

31/11/23

1. Database System Architectures

Types of DBMS

- a) centralized and client-server system
- b) server system architectures.
- c) Parallel Systems
- d) Distributed systems
- e) Network Types

a) centralized database system:- ^{one sys, one db} ^{one appn every thing} ^{on a single machine}

- centralized database system consists of single data server and multiple users connected to the server.
- Database is kept at single physical location.
- Different users can access this database from local machines or from remote machine.
- Management of the system and data are controlled centrally. Eg:- MS-office.

Advantages :- • operations such as insert, update, delete, backup, recovery are easier to perform.

• security in terms of authorization., Reliable (Security)

Disadvantages :- • when central computer or database system goes down, no one can access the database.

• All terminals must have connection with the central computer and this increase cost in networking.

Notes of DBMS

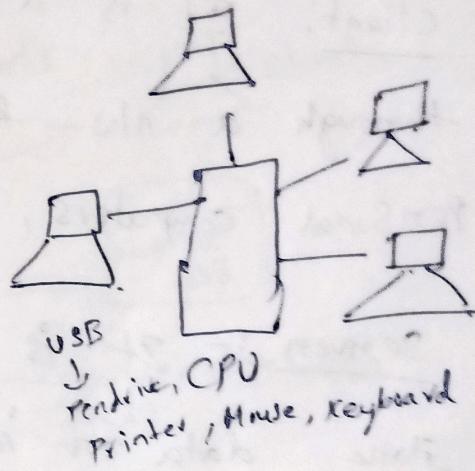
→ on the basis of Single user → single user
→ multi user

→ on the basis of location of server/site

- centralized
- distributed
- parallel
- client-server

centralized :- single resources, processor, storage devices.

- single location (I/O devices)
- support multiple users.
- one administrator.

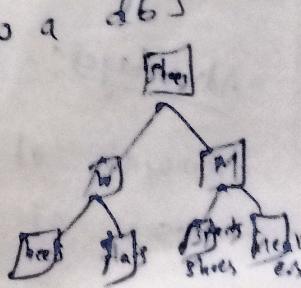


→ What is an advanced database?

a) Advanced database systems try to meet the requirements of present-day database app's by offering advanced functionality in terms of data modeling, multimedia datatype support, data integration capabilities, query languages, system features, and interfaces to other worlds.

Data model :- Data models describe how a db's

logical structure is represented.



Most advanced databases are

- MySQL
- PostgreSQL
- Oracle
- Microsoft SQL Server.

client - server systems :- (open)

1. client 2. Server

client: It is a device which request online for data through a network from a server. This may be personal computers, laptops, smartphone. (Front end)

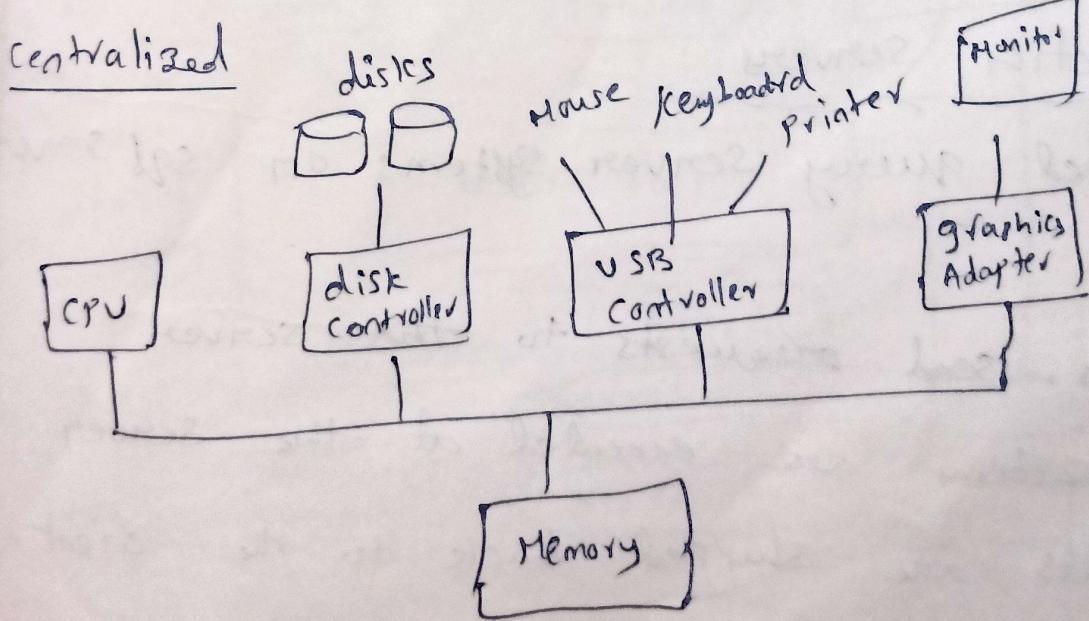
Server:- It is a large workstation which holds or store data on info and send to the client on request through a network. Servers are also called backend.

In client server db system, the DBMS app's and tools are available to all clients, but the actual database management system and physical database is available on servers.

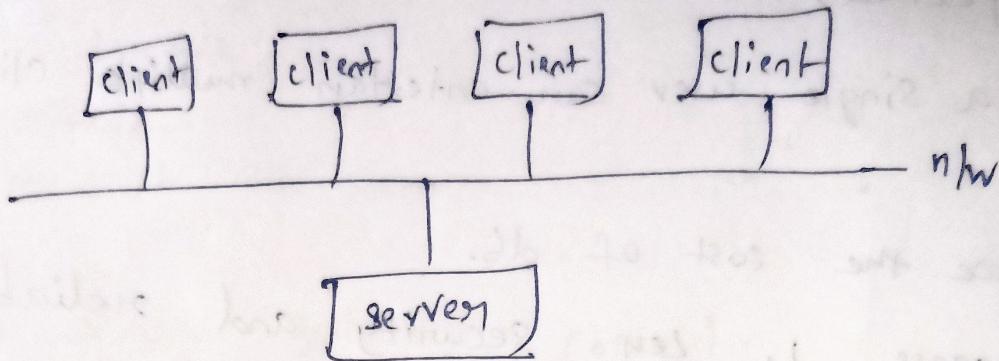
Ex: Artificial intelligence, online transaction processing
(Govt jobs)

Advantages:- 1. This system helps to manage a large amount of data to find out all possible successful uses of existing data.

2. It is more comfortable and lenient (tolerant than expected)
3. Even a single user can entertain multiple client at once.
4. Reduce the cost of db.
- disadvantages:
1. Len security and reliability
 2. The architecture can not provide complete range of systems tools to manage the n/w and associated.
-
- client1 clients client2
- client3 client4
- Back end Front end
- physical DB
- communication n/w (Cloud)
- online form ← Ex



client - server :-



→ 1/2/23

Server system: Architecture.

Server systems can be broadly categorized into two kinds

1. Transaction servers: which are widely used in relational database systems
2. data servers: used in object-oriented db sys

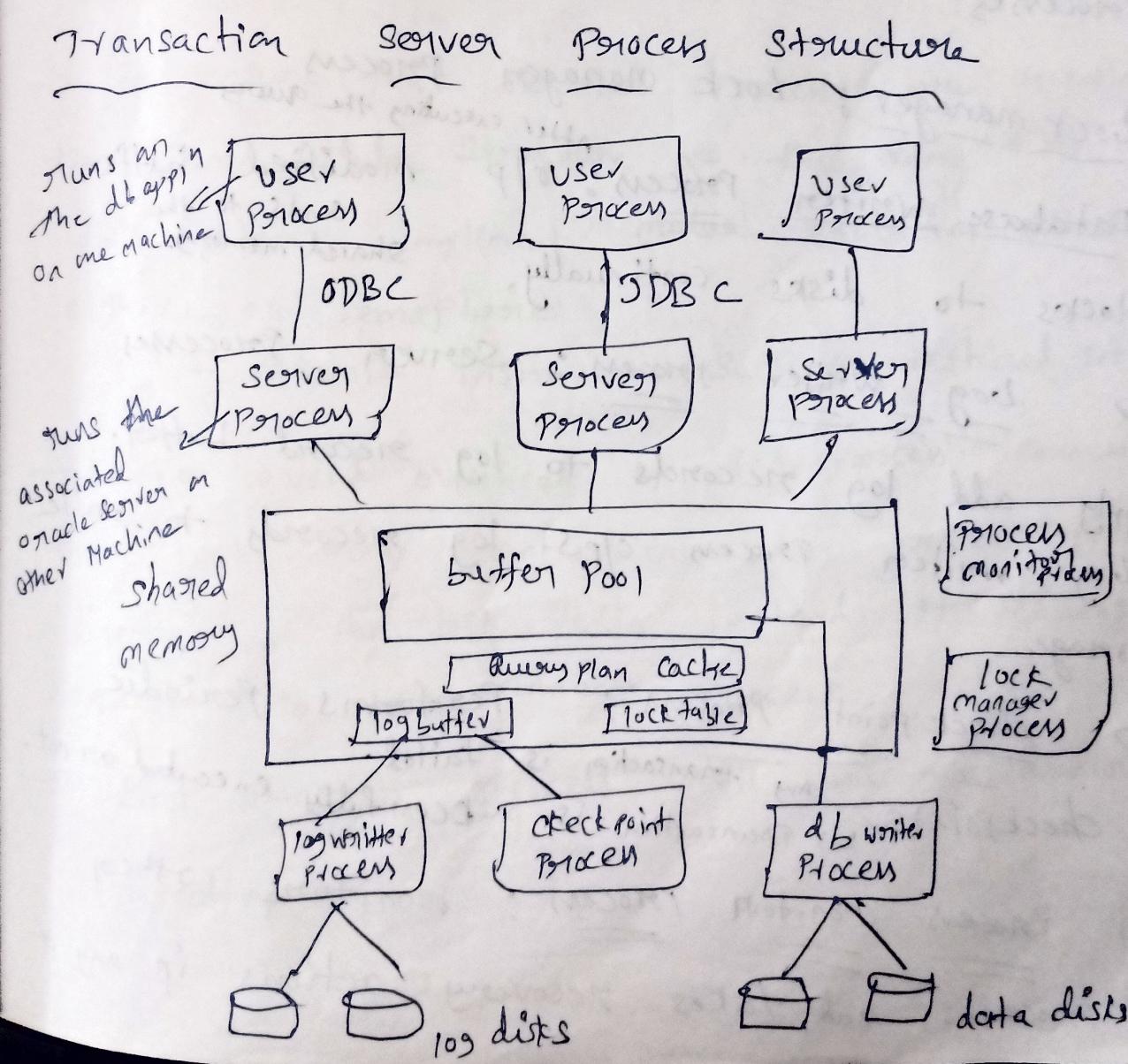
Transaction servers

Also called query server systems or sql server systems

- clients send requests to the server
- Transactions are executed at the server
- Results are shipped back to the client.

- Requests are specified in SQL and communicated to the server through a remote procedure call (RPC) mechanism.
- Transactional RPC allows RPC calls to form a transaction.
- Open database connectivity (ODBC) is a C language application program interface standard from Microsoft for connecting to a server, sending SQL requests and receiving results.

→ JDBC



- A typical transaction Server consists of multiple processes accessing data in shared memory.
- Server Processes: They receive user queries (transactions), execute them and send results back.
- Processes may be multithreaded, allowing single process to execute several user queries concurrently.
- Typically multiple multithreaded server processes.

Lock manager: Lock manager processes after executing the query

Database writer Process: O/P modified buffer sends to the shared memory.

blocks to disks continually.

- Log writer Process: Server processes simply add log records to log record buffer. Log writer processes o/p's log records to stable storage.
- check point Process: Performs periodic checkpoints. Any transaction is failed. Transaction is successfully executed or not
- Process monitor Process: monitors other processes, and takes recovery actions, if any of

other processes fail.

Eg: aborting any transactions being executed by a server process and restarting it.

Shared memory: contains shared data

→ Buffer Pool -

→ Lock table - it is a table who are granting and releasing all the locks which are available into the lock table

→ Log buffer

→ cached query plans (reused if same query submitted again)
L same data

All db processes can access shared memory

→ To ensure that no two processes are accessing the same data structure at the same time,

db's systems implement mutual exclusion using either

→ OS semaphores.

→ Atomic instructions such as test and set

→ To avoid overhead of interprocess communication for lock request, each db process operates directly on the lock table instead of sending requests to lock manager process.

→ Lock manager process still used for deadlock detection.

What is Database Architecture

A database Architecture is a representation of DBMS design. It helps to design, develop, implement, and maintain the db management System.

→ lock manager Process is responsible for granting and releasing the locks. for avoiding the Concurrency control.

Data Servers :- used in object-oriented db system. because now-a-days instead of using the relational db system.

→ used in high-speed LANs. in cases where the clients are comparable in processing power to the server. the tasks to be executed are compute intensive.

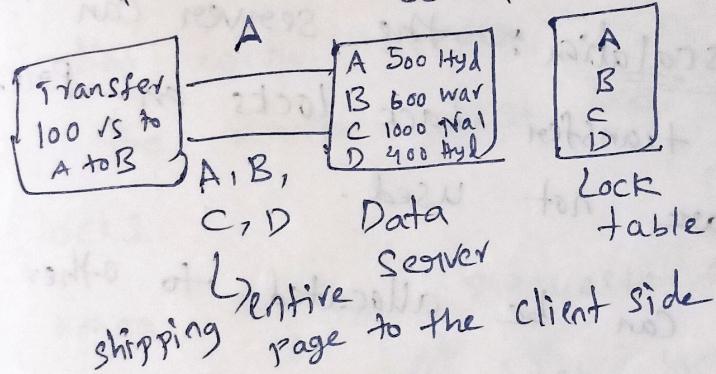
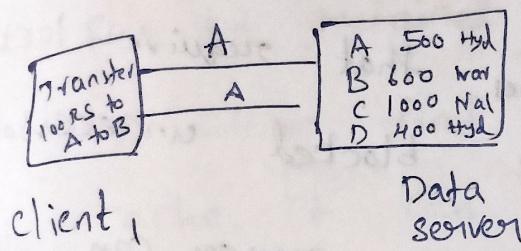
** In most of the time server will be busy so we can execute the query in client side also (transaction capable)

→ client can ask for data at server side

- The tasks to be executed are compute intensive.
- This architecture requires full back-end functionality at the clients.

→ issues:-

1) - Page-shipping versus Item-shipping (how much data shifted to the server side) - entire page shifted to client side.



- Fetching items even before they are requested is called prefetching.

↳ page shipping can be thought of as a form of prefetching.

- Unit of communication for data can be of greater granularity : page
- Coarse granularity : page
- Fine granularity : Tuple (referred as item)
- ↳ smaller (1 tuple or 1 record)
- Item shipping
- overhead of message passing is high compared to data transmit.

→ Instead, when a particular item is requested, other items are also sent which may be used in future.

Adaptive lock Granularity :-

→ locks are usually granted by the server for the data items that it ships to the client MC's.

Disadvantage of Page Shifting

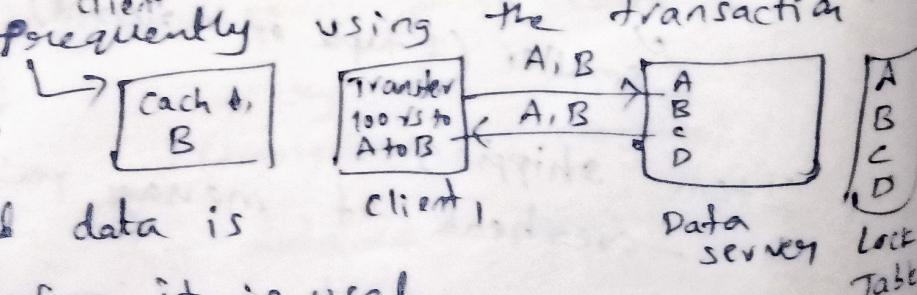
- client machines may be granted lock on a page.
- Locks all items contained in the page, even if the client is not accessing some items in the page.
- other client machines that require locks on those items may be blocked unnecessarily.
- Lock de-escalation: the server can request its clients to transfer back locks on prefetched items that are not used.
- Those locks can be allocated to other clients.

b) Data Caching:-

→ Data that are shipped to the client can be cached at client even after completing the transaction. (frequently using the transaction ^{client} _{A,B})

→ Cache Coherency:
check that the cached data is up-to-date before it is used

→ ensure that the data is not updated by diff client after caching.



→ Message must still be exchanged with the server to check its validity and to acquire lock on the data again.

Lock caching

- If the use of data is partitioned among the clients, with clients rarely requesting data that are also requested by other clients.
- Locks can be retained by client machine.
- If the client finds data and lock permission in the cache, it can access the data without communicating with server.
- But the server must keep track of all cached locks.
- If other client requests a lock on data item, the server calls back locks from clients when it receives conflicting lock request.
- Client returns lock once no local transaction is using it.
- Similar to desescalation, but across transactions.

Parallel Systems: Parallel db systems consists of multiple processors and multiple disks connected by a fast interconnection n/w.

→ A Coarse-grain Parallel machine consists of a small of Powerful Processors.

→ A massively parallel or fine grain Parallel machine utilizes thousands of smaller processor.

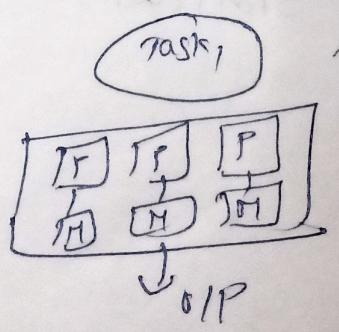
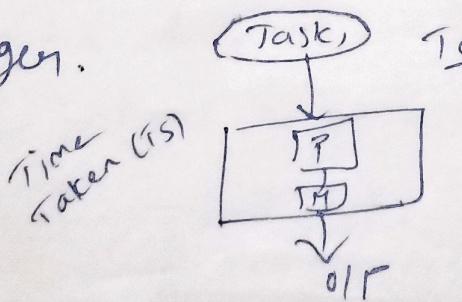
⇒ Two main Performance measures

throughput: the number of tasks that can be completed in a given time interval

response time: the amount of time it takes to complete a single task from the time it is submitted.

Speed-up and scale-up

Speed up: a fixed sized Problem executing on a small system is given to a system which is N -times larger.



Measured by

$$\text{Speedup} = \frac{\text{Small system elapsed time}}{\text{Large system elapsed time}} = \frac{T_S}{T_L}$$

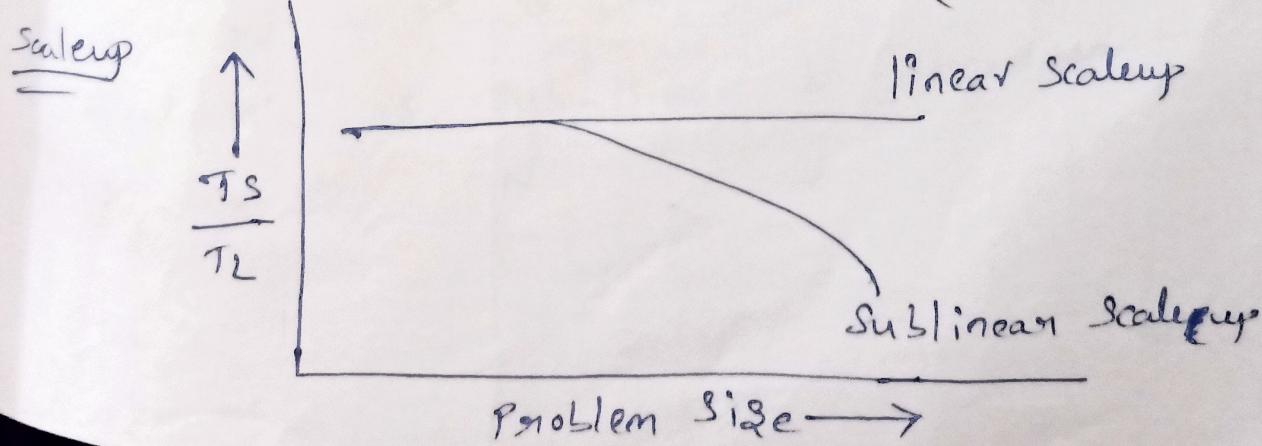
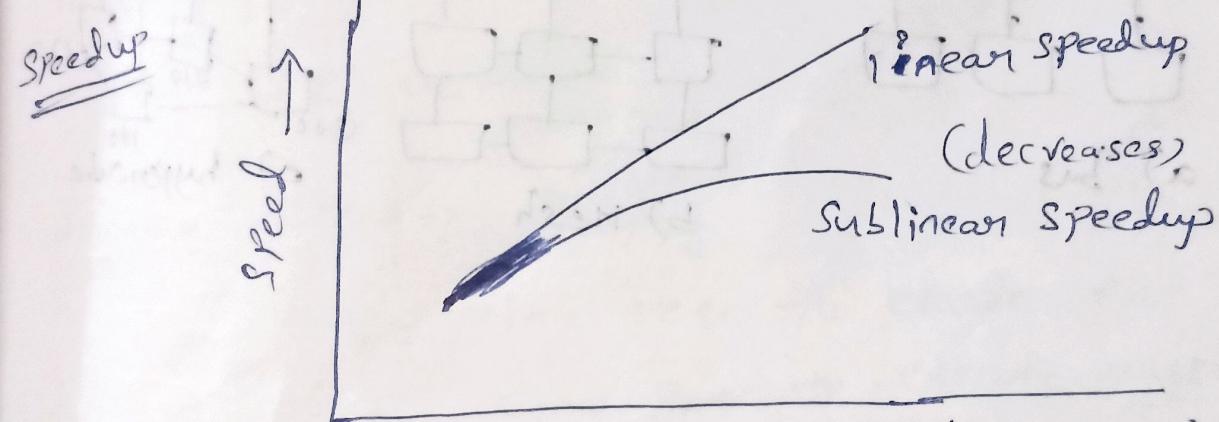
Speed up is linear if equation equals N.

Scaleup increase the size of both the problem and the system N-times : large system used to perform N-time larger job

Measured by:

$$\text{Scaleup} = \frac{\text{Small system small problem elapsed time}}{\text{Big system big problem elapsed time}}$$

Scaleup. is linear in equation equals 1



Interconnection

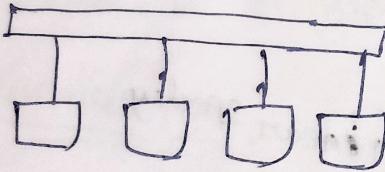
N/W Architecture

Bus: system components send data on and receive data from a single communication bus;

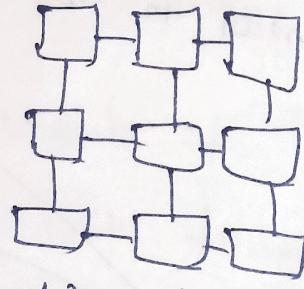
Does not scale well with increasing parallelism.

Mesh: components are arranged as nodes in grid, and each component is connected to all adjacent components.

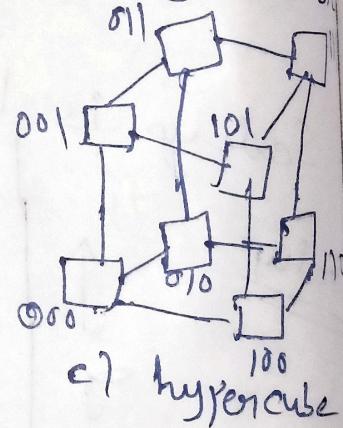
Hypercube: components are numbered in binary. Components are connected to one another if their binary representations differ in exactly one bit.



a) bus



b) mesh



c) hypercube

Parameters for Parallel Databases

To judge the performance of parallel database we are using the parameters

③ Response time: It is the time taken to complete a single task for given time.

2) Speedup in Parallel db

→ speedup is the process of increasing degree of (resource) parallelism to complete a running task in less time.

→ The time required for running task is inversely proportional to number of resources.

$$T(R) \propto \frac{1}{\text{no. of reses}}$$

formula Speedup = T_S/T_L

T_S = Time req to execute task of size Q

T_L = Time req to execute task of size $N+Q$

linear speed up is \sqrt{N} (no. of resources)

speed-up is sub-linear if speed up is less than N.

③ Scaling in Parallel database

Scale-up is the ability to keep performance constant, when no. of Process and resources increase proportionally.

Formula :

Let Q be the task & Q_N the task where N is greater than Q .

T_S = Execution time of task Q on smaller MC M_S

T_L = Execution time of task Q on smaller MC M_L

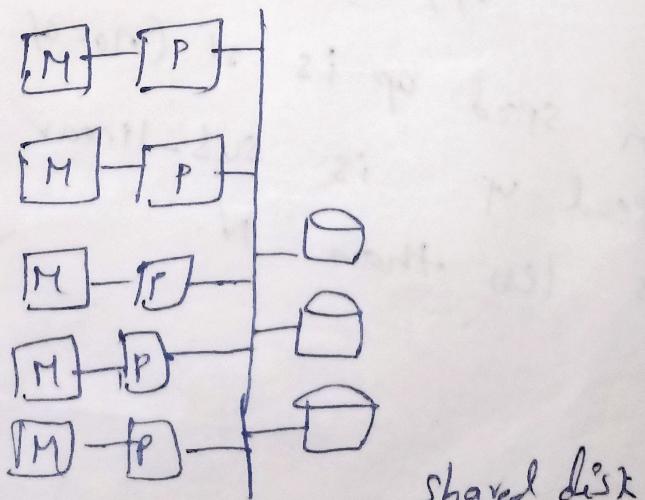
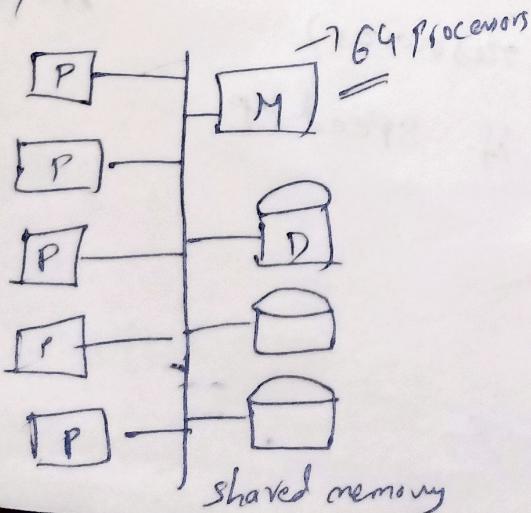
$$\text{Scale up} = \frac{T_S}{T_L}$$

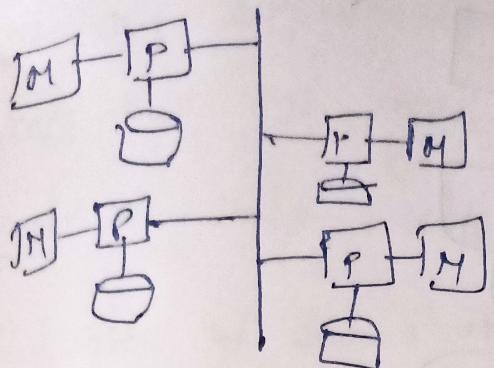
nos of Large MC
vol of smaller MC

Parallel Database Architectures

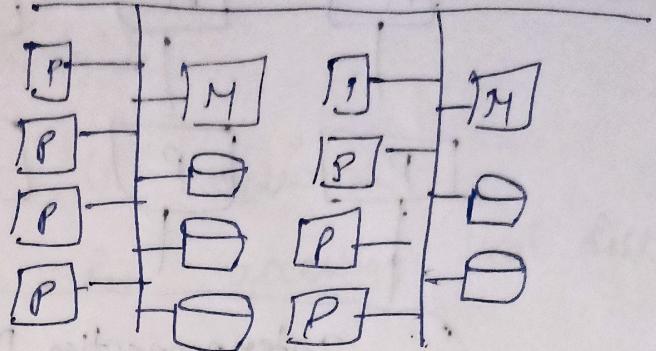
- 1) Shared memory - processors share a common memory
- 2) Shared disk - Processors share a common disk
- 3) shared nothing - Processors share neither a common memory nor common disk.

- 4) Hierarchical - hybrid of the above architecture





shared nothing



Hierarchical

Shared Advantage

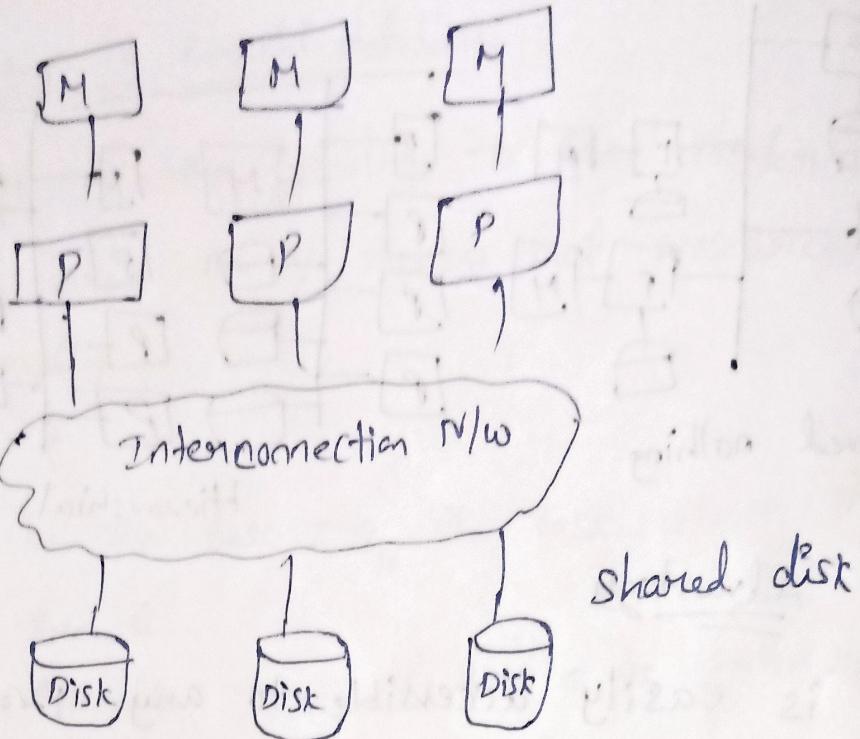
- Data is easily accessible to any processors
- one processor can send msg to other efficiently.

Disadv

- Waiting time of Processors is increased due to more no of Processors
- Bandwidth Problem.

Shared disk:

- Shared disk system uses multiple processor which are accessible to multiple disk via intercommunication channel and every processor has local memory.
- Each Processor has its own memory so that data sharing is efficient.
- The system built around this system are called as clusters.



Adv

i) Fault tolerance is achieved using shared disk system.
Fault tolerance: If a Processor or its memory fails the other Processor can complete the task. This is called as fault tolerance.

disadv:

→ Shared disk sys has limited scalability as large amount of data travels through the interconnection channel.
 → If more Processors are added the existing Processors are slowed down.

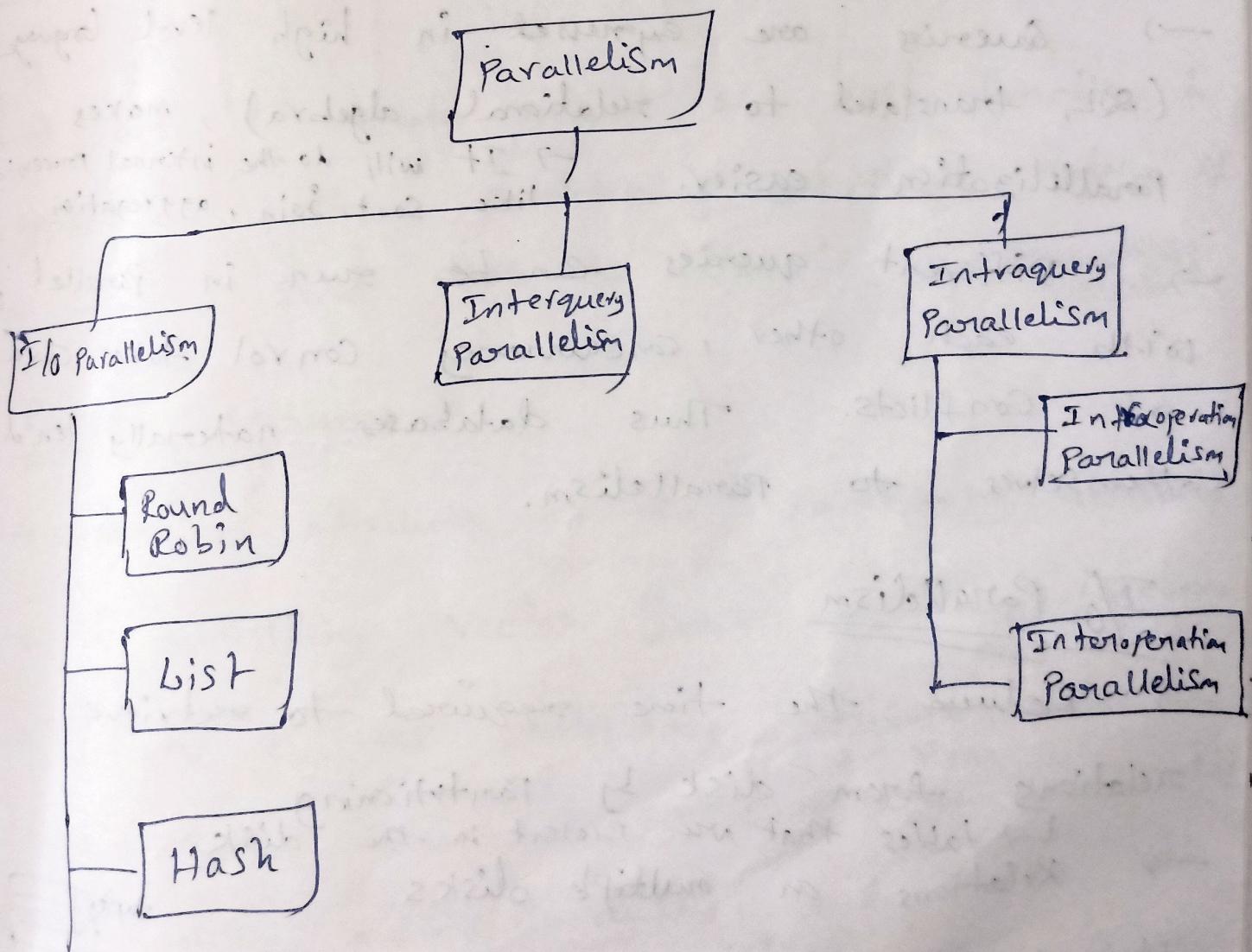
Shared nothing:- Adv i) No. of Processors and disk can be connected as per the requirement in shared nothing disk sys.

→ It can support for many Processors, which makes the sys more scalable.

Disadv 1) Data Partitioning is required

2) cost of communication for accessing local disk is much higher.

→ How Parallelism is achieved



Parallelism in Databases

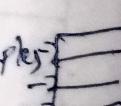
- Data can be Partitioned across multiple for Parallel I/O. → More no. of disks we can store more data
- Data can be partitioned and each process can work independently on its own partition.
- Individual relational operations (e.g. select, join, aggregation) executed in parallel.
- Queries are expressed in high level language (SQL, translated to relational algebra), makes parallelization easier. ↳ It will do the internal processing like sort, join, aggregation
- Different queries can be run in parallel with each other, concurrency control takes care of conflicts. Thus databases naturally lead themselves to parallelism.

I/O Parallelism

- Reduce the time required to retrieve relations from disk by Partitioning
 - ↳ tables that are present in the disk.
- Relations on multiple disks.

Horizontal Partitioning - tuples of a relation

tuples



are divided among many disks such that each tuple resides on one disk. (number of disks = 4)

Partitioning techniques (number of disk = n)

Round-Robin :

Send the i^{th} tuple inserted in the relation to disk $i \bmod n$.

Hash Partitioning

→ choose one or more attributes as the Partitioning attributes.

→ choose the hash function h with range $0 \dots n-1$.

→ Let h denote result of hash function h applied to the partitioning attribute value of a tuple. Send tuple to disk i . 

Range Partitioning

→ choose an attribute as the partitioning attribute

→ A partitioning vector $[v_0, v_1, \dots, v_{n-2}]$ is chosen $n=3$

→ Let v be the partitioning attribute value of tuple i . Tuples such that $v_i \leq v \leq v_{i+1}$ go to disk $i+1$.

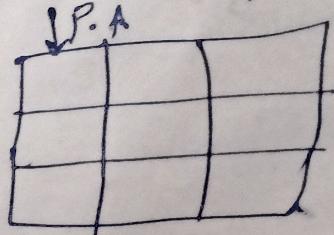
2) Tuples with $v < v_0$ to disk 0 and tuples with

$v \geq v_{n-2}$ go to disk $n-1$.

Vector $[v_0, v_1]$ tuple with $v < v_i$

go to disk 0

tuple with $v > v_{n-2}$ go to disk $n-1$



Eg: with a Partitioning vector [5, 11], with Partitioning attribute value of 2 will disk 0, a tuple with value 8 will go to while a tuple with value 20 will go to

Round Robin

i - record no

n - no. of disks

$i \bmod n$ is used for splitting records

$1 \bmod 3 = 1$ disk₁

$2 \bmod 3 = 2$ (disk₂)

Disk₀

Disk₁

2 Ram war

Record	id	Name	Branch
1	1	Sam	Hyd
2	2	Ram	War
3	3	Kiran	Gval
4	4	Xizhi	Hyd
5	5	Athore	Nai

no. of disks = 3

$i \bmod 3$

$i \bmod 3 = 1$

$i \bmod 3 = 2$

$i \bmod 3 = 0$

$2 \bmod 3 = 2$

$3 \bmod 3 = 0$

$4 \bmod 3 = 1$

$5 \bmod 3 = 2$

$6 \bmod 3 = 0$

Round

Robin

disadv
→ only

suitable for full table scans

→ not suitable for point queries or range queries.

L) select * from employee where name = 'Sam';

range

Select * from employee where

id between

3 and 5

Range: Range Partitioning strategy Partitions the data based on the partitioning attribute values.

→ We need to find set of range vectors on which we are about to partition.
Eg : "the records with salary range 100 to 5000 will be in disk 1, 5001 to 10000 in disk 2

id	Name	Branch	Salary	Salary Key
----	------	--------	--------	------------

Sal < 10000 - P0

Sal - 10000 < 30000 - P1

Sal > 30000 - P2

List Partitioning

chennai, vellore, hyd, tamilnadu, mumbai
↳ based on Branch.

Tamilnadu Partition	Maharashtra Partition
chennai	Mumbai
vellore	Pune

→ List Partitioning enables you to explicit control how rows map to partitions by specifying a list of discrete values for the Partitioning key in the description for each partition.

→ for a table with a Branch column by the Partitioning key, the Tamilnadu Partition

might contain values chennai and vellore.
the Maharashtra Partition might contain Mumbai

Hash:- we need to choose one or more
J.P.A

attributes from the table.

we need to apply hash function h

R.No
1
2
3

0 - 2
 $\rightarrow 0, 1, 2$

if h with orange 0 --- $n-1$
is zero is the result of \hookrightarrow no. of disks
hash function the tuple will move to
disk 0.

Comparison of Partitioning Techniques

1) Scanning the entire relation - Scan queries
 \hookrightarrow table

2) Locating a tuple associatively - Point queries
 \hookrightarrow A particular E.g. $A = 25$
tuple

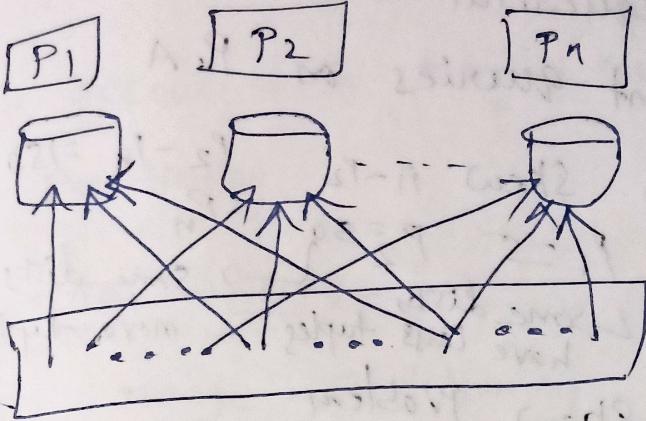
3) Range queries
Sal b/w 10000 & 20,000
 $\underbrace{10 \text{ and } 25}_{\text{Range}}$

Round robin Partitioning

→ Best suited for sequential scan of entire relation on each query.
• All disks have almost an equal no. of tuples.
• Retrieval work for entire relation is thus well balanced b/w disks (equal distribution)

→ Point queries and range queries are difficult to process.

- No clustering - tuples are scattered across all disks.
↳ no grouping



Given data.

Hash Partitioning

→ Good for sequential access

- Assuming hash function is good and Partitioning attributes from a key tuples will be equally distributed b/w disks. (Primary key)

• Retrieval work for entire relation is then well balanced b/w disks.

→ Good for point queries as P.A is a key

- Can lookup single disk, leaving others available for answering other queries.

• Index on P.A can be local to disk, making lookup and update more efficient.

→ No clustering, so difficult to answer range queries.

Range: $\{ \dots \} \{ 5 \dots \}$ $\{ \dots \}$ another group
another 2 group

→ Provides data clustering by Partitioning attribute value.

→ Good for sequential access

→ Good for Point queries on P.A

Problem Execution skew $P_1 - P_2$ $P_2 - P_n = 150$

(imbalance b/w disks) $P_1 = 5$ $P_2 = 50$ P_n → some disks have less tuples, more tuples

Handling of skew problem

→ Partitioning a Relation across disks

1) If a relation has only a few tuples which will fit into a single disk block, then assign the relation to a single disk.

2) Large relations are preferably partitioned across all the available disks.

3) If a relation consists of m disk blocks and there are n disks available in the system then the relation should be allocated min(m, n) disks.

→ The distribution of tuples to disks may be skewed → some disks have many tuples while others may have fewer tuples

Types of Skew:

Attribute - Value Skew:

- 1) All the tuples with the same value for the partitioning attribute end up in the same Partition.
- a) can occur with range-partitioning and hash-partitioning.

Partition Skew:-

- 2) With range-partitioning, badly chosen Partition Vector may assign too many tuples to some partitions and too few to others.
- b) less likely with hash-partitioning if a good hash-function is chosen.

Handling Skew in Range-Partitioning

→ To create a balanced Partitioning Vector (assuming attribute forms a key of the relation)

- a) sort the relation on the partitioning attribute.
- b) construct a Partition vector by scanning the relation in sorted order as follows

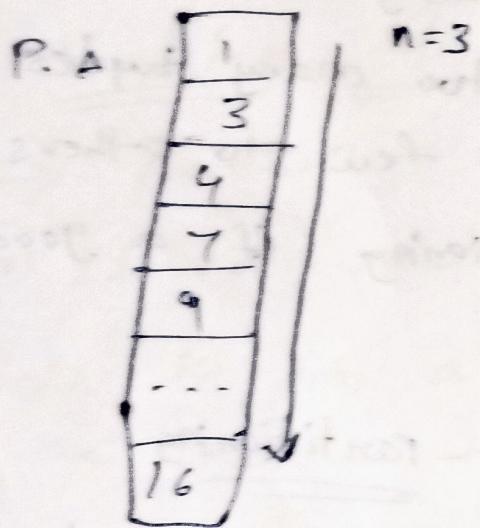
After every $\frac{1}{n^{th}}$ of the relation has been read, the value of the partitioning attribute

of the next tuple is added to the Partition vector.

b) n denotes the the number of Partitions to constructed

c) Duplicate entries or imbalances can result if duplicates are present in P.A.

Perform Partitioning using the Balanced Partition vector.



Balanced Partitioning Vector

1...4	7	[7, 12, 16]
7...11	12	(7, 12, 16)
12...16	16	<7 = 7 to 12 = 12 to 16

→ Handling skew using Histograms

> 16

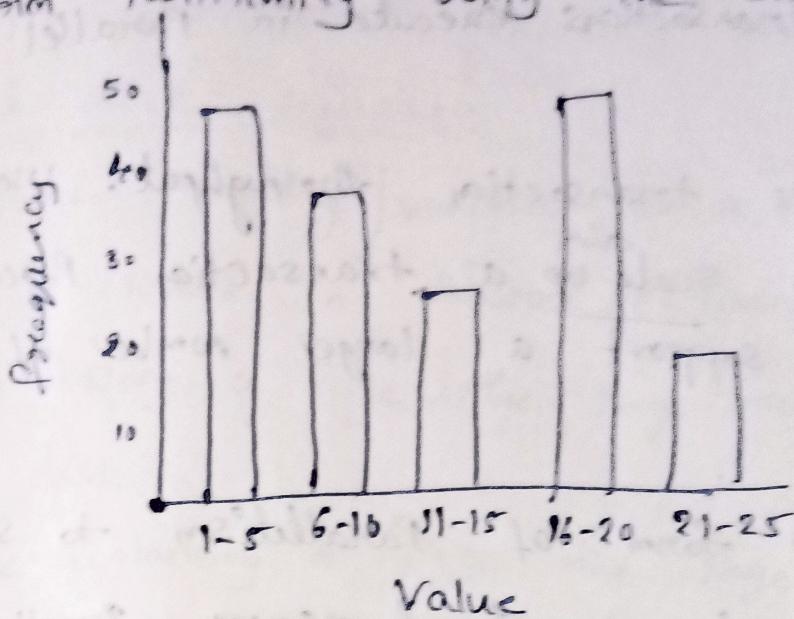
→ Balanced Partitioning Vector can be constructed from histogram in a relatively straightforward fashion.

→ Assume uniform distribution within each bin of the histogram.

→ Histogram can be constructed by Scanning relation, or Sampling (blocks containing tuples of relation)

→ Histograms can be stored in the log catalog.

Perform Partitioning using the Balanced P-V

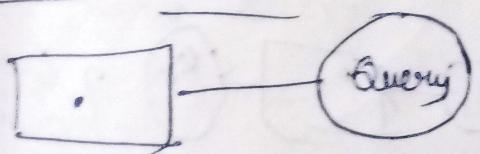


Handling Skew using virtual Processor Partitioning
skew in orange partitioning can be handled elegantly
using virtual Processor Partitioning

→ Create a large no. of Partitions (say 10 to 20
times the no. of processors)

→ Assign virtual Processors either in round-robin
fashion or based on estimated cost of processing
each virtual Partition.

Query Parallelism



1 Processor



$$\text{Speedup} = \frac{4 \text{ min}}{1 \text{ min}} = 4$$

Time: 4 min → In parallel db (S)
to improve the performance of the S.

Time: 1 min → It is achieved through query parallelism.

Interquery Parallelism

→ ^{Different} queries / transactions execute in parallel with one another.

→ Increases transaction throughput! used primarily to ^{aim} scale up a transaction processing system to support a larger number of trans per second.

→ Easiest form of "Parallelism" to support particularly in a shared memory parallel db because sequential db systems support concurrent processing.

→ More complicated to implement to on shared disk or shared-nothing architectures.

1) Locking and logging must be coordinated by passing msgs b/w processors.

2) Data in a local buffer may have been updated at another processor.

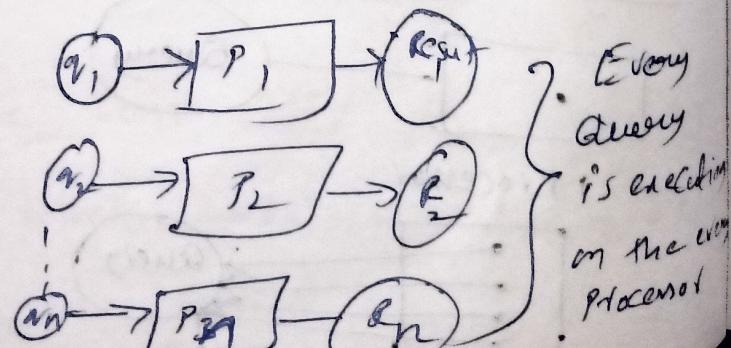
→ check that the cached data is up-to-date
cache-coherency has to be maintained - reads

and writes of data in buffer must find latest version of data.

→ Select * from emp;

→ Select * from Dept where mgrname = 'Steve'

→ Select Furniture_Name, Cost
from Furniture;



Cache Coherency Protocol

Example of Cache Coherency protocol for shared disk systems:

- Before reading / writing to a Page, the Page must be locked in Shared / exclusive mode.
- On locking a Page, the Page must be read from disk.
- Before unlocking a Page, the Page must be written to disk if it was modified.
- More complex Protocols with fewer disks read / writes exist.
- Cache Coherency Protocols for Shared Systems are similar. Each db. Page is assigned a home Processor. Requests to fetch the Page on write it to disk are sent to the home Processor.
→ nothing
Most common

IntraQuery Parallelism

The same query is being Parallelized — Means Part of same query is given to one Processor and another part of same query is

given to other processor/disk

→ Select * from Emp Dept where Emp.dno = Dept.dno;

- execution of single query in parallel multiple Processors / disks. ^{purpose} ^{aim}
- Important to speeding up long-running queries
- two complementary forms of intra-query Parallelism.

1) Intra-operation Parallelism - Parallelize the execution of each individual operation in the query. (Same operation is being parallelized).

Sort - 1 lakh

$\frac{1}{50,000} = \frac{50,000}{1P}$ $\frac{1}{2P}$

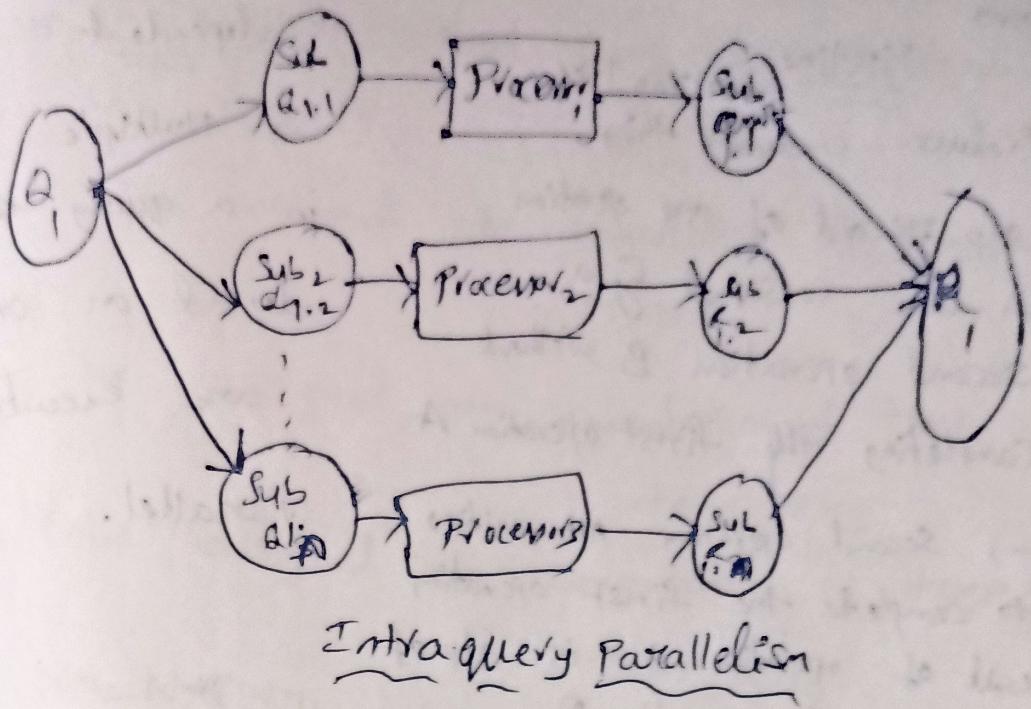
Eg: Parallel Sort,
Parallel Search

→ This form scales better with increasing parallelism because the number of tuples processed by each operation is typically more than the no. of operations in a query.

Inter-operation Parallelism - Execute the different operations in a query expression in parallel.

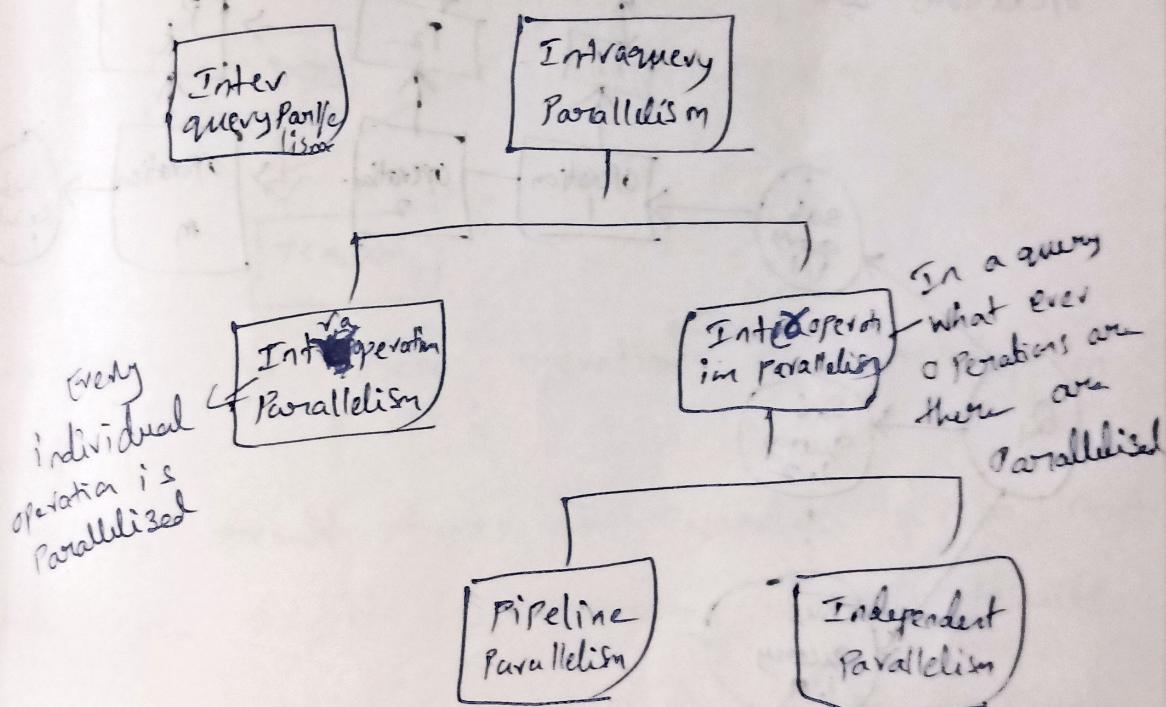
Example: simultaneous sorting or searching

→ The operations Sorting and Searching are parallelized.



Intraquery Parallelism

→ Forms of query Parallelism



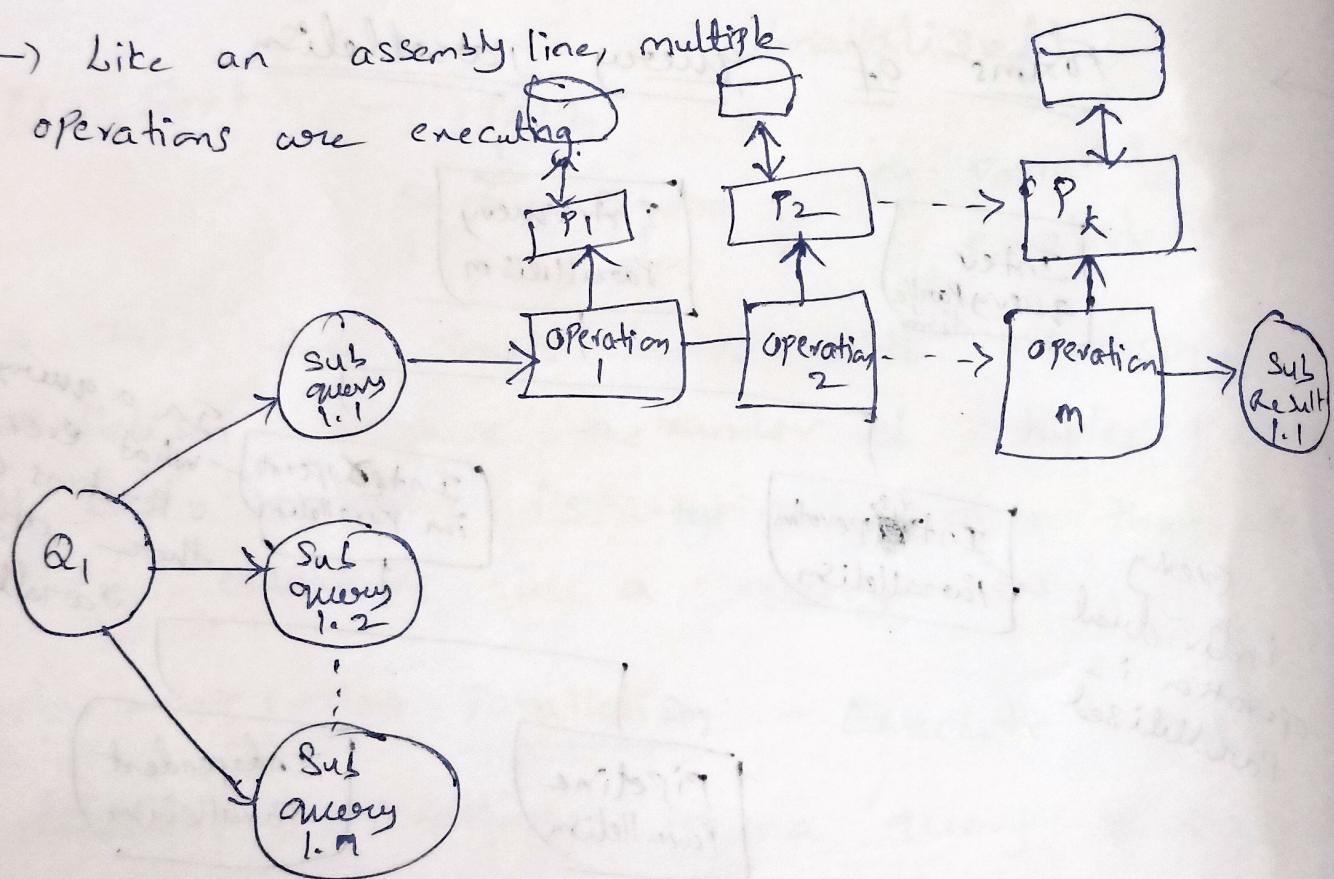
→ Interoperation Parallelism: Within a same query or transaction different operations are concurrently executing.

Pipeline working like
producer-consumer thing.
O/p record of one operation
A will be consumed by a
second operation B without
completing the first operation A.

→ Second operation not waiting
to complete the first operation
Result of operation A is I/p of
Second operation B.

→ Like an assembly line, multiple
operations are executing

Independent
→ Multiple operations
in a query that do
depend on one another
are executed in parallel.



Pipeline Parallelism

Consider a Join of four relations

\$R_1 \bowtie R_2 \bowtie R_3 \bowtie R_4

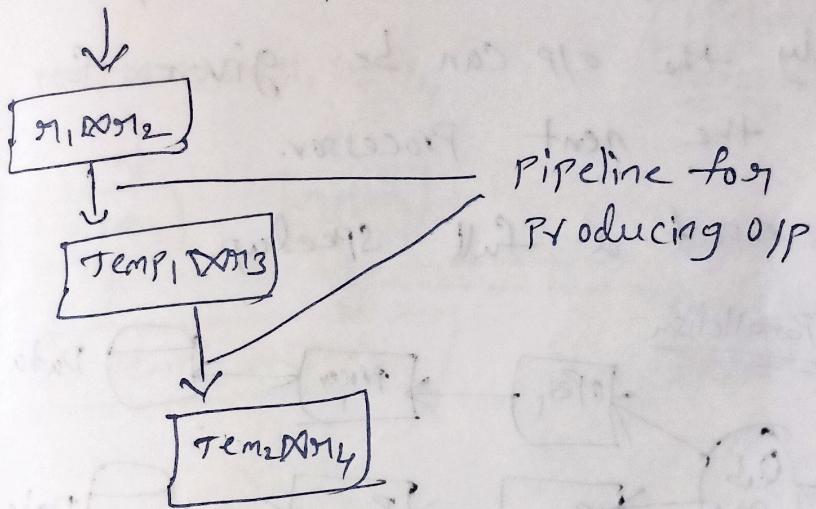
Set up a Pipeline that computes the three joins in parallel

Let P_1 be assigned to computation of $\text{temp}_1 = \pi_1, \Delta \pi_2$

and P_2 be assigned to computation of $\text{temp}_2 = \text{temp}_1, \Delta \pi_3$

P_3 be assigned to computation of $\text{temp}_3 = \text{temp}_2, \Delta \pi_4$

$\pi_1, \Delta \pi_2, \Delta \pi_3, \Delta \pi_4$



→ Each of these operations can execute in parallel.

→ sending result tuples it computes to the next operation even as it is consuming further results.

→ Provided a pipelineable join evaluation algorithm (e.g. indexed nested loop join) is used.

→ Factors Limiting utility of Pipeline Parallelism

1. Pipelined parallelism is not the good choice if degree of parallelism is high.
2. Useful with small number of Processors.
3. Not all operation can be Pipelined.

→ Example: Select AVG(salary) FROM Employee

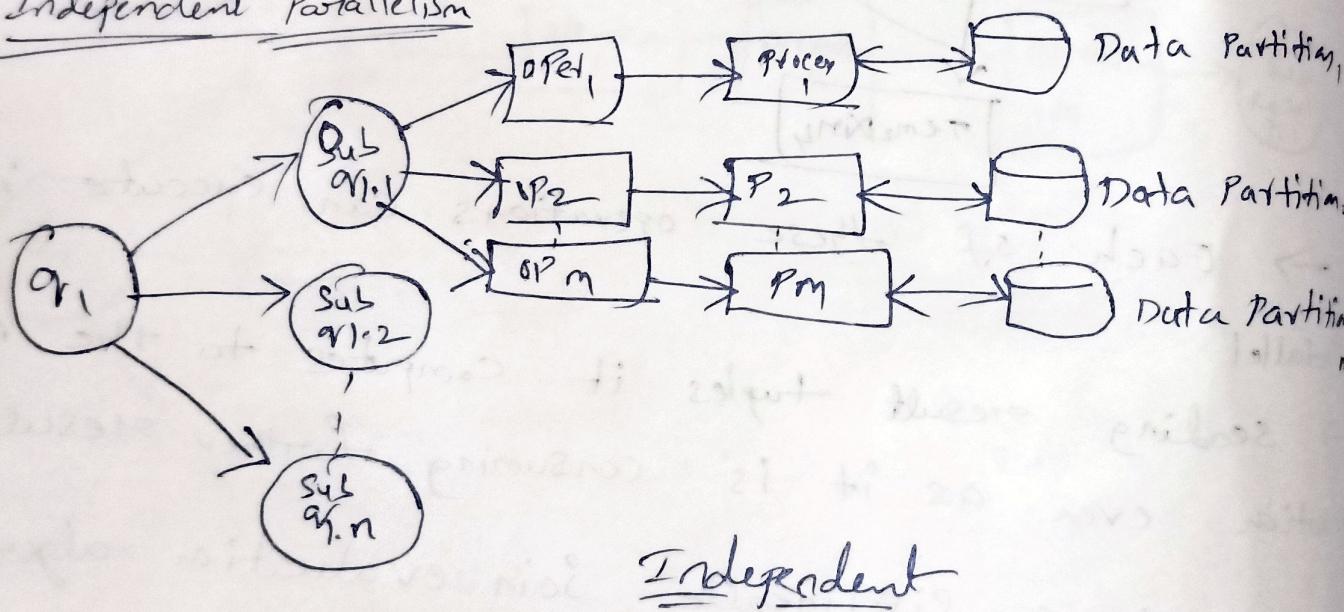
GROUPBy Dept_Id;

→ need to group at least one department

→ Then only the op can be given for aggregated operation on the next Processor.

4. cannot expect a full speedup.

→ Independent Parallelism



Consider a Join four relations

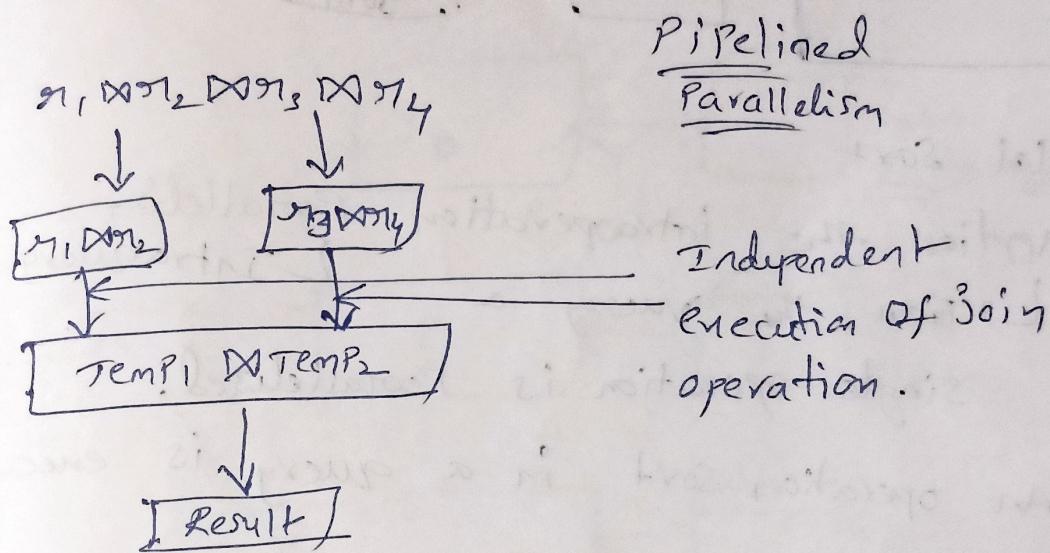
$$\pi_1 \bowtie \pi_2 \bowtie \pi_3 \bowtie \pi_4$$

Let P_i be the assigned the computation of π_i

$$= \pi_1 \bowtie \pi_2$$

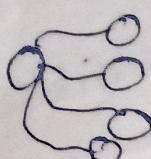
Let P_2 be assigned the computation of $\text{temp}_2 = \pi_3 \text{D}_{\text{avg}}$,
 Let P_3 be assigned the computation of $\text{temp}_1 \Delta \text{temp}_2$
 P_1 and P_2 can work independently in parallel.
 P_3 has to wait for IIP from P_1 & P_2 .

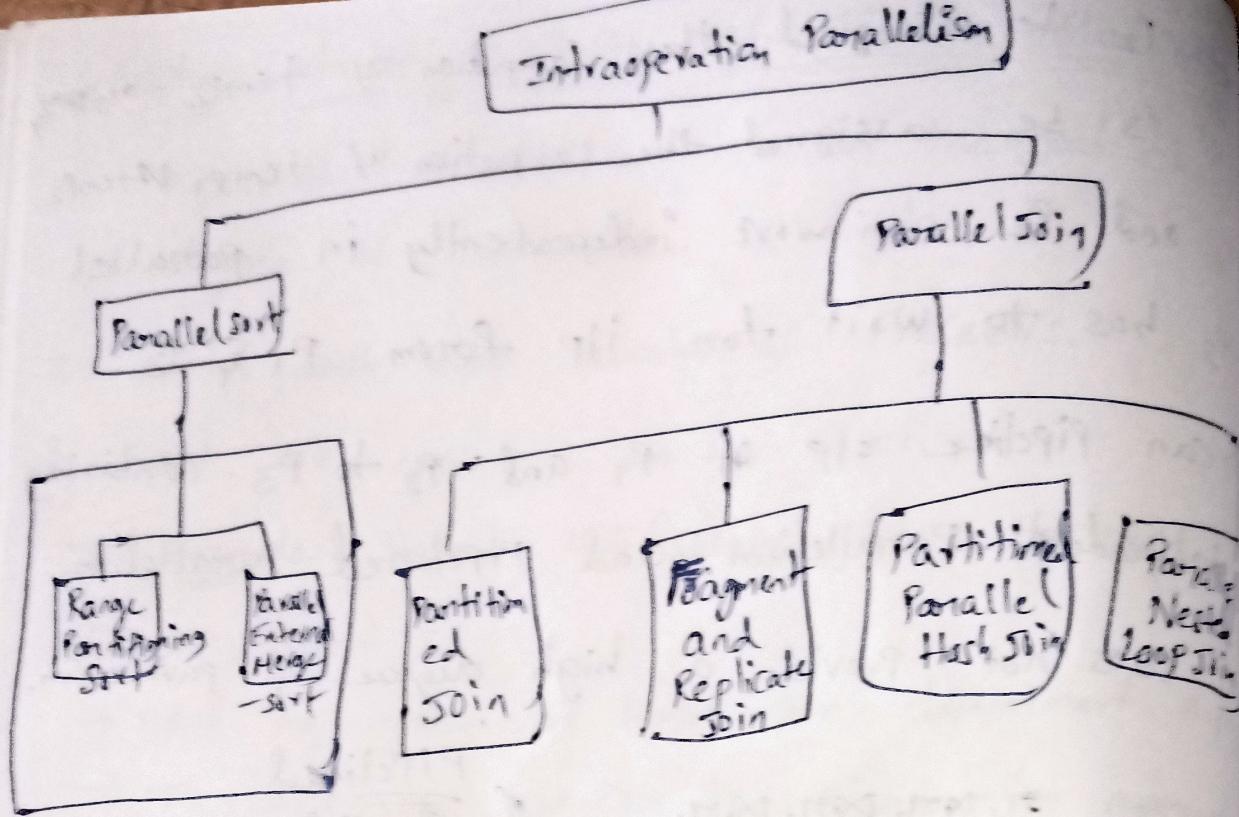
- Can pipeline IIP of P_1 and P_2 to P_3 . Combining independent Parallelism, and Pipelined Parallelism.
- Does not provide a high degree of parallelism.



→ Intraoperation Parallelism

- Each individual operation in a query is parallelized.
- Example: Parallel sort, parallel search, parallel join.
- A single operation joining is parallelized.





Parallel Sort

- Applies the intraoperation Parallelism
- ↳ From the query a single operation is parallelized.
- An operation sort in a query is executed in parallel.

Parallel Sorting Techniques

1. Range Partitioning Sort
2. External Sort-Merge.

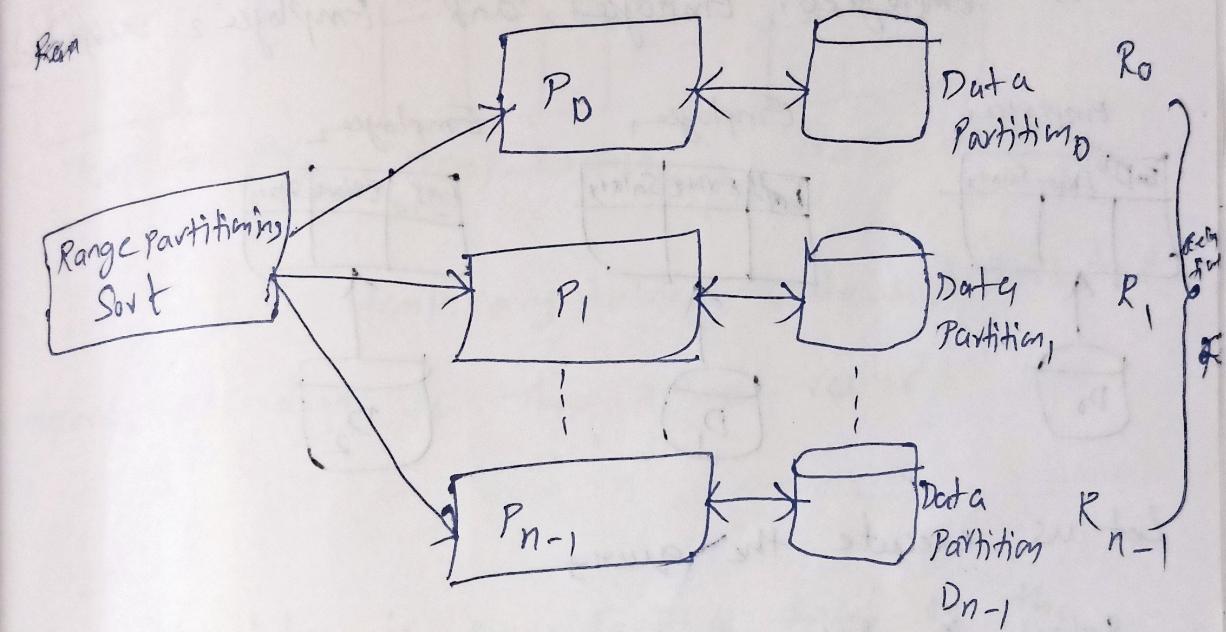
Assumptions:

- Assume n Processors $P_0, P_1, P_2, \dots, P_{n-1}$ and n disks $D_0, D_1, D_2, \dots, D_{n-1}$
- Disk D_i is associated with Processor P_i .

→ Relation R is Partitioned into R_0, R_1, R_2, \dots
using Round-Robin technique or Hash-partitioning
technique or Range Partitioning technique.

objective
our objective is to sort a relation (table) R that
resides on n disks on an attribute A in Parallel
sort

Ex: Sort on a table student records by roll no $\rightarrow A$



STEPS:

STEP1: Partition the relation R on the sorting attribute A at every processor using a range vector v.

Send the Partitioned records which fall in the i^{th} range to processor P_i . Where they are temporarily sorted in D_i .

Step 1: Sort each partition locally at each processor.

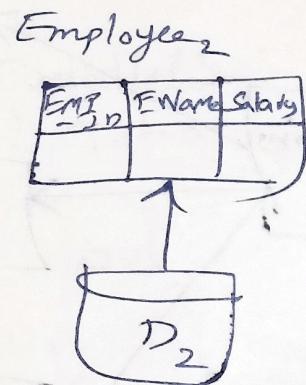
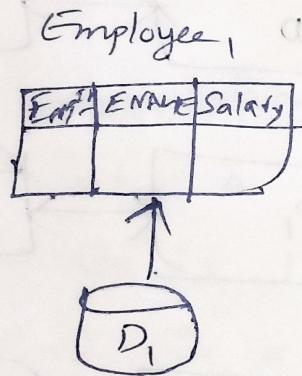
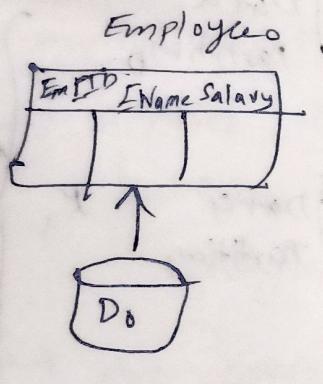
P; in parallel merge the sorted results.

Eg

Employee (Emp-ID, EName, Salary)

→ Assume that relation Employee is Partitioned using round robin technique into 3 disks D₀, D₁, D₂ which are associated with Processors P₀, P₁, P₂.

→ At Processors P₀, P₁, P₂, the relations are named Employee₀, Employee₁, and Employee₂ respectively.



Let us execute the query

Select * From

Employee

ORDER By Salary

→ Here the table is not partitioned according to the salary attribute.

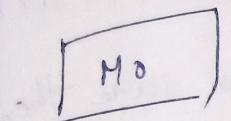
→ It is Partitioned by Emp-ID.

Step 1 (Range Partitioning) (Redistribution)

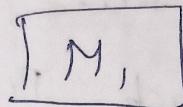
→ Identifying a range vector v on the salary attribute as v {v₀, v₁, ..., v_{n-1}}

In our Example, the tuples are distributed in 3 ranges are

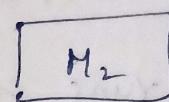
- Employee₀ in disk₀ with Range₀ (≤ 14000 and less)
- Employee₁ in disk₁ with Range₁ (14000 to 24000)
- Employee₂ in disk₂ with Range₂ (24001 to more)



Range: ≤ 14000



Range: > 14000



Range: > 24000

EmpID	ENAME	Salary

Temp_Employee₀

EmpID	ENAME	Salary

Temp_Employee₁

EmpID	ENAME	Salary

Temp_Employee₂

These temporary tables contains distributed records according to the range vector

Step 2:

- Sort the data in every Partition based on the Salary attribute.
- Every processor sorts its disk in parallel on different data sets. It is called Data Parallelism.

Final Result:

- After the processors completed the sorting, we can simply collect the data from different Processors & merge them.

Parallel Sort - Parallel External Sort - Merge

- Assume the relation has already been Partitioned among disks D_0, D_1, \dots, D_{n-1} (using any technique)
- In Parallel External Sort-merge
 - Each Processor P_i locally sorts the data on disk.
 - The sorted runs on each processor are then merged to get the final sorted O/P.
- Merging of sorted runs in step 2 can be parallelized as follows:
 - The sorted partitions at each processor P_i are range partitioned across the processor P_0, \dots, P_{m-1} .
 - Each processor P_i performs a merge on the stream as they are received, to get a single sorted partition.
 - The sorted runs on processors P_0, \dots, P_{m-1} are concatenated to get the final result.

Example:

- Employee relation is Partitioned using Round robin
 - Select * from Employee order by Salary;
 - Tables are not Partitioned on sorting attribute Salary.

D

ID	Name	Dept_name	Salary
6			10000
3			5000
9			6000
12			15000

P₁

ID	Name	Dept_name	Salary
1			15000
4			35000
7			10500

1 mode

2

P₂

ID	Name	Dept_name	Salary
2			25000
125			5000
8			7500

Step 1 Each processor P_i locally sorts the data on disk D_i on sorting attribute salary

D₀

ID	Name	Dept_name	Salary
3			5000
9			6000
6			10000
12			15000

D₁

ID	Name	Dept_name	Salary
7			10500
1			15500
4			35000

P₀ P₁

ID	Name	Dept_name	Salary
5			5000
8			7500
2			20000

- step 2
- Partition vector on Salary attribute $\{14000, 25000\}$
 - Redistribute every partition using merge Vector on three disks.

P_0 processor P_0 merges incoming data at Partition D_0

ID	Name	Dept name	Salary
3			5000
5			5000
9			6000
8			75000
6			10000
7			10500

P_1 $\geq 14000 < 25000$

ID	Name	Dept name	Salary
12			15000
11			15500
10			16000
9			16500

Processor P_1 merges incoming data at Partition D_1

P_2

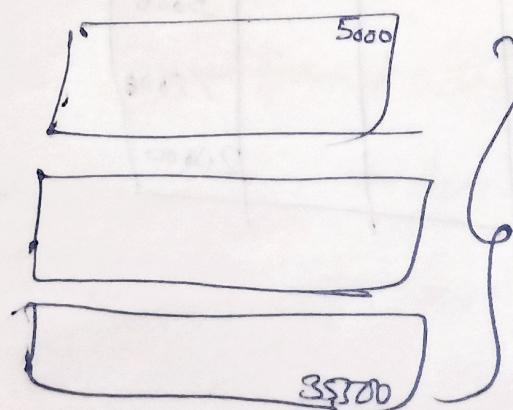
ID	Name	Dept name	Salary
2			25000
4			35000

Processor P_2 merges incoming data at Partition

Final merge

concatenate of sorted data from all the disks

D_0, D_1, D_2



The data is sorted.

Parallel Join

- The join operation requires pairs of tuples to be tested to see if they satisfy the join condition, and if they do, the pair is added to the join op.
- Parallel Join algorithms attempt to split the pairs to be tested over several processors. Each processor then computes part of the join result locally.
- In a final step, the result from each processor can be collected together to produce the final result.

Let r and s be the input relations and we want to compute $r \bowtie s = r \cdot s$.
→ Equi Join and s each are partitioned into n partitions denoted r_0, r_1, \dots, r_{n-1} and s_0, s_1, \dots, s_{n-1} . Equi join is a special type of join in which we use only equality operator.

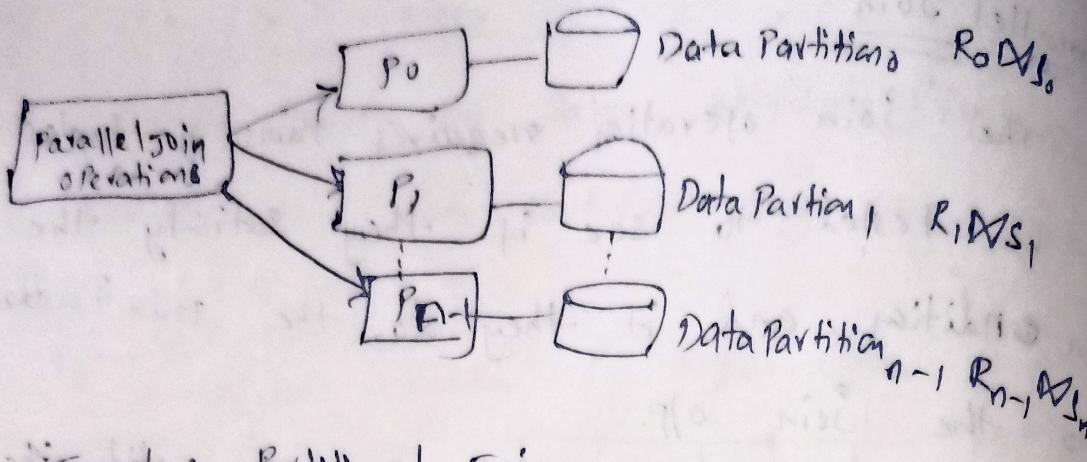
`SELECT * FROM Emp INNER JOIN Dept`

`ON Emp. DEPTID = Dept. DEPTID;`

Natural Join

Natural join is a type of equi join which occurs implicitly by comparing all the names columns in both tables.

→ `SELECT * FROM Emp NATURAL JOIN Dept`



Example :- Partitioned Join

Student			
RollNo	Name	Gen	Phone
1	Aman	M	---
3	Anand	M	---
4	Rupa	F	---
2	Ram	M	---

common attribute
CoA : RollNo

Course	
RollNo	Course Name
4	AD
2	AI
3	DM
1	AI

Range (or)
Hash Partition

Relations: Student and Course

Joining attribute : Roll No

No. of disks = 2 hence No. of Partitions = 2

Partitioning technique : Hash Partitioning

P. A : Roll No.

$$\text{Hash Function: } h(\text{RollNo}) = (\text{RollNo} \bmod n) = (\text{RollNo} \bmod 2)$$

RollNo	Name	Gen	Phone
2			
4			

Student

RollNo	Name	Gen	Phone
1			
3			

Student

$$\begin{aligned} \text{Partition 0} &= 1 \bmod 2 = 1 \\ &2 \bmod 2 = 0 \\ &3 \bmod 2 = 1 \\ &4 \bmod 2 = 0 \end{aligned}$$

Student course.

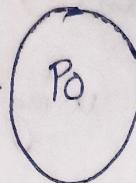
RollNo	Course Name
4	
2	

Student course

RollNo	Course
3	
1	

Students

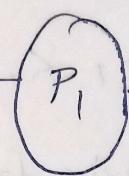
Roll	Name	Gen	Phone
4			
2			



Courses

RollNo	Course Name
4	
2	

Roll			
1			
3			



Roll	
3	
1	

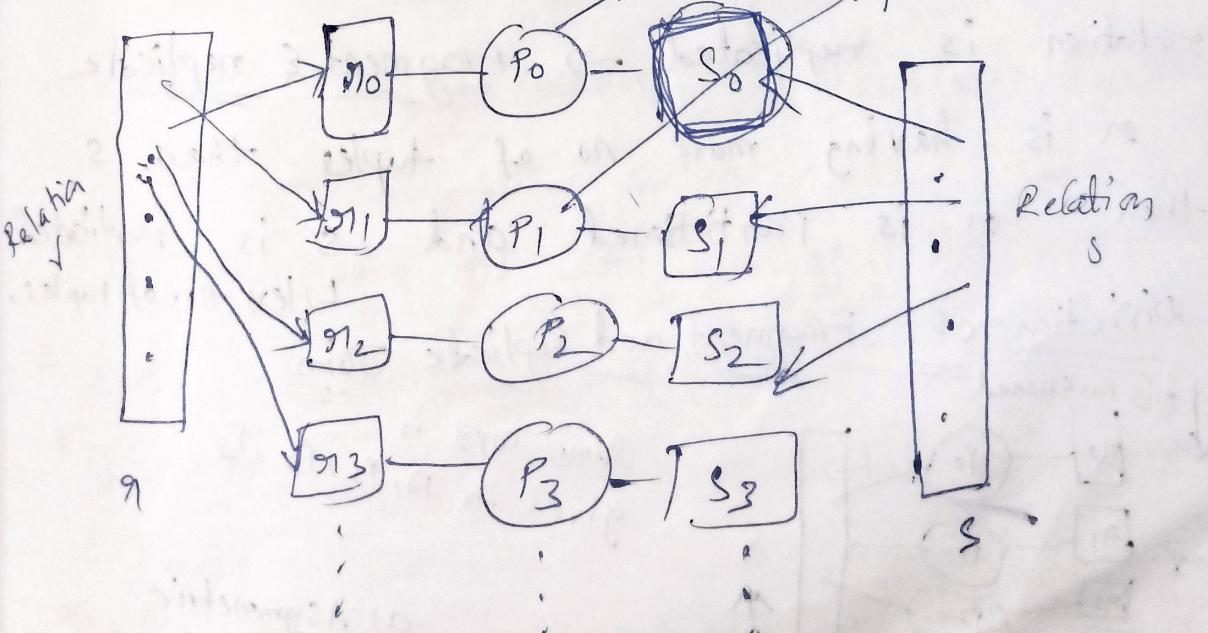
Result of Parallel Join operation: Student & Course

RollNo	Name	Gen	Phone	Course
4	Rupa	F	---	AD
2	Ram	M	---	AI
1	Aman	M	---	AI
3	Anand	M	---	DM

{Parallel}

Roll No:

R, NS, Execution



relation R

relation S

→ Fragment and Replicate Join

→ Partitioning not Possible for Some Join

e.g.: non-equiJoin conditions, such as $r.A > s.B$

- condition is established using all comparison operators except the equal ($=$) operator.

like $!=, >, \leq, =, <, >$.

• Special case - Asymmetric fragment and replicate

- one of the relations, say r , is partitioned

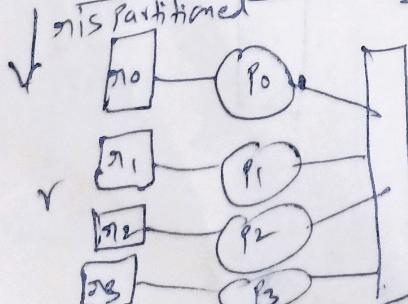
- any Partitioning technique can be used.

- The other relation, s , is replicated across all the processors.

- Processor P_i , then locally computes the join of r_i with all of s using any join technique.

→ one of the relation is partitioned and other relation is replicated → Fragment & replicate
 or is having more no. of tuples than s
 then r is partitioned and s is replicated

Dipiction of Fragment and Replicate Joins

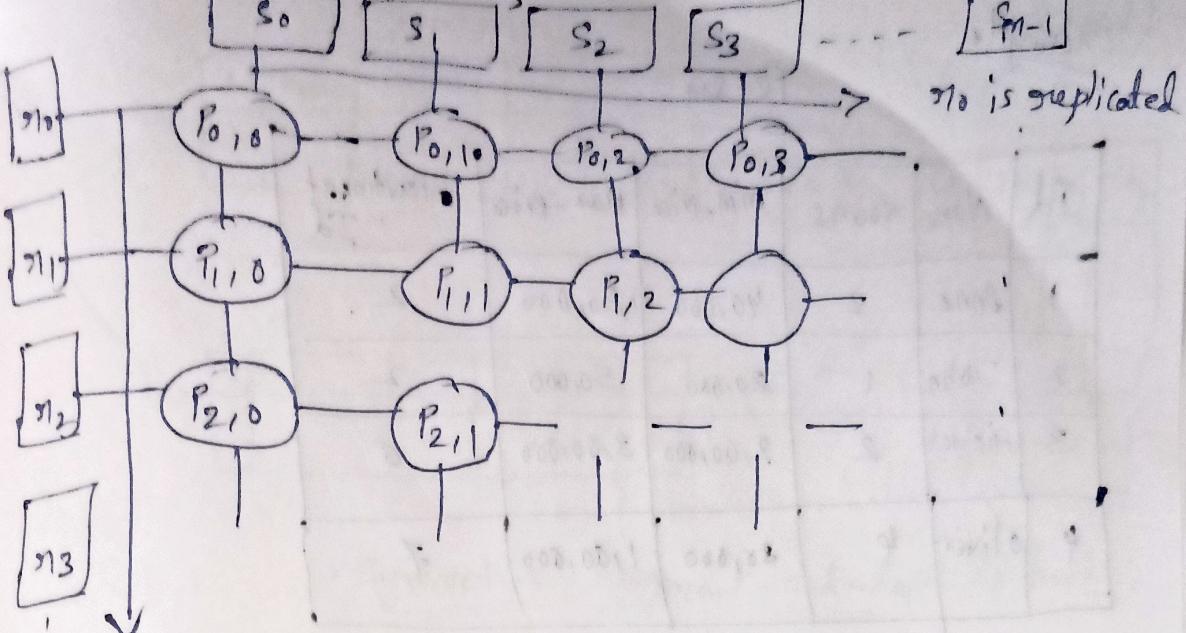


L less no. of tuple

Same copy is given to p_0, p_1, p_2, p_3

s is replicated

a) Asymmetric fragment and replicate



s_0 is

replicated

$r_{1,n-1}$

b) Fragment and replicate.

→ r is partitioned into n partitions r_0, r_1, \dots, r_{n-1}

→ s is partitioned into m partitions $s_0, s_1, s_2, \dots, s_{m-1}$

→ Any partitioning technique may be used.

→ There must be at least $m \times n$ Processors

→ Label the Processors as $P_{0,0}, P_{0,1}, \dots, P_{0,m-1}, P_{1,0}, P_{1,1}, \dots, P_{n-1,m-1}$

$P_{n-1,m-1}$

→ P_{ij} the Join of r_i with s_j . In order to do

s_0, r_i is replicated to $P_{i,0}$

$P_{i,1} \dots P_{i,m-1}$, while s_i is replicated to $P_{0,i}, P_{1,i}, \dots, P_{n-1,i}$

$P_{n-1,i}$

→ Any join technique can be used at each processor P_{ij} .

example

Person

Id	Name	rooms	min-price	max-price	Apartment - Id
1	Anne	2	40,000	150,000	2
2	John	1	20,000	50,000	2
3	Michael	2	2,00,000	3,00,000	6
4	Oliver	4	85,000	1,00,000	7

Apartment

Id	rooms	Price	city
1	2	30,000	Houston
2	2	45,000	Dallas
3	3	1,25,000	Chicago
4	5	2,45,000	Los Angeles
5	4	3,40,000	San Jose
6	4	2,20,000	San Diego
7	1	36,000	Cleveland

SELECT Name, min_price, max_price, Price, city
 FROM Person JOIN apartment ON apartment.id = person.apartment_id
 AND Price BETWEEN min_price AND max_price;

Person Fragmentation

id	name	rooms	min	Max	Offering
2					
4					

Person table

(Copied) Replication

id				
1				
3				

P0

id	Room	Price	City

P1

Symmetric Fragment-and-Replicate Join

Apartment table

id	rooms	Price	city
2			
4			

id			
1			
3			
5			

Person table

Apartment table

id			
2			
4			

P00

P01

id			
1			
3			

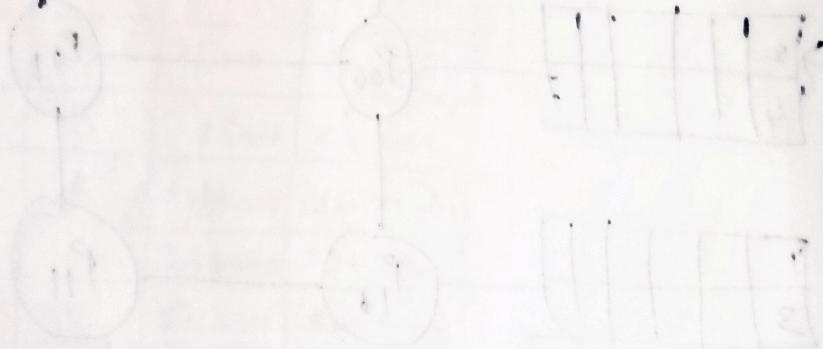
P10

P11

Fragment and Replicate Join

Name	Min-Price	Max-Price	Price	City
John	20,000	50,000	30,000	Houston
John	20,000	50,000	36,000	Cleveland
Anne	40,000	1,50,000	1,25,000	Chicago
Michael	2,00,000	3,00,000	2,45,000	Los Angeles
Oliver	30,000	1,00,000	45,000	Dallas
Oliver	30,000	1,00,000	30,000	Houston

- Both versions work with any join condition: every tuple in σ can be tested with every tuple in τ .
- usually has a higher cost than Partitioning:
since one of the relations (for asymmetric fragment-and-replicate) or both relations (for general fragment-and-replicate) have to be replicated
- sometimes asymmetric f-a- σ is preferable even though Partitioning could be used.



more elaborate form to implement

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----

Parallel Join : Partitioned Parallel Hash Join

Assume that we have

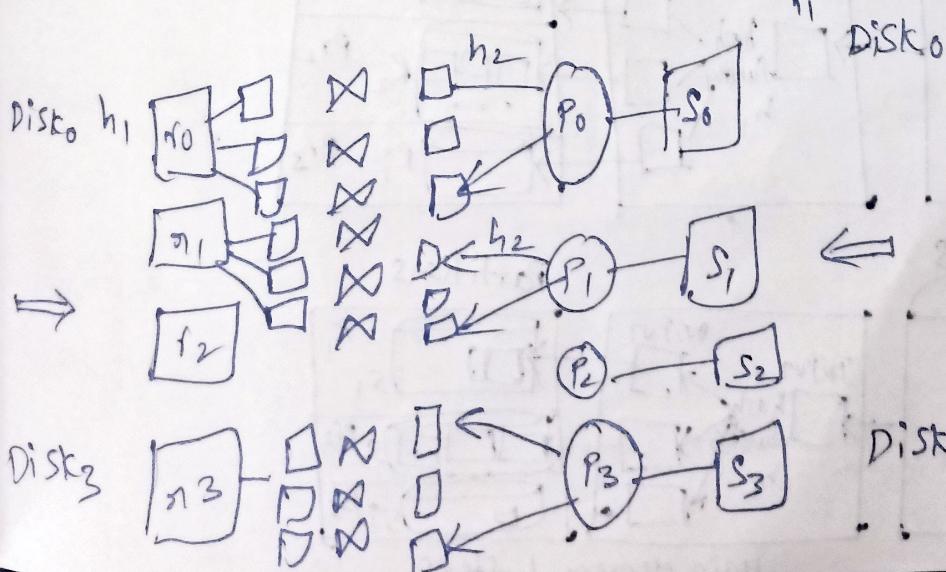
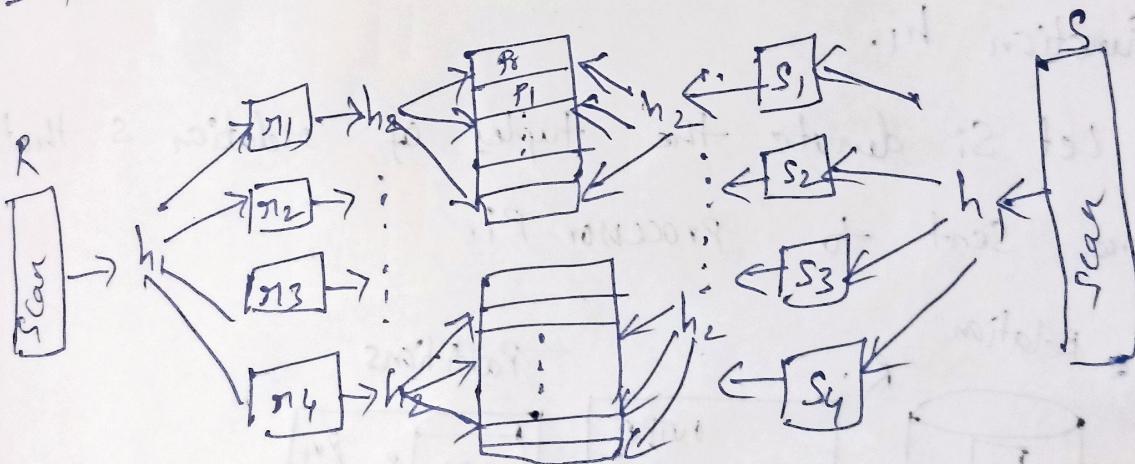
processors P_0, P_1, \dots, P_{n-1}

Two tables R and S for join operation
 (R and S are already partitioned into dists
 of n processors)

Assume S is smaller than R and therefore

S is chosen as the build relation.

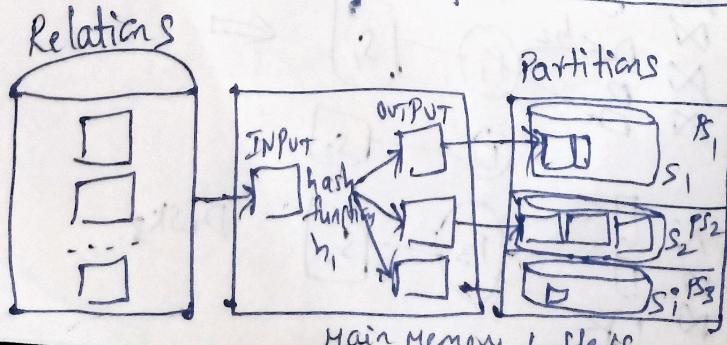
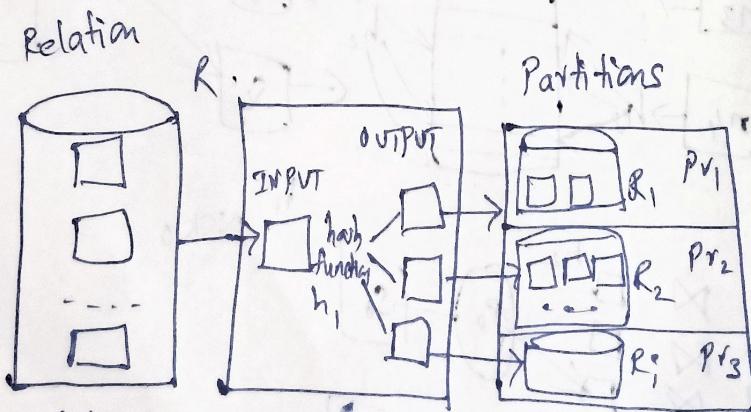
Hash functions : h_1 & h_2 .



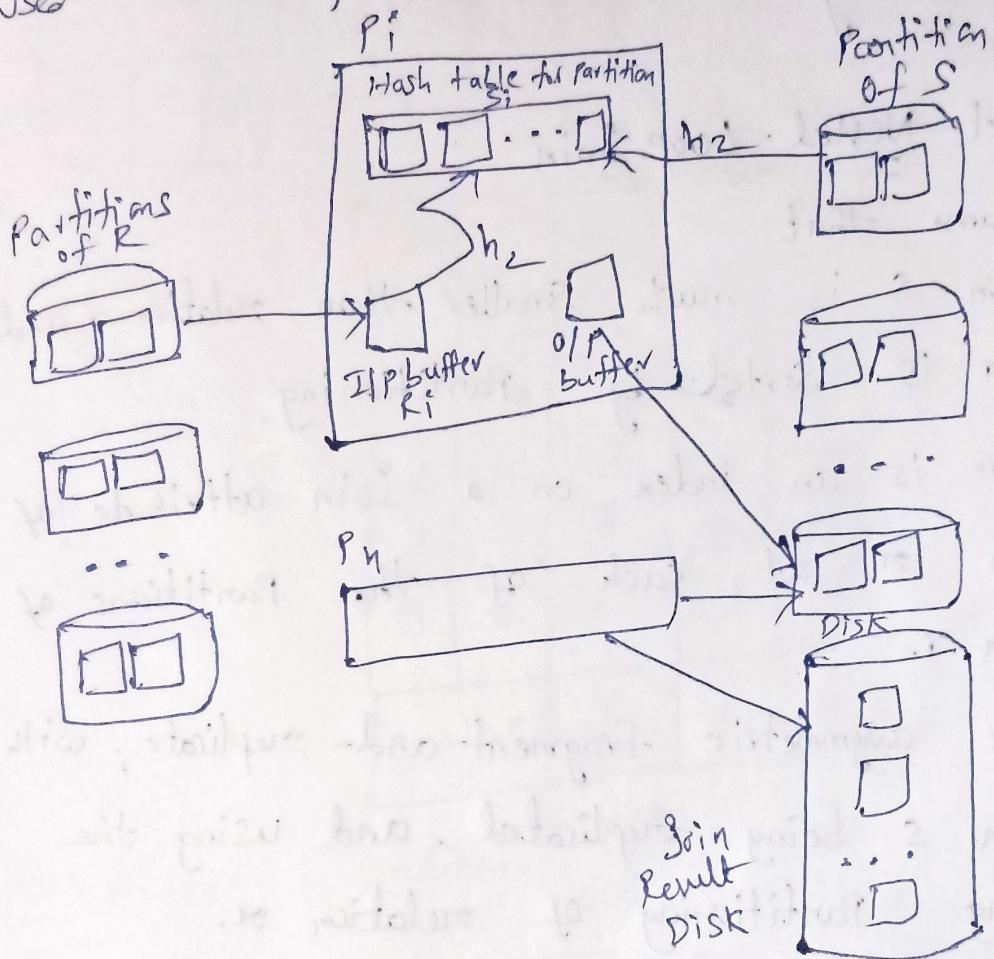
- h_1 is used for first level Partitioning
- h_2 is used for Second level Partitioning
- Later join operation takes place in P

Hash function h_1

- takes the join attribute value of each tuple in s_1 and maps the tuple to one of the n processors.
- Each Processor P_i stores the tuples of s_1 which are on its disk D_i , and sends each tuple to the appropriate Processor based on hash function h_1 .
- Let S_i denote the tuples of relation s_1 which are sent to processor P_i .



→ Hash function h_2
 → As tuples of relation S are received at the destination Processors, they are partitioned further using another hash function h_2 , which is used to compute the hash-join locally.



- Each Processor P_i executes the build and probe phases of the hash-join algorithm on the local partitions r_i and s_i of R and S to produce a partition of the final result of the hash-join.
- Note: Hash-join optimizations can be applied to the parallel case.

→ e.g. - the hybrid hash-join algorithm can used to cache some of the incoming tuples in memory and avoid the cost of writing and reading them back in.

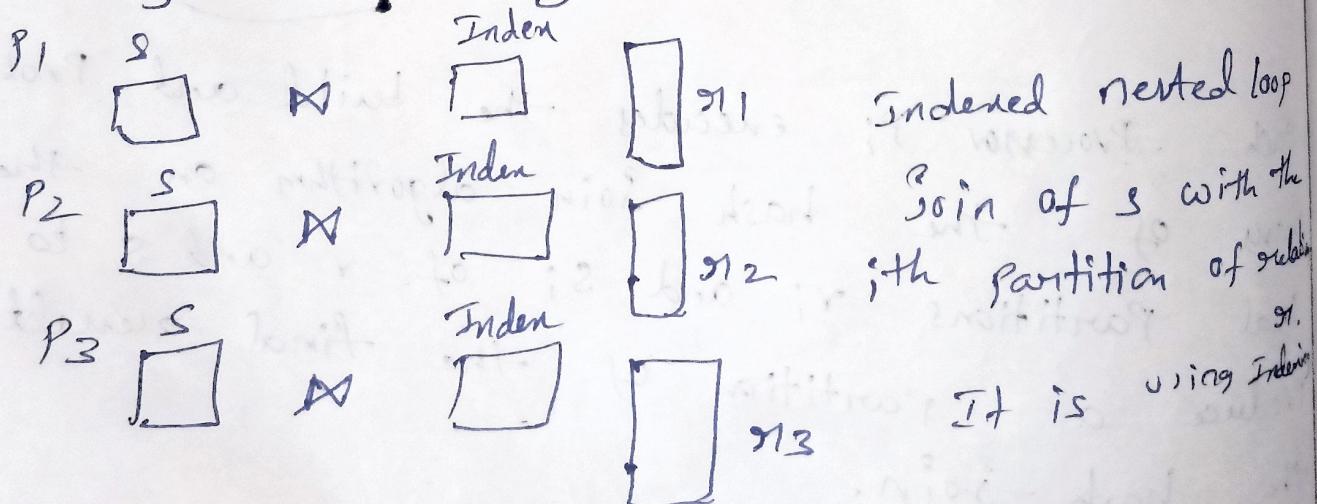
Parallel Nested-Loop Join

→ Assume that

- Relation S is much smaller than relation r and that r is sorted by Partitioning.

→ There is an index on a join attribute of relation r at each of the partitions of relation r .

→ Use asymmetric fragment and replicate, with relation S being replicated, and using the existing Partitioning of relation r .



Join :-

Cartesian product

P

A	B	C
1	a	b
2	c	d
3	e	f

Q

D	E
1	c
2	d

$P \times Q$

A	B	C	D	E
1	a	b	1	c
1	a	b	2	d
2	c	d	1	c
2	c	d	2	d
3	e	f	1	c
3	e	f	2	d

Natural Join :- condition.

P

(A)	B	C
1	a	b
2	c	d

(A)	D
1	d
3	e

Same name
and datatype

P Natural Join Q Attribute

Record
of
tuple

A	B	C	D
1	a	b	d

P	A	B
1	a	
2	c	

Q	C	D
1	b	
3	d	

P Left outer join Q
 $A = c$

A	B	C	D
1	a	b	b
2	c	NULL	NULL

P Right outer join Q
 $A = c$

A	B	C	D
1	a	1	b
3	NULL	3	d

Right Full outer join

A	B	C	D
1	a	b	b
2	c	NULL	NULL
NULL	NULL	3	d

INNER

P INNER Join Q $A = c$

A	B	C	D
1	a	1	b

Based on condition

SELECT * FROM P JOIN Q ON P.A = Q.C

→ Inter operation :-

SELECT * FROM Emp WHERE Salary > 5000

ORDER BY NAME;

SCAN Emp table for Salary > 5000

Sort all the Emp records on Name attribute

Two processes come to scan table, and the other to perform sorting.

Interoperation

Parallelism

It is about executing different operations of a query in parallel. A single query may involve multiple operations at once. We may exploit parallelism to achieve better performance of such queries. Consider the example query given below

SELECT AVG(Salary) FROM Employee GROUP BY Dept_Id;

It involves two operations.

First one is an aggregation and second is grouping.

For executing this query we need to group all the emp. records based on the attribute Dept_Id first.

Then, for every group we can apply the AVG aggregate function to get the final result.

We can use Interoperation Parallelism concept to parallelize these two operations

1. Pipelined
- 2) Independent

In Pipelined Parallelism, the idea is to consume the result produced by one operation by the next operation in the pipeline.

Example - Hash Partitioning

- Assume that there are $n=3$ disks, D_0, D_1, D_2 , among which the relation is to be distributed and partitioning attribute is ID.
- The hash function chosen is $H(x) = x \bmod n$ where x is value of a P.A ID.
- if $x=1$, the tuple is sent to $H(1) = 1 \bmod 3 = 1$, that is, D_1 .
- if $x=2$, the tuple is sent to $H(2) = 2 \bmod 3 = 2$, that is, D_2 .
- if $x=3$, , , , , $H(3) = 3 \bmod 3 = 0$, that is, D_0 .
- ;
- ;
- if $x=10$ / / / / / $H(10) = 10 \bmod 3 = 1$, that is, D_1 ,
- $x=11$ / / / / / $H(11) = 11 \bmod 3 = 2$, that is, D_2 ,
- $x=12$ 1 1 1 1 1 $H(12) = 12 \bmod 3 = 0$, that is, D_0 ,

ID	Name	Dept_name	Salary	RR ID	Name
1	Srinivas	CSE	40000	D_0	3
2	Katz	ENTC	35000	D_2	12
3	Crick	ENTC	49000	D_1	10
;					
10	Singh	CSE	75000	D_2	2
11	Brooks	MECH	60000		11
12	Tanaka	ENTC	65000		

Range: Assume that there are $n=5$ disks D_0, D_1, D_2, D_3, D_4

P.A = ID

Partition vector is chosen as $[20, 40, 60, 80]$

ID	Name	Dept-name	Salary	D ₀	ID
1				19	
2				D ₁	20
3					39
4				D ₂	40
5					59
6					
7					
8					
9					
10					

It chooses a P.A A

It chooses a P.V [v₀, v₁, ..., v_{n-2}]

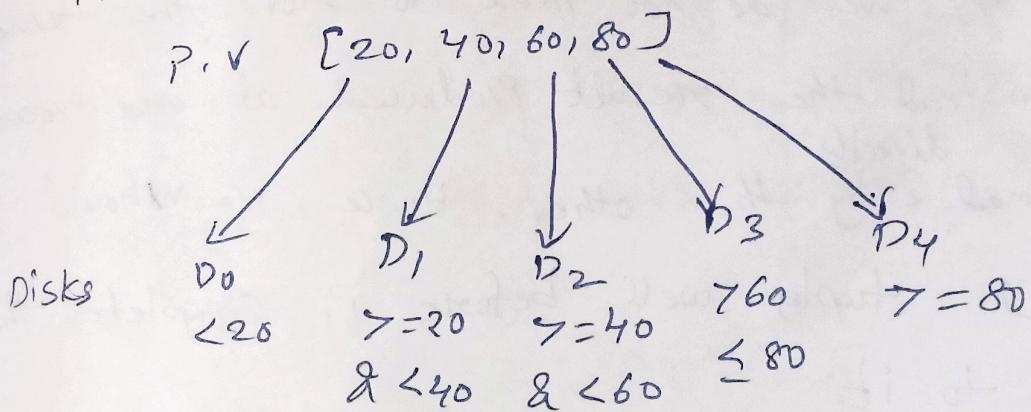
such that if i < j then v_i < v_j
consider a tuple t such that t[A] = x

if x < v₀, then t goes to disk 0

if x > v_{n-2}, then t goes to disk n-1

if v_i ≤ x < v_{i+1} then goes to disk D_{i+1}

P.V [20, 40, 60, 80]



For example, consider the following operation

$\sigma_1, \sigma_2, \sigma_3, \sigma_4$

The above expression shows natural join operation.
This actually joins four tables. This operation can be pipelined as follows.

Perform $\text{temp}_1 \leftarrow \sigma_1, \sigma_2$ at processor P_1 and

Send result temp_1 to processor P_2 to perform

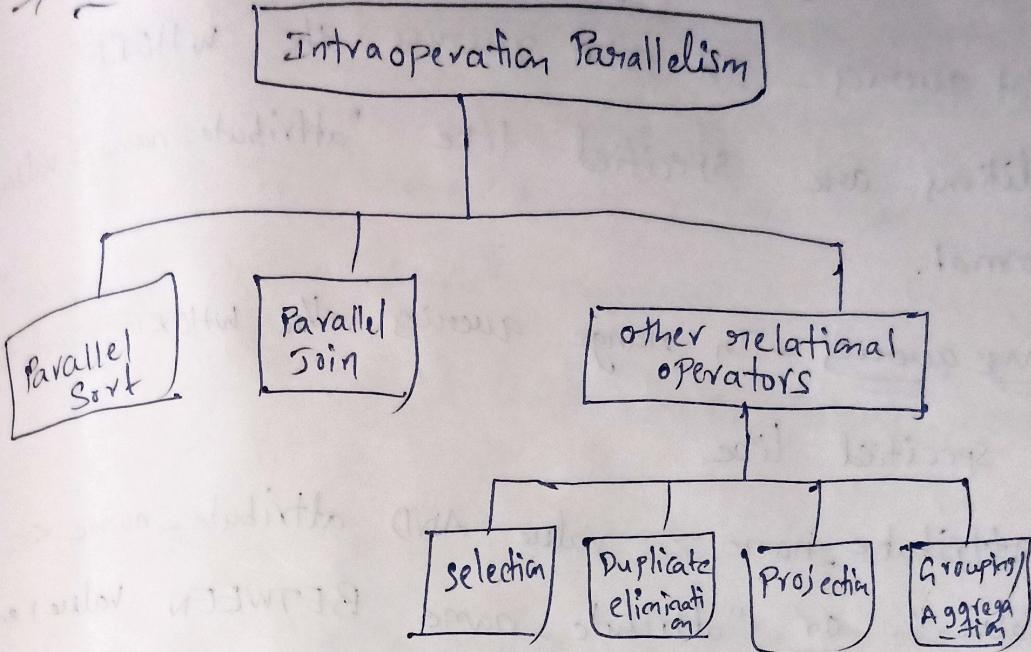
$\text{temp}_2 \leftarrow \text{temp}_1, \sigma_3$ and send the result temp_2 to

processor P_3 to perform $\text{result} \leftarrow \text{temp}_2, \sigma_4$.

The advantage is, we do not need to store the intermediate results, and instead the result produced at one processor can be consumed directly by the other. Hence, we should start receiving tuples well before P_1 completes the join assigned to it.

Join statement is mainly used to combine two tables based on a specified common field b/w them.

→ other relational operators



→ Relation: Employee

- Assume that table is Partitioned and stored in disk D_0, D_1, \dots, D_{n-1} with Processors P_0, P_1, \dots, P_{n-1} .
- Employee is Partitioned using Range Partitioning technique.

EMPID	EName	Salary
E101		27000
E102		22500
E103		80000
E106		270080
E113		20000
E111		31500
E109		26500
E110		22500
E212		40500

D_0

EMPID	EName	Salary
E112		22500
E113		20000
E110		22500

D_1

EMPID	EName	Salary
E101		27000
E106		27000
E109		26500

D_2

EMPID	EName	Salary
E103		80000
E111		31500
E112		40500

Recall the 3 types of queries based on data access

- Point queries - In point queries, the WHERE conditions are specified like "attribute-name = value" format.
- Range queries - In range queries, the WHERE conditions are specified like "attribute-name >= value AND attribute-name <= value". or "attribute-name BETWEEN value1 and value2."
- Scanning the entire relation - In queries where the WHERE condition involves non-key attributes, then the system has to look for all the data in the specified table. This is called the relation scan.

Parallelizing the Selection operation

Consider the following general syntax of any SQL select query,

```
SELECT list_of_attributes FROM table_name  
WHERE condition Ø;
```

This query can be parallelized based on the condition Ø given in the WHERE clause.

selection $\sigma_Q(r)$

- Case 1: If Q is of the form $a_i = v$, where a_i is an attribute and v is a value.
- If Q is Partitioned on a_i → the Selection is performed at a single Processor.

SELECT * FROM emp WHERE Salary = 2000

- Case 2: If Q is of the form $l \leq a_i \leq u$ (i.e., is a range selection) and the relation has been range-partitioned on a_i :
- Selection is performed at each Processor whose Partition overlaps with the specified range of values.
- SELECT * FROM emp WHERE Salary BETWEEN 25000 and 30000.

- Case 3: other cases: the selection is performed in parallel at all the Processors.

SELECT * FROM emp WHERE Ename = 'Sawyer'

Duplicate elimination

- Perform by using either of the Parallel sort techniques
 - Eliminate duplicates as soon as they are found during sorting.

→ can also partition the tuples (using either range or hash Partitioning) and perform duplicate elimination locally at each processor.

Projection:-

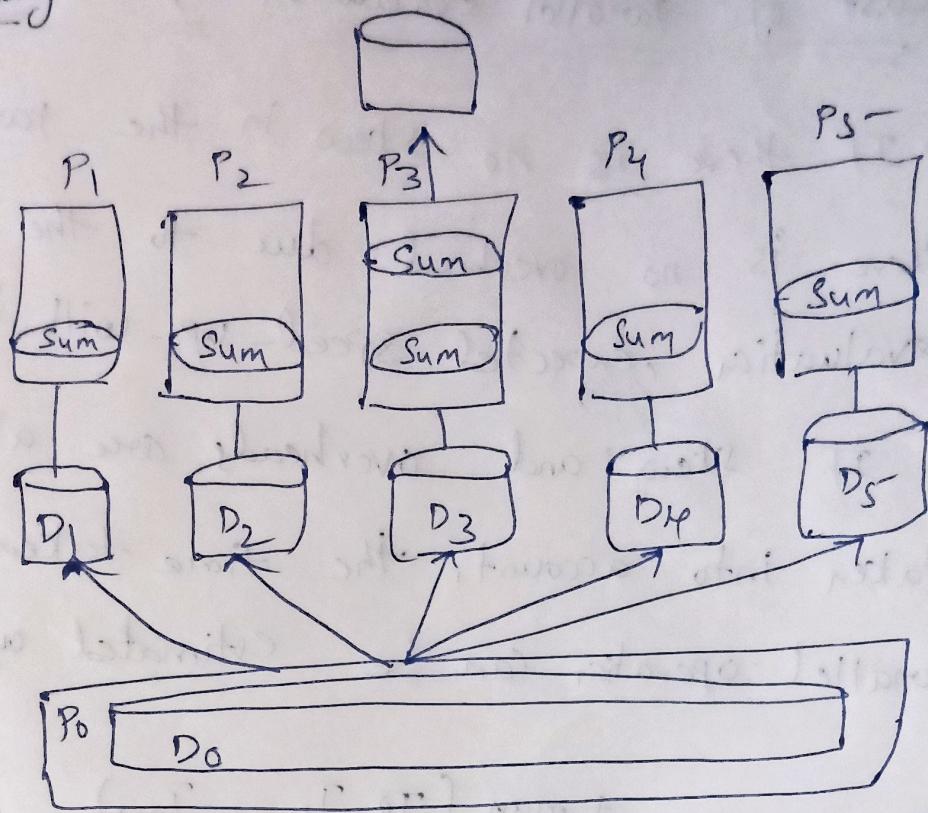
- Projection without duplicate elimination can be performed as tuples are read in form dist in parallel.
- If duplicate elimination is required, any of the above duplicate elimination techniques can be used.

→ Grouping / Aggregation

- Partition the relation on the grouping attributes and then compute the aggregate values locally at each processor.
- can reduce cost of transferring tuples during Partitioning by partly computing aggregate values before Partitioning.

Partly computing

Partly computing



Grouping
partitioning

→ Consider the 'sum' aggregation operation:

→ Perform aggregation operation at each processor

P_i on those tuples stored on disk D_i.

→ results in tuples with Partial sums at each processor.

→ Results of the local aggregation is partitioned on the grouping attributes.

→ Aggregation performed at each processor P_i to get the final result.

→ Fewer tuples need to be sent to other processors during partitioning.

Cost of Parallel evaluation of operations

- If there is no skew in the partitioning and there is no overhead due to the parallel evaluation, expected speed-up will be $\frac{1}{n}$.
- If skew and overheads are also to be taken into account, the time taken by a parallel operation can be estimated by

$$T_{\text{part}} + T_{\text{asm}} + \max(T_0, T_1, \dots, T_{n-1})$$

- T_{part} is the time for partitioning the relation
 - T_{asm} is the time for assembling the results
 - T_i is the time taken for the operation at processor p_i
- this needs to be estimated taking into account the skew and the time wasted on contentions.