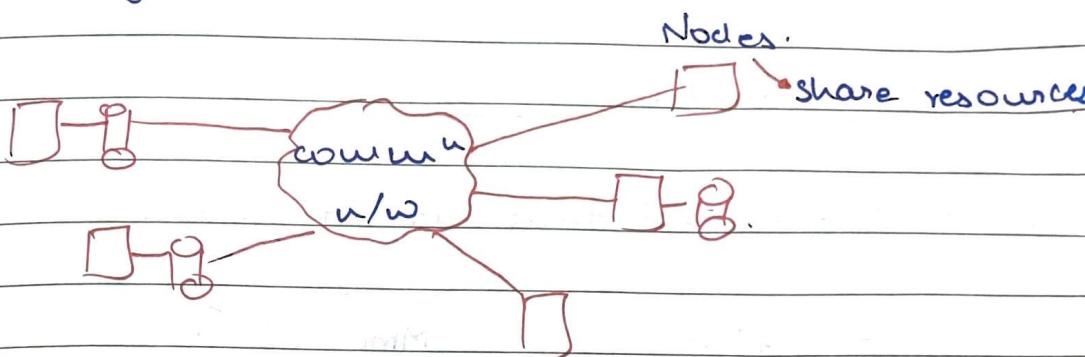


## Distributed Systems

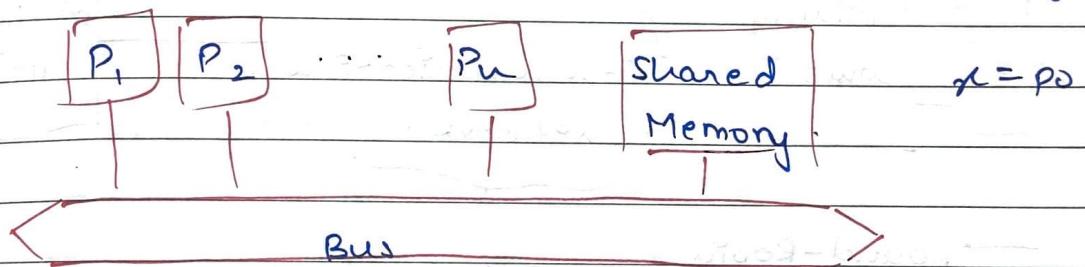
→ It is a collection of independent computers that appears to the users of the system as a single coherent system.

Eg.



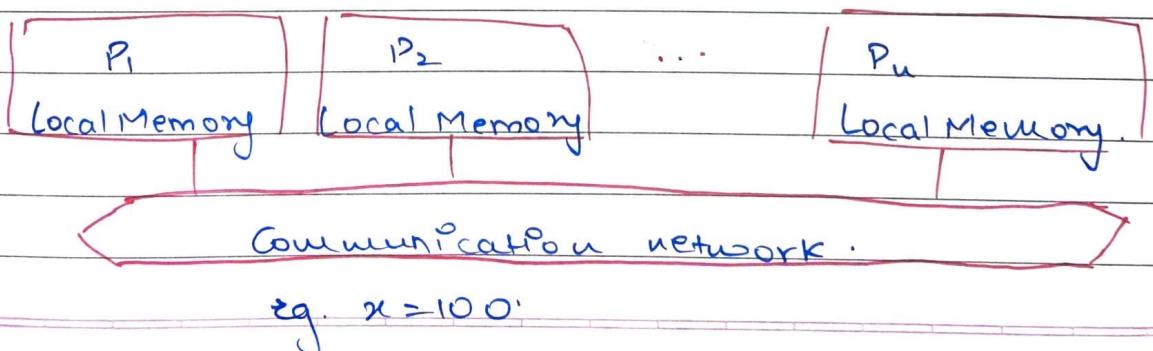
## Types of architecture

1) Tightly coupled architecture → Shared memory.



→ Not scalable due to common bus.

2) Loosely coupled architecture → shared Nothing.



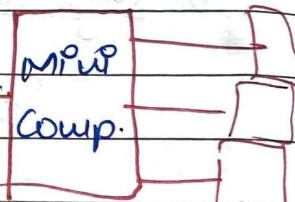
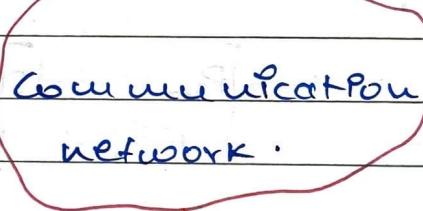
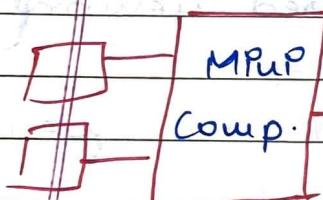
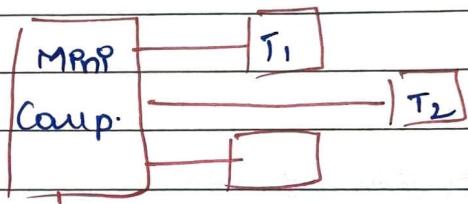
## Distributed computing system models

- 1) Minicomputer model
- 2) Workstation model
- 3) Workstation server model
- 4) Processor model.
- 5) Hybrid model.

## MPU computer model

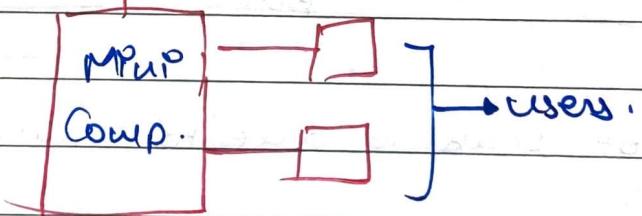
→ Centralized time sharing system.

+ → +time

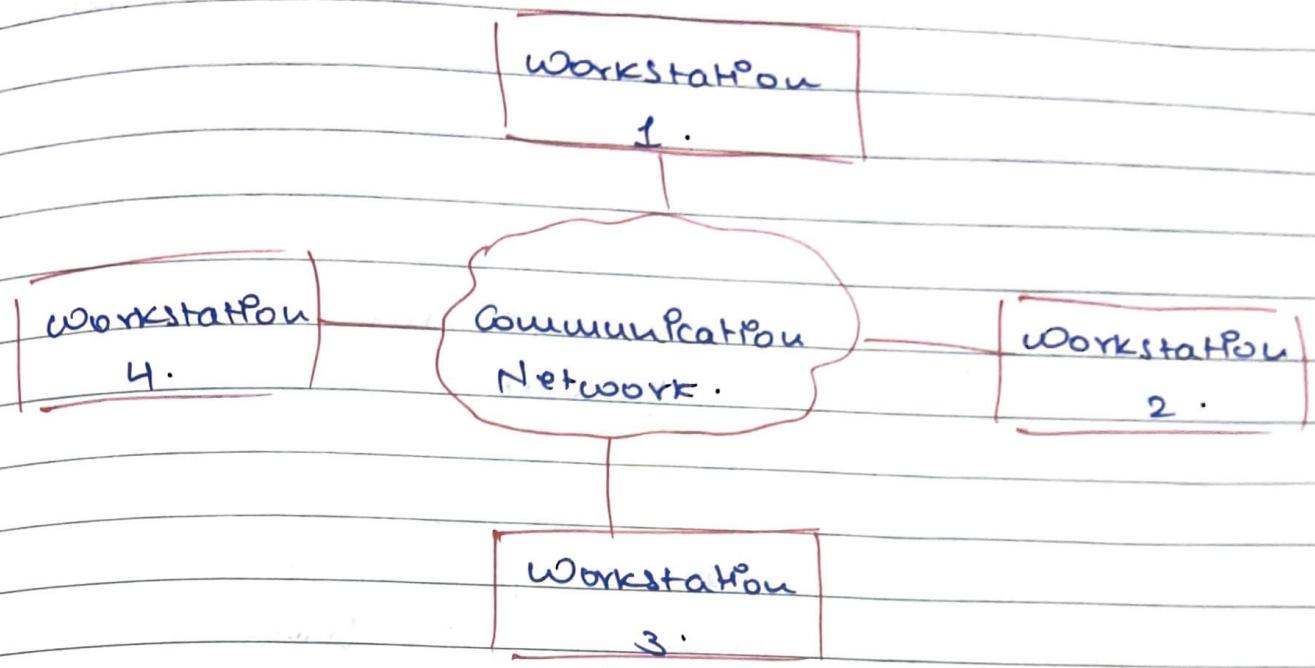


→ Round-Robin

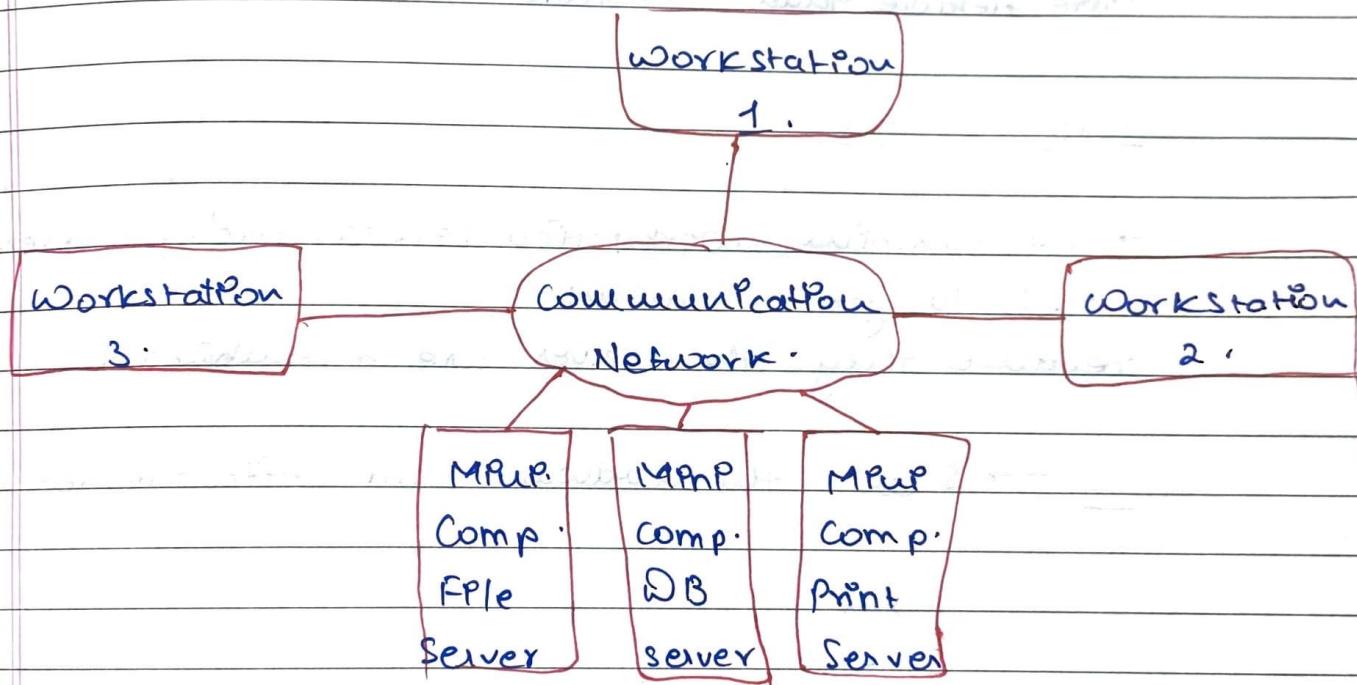
Algorithm.



## Workstation Model

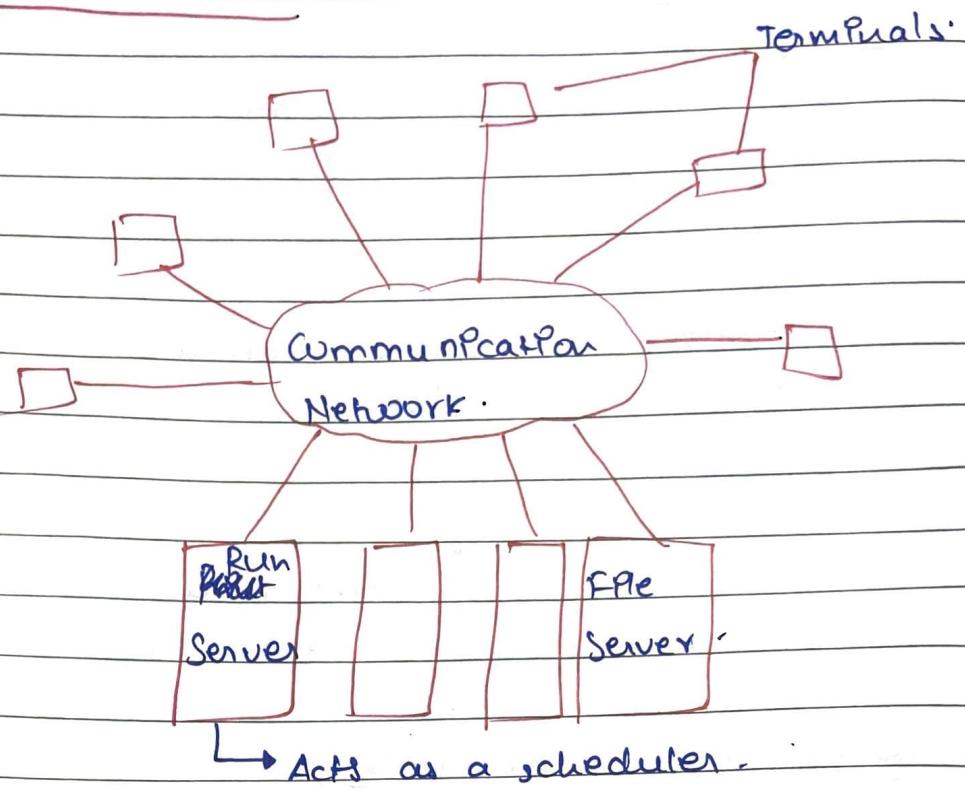


## Workstation Server model



- ① Simple planned server model.
- ② Works as a client server model.
- ③ Easy to manage resources

### Processor Pool model



\* More flexible than workstation model.

### Hybrid model

- We combine workstation server and processor pool models to get this model
- Because they themselves are a combination.

use → Sharing the hardware and software resources

## Advantages of DC

- 1) inherently distributed applications → locally as well as globally the data is distributed.
- 2) Information sharing among distributed users.
- 3) Resource sharing → e.g. Printer.
- 4) Better price performance ratio.
- 5) Shorter response time and higher throughput.
  - In workstation model, ideal resources can be shared.
  - Can distribute the tasks to lightly loaded machine according to requirements.
- 6) Higher reliability and availability.
  - because of replication
- 7) Extensibility & Incremental growth.
- 8) Better flexibility.

## Disadvantages of DC

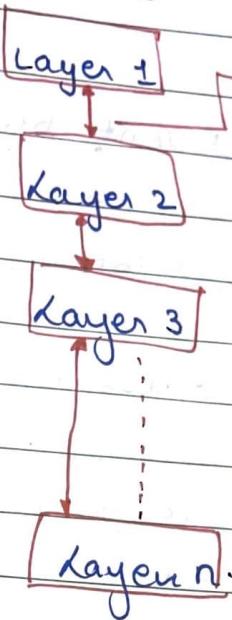
- 1) Higher maintenance if the systems are not high-end.
- 2) Security breach.

## Issues in Designing the distributed systems

- 1) Transparency → <sup>underlying</sup> Hiding details of the system.  
→ Why? → It appears as a single system.
- a) Access → Hide details of data coming from various nodes
- b) Location → Hide details of location
  - Naming → PName given to resource is not known
  - Migration → Resources moving from one location to another is not known.
- c) Relocation → Hide details of why the resources are relocated from one server to another.
- d) Concurrency → Hiding the no. of users using the servers.
- e) Failure → Hide failure & recovery of resources
- f) Persistence → Hide the software resource in memory or disk. (Permanent)
- Q. Identify and explain transparency for the given scenario.

## Architecture Styles

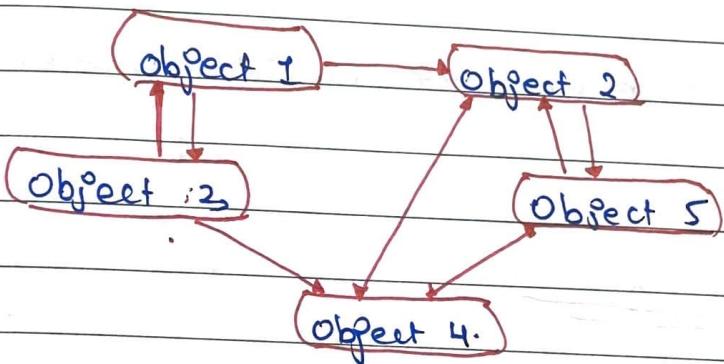
1) Layered Architecture → e.g. OSI model.



W<sup>o</sup>s advantage:

Each layer can communicate with the next layer only

2) Object-based style.

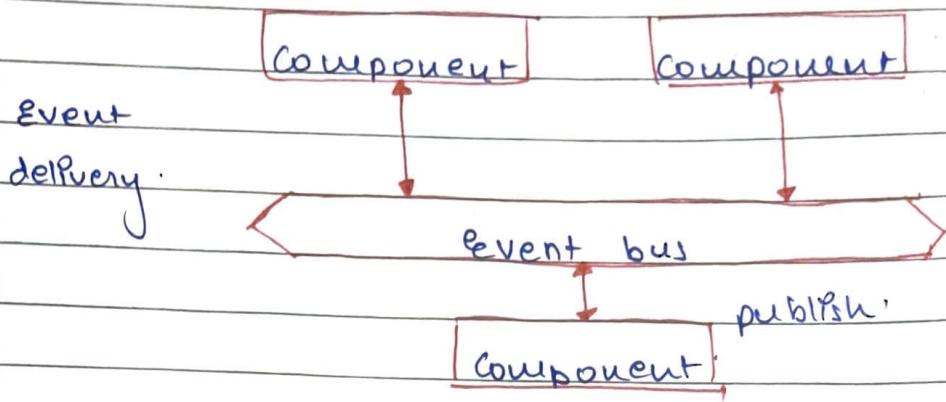


- No restrictions.

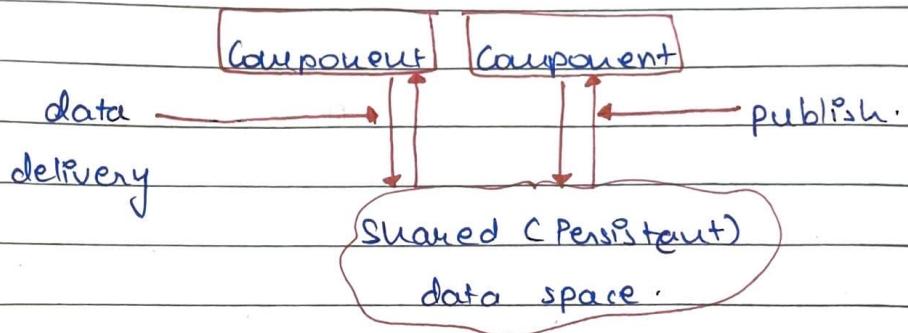
e.g. Client server system.

3) Event based.

- 3) Event based can be used at a particular time only.



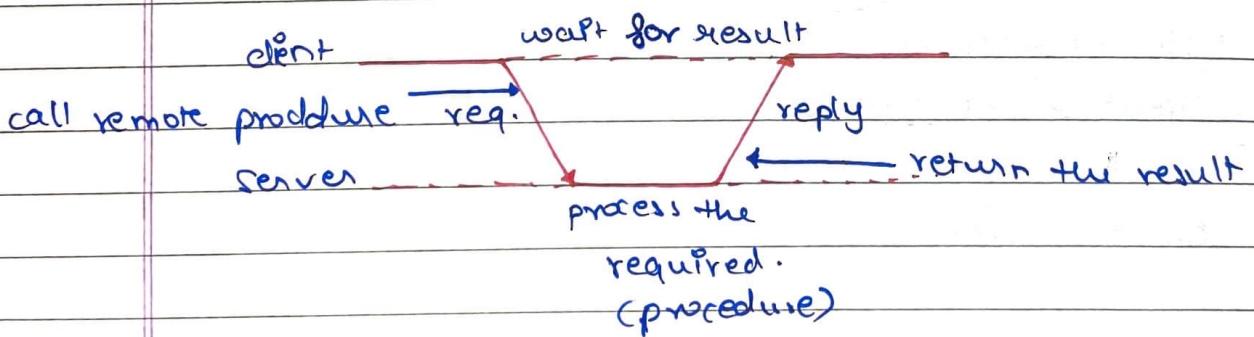
- 4) Shared data space → Anybody can use anything.



Bulletin board / Blackboard architecture → another name.

Eg. Blackboard.

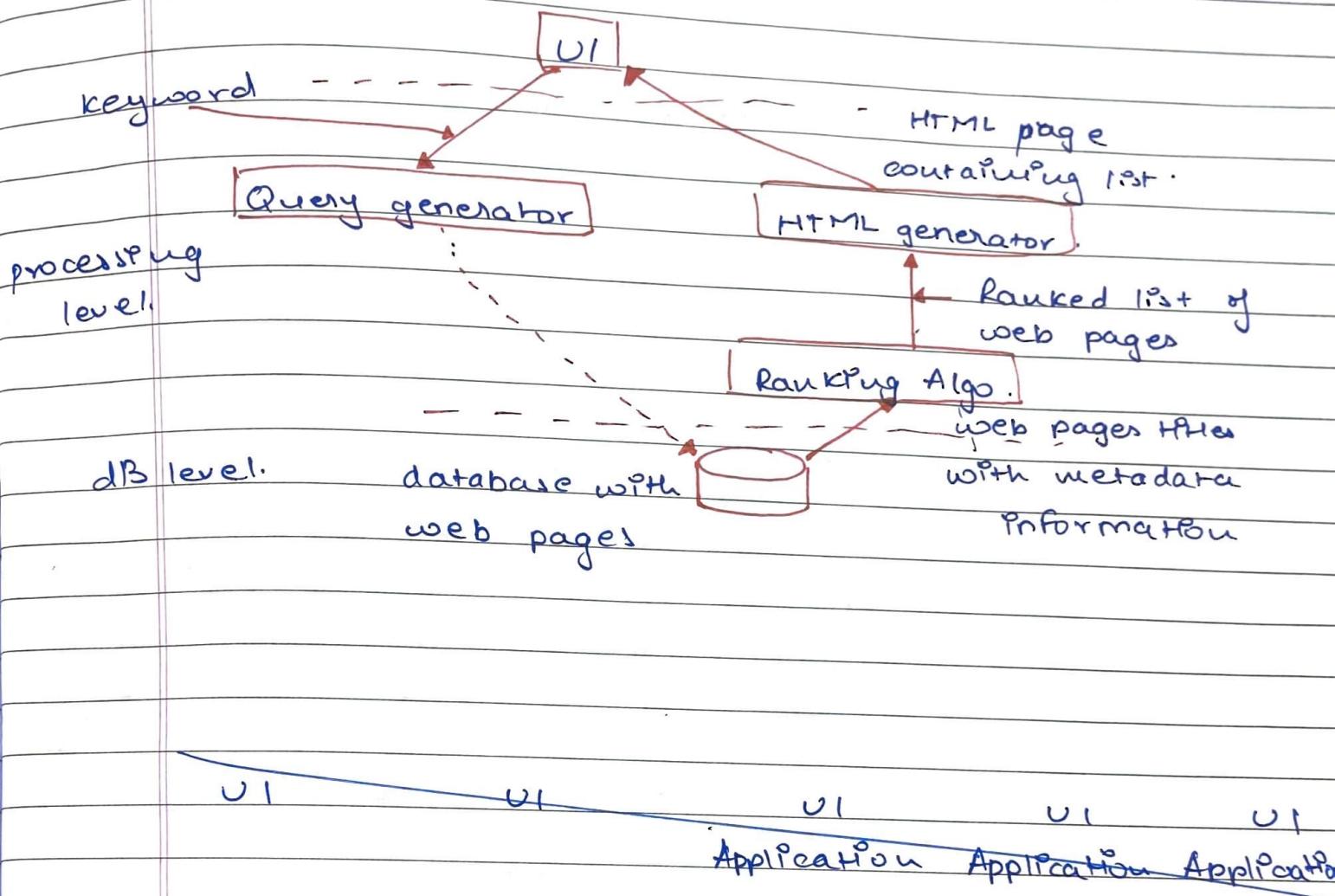
### Client server Architecture



## Application Layering

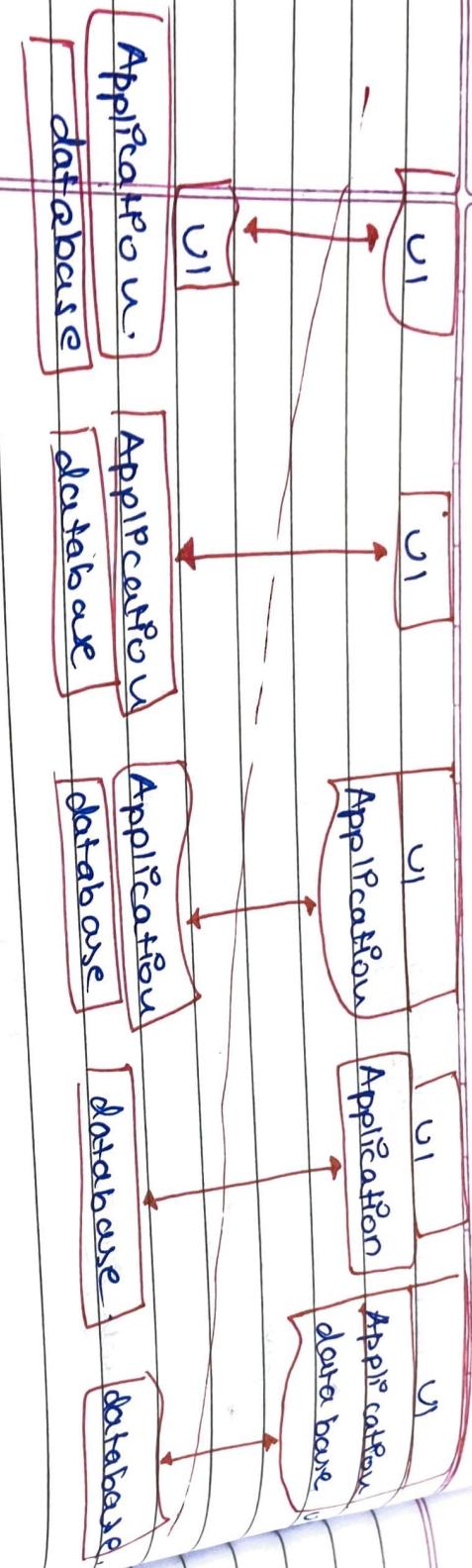
- UI
- processing
- database

eg. Search Engine



## Client machine

PAGE NO. / /  
DATE / /



## Server machine

remote Procedure call  
So pt looks like a  
coherent system

warn()

procedure  
call

value

machine

to exec

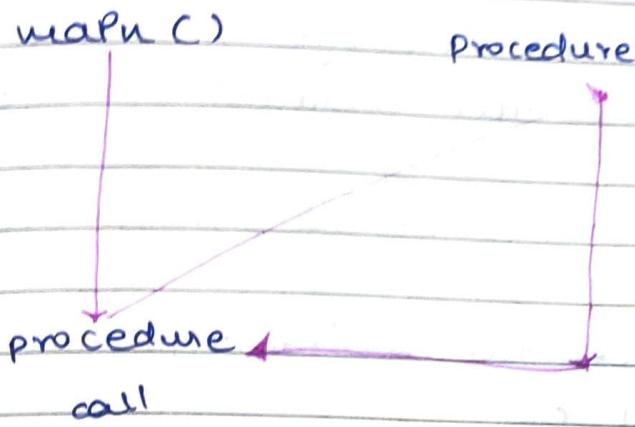
Inform

ru

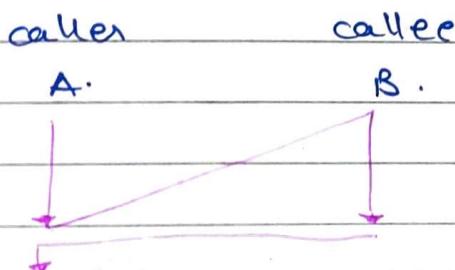
proc

## Remote Procedure Call (RPC)

- So Pt looks like centralized computing (single coherent system).



► When process on machine A calls a procedure on machine B, the calling process on A is suspended, & execution of called procedure takes place on B. Information can be transported from caller to callee in the parameters & can come back in the procedure result. No message passing is visible to the programmer. This is known as RPC.



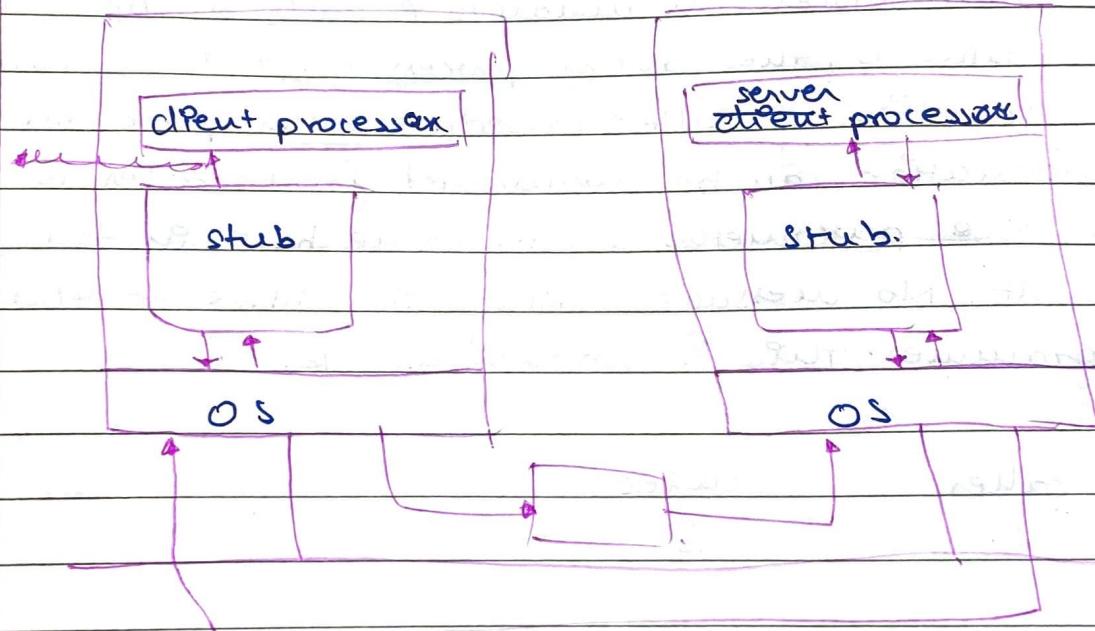
Stubs → creation of packets.  
 Runtime RPC → Local compiler.  
 parameter marshalling → packing the parameters.

PAGE NO. / / /

DATE / / /

### Steps of RPC

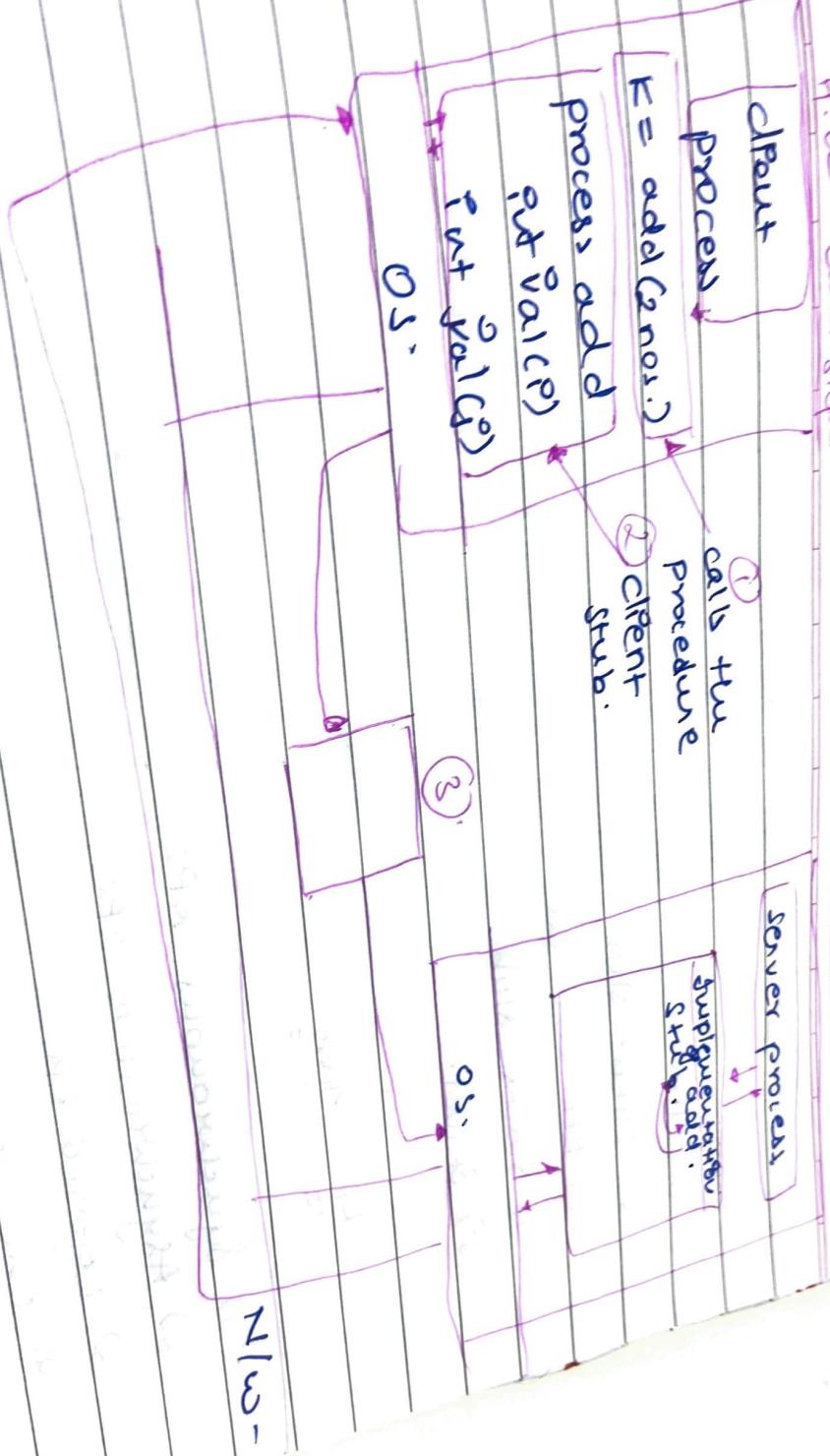
- 1) Client procedure calls the client stub.
- 2) The client stub builds the message and calls the local OS.
- 3) The client OS sends a message to remote OS.
- 4) The remote OS gives message to server stub.
- 5) Server stub unpacks the parameter P calls the server.
- 6) The server does the work & returns the result to the stub.
- 7) The server stub packs the <sup>result</sup> message & calls the local OS.
- 8) Server OS sends the message to client OS.



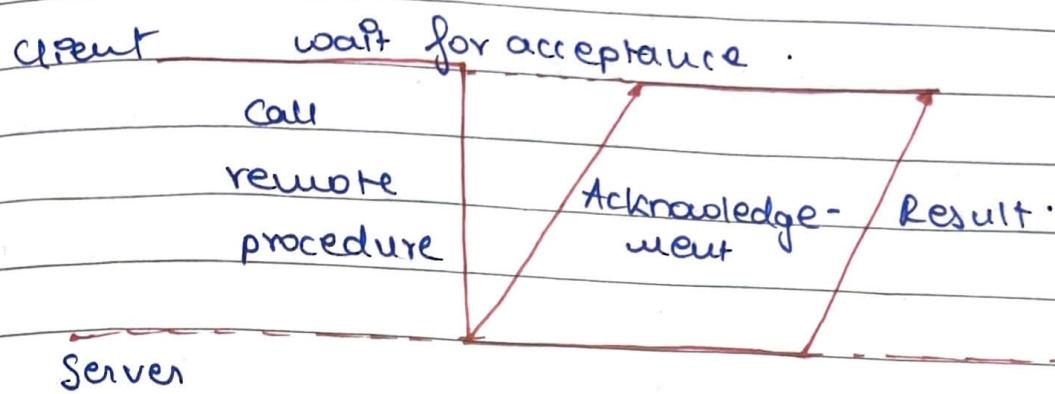
- 9) Client OS gives the message to the client stub.
- 10) Stub unpacks the result & returns to the client.

8. Add 2 numbers.

pp. 10-11 add ships

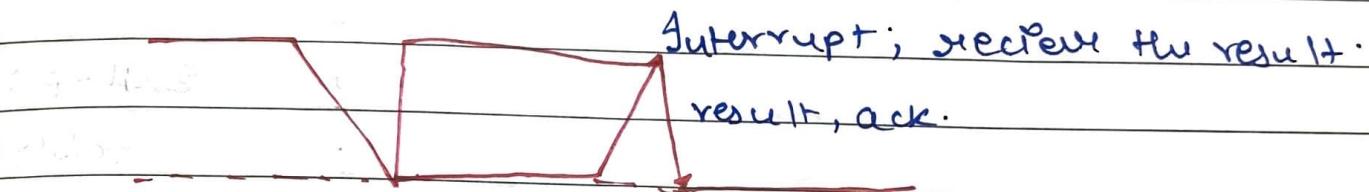


Asynchronous RPC → waiting for client



Deferred RPC → no wait time

→ combining two asynchronous RPCs.



Eg. OTP verification.

One-way RPC

→ Server does not wait for acknowledgement.

① Port mapping services → to map port to the server.

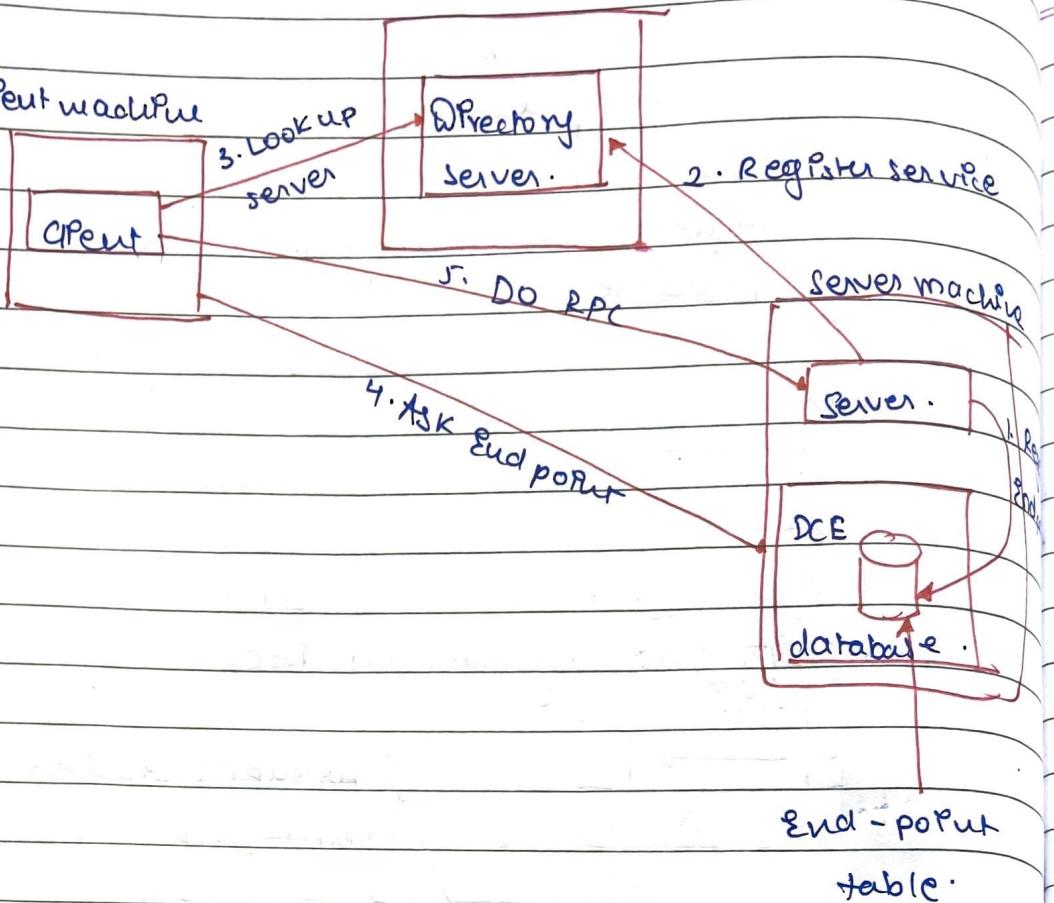
PAGE NO.	/ /
DATE	/ /

Q.

## Binding Client & Server

Directory machine

Client machine



## Failure semantics

### 1) Server failure semantics

→ At least one semantic → If the client received the reply from the server, it means that all has been executed at least once on the server.

→ At most once → If the client gets the reply from the server it means that the RPC call have been executed at most once on the server.

→ Exactly once → Same exactly once.

e.g. Gmail.

## Client failure semantics

- If server executed RPC call, but meanwhile client crashed, server will discard the response
- If server was executing long running RPC & in both executions it comes to know that client has crashed then such calls are called Orphaned RPC calls
- In this, server can do the following:-
  - 1) Extermination → when client gets crashed, server can eliminate the process.
  - 2) garnation. → when the server broadcast message to know if the client is alive.

## Light-weight RPC

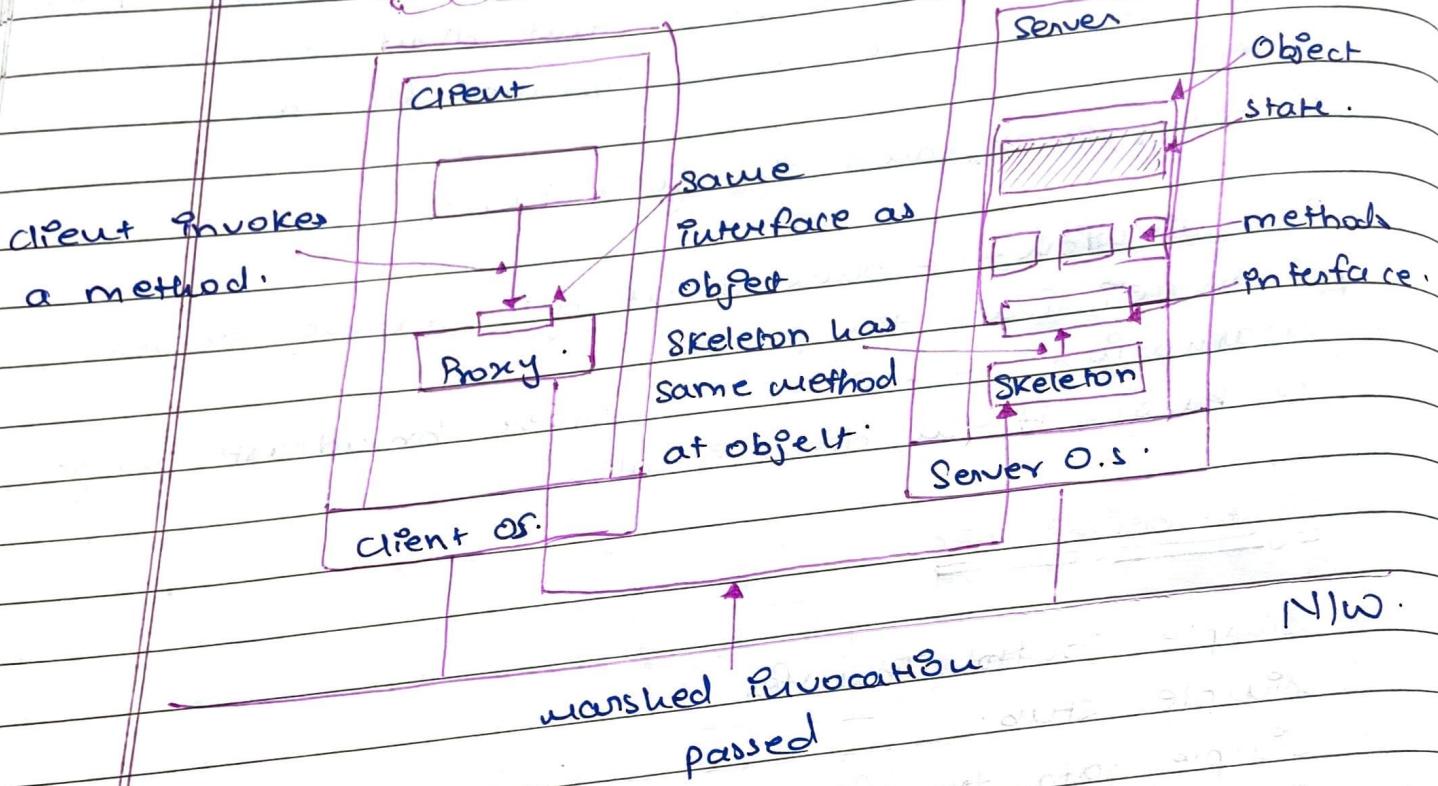
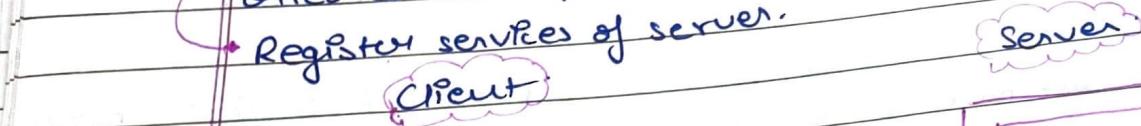
- 1) Simple control Transfer
- 2) Simple stub.
- 3) Simple data transfer
- 4) Design and concurrency

How to make RPC flexible?  
 → Layered approach → by making it reliable  
 Problems in RPC → accessing global variables

replication.  
 PAGE NO. / /  
 DATE / /

## Remote Object Invocation

- RPC → procedure oriented concept
- Register services of server.



- Skeleton → marshalling  
 Proxy → unmarshalling.
- Persistent object → requires secondary storage.  
 Transient object → exists during runtime only  
 → temporary.

How to get the server up-to-date → clock-synchronisation

Binding of object :- 1) Implicit.  
 2) Explicit



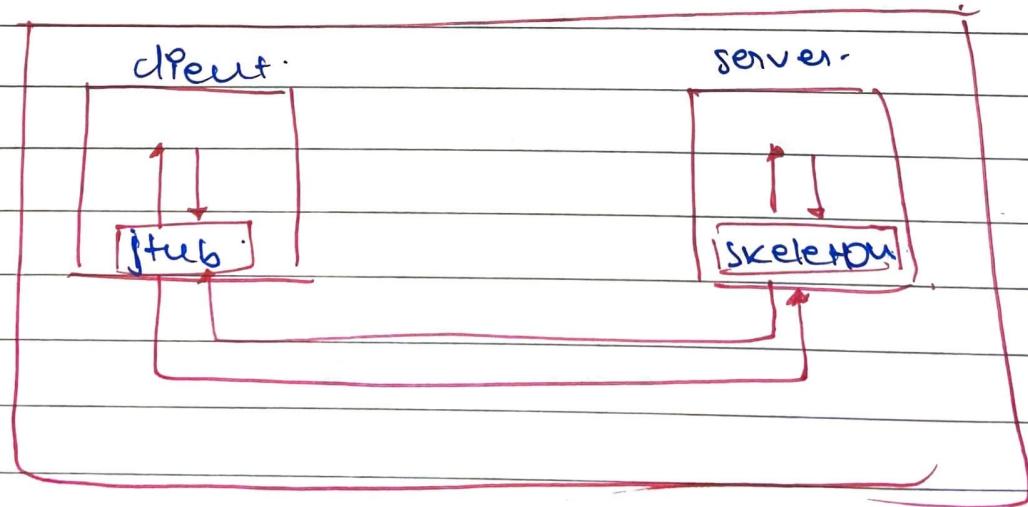
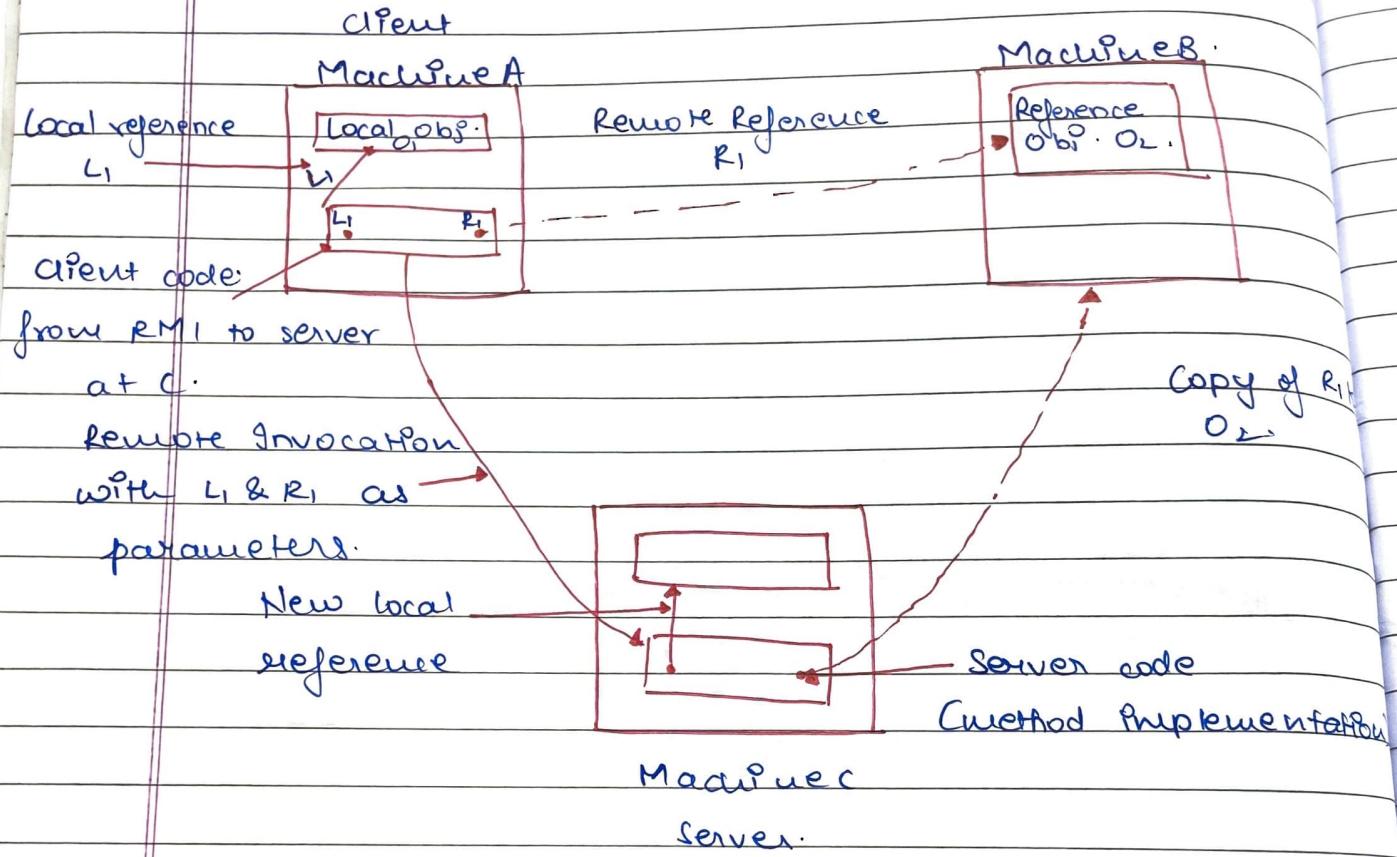
FCE → Distributed computing environment → Locating the reference object.

PAGE NO. / / / /  
 DATE / / / /  
 Java implementation of RMI.  
 objects in different programs can communicate

## RMI (Remote Method Invocation)

Parameter passing

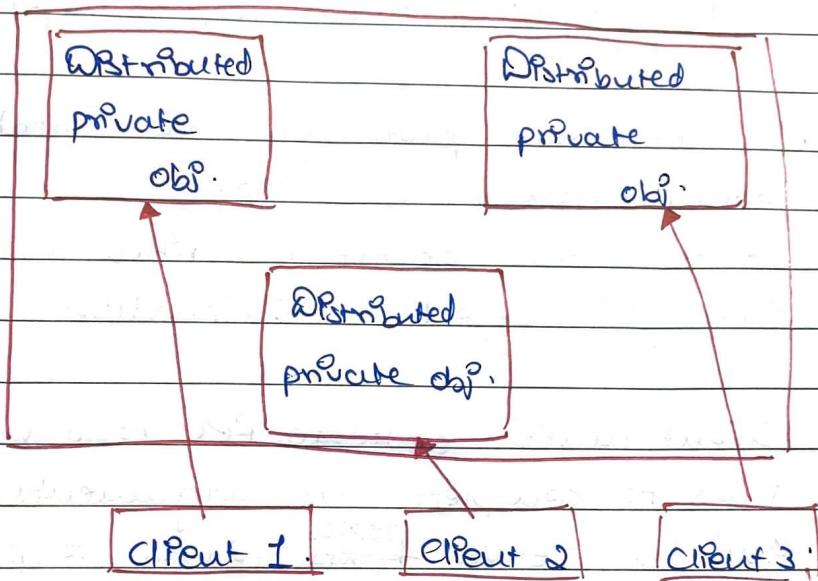
→ System-wide object reference.



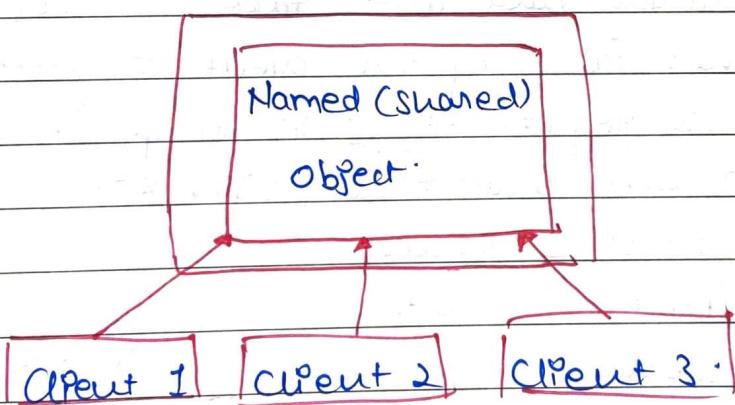
## The DCE Distributed Object Model

- Distributed of dynamic object → server creates on behalf client
- Distributed named object

Server machine



Server machine



PAGE NO. / / /  
DATE / / /

## Method

### Remote Invocation

#### RPC vs RMI

- RPC → procedural programming.
- RMI → object based.
- RPC → parameters are ordinary data structures.
- RMI → possible to pass objects as parameters
- If marshalled parameters are local → object serialization is used. (Copy).
- State of the object is written to a byte stream through object serialization.

Q. Client makes a remote RPC to a server. Client takes 5ms to complete the arguments for each request. And the server takes 10ms to process each request. Local OS's processing time for each send or receive operation is 0.5 ms and the network time to transmit each msg or reply each msg is 3ms. Marshalling and unmarshalling takes 0.5ms per msg. Calculate the time taken by the client to generate a return from 2 requests if it is single threaded. You can ignore the context switching time.

Request → 5ms.

Marshalling → 0.5"

Client → 0.5"

Network → 3 → 0.5 (Server)  
↓ (Delete)

↓ 0.5 (Server)  
↓ (n)

↓ 0.5  
↓ (n)

↓ 0.5

↓ 0.5



## Steps

- 1) Create the
- 2) Provide the
- 3) Compile & skeleton
- 4) Start the
- 5) Create
- 6) Create

1)

2)

## Steps for writing the RMI program.

- 1) Create the remote interface.
- 2) Provide the implementation of remote interface.
- 3) Compile the implementation class & create the proxy & skeleton objects using rmic tools.
- 4) Start the registry services by rmiregistry tools.
- 5) Create & start the remote application.
- 6) Create & start the client application.

1)

```
import java.rmi.*;  
public interface Adder extends Remote  
{  
    public int add (int x, int y).  
    throws RemoteException;  
}
```

2)

```
import java.rmi.*;  
import java.rmi.server;  
public class AdderRemote  
extends UnicastRemoteObject  
implements Adder{  
    AdderRemote () throws  
    RemoteException {}  
    super ();
```

3 .

```
public int add (int x, int y).  
{ return (x+y); }  
}
```

3)

rmic Adder Remote

4)

rmiregistry 5000

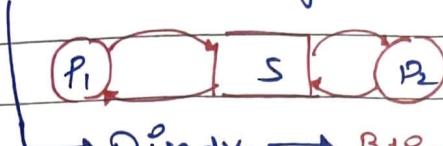
```
import java.rmi.*;  
import java.rmi.registry.*;  
public class MyServer {  
    public static void main (String args[]) {  
        try {  
            Adder stub = new AdderRemote();  
            Naming.rebind ("rmi://localhost:  
                5000/abc", stub);  
        }  
        catch (Exception e) {}  
    }  
}
```

```
import java.rmi.*;  
public class MyClient {  
    public static void main (String args[]) {  
        try {  
            Adder stub = (Adder) Naming.lookup  
                ("rmi://localhost:5000/abc");  
            System.out.println (stub.add (3,4));  
        }  
        catch (Exception e) {}  
    }  
}
```

## Message Passing

→ Interprocess communication (IPC)

→ Shared memory



→ Disadv. → Blind write.

Message format is same

## Features of good message passing

1) Simplicity.

→ Send & receive.

2) Uniform semantics. → Local → same machine.

↳ remote → different machine

sending data and ack. packed. ↳ Reliable data transfer.

3) Efficiency → Reducing no. of messages exchanged.  
 ↳ faster communication. ↳ Reducing no. of acknowledgement.  
 ↳ minimizing cost of communication. ↳ sending one for multiple.

4) Reliability → no packet loss

↳ in-time delivery

↳ no errors

↳ no duplication.

5) Correctness → can be used for group communication

↳ atomicity → all or nothing.

↳ ordered delivery → sequence no.

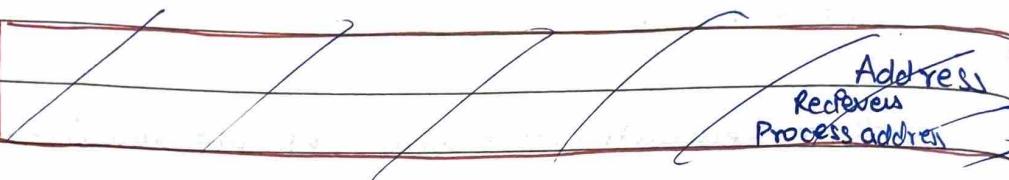
↳ survivability

6) Flexibility → broadcast & protocol  
 multi cast protocol

7) Security → Encryption  
 ↳ Authenticate user.

8) Portability → Transferable to any system.  
 ↳ Should be visible to all users

### Issues in IPC by Message Passing



Actual data or pointer to data.	Structural information	Sequence no. or message ID	Address
	no. of bytes	type	Receiver's process address Sender's process address

variable size → fixed length header.

### Synchronisation

→ Blocking → Synchronous

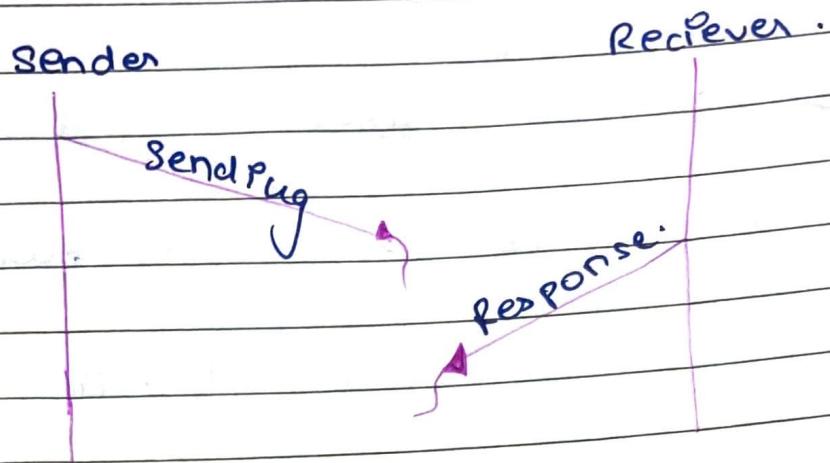
→ Non-blocking → Asynchronous

↳ Receiver Knows? → through notification.

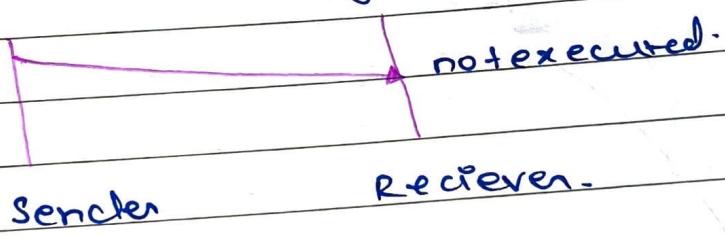
## Buffering

- 1) Null buffer → normal phone calls.
- ~~Advantage~~ → Sender & receiver should be active.
- 2) Single message buffer → synchronous communication  
→ OTP  
→ Message will be in sender as and receiver address space.
- 3) Unbounded capacity buffer → unlimited usage.
- 4) Bounded capacity buffer → limited storage

## Failure handling techniques in message passing.



- 1) Loss of request → receiver is crashed.  
↳ lack of buffer.
- 2) Loss of response.
- 3) Unsuccessful admission of request.



- couple execution, server crashed.

Page No.	
Date	/ /

Protocol:

Sender

Receiver

(1) → blocked state.

Reply

Ack.

processes the  
request

User:

useless because  
process is in blocked  
state for too long.

request

process the request

Reply.

ack.

↳ If acknowledgement is removed it does not remain reliable.

- Idempotency means reapplicability.
- Idempotent operation produces same result without any side effects no matter how many times it is performed with the same arguments. (Consistency).
- Operations that do not necessarily produce same results when executed repeatedly with the same arguments are said to be non-idempotent.

Eg. Loading the website.

Authentication while login.

max, sqrt functions used.

debit (amount).

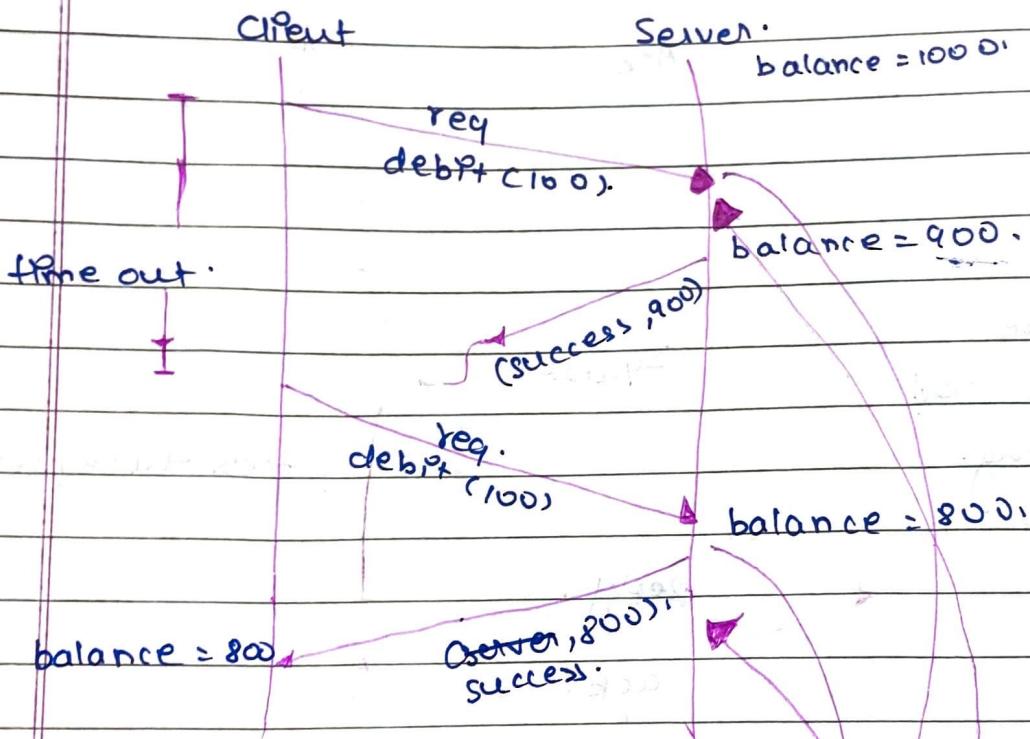
if (balance > amount)

    balance = balance - amount;

    return ("success", balance);

else return ("failure", balance);

end.



To prevent this from happening, cache table is maintained

req. id.	reply to be sent
1.	900.

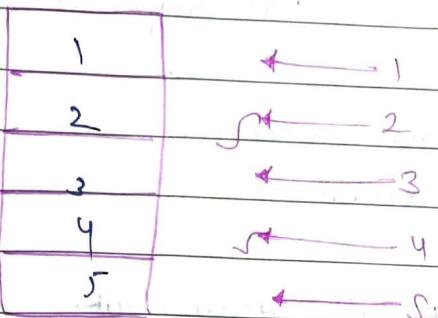
Multipacket flow → Single message with multiple packets.  
Issues → would not be in sequence, packet loss, waiting increased when stop and wait is used.

Reducing ack. → using blast protocol → acknowledgement for all packets

Stop and wait → synchronous communication

Conditions required → packets received in sequence  
Receiver has cache to store the sequence id.

Receiver.



what to do in this case?

The receiver already has a count of the no. of packets.

We have to use selective repeat after timeout.

## Group Communication

- 1) 1 - M
- 2) M - 1
- 3) M - M
- 4) 1 - 1

\* Broadcast → maximum utilization of server, time saving

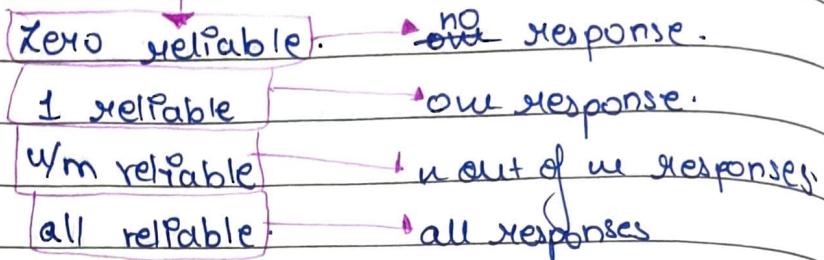
## Group management:

- open group
- close group

## Group addressing

- multicast and broadcast addressing.
- \* multicast → asynchronous communication → since broadcast sender does not wait for acknowledgement.
- both buffered and unbuffered
- shared-space or bulletin board architecture is used.

## \* Flexibility and reliability in multicast communication



2) u to 1 → many servers one client  
→ Eg. moodle server.

3) u to u → Eg. Conference call, group chat.  
→ Order should be maintained.

Assignment → Failure handling, reliability, buffering  
of. IPC, group communication  
for diagrams, use cases, protocols, any application (social media) → study and write in your own words.

Any feature required to add. (Eg. last seen, status should not be seen). 5-6 pages. (10 marks). 5th October.

First week of October.

## Message Oriented Communication

Message queue  
Local buffer

Sending host

Communication  
server

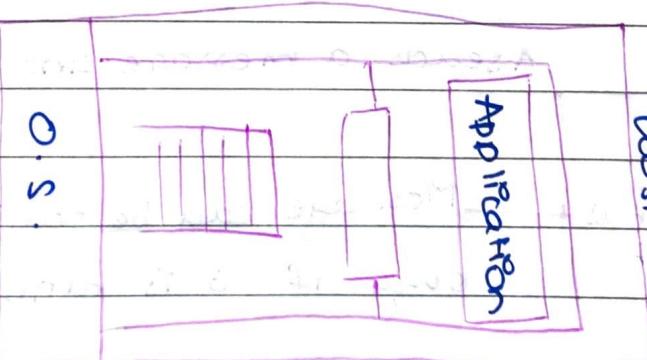
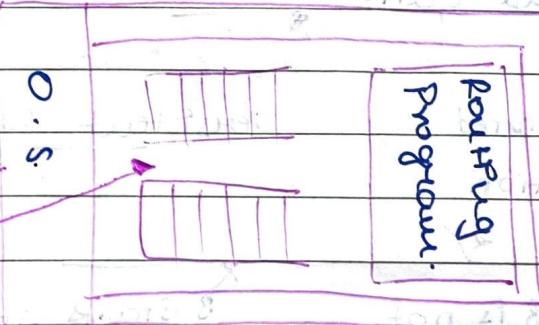
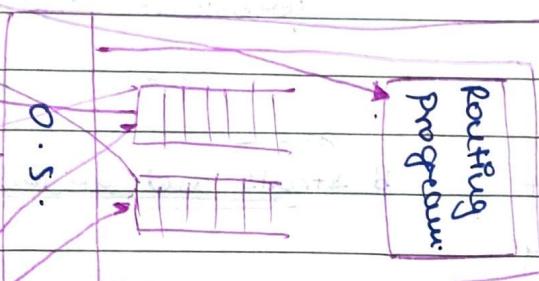
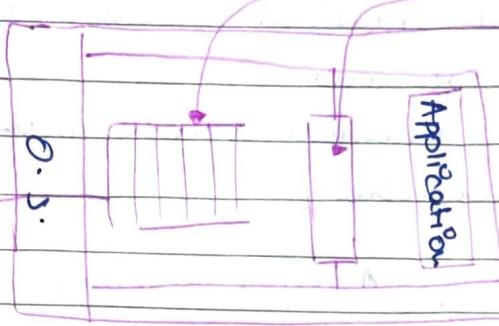
Communication  
server

Receiving  
host

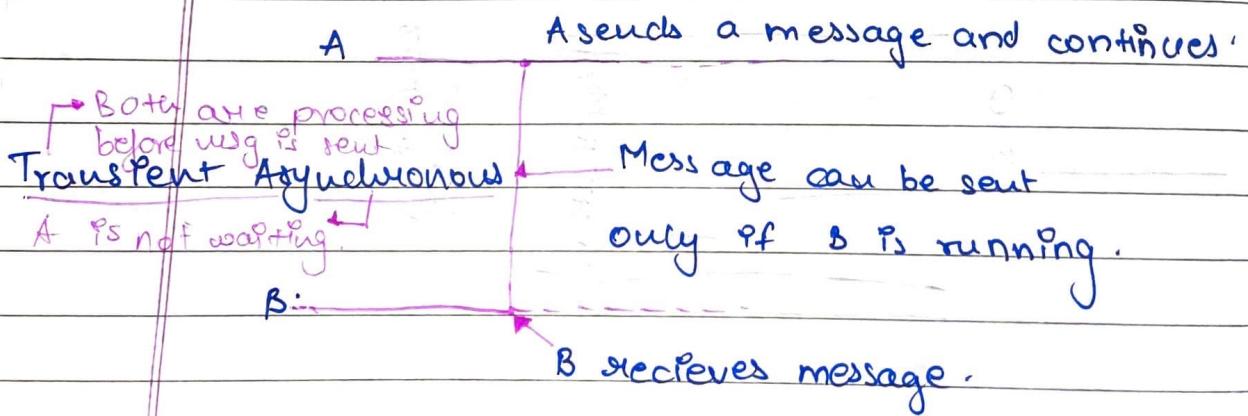
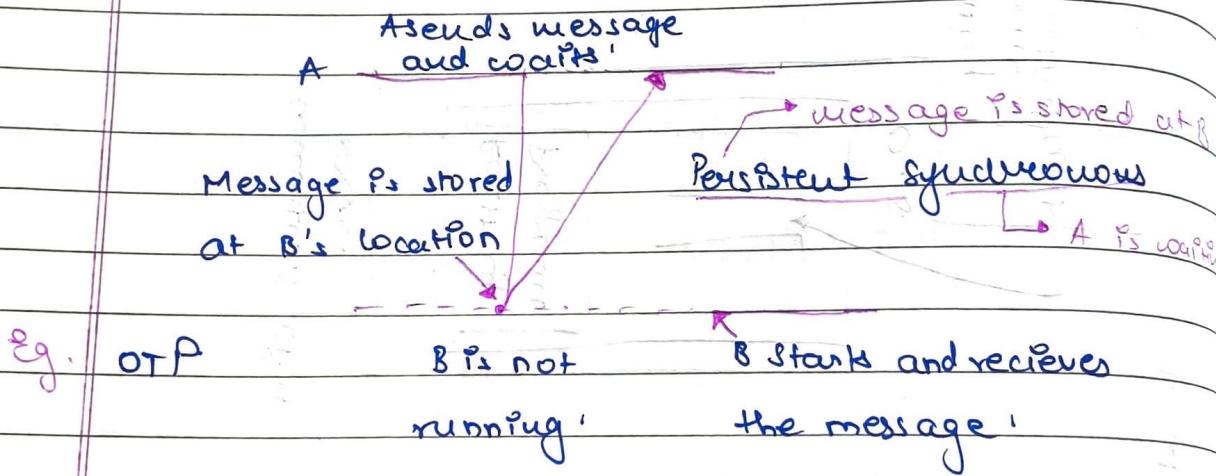
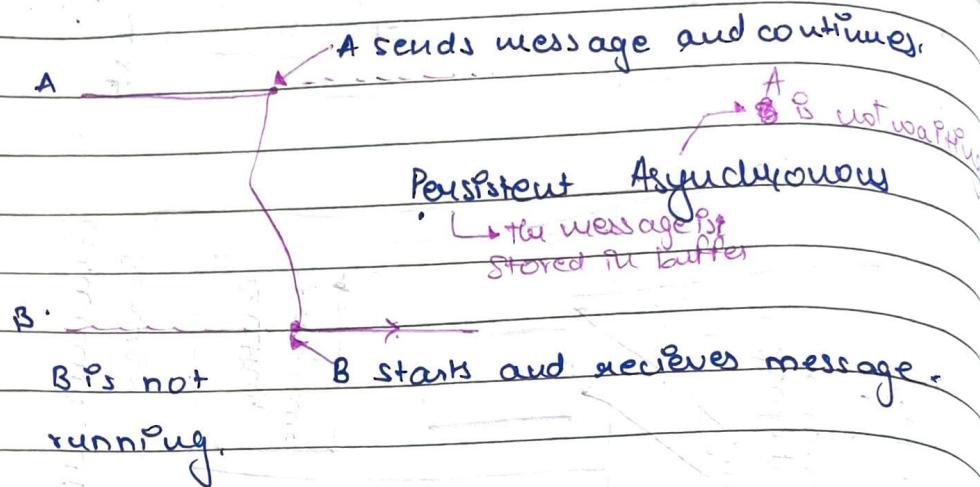
local N/W.

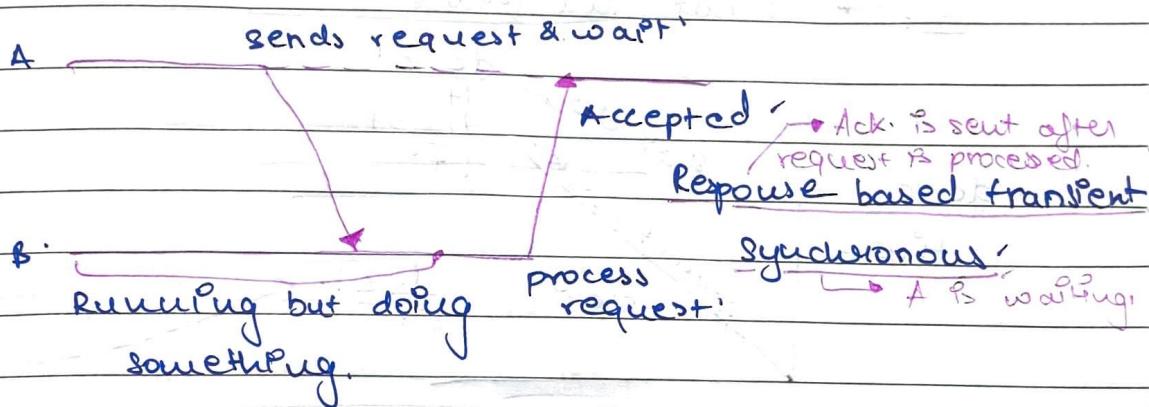
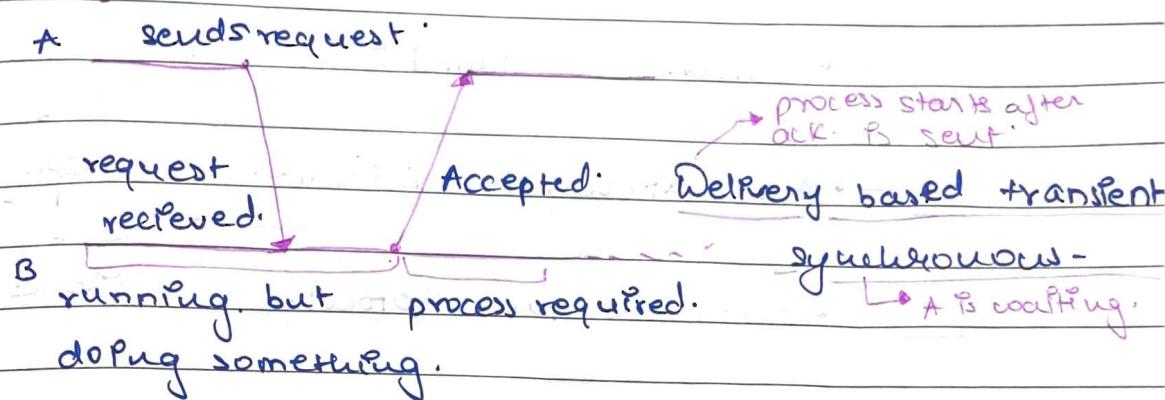
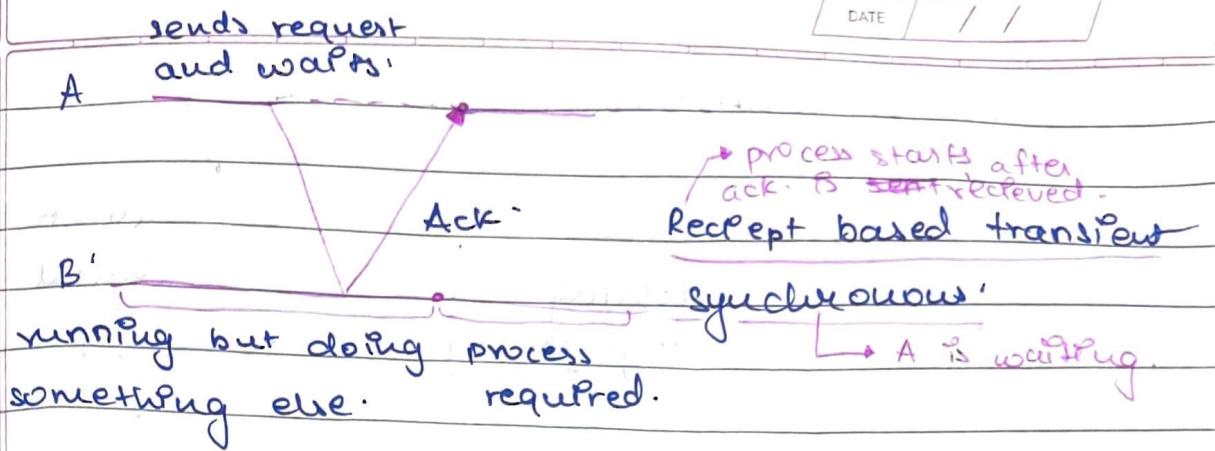
local N/W.

Buffering  
and dependent  
of communication host



- (1) Persistent synchronous  
 (2) Receipt based transient synchronous.  
 (3) Response based transient synchronous.  
 (4) Persistent synchronous  
 (5) Delivery based transient synchronous  
 (6) Transient synchronous.
- (some)*





Q. Explain the different forms of communication with example.

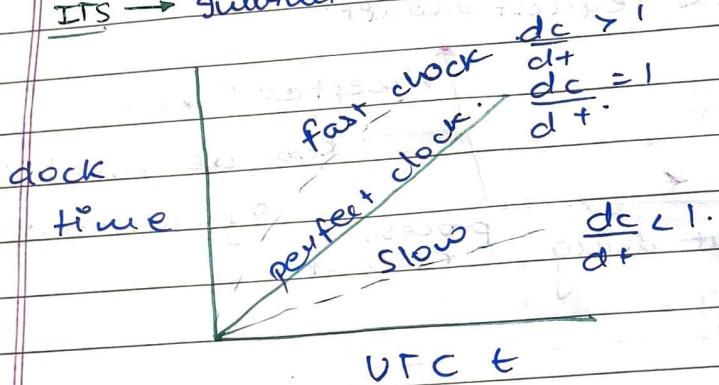
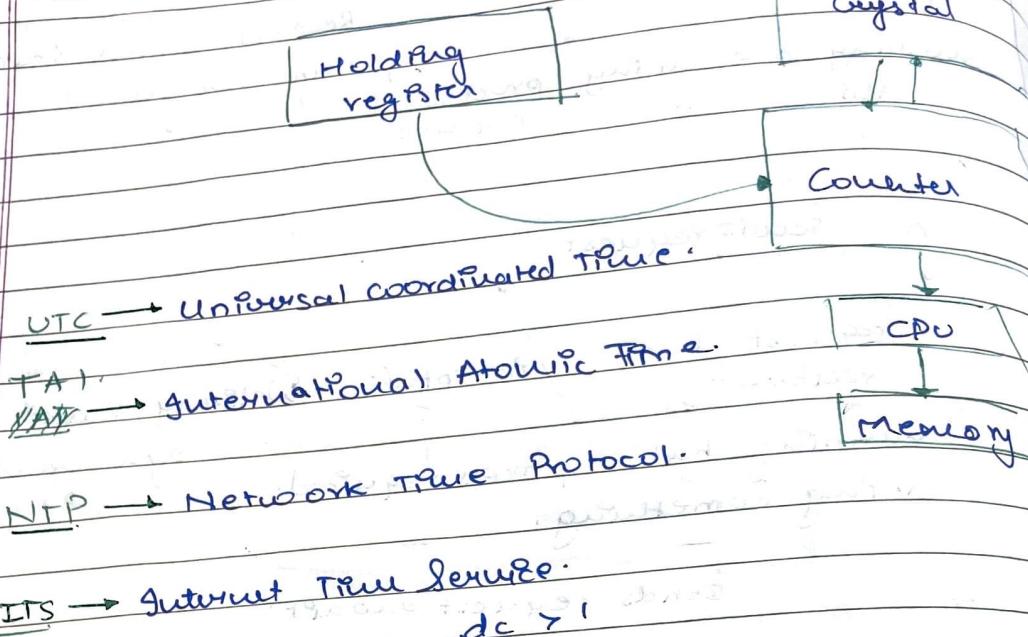
(\*) Persistent → message is stored.

transient → Both sides are active

Synchronous → waits

Asynchronous → does not wait.

## Clock Synchronisation



- ① Cannot adjust manually → the transaction would not stay valid.

② OS provides with the date & time.

$$1 - \delta \leq \frac{dc}{dt} \leq 1 + \delta$$

Centralized → where process can share clock & memory.  
It relies on shared memory.

→ Centralized clock synchronization algorithms are passive time server & active time server.

① Christian's Algorithm

② Berkeley's Algorithm.

③ Christian's Algorithm.

→ Passive time server algorithm.

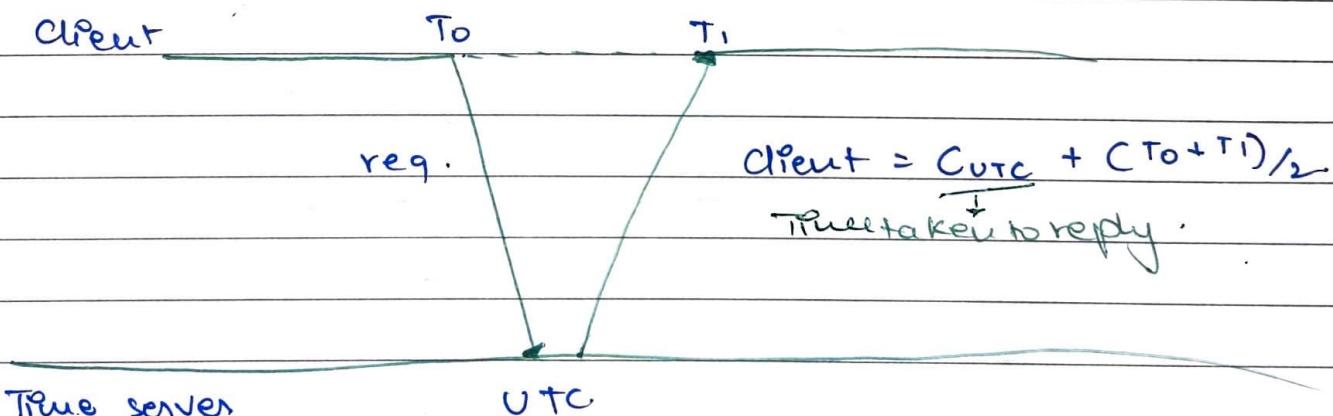
→ users RPC.

→ There is a machine with ~~now~~ receiver which receives precise UTC.

→ Machine sends a request to time server at least  $\delta/2\delta$ , where  $\delta$  is a maximum difference allowed between clock & UTC.

→ The time server sends reply message with current UTC where req receives required request.

→ The machine measures the delay between the server's sending message & machines receiving message. Then it uses the measure to adjust the clock.



clock drift  $\rightarrow$  clocks running at different rates.

PAGE NO.	
DATE	/ /

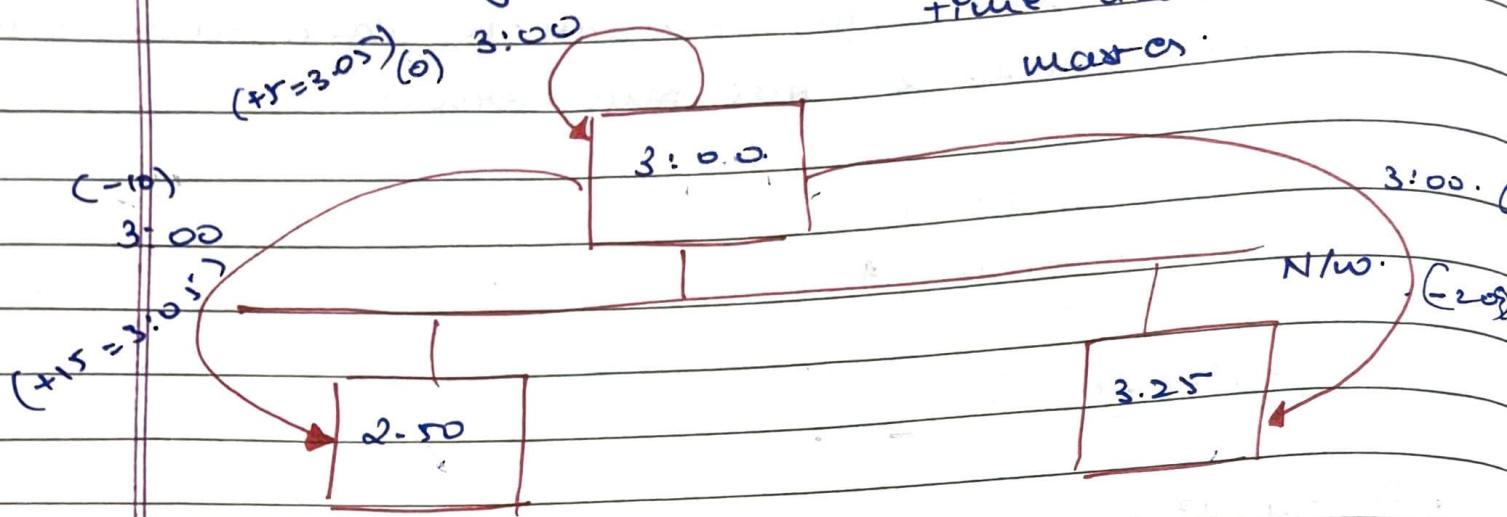
## ② Berkeley's Algorithm

→ used in systems without UTC server.

green time.

+ time daemon

master.



$$25 - 10 + 0 = 15 \quad \frac{15}{3} = 5^{\circ}$$

\* If master crashes, new one is appointed with election algorithm.

→ Bound the higher and lower values to take away

\* Clock synchronisation is done to maintain consistency which is required in any system.

→ real time is not required in clock synchronization

### Logical clock

P<sub>1</sub> 0 1 2 3 4 .

P<sub>2</sub> 0 1 2 3 4

Ordered. → event ordering is required.

\* Lamport clock is used.

\* Concurrent events → same timing events.

### Happened before Relation

- If A & B are the events in the same process, and A is executed before B, then A → B.
- If A represents sending of message & B is a receipt of message, then A → B.
- Relation is transitive
- A → B, B → C then A → C.
- Relation is undefined across the processes that don't exchange messages.

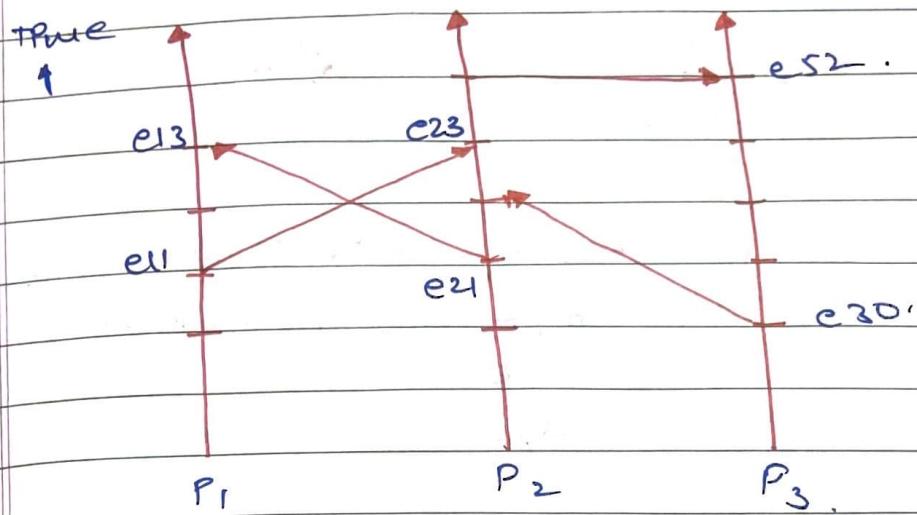
### Event ordering using Happened Before

- Define the notion of time:  
If A → B, then c(A) < c(B).
- If A & B are concurrent,  
then c(A) <, >, =, c(B).
- Each process maintains the logical clock (c<sub>p</sub>)
- whenever event occurs locally at i,  
 $c_i^o = c_i^o + I$

- When I sends message to J,  
piggyback  $L_{CJ}^P$ .
- When J receives message from I,  
if  $L_{CI}^P < L_{CJ}^P$ , then  $L_{CJ}^P = L_{CI}^P + 1$ .  
else, do nothing.

## Lamport Clock

→ If absolute time is not important, this can be used  
or global time.



Eg:  $e_{11} \rightarrow e_{20}$   
 $e_{21} \rightarrow e_{30}$  ] → concurrent events.

→ Used to store in the database and fetch it.

## Vector clock

→ Each process  $i$  maintains a vector  $V_i$

- $V_i[i]$ : number of events that have occurred at  $i$ .
- $V_i[j]$ : number of events we know that have occurred at  $j$ .

→ Update the vector clock as follows

- Local event: increment  $V_i[i]$ .

- Send message: piggyback entire vector  $V$ .

- Receipt of message  $V_j[k] = \max(V_j[k], V_i[k])$

- a) Receiver  $i$  is told about how many events the sender knows occurred at another process  $k$ .

Also,  $V_j C_i^j = V_0 C_i^0 + 1$ .

$(1, 0, 0)$   $(0, 2, 0, 0)$

$(0, 0, 0)$   $P_1$  a b

$(2, 1, 0)$

$(0, 0, 0)$   $P_2$

c d  $(2, 2, 0)$

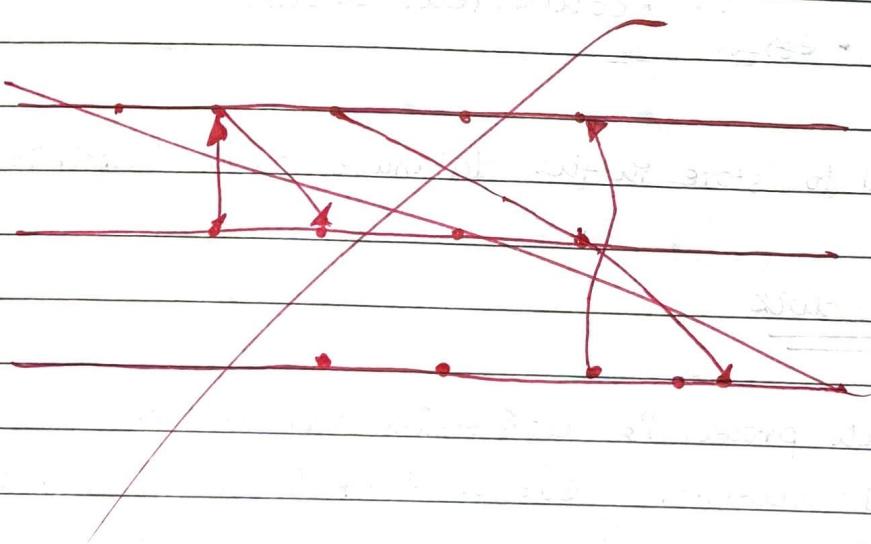
$(0, 0, 0)$   $P_3$

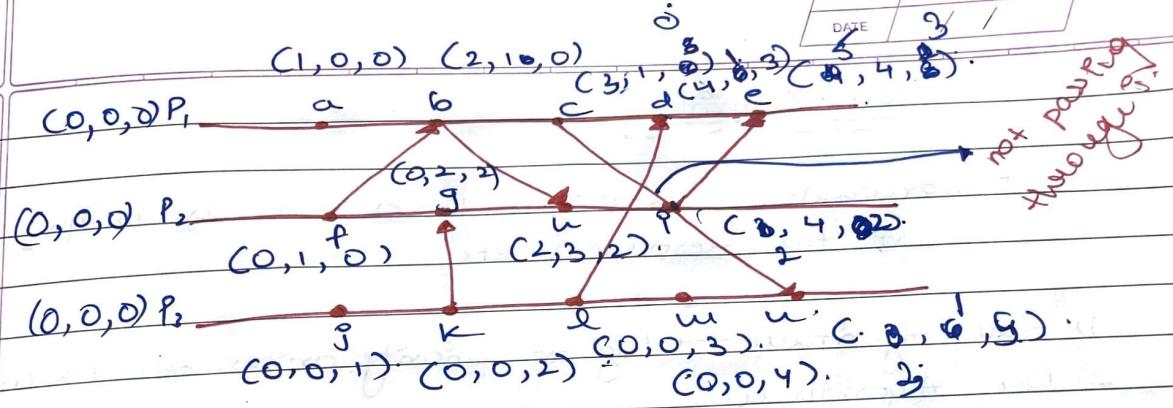
e

$(0, 0, 1)$

f  $(2, 2, 2)$

Concurrent events  $\rightarrow a, e \rightarrow$  If there is no  
b, e communication.  
c, e  
d, e.

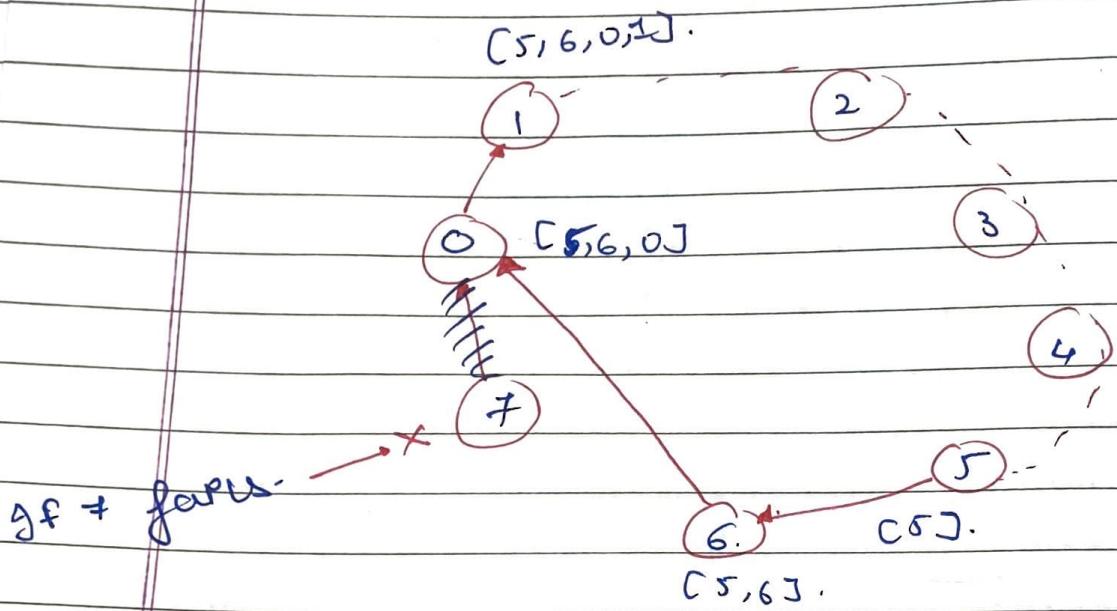




## Election Algorithm.

→ Coordinator → assigns tasks such as assign or request the resources etc.

- 1) Bully Algorithm → highest priority process becomes coord.
- 2) Ring Algorithm.



Q: which Algorithm is better? why??

\* It sends (IDt message).

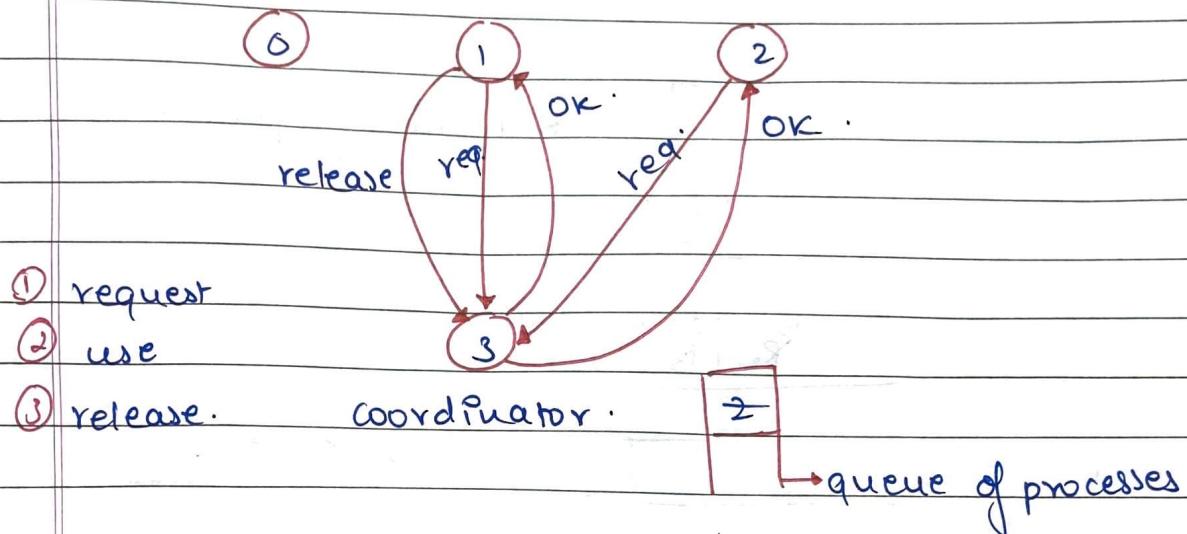
\* When it again receives (5) then it looks for higher priority in the array and chooses that as coordinator.

\* How will we know (7) has failed → ping count.

## Mutual Exclusion Algorithm

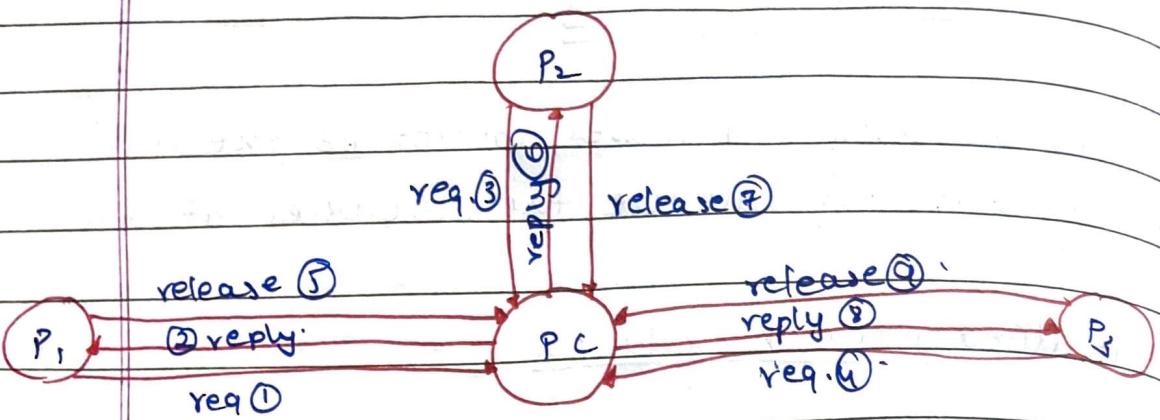
- Simultaneously 2 processes cannot be executed.
- Only 1 process can use the resources at a time.

### ① Centralised Mutual Exclusion



### Disadvantage

- 1) Centralised Coordinator Crashed
- 2) Starvation.



[ ] initial status of queue .

[ ]  $P_2$

[ ]  $P_2$  |  $P_3$

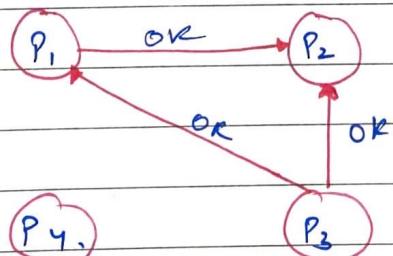
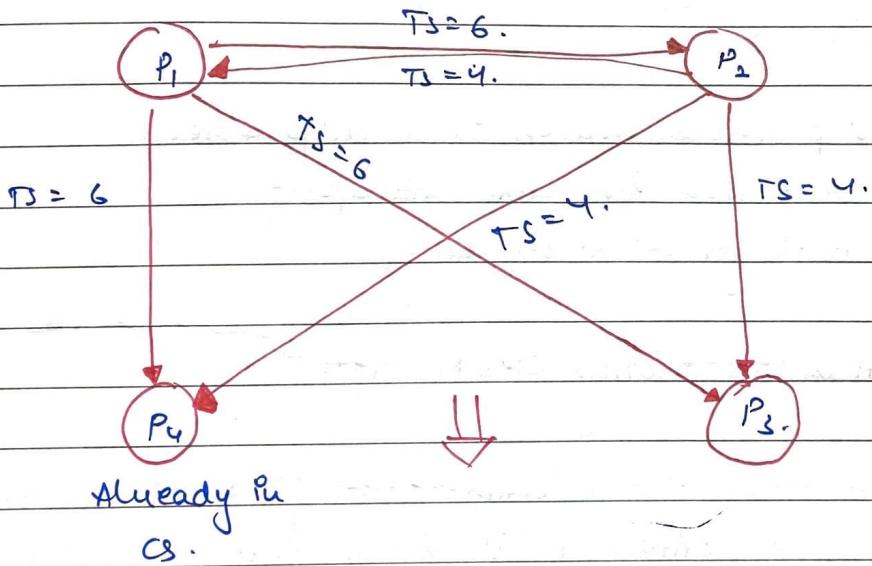
[ ]

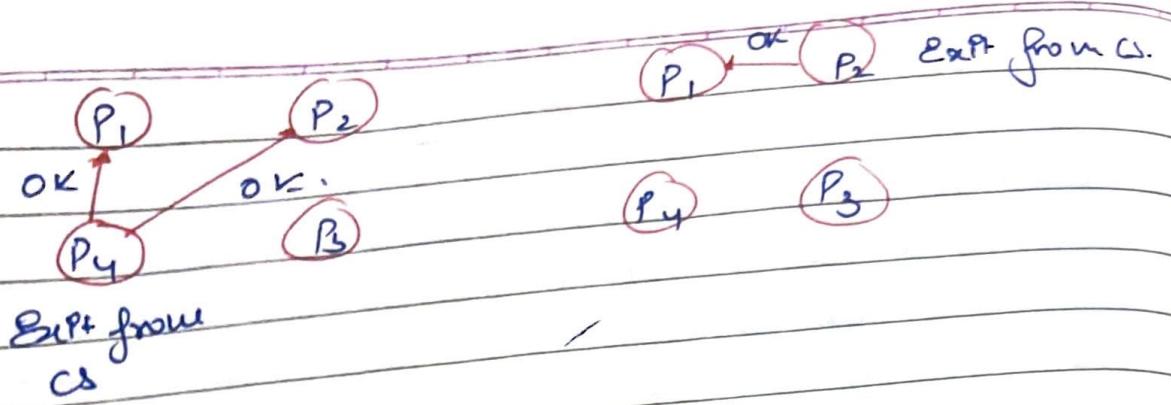
For the arrival of new event, timestamp has to increase by 1.

PAGE NO.	
DATE	/ /

## (2) Distributed Mutual Exclusion

- Based on event ordering & time stamps.
- Process K enters critical section as follows:
  - Generate new time stamp  $T_{SK} = T_{SK} + 1$
  - Send request  $(K, T_{SK})$  all other n-1 processes
  - Enter critical section
  - Upon receiving a request message, process P;
    - sends reply if no contention.
    - If already in critical section, does not send reply, queue required.
    - If wants to enter, compare  $T_{SP}$  with  $T_{SK}$  & sends reply if  $T_{SK} < T_{SP}$ , else queue.





### Difference

- No coordinator in distributed.
- Single point of failure in 'n' point of failure in centralized and distributed respectively.

\* Timeout mechanism used to solve single 'n' point of failure.

### Disadvantage

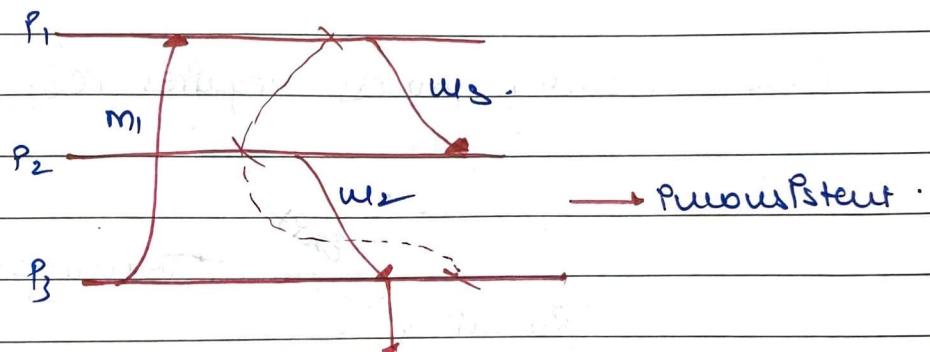
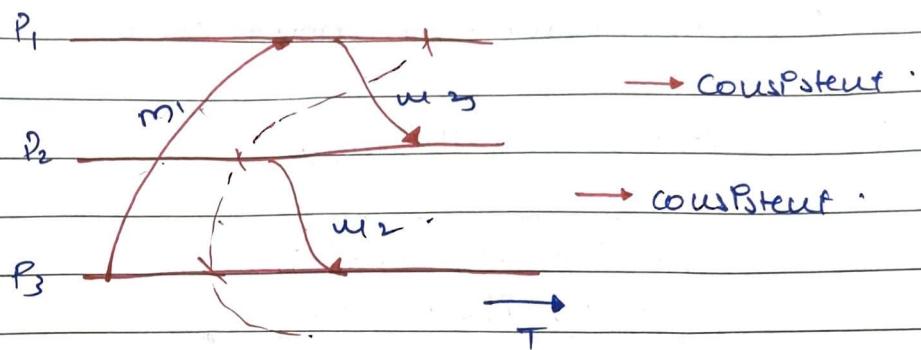
- All processes know id of all processes.
- Permission for time stamp  
(waits for other processes)

\* Solution → Voting (Majority)

- Q. Either draw the diagram through the time stamp given or write the algorithm through the diagram given.

## Global State

- Global state of D.S.
- Local state of each process
- Messages sent but not received.
- Many appln need to know the state of the system.
- Failure recovery, distributed deadlock detection.

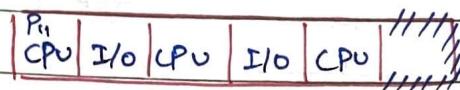


Because it cannot be received before checkpoint.  
It is received without being sent.

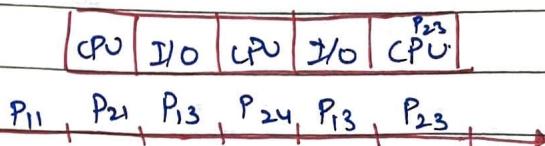
## Processes

- 1) Multiprogramming → Single Processor
- 2) Multiprocessing
  - CPU bound
  - I/O bound
- Q. why all processes are not either CPU or I/O bound  
 → Multiprogramming won't occur.

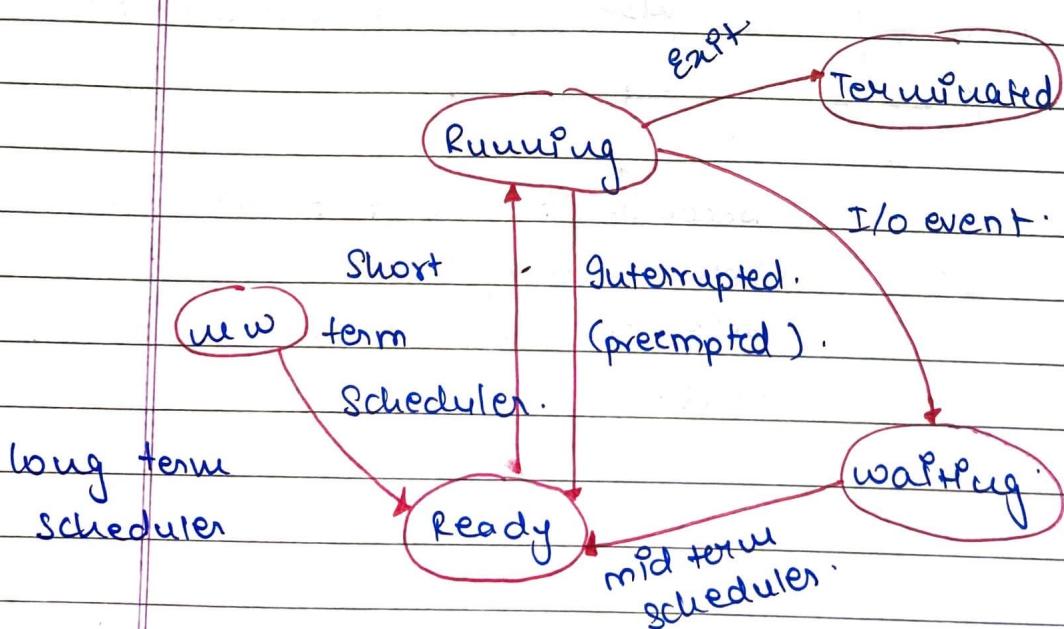
prog 1



Prog 2.



- \* Each and every process requires PCB.



- \* System call → interrupt generated by the system.
- \* Ready queue and long term scheduler is required because we need to balance I/O and CPU bound.

### Objectives of scheduling

- 1) Maximum CPU utilization
- 2) Throughput becomes maximum.
- 3) Minimum waiting time
- 4) Response time.
- 5) Minimum turn around time.

### CPU scheduling Algorithms

#### Non-preemptive

- 1) FCFS
- 2) SJF
- 3) Priority

#### Preemptive

- 1) SJF (SRPT)
- 2) Priority
- 3) RR.
- 4) EDF

Q. How to know if it is pre-emptive or not?

→ Starting time of all processes is same → non-preemptive  
Hence, why RR is non-preemptive.