



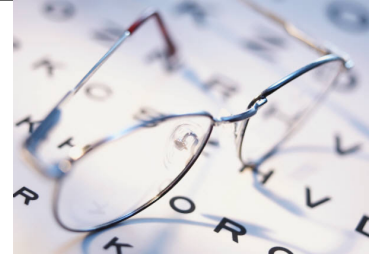
UNIVERSITY OF COPENHAGEN



# Data Processing: External Sorting using Index, Implementation of Relational Operators

ACS, Yongluan Zhou

# What should we learn today?



- External Sorting using Index
- Implementation of relational operators,
  - including selections, projections, joins, set operations, and aggregation
- Loop-based implementations
  - and techniques such as use of blocks and indices to improve their performance
- Hashing- and sorting-based implementation
  - durability
- Operator Interface

## Can we speed up sorting with Index?

- **B+ Trees:** What is a B+ Tree? What is the difference between a B+ Tree and a binary tree?

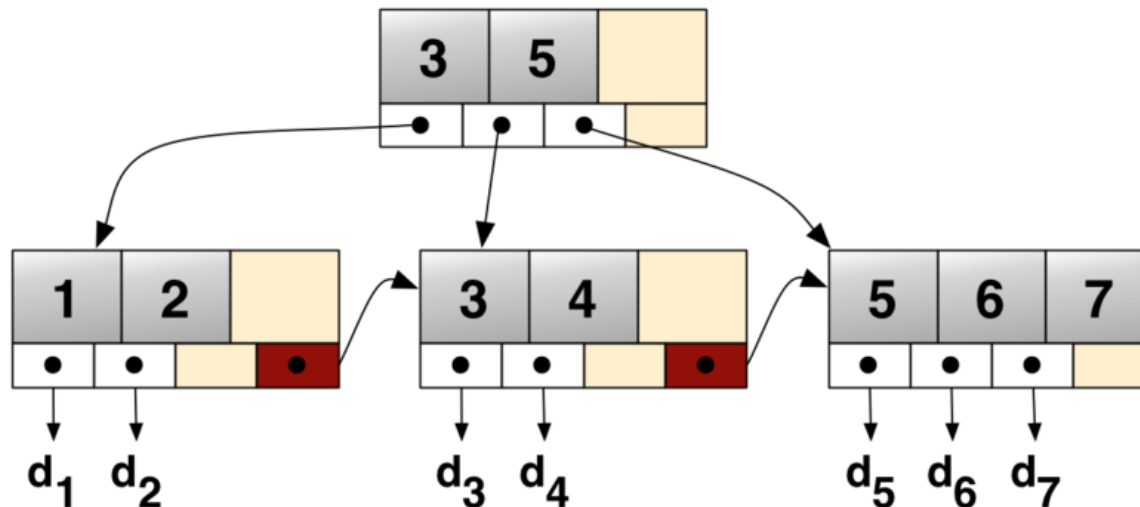
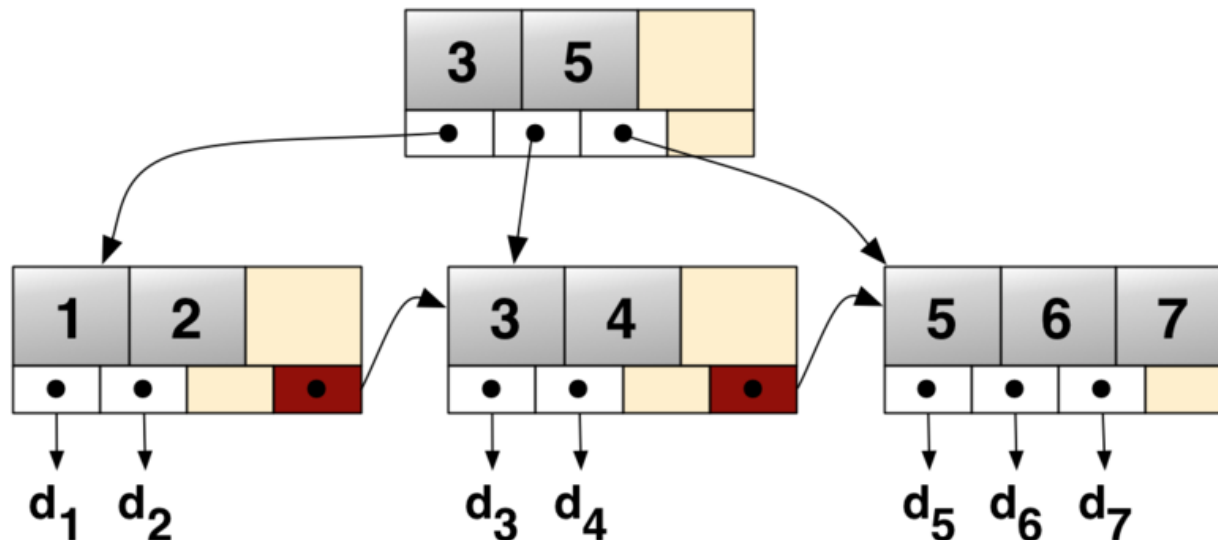


Image source: Wikipedia



## Basic concepts: B+ trees

- Each **node** in the tree occupies a **page**
- Entries in non-leave nodes → called **index entries**:  
**<key value, page\_id>**
- Entries in leaf nodes → called **data entries**:
  - either containing actual data (direct index)
  - or pointer to them (indirect index)



# Using B+ Trees for Sorting

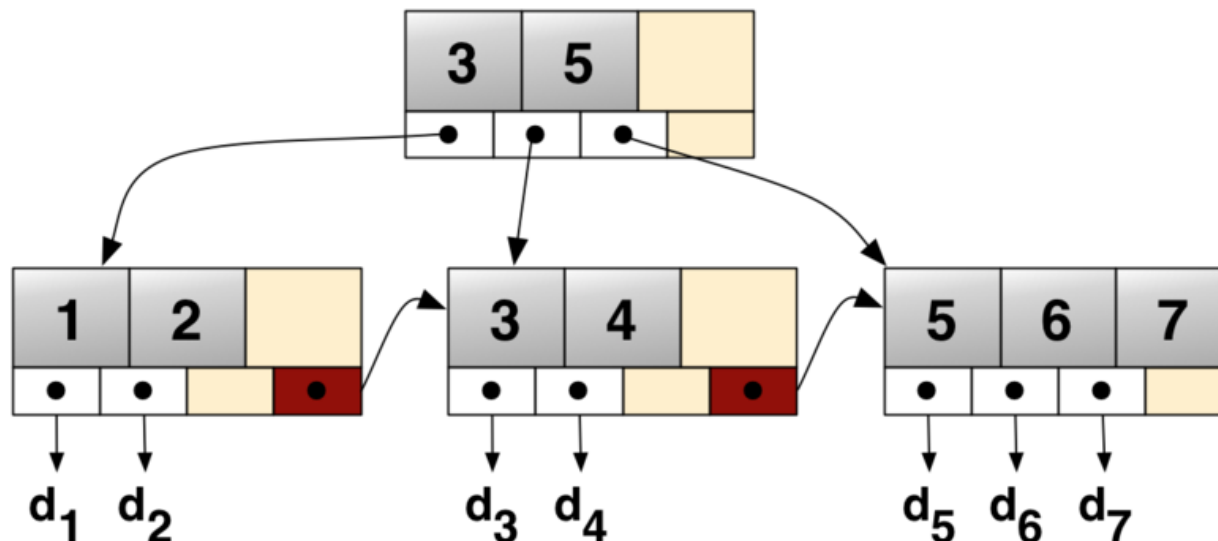
## Scenario:

- Table to be sorted has B+ tree index on sorting column(s).

## Idea:

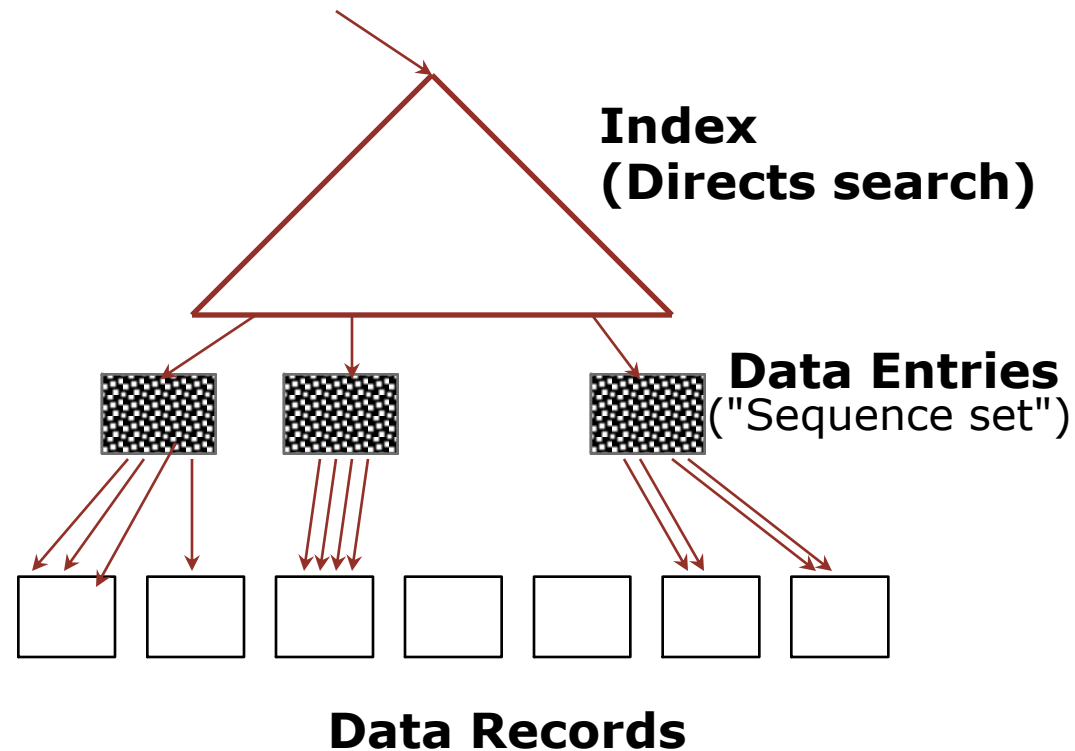
- Can retrieve records in order by traversing leaf pages.

Is it always a good idea?



## Using B+ Trees for Sorting

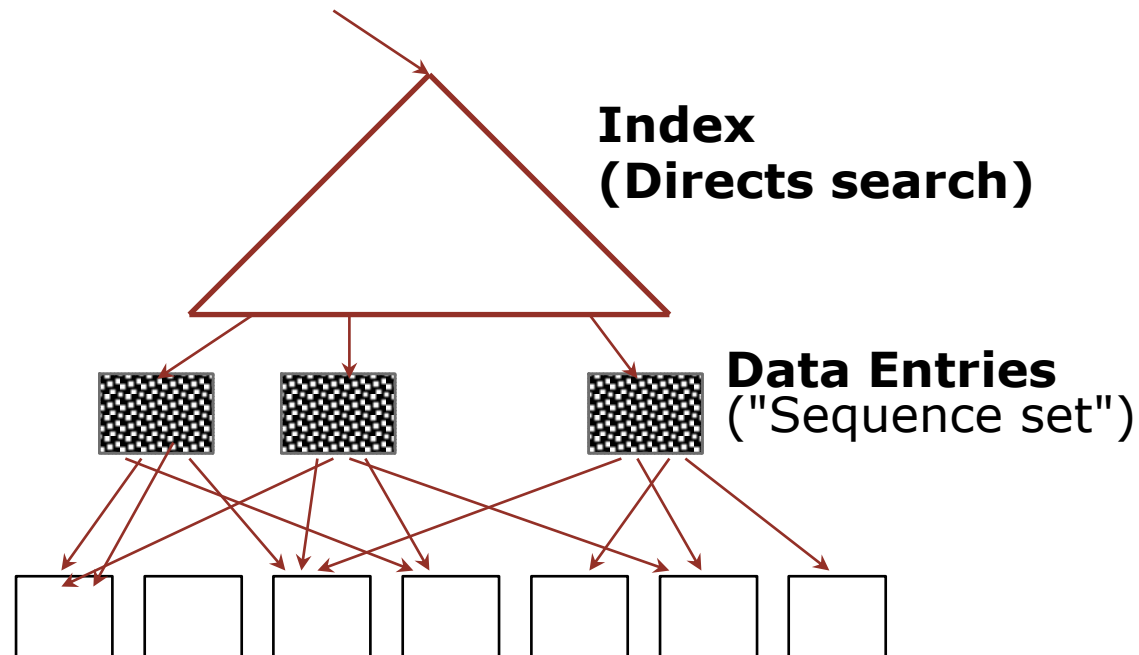
- **Cost:** root to the left-most leaf, then retrieve all leaf pages (direct index)
- If indirect index is used?  
Additional cost of retrieving data records:  
each page fetched just once.



➡ *Always better than external sorting!*

# Using B+ Trees for Sorting

- Unclustered (and indirect) index:
  - Each data entry contains pointer to a data record.
  - Data records are not stored in the sorted order of the index
- In general, **one I/O per data record!**



# Implementation of Relational Operators



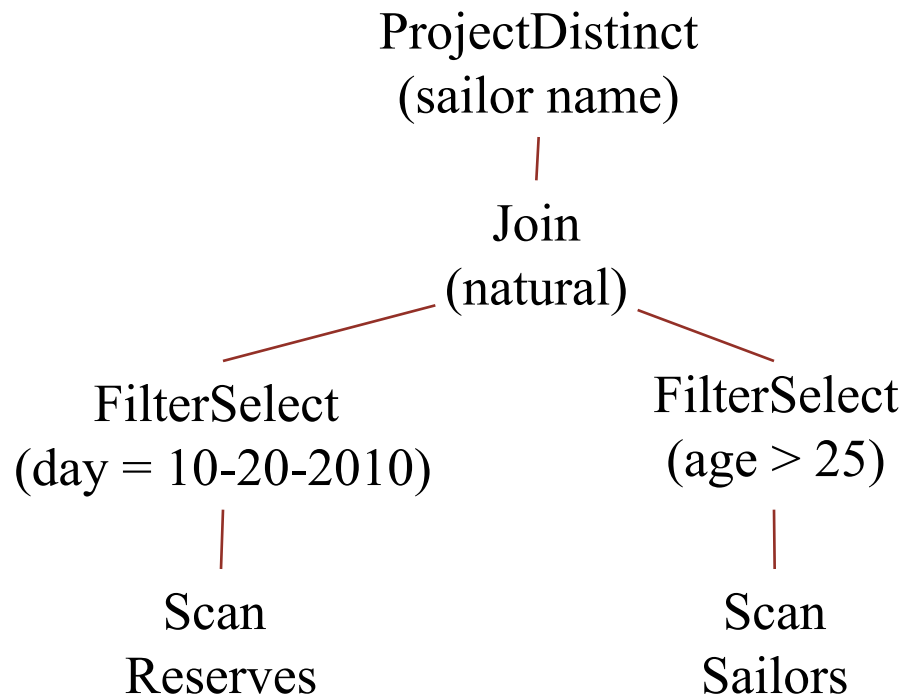


## Translating SQL Query to Relational Operators

SELECT DISTINCT S.name

FROM R JOIN S

WHERE R.day = 10-20-2010 AND S.age > 25



# Relational Operator Implementation

Implementing relational operators is challenging because:

- Relational queries are declarative
  - There is no efficient predefined strategy
- The data sets are typically very large



# Relational Operators

We now study implementation alternatives

- Select
- Project
- Join
- Set operations (union, intersect, except)
- Aggregation



## Select Operator

```
SELECT *  
FROM   Sailor S  
WHERE  S.Age = 25 AND S.Salary > 100K
```

- How best to perform? Depends on:
  - what indexes are available
  - expected size of result
- Case 1: No index on any selection attribute
- Case 2: Have “matching” index on all selection attributes
- Case 3: Have “matching” index on some (but not all) selection attributes



## Case 1: No index on any selection attribute

```
SELECT *  
FROM   Sailor S  
WHERE  S.Age = 25 AND S.Salary > 100K
```

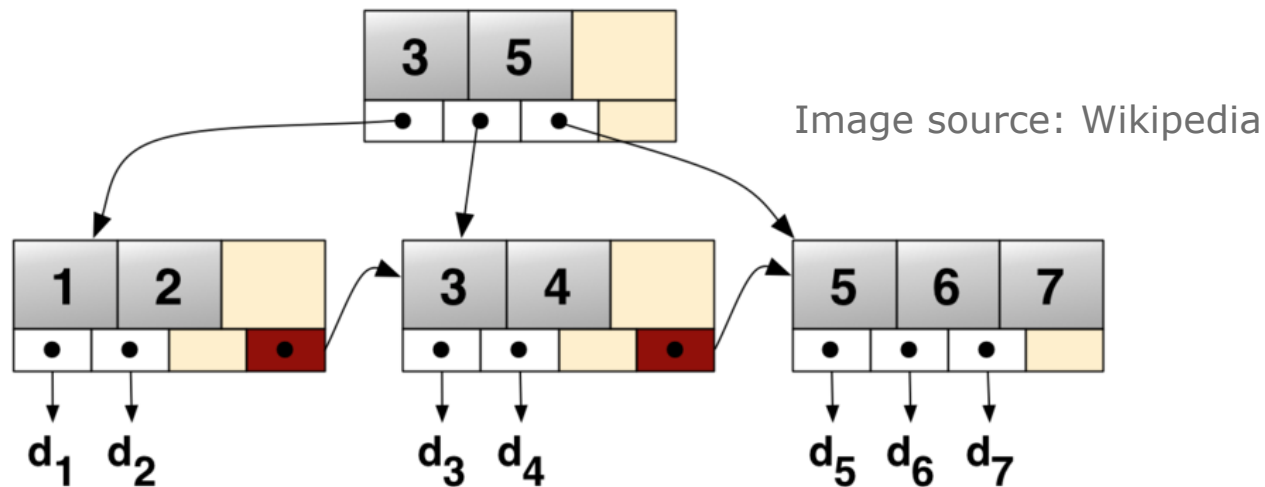
- Single loop: Just scan and filter!
- If relation has N pages, cost = N
  - Assume  $|S| = 1000$  pages, cost = 1000 pages



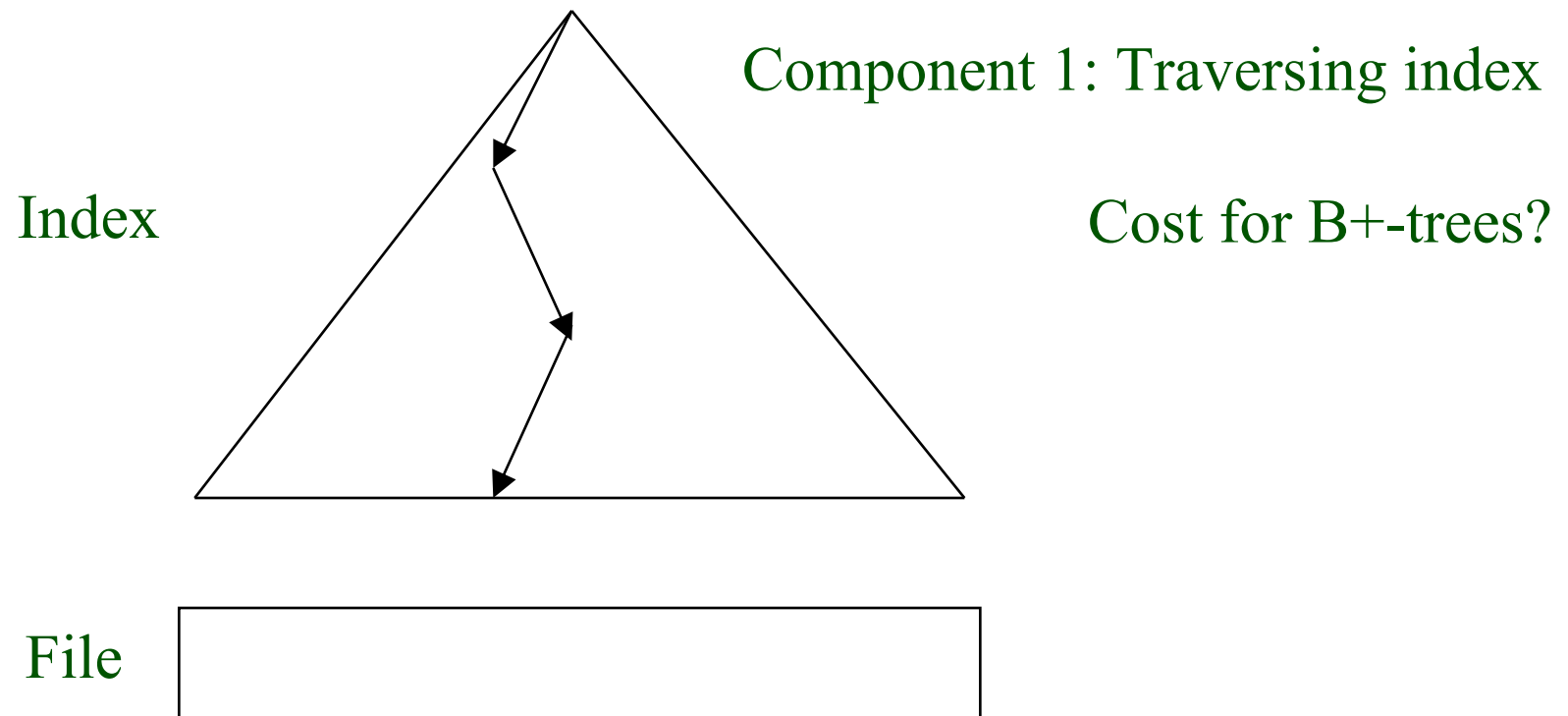
## Case 2: "Matching" index on all selection attributes

```
SELECT *  
FROM   Sailor S  
WHERE  S.Age = 25 AND S.Salary > 100K
```

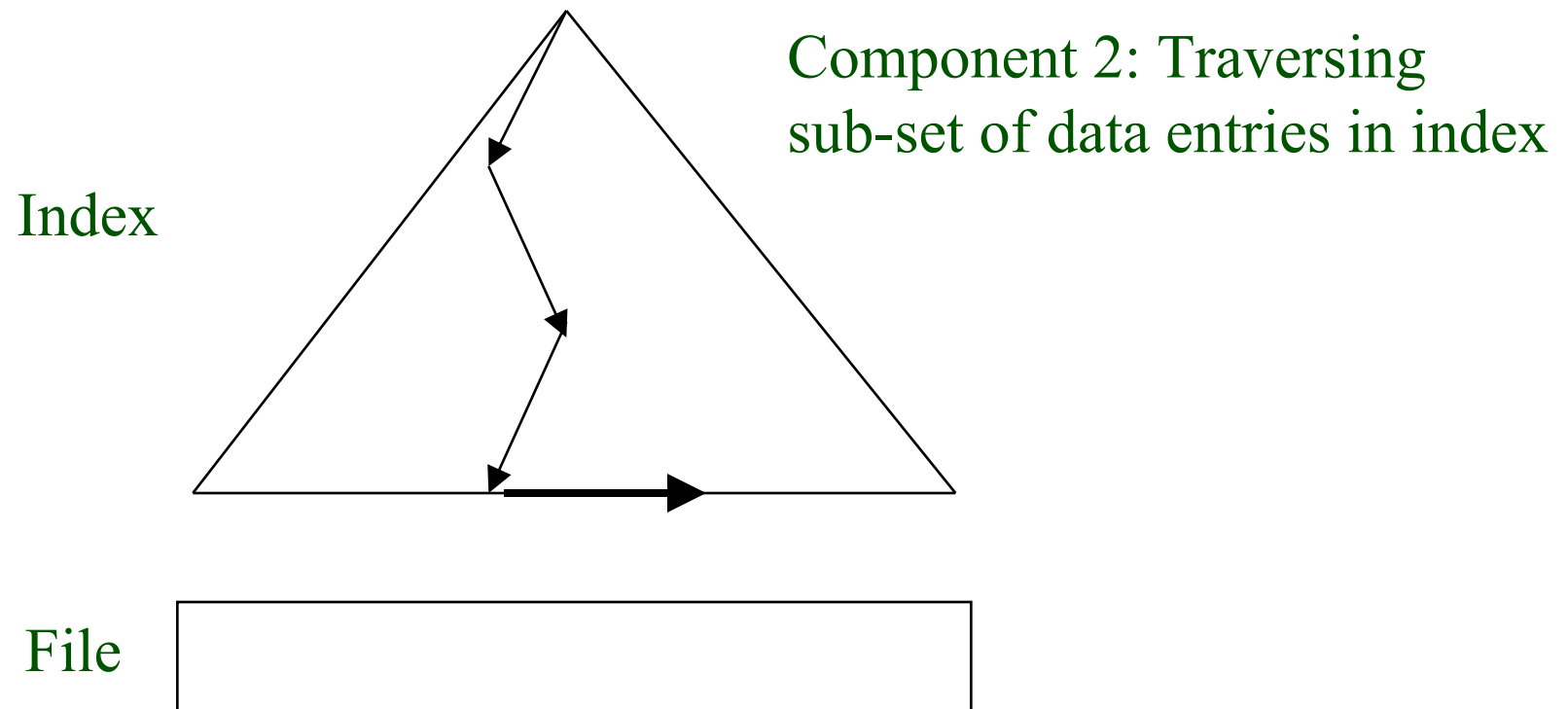
- Assume index on (Age, Salary)



## Case 2: Cost Components

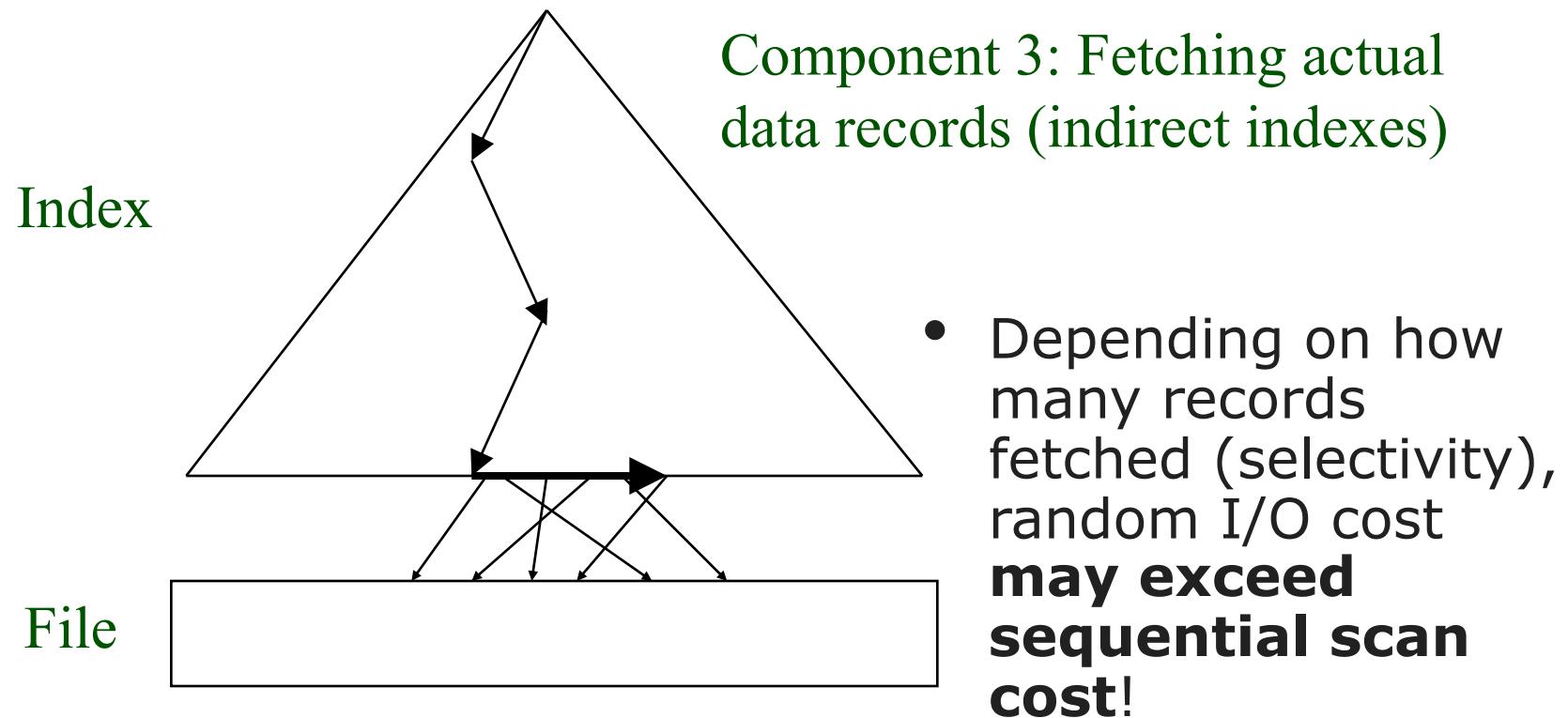


## Case 2: Cost Components



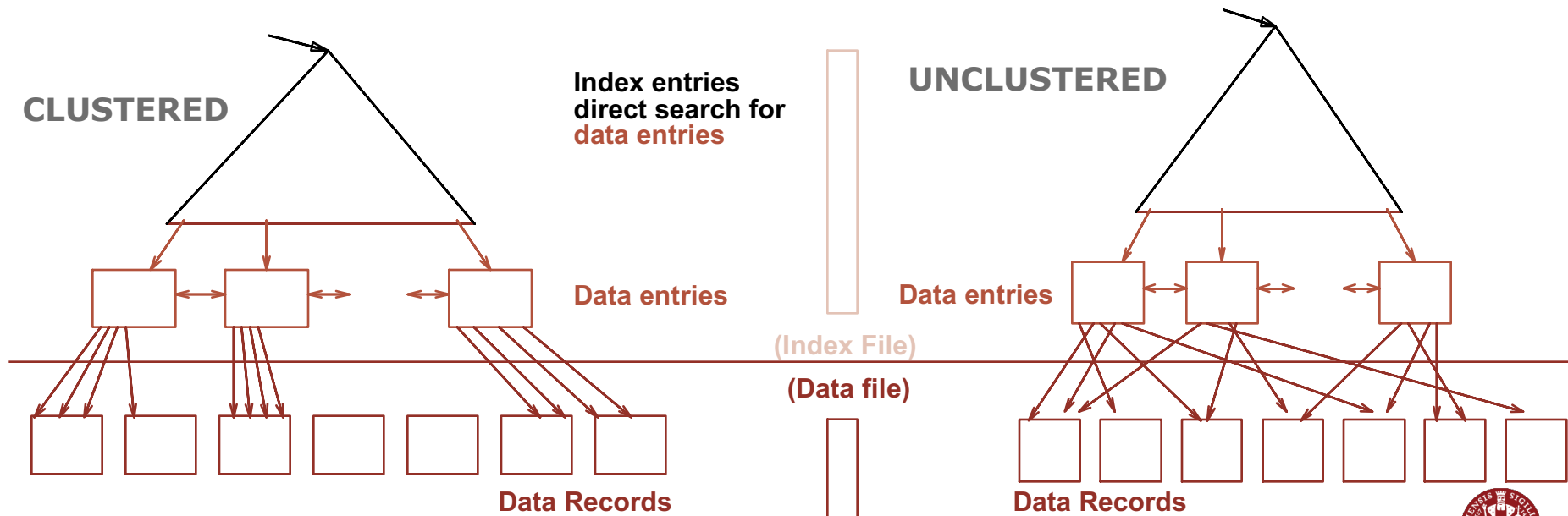


## Case 2: Cost Components



## Case 2: Cost Component 3

- Assume selectivity = 10% (100 pages, 10000 tuples):
  - If clustered index, cost is 100 I/Os;
  - If unclustered, could be up to 10000 I/Os!



## Case 3: “Matching” index on some attributes

```
SELECT *  
FROM   Sailor S  
WHERE  S.Age = 25 AND S.Salary > 100K
```

- Assume index on Age only



## Case 3: Evaluation Alternatives

- Alternative 1
  - Use available index (on Age) to get superset of relevant data entries
  - Retrieve the tuples corresponding to the set of data entries
  - Apply remaining predicates on retrieved tuples
  - Return those tuples that satisfy all predicates
- Alternative 2
  - Sequential scan! (always available)
  - May be again better depending on selectivity



## Case 3: “Matching” index on some attributes

```
SELECT *  
FROM   Sailor S  
WHERE  S.Age = 25 AND S.Salary > 100K
```

- Assume separate indices on Age and on Salary



## Case 3: Evaluation Alternatives

- Alternative 1
  - Choose most **selective** access path (index)
    - Could be index on Age or Salary, depending on selectivity of the corresponding predicates
  - Use this index to get **superset** of relevant data entries
  - Retrieve the tuples corresponding to the set
  - Apply remaining predicates on retrieved tuples
  - Return those tuples that satisfy all predicates
- Alternative 2
  - Get rids of data records using each index
    - Use index on Age and index on Salary
  - **Intersect** the rids
  - Retrieve the tuples corresponding to the rids
  - Apply remaining predicates on retrieved tuples
  - Return those tuples that satisfy all predicates
- Alternative 3
  - Sequential scan!



# Questions?



# Relational Operators

- We now study implementation alternatives
- Select
- Project
- Join
- Set operations (union, intersect, except)
- Aggregation





## Projection

```
SELECT DISTINCT S.Name, S.Age  
FROM   Sailor S
```

Main issue is duplicate elimination.

- Assume we do *not* have any indices

How would you implement  
duplicate elimination?  
Can you use sorting or hashing?

## Projection without Indices

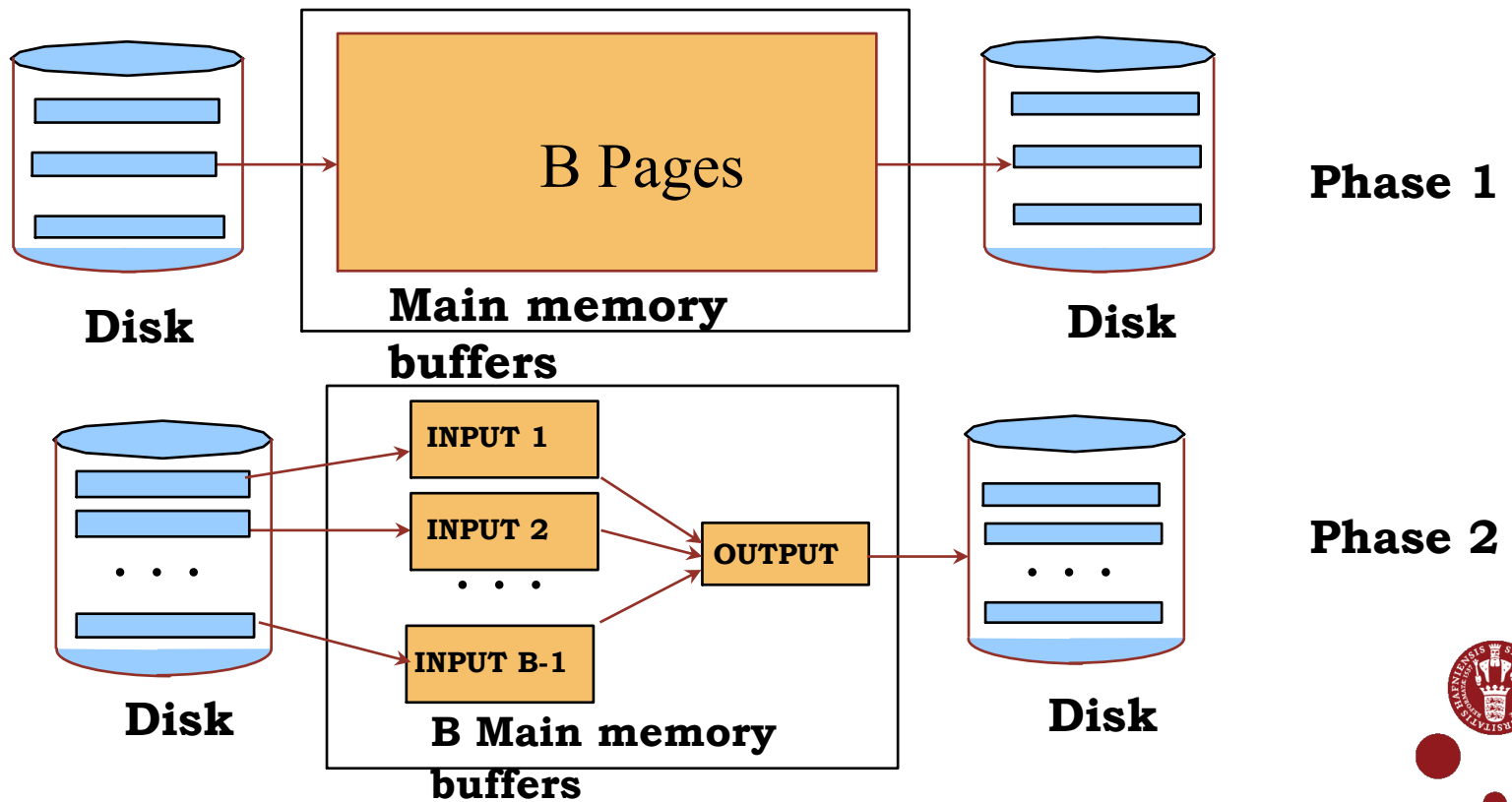
```
SELECT DISTINCT S.Name, S.Age  
FROM   Sailor S
```

- We have *no* indices
- What strategies can we use?
  - **Sorting:** Duplicates adjacent after sorting
  - **Hashing:** Duplicates hash to same buckets (in disk and memory)



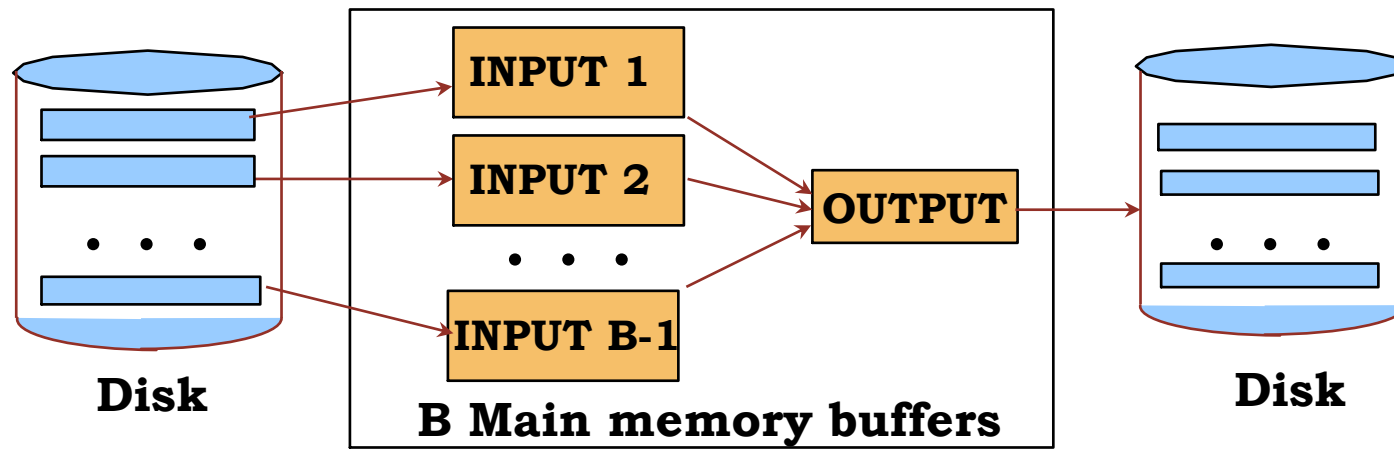
## Do-it-yourself Recap: External Sorting

- What were the two phases of multi-way external sorting and how did they work?
- What optimizations could you apply to external sorting?



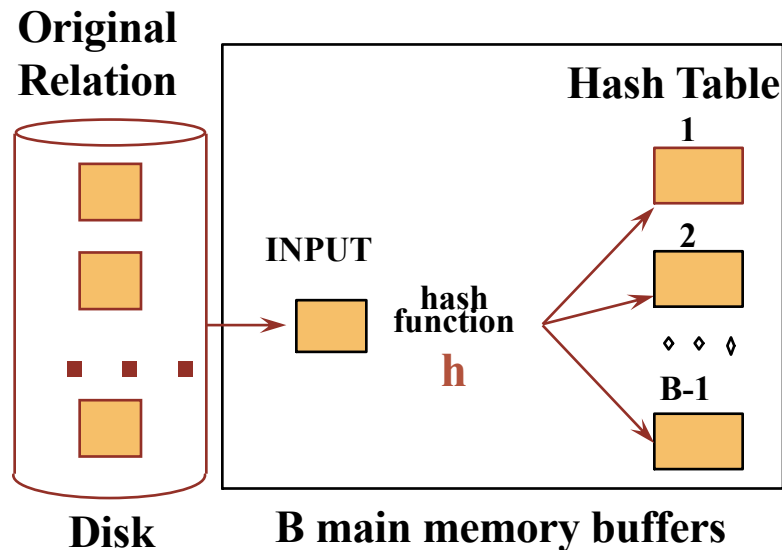
# Projection with External Sorting

- Phase 1
  - Project out unwanted columns
  - Still produce runs of length B pages
  - But tuples in runs are smaller than input tuples
- Phase 2
  - Eliminate duplicates during merge



## Duplicate Elimination with Hashing

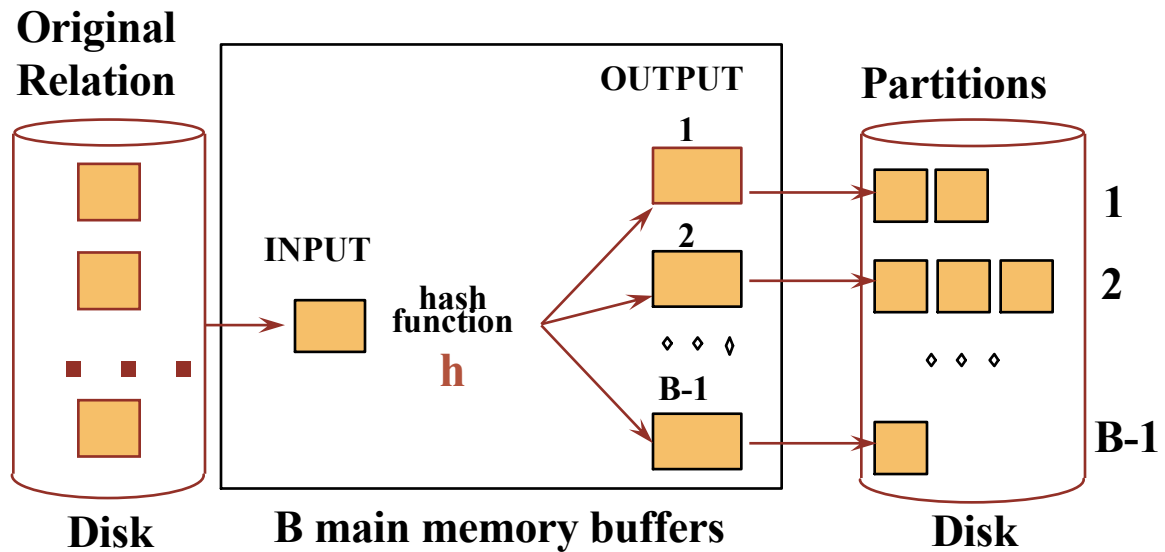
1. Apply hash function
2. Look for duplicates in the corresponding bucket
3. If the input buffer is empty, then read in a page and goto 1.



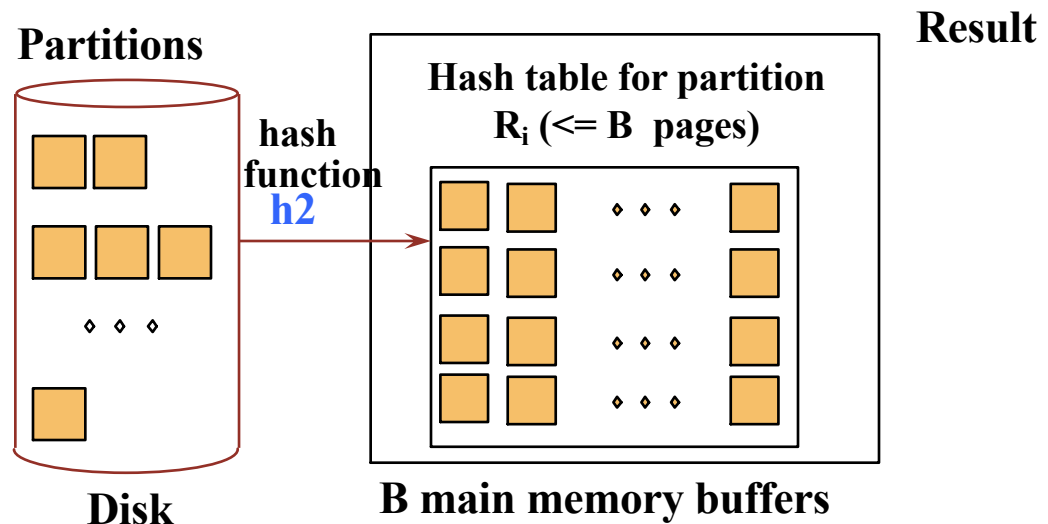
**How to build hash table for data larger than the memory?**

# Duplicate Elimination using Hashing

- **Partition:**



- **Rehash:**

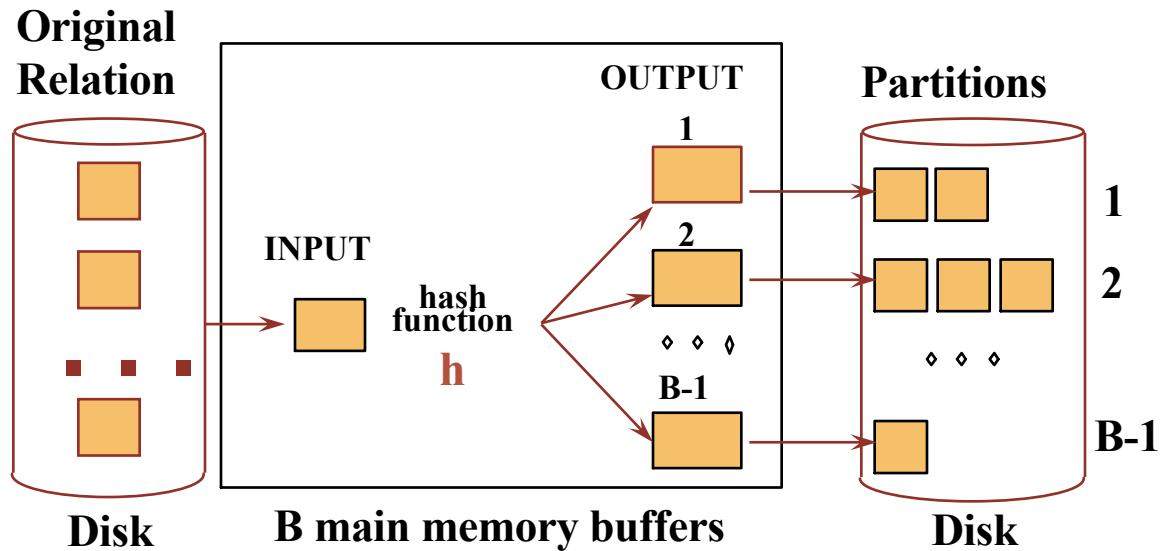


# General Idea

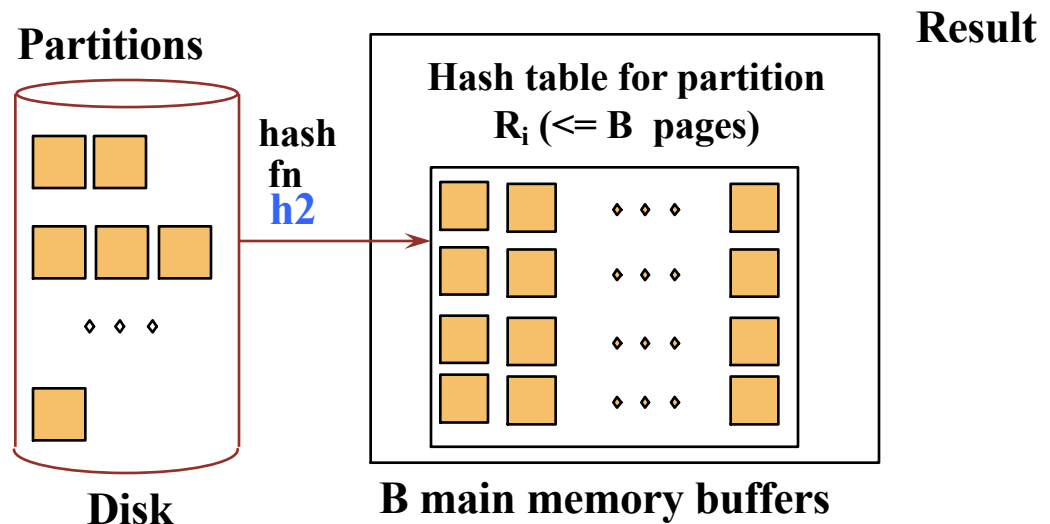
- Two phases:
  - **Partition**: use a hash function  $h$  to split tuples into partitions on disk.
    - Key property: all matches live in the same partition.
  - **ReHash**: for each partition on disk, build a main-memory hash table using a hash function  $h_2$



# Duplicate Elimination using Hashing



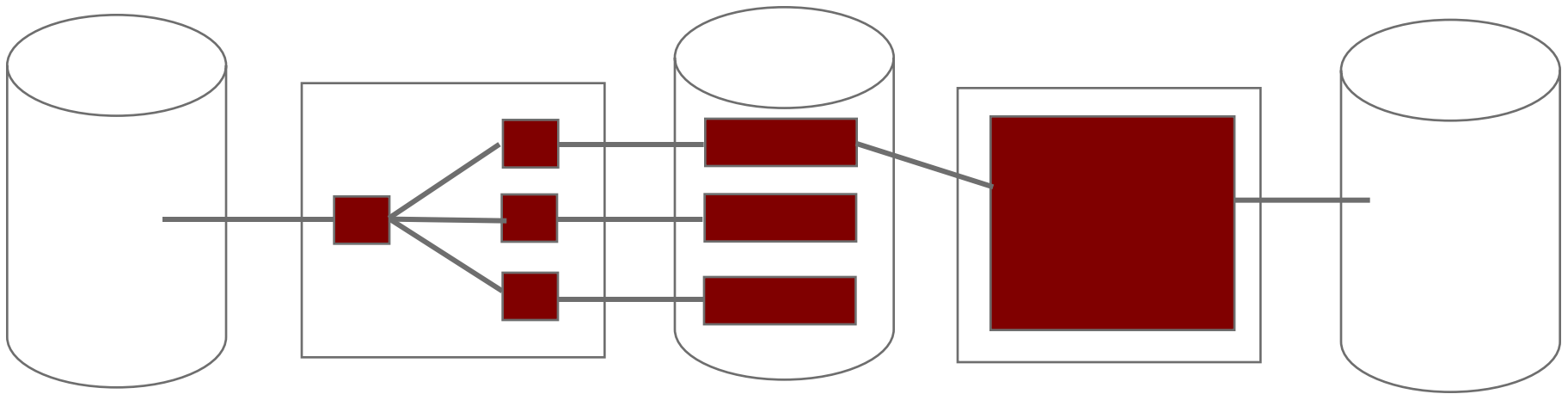
**What if a partition  
still cannot fit into  
memory?**





## Cost of External Hashing

- Suppose it can be done in 2 passes



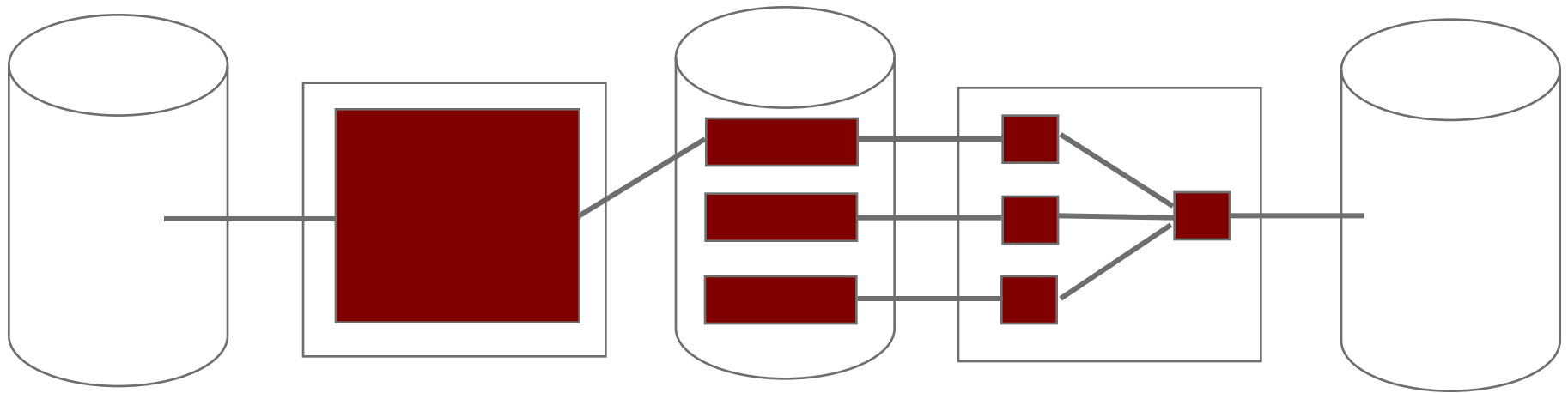
$$\text{cost} = 3 * |R| \text{ IO's}$$

How does this compare with **external sorting**?



## Cost of External Sorting

- Suppose sorting can be done in 2 passes

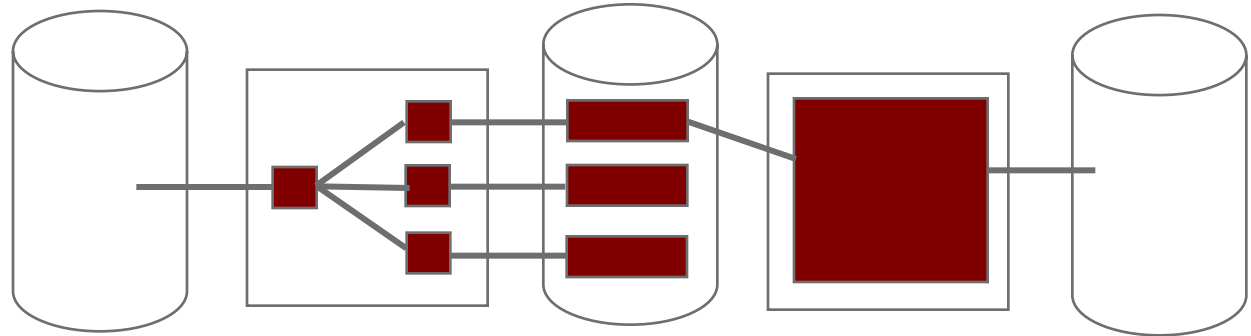


$$\text{cost} = 3 * |R| \text{ IO's}$$

# Duality of External Hashing and External Sorting

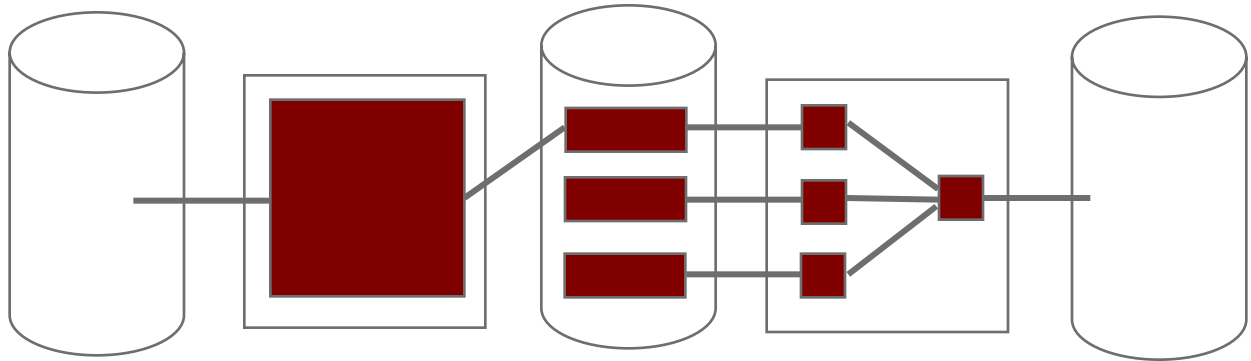
External  
Hashing

cost =  $3 * |R|$  IO's



External  
Sorting

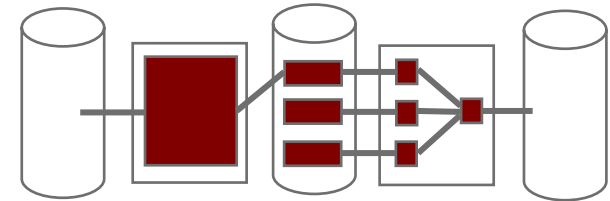
cost =  $3 * |R|$   
IO's



# Duality of Sorting and Hashing

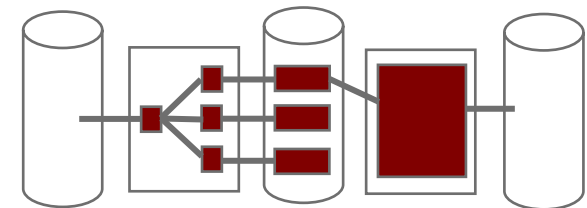
- **Sorting**

- Physical division, logical combination
  - Split followed by merge
  - Recurse on merging
- Sequential write (phase 1), random read (phase 2)
- Fan-in
- If pipelining and  $\sqrt{N} < B < N$ , Total Cost =  $3N$



- **Hashing**

- Logical division, physical combination
  - Partition followed by concatenate
  - Recurse on partitioning
- Random write (phase 1), sequential read (phase 2)
- Fan-out
- If pipelining and  $\sqrt{N} < B < N$ , Total Cost =  $3N$



Source: G. Graefe, Query Evaluation Techniques for Large Databases.  
ACM Computing Surveys, Vol 25, No. 2, June 1993 (partial)

## So which is better??

- **Sorting pros:**
  - Great if input already sorted (or *almost* sorted)
  - Great if need output to be sorted anyway
  - Not sensitive to “data skew” or “bad” hash functions
- **Hashing pros:**
  - Highly parallelizable
  - Can exploit extra memory to reduce # IOs with hybrid hashing (*will not covered in this course*)



# Relational Operators

- We now study implementation alternatives
- Select
- Project
- Join (equi-joins only)
- Set operations (union, intersect, except)
- Aggregation



# Joins

- Joins are very common.

```
SELECT *  
FROM   Reserves R1, Sailors S1  
WHERE  R1.sid=S1.sid
```

- Join techniques we will cover today:
  1. Nested-loops join
  2. Index-nested loops join





## Simple Nested Loops Join

$R \bowtie S$ :

```
foreach tuple r in R do
    foreach tuple s in S do
        if r.sid == s.sid then add <r, s>
to result
```

- Suppose R has 1000 pages, S has 500 pages, and each page (of R and S) has 100 tuples
- Cost =  $|R| + (pR * |R|) * |S| = 1000 + 100 * 1000 * 500$  IOs
  - At 10ms/IO, Total time: ???
    - ~ 6 days!
- What if smaller relation (S) was "outer"?
- What if one relation can fit entirely in memory?



## Page-Oriented Nested Loops Join

$R \bowtie S$ :

```
foreach page  $b_R$  in  $R$  do
  foreach page  $b_S$  in  $S$  do
    foreach tuple  $r$  in  $b_R$  do
      foreach tuple  $s$  in  $b_S$  do
        if  $r_i == s_j$  then add  $\langle r, s \rangle$  to result
```

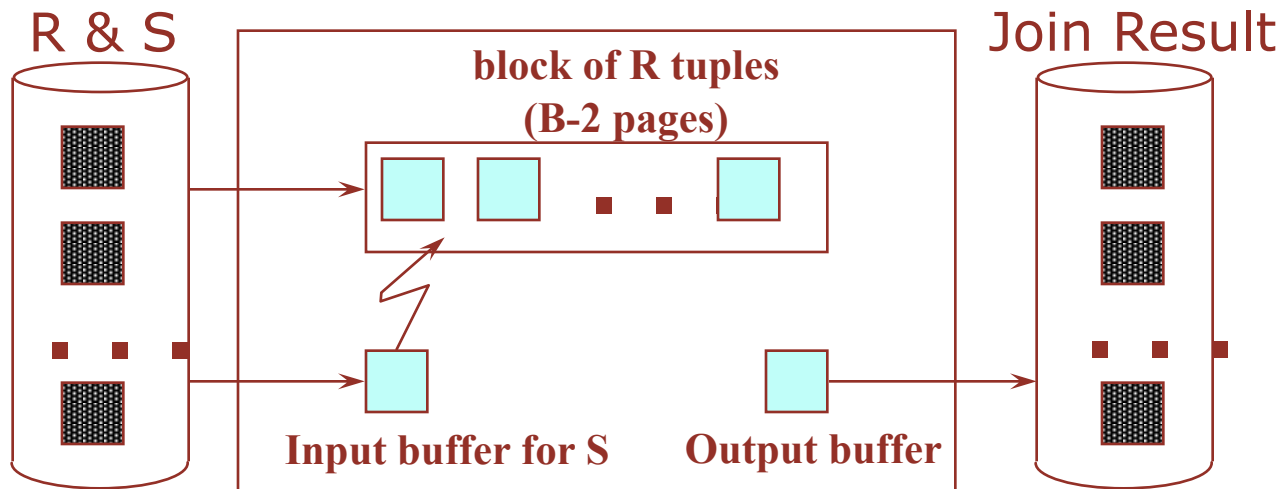
- Cost =  $|R| * |S| + |R| = 1000 * 500 + 1000$
- If smaller relation ( $S$ ) is outer, cost =  $500 * 1000 + 500$
- Much better than naïve per-tuple approach!
  - At 10ms/IO, total time  $\sim 1.4$  hour
- The trick is to reduce the # complete reads of the inner table

Can we reduce it  
even further?



## Block Nested Loops Join

- Page-oriented NL doesn't exploit extra buffers :(
- Idea to use memory efficiently:



**Cost: Scan outer + (#outer blocks \* scan inner)**

$\#outer\ blocks = \lceil \# of\ pages\ of\ outer / blocksize \rceil$



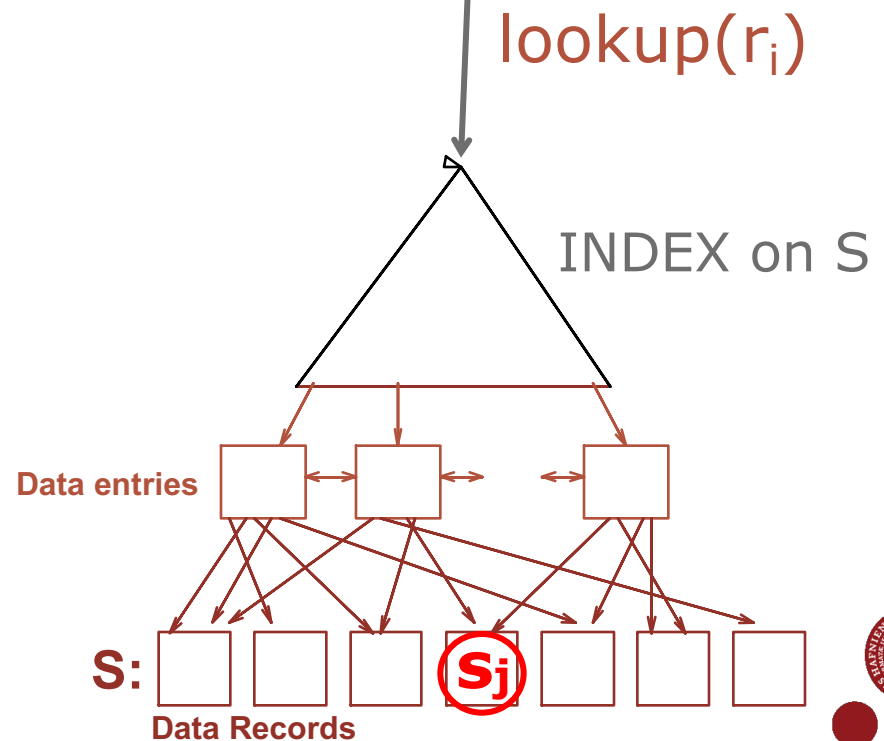
## Examples of Block Nested Loops Join

- Say we have  $B = 100+2$  memory buffers
- Join cost =  $|outer| + (\#outer\ blocks * |inner|)$ 
  - $\#outer\ blocks = |outer| / 100$
- With R as outer ( $|R| = 1000$ ):
  - Scanning R costs 1000 IO's (done in 10 blocks)
  - Per block of R, we scan S; costs  $10 * 500$  I/Os
  - Total =  $1000 + 10 * 500$ .
  - At 10ms/IO, total time:  $\sim 1$  minute
- With S as outer ( $|S| = 500$ ):
  - Scanning S costs 500 IO's (done in 5 blocks)
  - Per block of S, we scan R; costs  $5 * 1000$  IO's
  - Total =  $500 + 5 * 1000$ .
  - At 10ms/IO, total time:  $\sim 55$  seconds



# Index Nested Loops Join

$R \bowtie S$ :  
 foreach tuple  $r$  in  $R$  do  
~~foreach tuple  $s$  in  $S$  where  $r_i == s_j$  do~~  
 add  $\langle r, s \rangle$  to result



## What should we learn today?



- Discuss the design of a **pull-based interface** for data processing operators, and how such an organization helps with composability and pipelining
- Explain and reason about the **implementation of physical relational operators**, including selections, projections, joins, set operations, and aggregation
- Explain simple **loop-based implementations** to relational operators and techniques such as use of blocks and indices to improve their performance
- Discuss the **duality of hashing and sorting** and how these algorithmic approaches apply to the implementation of relational operators