

Fundamental Operators

Let **r** and **s** be relations with **schemas R and S**

union	$r \cup s = \{ t \mid t \in r \vee t \in s \}$
difference	$r - s = \{ t \mid t \in r \wedge t \notin s \}$
cartesian_product	$r \times s = \{ t \mid t = t_r t_s \text{ where } t_r \in r \wedge t_s \in s \}$
selection	$\sigma_p(r)$
projection	$\pi_A(r)$

relational algebra

set operations



set union



set intersection



set difference



cartesian product

relational database
specific operations



selection



projection



join



set division

set functions



sum

avg

count

any

max

min

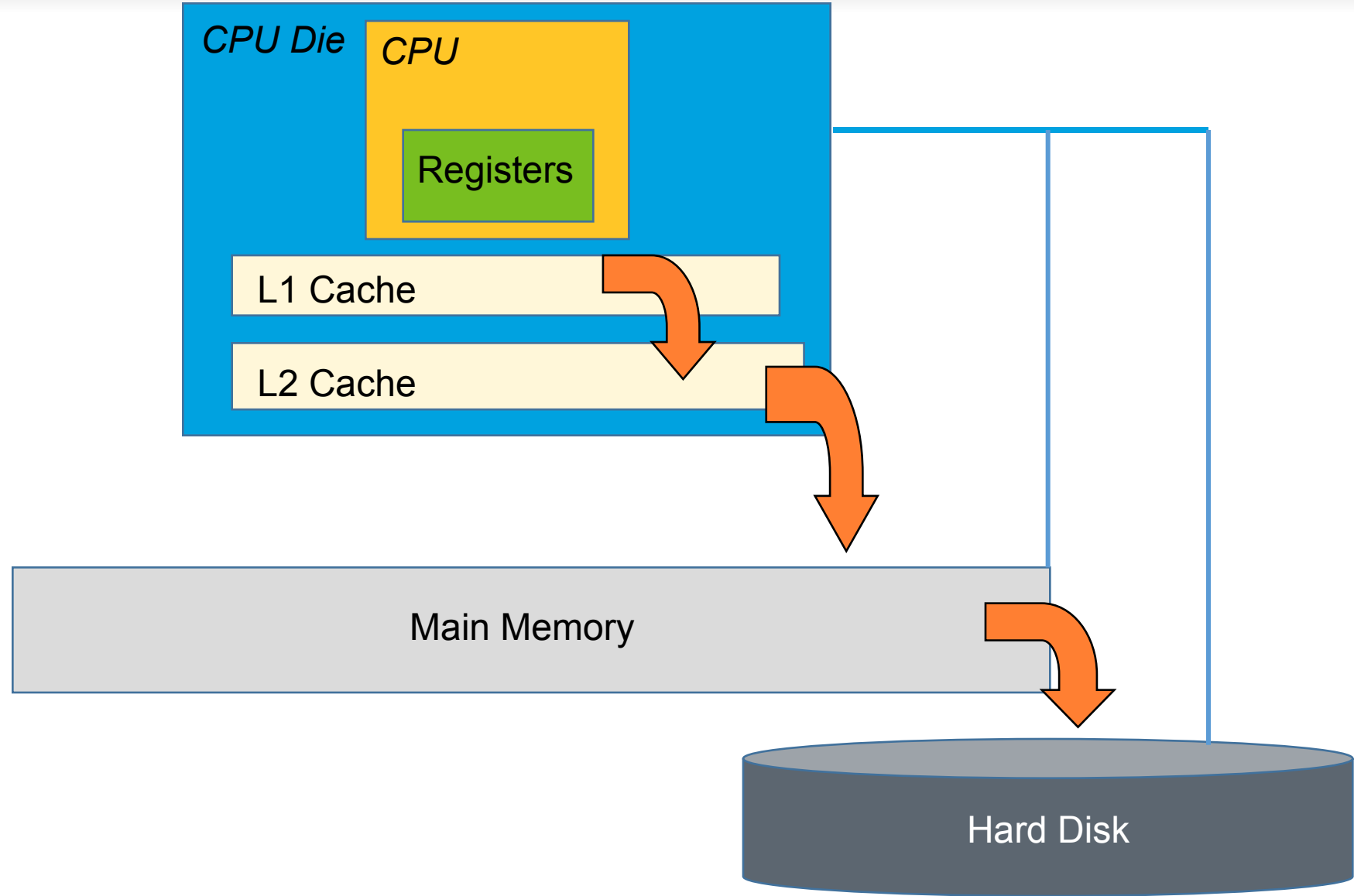
The Memory Hierarchy

Example:

Intel PIII

CPU: 450MHz

Memory: 512MB



Internal Data Storage

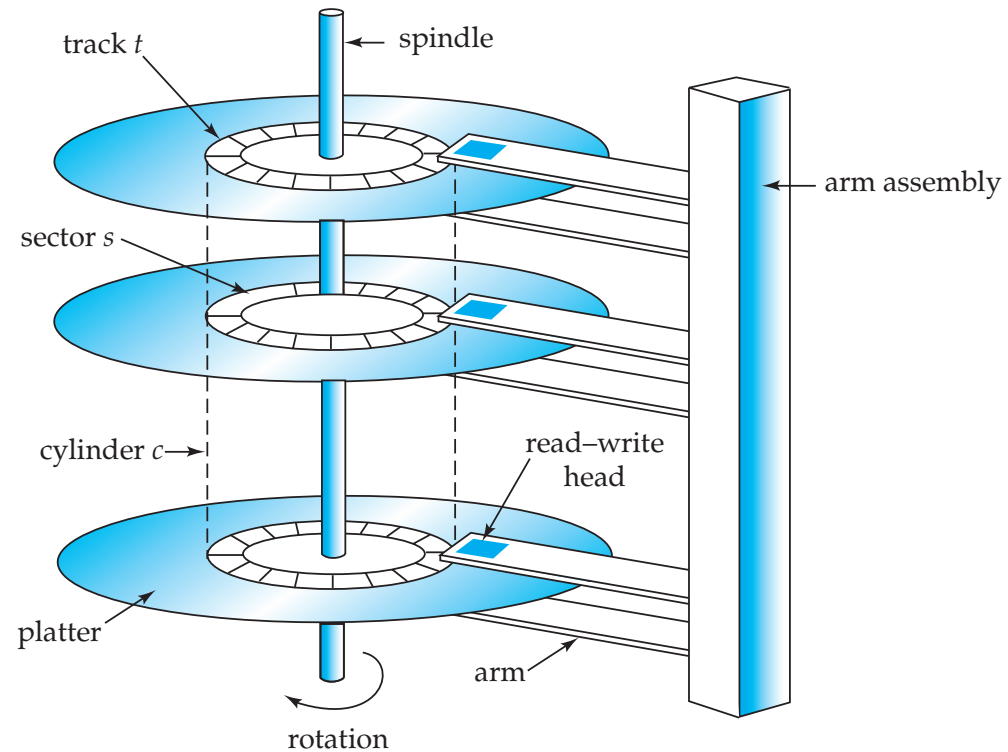
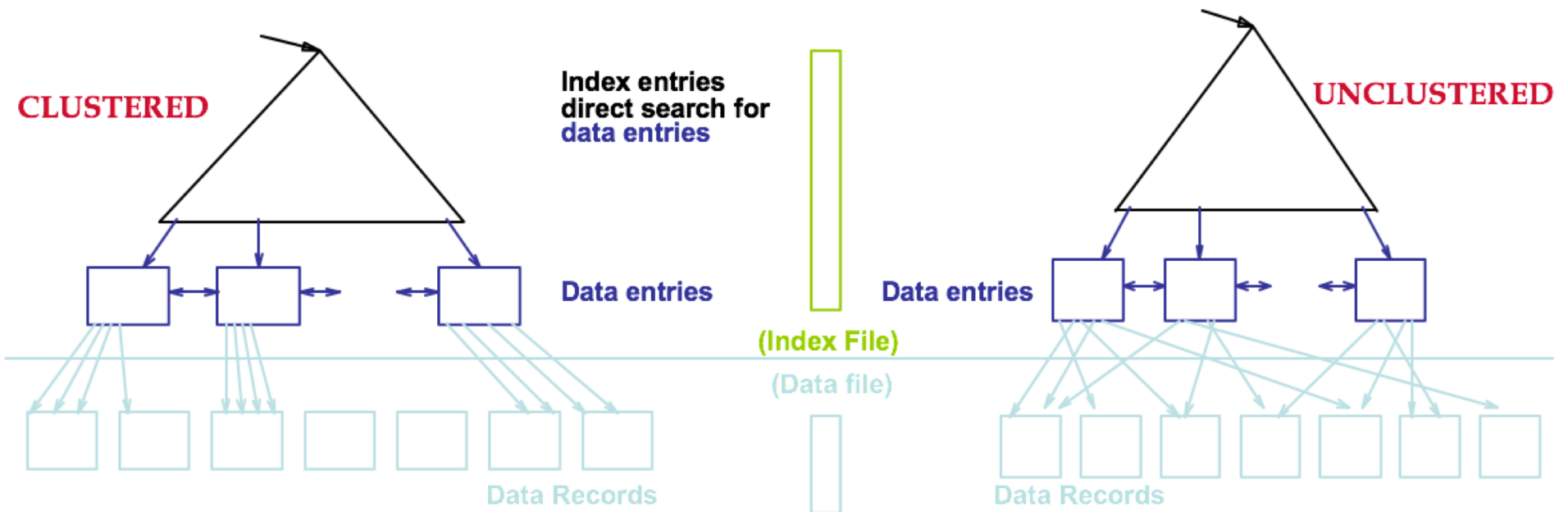
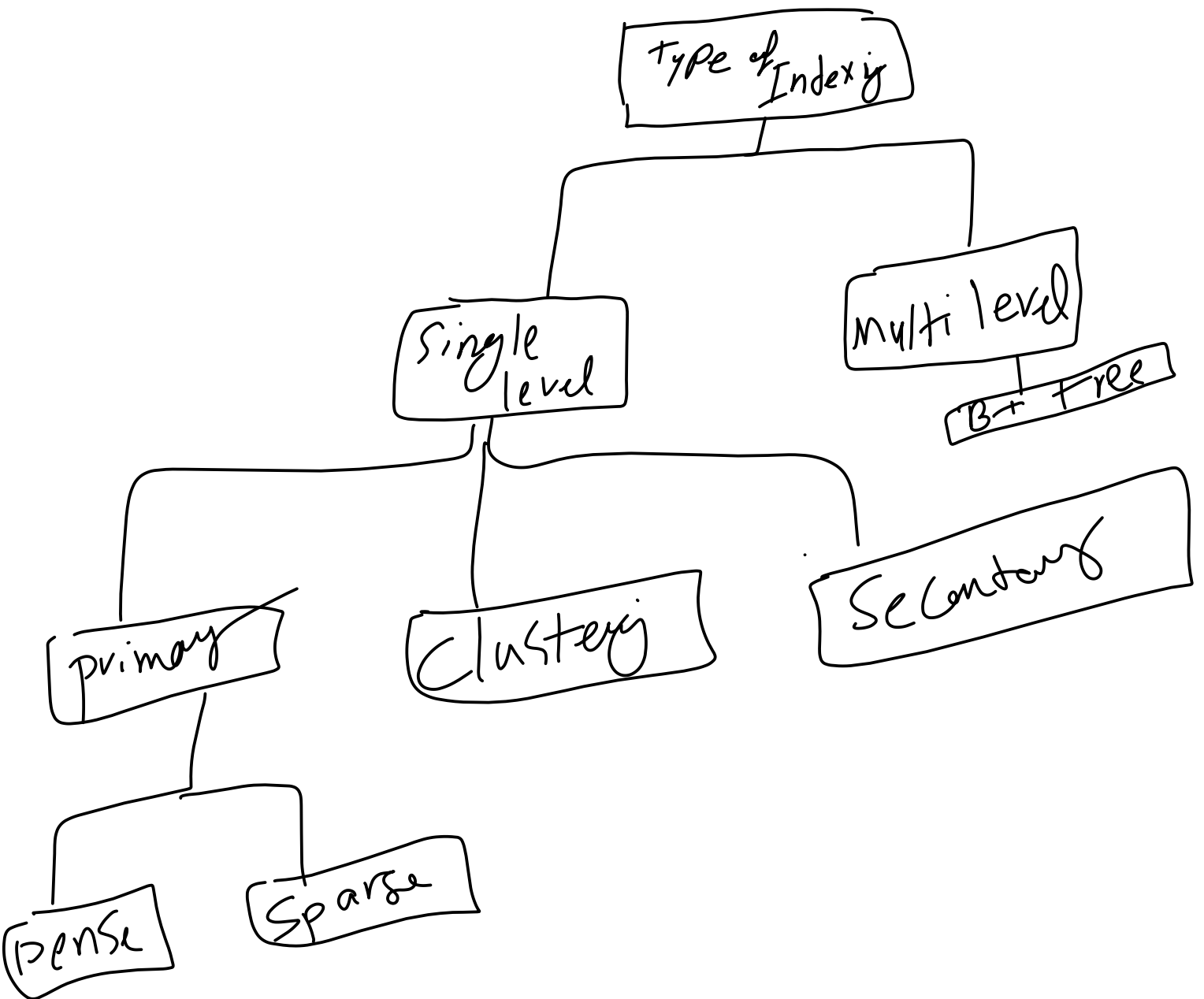


Image source: <https://www.snia.org>

Clustered vs. Unclustered Index



Types of Indexing



1-single level

a- Primary Indexing

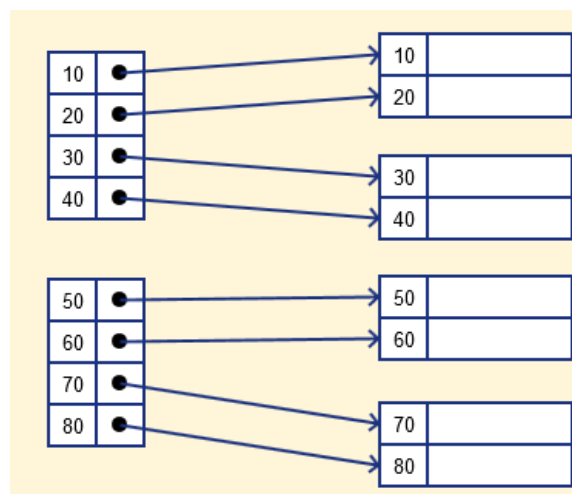
Primary Index is an ordered file which is fixed length size with two fields. The first field is the same as a primary key and second, field is pointed to that specific data block. In the primary Index, there is always one to one relationship between the entries in the index table.

The primary Indexing is also further divided into two types.

- Dense Index
- Sparse Index

a-1- Dense Index

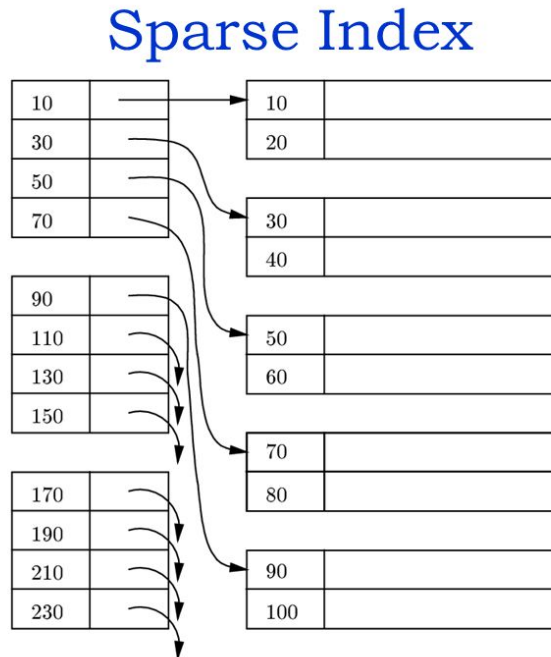
- The dense index contains an index record for every search key value in the data file. It makes searching faster.
- In this, the number of records in the index table is same as the number of records in the disk.
- It needs more space to store index record itself. The index records have the search key and a pointer to the actual record on the disk.



a-2-Sparse Index

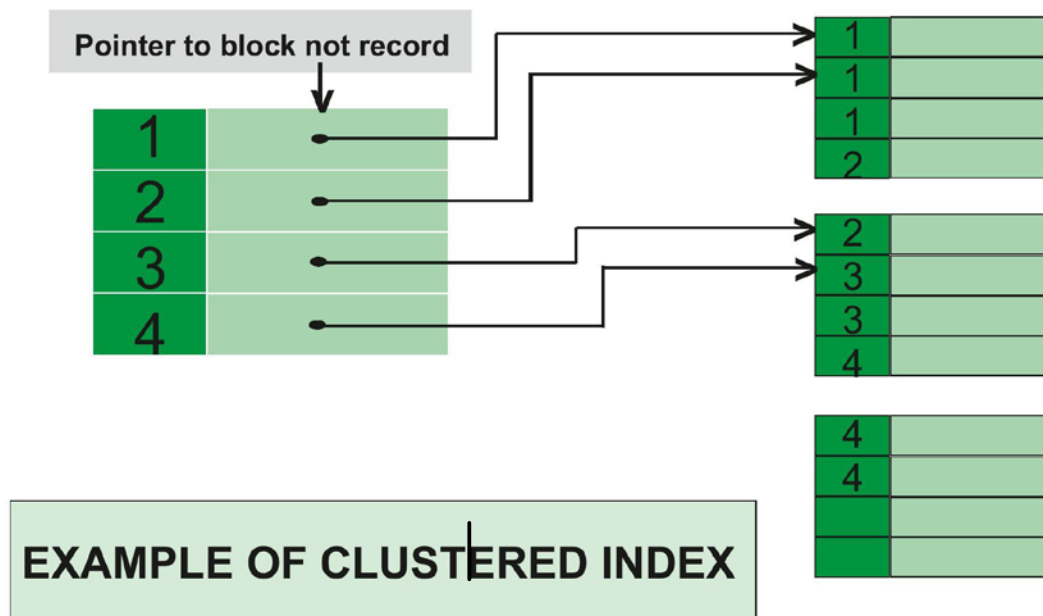
- In the data file, index record appears only for a few items. Each item points to a block.
- In this, instead of pointing to each record in the main table, the index points to the records in the main table in a gap.

Example of Sparse Index



b- Clustering Index

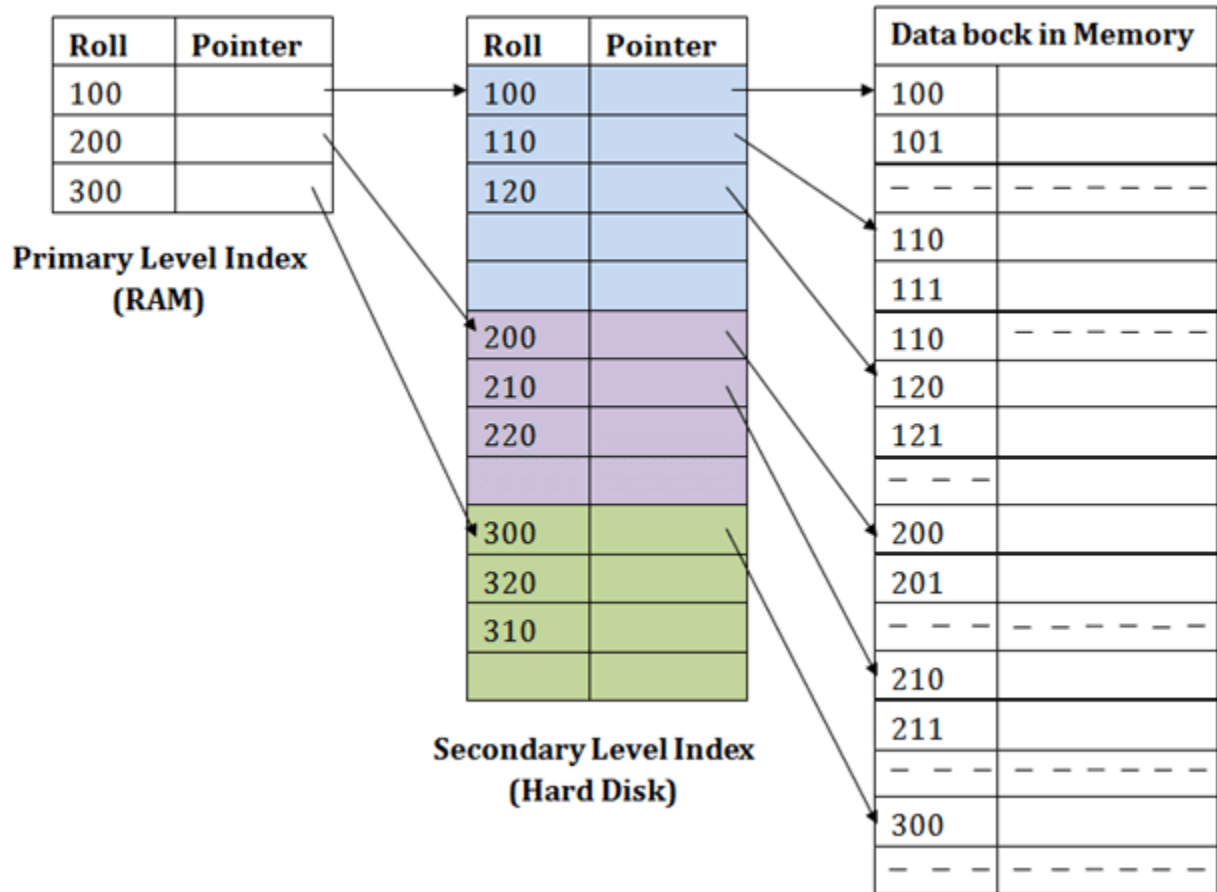
- A clustered index can be defined as an ordered data file. Sometimes the index is created on non-primary key columns which may not be unique for each record.



c- Secondary Index(non-clustering index)

In the sparse indexing, as the size of the table grows, the size of mapping also grows. These mappings are usually kept in the primary memory so that address fetch should be faster. Then the secondary memory searches the actual data based on the address got from mapping. If the mapping size grows then fetching the address itself becomes slower. In this case, the sparse index will not be efficient. To overcome this problem, secondary indexing is introduced.

In secondary indexing, to reduce the size of mapping, another level of indexing is introduced. In this method, the huge range for the columns is selected initially so that the mapping size of the first level becomes small. Then each range is further divided into smaller ranges. The mapping of the first level is stored in the primary memory, so that address fetch is faster. The mapping of the second level and actual data are stored in the secondary memory (hard disk).



2- Multi-level:

B+-tree

Transaction ACID Properties

ACID

Atomic

"ALL OR NOTHING"

Transaction cannot be subdivided

Consistent

Transaction → transform database from one consistent state to another consistent state

Isolated

Transactions execute independently of one another

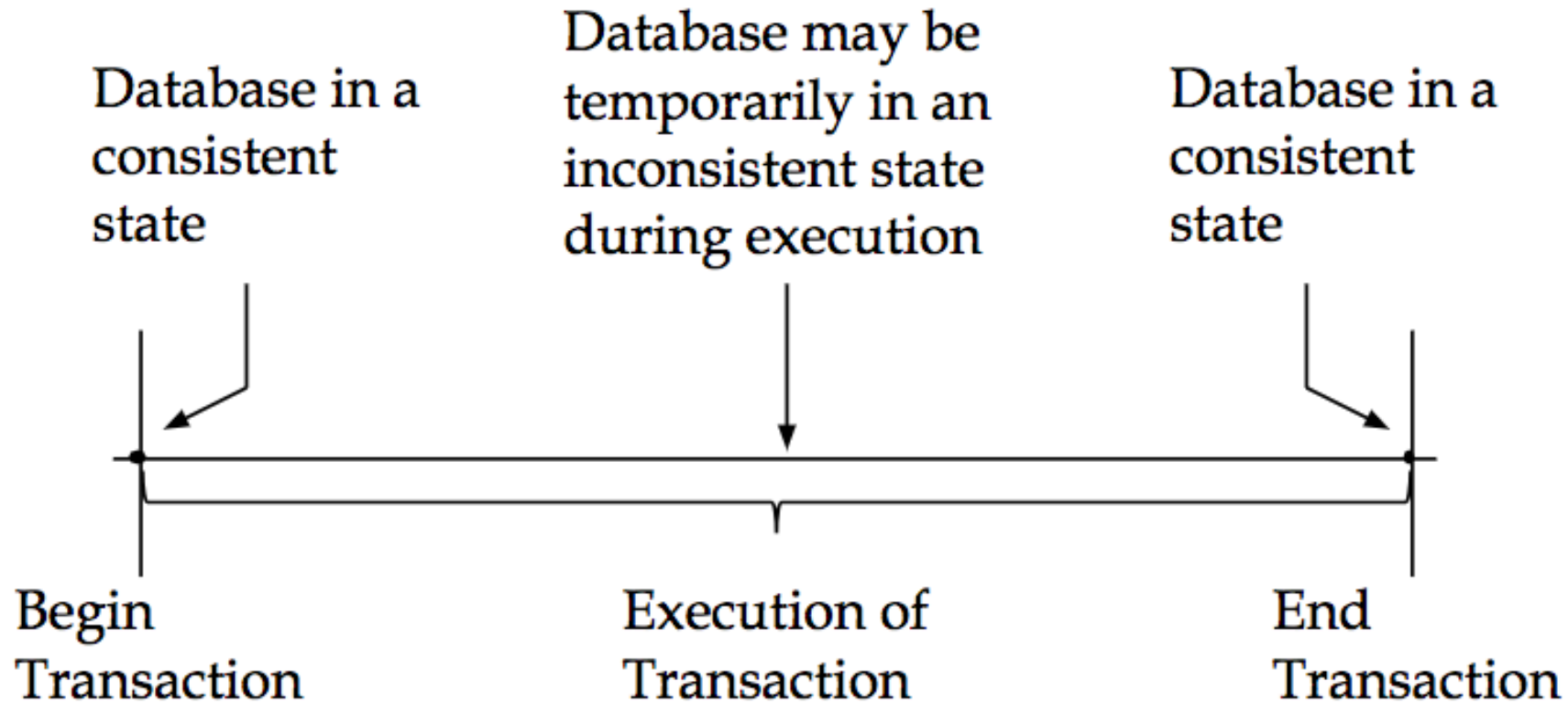
Database changes not revealed to users until after transaction has completed

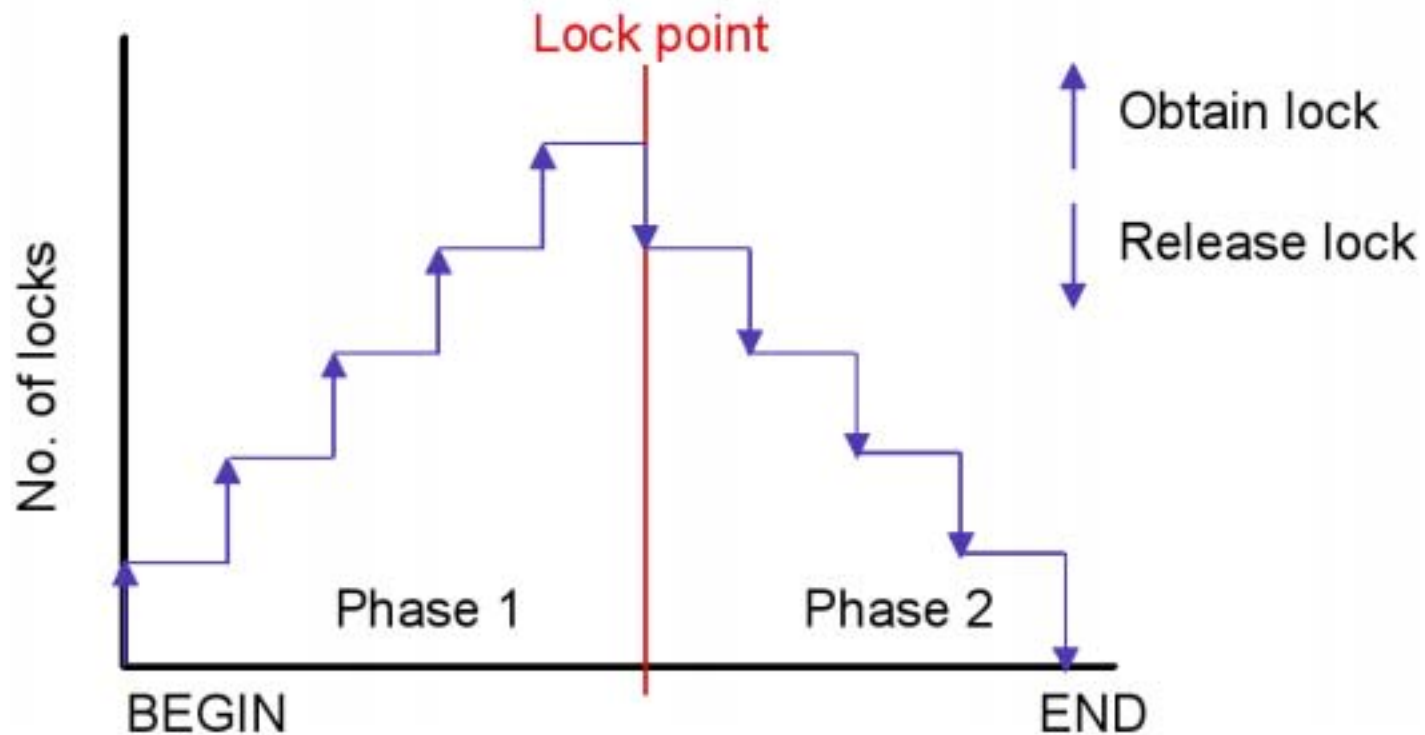
Durable

Database changes are permanent

The permanence of the database's consistent state

Transactions





Initial State



begin



Transaction

commit



roll back



**Transaction
completed**



**Transaction
failed**

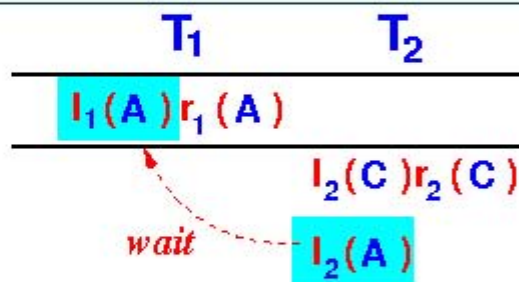
- **Wait-for graph** is a **graph** where:

- **Node** represents a **transaction**

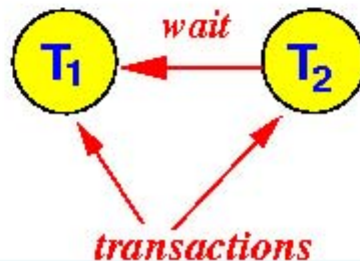
- **Edge** $i \Rightarrow j$ represents the fact that:

- The **transaction** i is **waiting** for a **lock** held by the **transaction** j

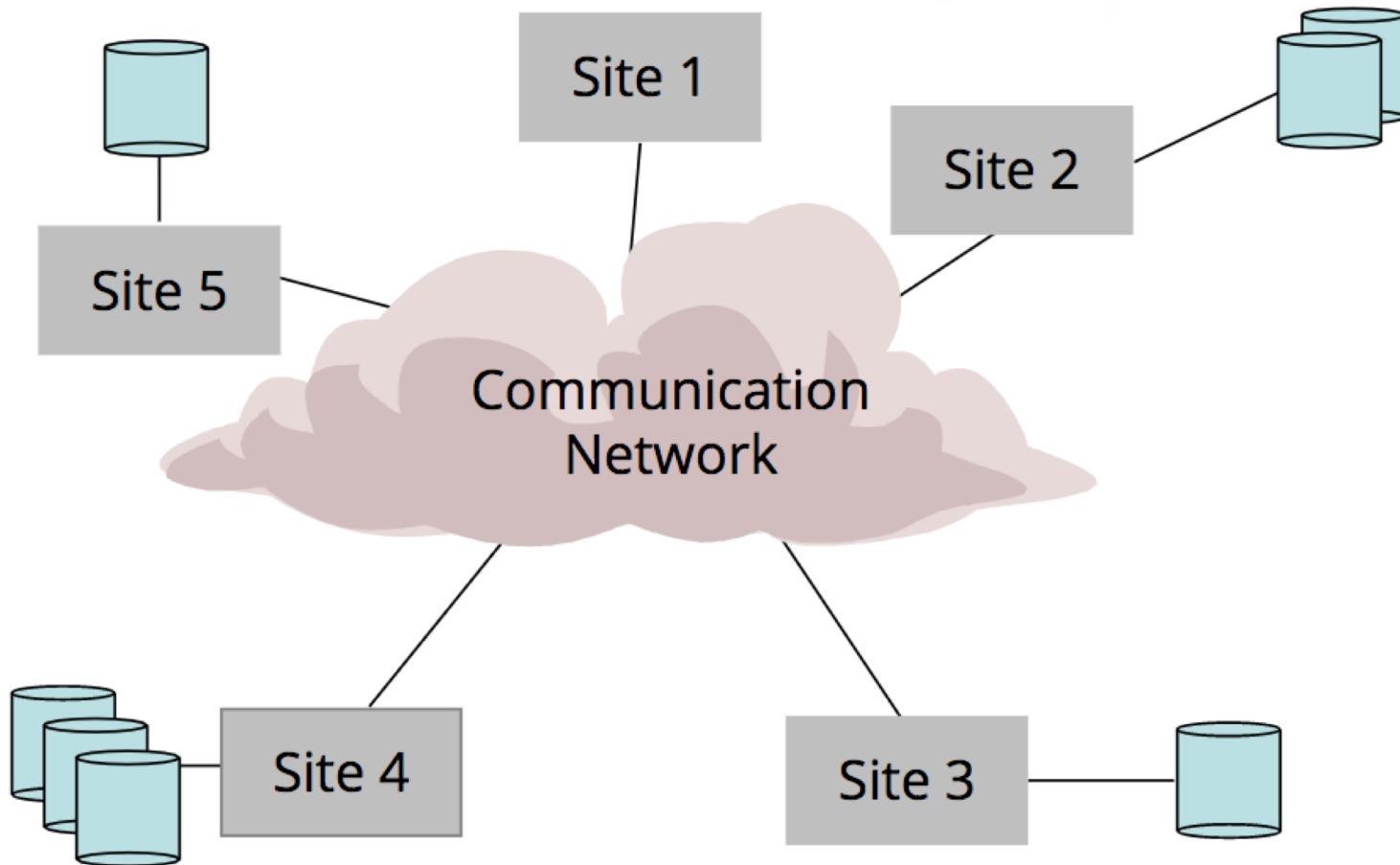
Example:



Wait-for graph:



Distributed DBMS Environment



Definition

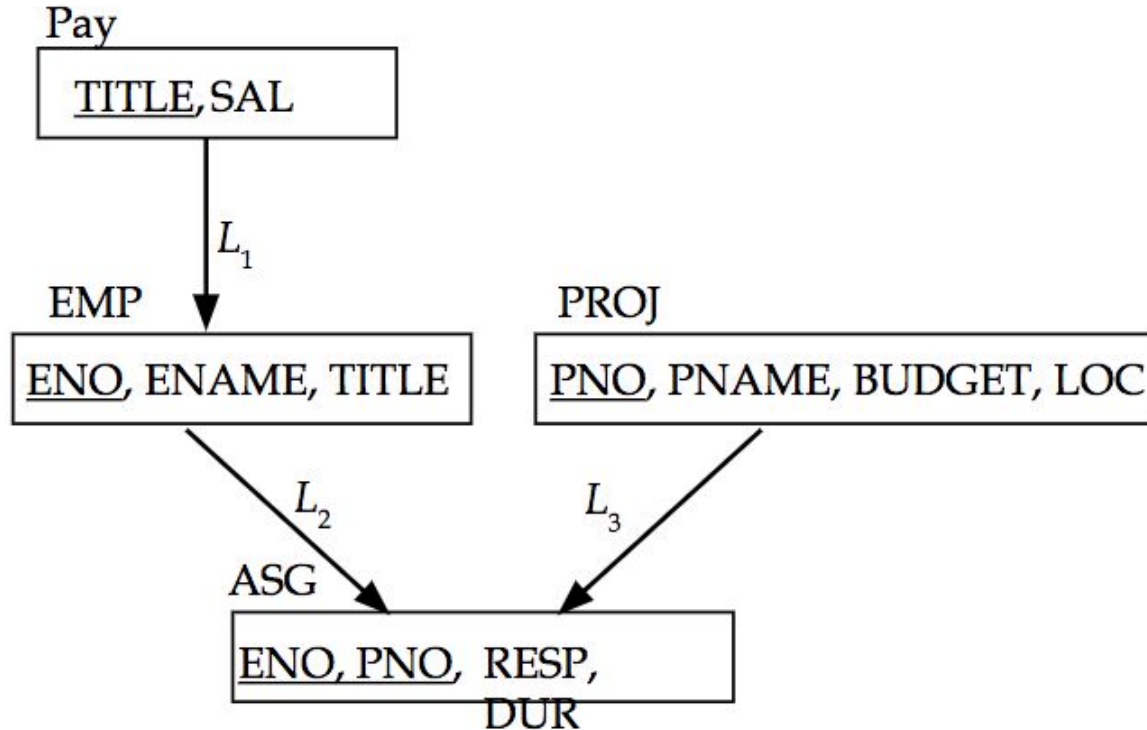
Primary Horizontal Fragmentation is:

$$R_j = \sigma_{F_j}(R), \quad 1 \leq j \leq w$$

| A horizontal fragment R_i of relation R consists of all the tuples of R which satisfy a minterm predicate m_i .

| Given a set of minterm predicates M , there are as many horizontal fragments of relation R as there are minterm predicates.

Derived Horizontal Fragmentation (DHF)



Definition

| Given a link L where $owner(L)=S$ and $member(L)=R$, the derived horizontal fragments of R are defined as

$$R_i = R \bowtie_F S_i, 1 \leq i \leq w$$

| where w is the maximum number of fragments that will be defined on R and

$$S_i = \sigma_{F_i}(S)$$

| where F_i is the formula according to which the primary horizontal fragment S_i is defined.

DHF Example

| Given link L_1 where $\text{owner}(L_1)=\text{PAY}$ and $\text{member}(L_1)=\text{EMP}$

- $\text{EMP}_1 = \text{EMP} \times \text{PAY}_1$
- $\text{EMP}_2 = \text{EMP} \times \text{PAY}_2$

| Where

- $\text{PAY}_1 = \sigma_{\text{SAL} \leq 30000}(\text{PAY})$
- $\text{PAY}_2 = \sigma_{\text{SAL} > 30000}(\text{PAY})$

EMP_1

ENO	ENAME	TITLE
E3	A. Lee	Mech. Eng.
E4	J. Miller	Programmer
E7	R. Davis	Mech. Eng.

DHF Example

| Given link L_1 where $\text{owner}(L_1)=\text{PAY}$ and $\text{member}(L_1)=\text{EMP}$

- $\text{EMP}_1 = \text{EMP} \times \text{PAY}_1$
- $\text{EMP}_2 = \text{EMP} \times \text{PAY}_2$

| Where

- $\text{PAY}_1 = \sigma_{\text{SAL} \leq 30000}(\text{PAY})$
- $\text{PAY}_2 = \sigma_{\text{SAL} > 30000}(\text{PAY})$

EMP_2

ENO	ENAME	TITLE
E1	J. Doe	Elect. Eng.
E2	M. Smith	Syst. Anal.
E5	B. Casey	Syst. Anal.
E6	L. Chu	Elect. Eng.
E8	J. Jones	Syst. Anal.

VF: Information Requirements

Application Information:

| Attribute affinities

- a measure that indicates how closely related the attributes are
- This is obtained from more primitive usage data

| Attribute usage values

- Given a set of queries $Q = \{q_1, q_2, \dots, q_q\}$ that will run on the relation: $R[A_1, A_2, \dots, A_n]$,

$$use(q_i, A_j) = \begin{cases} 1 & \text{if attribute } A_j \text{ is referenced by query } q_i \\ 0 & \text{otherwise} \end{cases}$$

$use(q_i, \bullet)$ can be defined accordingly

VF – Definition of $use(q_i, A_j)$

Consider the following 4 queries for relation PROJ:

q_1 : SELECT BUDGET
FROM PROJ
WHERE PNO=Value

q_2 : SELECT PNAME, BUDGET
FROM PROJ

q_3 : SELECT PNAME
FROM PROJ
WHERE LOC=Value

q_4 : SELECT SUM(BUDGET)
FROM PROJ
WHERE LOC=Value

	A_1	A_2	A_3	A_4
q_1	1	0	1	0
q_2	0	1	1	0
q_3	0	1	0	1
q_4	0	0	1	1

Attribute Usage Matrix

Let $A_1 = \text{PNO}$, $A_2 = \text{PNAME}$, $A_3 = \text{BUDGET}$, $A_4 = \text{LOC}$

VF – Affinity Measure $aff(A_i, A_j)$

| The **attribute affinity measure** between two attributes A_i and A_j of a relation $R[A_1, A_2, \dots, A_n]$ with respect to the set of applications $Q = (q_1, q_2, \dots, q_q)$ is defined as:

$$aff(A_i, A_j) = \sum_{\text{all queries that access } A_i \text{ and } A_j} (\text{query access})$$

$$\text{query access} = \sum_{\text{all sites}} \text{access frequency of a query} \times \frac{\text{access}}{\text{execution}}$$

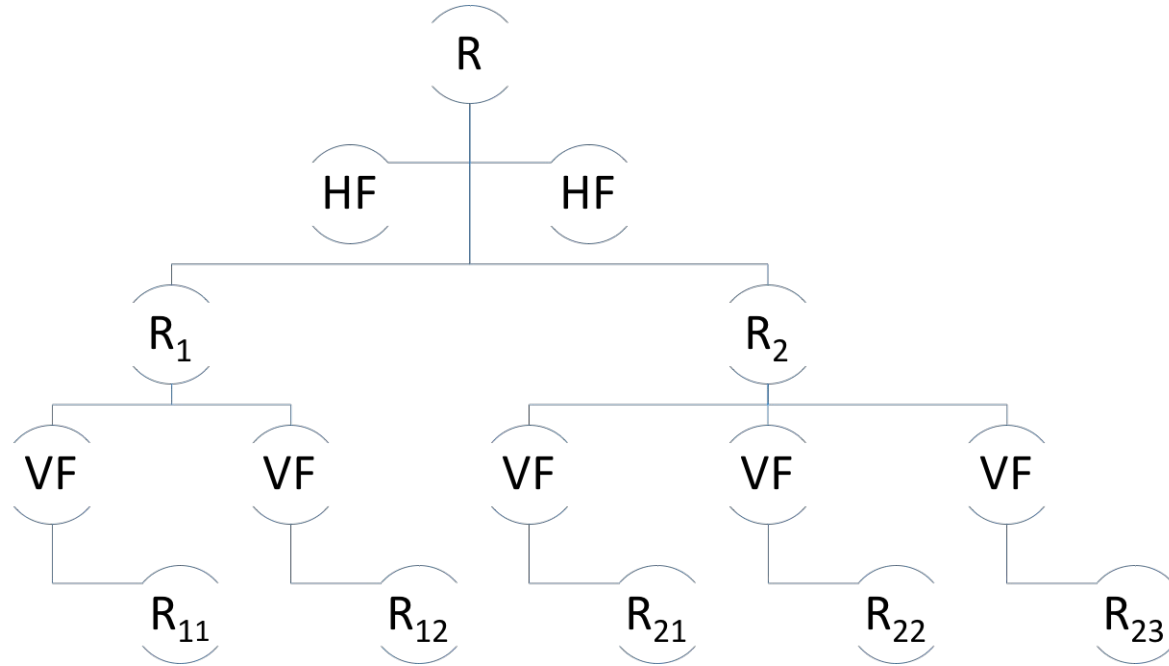
Attribute Affinity Matrix

	A1	A2	A3	A4	A5
A1					
A2	50				
A3	45	48			
A4	1	1	0		
A5	0	0	4	75	

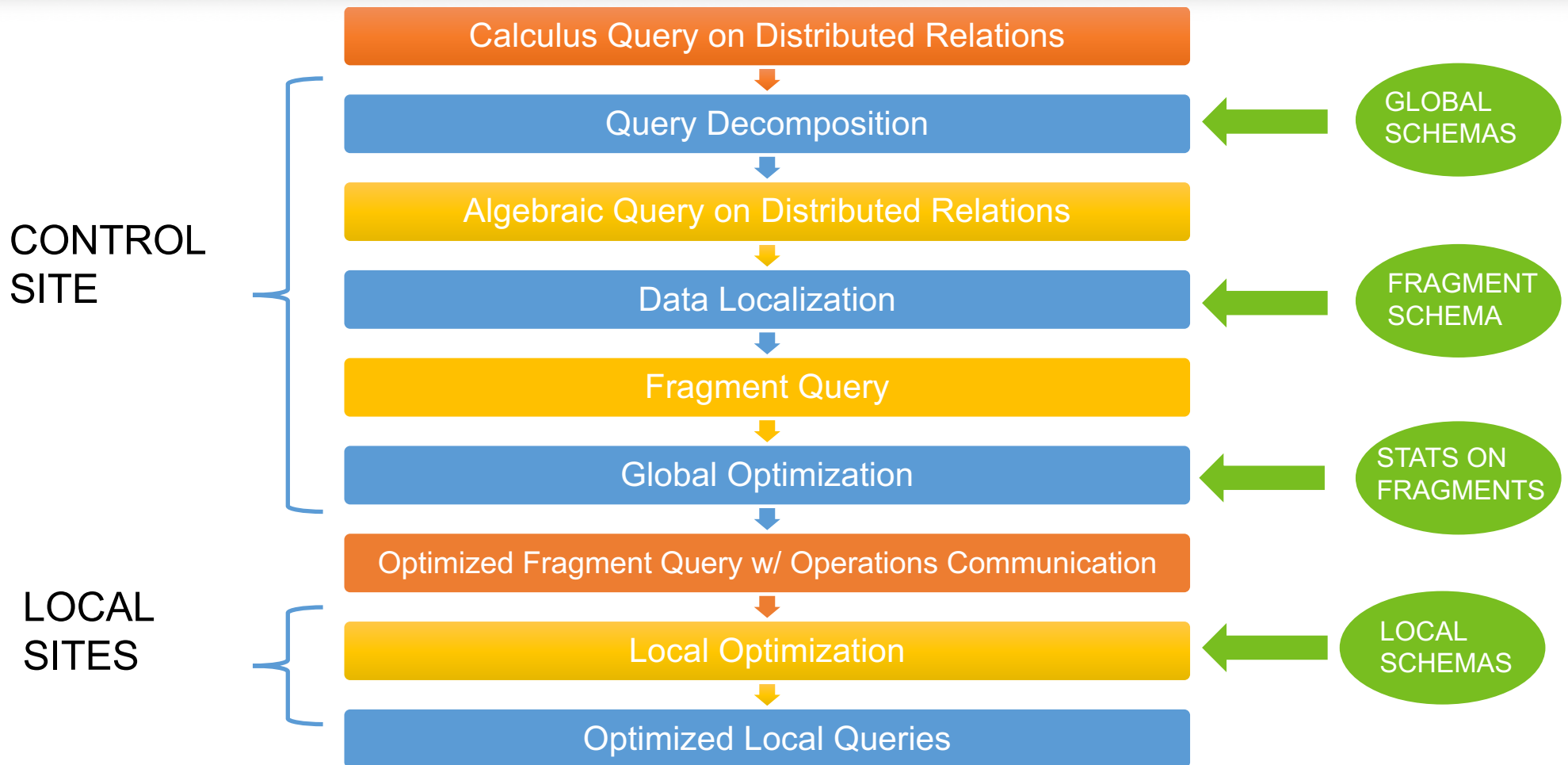
$R_1[K, A_1, A_2, A_3]$

$R_2[K, A_4, A_5]$

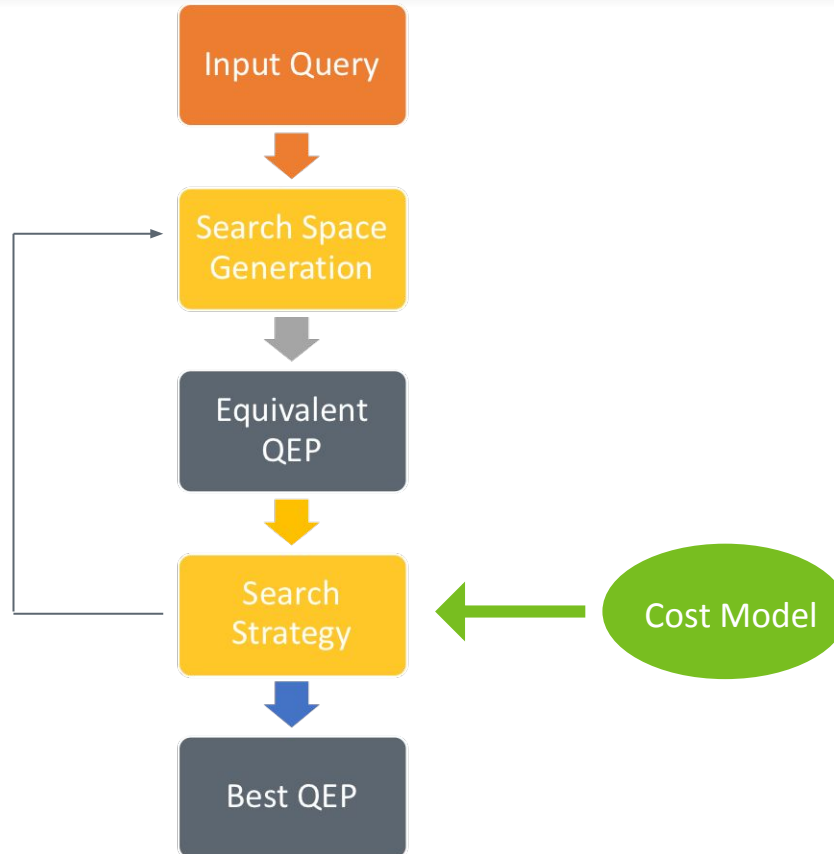
Hybrid Fragmentation



Distributed Query Processing



Query Optimization Process



Total Cost

Total cost= CPU cost + I/O cost + communication cost

CPU cost= unit instruction cost * no. of instructions

I/O cost= unit disk I/O cost * no. of disk I/Os

communication cost = message initiation + transmission

Response Time

Response time = CPU time + I/O time + communication time

CPU time = unit instruction time * no. of sequential instructions

I/O time = unit I/O time * no. of sequential I/Os

communication time = unit msg initiation time * no. of sequential msg

+ unit transmission time * no. of sequential bytes

Parallel DBMS

| Pipeline Parallelism

- many machines each doing one step in a multi-step process.



| Partition Parallelism

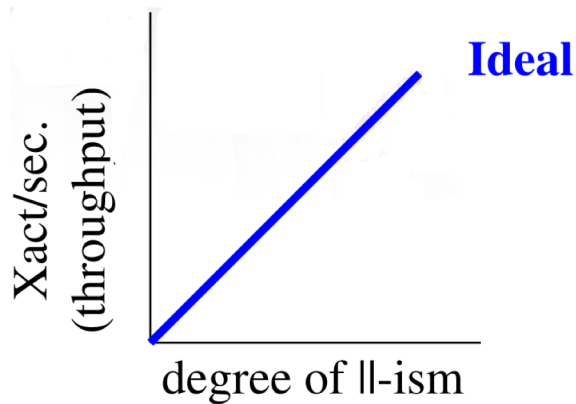
- many machines doing the same thing to different pieces of data.



Terminology

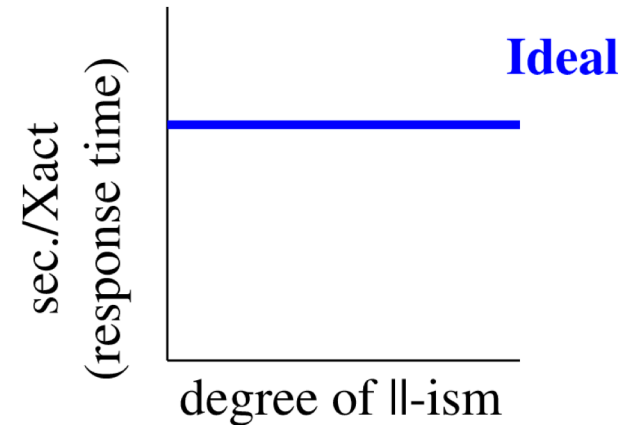
Speed-Up

- More resources means proportionally less time for given amount of data.



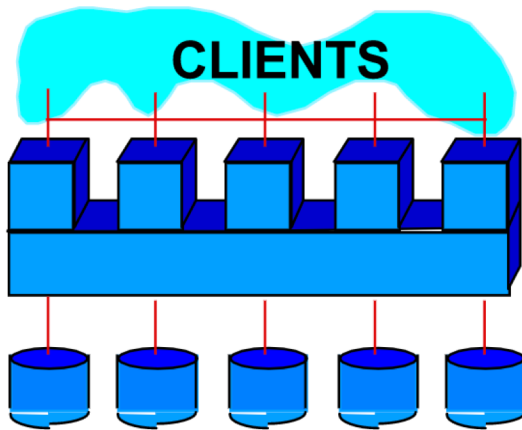
Scale-Up

- If resources increased in proportion to increase in data size, time is constant.



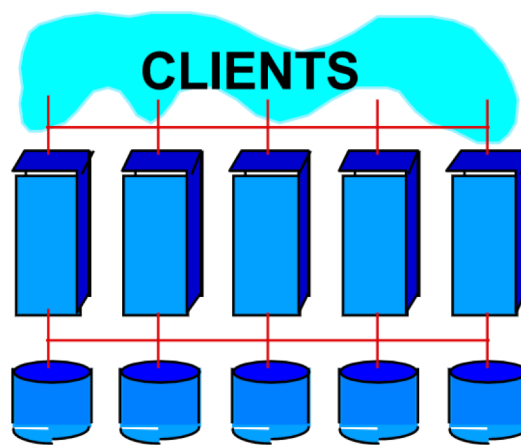
Architecture Issues: Shared What?

**Shared Memory
(SMP)**

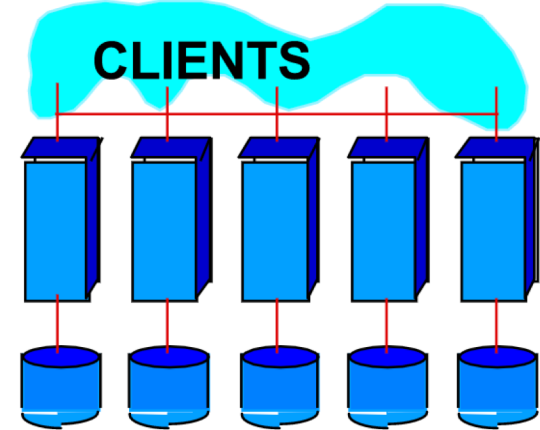


**Easy to program
Expensive to build
Difficult to scaleup**

Shared Disk



**Shared Nothing
(network)**



**Hard to program
Cheap to build
Easy to scaleup**

CAP Theorem

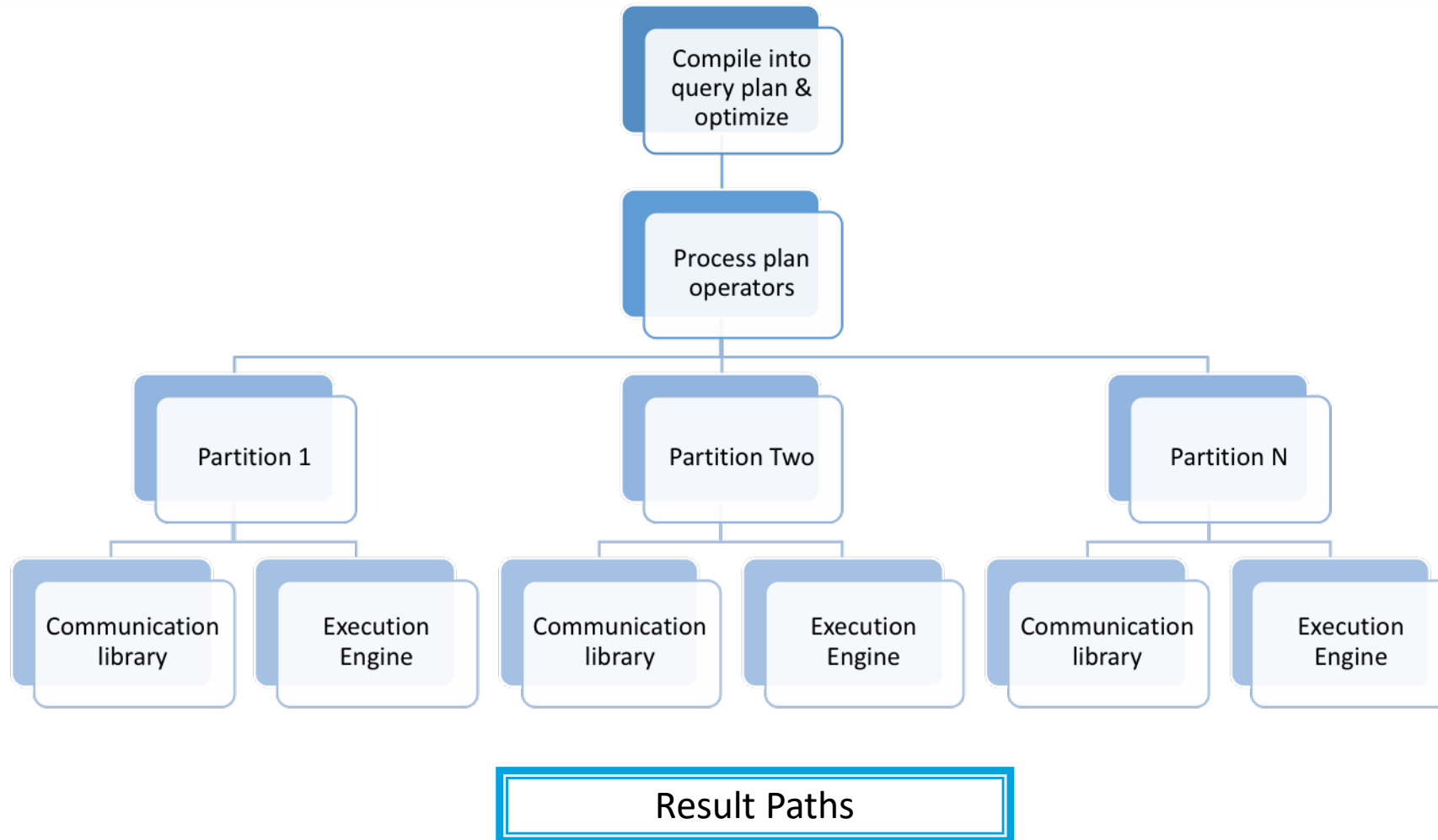
| ACID

- A DBMS is expected to support “ACID transactions,” processes that are:
 - **Atomicity**: either the whole process is done or none is
 - **Consistency**: only valid data are written
 - **Isolation**: one operation at a time
 - **Durability**: once committed, it stays that way

| CAP

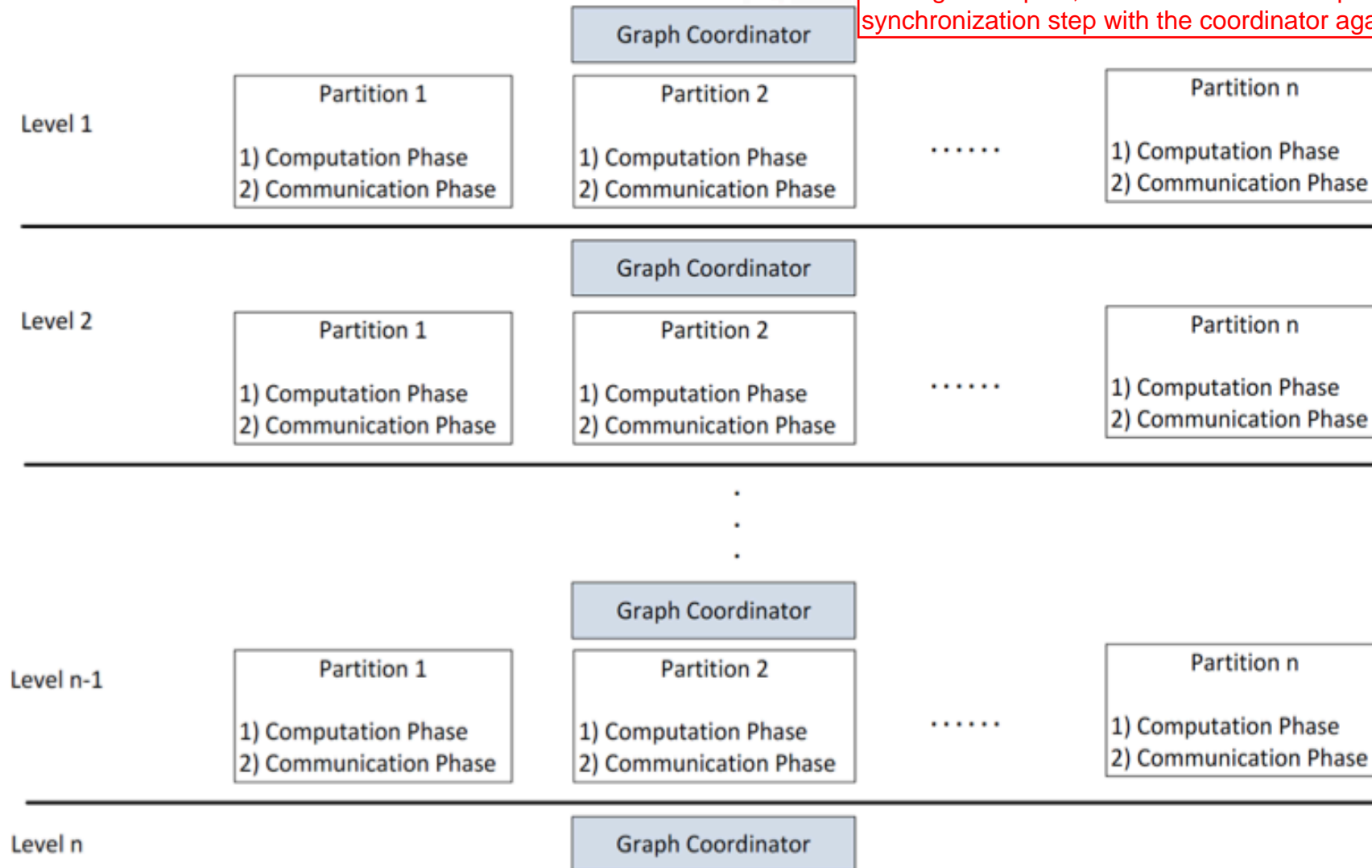
- **Consistency**: all data on cluster has the same copies
- **Availability**: cluster always accepts reads and writes
- **Partition tolerance**: guaranteed properties are maintained even when network failures prevent some machines from communicating with others

Distributed Query Execution



Distributed BFS Traversal

The graph coordinator optimizes the query, gets the finite state machine and algebraic plan, and sends it to all the partitions. There's a synchronization step with the coordinator again before we proceed.

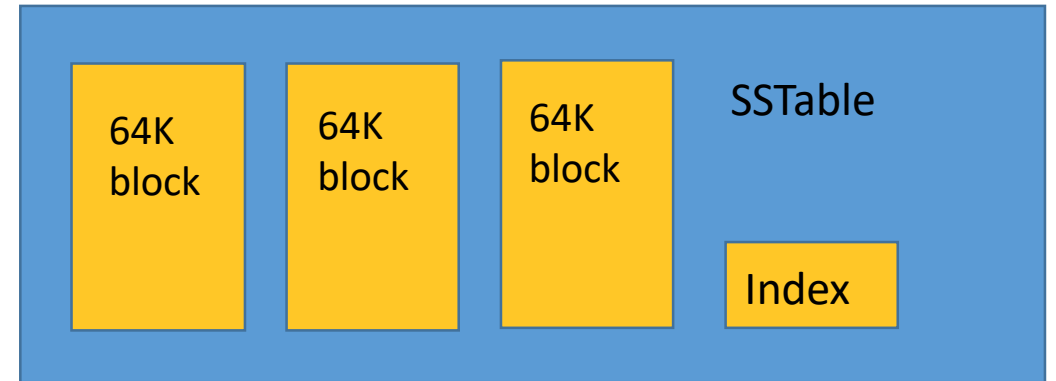


SSTable

| Immutable, sorted file of key-value pairs

| Chunks of data plus an index

- Index is of block ranges, not values

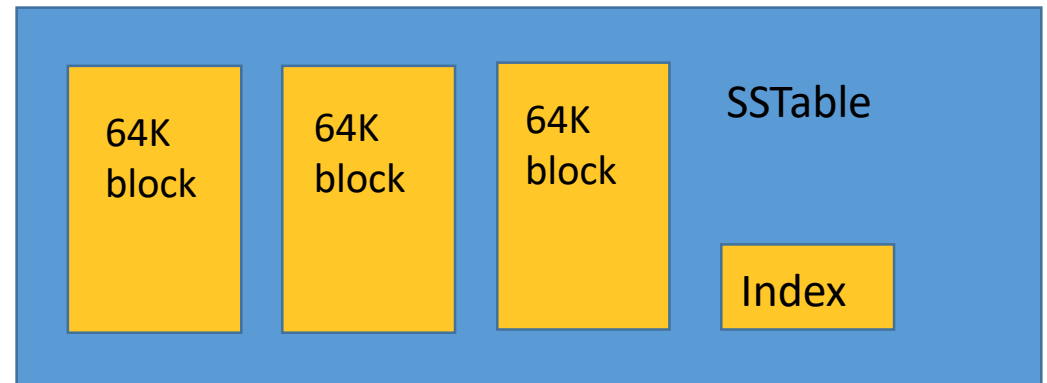
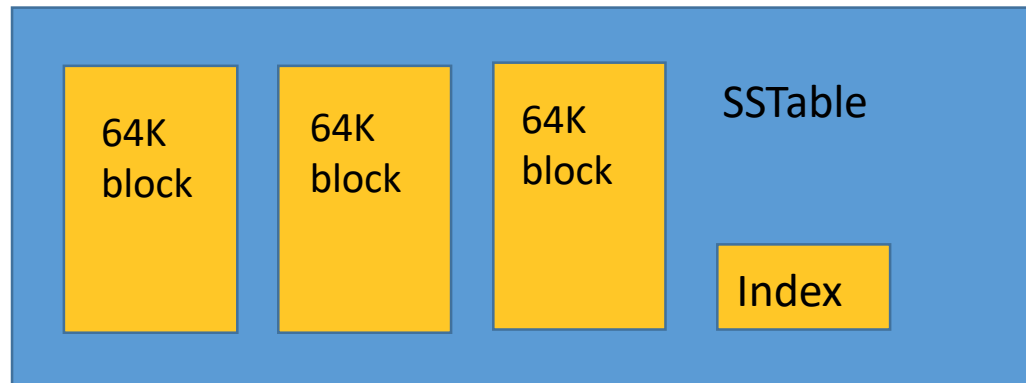


Tablet

- Contains some range of rows of the table

- Built out of multiple SSTables

Tablet Start:aardvark End:apple



Table

- | Multiple tablets make up the table
- | SSTables can be shared
- | Tablets do not overlap, SSTables can overlap

