

# **07- File Organization and Indexing**

**School of Computer Science  
University of Windsor**

**Dr. Shafaq Khan**

---

# Agenda

## ➤ Lecture

- Organizing the data on the disk
- Introduction to Indexing
- Single-Level Ordered Indexes
  - Primary Indexes
  - Clustering Indexes
  - Secondary Indexes
- Multilevel Indexes
  - Two-Level Primary Indexing
- Dynamic Multilevel Indexes
  - B Trees
  - B+ Trees
- CREATE INDEX Command in SQL

## ➤ Assignment-2 Quiz

# Announcements

## **P2: Milestone report submission**

(Sec 2: Feb 27; Sec 3: Feb 28; Sec 4: Mar 1)

**P2: Milestone presentation: Next week (Feb 28 – Mar 2)**



# Introductory Questions

Why is indexing important?

What is the aim of having an index?

What are the downsides of indexes?

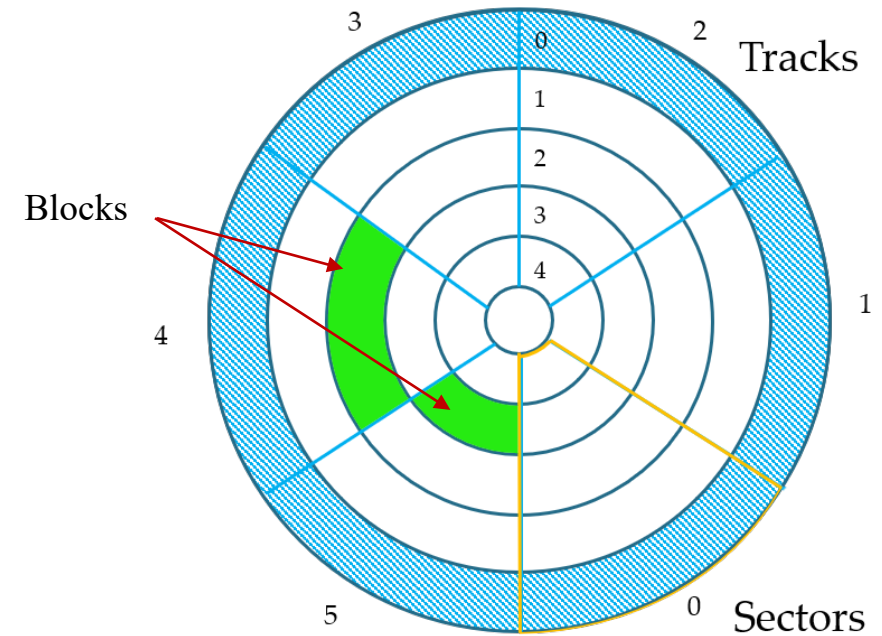
What are the different types of indexes used for query optimizations?

What to choose for creating indexes?

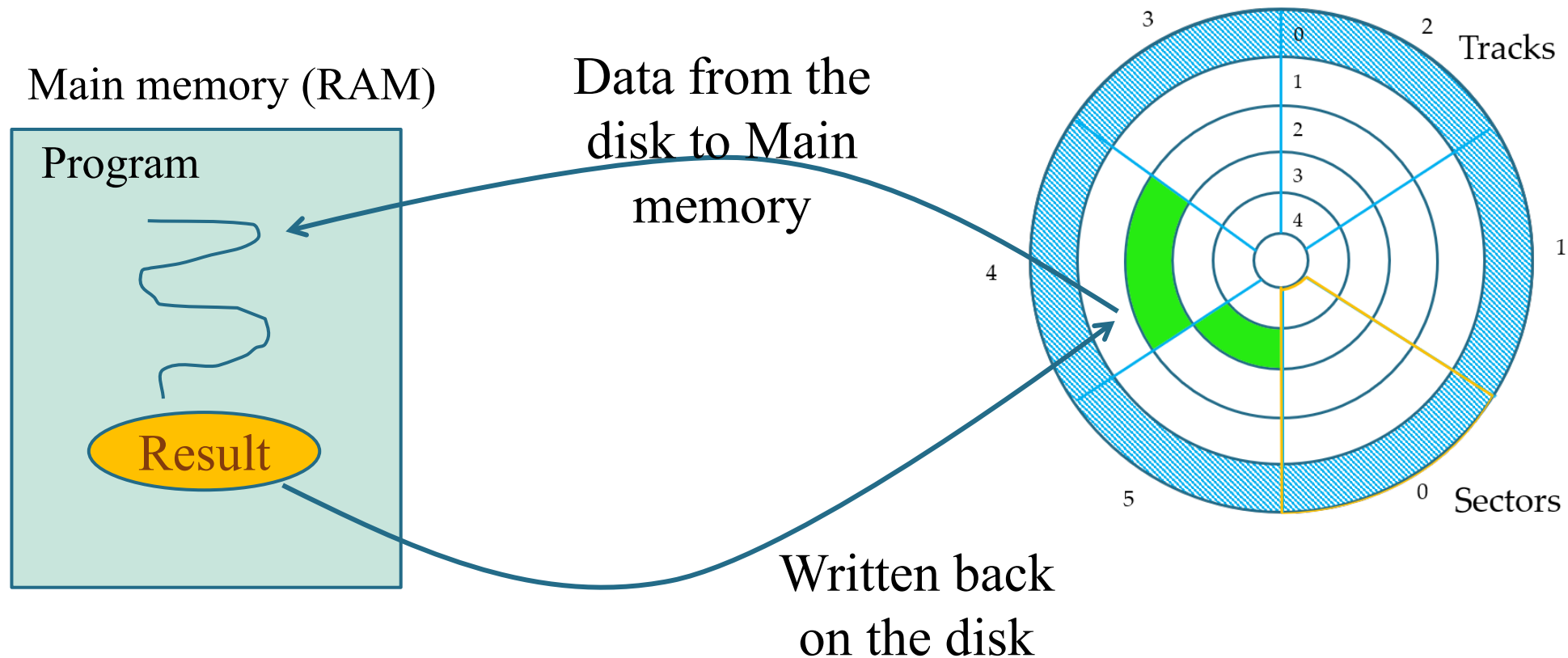
# Disk Structure

*Block Address*  
 $= \langle \text{Track Num}, \text{Sector Num} \rangle$

Let's assume a block size is 512 bytes



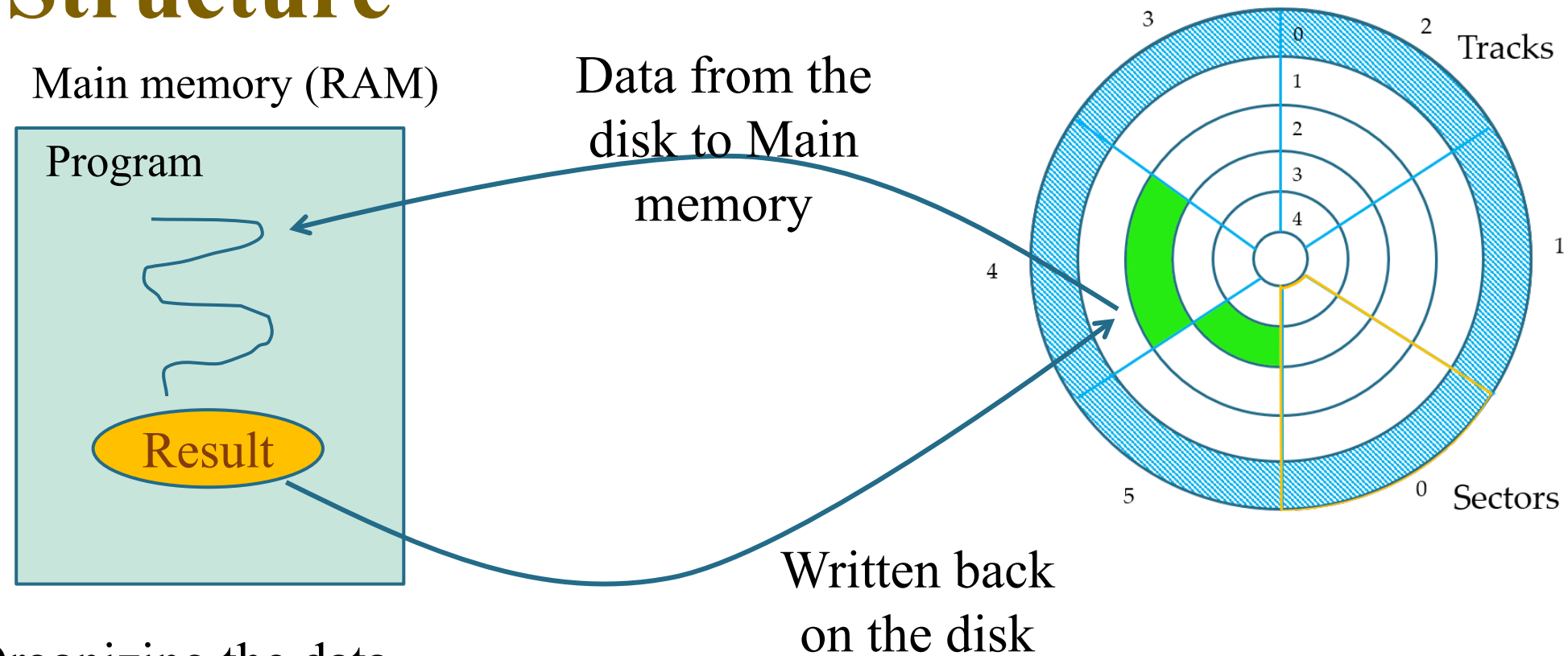
# Disk Structure



So, the data cannot be directly processed upon the disk it has to be brought into the main memory and then access .



# Disk Structure



Organizing the data inside the main memory that is directly used by the program is **Data Structures**.

Organizing the data on the disk efficiently so that it can be easily utilized that is **DBMS**.

7



# How is data stored on Disk?

Employees

Fields	Size
Eid	10
Name	50
Depart	30
Gender	8
Salary	30
Total size	128 bytes

In each block how many rows can be stored?

Number of Records/ Block = 512/128 = 4

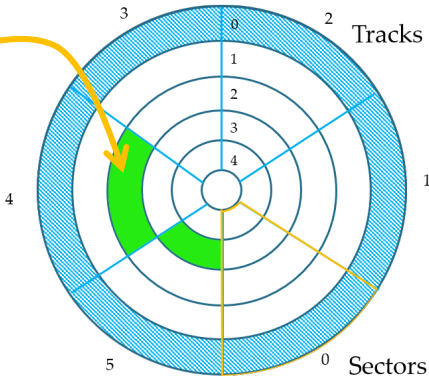
128 bytes

Employees

Eid	Name	Depart
1	Jenkins	Manager
2	Williams	Sales Rep
3	Smith	Sales Rep
4	Crosby	Manager
5	Albright	Secretary
6	Sawyer	Sales Rep
7	Thomas	Secretary
8	Albright	Worker
9	Crawford	Manager

100 records

What is the size of a record?



Block Size = 512 bytes





# How is data stored on Disk?

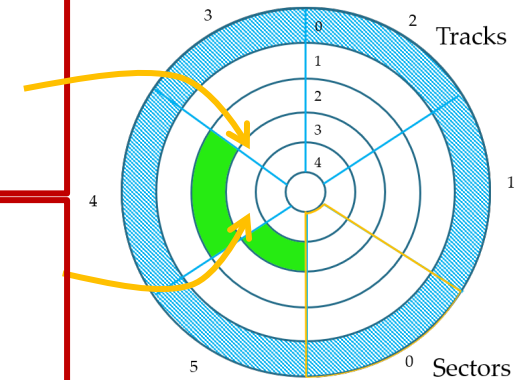
Employees

Fields	Size
Eid	10
Name	50
Depart	30
Gender	8
Salary	30
Total size	<b>128 bytes</b>

**SELECT** Eid=7  
**FROM** Employees

Employees

Eid	Name	Depart
1	Jenkins	Manager
2	Williams	Sales Rep
3	Smith	Sales Rep
4	Crosby	Manager
5	Albright	Secretary
6	Sawyer	Sales Rep
7	Thomas	Secretary
8	Albright	Worker
9	Crawford	Manager



Block Size = 512 bytes

$$\text{Number of Records/ Block} = 512/128 \\ = 4$$

How many blocks  
are required for  
100 records?

$$100 \text{ Records can be stored in} = 100/4 \text{ or} \\ 12800/512 \\ = \mathbf{25 \text{ blocks}}$$

100 records



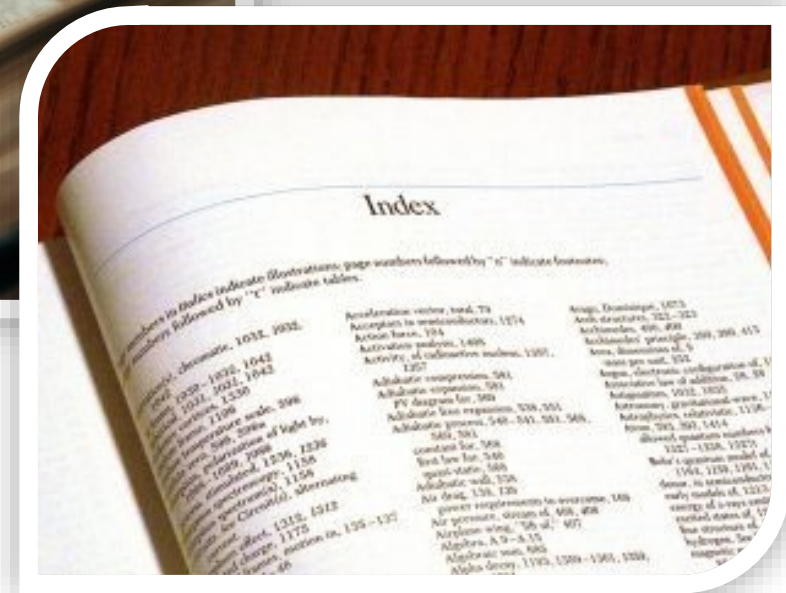
# What is Indexing?

An index is a structure that provides accelerated access to the rows of a table based on the values of one or more columns.



It makes our search simpler and quicker.

- ✓ How do we create the indexes?
- ✓ How these indexes help to access the data?



# Motivation- Searching for a Key Value

- Ex: Unsorted List
  - $1024 \times 1024 \times 1024 = 1,073,741,824$  (I Billion elements –Approx.)
  - Takes 1 Billion comparisons in the worst case
- Ex|: Sorted List
  - $1024 \times 1024 \times 1024 = 1,073,741,824$  (I Billion sorted elements –Approx.)
  - Takes 30 comparisons in the worst case
- 30 comparisons vs. 1 Billion comparisons
- Overhead involved in sorting is much lesser than that in searching?



# What is Indexing?

Where do we store the index?

How many blocks for storing index?

Index

Key	Pointer
1	
2	
3	
4	
5	

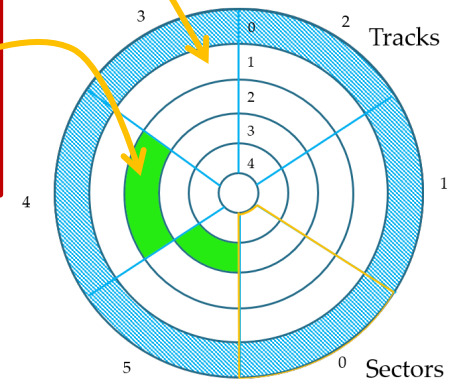
100 records

**Dense Index**

Employees

Eid	Name	Depart
1	Jenkins	Manager
2	Williams	Sales Rep
3	Smith	Sales Rep
4	Crosby	Manager
5	Albright	Secretary
6	Sawyer	Sales Rep
7	Thomas	Secretary
8	Albright	Worker
9	Crawford	Manager

100 records



Block Size = 512 bytes



# What is Indexing?

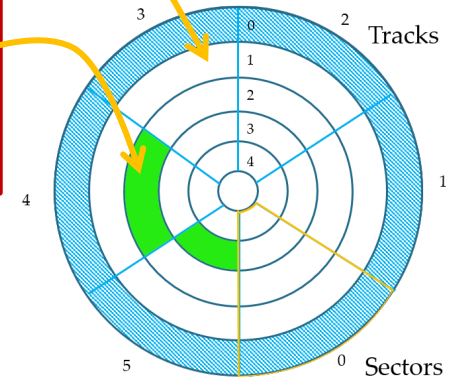
Index	
Fields	Size
Eid	10
Pointer	6
Total size	<b>16 bytes</b>

Index	
Key	Pointer
1	→
2	→
3	→
4	→
5	

100 entries

Employees		
Eid	Name	Depart
1	Jenkins	Manager
2	Williams	Sales Rep
3	Smith	Sales Rep
4	Crosby	Manager
5	Albright	Secretary
6	Sawyer	Sales Rep
7	Thomas	Secretary
8	Albright	Worker
9	Crawford	Manager

100 records



Block Size = 512 bytes

1 index entry requires : 16 bytes  
100 entries require: 1600 bytes  
100 entries can be stored in =  $1600/512 = 3.1 \approx 4$  blocks

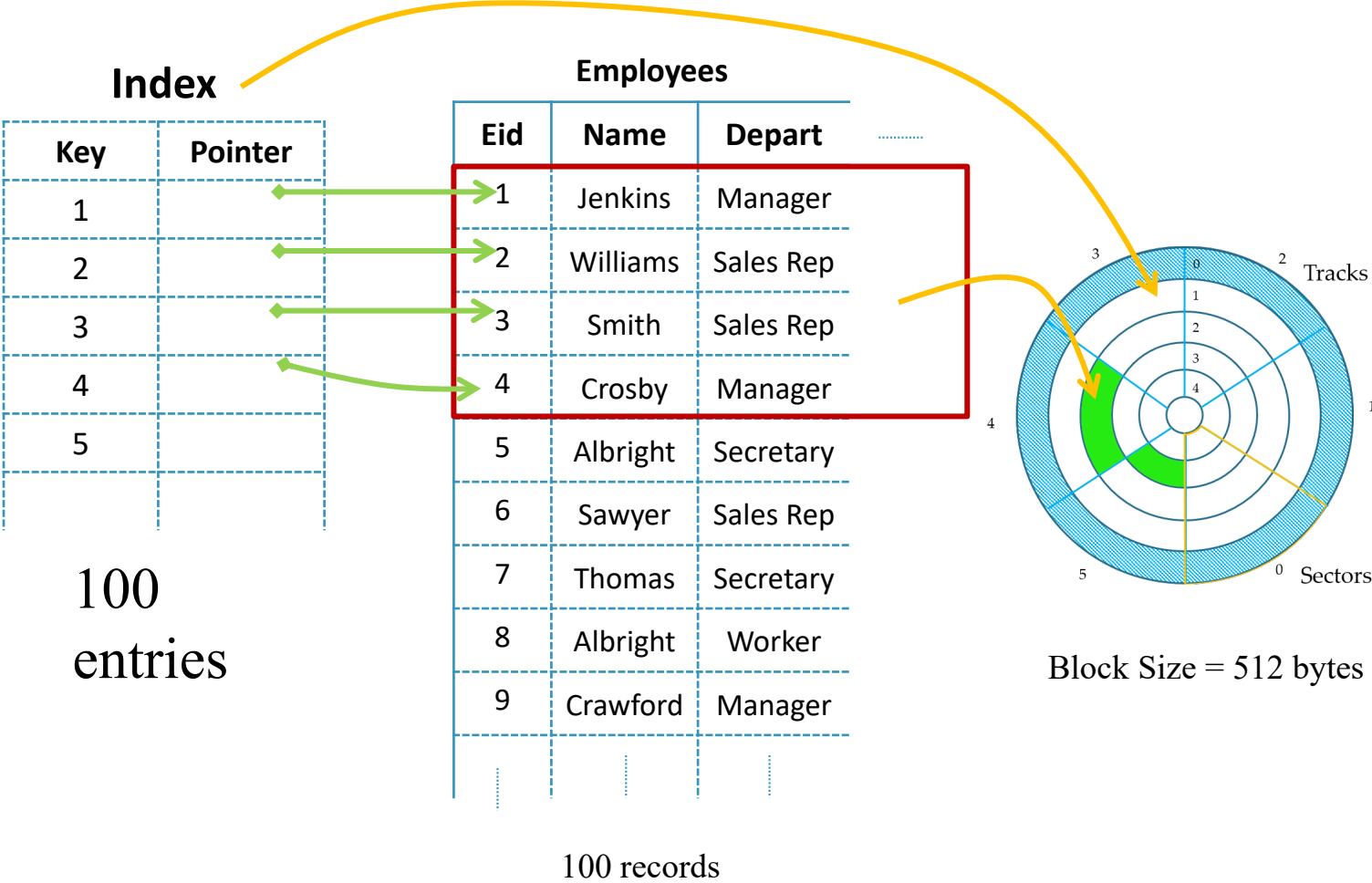


# What is Indexing?

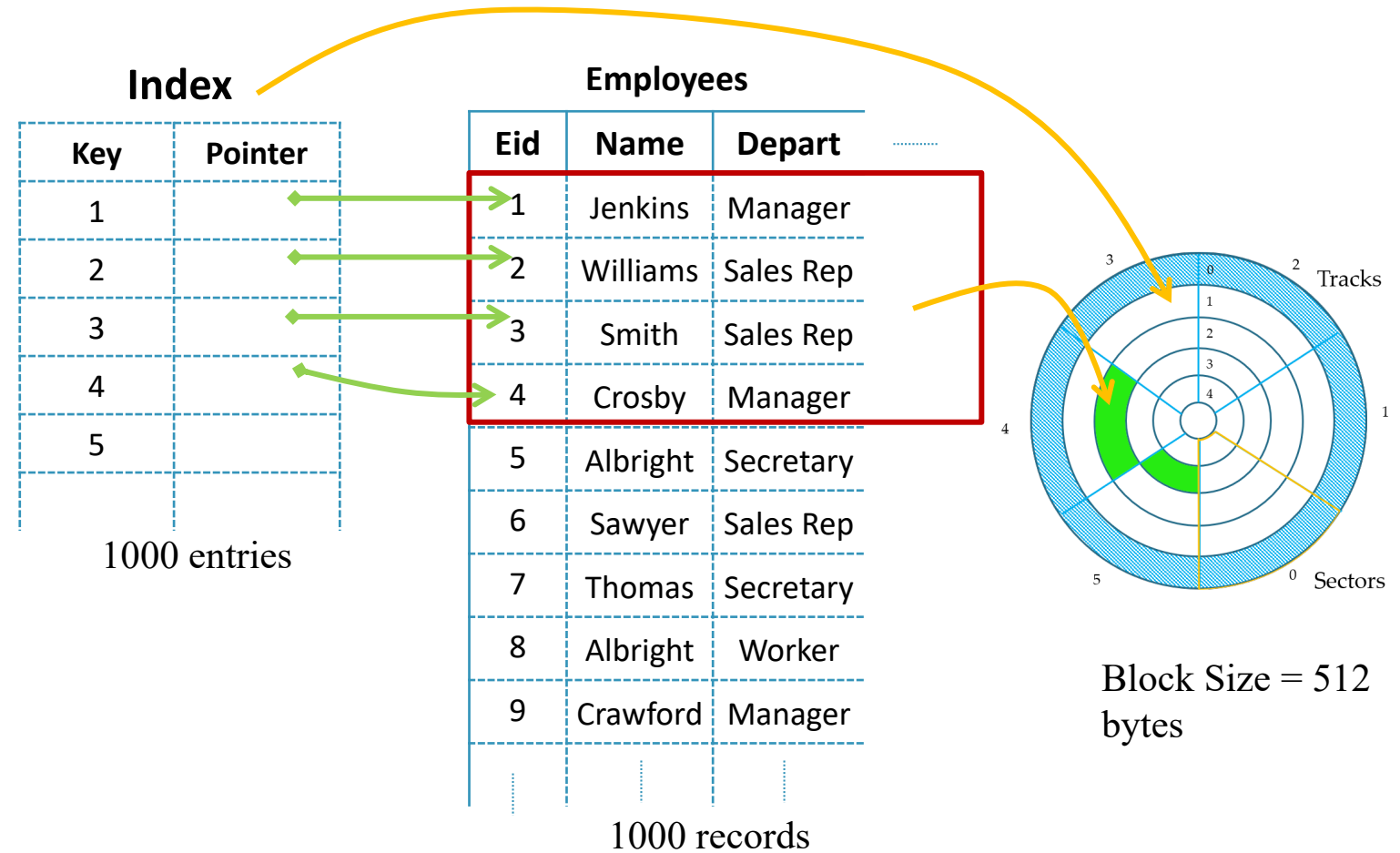
**SELECT** Eid=7  
**FROM** Employees

100 entries can be stored in = 1600/512  
≈ **4 blocks**

Total number of blocks required = 4 +1  
= **5 blocks**

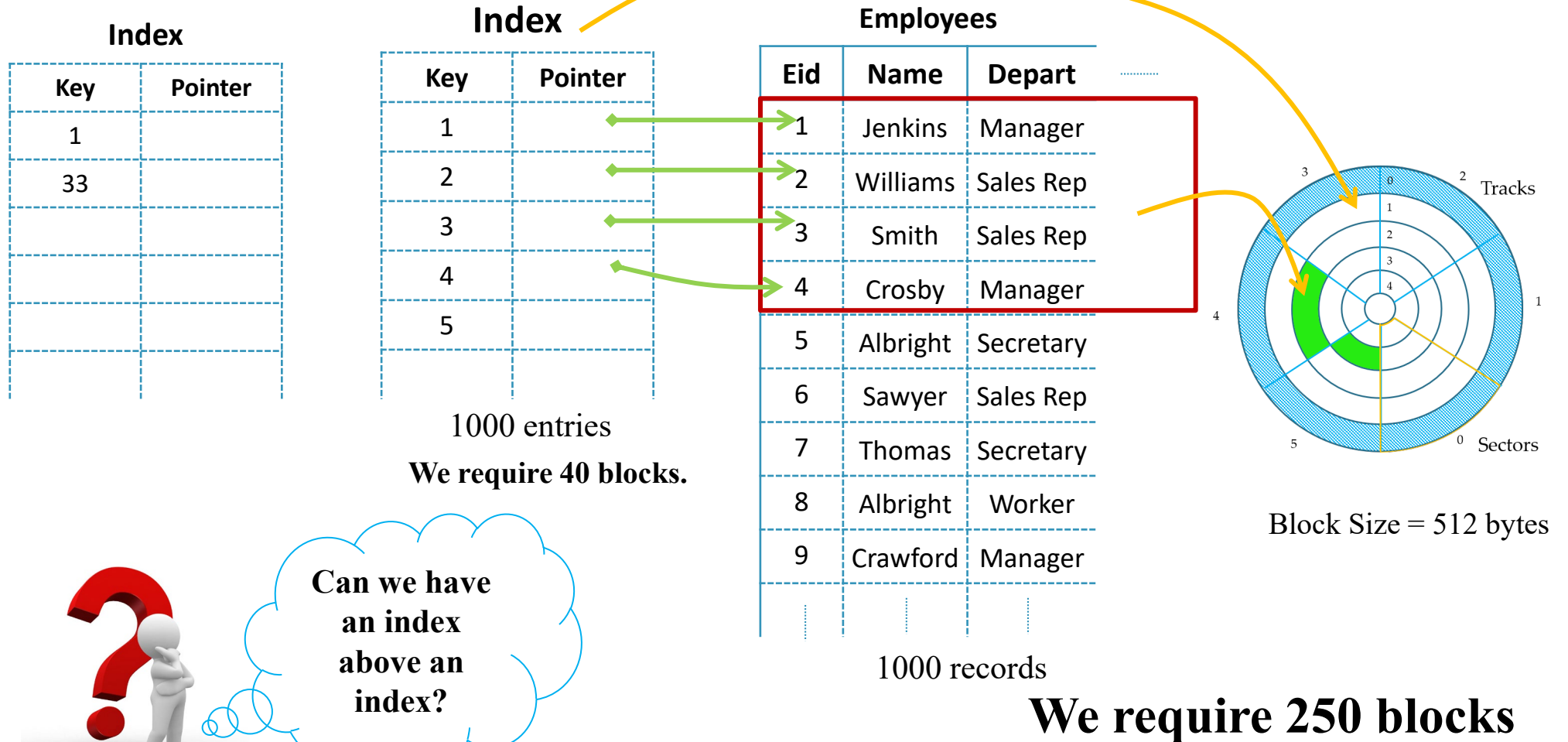


# What is Indexing?





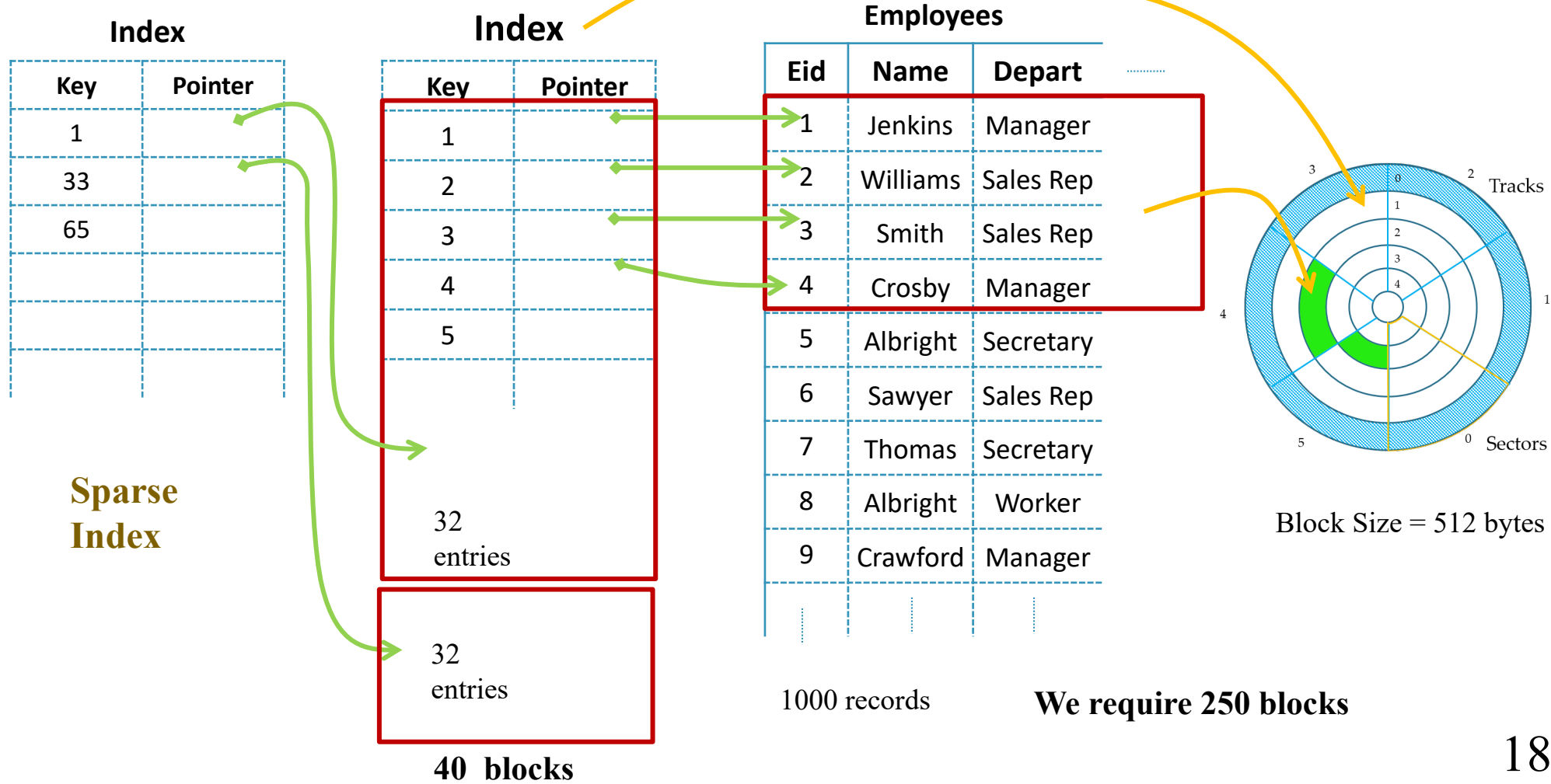
# What is multi-level indexing?



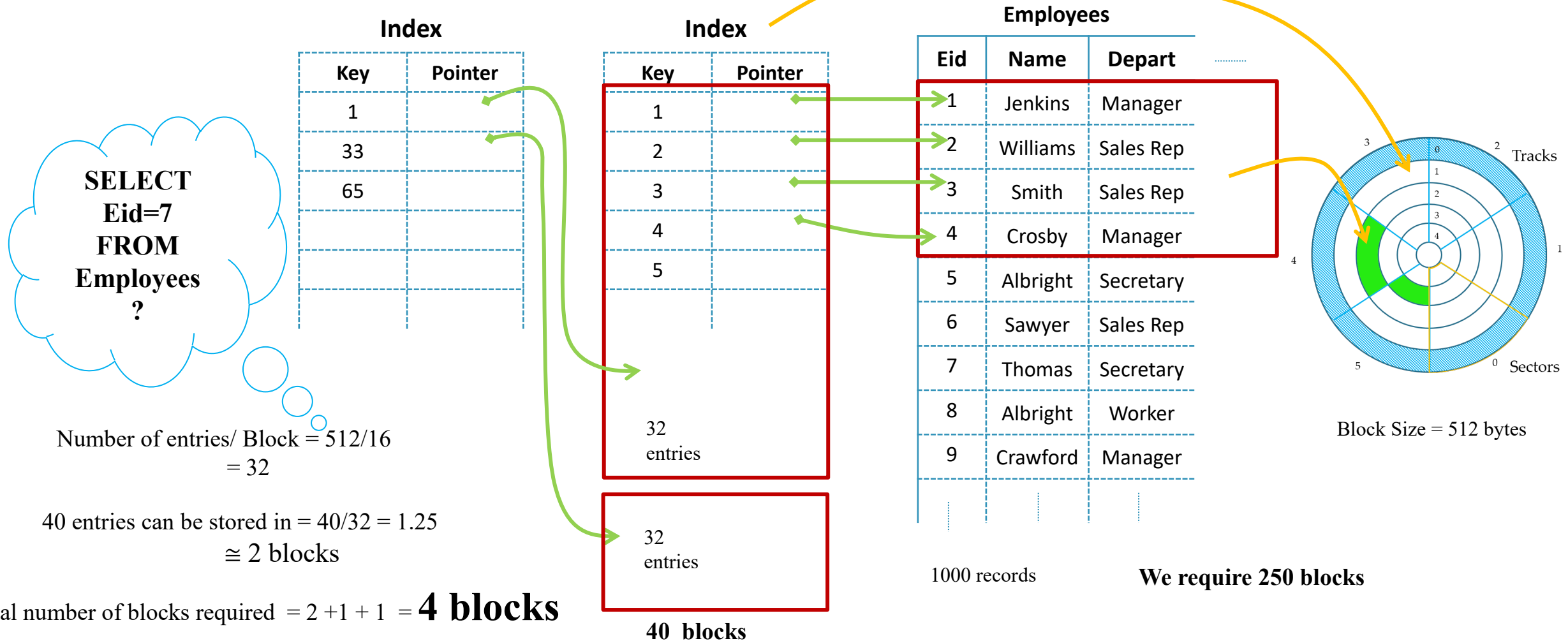




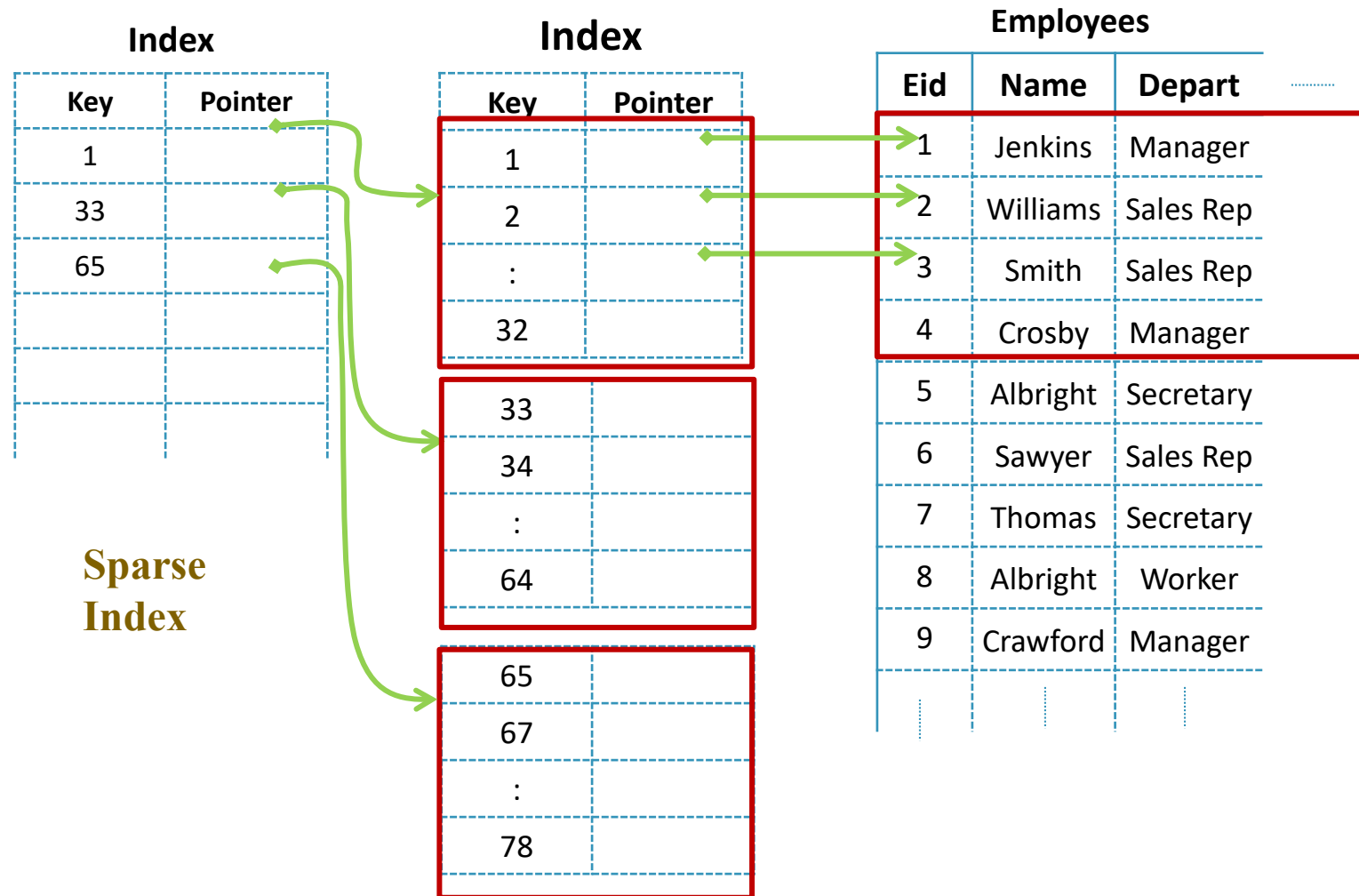
# What is multi-level indexing?

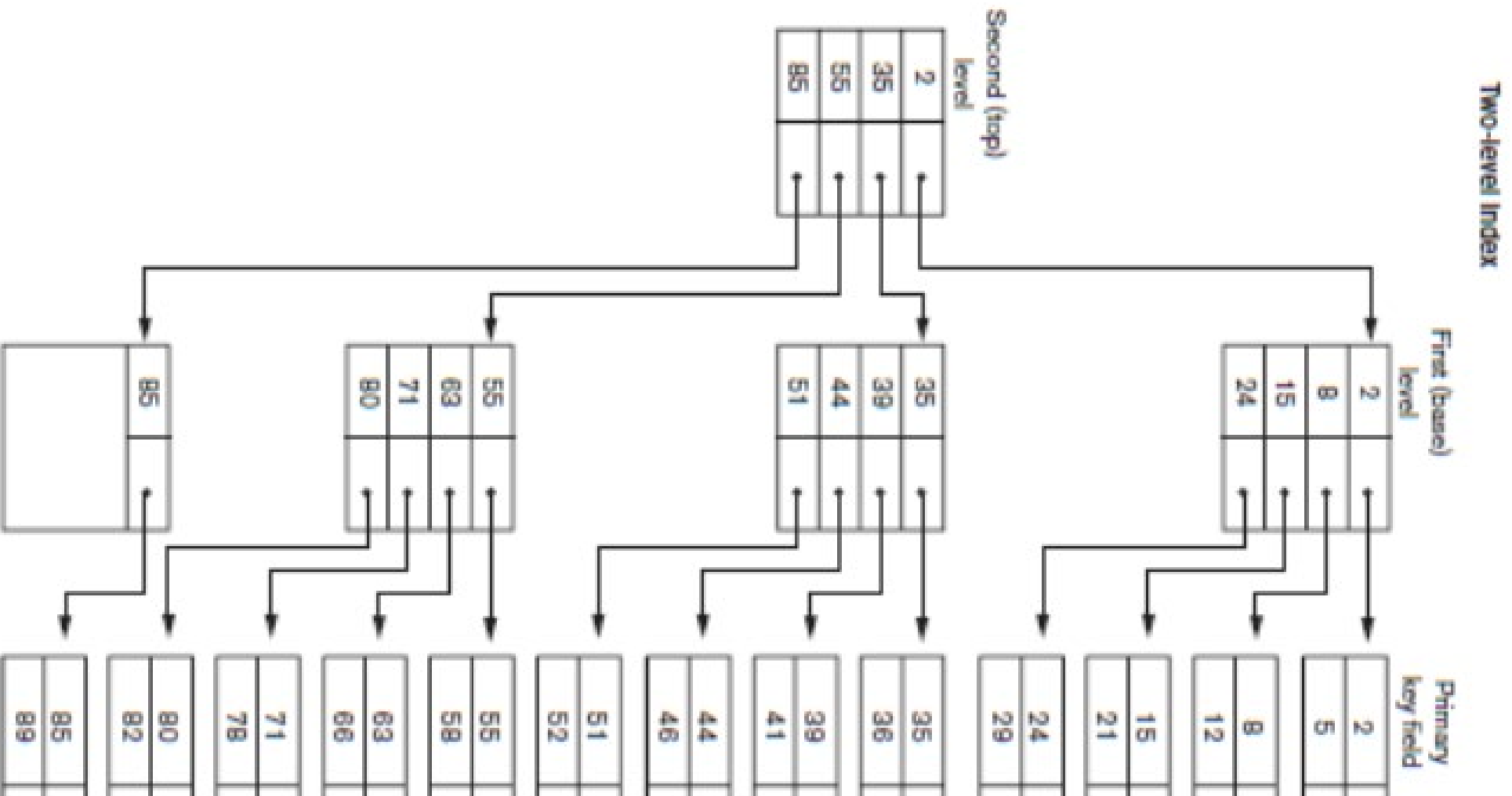


# What is multi-level indexing?



# What is multi-level indexing?





# Indexing cont.

- Database Indexes are special files which are stored along with the Data files of the Database
- Database indexes help the query processor to **expedite** search and retrieval of records from Blocks
- By default, the DBMS typically creates an index based on the PK of each table
  - However, depending on the type and frequency of queries, users can create their own indexes for better overall performance
- When queries corresponding to particular tables are executed, the corresponding index files (if any) are transferred to the main memory for the **query optimization** module to work on



# Single-Level Ordered Indexes



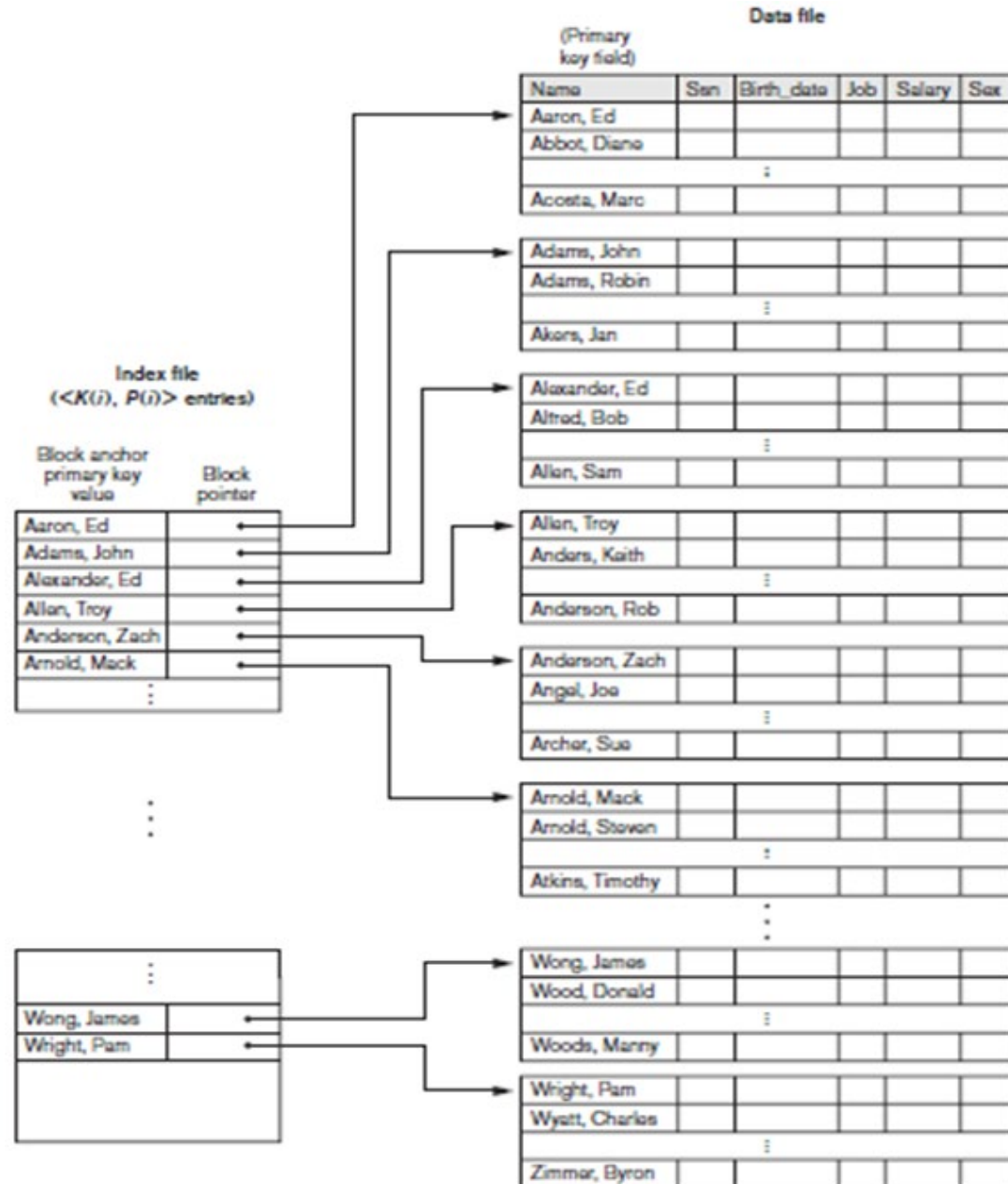
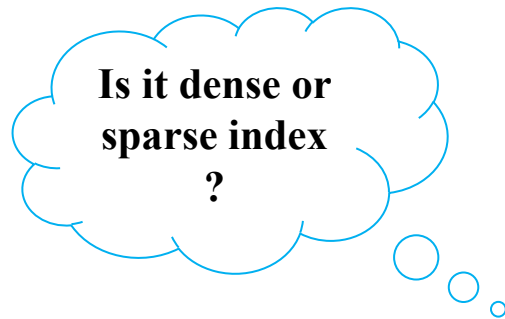
# Single-Level Ordered Index

- Primary Index
- Clustered Index
- Secondary Index



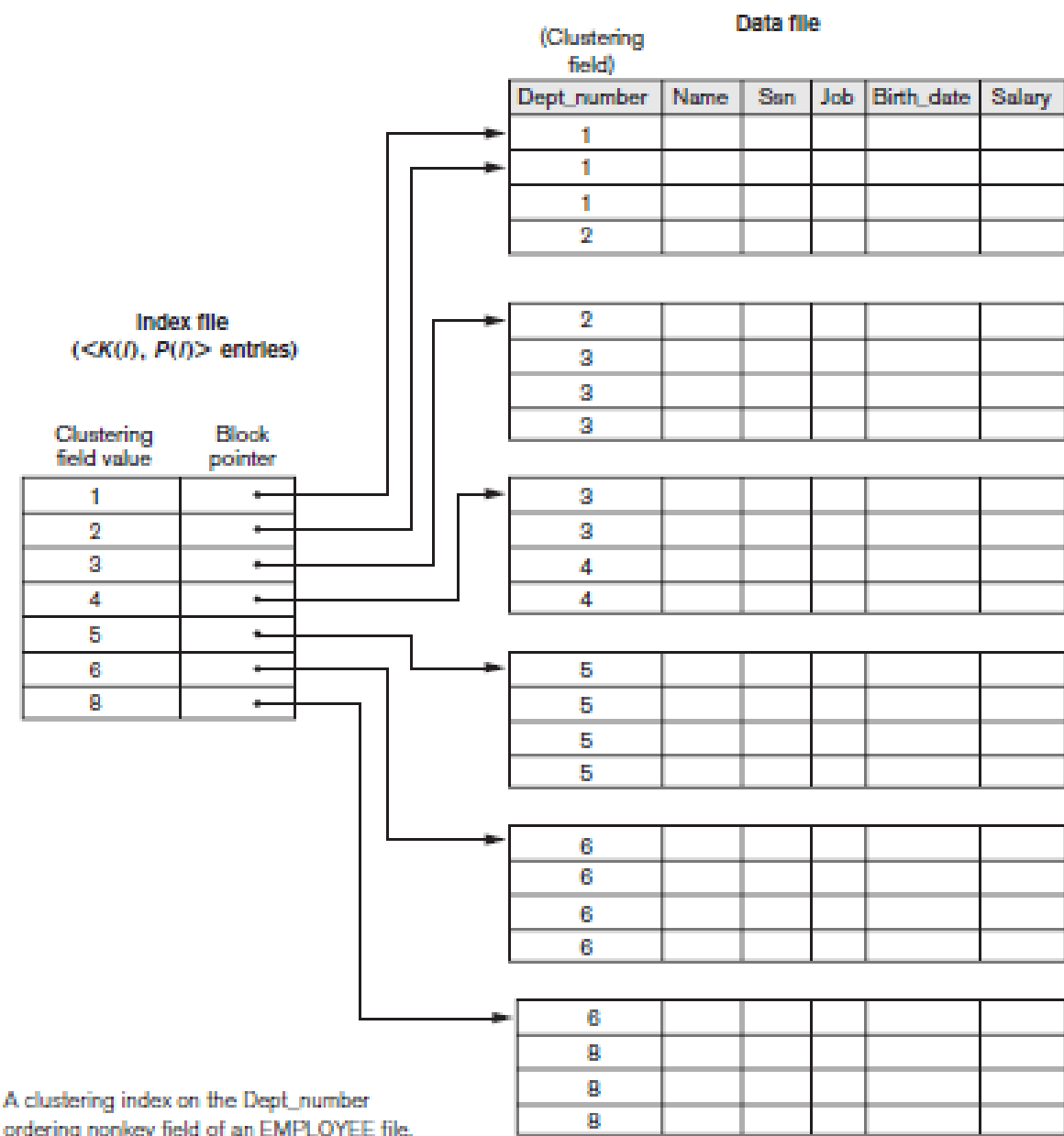
# Primary Indexes

- Records are **sorted and stored** across blocks based on a key
- The first Record in each block is called the **Block anchor**
- Each block anchor has an entry in the index file that has a **pointer** to the block



# Clustering Index

- Records are **sorted and stored** across blocks based on a **Clustering Field value** (non-key value)
- Each Clustering Field value has an entry in the index file that has a **pointer** to the **first block** that contains the Clustering Field value

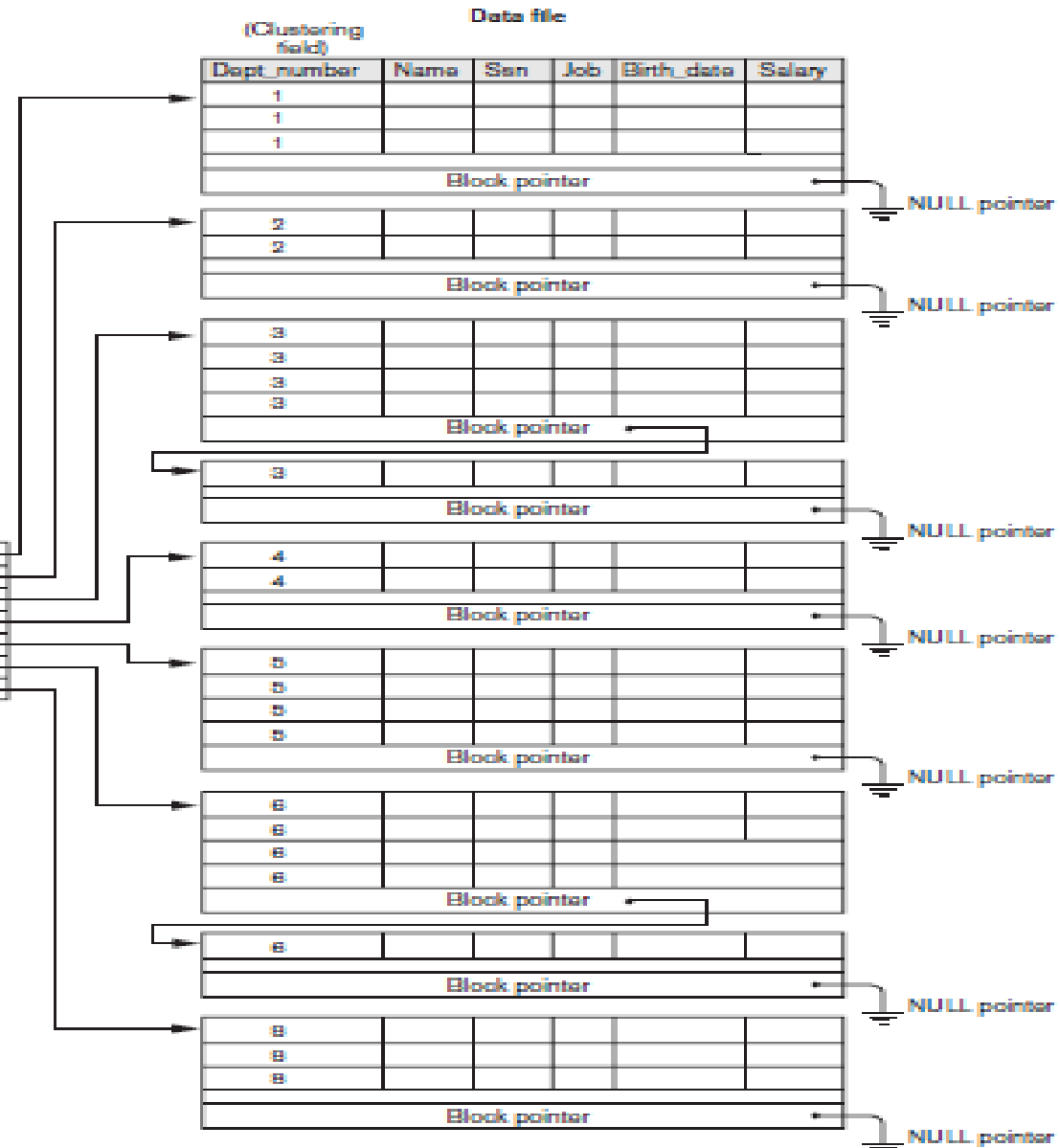


# Clustering Indexes with Separate Block Clusters

Clustering index with a separate block cluster for each group of records that share the same value for the clustering field.

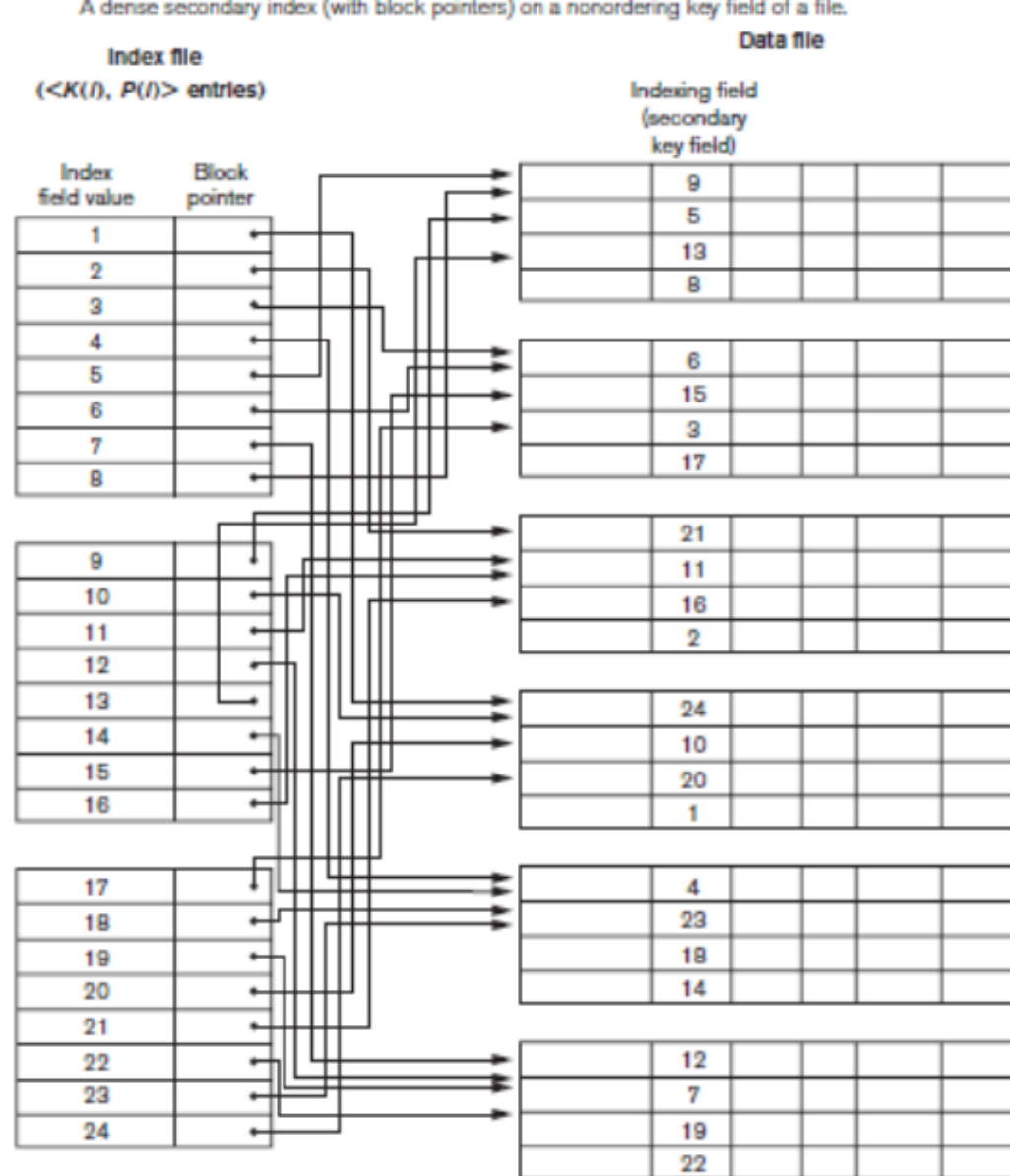
Index file  
( $\langle K(i), P(i) \rangle$  entries)

Clustering field value	Block pointer
1	→
2	→
3	→
4	→
5	→
6	→
8	→



# Secondary Index

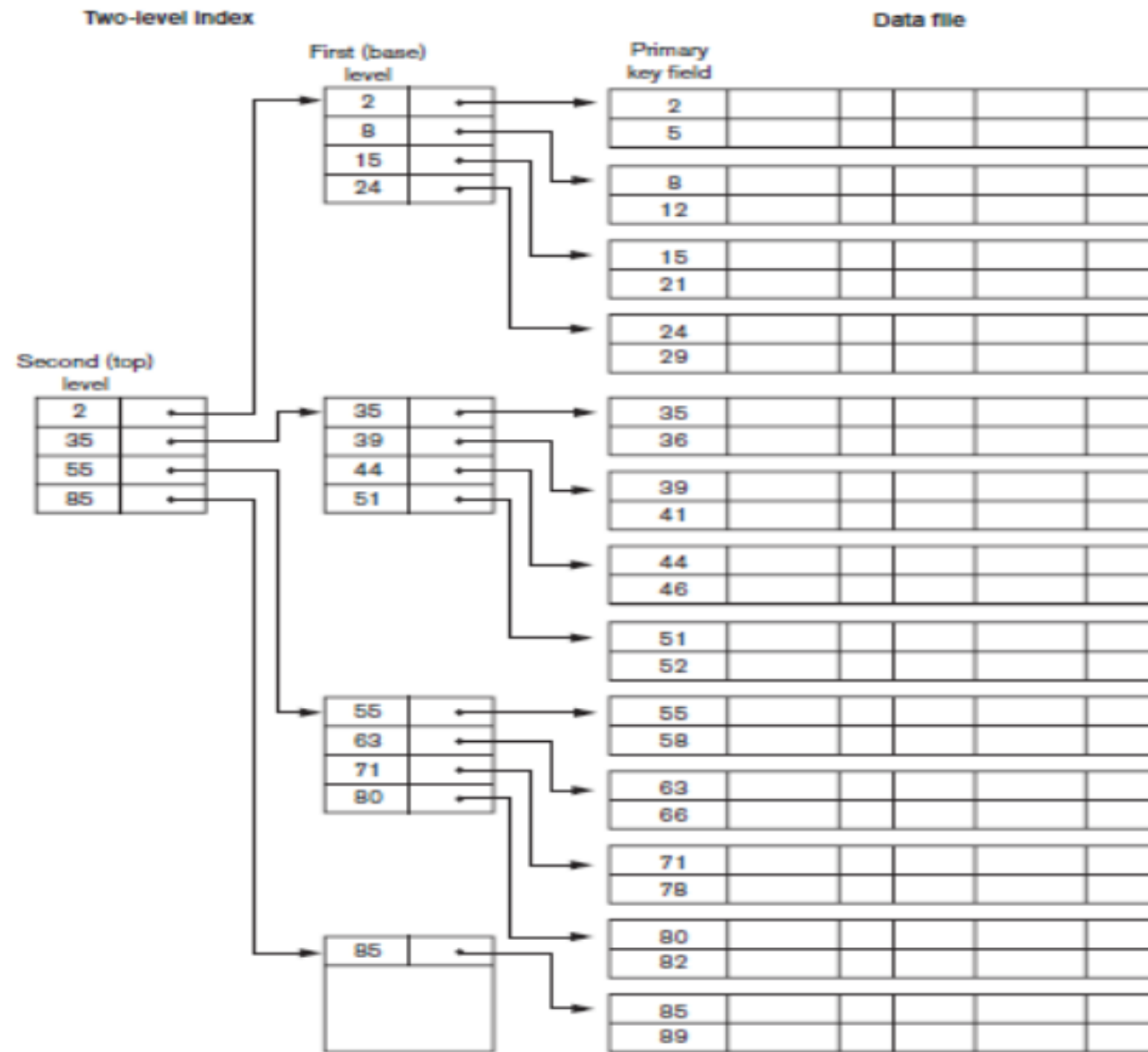
- Indexes based on **other key values** (which are not the PK)
- The records in this case are **not sorted** across the blocks based on the secondary key fields (Since they are sorted using the PK)
- Anchoring is not possible (Like in Primary Index)



# Multilevel Indexes



## Two-Level Primary Index

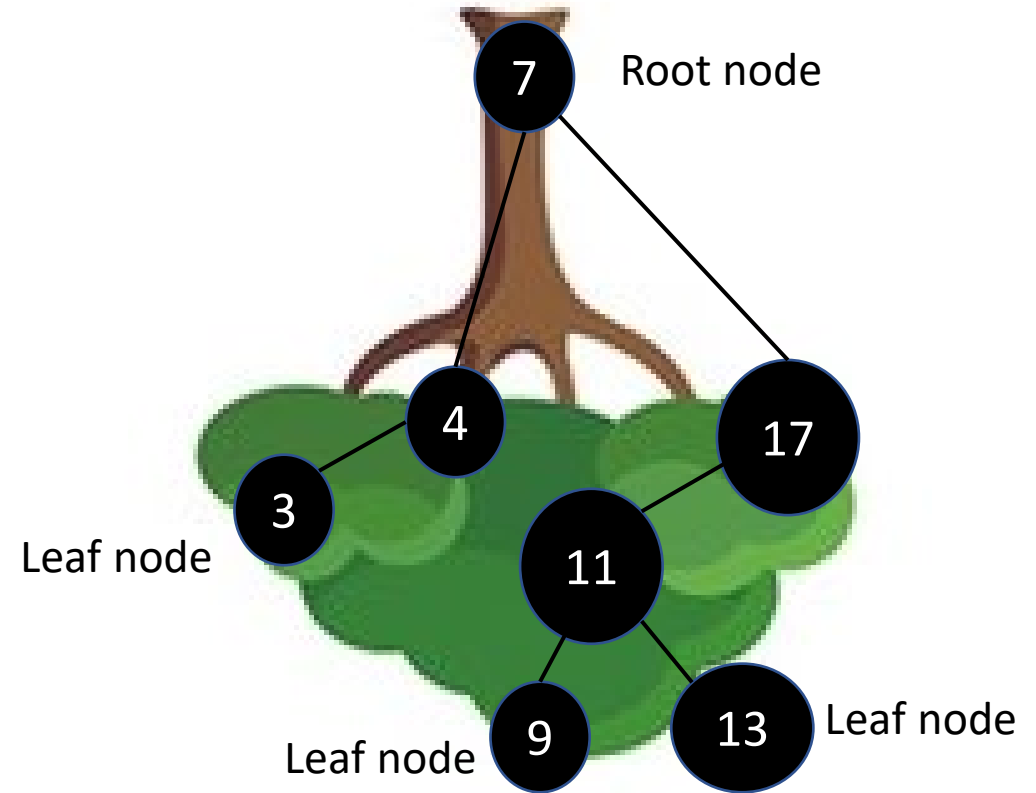


# Dynamic Multilevel Indexes using B Trees and B+ Trees



# TREE

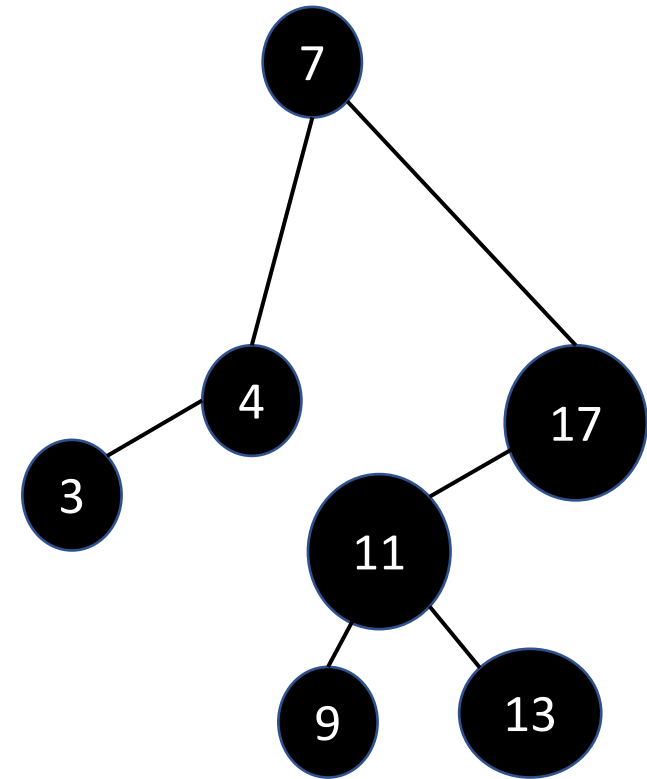
Tree: An organizational structure for the storage and retrieval of data



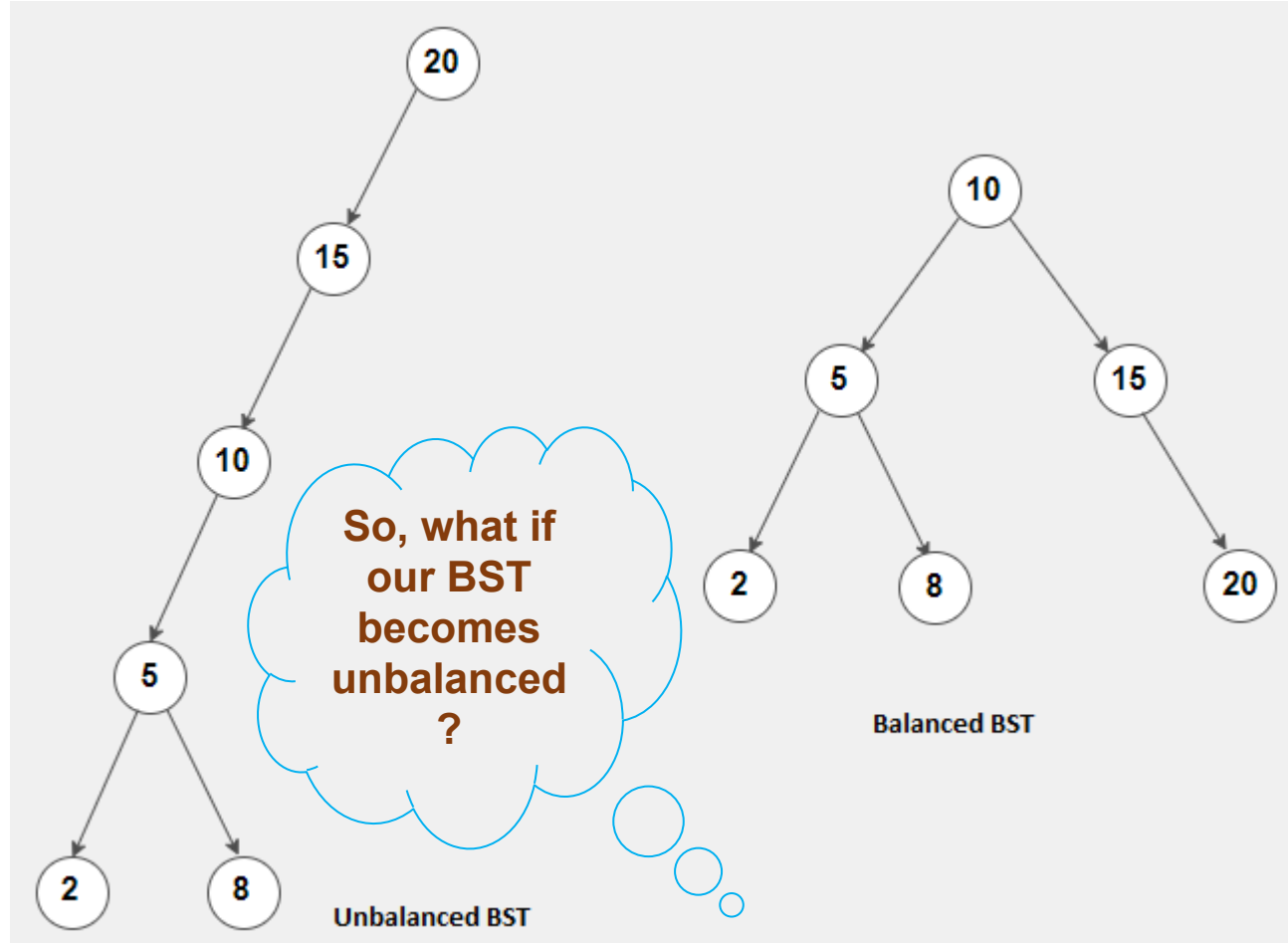


# Binary Search Tree

- Binary search trees allow binary search for fast lookup, addition, and removal of data items.
- A node can have at most two children.



# Binary Search Tree BST (Balanced vs. Unbalanced)

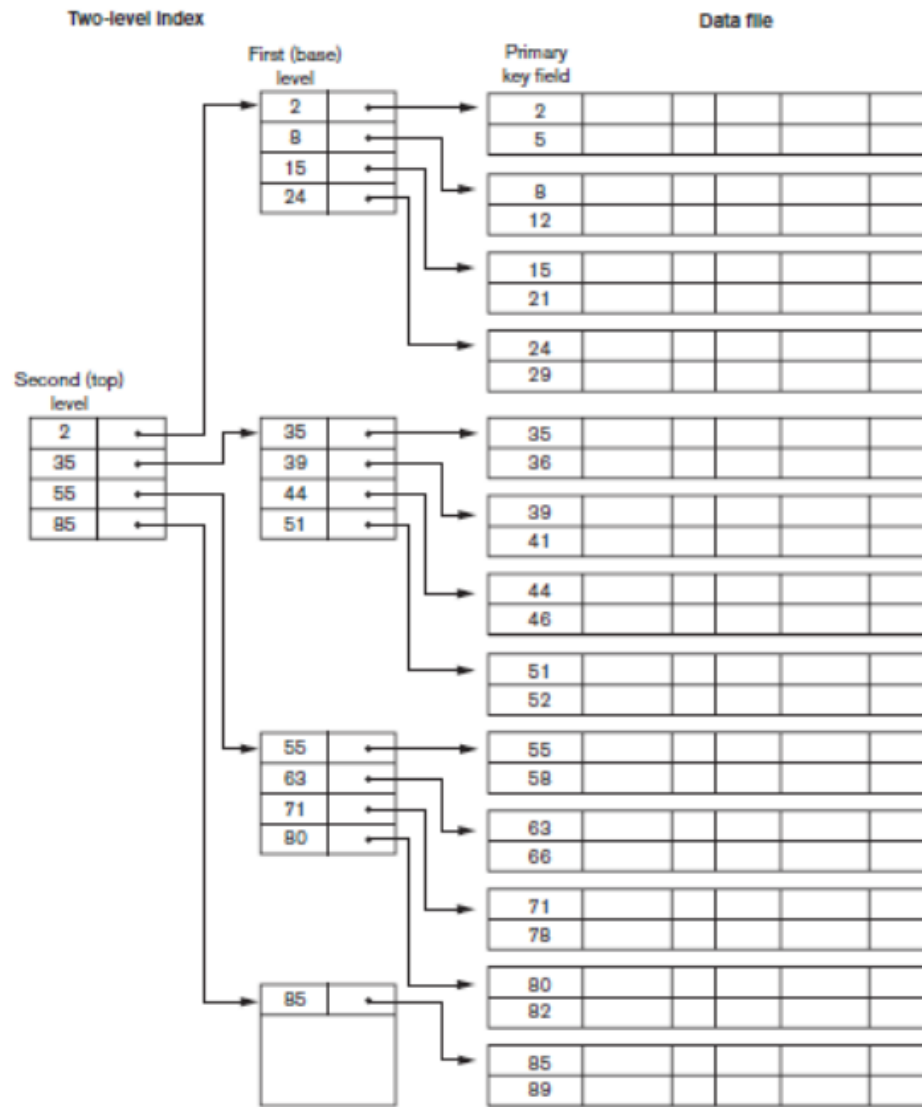
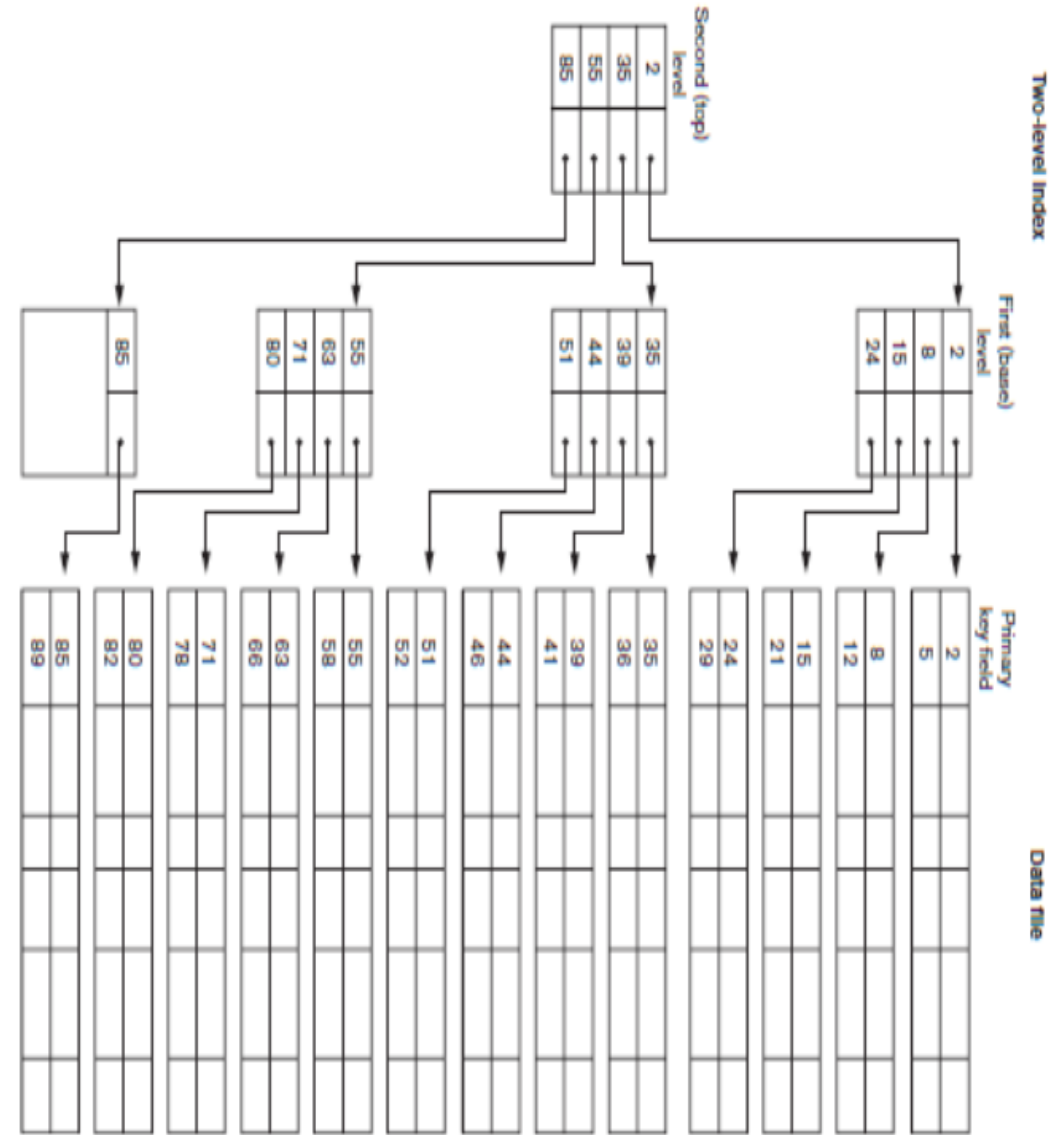


- **Balanced BST:** The difference between the height of the left sub-tree and the right sub-tree of **every node** is not more than 1
- The time complexity of operations on the binary search tree is directly proportional to the height of the tree.  
**Balanced BST:  $O(\log(n))$**   
**Unbalanced BST:  $O(n)$**
- It can also be expressed in terms of its height ( $h$ ) as  $O(h)$
- A balanced BST minimizes the search time by minimizing the value of  $h$

# B-Tree

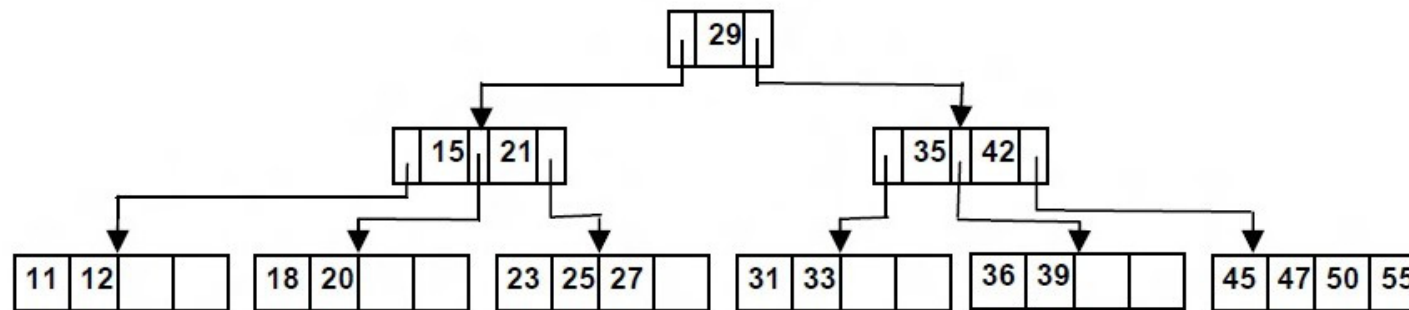
- A B-tree is a self-balancing tree data structure that maintains sorted data and allows searches, sequential access, insertions, and deletions in logarithmic time.
- Invented by Rudolf Bayer and Ed McCreight in 1972
- The B-tree **generalizes** the binary search tree
  - Allows nodes with more than two children (Whereas a node in a BST cannot have more than two children)
- Useful for implementing multi-level indexing



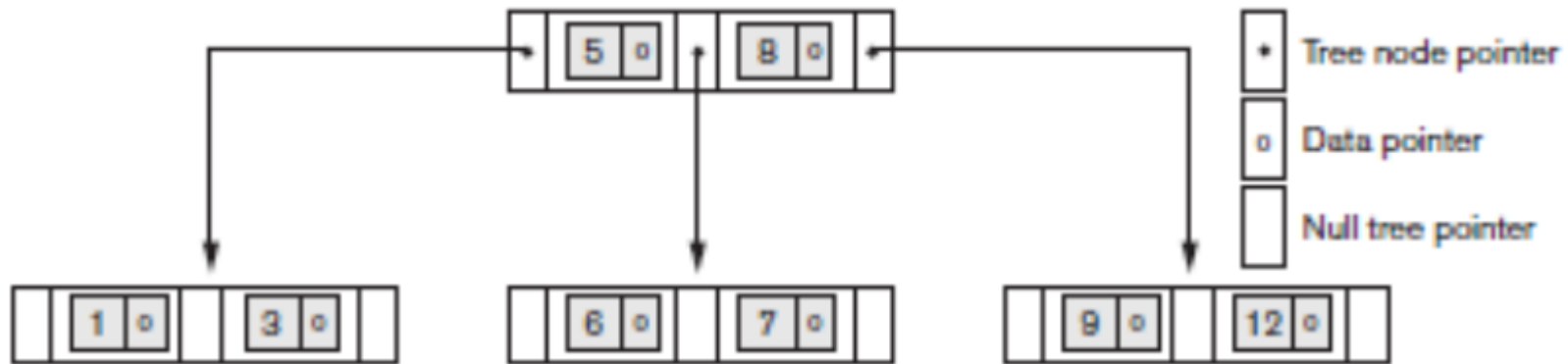


# B Tree Rules

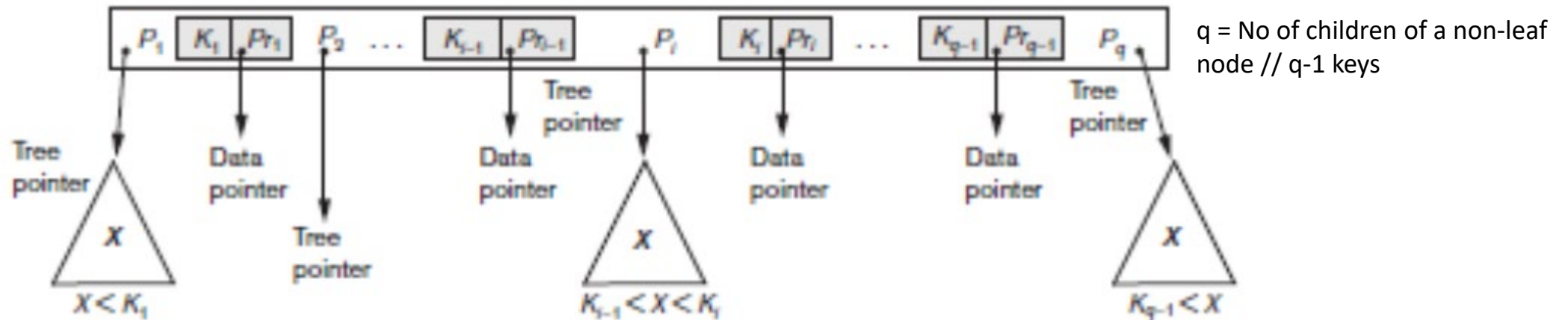
- A B-tree of order  $m$  is a tree which satisfies the following properties:
  - Every node has at most  $m$  children.
  - Every non-leaf node (except root) has at least  $\lceil m/2 \rceil$  child nodes.
  - The root has at least two children if it is not a leaf node.
  - A non-leaf node with  $k$  children contains  $k - 1$  keys.



B-Tree of order 5



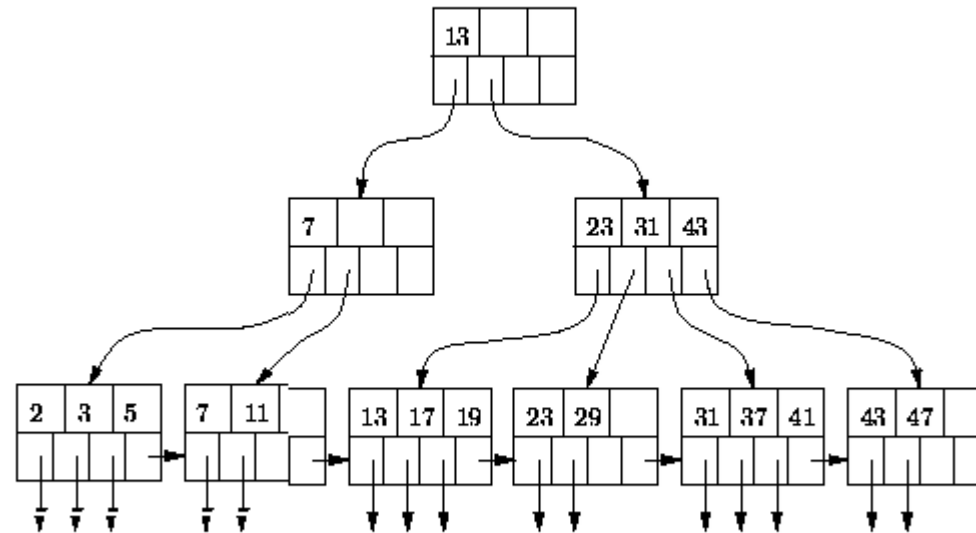
**B-Tree of order 3 (With Tree-Node Pointer and Data Pointer)**



**A node in a B-Tree with  $q-1$  search values**

# B+ Trees

- Data pointers only at the leaf node to facilitate faster search



# B TREE DEMO





# INDEX CREATION

## The general form

**CREATE** [ **UNIQUE** ] **INDEX** <index name>

**ON** <table name> ( <column name> [ <order> ] { , <column name> [ <order> ] } )

[ **CLUSTER** ] ;

- **CREATE INDEX** DnoIndex  
ON EMPLOYEE (Dno)  
CLUSTER ; //Clustered Index
- **CREATE INDEX** EmpAgeIndex  
ON EMPLOYEE(Age DESC)  
CLUSTER; //Clustered Index
- **CREATE INDEX** EmpEmailIndex  
ON EMPLOYEE (Email) // Creates a secondary index
- **CREATE INDEX** EmpAgeService  
ON EMPLOYEE(Age,YearsOfService); //Dense Index

## Choose indexes – Guidelines for choosing ‘wish-list’

1. Do not index small relations.
2. Add secondary index to a FK if it is frequently accessed.
3. Add secondary index to any attribute heavily used as a secondary key.
4. Add secondary index on attributes involved in: selection or join criteria; ORDER BY; GROUP BY; and other operations involving sorting (such as UNION or DISTINCT).
5. Add secondary index on attributes involved in built-in functions.
6. Avoid indexing an attribute or relation that is frequently updated.
7. Avoid indexing an attribute if the query will retrieve a significant proportion of the relation.
8. Avoid indexing attributes that consist of long character strings.



# Choose indexes

Suppose, we have a table named buyers where the SELECT query uses indexes like below:

```
SELECT buyer_id          /* no need to index */
FROM buyers
WHERE first_name='Carolyn' /* consider to use index */
AND last_name='Begg'      /* consider to use index */
```

Since "buyer\_id" is referenced in the SELECT portion, MySQL will not use it to limit the chosen rows. Hence, there is no great need to index it

The below is another example little different from the above one:

```
SELECT buyers.buyer_id, country.name /* no need to index */
FROM buyers LEFT JOIN country ON buyers.country_id=country.country_id /* consider to use index */
WHERE first_name='Carolyn' /* consider to use index */
AND last_name='Begg' /* consider to use index */
```

According to the above queries first\_name, last\_name columns can be indexed as they are located in the WHERE clause. Also an additional field, country\_id from country table, can be considered for indexing because it is in a JOIN clause. So indexing can be considered on every field in the WHERE clause or a JOIN clause.

# Question

```
Select * FROM Student, Apply, University  
WHERE Student.sID=Apply.sID AND Apply.uName=University.uName  
AND Student.GPA>1.5 AND University.uName<'Uwindsor'
```

Suppose we are allowed to create two indexes, and assume all indexes are tree based. Which two indexes do you think would be most useful for speeding up query execution?

1. Student.sID, University.uName
2. Student.sID, Student.GPA
3. Apply.uName, University.uName
4. Apply.sID, Student.GPA



# Downsides of Indexes

- ✓ Extra space
- ✓ Index creation
- ✓ Index maintenance



# Summary and Conclusion

- File Organization
- Introduction to Indexing
- Single-Level Ordered Indexes
  - Primary Indexes //SPARSE
  - Clustering Indexes // SPARSE (Optionally with Block Pointer)
  - Secondary Indexes //DENSE
- Multilevel Indexes
  - Two-Level Primary Indexing //SPARSE
- Dynamic Multilevel Indexes
  - B Trees //Generalized version of Balances Search Tree- Allows more than 2 children per node
  - B+ Trees // SELF BALANCING – Data Pointers only at the leaf level
- CREATE INDEX Command in SQL
- Choosing an INDEX



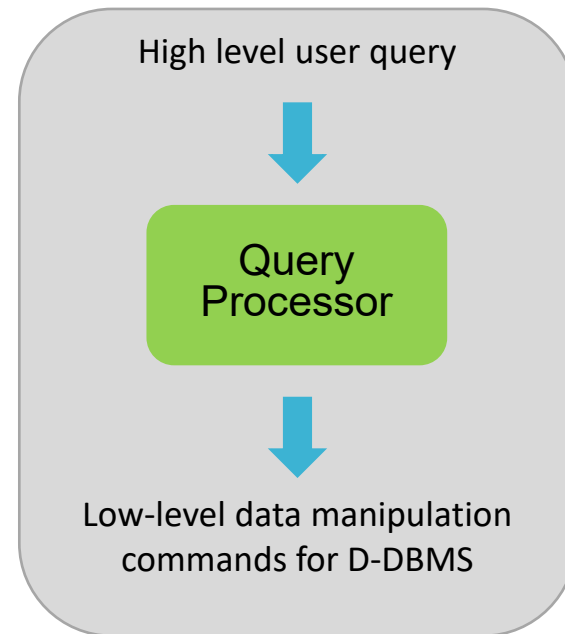
# Extra Material

The image features a minimalist design. On the left is a large white rectangle. To its right is a light gray rectangle. Within the gray rectangle, the text "Extra Material" is written in a bold, brown, serif font. In the bottom right corner of the gray rectangle, there is a yellow triangle pointing towards the bottom right corner of the overall image.

# Query Processing

- **Aim:**

Transform a query written in a high-level language, typically SQL, into a correct and efficient execution strategy expressed in a low-level language (implementing the relational algebra), and to execute the strategy to retrieve the required data.

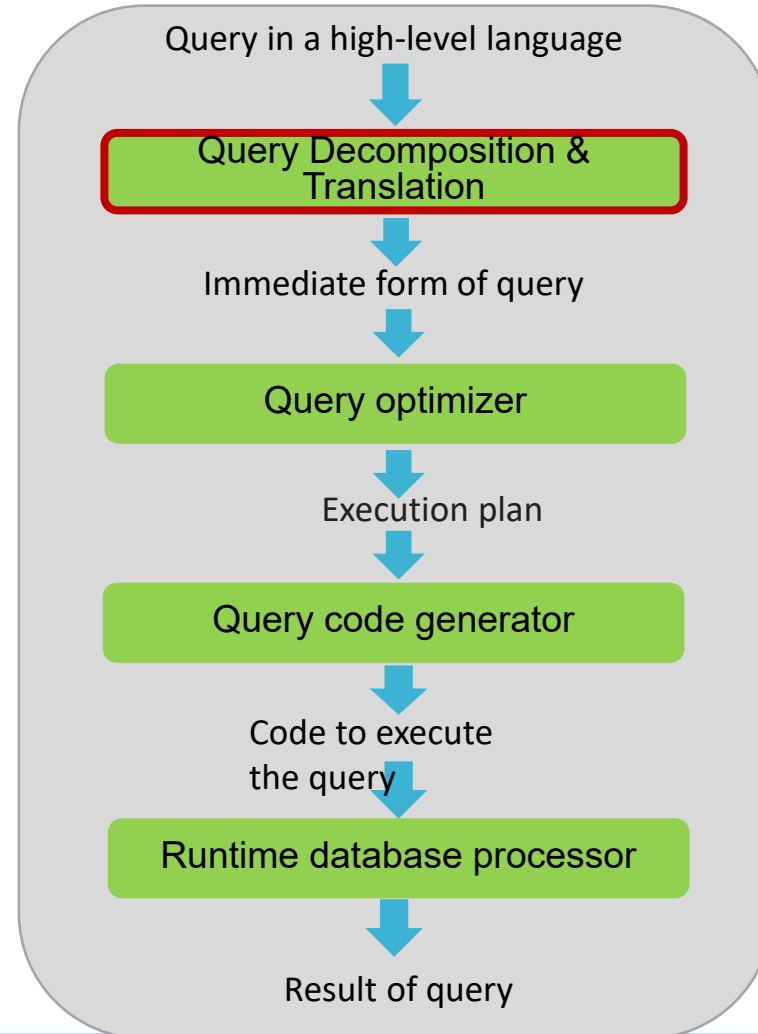


An important aspect of query processing is **query optimization**.





# Query Processing



Transform a high-level query into a relational algebra query and to check whether the query is syntactically and semantically correct.

The typical stages:

- ✓ Analysis,
- ✓ Normalization,
- ✓ Semantic analysis,
- ✓ Simplification, and
- ✓ Query restructuring.



# Query Decomposition

- **Analysis:**

```
Select staffNumber  
From Staff  
Where position > 10;
```

Staff

StaffNo	fName	lName	Position	Sex	DOB	Salary	BranchNo
SL21	John	White	Manager	M	1-Oct-45	30000	B005
SG37	Ann	Beech	Assistant	F	10-Nov-60	12000	B003
SG14	David	Ford	Supervisor	M	24-Mar-58	18000	B003
SA9	Mary	Howe	Assistant	F	19-Feb-70	9000	B007
SG5	Susan	Brand	Manager	F	3-Jun-40	24000	B003
SL41	Julie	Lee	Assistant	F	13-Jun-65	9000	B005

This query would be rejected on two grounds:

- (1) In the select list, the attribute staffNumber is not defined for the Staff relation (should be staffNo).
- (2) In the WHERE clause, the comparison “>10” is incompatible with the data type position, which is a variable character string.



# Query Decomposition

## Normalization:

- ✓ The normalization stage of query processing converts the query into a normalized form that can be more easily manipulated.
- ✓ Predicate can be converted into one of two forms:
  - **Conjunctive normal form:**  $(\text{position} = \text{'Manager'} \vee \text{salary} > 20000) \wedge \text{branchNo} = \text{'B003'}$
  - **Disjunctive normal form:**  $(\text{position} = \text{'Manager'} \vee \text{branchNo} = \text{'B003'}) \vee (\text{salary} > 20000 \vee \text{branchNo} = \text{'B003'})$

## Semantic analysis:

- ✓ The objective of semantic analysis is to reject normalized queries that are incorrectly formulated or contradictory.
- ✓ For example:
  - the predicate  $(\text{position} = \text{'Manager'} \vee \text{position} = \text{'Assistant'})$  on the Staff relation  $\rightarrow$  **contradictory**



# Query Decomposition

## Simplification:

- ✓ detect redundant qualifications,
- ✓ eliminate common subexpressions, and
- ✓ transform the query to a semantically equivalent but more easily and efficiently computed form.
- ✓ Access restrictions, view definitions, and integrity constraints are considered at this stage.
  - introduce redundancy.
- ✓ Assuming user has appropriate access privileges, first apply well-known idempotency rules of Boolean algebra.

$$p \dot{\cup} (p) \equiv p$$

$$p \dot{\cup} \text{false} \equiv \text{false}$$

$$p \dot{\cup} \text{true} \equiv p$$

$$p \dot{\cup} (\sim p) \equiv \text{false}$$

$$p \dot{\cup} (p \dot{\cup} q) \equiv p$$

$$p \dot{\cup} (p) \equiv p$$

$$p \dot{\cup} \text{false} \equiv p$$

$$p \dot{\cup} \text{true} \equiv \text{true}$$

$$p \dot{\cup} (\sim p) \equiv \text{true}$$

$$p \dot{\cup} (p \dot{\cup} q) \equiv p$$



# Query Decomposition

For example, consider the following integrity constraint,:

```
CREATE ASSERTION OnlyManagerSalaryHigh  
CHECK ((position <> 'Manager' AND salary < 20000)  
OR (position = 'Manager' AND salary > 20000));
```

and consider the effect on the query:

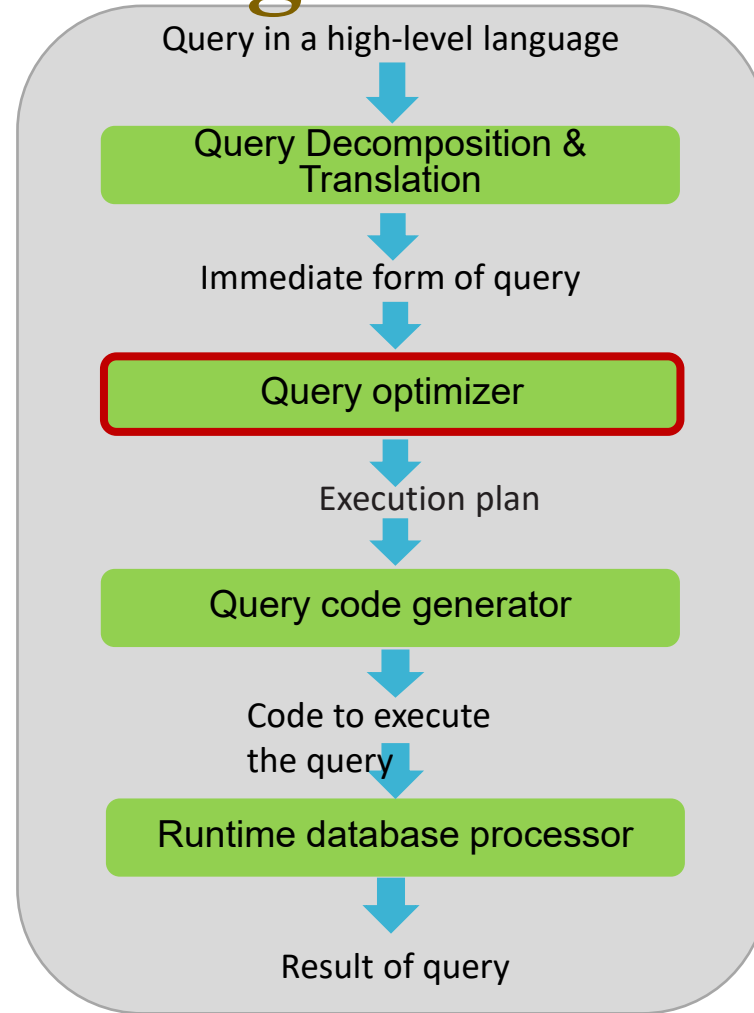
```
SELECT *  
FROM Staff  
WHERE (position = 'Manager' AND salary = 15000);
```

A contradiction of the integrity constraint so there can be no tuples that satisfy this predicate.

**Query restructuring:** the query is restructured to provide a more efficient implementation



# Query Processing



## Aim:

As there are many equivalent transformations of the same high-level query, choose the one that minimizes resource usage.

There are two main techniques for query optimization.

- ✓ **Heuristic rules**
- ✓ Systematically estimating



# Query optimization

## Heuristic rules:

- ✓ Uses **transformation rules** to convert one relational algebra expression into an equivalent form that is known to be more efficient.



# Transformation Rules for the Relational Algebra Operations

1. **Conjunctive Selection operations can cascade into individual Selection operations (and vice versa).**

$$\sigma_{p \wedge q \wedge r}(\mathbf{R}) = \sigma_p(\sigma_q(\sigma_r(\mathbf{R})))$$

Sometimes referred to as cascade of Selection.

$$\sigma_{\text{branchNo}='B003' \wedge \text{salary}>15000}(\mathbf{Staff}) = \sigma_{\text{branchNo}='B003'}(\sigma_{\text{salary}>15000}(\mathbf{Staff}))$$



# Transformation Rules for the Relational Algebra Operations

## 2. Commutativity of Selection operations.

$$\sigma_p(\sigma_q(R)) = \sigma_q(\sigma_p(R))$$

For example:

$$\sigma_{\text{branchNo}='B003'}(\sigma_{\text{salary}>15000}(\text{Staff})) = \sigma_{\text{salary}>15000}(\sigma_{\text{branchNo}='B003'}(\text{Staff}))$$



# Transformation Rules for the Relational Algebra Operations

3. In a sequence of Projection operations, only the last in the sequence is required.

$$\Pi_L \Pi_M \dots \Pi_N(R) = \Pi_L(R)$$

For example:

$$\Pi_{\text{IName}} \Pi_{\text{branchNo, IName}}(\text{Staff}) = \Pi_{\text{IName}}(\text{Staff})$$



# Transformation Rules for the Relational Algebra Operations

## 4. Commutativity of Selection and Projection.

If predicate  $p$  involves only attributes in projection list, Selection and Projection operations commute:

$$\Pi_{A_1, \dots, A_m}(\sigma_p(\mathbf{R})) = \sigma_p(\Pi_{A_1, \dots, A_m}(\mathbf{R})) \quad \text{where } p \in \{A_1, A_2, \dots, A_m\}$$

For example:

$$\Pi_{fName, lName}(\sigma_{lName='Beech'}(\mathbf{Staff})) = \sigma_{lName='Beech'}(\Pi_{fName, lName}(\mathbf{Staff}))$$



# Transformation Rules for the Relational Algebra Operations

## 5. Commutativity of Theta join (and Cartesian product).

$$R \bowtie_p S = S \bowtie_p R$$

$$R \times S = S \times R$$

Rule also applies to Equijoin and Natural join

For example:

$$\text{Staff} \bowtie_{\text{staff.branchNo}=\text{branch.branchNo}} \text{Branch} = \text{Branch} \bowtie_{\text{staff.branchNo}=\text{branch.branchNo}} \text{Staff}$$



# Transformation Rules for the Relational Algebra Operations

## 6. Commutativity of Selection and Theta join (or Cartesian product).

If the selection predicate involves only attributes of one of the relations being joined, then the Selection and Join (or Cartesian product) operations commute:

$$\begin{aligned}\sigma_p(\mathbf{R} \bowtie_r \mathbf{S}) &= (\sigma_p(\mathbf{R})) \bowtie_r \mathbf{S} \\ \sigma_p(\mathbf{R} \times \mathbf{S}) &= (\sigma_p(\mathbf{R})) \times \mathbf{S} \quad \text{where } p \in \{A_1, A_2, \dots, A_n\}\end{aligned}$$

If selection predicate is conjunctive predicate having form  $(p \wedge q)$ , where  $p$  only involves attributes of  $R$ , and  $q$  only attributes of  $S$ , Selection and Theta join operations commute as:

$$\begin{aligned}\sigma_{p \wedge q}(\mathbf{R} \bowtie_r \mathbf{S}) &= (\sigma_p(\mathbf{R})) \bowtie_r (\sigma_q(\mathbf{S})) \\ \sigma_{p \wedge q}(\mathbf{R} \times \mathbf{S}) &= (\sigma_p(\mathbf{R})) \times (\sigma_q(\mathbf{S}))\end{aligned}$$

For example:

$$\sigma_{\text{position}='Manager' \wedge \text{city}='London'}(\mathbf{Staff} \bowtie_{\text{Staff.branchNo}=\text{Branch.branchNo}} \mathbf{Branch}) = (\sigma_{\text{position}='Manager'}(\mathbf{Staff})) \bowtie_{\text{Staff.branchNo}=\text{Branch.branchNo}} (\sigma_{\text{city}='London'}(\mathbf{Branch}))$$

BranchNo	Street	City	Postcode
B005	22 Deer Rd	London	SW14EH
B007	16 Argyll St	Aberdeen	AB23SU
B003	163 Main St	Glasgow	G11 9XX

StaffNo	fName	IName	Position	Sex	DOB	Salary	BranchNo
SL21	John	White	Manager	M	1-Oct-45	30000	B005



# Transformation Rules for the Relational Algebra Operations

## 7. Commutativity of Projection and Theta join (or Cartesian product).

If projection list is of form  $L = L_1 \cup L_2$ , where  $L_1$  only has attributes of  $R$ , and  $L_2$  only has attributes of  $S$ , provided join condition only contains attributes of  $L$ , Projection and Theta join commute:

$$\Pi_{L_1 \cup L_2}(R \bowtie_r S) = (\Pi_{L_1}(R)) \bowtie_r (\Pi_{L_2}(S))$$

If join condition contains additional attributes not in  $L$  ( $M = M_1 \cup M_2$  where  $M_1$  only has attributes of  $R$ , and  $M_2$  only has attributes of  $S$ ), a final projection operation is required:

$$\Pi_{L_1 \cup L_2}(R \bowtie_r S) = \Pi_{L_1 \cup L_2}((\Pi_{L_1 \cup M_1}(R)) \bowtie_r (\Pi_{L_2 \cup M_2}(S)))$$

For example:

$$\Pi_{\text{position, city, branchNo}}(\text{Staff} \bowtie_{\text{Staff.branchNo=Branch.branchNo}} \text{Branch}) = (\Pi_{\text{position, branchNo}}(\text{Staff})) \bowtie_{\text{Staff.branchNo=Branch.branchNo}} (\Pi_{\text{city, branchNo}}(\text{Branch}))$$

and using the latter rule:

$$\Pi_{\text{position, city}}(\text{Staff} \bowtie_{\text{Staff.branchNo=Branch.branchNo}} \text{Branch}) = \Pi_{\text{position, city}}((\Pi_{\text{position, branchNo}}(\text{Staff})) \bowtie_{\text{Staff.branchNo=Branch.branchNo}} (\Pi_{\text{city, branchNo}}(\text{Branch})))$$



# Transformation Rules for the Relational Algebra Operations

8. **Commutativity of Union and Intersection (but not set difference).**

$$\mathbf{R \cup S = S \cup R}$$

$$\mathbf{R \cap S = S \cap R}$$

9. **Commutativity of Selection and set operations (Union, Intersection, and Set difference).**

$$\sigma_p(\mathbf{R \cup S}) = \sigma_p(\mathbf{S}) \cup \sigma_p(\mathbf{R})$$

$$\sigma_p(\mathbf{R \cap S}) = \sigma_p(\mathbf{S}) \cap \sigma_p(\mathbf{R})$$

$$\sigma_p(\mathbf{R - S}) = \sigma_p(\mathbf{S}) - \sigma_p(\mathbf{R})$$

10. **Commutativity of Projection and Union.**

$$\Pi_L(\mathbf{R \cup S}) = \Pi_L(\mathbf{S}) \cup \Pi_L(\mathbf{R})$$



# Transformation Rules for the Relational Algebra Operations

## 11. Associativity of Theta join (and Cartesian product).

Cartesian product and Natural join are always associative:

$$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$$

$$(R \times S) \times T = R \times (S \times T)$$

If join condition  $q$  involves attributes only from  $S$  and  $T$ , then Theta join is associative:

$$(R \bowtie_p S) \bowtie_{q \wedge r} T = R \bowtie_{p \wedge r} (S \bowtie_q T)$$

For example:

$$(\text{Staff} \bowtie_{\text{Staff.staffNo}=\text{PropertyForRent.staffNo}} \text{PropertyForRent}) \bowtie_{\text{ownerNo}=\text{Owner.ownerNo} \wedge \text{staff.IName}=\text{Owner.IName}}$$

$$\text{Owner} = \text{Staff} \bowtie_{\text{staff.staffNo} = \text{PropertyForRent.staffNo} \wedge \text{staff.IName}=\text{IName}} (\text{PropertyForRent} \bowtie_{\text{ownerNo}} \text{Owner})$$





# Transformation Rules for the Relational Algebra Operations

## 12. Associativity of Union and Intersection (but not Set difference).

$$(R \cup S) \cup T = S \cup (R \cup T)$$

$$(R \cap S) \cap T = S \cap (R \cap T)$$



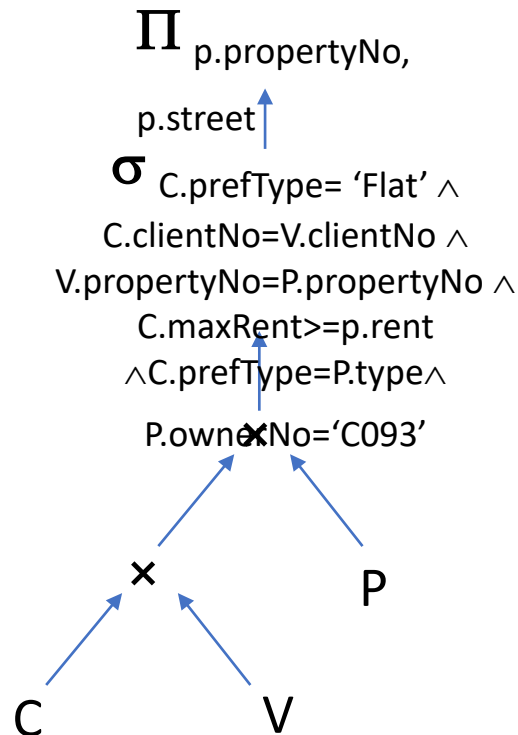
# Example : Use of Transformation Rules

For prospective renters of flats, find properties that match requirements and owned by CO93.

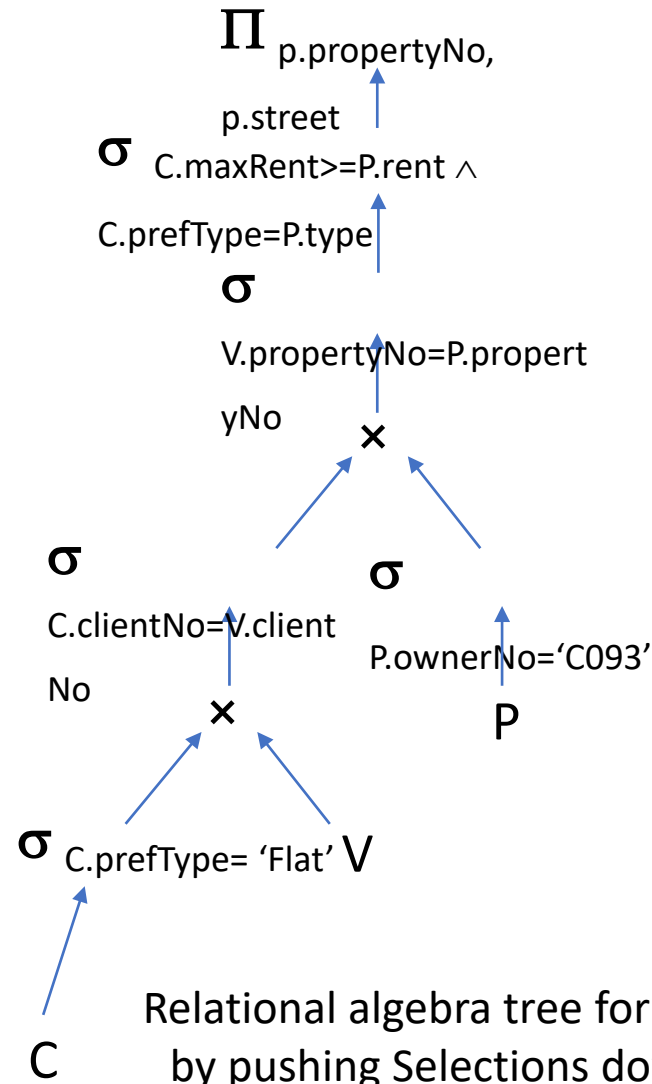
```
SELECT p.propertyNo, p.street
FROM Client c, Viewing v, PropertyForRent p
WHERE  c.prefType = 'Flat' AND
       c.clientNo = v.clientNo AND
       v.propertyNo = p.propertyNo AND
       c.maxRent >= p.rent AND
       c.prefType = p.type AND
       p.ownerNo = 'CO93';
```


$$\Pi_{p.propertyNo, p.street} (\sigma_{c.prefType='Flat' \wedge c.clientNo=v.clientNo \wedge v.propertyNo=p.propertyNo \wedge c.maxRent \geq p.rent \wedge c.prefType=p.type \wedge p.ownerNo='CO93'} ((c \times v) \times p))$$

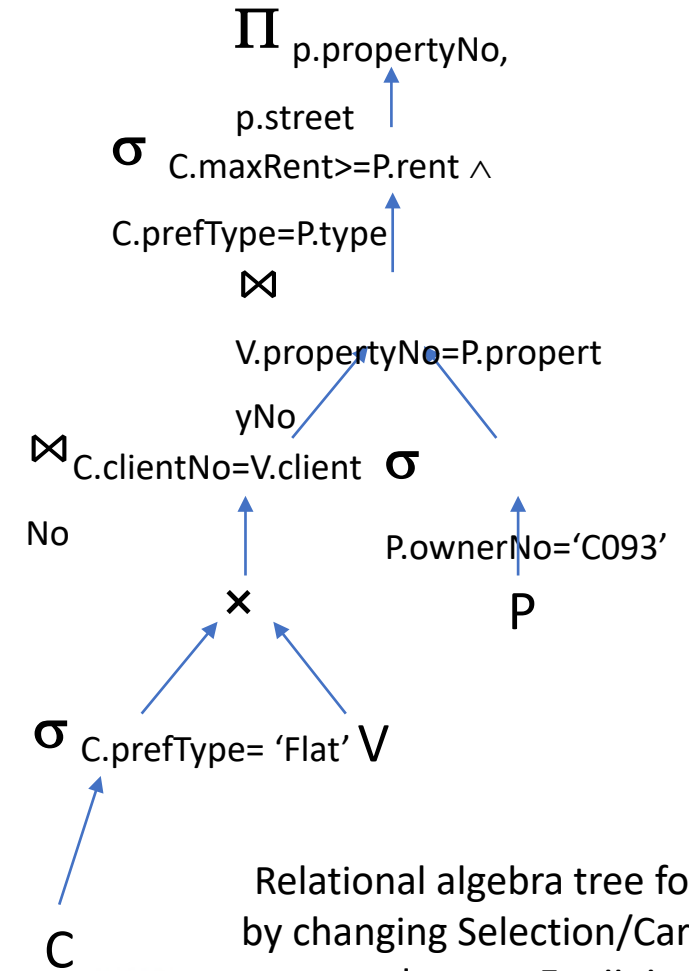

$\Pi_{p.propertyNo, p.street} (\sigma_{c.prefType = 'Flat' \wedge c.clientNo = v.clientNo \wedge v.propertyNo = p.propertyNo \wedge c.maxRent \geq p.rent \wedge c.prefType = p.type \wedge p.ownerNo = 'C093'} ((c \times v) \times p))$



Canonical relational algebra tree



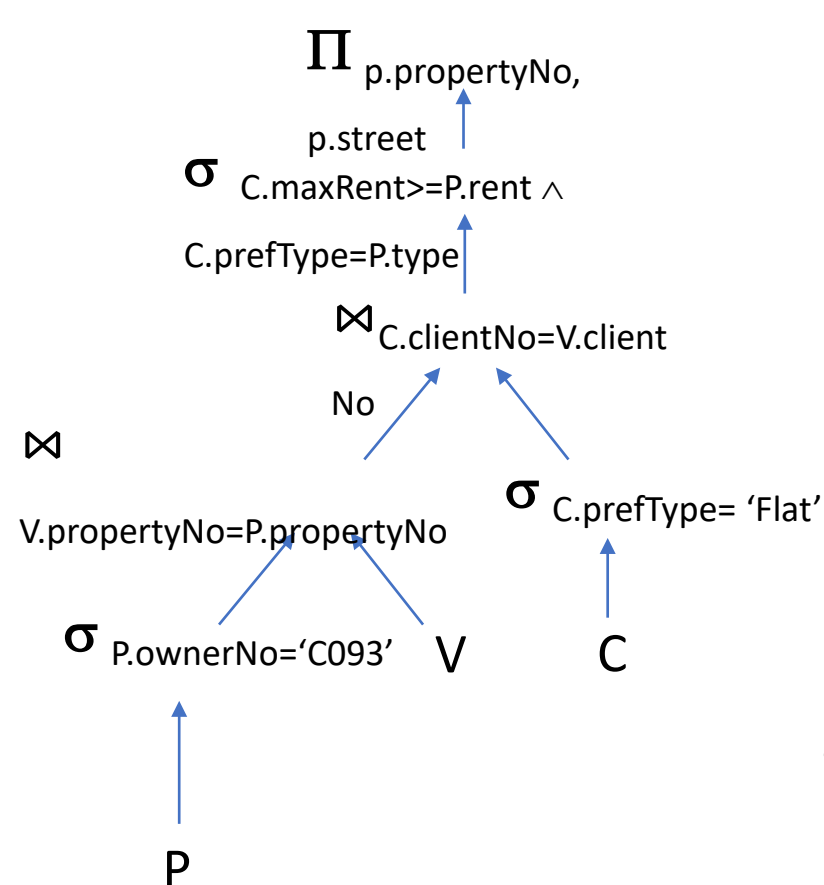
Relational algebra tree formed by pushing Selections down



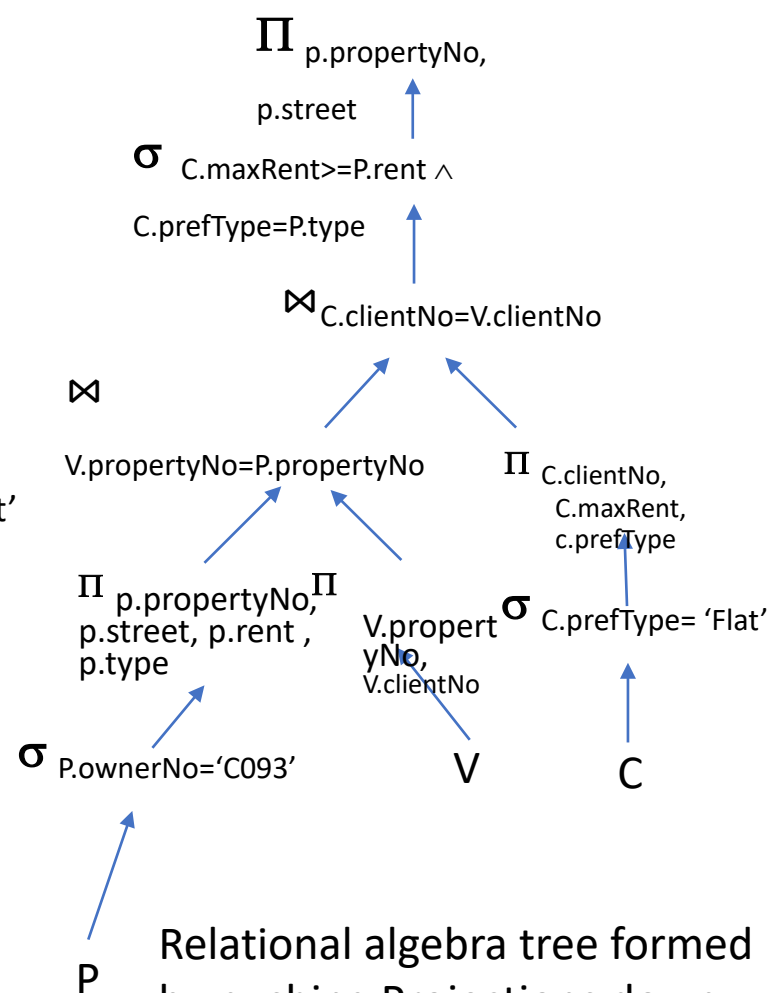
Relational algebra tree formed by changing Selection/Cartesian products to Equijoins;



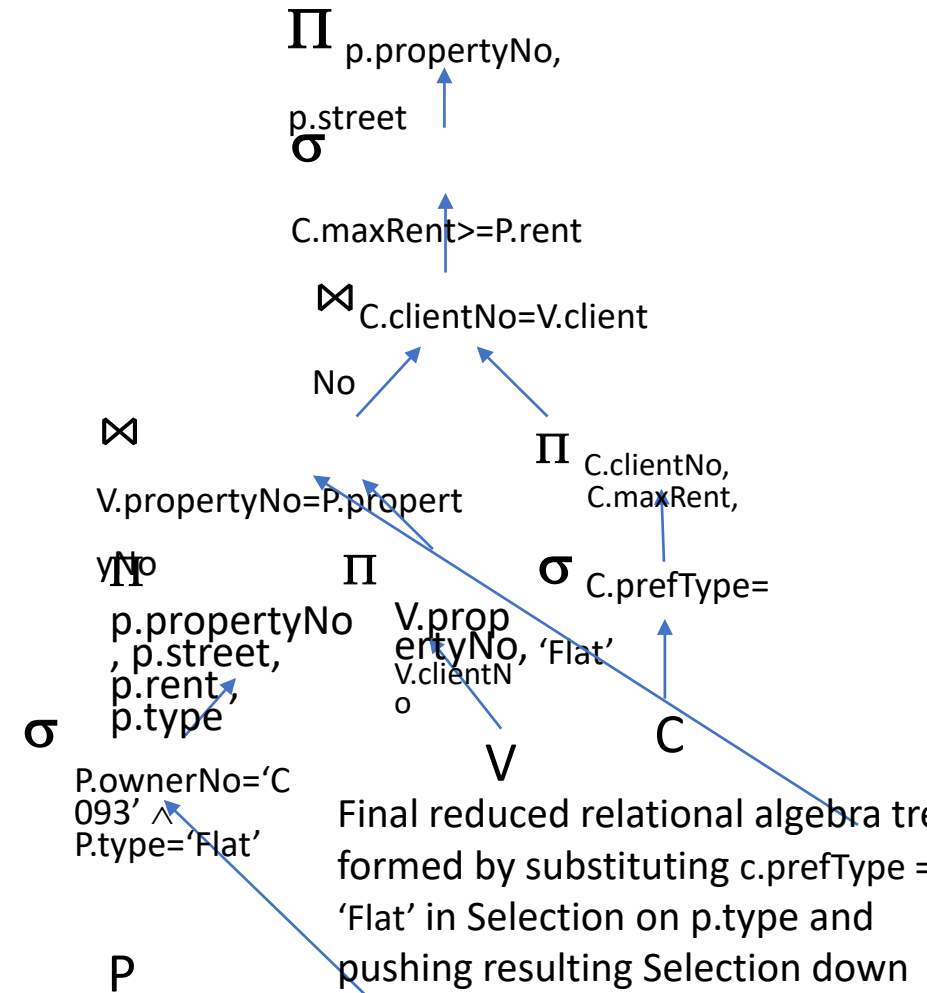
$\Pi_{p.propertyNo, p.street} (\sigma_{c.prefType = 'Flat' \wedge c.clientNo = v.clientNo \wedge v.propertyNo = p.propertyNo \wedge c.maxRent \geq p.rent \wedge c.prefType = p.type} (p \bowtie v \bowtie c))$



Relational algebra tree formed using associativity of Equijoins;



Relational algebra tree formed by pushing Projections down;



Final reduced relational algebra tree formed by substituting  $c.prefType = 'Flat'$  in Selection on  $p.type$  and pushing resulting Selection down tree.



# Heuristical Processing Strategies

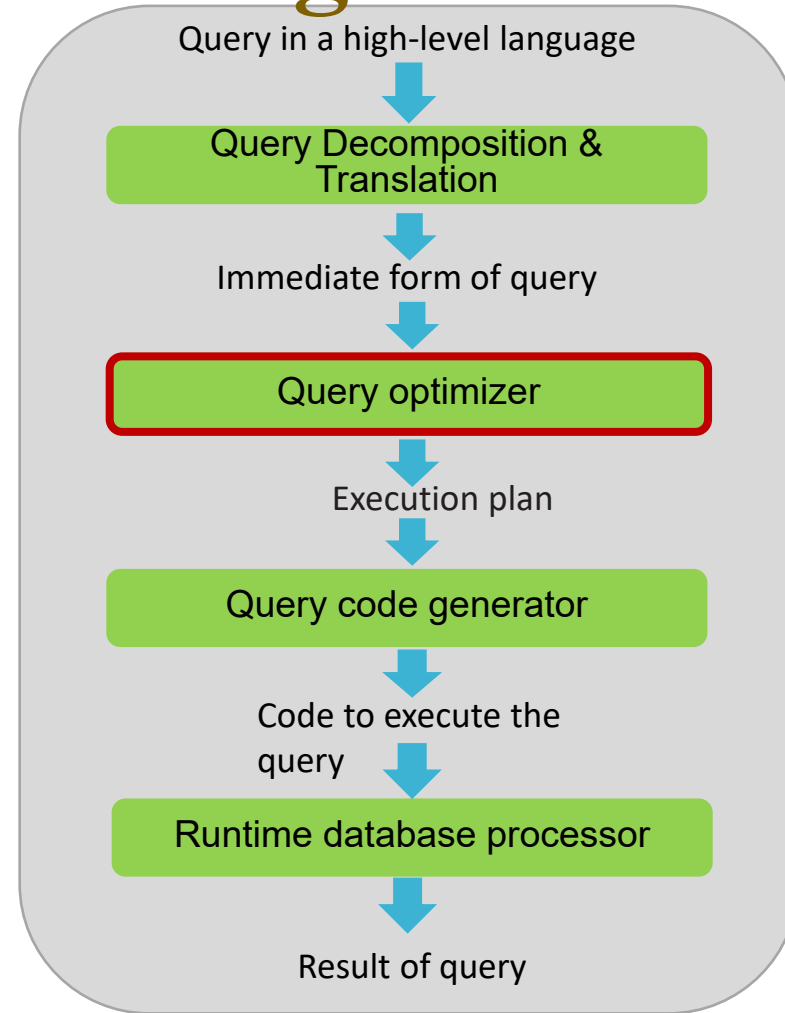
1. Perform Selection operations as early as possible.
  - ✓ Keep predicates on same relation together.
2. Combine Cartesian product with subsequent Selection whose predicate represents join condition into a Join operation.

$$\sigma_{R.a \theta S.b}(RXS) = R \bowtie_{R.a \theta S.b} (S)$$
$$(R \bowtie_{R.a \theta S.b} S) \bowtie_{S.c \theta T.d} T$$

3. Use associativity of binary operations to rearrange leaf nodes so leaf nodes with most restrictive Selection operations executed first.
4. Perform Projection as early as possible.
  - ✓ Keep projection attributes on same relation together.
5. Compute common expressions once.
  - ✓ If common expression appears more than once, and result not too large, store result and reuse it when required.
  - ✓ Useful when querying views, as same expression is used to construct view each time.



# Query Processing



## Aim:

As there are many equivalent transformations of the same high-level query, choose the one that minimizes resource usage.

There are two main techniques for query optimization.

- ✓ Heuristic rules
- ✓ **Systematically estimating**



# Cost Estimation for the Relational Algebra Operations

- Many different ways of implementing Relational Algebra (RA) operations.
- Aim of Query Optimization(QO) is to choose most efficient one.
  - Use formulae that estimate costs for a number of options and select one with lowest cost.
- Consider only cost of disk access, which is usually dominant cost in QP.
- Many estimates are based on cardinality of the relation, so need to be able to estimate this.



# Database Statistics

The success of estimating the size and cost of intermediate relational algebra operations depends on the amount and currency of the statistical information that the DBMS holds.

For each base relation R:

- ✓  $nTuples(R)$ - the number of tuples (records) in relation R (that is, its cardinality).
- ✓  $bFactor(R)$ - the blocking factor of R (that is, the number of tuples of R that fit into one block).
- ✓  $nBlocks(R)$ - the number of blocks required to store R.

For each attribute A of base relation R:

- ✓  $nDistinct_A(R)$ - the number of distinct values that appear for attribute A in relation R.
- ✓  $min_A(R)$ ,  $max_A(R)$ - the minimum and maximum possible values for the attribute A in relation R.
- ✓  $SC_A(R)$ —the selection cardinality of attribute A in relation R.

For each multilevel index I on attribute set A:

- ✓  $nLevels_A(I)$ —the number of levels in I.
- ✓  $nLfBlocks_A(I)$ —the number of leaf blocks in I.





# Selection Operation ( $S = \sigma_p(R)$ )

Predicate may be simple or composite.

Number of different implementations, depending on file structure, and whether attribute(s) involved are indexed/hashed.

Main strategies are:

1. Linear Search (Unordered file, no index):  $\lceil n\text{Blocks}(R)/2 \rceil$ , for equality condition on key attribute  $n\text{Blocks}(R)$ , otherwise
2. Binary Search (Ordered file, no index):  $\lceil \log_2 (n\text{Blocks}(R)) \rceil$ , for equality condition on ordered attribute  $\lceil \log_2 (n\text{Blocks}(R)) \rceil + \lceil \text{SC}_A(R)/b\text{Factor}(R) \rceil - 1$ , otherwise
3. Equality on hash key: 1, assuming no overflow
4. Equality condition on primary key:  $n\text{Levels}_A(I) + 1$
5. Inequality condition on primary key:  $n\text{Levels}_A(I) + \lceil n\text{Blocks}(R)/2 \rceil$
6. Equality condition on clustering (secondary) index:  $n\text{Levels}_A(I) + \lceil \text{SC}_A(R)/b\text{Factor}(R) \rceil$
7. Equality condition on a non-clustering (secondary) index:  $n\text{Levels}_A(I) + \lceil \text{SC}_A(R) \rceil$
8. Inequality condition on a secondary B+-tree index:  $n\text{Levels}_A(I) + \lceil n\text{LfBlocks}_A(I)/2 + n\text{Tuples}(R)/2 \rceil n\text{Levels}_A(I) + 1$



# Selection Operation ( $S = \sigma_p(R)$ )

Cost estimation for Selection operation:

We make the following assumptions about the Staff relation:

There is a hash index with no overflow on the primary key attribute staffNo.

There is a clustering index on the foreign key attribute branchNo.

There is a B+-tree index on the salary attribute.

The Staff relation has the following statistics stored in the system catalog:

nTuples(Staff) = 3000	→ nBlocks(Staff) = 100
bFactor(Staff) = 30	→ $SC_{branchNo}(Staff) = 6$
nDistinct <sub>branchNo</sub> (Staff) = 500	→ $SC_{position}(Staff) = 300$
nDistinct <sub>position</sub> (Staff) = 10	→ $SC_{salary}(Staff) = 6$
nDistinct <sub>salary</sub> (Staff) = 500	max <sub>salary</sub> (Staff) = 50,000
min <sub>salary</sub> (Staff) = 10,000	
nLevels <sub>branchNo</sub> (I) = 2	
nLevels <sub>salary</sub> (I) = 2	nLfBlocks <sub>salary</sub> (I) = 50

The estimated cost of a linear search on the key attribute staffNo is 50 blocks,

The cost of a linear search on a non-key attribute is 100 blocks.

Consider the following Selection operations:

S1:  $\sigma_{staffNo='SG5'}(Staff)$

Equality condition on the primary key. the attribute staffNo is hashed, estimate the cost as 1 block. The estimated cardinality of the result relation is  $SC_{staffNo}(Staff) = 1$ .

S2:  $\sigma_{position='manager'}(Staff)$

The attribute in the predicate is a non-key, non-indexed attribute, so we cannot improve on the linear search method, giving an estimated cost of 100 blocks. The estimated cardinality of the result relation is  $SC_{position}(Staff) = 300$



# Join Operation ( $T = (R \bowtie_F S)$ )

The most time-consuming operation to process.

The main strategies for implementing the Join operation.

✓ Block Nested Loop Join:

$nBlocks(R) + 1 (nBlocks(R) * nBlocks(S))$ , if buffer has only one block for R and S  
 $nBlocks(R) + 1 [nBlocks(S) * (nBlocks(R) / (nBuffer - 2))]$ , if  $(nBuffer - 2)$  blocks for R  
 $nBlocks(R) + 1 nBlocks(S)$ , if all blocks of R can be read into database buffer

✓ Indexed Nested Loop Join:

Depends on indexing method; for example:

$nBlocks(R) + 1 nTuples(R) * (nLevels_A(I) + 1)$ , if join attribute A in S is the primary key  
 $nBlocks(R) + 1 nTuples(R) * (nLevels_A(I) + 1 [SCA(R) / bFactor(R)])$ , for clustering index I on attribute A

✓ Sort-Merge Join:

$nBlocks(R) * [\log_2(nBlocks(R)) + 1] + nBlocks(S) * [\log_2(nBlocks(S))]$ , for sorts  
 $nBlocks(R) + 1 nBlocks(S)$ , for merge

✓ Hash Join:

$3(nBlocks(R) + 1 nBlocks(S))$ , if hash index is held in memory  
 $2(nBlocks(R) + 1 nBlocks(S)) * [\log_{nBuffer-1}(nBlocks(S)) + 1] + 1 nBlocks(R) + 1 nBlocks(S)$ , otherwise



# Projection Operation( $S = \Pi_{A_1, A_2, \dots, A_m}(R)$ )

To implement projection, need to:

- ✓ Remove attributes that are not required;
- ✓ Eliminate any duplicate tuples produced from previous step.

Estimating the cardinality of the Projection operation:

When the Projection contains a key attribute: the cardinality of the Projection is:  $nTuples(S) = nTuples(R)$

If the Projection consists of a single non-key attribute ( $S = \Pi_A(R)$ ), we can estimate the cardinality of the Projection as:  $nTuples(S) = SC_A(R)$

Two main approaches to eliminating duplicates:

- ✓ Sorting;
- ✓ Hashing.



# The Relational Algebra Set Operations( $T = R \cup S$ , $T = R \cap S$ , $T = R - S$ )

Implemented by

- ✓ sorting both relations on same attributes, and
- ✓ then scanning through each of sorted relations once to obtain desired result.

For all these operations, we could develop an algorithm using the sort–merge join algorithm as a basis.

The estimated cost in all cases is simply:  $nBlocks(R) + nBlocks(S) + nBlocks(R) * [\log_2 (nBlocks(R))]$   
 $+ nBlocks(S) * [\log_2 (nBlocks(S))]$



# Any Questions

