMPLE

---

Determine how many passes it takes to sort 108 pages with 5 buffer pool pages:  **$N=108$ ,  $B=5$**

- **Pass #0:**  **$\lceil N / B \rceil$**  =  $\lceil 108 / 5 \rceil$  = 22 sorted runs of 5 pages each (last run is only 3 pages).
- **Pass #1:**  **$\lceil N' / B-1 \rceil$**  =  $\lceil 22 / 4 \rceil$  = 6 sorted runs of 20 pages each (last run is only 8 pages).



## EXAMPLE

---

Determine how many passes it takes to sort 108 pages with 5 buffer pool pages:  **$N=108$ ,  $B=5$**

- **Pass #0:**  **$\lceil N / B \rceil$**  =  $\lceil 108 / 5 \rceil$  = 22 sorted runs of 5 pages each (last run is only 3 pages).
- **Pass #1:**  **$\lceil N' / B-1 \rceil$**  =  $\lceil 22 / 4 \rceil$  = 6 sorted runs of 20 pages each (last run is only 8 pages).
- **Pass #2:**  **$\lceil N'' / B-1 \rceil$**  =  $\lceil 6 / 4 \rceil$  = 2 sorted runs, first one has 80 pages and second one has 28 pages.





## EXAMPLE

---

Determine how many passes it takes to sort 108 pages with 5 buffer pool pages:  **$N=108$ ,  $B=5$**

- **Pass #0:**  **$\lceil N / B \rceil$**  =  $\lceil 108 / 5 \rceil$  = 22 sorted runs of 5 pages each (last run is only 3 pages).
- **Pass #1:**  **$\lceil N' / B-1 \rceil$**  =  $\lceil 22 / 4 \rceil$  = 6 sorted runs of 20 pages each (last run is only 8 pages).
- **Pass #2:**  **$\lceil N'' / B-1 \rceil$**  =  $\lceil 6 / 4 \rceil$  = 2 sorted runs, first one has 80 pages and second one has 28 pages.
- **Pass #3:** Sorted file of 108 pages.



## EXAMPLE

Determine how many passes it takes to sort 108 pages with 5 buffer pool pages:  $N=108$ ,  $B=5$

- **Pass #0:**  $\lceil N / B \rceil = \lceil 108 / 5 \rceil = 22$  sorted runs of 5 pages each (last run is only 3 pages).
- **Pass #1:**  $\lceil N' / B-1 \rceil = \lceil 22 / 4 \rceil = 6$  sorted runs of 20 pages each (last run is only 8 pages).
- **Pass #2:**  $\lceil N'' / B-1 \rceil = \lceil 6 / 4 \rceil = 2$  sorted runs, first one has 80 pages and second one has 28 pages.
- **Pass #3:** Sorted file of 108 pages.

$$\begin{aligned} 1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil &= 1 + \lceil \log_4 22 \rceil = 1 + \lceil 2.229... \rceil \\ &= 4 \text{ passes} \end{aligned}$$



## USING B+TREES FOR SORTING

---

If the table that must be sorted already has a B+Tree index on the sort attribute(s), then we can use that to accelerate sorting.

Retrieve tuples in desired sort order by simply traversing the leaf pages of the tree.

Cases to consider:

- Clustered B+Tree
- Unclustered B+Tree



## CASE #1 – CLUSTERED B+TREE

Traverse to the left-most leaf page,  
and then retrieve tuples from all leaf  
pages.

This is always better than external  
sorting because there is no  
computational cost, and all disk  
access is sequential.

*B+Tree Index*

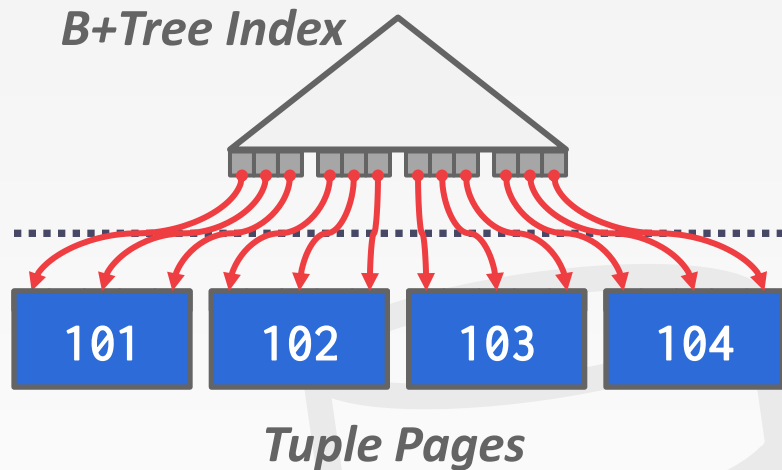


*Tuple Pages*

## CASE #1 – CLUSTERED B+TREE

Traverse to the left-most leaf page,  
and then retrieve tuples from all leaf  
pages.

This is always better than external  
sorting because there is no  
computational cost, and all disk  
access is sequential.



## CASE #2 – UNCLUSTERED B+TREE

Chase each pointer to the page that contains the data.

This is almost always a bad idea.  
In general, one I/O per data record.

*B+Tree Index*

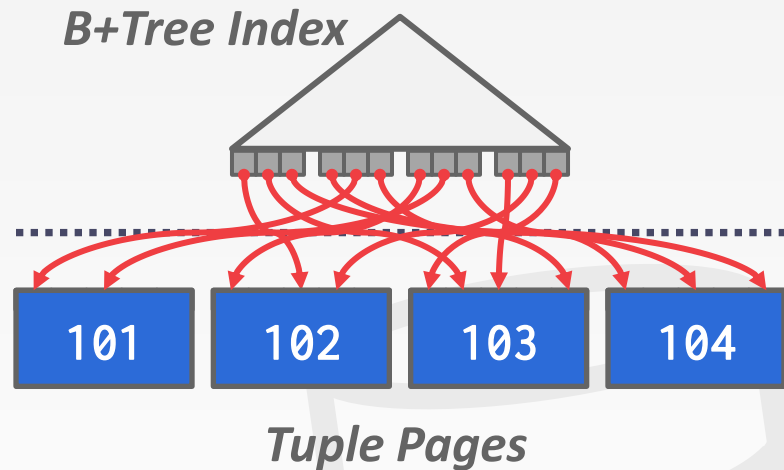


*Tuple Pages*

## CASE #2 – UNCLUSTERED B+TREE

Chase each pointer to the page that contains the data.

This is almost always a bad idea.  
In general, one I/O per data record.



# AGGREGATIONS

---

Collapse values for a single attribute from multiple tuples into a single scalar value.

Two implementation choices:

- Sorting
- Hashing





## SORTING AGGREGATION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
 ORDER BY cid
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

# SORTING AGGREGATION

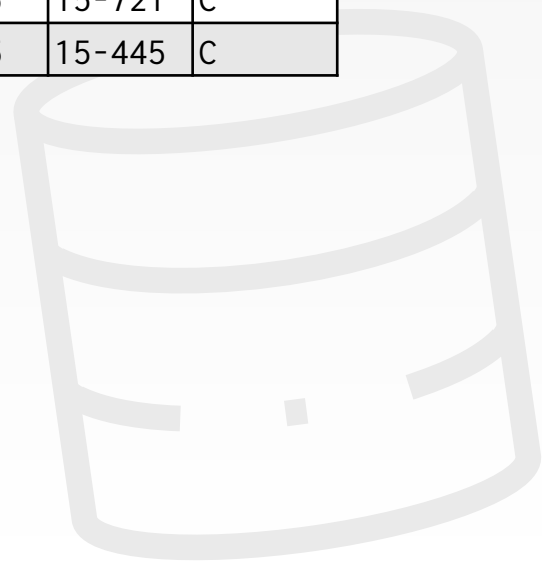
```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
 ORDER BY cid
```



sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



# SORTING AGGREGATION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
 ORDER BY cid
```

enrolled(sid,cid,grade)

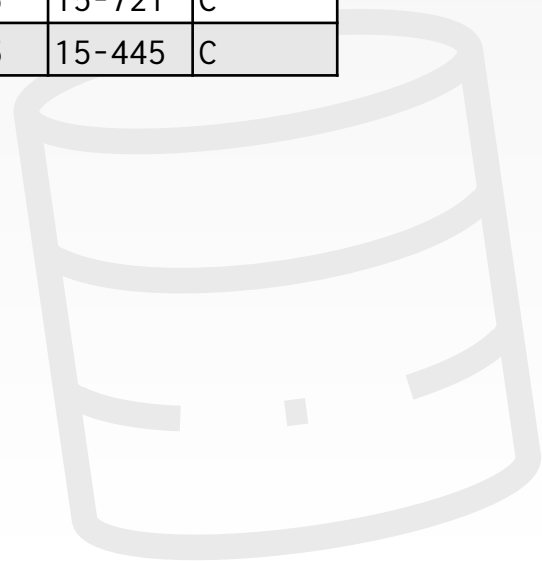
sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C



cid
15-445
15-826
15-721
15-445



# SORTING AGGREGATION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
 ORDER BY cid
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

  
**Filter**

sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C

  
**Remove  
Columns**

cid
15-445
15-826
15-721
15-445

  
**Sort**

cid
15-445
15-445
15-721
15-826

# SORTING AGGREGATION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
 ORDER BY cid
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

  
**Filter**

sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C

  
**Remove  
Columns**

cid
15-445
15-826
15-721
15-445

  
**Sort**

cid
15-445
15-445
15-721
15-826

  
**Eliminate  
Dupes**

# SORTING AGGREGATION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
 ORDER BY cid
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C



cid
15-445
15-826
15-721
15-445



cid
15-445
<del>15-445</del>
15-721
15-826



**Eliminate  
Duples**

# ALTERNATIVES TO SORTING

---

What if we do not need the data to be ordered?

- Forming groups in **GROUP BY** (no ordering)
- Removing duplicates in **DISTINCT** (no ordering)



# ALTERNATIVES TO SORTING

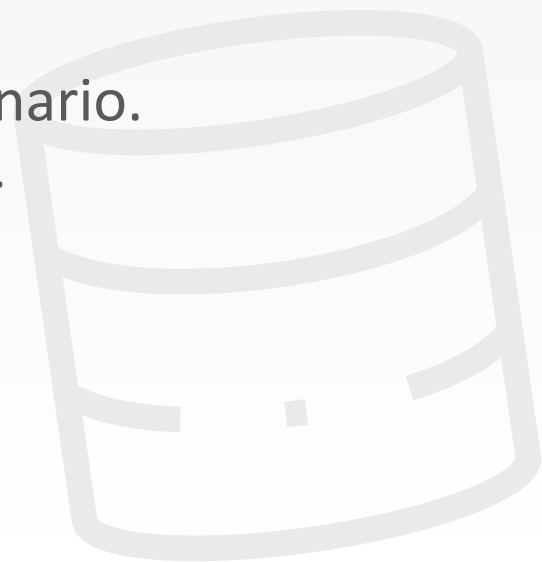
---

What if we do not need the data to be ordered?

- Forming groups in **GROUP BY** (no ordering)
- Removing duplicates in **DISTINCT** (no ordering)

Hashing is a better alternative in this scenario.

- Can be computationally cheaper than sorting.





# HASHING AGGREGATE

---

Populate an ephemeral hash table as the DBMS scans the table. For each record, check whether there is already an entry in the hash table:

- **DISTINCT**: Discard duplicate
- **GROUP BY**: Perform aggregate computation

If everything fits in memory, then this is easy.

If the DBMS must spill data to disk, then we need to be smarter...



# EXTERNAL HASHING AGGREGATE

---

## Phase #1 – Partition

- Divide tuples into buckets based on hash key
- Write them out to disk when they get full

## Phase #2 – ReHash

- Build in-memory hash table for each partition and compute the aggregation



## PHASE #1 – PARTITION

---

Use a hash function  $h_1$  to split tuples into partitions on disk.

- A partition is one or more pages that contain the set of keys with the same hash value.
- Partitions are “spilled” to disk via output buffers.

Assume that we have  $B$  buffers.

We will use  $B-1$  buffers for the partitions and  $1$  buffer for the input data.

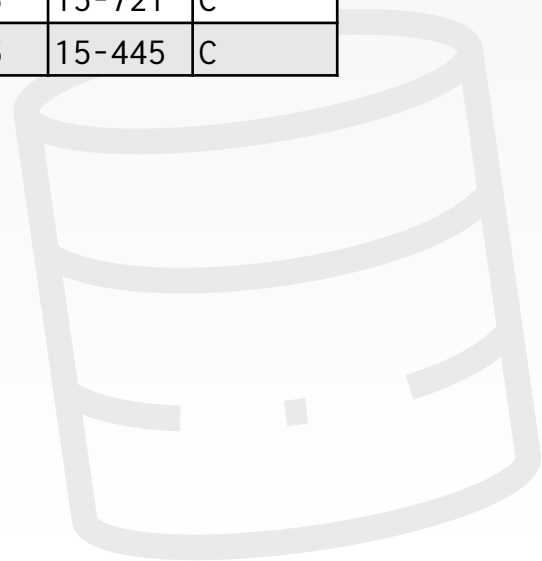


## PHASE #1 – PARTITION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



## PHASE #1 – PARTITION

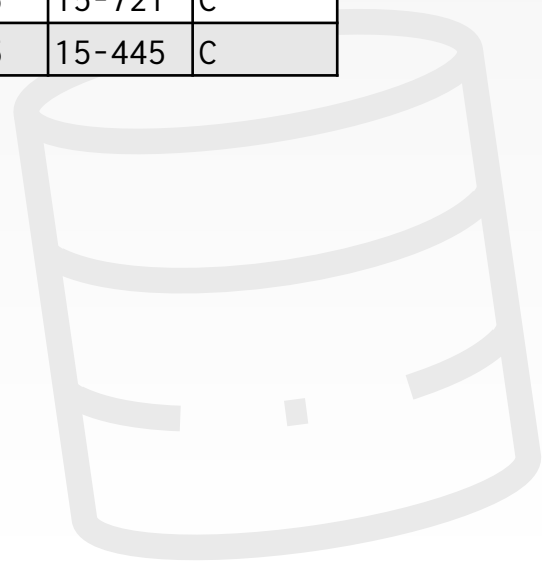
```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

  
*Filter*

sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C

*enrolled(sid,cid,grade)*

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



# PHASE #1 – PARTITION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

  
**Filter**

sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C

  
**Remove  
Columns**

cid
15-445
15-826
15-721
15-445

⋮

# PHASE #1 – PARTITION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

**Filter**

sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C

**Remove Columns**

cid
15-445
15-826
15-721
15-445
⋮

**enrolled(sid,cid,grade)**

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

**B-1 partitions**

$h_1$

15-445 15-826  
15-445 15-445  
15-826 15-445

15-721  
15-721

⋮

# PHASE #1 – PARTITION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

enrolled(sid,cid,grade)

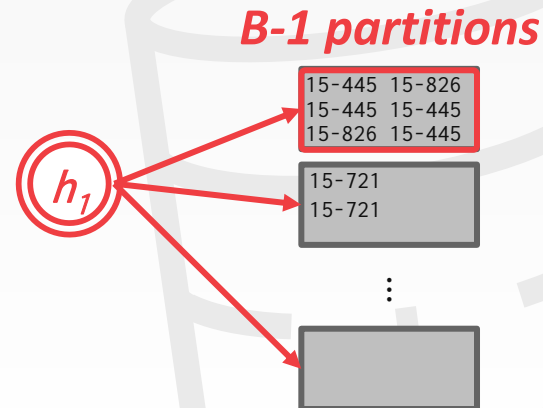
sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C



cid
15-445
15-826
15-721
15-445
⋮





# PHASE #1 – PARTITION

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

**Filter**

sid	cid	grade
53666	15-445	C
53688	15-826	B
53666	15-721	C
53655	15-445	C

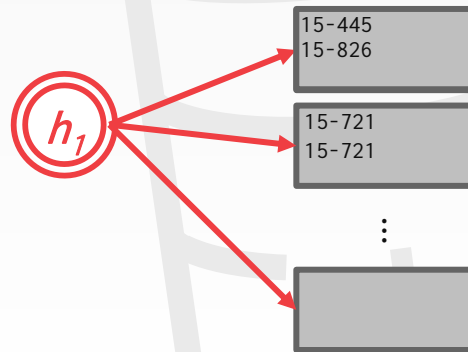
**Remove Columns**

cid
15-445
15-826
15-721
15-445

**enrolled(sid,cid,grade)**

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

**B-1 partitions**



## PHASE #2 – REHASH

---

For each partition on disk:

- Read it into memory and build an in-memory hash table based on a second hash function  $h_2$ .
- Then go through each bucket of this hash table to bring together matching tuples.

This assumes that each partition fits in memory.

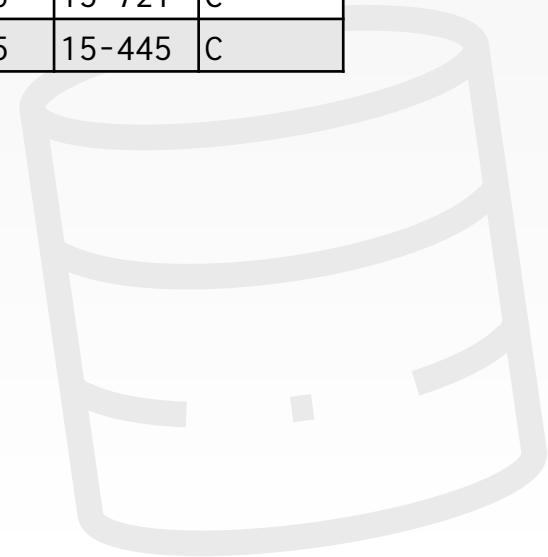


## PHASE #2 – REHASH

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



## PHASE #2 – REHASH

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

### *Phase #1 Buckets*

15-445 15-826  
15-445 15-445  
15-826 15-445

15-445  
15-826

⋮

**enrolled(sid,cid,grade)**

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

## PHASE #2 – REHASH

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

### Phase #1 Buckets

**B-1  
Partitions**

15-445 15-826  
15-445 15-445  
15-826 15-445

15-445  
15-826

⋮

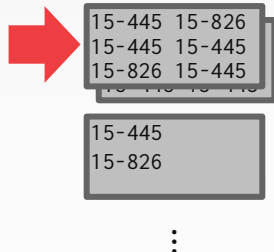
**enrolled(sid,cid,grade)**

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

## PHASE #2 – REHASH

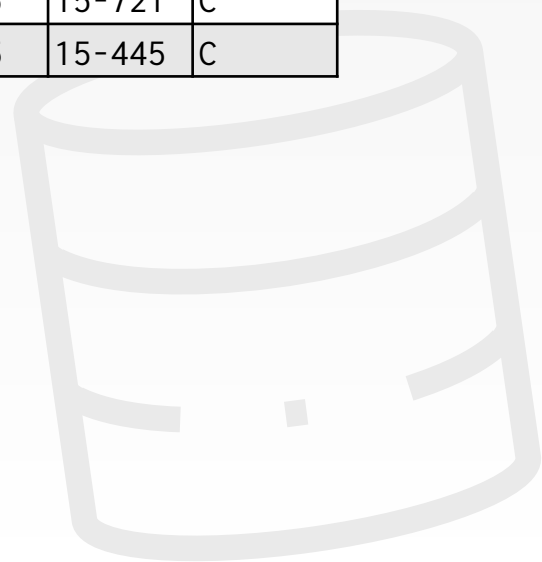
```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

### *Phase #1 Buckets*



**enrolled(sid,cid,grade)**

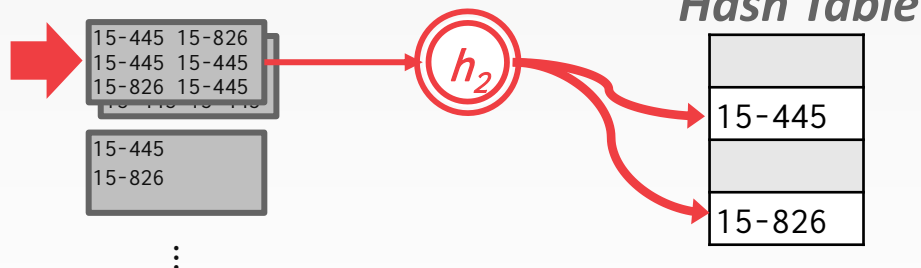
sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



## PHASE #2 – REHASH

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

### Phase #1 Buckets



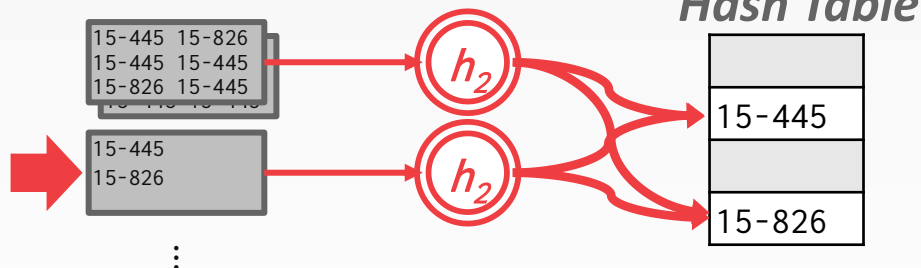
enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

## PHASE #2 – REHASH

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

### Phase #1 Buckets



enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C



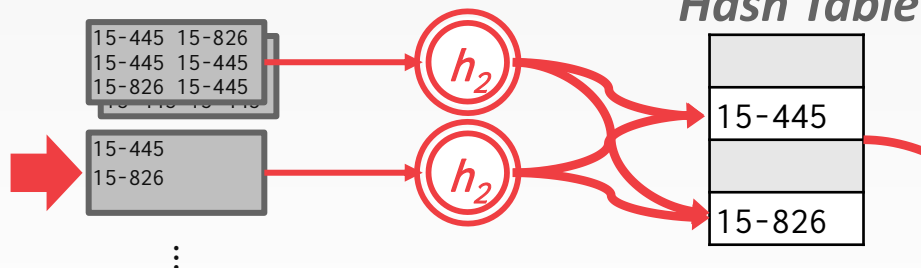
## PHASE #2 – REHASH

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

*Phase #1 Buckets*



*Final Result*

cid
15-445
15-826

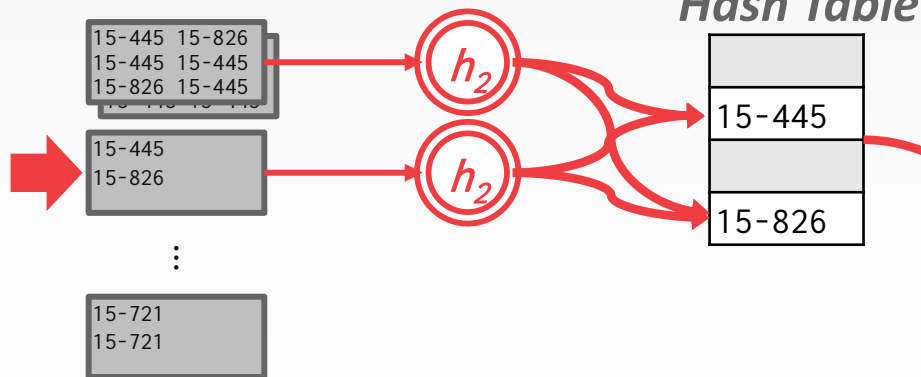
## PHASE #2 – REHASH

```
SELECT DISTINCT cid
FROM enrolled
WHERE grade IN ('B','C')
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

### Phase #1 Buckets



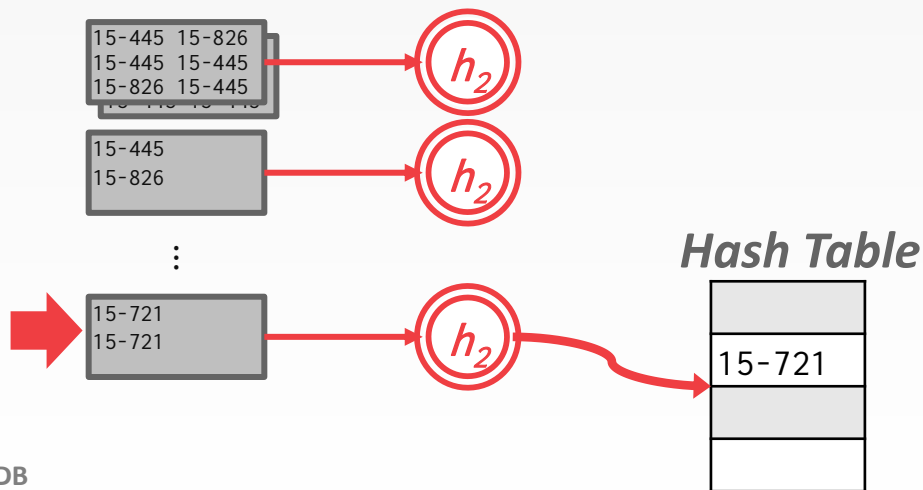
### Final Result

cid
15-445
15-826

## PHASE #2 – REHASH

```
SELECT DISTINCT cid
  FROM enrolled
 WHERE grade IN ('B','C')
```

### Phase #1 Buckets



### enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

### Final Result

cid
15-445
15-826

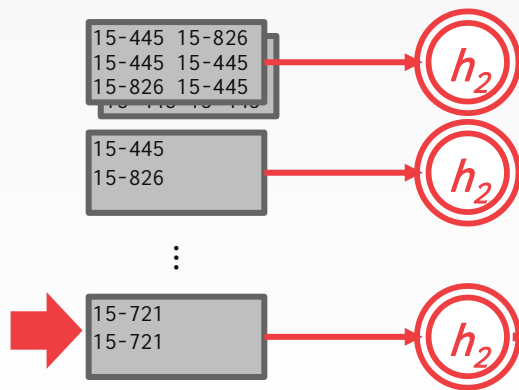
## PHASE #2 – REHASH

```
SELECT DISTINCT cid
FROM enrolled
WHERE grade IN ('B','C')
```

enrolled(sid,cid,grade)

sid	cid	grade
53666	15-445	C
53688	15-721	A
53688	15-826	B
53666	15-721	C
53655	15-445	C

### Phase #1 Buckets



### Hash Table

15-721

### Final Result

cid
15-445
15-826
15-721

# HASHING SUMMARIZATION

---

During the ReHash phase, store pairs of the form  
(**GroupKey**→**RunningVal**)

When we want to insert a new tuple into the hash table:

- If we find a matching **GroupKey**, just update the **RunningVal** appropriately
- Else insert a new **GroupKey**→**RunningVal**



# HASHING SUMMARIZATION

```
SELECT cid, AVG(s.gpa)
  FROM student AS s, enrolled AS e
 WHERE s.sid = e.sid
GROUP BY cid
```

*Phase #1  
Buckets*

15-445  
15-445

15-826

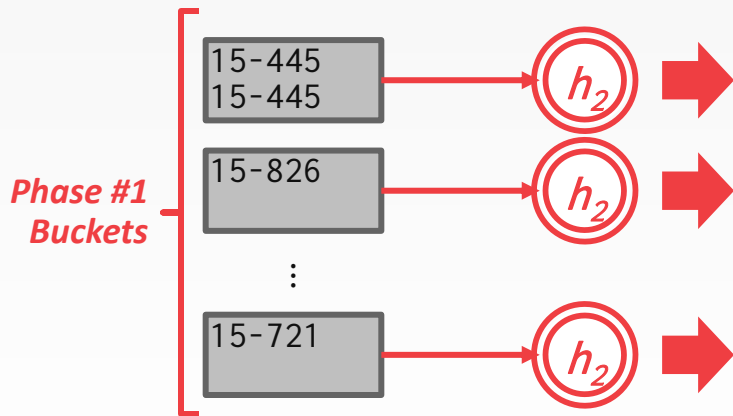
⋮

15-721



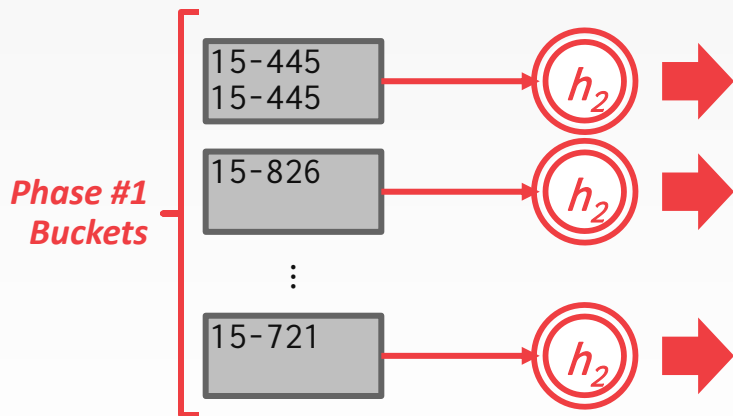
# HASHING SUMMARIZATION

```
SELECT cid, AVG(s.gpa)
  FROM student AS s, enrolled AS e
 WHERE s.sid = e.sid
GROUP BY cid
```



# HASHING SUMMARIZATION

```
SELECT cid, AVG(s.gpa)
  FROM student AS s, enrolled AS e
 WHERE s.sid = e.sid
 GROUP BY cid
```



*Hash Table*

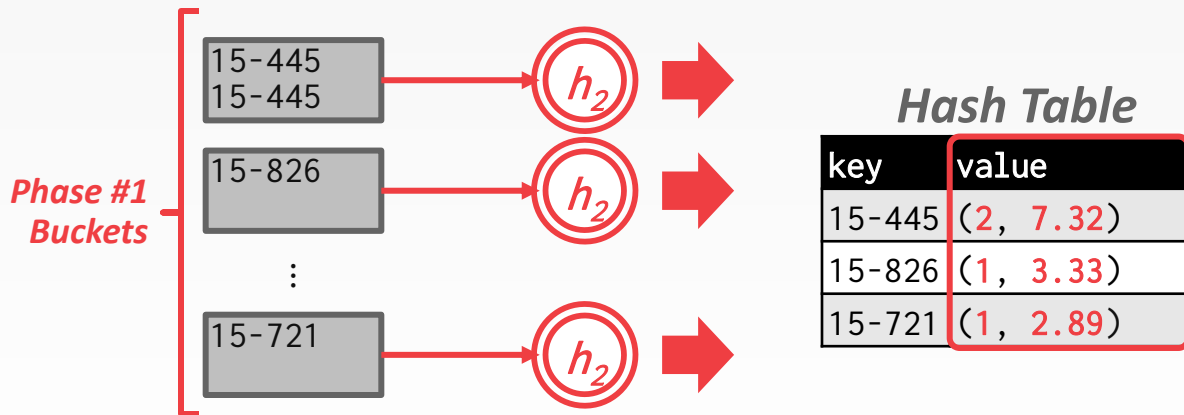
key	value
15-445	(2, 7.32)
15-826	(1, 3.33)
15-721	(1, 2.89)





# HASHING SUMMARIZATION

```
SELECT cid, AVG(s.gpa)
  FROM student AS s, enrolled AS e
 WHERE s.sid = e.sid
 GROUP BY cid
```

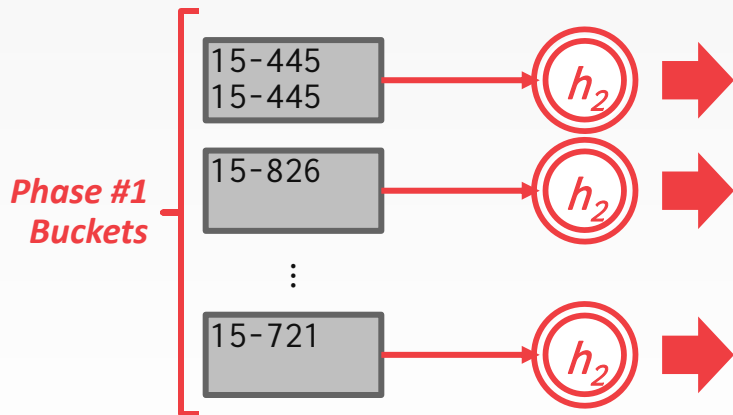


# HASHING SUMMARIZATION

```
SELECT cid, AVG(s.gpa)
  FROM student AS s, enrolled AS e
 WHERE s.sid = e.sid
 GROUP BY cid
```

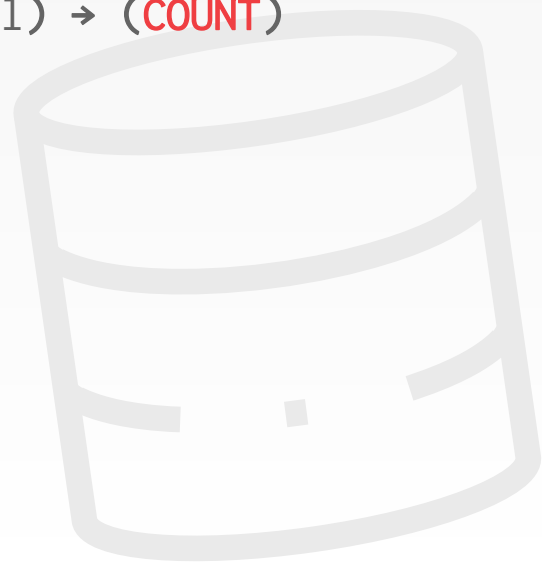
## Running Totals

AVG(col) → (COUNT, SUM)  
 MIN(col) → (MIN)  
 MAX(col) → (MAX)  
 SUM(col) → (SUM)  
 COUNT(col) → (COUNT)



## Hash Table

key	value
15-445	(2, 7.32)
15-826	(1, 3.33)
15-721	(1, 2.89)

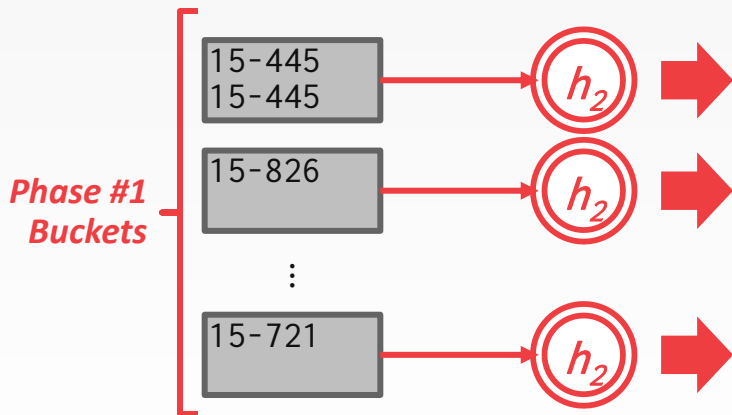


# HASHING SUMMARIZATION

```
SELECT cid, AVG(s.gpa)
  FROM student AS s, enrolled AS e
 WHERE s.sid = e.sid
 GROUP BY cid
```

## Running Totals

AVG(col) → (COUNT, SUM)  
 MIN(col) → (MIN)  
 MAX(col) → (MAX)  
 SUM(col) → (SUM)  
 COUNT(col) → (COUNT)



## Hash Table

key	value
15-445	(2, 7.32)
15-826	(1, 3.33)
15-721	(1, 2.89)

## Final Result

cid	AVG(gpa)
15-445	3.66
15-826	3.33
15-721	2.89

## CONCLUSION

---

Choice of sorting vs. hashing is subtle and depends on optimizations done in each case.

High-level principle to handle data that does not fit into memory:

→ Chunk data into partitions that fit into memory to handle separately



## NEXT CLASS

---

Nested Loop Join

Sort-Merge Join

Hash Join

