



INFO20003 Database Systems

Dr Renata Borovica-Gajic

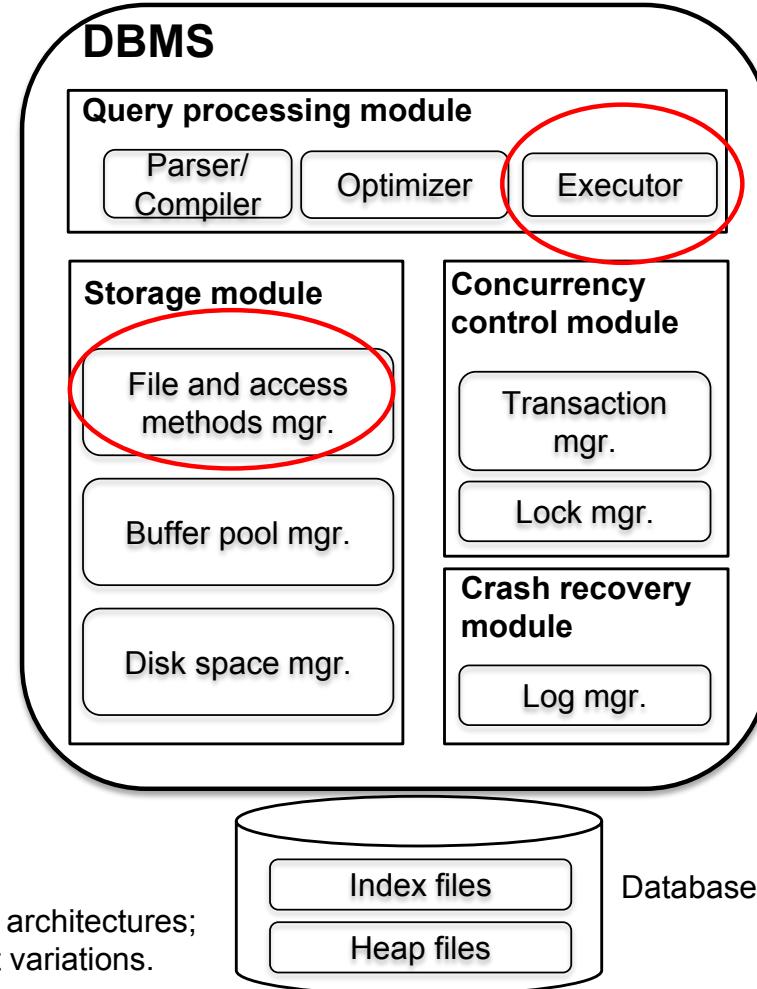
Lecture 11
Query Processing Part I

Week 6



Remember this? Components of a DBMS

Will briefly
touch upon ...



This is one of several possible architectures;
each system has its own slight variations.

- Query Processing Overview
- Selections
- Projections

Readings: Chapter 12 and 14, Ramakrishnan & Gehrke, Database Systems



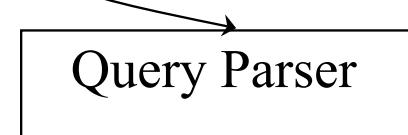
- Some database operations are **EXPENSIVE**
 - one operation on a single record is not so expensive
- DBMSs can greatly improve performance by being ‘smart’
 - cumulative → every operation is optimized
 - e.g., can speed up 1,000,000x over naïve approach
- Main weapons are:
 1. clever implementation techniques for operators
 2. exploiting ‘equivalencies’ of relational operators
 3. using cost models to choose among alternatives



MELBOURNE

Query

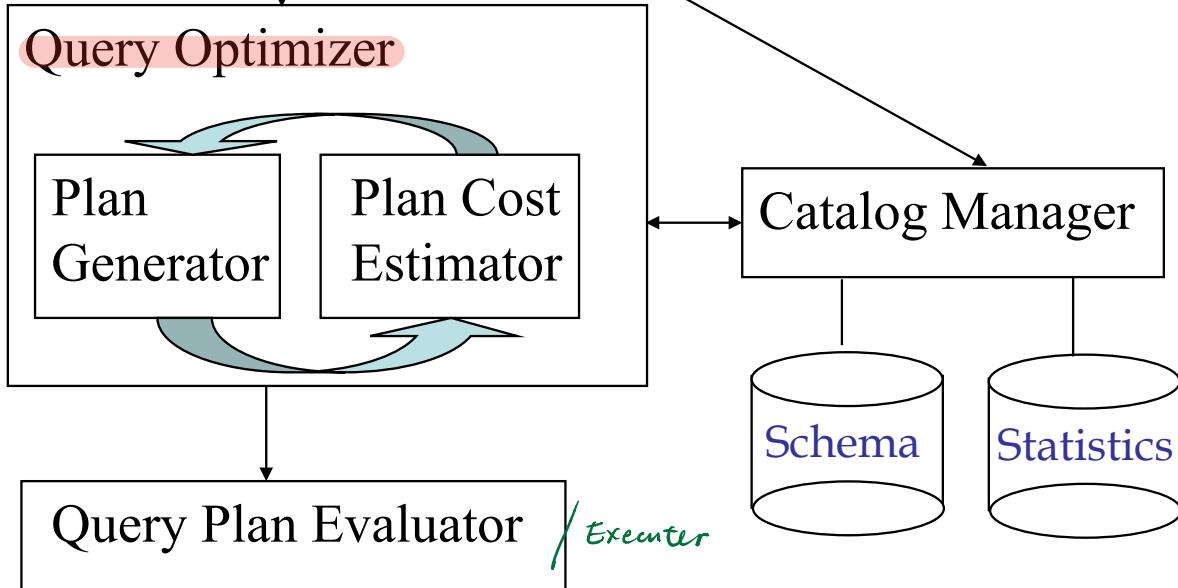
```
Select *\nFrom Blah B\nWhere B.blah = "foo"
```



check syntax

Usually there is a heuristics-based **rewriting step** before the cost-based steps.

Next week



MELBOURNE

- We will consider how to implement:
 - Selection (σ) Selects a subset of rows from relation
 - Projection (π) Deletes unwanted columns from relation
 - Join (\bowtie) Allows us to combine two relations
- Operators can be then be *composed* creating *query plans*



- Query Processing Overview
- Selections
- Projections

Readings: Chapter 14, Ramakrishnan & Gehrke, Database Systems



Sailors (sid: integer, sname: string, rating: integer, age: real)
Reserves (sid: integer, bid: integer, day: dates, rname: string)

- **Sailors (S):**

- Each tuple is 50 bytes long, 80 tuples per page, **500 pages**
- $N = NPages(S) = 500$, $p_S = NTuplesPerPage(S) = 80$
- $NTuples(S) = 500 * 80 = 40000$

- **Reserves (R):**

- Each tuple is 40 bytes long, 100 tuples per page, **1000 pages**
- $M = NPages(R) = 1000$, $p_R = NTuplesPerPage(R) = 100$
- $NTuples(R) = 100000$



MELBOURNE

- Of the form $\sigma_{R.attr \ op \ value}(R)$
- Example:

```
SELECT *
FROM Reserves R
WHERE R.BID > 20;
```

- The best way to perform a selection depends on:
 1. available indexes/access paths
 2. **expected size of the result** (number of tuples and/or number of pages)



- **Size of result** approximated as:

*size of relation * \prod (reduction factors)*

- **Reduction factor** is usually called **selectivity**. It estimates what portion of the relation will qualify for the given predicate, i.e. satisfy the given condition.
 - This is estimated by the optimizer (will be taught next week)
 - E.g. 30% of records qualify, or 5% of records qualify



1. With no index, unsorted:

—Must scan the whole relation, i.e. perform **Heap Scan**

—**Cost = Number of Pages of Relation, i.e. NPages(R)** *traverse the whole*

—**Example:** Reserves cost(R) = 1000 IO (1000 pages)

2. With no index, but file is sorted:

find the first record which qualify

—cost = **binary search cost + number of pages containing results**

—**Cost = $\log_2(NPages(R)) + (RF \times NPages(R))$**

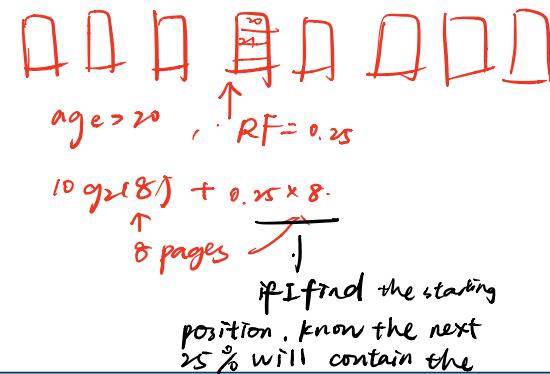
—**Example:** Reserves cost(R) = 10 I/O + ($RF \times NPages(R)$)

3. With an index on selection attribute:

—Use index to find qualifying data entries,

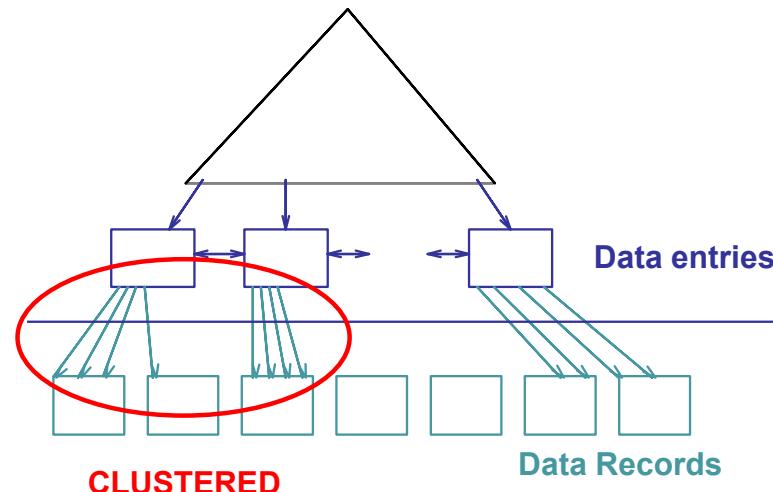
—Then retrieve corresponding data records

—Discussed next....

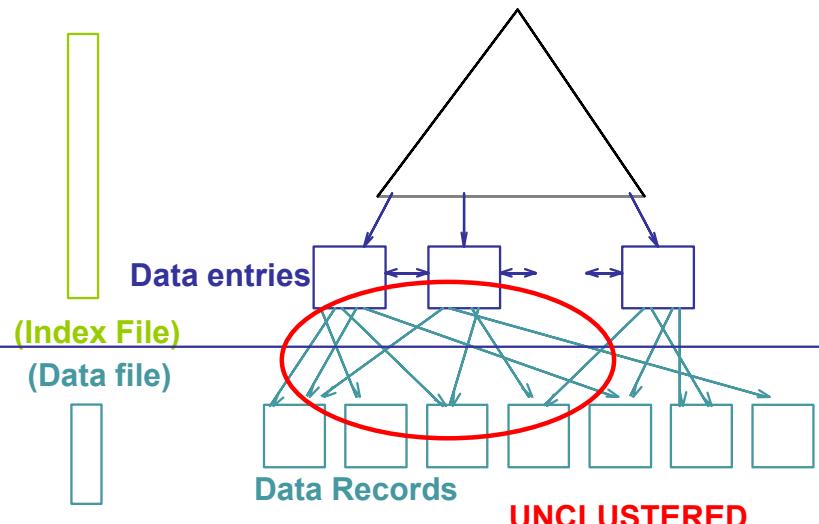




Clustered vs. unclustered



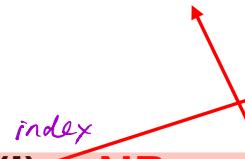
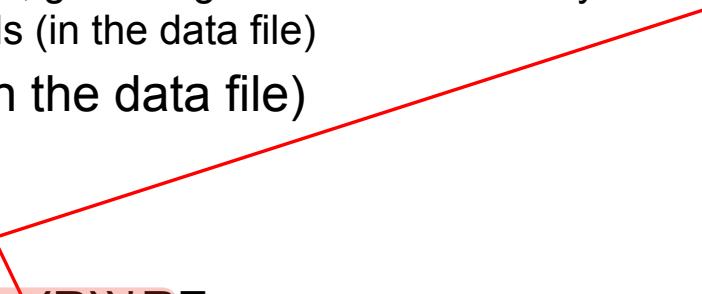
① does this mean
the data file is sorted?



for each record, one I/O operation
(records in different pages)



MELBOURNE

- Cost depends on the number of qualifying tuples
 - Clustering is important when calculating the total cost
 - Steps to perform:
 1. Find qualifying data entries:
 - Go through the index: height typically small, 2-4 I/O in case of B+tree, 1.2 I/O in case of hash index (*negligible* if many records retrieved)
 - Once data entries are reached, go through data entries one by one and look up corresponding data records (in the data file)
 2. Retrieve data records (in the data file)
 - **Cost:**
 1. Clustered index:
 $\text{Cost} = (\text{NPages}(I) + \text{NPages}(R)) * \text{RF}$
 2. Unclustered index:
 $\text{Cost} = (\text{NPages}(I) + \text{NTuples}(R)) * \text{RF}$
- 
- 

MELBOURNE

- **Example:** Let's say that 10% of Reserves tuples qualify, and let's say that index occupies 50 pages
- $RF = 10\% = 0.1$, $NPages(I) = 50$, $NPages(R) = 1000$, $NTuplesPerPage(R) = 100$

• Cost:

1. Clustered index:

$$\text{Cost} = (NPages(I) + NPages(R)) * RF$$

$$\text{Cost} = (50 + 1000) * 0.1 = \textcolor{red}{105 \text{ (I/O)}}$$

clustered

$$(50 + 1000) * 0.1 = 105$$

unclustered

$$(50 + 100000) * 0.1 = 10005$$

Cheapest access path

2. Unclustered index:

$$\text{Cost} = (NPages(I) + NTuples(R)) * RF$$

$$\text{Cost} = (50 + 100000) * 0.1 = \textcolor{red}{10005 \text{ (I/O)}}$$

3. Heap Scan:

$$\text{Cost} = NPages(R) = 1000 \text{ (I/O)}$$



- Typically queries have multiple predicates (conditions)
- **Example:** day<8/9/94 AND rname='Paul' AND bid=5 AND sid=3
- A B-tree index **matches** (a combination of) predicates that involve only attributes in a **prefix of the search key**
 - Index on $\langle a, b, c \rangle$ matches predicates on: (a,b,c) , (a,b) and (a)
 - Index on $\langle a, b, c \rangle$ matches $a=5 \text{ AND } b=3$, but will not be used to answer $b=3$ $\rightarrow \text{not a prefix}$
 - This implies that only reduction factors of predicates that are part of the prefix will be used to determine the cost (they are called matching predicates (or primary conjuncts))

a	b
1	1
1	2
1	7
2	2
2	8
3	1
7	7
8	6

I(a,b)

① a=3 ✓ RF(a)

② a=1 and b>2 ✓ RF(a) F(b)

③ b=1 X —

④ a=1 and b>3 and c=7

(suppose
table has c
but index don't have c)

c will be checked once
the data is already brought from index
it is not the cost of index

MELBOURNE

1. Find the **cheapest access path**
 - An index or file scan with the **least estimated page I/O**

2. Retrieve tuples using it
 - Predicates that match this index reduce the number of tuples *retrieved (and impact the cost)* *doesn't contribute to the cost*
[cost is in I/O operations].

3. Apply the predicates that **don't** match the index (if any) later on
 - These predicates are used to discard some retrieved tuples, but do not affect number of tuples/pages fetched (nor the total cost)
 - In this case selection over other predicates is said to be done “on-the-fly”



MELBOURNE

- **Example:** $\text{day} < 8/9/94 \text{ AND } \text{bid}=5 \text{ AND } \text{sid}=3$
- A **B+ tree index on day** can be used;
 - $\text{RF} = \text{RF}(\text{day})$
 - Then, $\text{bid}=5$ and $\text{sid}=3$ must be checked for each retrieved tuple *on the fly*
- Similarly, a **hash index on $\langle \text{bid}, \text{sid} \rangle$** could be used;
 - $\prod \text{RF} = \text{RF}(\text{bid}) * \text{RF}(\text{sid})$
 - Then, $\text{day} < 8/9/94$ must be checked *on the fly*
- How about a B+tree on $\langle \text{rname}, \text{day} \rangle$? (Y/N)
- How about a B+tree on $\langle \text{day}, \text{rname} \rangle$? (Y/N)
- How about a Hash index on $\langle \text{day}, \text{rname} \rangle$? (Y/N)

↑
use equality
condition

it can't speed up
 $\text{day} \rightarrow \text{range condition}$

$\text{RF}(\text{day}).$
 $\text{RF}(\text{day}).$

X
✓
X



- Overview
- Selections
- Projections

Readings: Chapter 14, Ramakrishnan & Gehrke, Database Systems

The Projection Operation

MELBOURNE

- Issue with projection is removing duplicates

① sort data
↓
duplicate will be adjacent

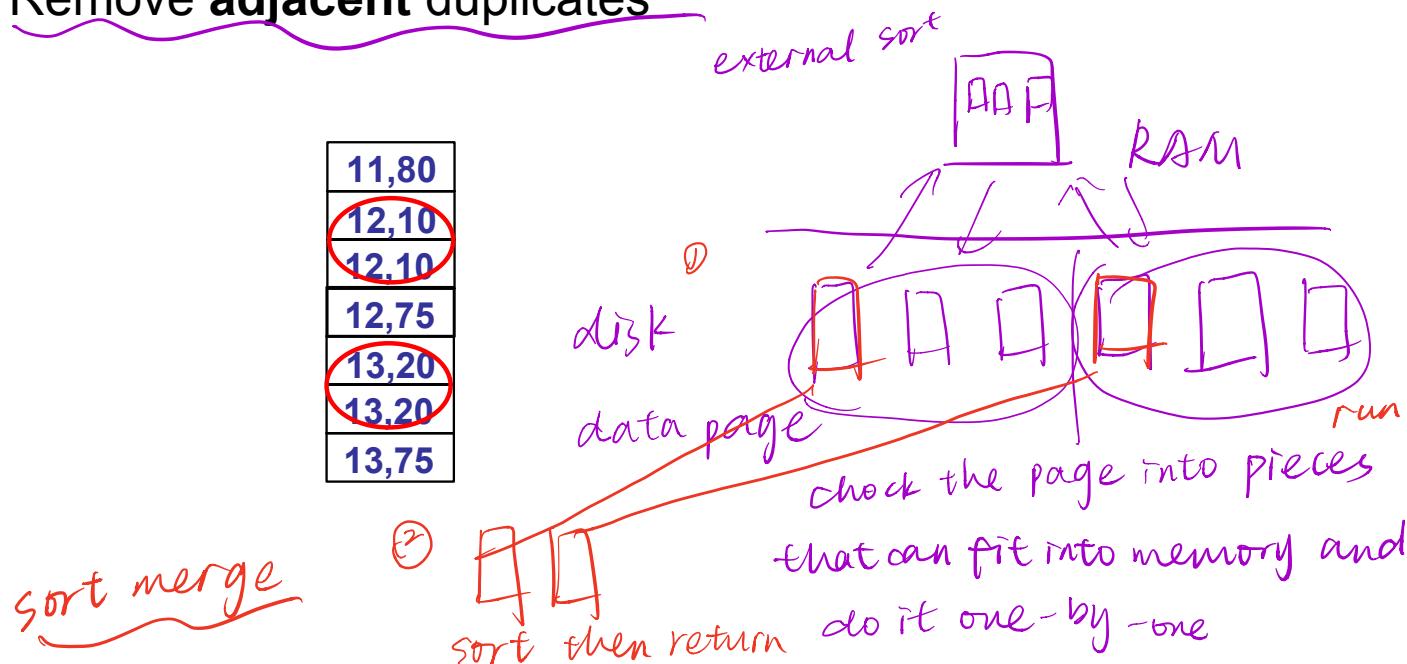
```
SELECT DISTINCT R.sid, R.bid  
FROM Reserves R
```

my sql doesn't remove duplicates
by default

② hash data
duplicate will in the
same bucket

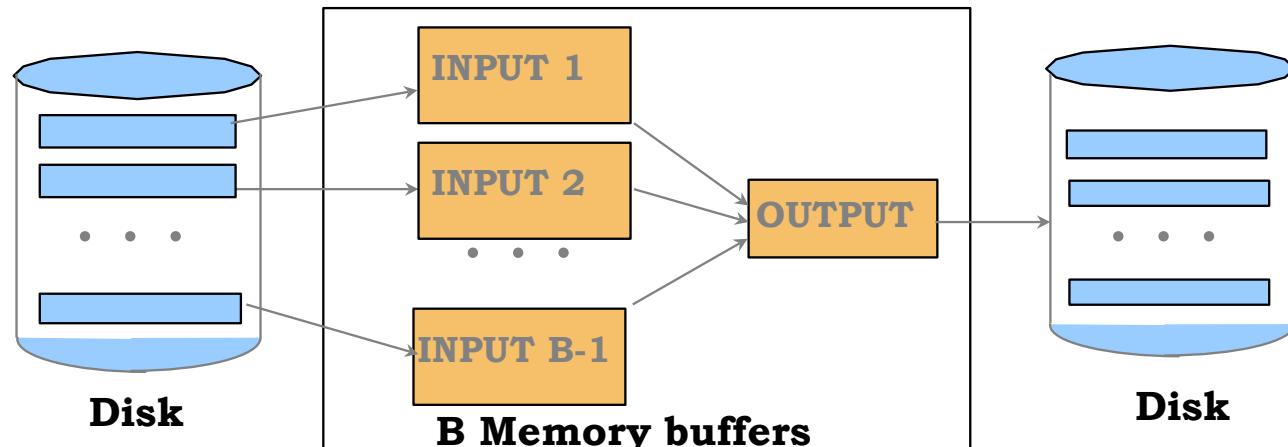
- Projection can be done based on **hashing** or **sorting**

- Basic approach is to use **sorting**
 - 1. Scan R, extract only the **needed** attributes
 - 2. Sort the result set (typically using external merge sort)
 - 3. Remove **adjacent** duplicates



MELBOURNE

- If data does not fit in memory do several passes
 - Sort runs: Make each B pages sorted (called runs)
 - Merge runs: Make multiple passes to merge runs
 - Pass 2: Produce runs of length $B(B-1)$ pages
 - Pass 3: Produce runs of length $B(B-1)^2$ pages
 - ...
 - Pass P: Produce runs of length $B(B-1)^P$ pages
- We will let you know
how many passes there are

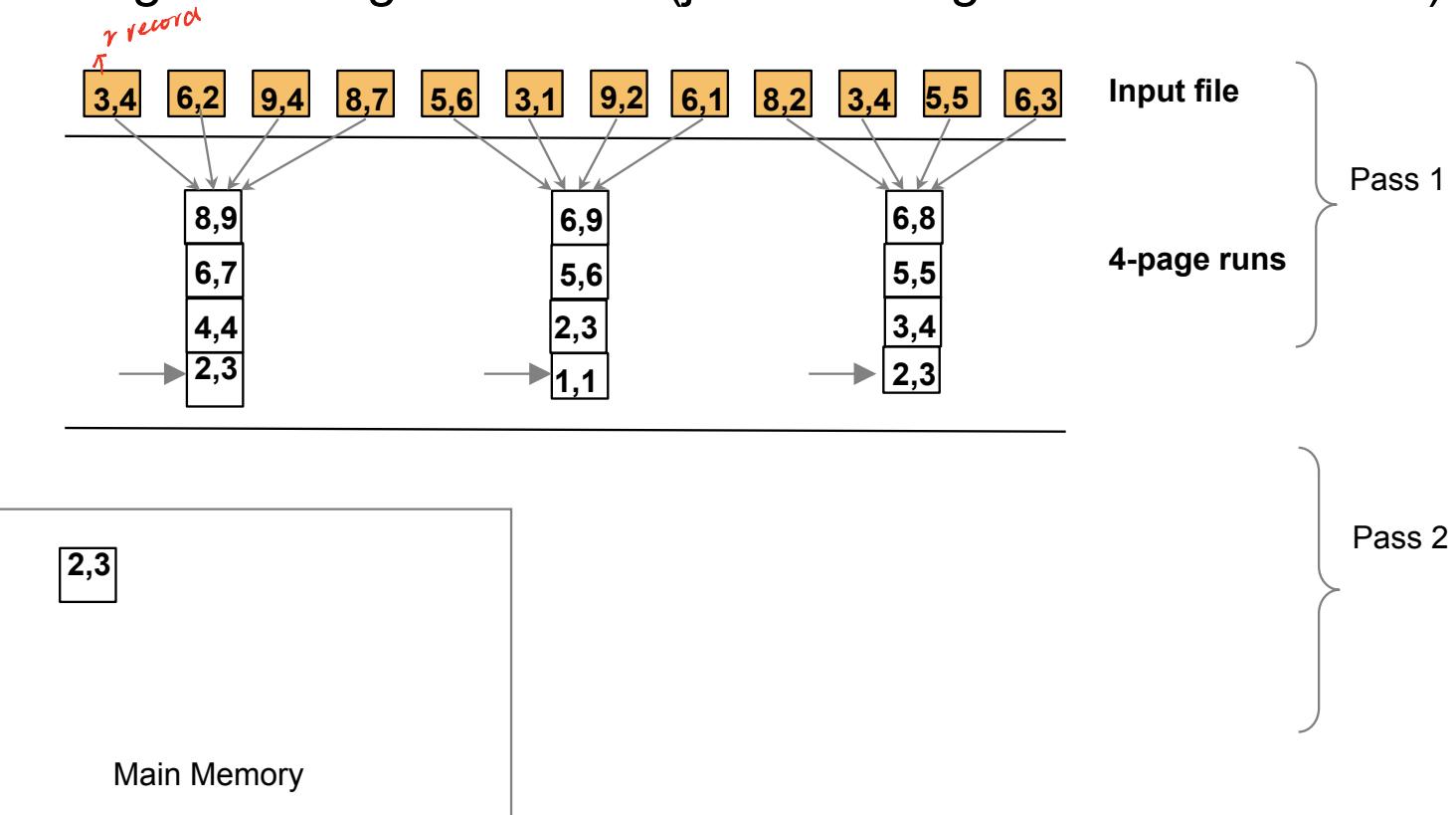


Readings: Chapter 13, Ramakrishnan & Gehrke, Database Systems

External Merge Sort: Example

MELBOURNE

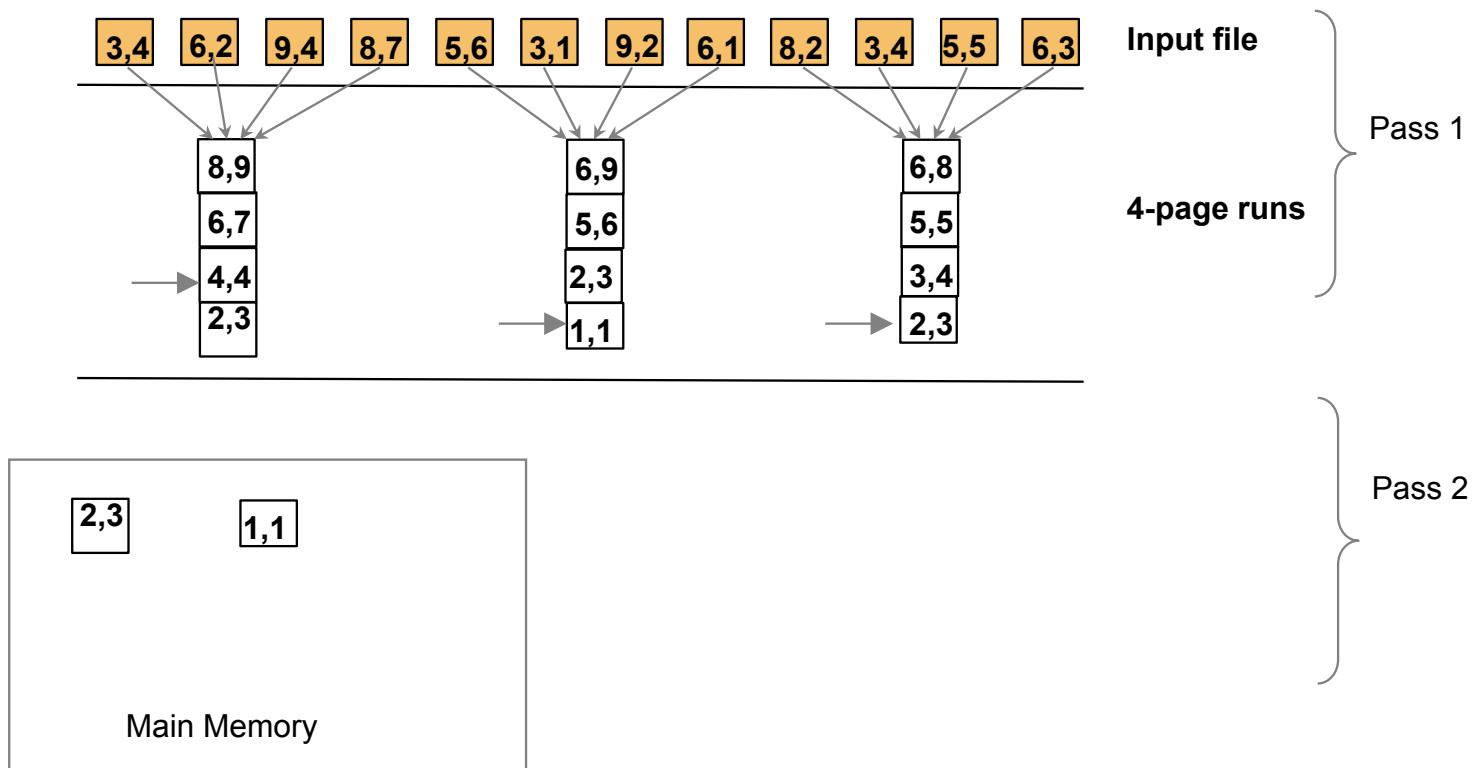
buffer pages in memory $B = 4$, each page 2 records,
sorting on a single attribute (just showing the attribute value)



External Merge Sort: Example

MELBOURNE

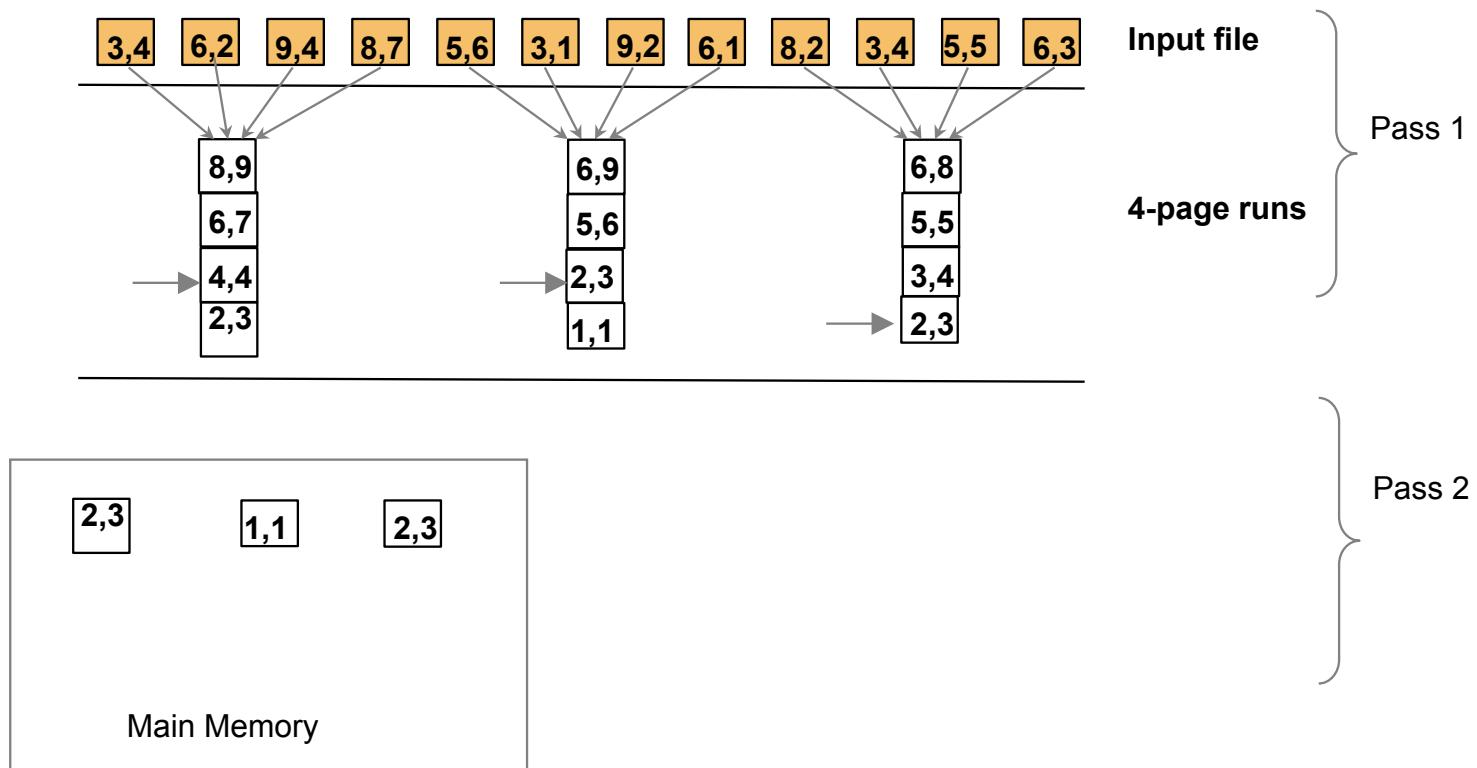
buffer pages in memory $B = 4$, each page 2 records





MELBOURNE

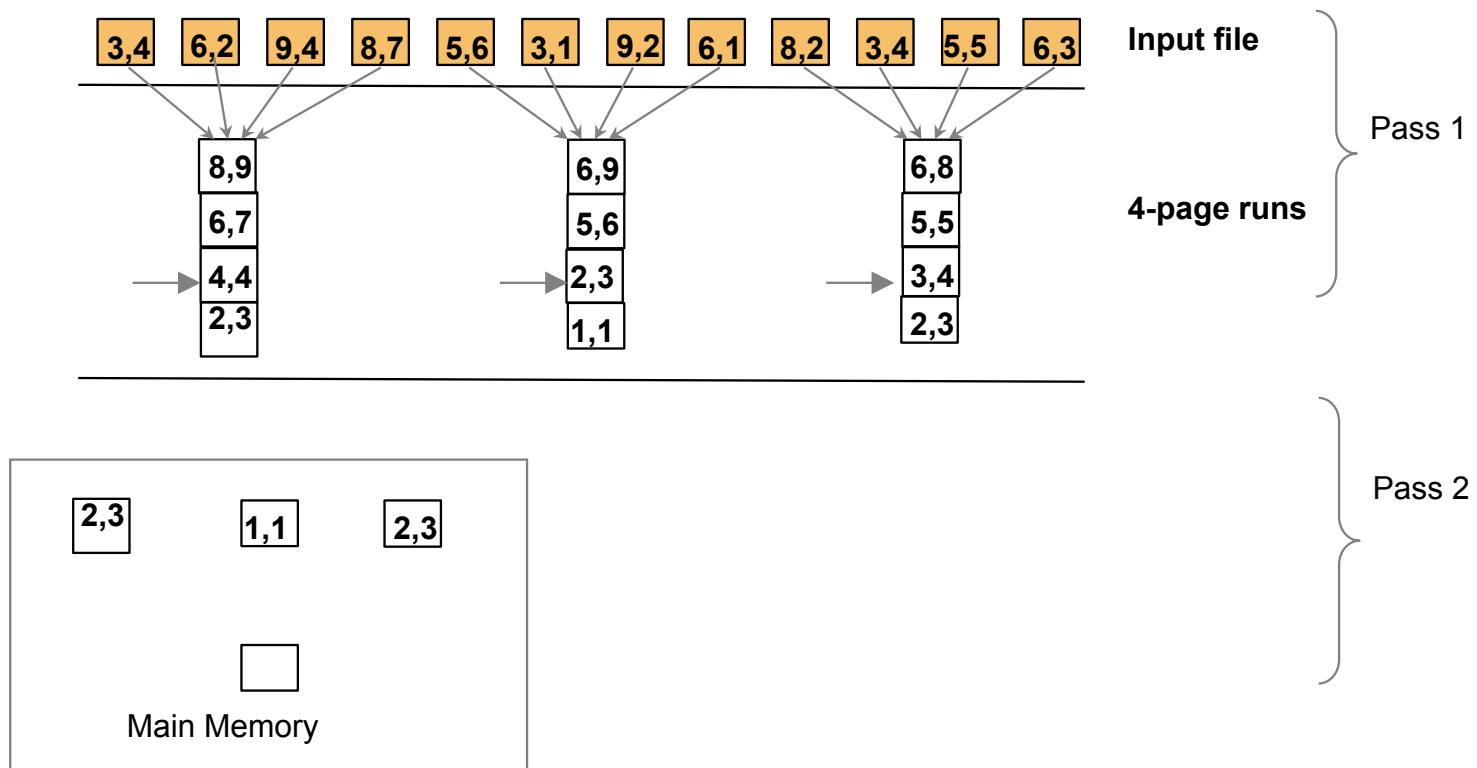
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

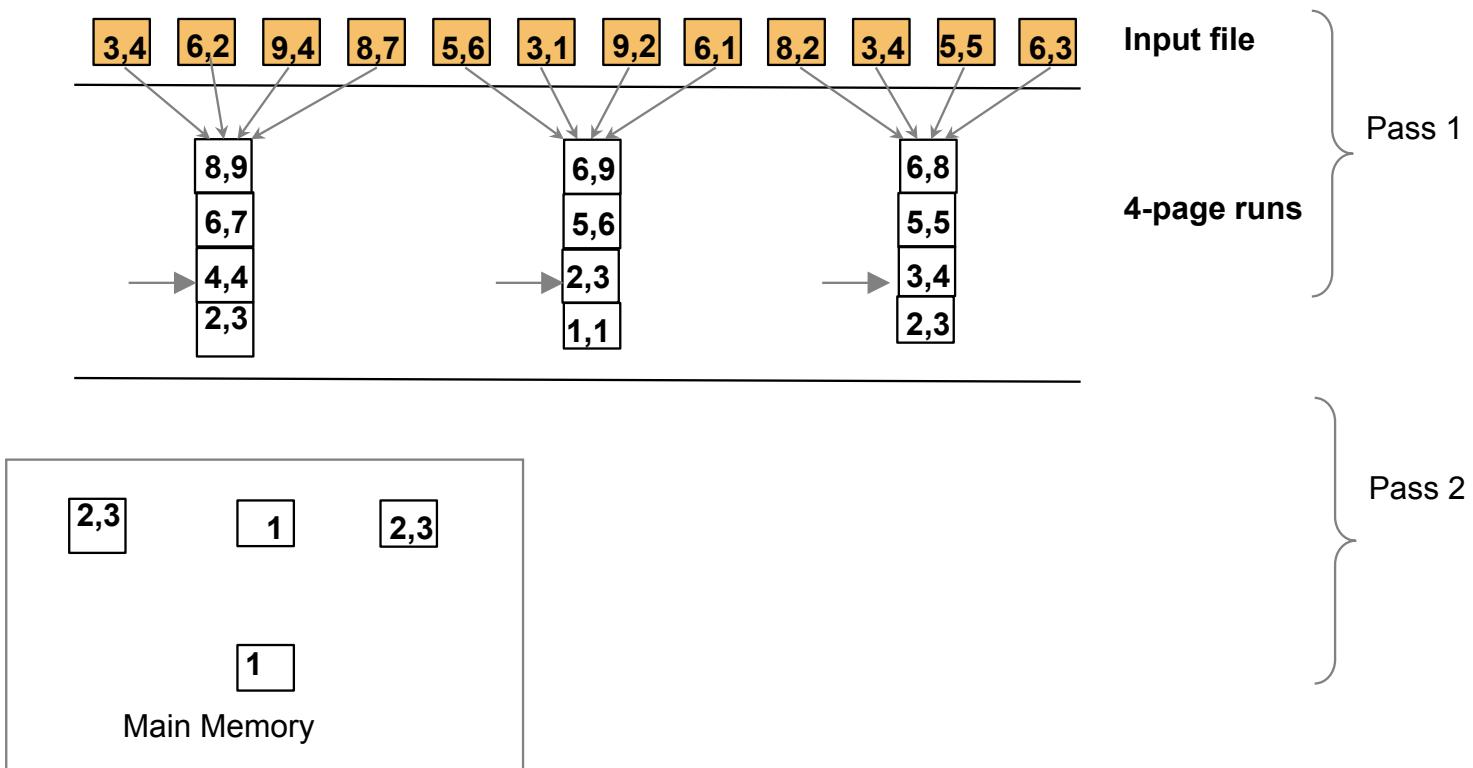
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

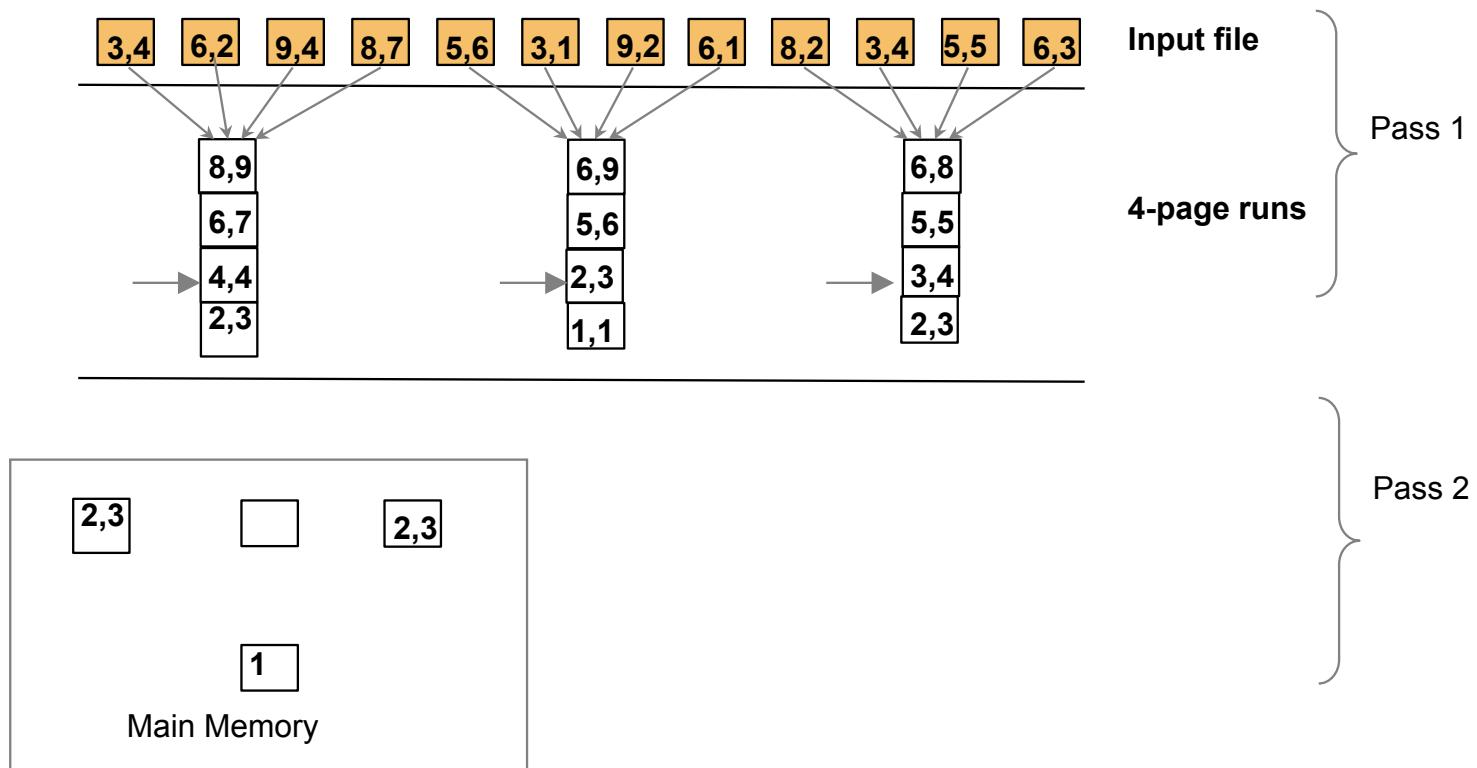
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

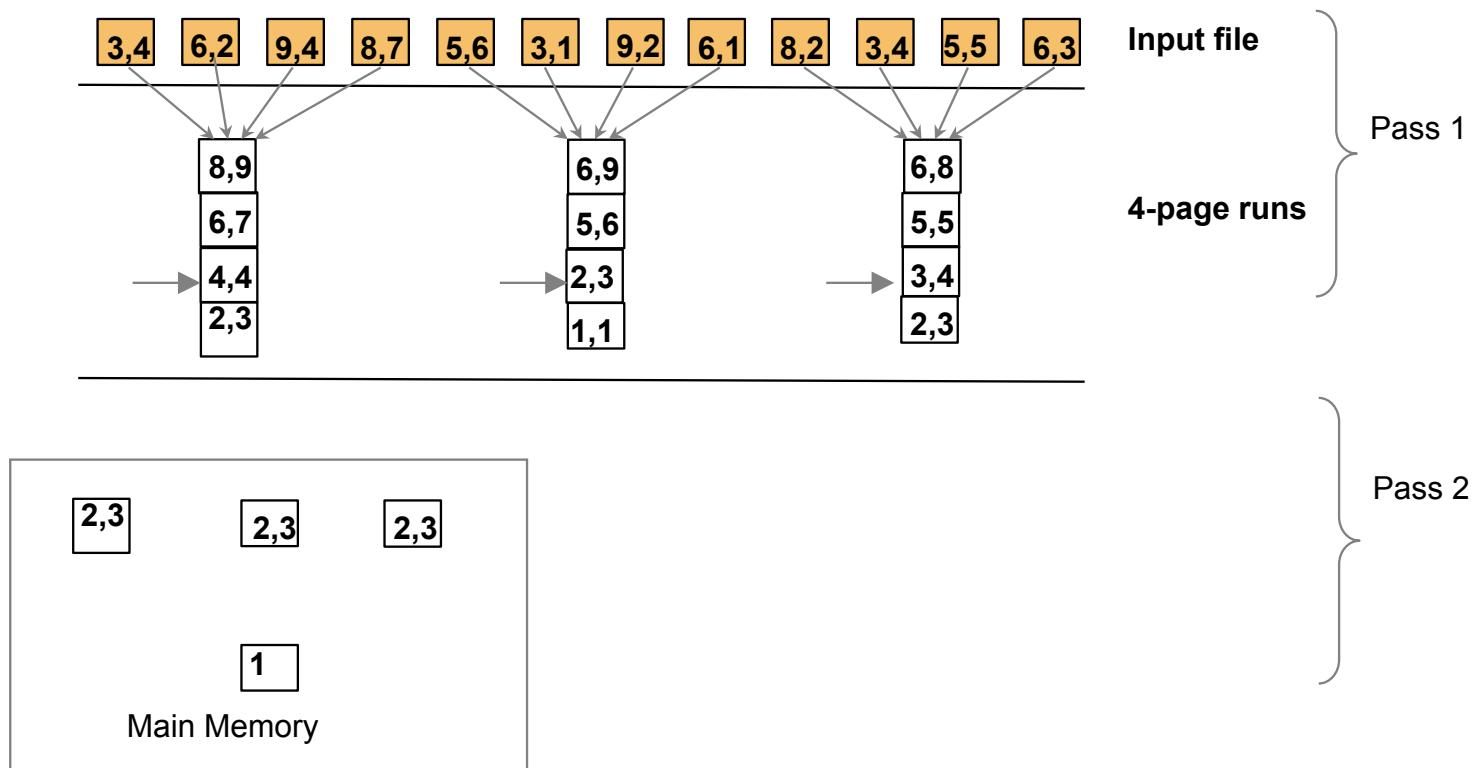
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

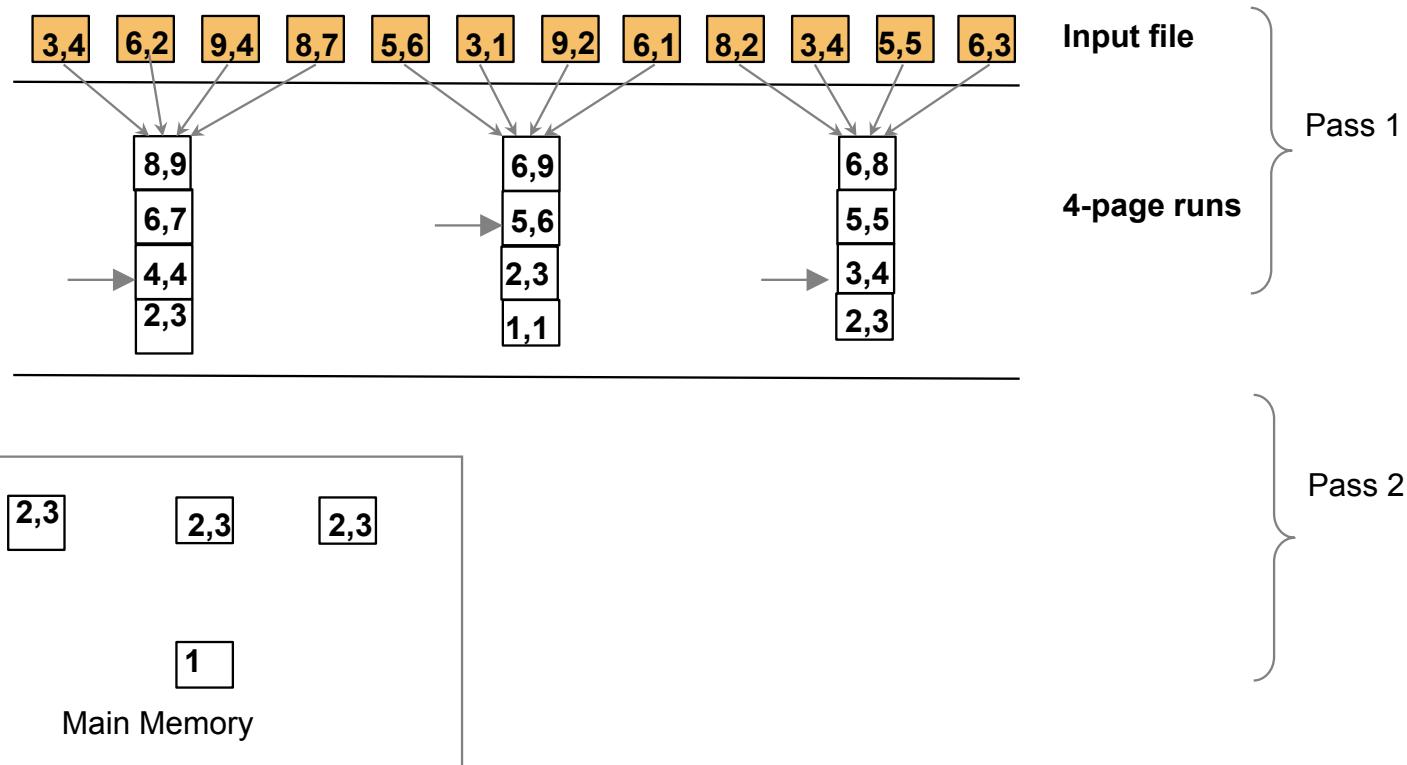
buffer pages in memory $B = 4$, each page 2 records





MELBOURNE

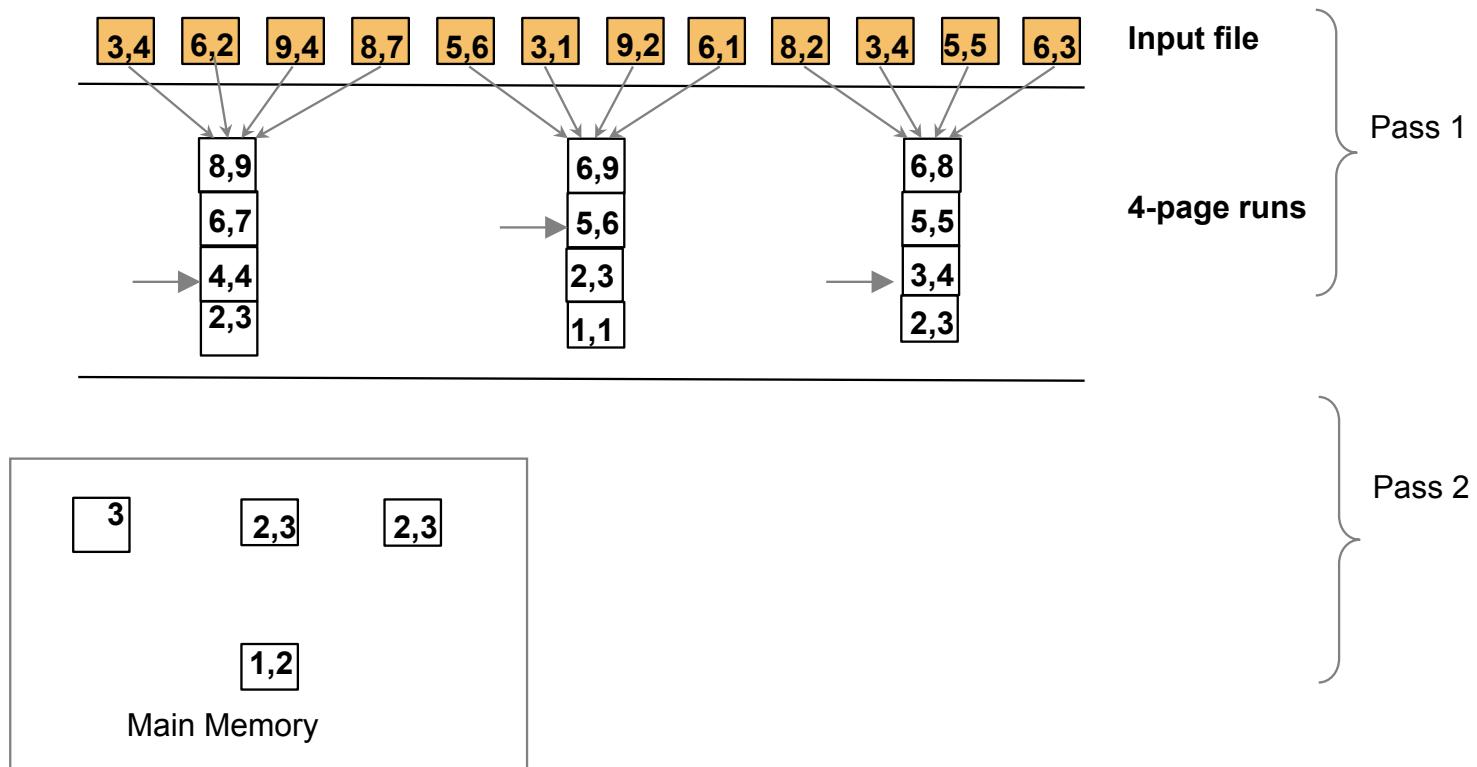
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

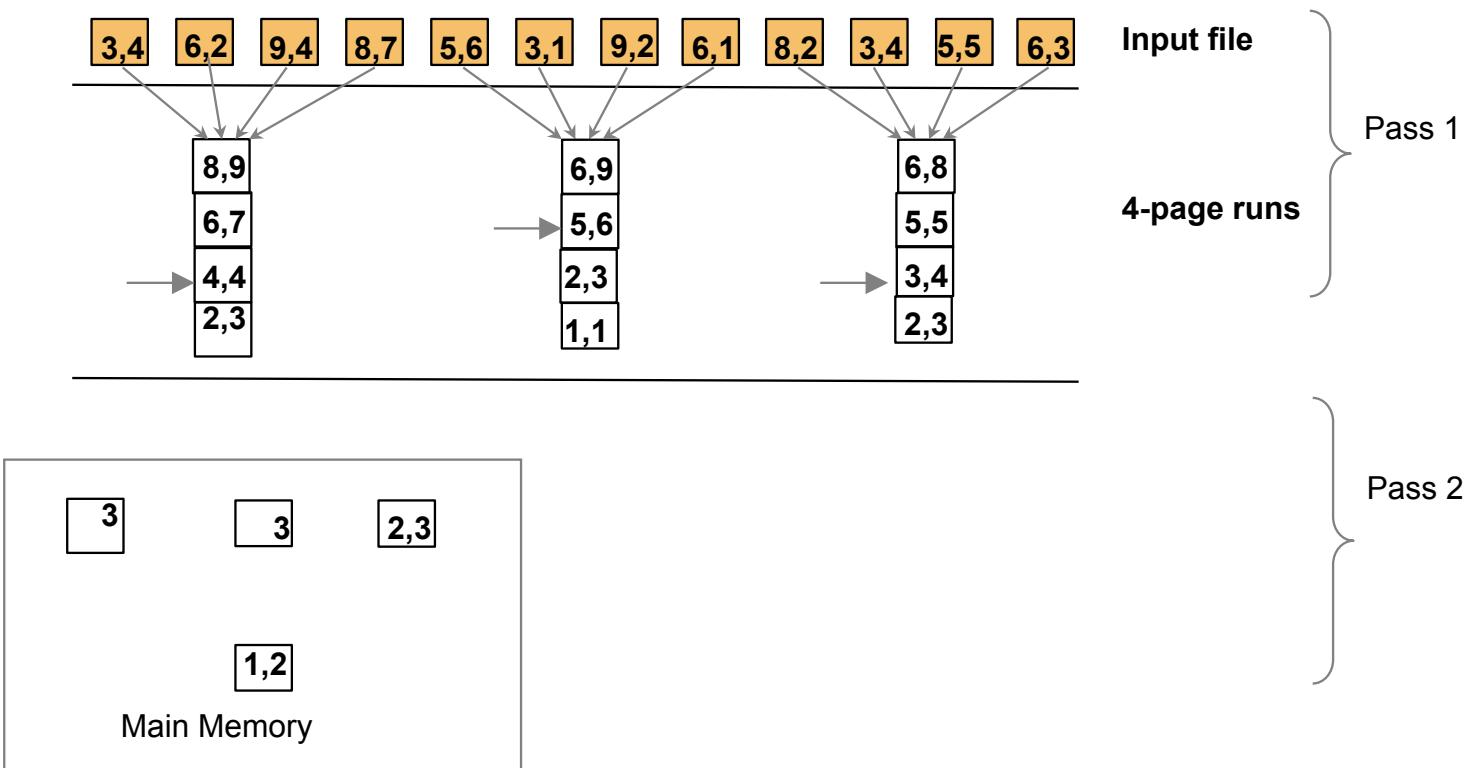
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

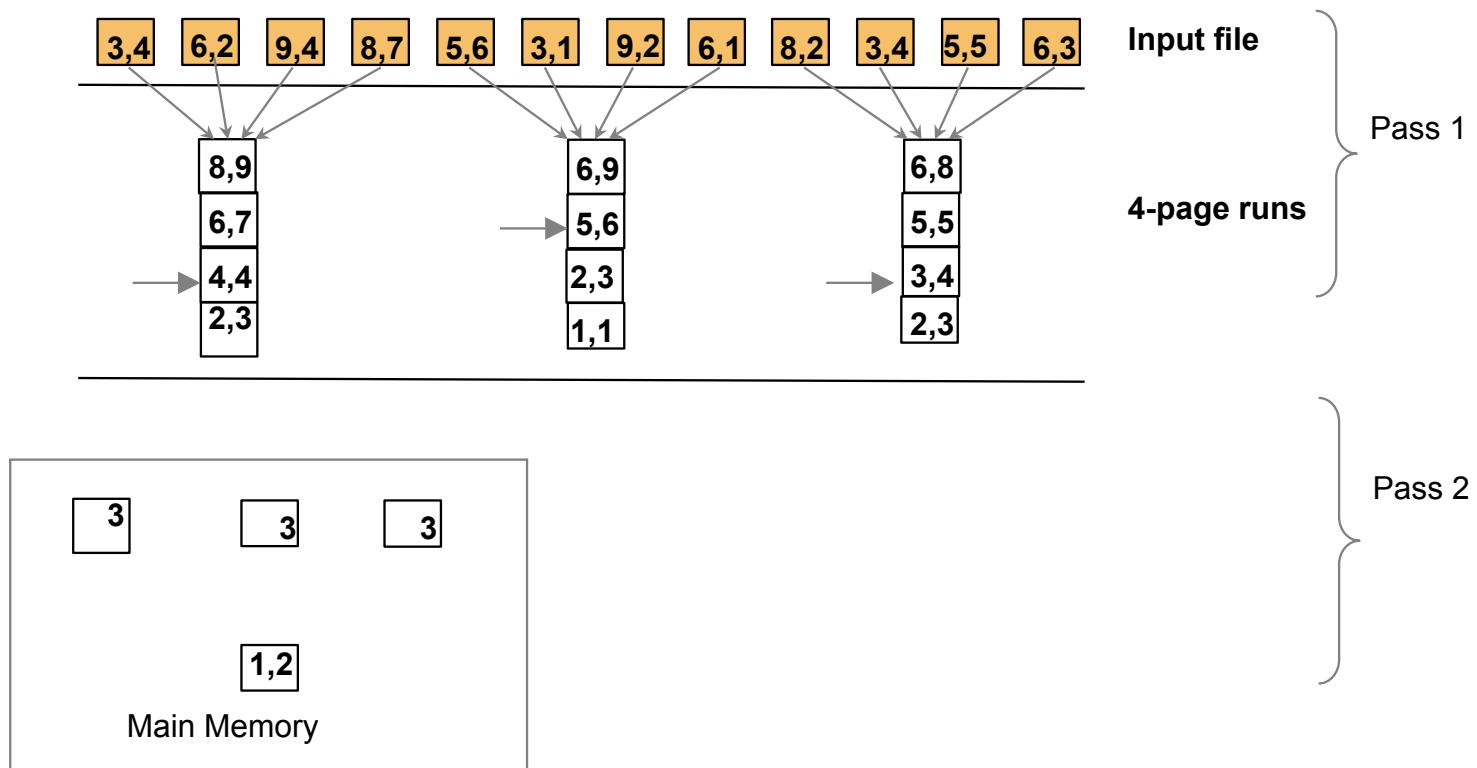
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

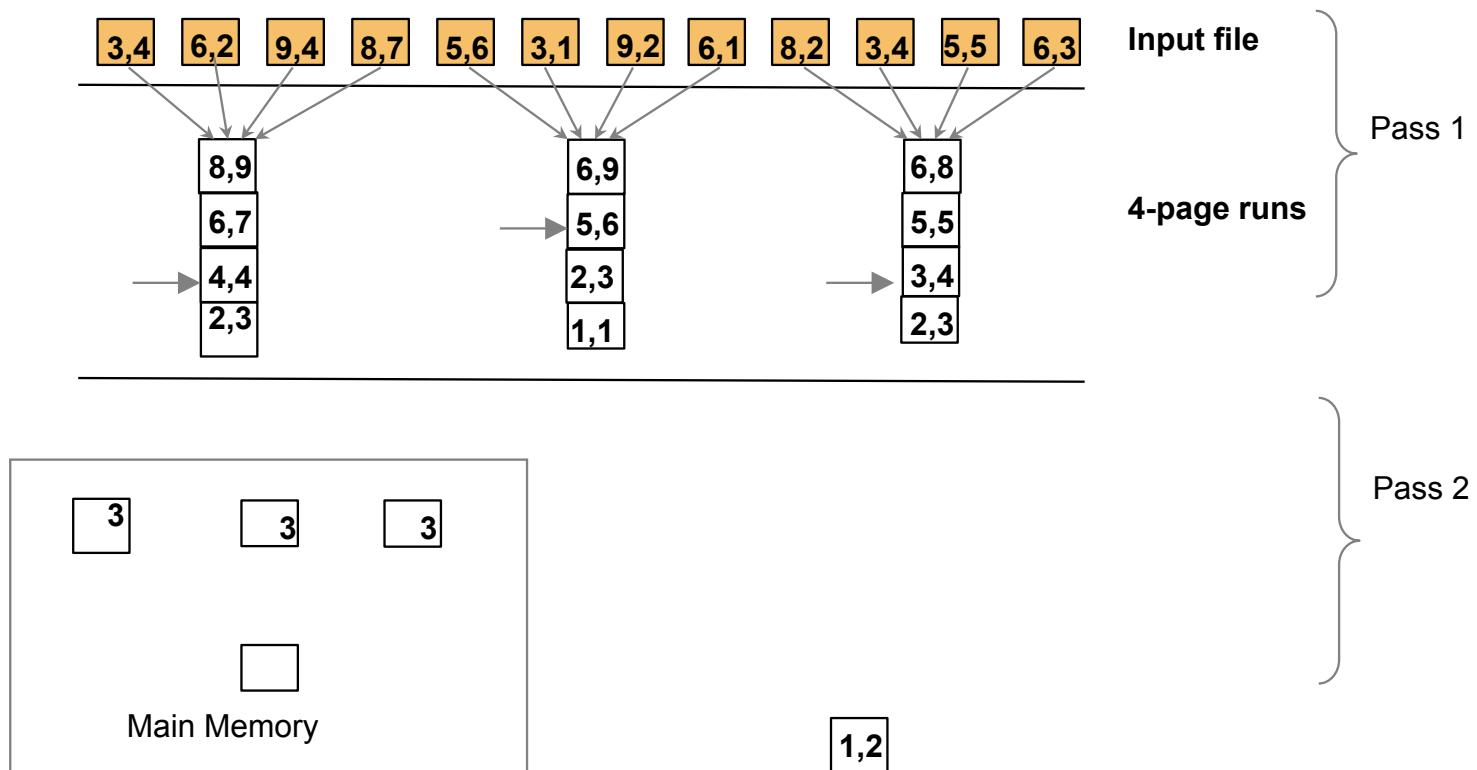
buffer pages in memory $B = 4$, each page 2 records



External Merge Sort: Example

MELBOURNE

buffer pages in memory $B = 4$, each page 2 records



MELBOURNE

- Sorting with **external sort**:
 - 1. Scan R, extract only the needed attributes
 - 2. Sort the result set using EXTERNAL SORT
 - 3. Remove adjacent duplicates

Cost =

ReadTable +	Read the entire table and keep only projected attributes
WriteProjectedPages +	Write pages with projected attributes to disk
SortingCost +	Sort pages with projected attributes with external sort
ReadProjectedPages	Read sorted projected pages to discard adjacent duplicates

WriteProjectedPages = NPages(R) * PF

how many columns we keep relative to # columns

PF: Projection Factor says how much are we projecting, ratio with respect to all attributes (e.g. keeping $\frac{1}{4}$ of attributes, or 10% of all attributes)

Every time we read and write

SortingCost = 2*NumPasses*ReadProjectedPages 

- **Example:** Let's say that we project $\frac{1}{4}$ of all attributes, and let's say that we have 20 pages in memory
- $PF = 1/4 = 0.25$, $NPages(R) = 1000$
- With 20 memory pages we can sort in 2 passes

$$(1000 + 1000 \times 0.25 + 2 \times 2 \times 250) \times 25$$

Cost = ReadTable +
WriteProjectedPages +
SortingCost +
ReadProjectedPages
 $= 1000 + 0.25 * 1000 + 2 * 2 * 250 + 250 = 2500$ (I/O)

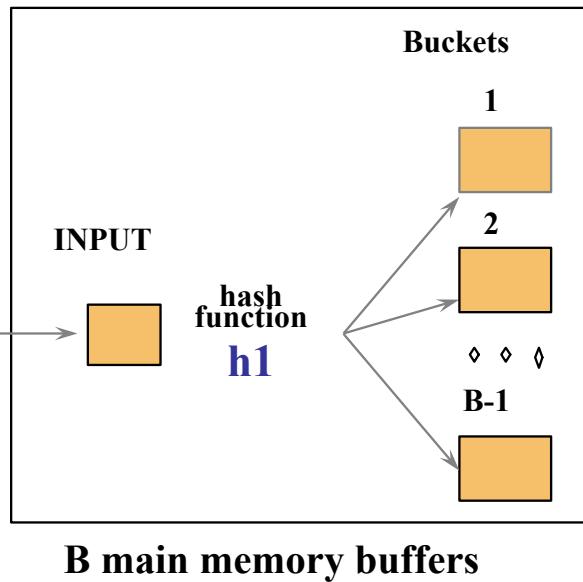
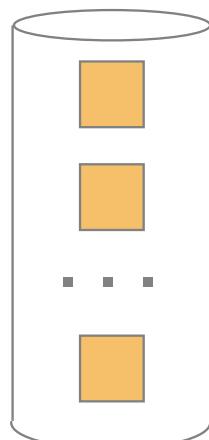


MELBOURNE

- Hashing-based projection
 - 1. Scan R, extract only the **needed** attributes
 - 2. Hash data into buckets
 - Apply hash function $h1$ to choose one of B output buffers
 - 3. Remove **adjacent** duplicates from a bucket
 - 2 tuples from different partitions guaranteed to be distinct

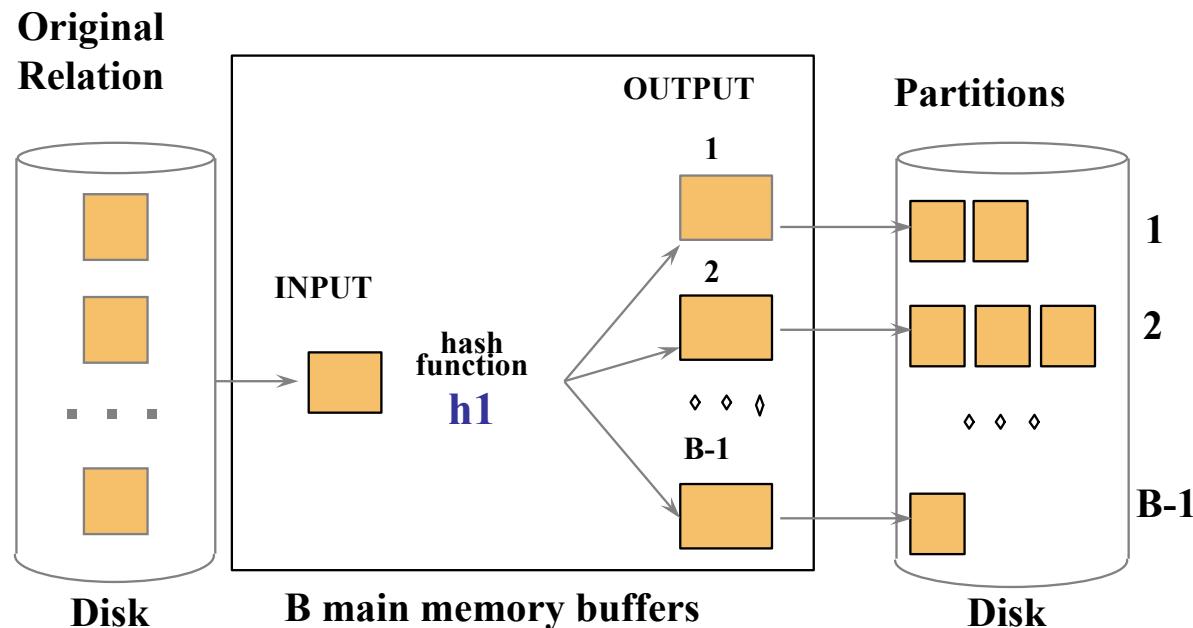


Original Relation



MELBOURNE

1. Partition data into B partitions with h_1 hash function
2. Load each partition, hash it with another hash function (h_2) and eliminate duplicates



1. Partitioning phase:

- Read R using one input buffer
- For each tuple:
 - Discard unwanted fields
 - Apply hash function $h1$ to choose one of $B-1$ output buffers
- Result is $B-1$ partitions (of tuples with no unwanted fields)
 - 2 tuples from different partitions guaranteed to be distinct

2. Duplicate elimination phase:

- For each partition
 - Read it and build an in-memory hash table
 - using hash function $h2 (<> h1)$ on all fields
 - while discarding duplicates *if same partition, copy into one bucket*
- If partition does not fit in memory
 - Apply hash-based projection algorithm recursively to this partition (we will not do this...)

Cost of External Hashing

Cost = ReadTable + Read the entire table and project attributes
WriteProjectedPages + Write projected pages into corresponding partitions
ReadProjectedPages Read partitions one by one, create another hash table and discard duplicates within a bucket

Our example:

Cost = ReadTable +
WriteProjectedPages +
ReadProjectedPages
 $= 1000 + 0.25 * 1000 + 250 = 1500 \text{ (I/O)}$

- Understand the logic behind relational operators
- Learn alternatives for selections and projections (for now)
 - Be able to calculate the cost of alternatives
- Important for Assignment 3 as well

MELBOURNE

- Query Processing Part II
 - Join alternatives