

LECTURE NOTES

ON

**DISTRIBUTED SYSTEMS
ACADEMIC YEAR 2021-22**

IV B.Tech.–II SEMESTER (R16)

G.K.HAVILAH, Assistant Professor



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

**V S M COLLEGE OF ENGINEERING
RAMCHANDRAPURAM
E.G DISTRICT
533255**

DISTRIBUTED SYSTEMS

OBJECTIVES:

- Provides an introduction to the fundamentals of distributed computer systems, assuming the availability of facilities for data transmission, IPC mechanisms in distributed systems, Remote procedure calls.
- Expose students to current technology used to build architectures to enhance distributed Computing infrastructures with various computing principles

UNIT-I:

Characterization of Distributed Systems: Introduction, Examples of Distributed Systems, Resource Sharing and the Web, Challenges.

System Models: Introduction, Architectural Models- Software Layers, System Architecture, Variations, Interface and Objects, Design Requirements for Distributed Architectures, Fundamental Models- Interaction Model, Failure Model, Security Model.

UNIT-II:

Interprocess Communication: Introduction, The API for the Internet Protocols- The Characteristics of Interprocess communication, Sockets, UDP Datagram Communication, TCP Stream Communication; External Data Representation and Marshalling; Client Server Communication; Group Communication- IP Multicast- an implementation of group communication, Reliability and Ordering of Multicast.

UNIT-III:

Distributed Objects and Remote Invocation: Introduction, Communication between Distributed Objects- Object Model, Distributed Object Model, Design Issues for RMI, Implementation of RMI, Distributed Garbage Collection; Remote Procedure Call, Events and Notifications, Case Study: JAVA RMI

UNIT-IV:

Operating System Support: Introduction, The Operating System Layer, Protection, Processes and Threads –Address Space, Creation of a New Process, Threads.

UNIT-V:

Distributed File Systems: Introduction, File Service Architecture; Peer-to-Peer Systems: Introduction, Napster and its Legacy, Peer-to-Peer Middleware, Routing Overlays.

Coordination and Agreement: Introduction, Distributed Mutual Exclusion, Elections, Multicast Communication.

UNIT-VI:

Transactions & Replications: Introduction, System Model and Group Communication, Concurrency Control in Distributed Transactions, Distributed Dead Locks, Transaction Recovery; Replication-Introduction, Passive (Primary) Replication, Active Replication.

OUTCOMES:

- Develop a familiarity with distributed file systems.
- Describe important characteristics of distributed systems and the salient architectural features of such systems.
- Describe the features and applications of important standard protocols which are used in distributed systems.
- Gaining practical experience of inter-process communication in a distributed environment

TEXT BOOKS:

1. Ajay D Kshemkalyani, Mukesh Sigal, "Distributed Computing, Principles, Algorithms and Systems", Cambridge
2. George Coulouris, Jean Dollimore, Tim Kindberg, "Distributed Systems- Concepts and Design", Fourth Edition, Pearson Publication

REFERENCE BOOKS

1. Distributed-Systems-Principles-Paradigms-Tanenbaum PHI

CHARACTERIZATION OF DISTRIBUTED SYSTEMS

* Introduction of Distributed System:

Distributed System: A distributed system is one in which the hardware and software components that are located at network computer are able to communicate and also coordinate their actions by passing messages.

Significant Consequences:

(i) Concurrency: The connection of network of computers is the concurrent execution. more than one user can work on their system at a time by sharing resources such as webpage or files. By adding more resources to the network, its capacity of handling shared resources can be increased.

(ii) No Global clock: The programs that coordinate their action by message exchanging when they need to cooperate. And this coordination may depend upon the idea of time at which the actions of program occur. But the accuracy with which the clocks of the computers in a network can be synchronized may have limits. i.e there is no global notion of the correct time. This is a consequence about the

22

-that Sending messages through the network is the communication.

(iii) Independent failures - Any Computer System may fail, and hence it is the system designers responsibility to plan for such consequences of possible failures. The computers that are connected to network can be separated as a result of failures but they might not stop running. The programs running on them may not know that about the failure or slowness of the network. Similarly, the failure or program termination of computer may not be known to the other computer with which it communicates. Hence, the other systems can still running even if one fails.

* Examples of Distributed Systems

(i) Internet :-

In simple terms, internet is a network of network that consists of millions of private, public, academic, business and government networks connected to each other with the help of widely available internet working devices such as routers, gateways, bridges etc.

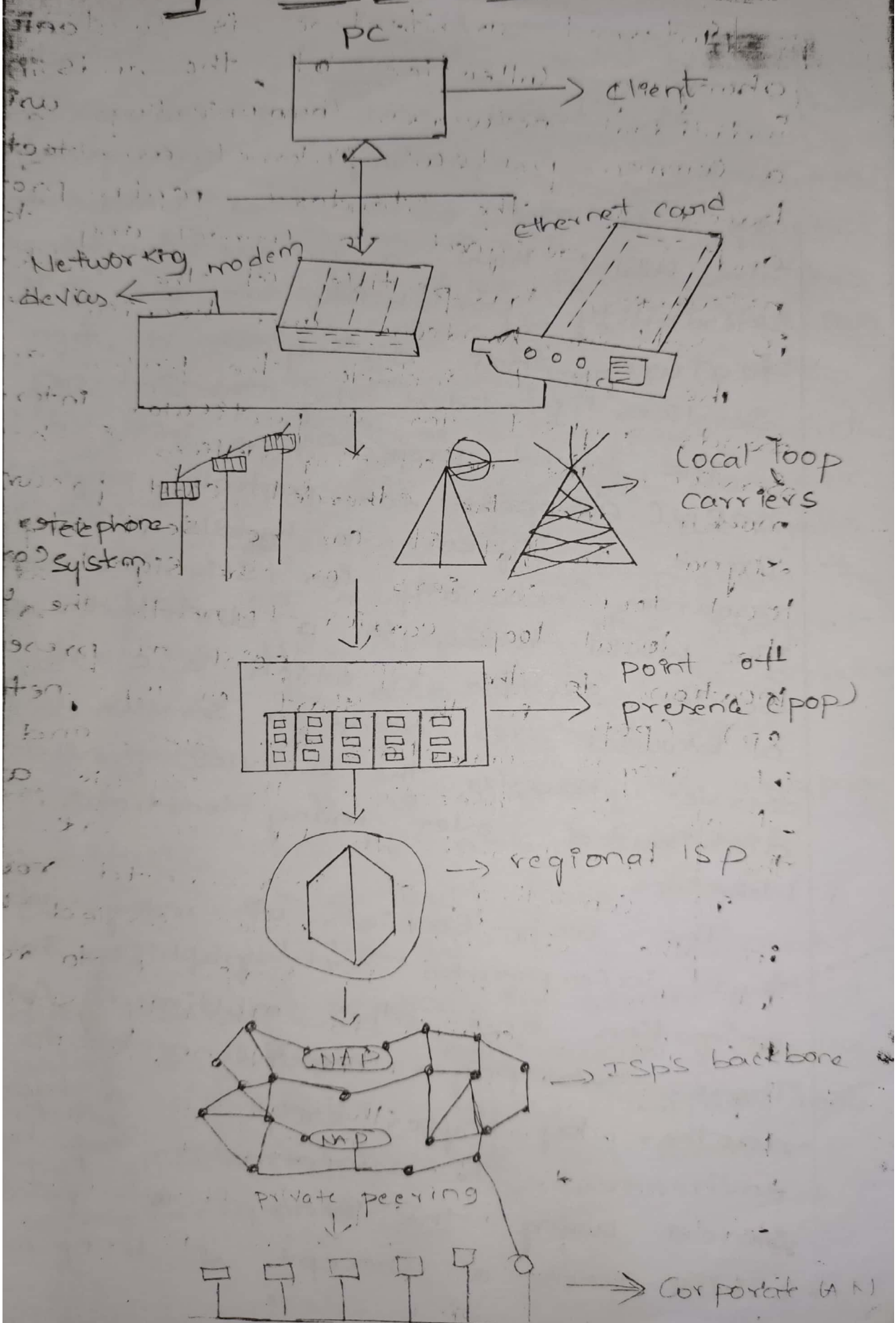
Architecture of Internet

Internet architecture is constantly changing collection of thousands of individual networks communicating using a common protocol. Internet architecture based on the standard TCP/IP protocol and is designed to connect any two networks, irrespective of the difference in software, hardware and technical design.

The client initiates the process and sends request for a particular internet service. The networking devices such as the modem and the ethernet card process the signal so that it can be sent over the local loop carriers or the signal carriers. The local loop carrier connects the user location to the ISP's point of presence (POP), (POP is the start of ISP's network. It accepts connections from users and authenticates the connections. The signals are then transferred to ISP's network.

It consists of interconnected routers in the different cities served by ISP. If the packet is supposed to reach the host, which is directly served by ISP, then it is delivered to the host. Otherwise, the packet is forwarded to reach the host which is directly served by ISP, ISP backbone operator. It interconnects the ISP's, POP's and also to other ISP's.

fig :- Internet architecture



If the packet received by the backbone is supposed to reach on ISP or company served by backbone then the packet is forwarded to the nearest router.

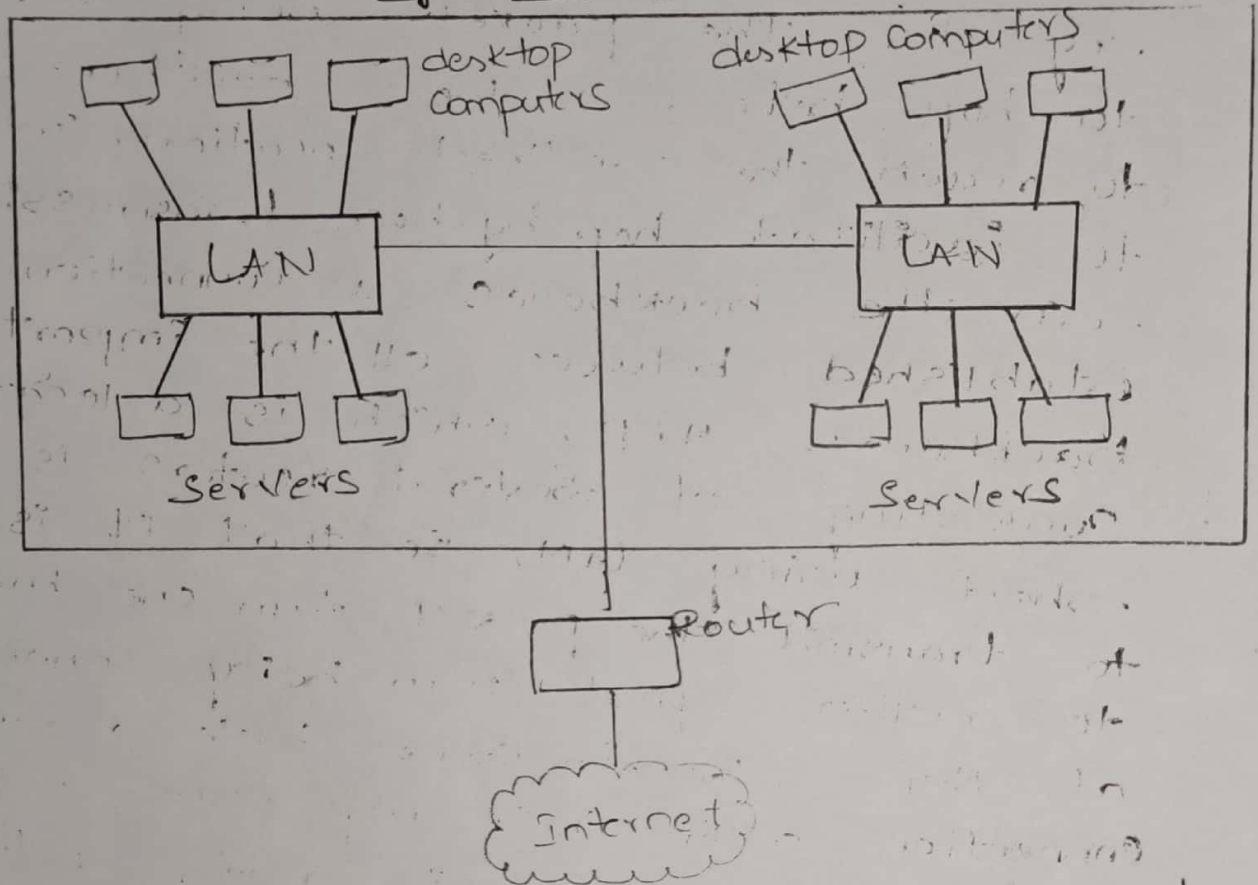
Since, there are many backbones present in the network, it is necessary for the packets being transmitted to hop between these backbones so as to reach the destined location. In order to facilitate hop-by-hop transmission between the backbones, a connection is established between all the important backbone's NAP, which is a location consisting of several routers is established using LAN, so that it is possible to transmit the packet from one backbone to another. A part from being connected at NAP's, the backbones have a direct connection with their respective router(s). Such a connection is referred to as "private peering".

(ii) Intranets:-

A portion of an internet that can be administered separately and which has boundary that is configured to impose security policies is known as intranet. Several 'LANs' are linked in its backbone connections. The organization which

administers the intranet is responsible for its network configuration and also vary from LAN on a single site to a set of connected LAN's of several branches of companies across different countries.

fig :- Intranet



An Intranet is connected to the internet through a router that allows the users of the intranet to use the services such as web or email. The services that it provides can also be accessed by the users of other intranets.

* Resource Sharing and the web:-

World Wide Web (WWW):-

The "world wide web (www) or the web, is a repository of information, spread all over the world and linked together. The WWW has a unique combination of flexibility, portability, and user friendly features that distinguish it from other services provided by the Internet.

The WWW project was initiated by CERN (European for particle physics) to create a system to handle distributed resources necessary for scientific research.

The WWW today is a distributed client-server service in which a client using a browser can access a service using a server. However the service provided is distributed over many locations called websites.

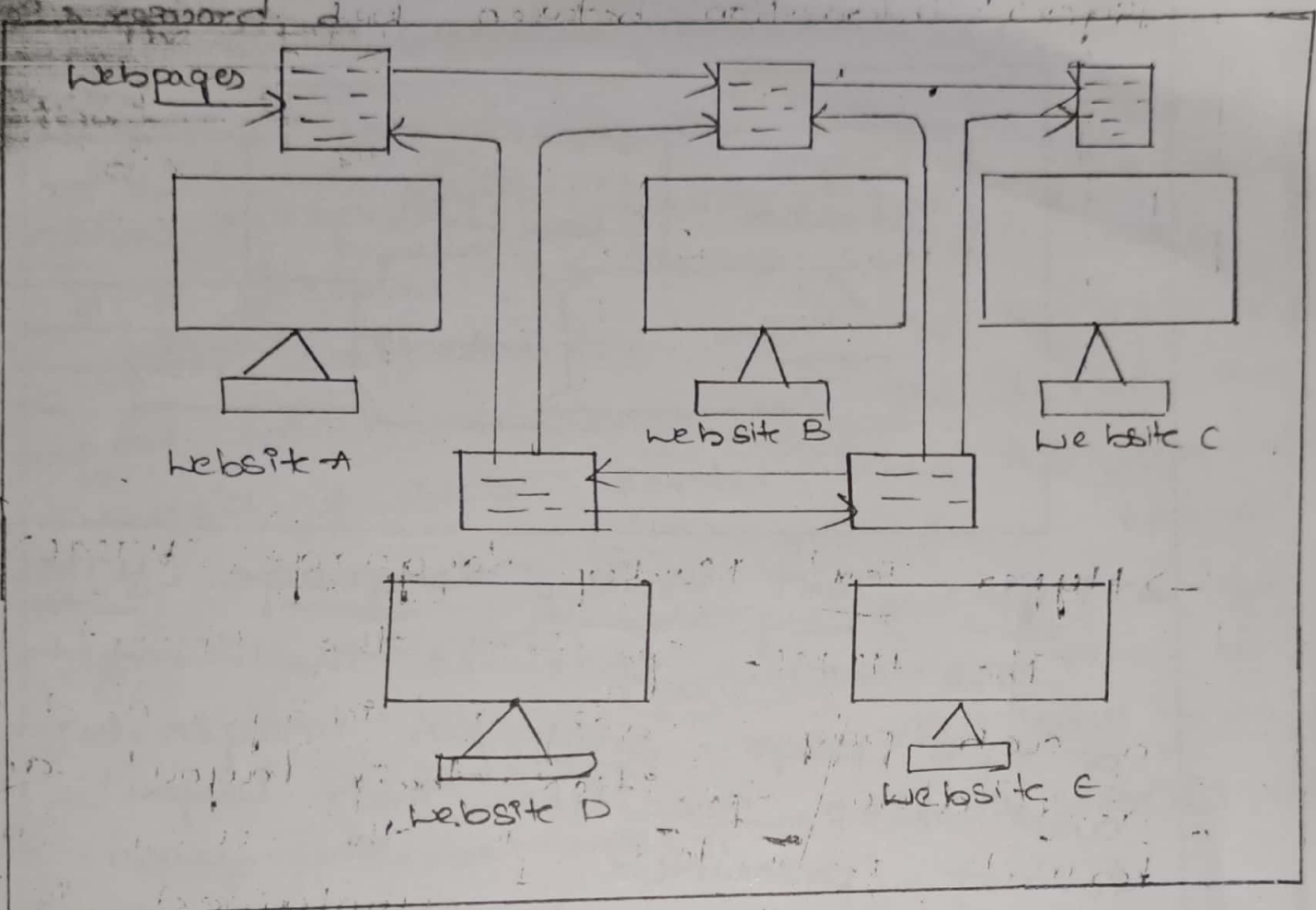
The web consists of many webpages that incorporate text, graphics, sound, animation and other multimedia components. These web pages are connected to one another by hyper text. In a hyper text environment the information is stored using the concept of pointers. WWW uses a concept of HTTP which

allows to communicate between web browser and web server. The web pages can be created by using a HTML, this language has some commands which are used to inform the browser about the way of displaying the text, graphics and multimedia files. HTML also has some commands through which we can give links to the web pages.

If we want to get a page from the web, we have to type URL for our desired pages, or otherwise we have to click on a link that provides the URL. The URL specifies the internet address of the web server, the directory and name of our desired page.

If there is no directory or web page specified, then the web server will provide a default home page.

The WWW today is a distributed client-server service in which a client using a browser can access a service using a server. Fig (1) illustrates how the different web site can communicate with each other.

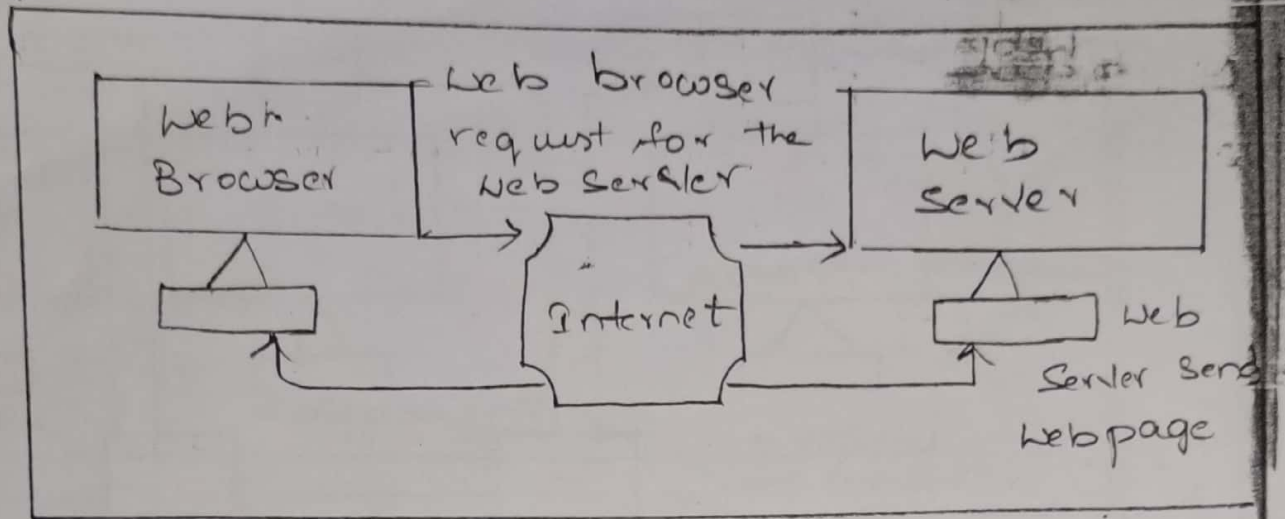


Working of Web:-

The web operates on a client-server model. A web browser acts as the client. In the WWW interaction, using this program, a user sends a request for a web page stored on a web server. The web server locates this web page and sends it back to the client computer. The web browser then interprets the web page written in the HTML language and then displays it on the client computer's screen.

The

Figure 1: Interaction between Web browser & Server



→ Hyper Text Markup Language (HTML) :-

The HTML specifies the contents of a web page such as images and text and also specifies their layout and format presented to the user. Such structured items are contained in the web pages as headings, paragraphs, tables and images. The links and its associated resources are also specified by the HTML.

The HTML can be produced by the users either by using standard text editor or by using an HTML-aware 'wysiwyg' editor that generates the HTML from the layout created by them. A typical HTML text would be as created by them. A typical HTML text would be as

```
<IMG SRC="http:
```

The HTML text is stored in a file that is accessed by the web server and later returned read and rendered by the browser into the formatted text and images laid on a web page. However, only the browser can interpret the HTML text. The server only informs the browser about the type of content it is returning, by inferring from file name extension in order to distinguish from other text files.

All the HTML directives known as tags are enclosed by angle brackets. The text between these tags is as it is presented on the web page of the link. The browser are also configured to show the text of the underlined links by default in paragraphs to the user.

If interested, also have a look at company details, the URL contained in a CA.

The association between the links displayed text and URL contained in the `<A HREF ...>` tag, when clicked, on the underlined text. The resource which is identified by the URL is identified and presented to user.

for example, a web page about the company details specified by the resource HTML file.

→ Uniform Resource Locator (URL) :-

The URL is the address of the system on the internet. It identifies the resources and hence can be examined by the browsers to access the corresponding resources. When the user clicks on a link or selects one of its 'book marks', the browser looks up the respective URL or fetches a resource embedded in a web page such as an image or text. The URL has two top-level components.

Scheme: Scheme-specific-identifier.

The first component 'Scheme' declares the type of the URL such as user's email, file to be retrieved using ftp, http and mid.

The URLs can be used of the form widget and the browsers must be given the capability of using the new 'widget' protocol by adding plug-in when a widget is invented with its own address. -ng Scheme and protocol to locate and

account them.

The resources can also be accessed by using HTTP URIs by using standard HTTP protocol. The HTTP URIs identify which web server maintains the resources and which resources at the server are needed. A web server files are maintained in a one or more sub trees of its file system. Every resource can be identified by a path name related to the server. The HTTP URI will be as, $http://\text{Servername}[:\text{port}][:\text{Path Name}][?query][\#fragment]$

The items which are enclosed in bracket are optional. A full HTTP URI begins with 'http://' and the server name follows it. The server's name is expressed as Domain name system and the "port" number that follows it is optional through which the request are retained by the server's and is 80 by default. Then followed by path name of the server's resource is optional and in its absence the server's default web page is needed.

request-reply protocol, the URL of the resource required is sent as a message by the client to the server. The server sends back the file contents as a reply to the client if the path name exists or else it sends an error response such as "404 Not Found".

(ii) Content Type

Every content type may not be handled by the browser. So, it adds the content types it prefers in the reply. For example, the browser may be applicable of displaying images of type 'png' format. The servers consider it and adds the content type in the reply to make the browser understand how to process it. The content types are indicated by the strings called MIME types. The set of actions taken by the browser for a given type are configurable.

(iii) One Resource per Request

One resource per request is specified by the client. Consider for example a web page contains nine images, then totally ten requests are issued.

by the browser in order to access the whole content of the page. Several requests are made concurrently by the browser so that the overall delay to the user can be reduced.

(iv) Simple Access Control

Any user who has network connectivity to a web server can access any of its resources by default. In case, if the user needs to restrict access to any resource, then the server can be configured to issue a challenge to the client who requests it. The client can access the resource for example by typing a password.

* Distributed System Challenges

The construction of distributed systems produce the challenges that are given below,

(a) Heterogeneity: The distributed system must be constructed from different networks, operating systems, computer hardware and programming languages. The difference in networks is marked by the internet communication protocol and the

Other differences are dealt by middle

-ware

(b) Openness:- They must be extensible by publishing the interfaces of the components first but the real challenges would be the integration of component & that is written by different programmers.

(c) Security:- Adequate protection of shared resources can be provided by using the encryption and also the sensitive information can be maintained confidential when transmitted over the network. The problem is still caused by denial of service attacks.

(d) Scalability:- If the cost of adding a user is constant in terms of resources that need to be added, then the distributed system is said to be scalable. The performance bottleneck should be avoided by the algorithms that are used to access the shared data and this data must be structured hierarchically to obtain best access times. The data that is frequently accessed can be replicated.

(e) Failure handling:- Any computer or network may fail independently without affecting others. Hence each of the component must be aware of possible

ways the components on which it depends fail and also must be designed in a way to be able to deal those failures.

(f) Concurrency: multiple users in distributed system send concurrent request to its resources. Hence, each of it must be designed to be safe in such environment.

(g) Transparency: The aim of transparency is to hide some aspects of the distributed system from the application program to make it available to the design of their particular application.

→ Heterogeneity:

Internet allows users to access services over a variety of different networks. It also allows to run applications over heterogeneous collection of computers and networks.

Heterogeneity is applied to different networks, programming languages, OS, computer hardware, implementation by different developers.

Integers are presented in different

forms on different hardware for example
there are two ways for byte ordering
of integers. When messages are exchanged
on different hardware these differences
of representation must be addressed
for example calls for exchanging
messages in windows is different
from calls in unix.

For characters and data structures
different programming languages use
different representation. Hence in order
to communicate with each other, these
differences must be dealt
for working in different environment
of network communication programming
language common standards has to be
adapted by developers so that they
can easily understand and can easily
communicate with each other using pro-
grams they have written.

(a) Middleware:-

Middleware refers to a software layer
whose purpose is to make heterogeneity
and provide programming abstraction. For
example Common object request broker
Architecture (CORBA). difference of under-
networks, operating system, computer h/w
is dealt by middleware by implement

internet protocols which masks these differences.

Java & remote method invocation is a middleware which support single program

-ming language

(b) mobile code and heterogeneity:-

mobile code is a term which applies to code that is sent from source computer to destination and runs at destination.

for example java applets, executable programs are specific to instruction set and host operating system. hence code which is suitable for running on one computer is not suitable for other.

example when an executable file is sent via email attachment by windows, macintosh users will not run on macos or x86 users will not run on x86 computer running linux operating system.

In virtual machine approach compiler generates code for virtual machine instead of particular hardware. for example java compiler generates code for virtual machine which is implemented once in each type of hardware to allow java programs to run. however this solution is not possible with other programming languages.

→ openness:-

Openness can be defined as a feature which determines whether the system can be re-implemented and extended in different ways. Openness of distributed system can be known by the level at which new resource sharing can be included and made available for use by client programs. Openness cannot be obtained until software developers are provided with documentation and specification of the key interface of the ports of system. In other words, until the key is published. This process is similar to standardizations of key interfaces but it usually ignore those standardizations because they are slow and complicated.

The internet protocol designers introduced series of specification and documents, series of documents are called as "Request for Comments" each of them is determined by a number specifications for internet protocols were published in the year 1980's which followed specifications for applications that run over them like email, telnet & file transfer.

The process of documentation and specification has been continued and copies of specification as well as discussion be obtained from www.ietf.org, for example

CORBA publishes documents and "Specifications" of the interfaces of its "Services" which can be obtained from www.omg.org.

System that are designed to serve resource sharing in this way are called open distributed system which are extendable at hardware level they are extended by including more number of computers attached to network and at the network and at the software level they can be extended by re-implementation of old ones and introduction of new services and by allowing application to share resources.

→ Security

Information resources that are of high value to users is maintained in distributed systems and their security is of high importance.

Security for information resources has three parts :-

1. Confidentiality
2. Availability
3. Integrity.

Security risks increases when one program in a computer interacts with a program on another computer because it allows sharing of resources freely in an intranet. Firewall can be used to stop unauthorized

users from entering and gaining access to important data in internet.

But firewalls do not ensure about the proper use of available resource in an intranet or internet.

Clients sends request to access data which is maintained by server in distributed systems and in reply server sends information via message over a network.

For example, a bank manager might request to access list of new customers from server.

In this the difficult task is to send the details of customers like account number

address in a message over a network in a secure manner. Hiding the contents of

message is not just important even the identity of the bank manager should be

checked and assured.

Some of the security challenges which has not been met fully are,

(i) Denial of Service

(ii) Security of mobile code.

(i) Denial of Service:- Some users try to disrupt service of target sites or services hosted on high profile web servers. Such as banks, offices by bombarding large number of pointless requests so that

Service becomes unavailable to the serious users of the site, that internet site stops functioning efficiently. This is called denial of service.

(ii) Security of mobile code: - mobile code has to be tackled very carefully. If a user receives an email attachment after loading that email locally on the computer and when it is being executed effects may be unpredictable because from above it may look a beautiful picture but in real without realizing we are being a part of denial of service attack.

→ Scalability:

Distributed systems work effectively and efficiently at different scales ranging from intranet to internet.

A distributed system can be defined as Scalable only if it works effectively and efficiently with an increase in the number of resources and users.

During the last 20 years there is a dramatic increase in the number of computers and services. From 1993-2003 there is significant increases in the no. of web servers and computers.

Some of the challenges encountered while designing scalable distributed

Systems are,

- 1. Controlling physical resources, Cost
- 2. Controlling the loss in performance
- 3. Preventing running out of SW resources
- 4. Avoiding performance delays

1. Controlling physical resources Cost:- The Cost of physical resources should be reasonable as the demand increases and should be possible to extend the system.

for example as number of user and computer increases the frequency at which files are being accessed also increases in an intranet

Server Computer should be possibly added to avoid performance delays, when a single server has to tackle n number of file system with N number of users to be scalable there should be at most $O(n)$ physical resources to support them.

for ex, if there are 20 users and a single server supports them, such 2 servers then should be capable of supporting 40 users.

(2) Controlling the loss in performance:-
 Let assume a data set whose size is equal to number of users in the system. for example, the table with domain names of computer and then IP addresses held

By domain name system, which is used for looking up domain name like www.slag-ups.com hierarchic structures scales better in algorithms than linear structures. But increased size in hierarchic structures will cause loss in performance. maximum performance loss should not be less than $O(\log)n$.

(3) Preventing Running out of Slw Resources

In 1970's 32 bits was used for assigning the numbers as Internet (IP) addresses but unavailability of numbers to assign as IP addresses have become the problem. To overcome this problem 128-bit addresses is used, there are solution for this problem correctly as said by early designers lot of changes will be required on slw components in order to make them work in 128 bits.

(4) Avoiding performance Delays

To avoid performance delay, algorithms should be decentralized. for example, in domain name system in which the name table was kept in a single master file and can be downloaded whenever needed this was not a problems until there were only few hundred computers that

needed it. (There was a delay in performance and administrative delays)

But problems arose when there was a dramatic increase in the number of computers which needed it. There was a delay in performance and administrative delays. This problem was solved by separating the name table between servers and administered locally.

The resources which are frequently used by users causes loss in performance. To overcome this problem catching and replication is used which improves performance of resources that are heavily used by users.

There need for modifications in system and software applications should be there even when the demand increases.

→ Techniques for dealing with failures!

When a failure occurs in hardware software proper work of systems fail which leads to wrong results and systems may stop working before the completion of computation.

There are different types of failure which occurs over the network and in the processes in distributed systems. In distributed systems, failures are partial.

which means some parts continues to work while the others fail hence it is difficult to tackle failures.

(1) Failure detection: Some failures can easily be detected. for example By using checksums corrupted data in a file or message can be detected. Some of the failures are impossible to detect like remote crashed server.

(2) Hidden Failures: Some detected failures can be hidden, for example (i) when a messages get fail to reach the destination it can be retransmitted.

(ii) A copy of file data is created on disks, if one gets corrupted other will be correct.

In worst cases, hiding failures techniques are not guaranteed, which means copy of data on the other file may also get corrupted. The message may not reach the desired destination. how often they may be transmitted.

(3) Tolerating failures: most of the internet services show failures and it will not be practical to detect and hide the failures in such a large network with so many parts.

for example, when a web browser tries to connect to web server and it doesn't get connected, the user has to tolerate the delay in service.

(4) Recovery from failures:- In recovery, the data can be rolled back or recovered after a server has been crashed. In some programs, a fault occurs when computations performed are incomplete and not consistent.

(5) Redundancy:- Some services can be tolerated by using redundant components. Every name table in domain name system should be copied in at least two servers. There should be at least two routers between two routers.

To ensure that data is accessible even after a failure, db should be copied in different servers, when a fault is detected in a server, clients are redirected to other servers without loss in performance. The data has to be updated by creating copies of the data.

⇒ Constructing Distributed System:-
In distributed systems, clients can share the resources provided by services, and

Applications. Hence, many clients will try to access the shared resources at the same time. For example, when applying for online scholarship and when deadline comes close, many clients will access frequently at the same time.

A program allows a shared resource to take one client request at a time but this limits the throughput hence services and applications allow many client requests concurrently. In order to make it more robust each resource is encapsulated as an object and are executed concurrently. In this case, many threads can be executed concurrently within an object.

In other cases, executing several threads concurrently on object may conflict and as a result it may produce inconsistent results.

For example, in a bank when two concurrent transaction is taking place with account number 502184, 502185 502184 is transferring 50,000 to some other account and 502185 has received 1 lakh rupees. If the correspon

...ding operations, are interleaved without
any control then they will become inconsistent

In distributed system, a shared resource
must ensure operations correctly in a
concurrent. Therefore, a programmer should
implement operations on objects and
server in a safe way and should
not be inconsistent.

By using standard techniques like
semaphores, operations on objects in a
concurrent environment can be safe and
consistent data can be achieved.

→ Transparency: It refers to component se-
-ration details of distributed system.
The user and the application programmer
This is done in order to deploy the
system as a whole rather than deploying
it into individual components.

There different eight forms of transp-
-arency are,

(1) performance Transparency: This type of
transparency allows system reconfiguration
-on so as to increase the performance
with change in loads.

(2) Scalability Transparency! This type of transparency allows extension of the system and application without effecting the system structure or the application algorithms.

(3) Failure Transparency! This type of transparency allows the hides the faults and failures of the system hardware and software components allowing the users and application programs to finish their tasks without any hindrance.

SYSTEM MODELS

Architectural Models:-

The structure is the architectural of a system. It is a collection of individually defined components. The complete objective is to make sure that the structure will satisfy the current and future demands on it. The vital responsibilities are to make the system cost effective, manageable, reliable and adaptable. The architectural design of building has some features like determining architectural style as well as its general structure. The architectural styles gives a constant frame of reference for the design.

The central architectural models are constructed almost on the notion of process and object the functioning of separate components of a distributed system is first simplified and then abstracted in an architectural model. This model later evaluates the following things.

- (a) The arrangement of components throughout a network of computers i.e. trying to specify beneficial patterns for the distribution of data and workload.
- (b) The interrelationships between the components i.e. the components functional roles and the pattern of communication.

The processes are divided into server processor, peer processes and client processes in initial simplification. In the latter simplification, the process that invite, unite and communicate in a proportional way are involved to perform a task. Through this division, processes recognize their concerns which is helpful in evaluating their workloads. It also makes easy to decide the effect of failures, in each process. To achieve the performance and reliability objectives for the resulting system, the arrangement of processes is important. This arrangement can be identified through the above analysis.

By altering client server model, one can construct few systems.

(a) A process can assign a task to another since it is possible the code from one process to another.

Example: The code from servers is downloaded by the client process and executed locally. To minimize access delays and communication traffic, the code and objects that access, the process can be transferred.

of the distributed systems are designed to allow the computers and other mobile devices to be added or removed smoothly, enabling to find, obtain services and also to provide their services to others.

There are many patterns that are used for the allocation of work in a distributed system. These patterns have a major effect on the performance of the resulting system.

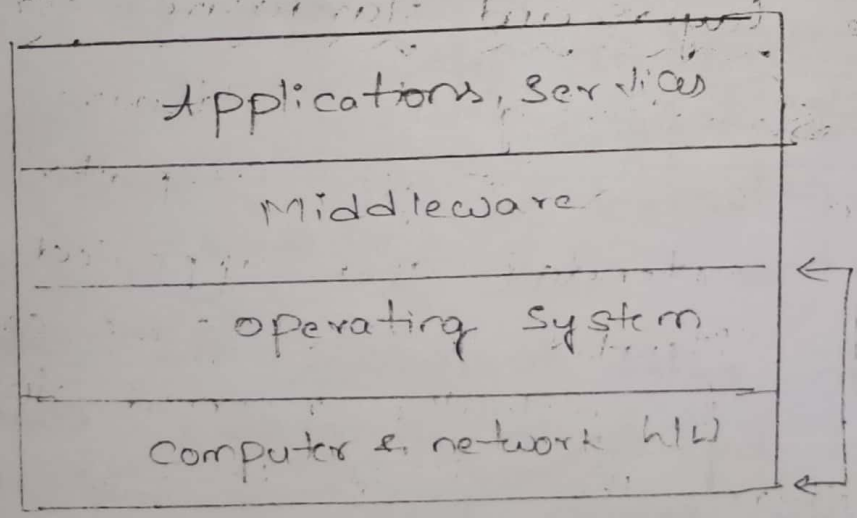
The particular issues of performance, reliability, cost and security influences arrangement of the processes in a network of computers. The above discussed architectural models give only a simple view of major patterns of distribution.

* Software layers and hardware service layer

The structuring of software as modules or layers in a single computer and the services extended and appealed between processes situated in the same or different computers is actually referred to as software architecture. The disclosed form of this process and service oriented view is a layer of services in distributed systems. This layer of services included both software and hardware service. The process which accepts request from other processes is a

a server. one or more servers can supply a distributed service to support a consistent system wide view of the service resources. these servers can communicate with one another and also with client processes.

AS an example, consider a network time service which is carried out by server programs on the internet located on the network time protocol (NTP). These server processes executing on hosts receive the request of client programs on the internet and provide current time to them as a result of interaction with each other. these processes alter their version of current time to existing reference



Software & hardware service layers

Platform:

The bottom layers (Operating System, Computer and network hardware layers) serve as platform for distributed applications. The services provided by these bottom layers is utilized by higher layers. In every computer, operating system layer and computer network layer are implemented separately. This implementation facilities communication and co-ordination between processes by transferring the system's programming interface to that level. Important examples of platform include Intel x86/Linux, Intel x86/Solaris, Intel x86/Windows, power pc/Macos.

2. Middle ware:

A piece of software placed between the application and operating system is called "middleware". Its intension is to enclose heterogeneity and to supply an accessible model to application programmers. The concept of middleware can be easily explained in the context of objects or processes in a set of computers that interact with each other to carryout communication and resource-sharing. Support for distributed application. Middleware supplies helpful building block in the construction of software components. These

Components can interact in distributed applications. In specific, it maintains abstract to increase the level of communication activities in a distributed system. These abstractions include the transfer of multimedia data in real time, the partitioning interactions between a group of processes, placement and extraction of shared data objects amongst operation systems, Remote method invocation (RMI), making multiple copies of shared data objects and notification of events.

The earliest instances of middleware include remote procedure calling packages (eg Sun RPC) and group communication systems (eg Isis). Object-oriented middleware products and standards are:

(a) CORBA

(b) Java RMI

(c) Microsoft's Distributed Component Object Model (DCOM)

(d) Web Services

(e) The ISO/ITU-T's Reference model for open distributed processing (RM-ODP)

The services provided by middleware products and standards, such as infrastructural services, can be used by applications programs. As an example, consider the services offered by CORBA. They are security, persistent storage, transactions, naming and event

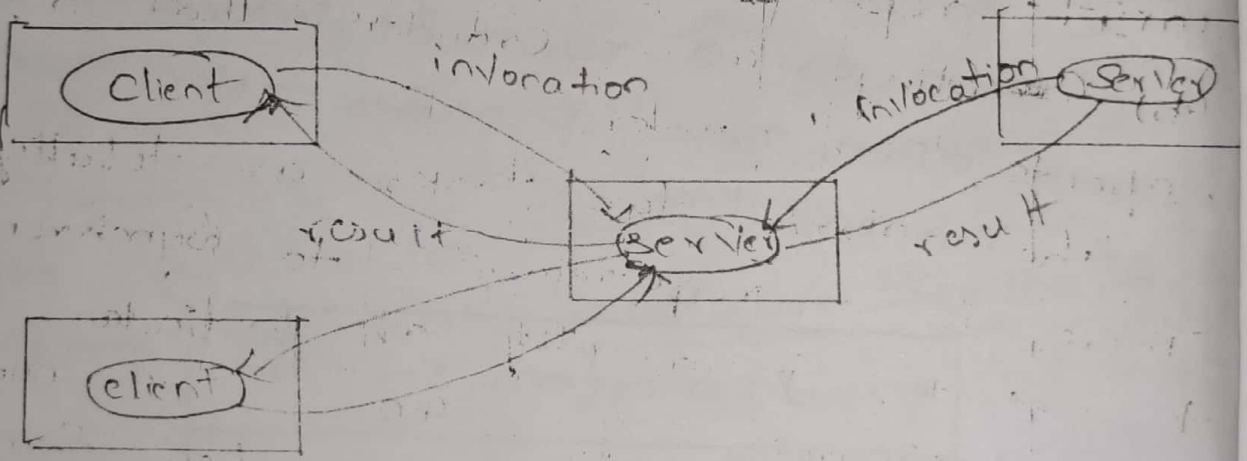
notification. The upper layer provides some specific services.

* System Architecture Models:-

(i) Client-Server Model:-

When dealing with distributed systems, the architecture very frequently referred to is Client-Server. This architecture is broadly utilized and significantly very important. The figure below shows the structure of client server model.

Fig. (i) Service provided by ^{single s.} multiple servers



In order to access the shared resources managed by server processes. Example, the local file server that controls the files consisting webpages has a common client called the files web servers and many other has many clients such as web servers and many other internet services search engines is other example needed to web

It allows users to search the sites all over the internet for synopsis of information available on web pages. These synopses are created by programs called web crawlers.

At search engine site, these web crawlers are executed in the background, using HTTP requests to access web servers all over the internet. A search engine thus plays the role of both a server and a client.

It performs the following two tasks:
(a) It replies to queries from browser clients.
(b) Executing web crawlers that act as clients web server.

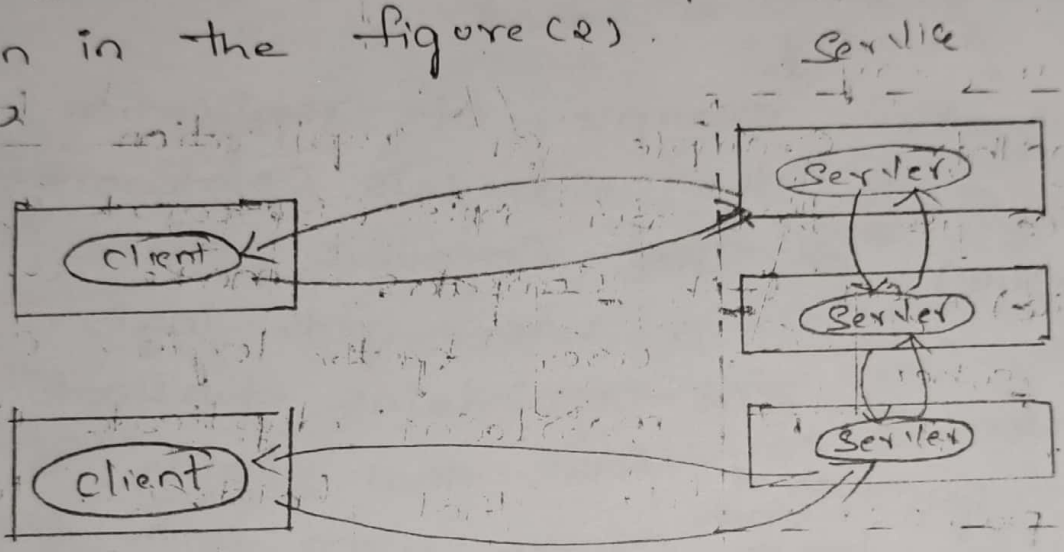
The above two tasks are totally separate. There is a requirement to synchronize them to execute concurrently. In particular, a normal search engine consists of numerous concurrent threads of execution. Some threads provide servers to client and others execute web crawlers.

(ii) Multiple Server model:-

In multiple server model, several servers implement the services as its processes in separate host computers, after implementation these computers interact with each other. If required, in order to provide a

Service to client processes. The servers may also provide a service to client processes. The servers may also provide the service either by partitioning the services into set of objects and distributing them among themselves or by maintaining replicate copies of these services on several hosts. The scenario of multiple server model is shown in the figure (2).

fig 2



Clients Invoke Multiple Individual servers

The concepts of this model are illustrated in the following examples.

1. On web each server maintains their own resources, called partitioned data. Therefore a user can access a resource of any one server.
2. Replication is an important concept of increasing availability and performance.

2) and for improving fault-tolerance, the processes that run on different computers can have multiple consistent copies of data due to replication. For an instance, data available at mittal, elec.com might also be provided by some other web services. It is due to mapping of service into several servers maintaining database replicated in memory.

3) Another example of replication based service is the Sun NIS (Network Information Service). On ATM, computers make use of NIS when a user gets login. Each server of NIS maintains duplicate copies of password file that holds a list of user's login names along with encrypted passwords.

→ Proxy Servers and Caches:-

A cache is a repository. The newly used data objects are stored in a cache. Whenever a new object is collected at a computer, it is integrated to the cache store. Some times, this integration requires restoring already available objects. When a client process requests an object, the caching service starts examining the cache for updated copy. If the updated copy is available then caching

Service provides the object from the repository. Otherwise, it retrieves an updated copy. In general, caches are situated in a proxy server that is shared by multiple clients or available with each client.

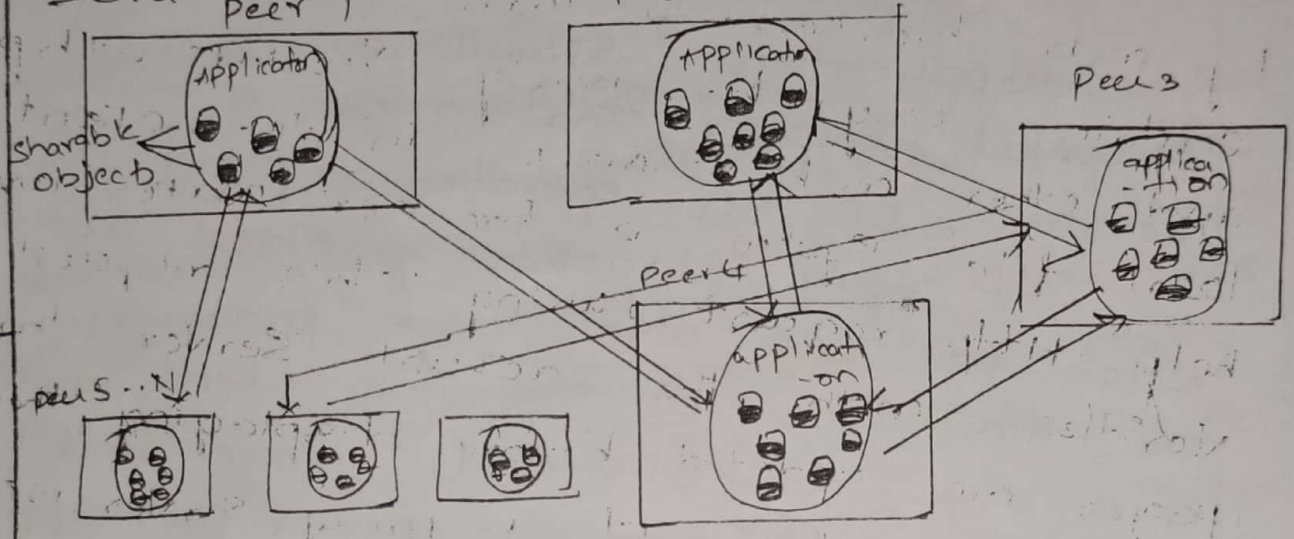
In reality, caches are employed widely by web browsers to maintain a cache that contains only the newly visited web pages and other web resources in the client's local file system. A browser makes use of this request to help HTTP request using this request to verify with the actual server that pages are updated before displaying them.

The objective of a proxy server is to extend the performance and availability of the service by decreasing the load on wide-area network and web servers through a firewall. For a client machine a site or over multiple sites, a shared cache is supplied by web proxy server.

→ peer processes:-

The peer processes architecture, all processes play akin roles and interact cooperatively with each other as peers, performing distributed activity or computation without making any differentiation among clients and servers.

In code in peer processes are responsible of maintaining consistency for application level resources and synchronizing application level resources actions based on the requirements. The scenario of peer processes architecture is shown in figure below.



Basically, a number of peer processes can take part in communication, whereas the communication pattern will be as per the need of application.

The elimination of server processes leads in the reduction of delays in inter-process communication that involves local objects accessing.

For example, assume an application of a distributed 'white board' that allows users to view and modify a picture shared between several computers. This can be made possible by implementing

as an application process at each site relying on middleware layers. So that event notification and group communication can be performed in order to notify the changes to the picture for all application processes. From this we can infer that users can have better interactive response from distributed shared objects than a server based architecture.

* Mobile Agents:

A mobile agent is a program under execution that goes from one system to another in a network. A network performs a task such as gathering information, providing results on behalf of someone. At each site visited by a mobile agent several calls are made to local resources. For example, to get single data base entries. This architecture when compared to static client that makes remote calls to some resources, perhaps carryout huge amount of data, with the substitution of remote calls with local ones there is a decrease in communication cost of time.

Mobile agents are possibly used in

cases such as,

(a) To install and maintain software on the computers within an organization

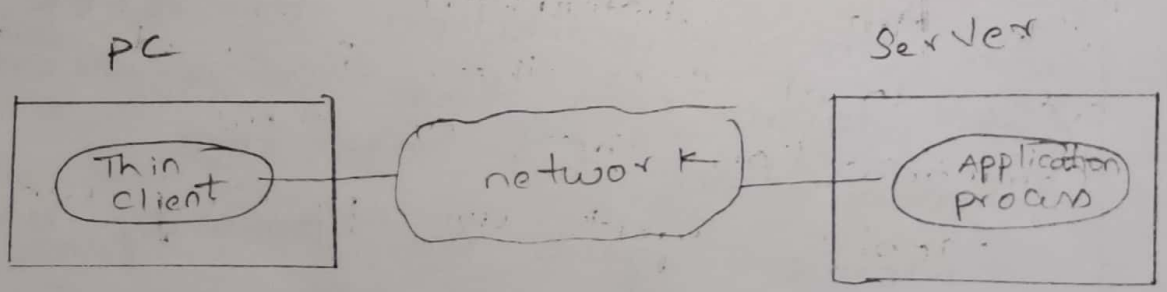
(b) To calculate the prices of products for different vendors by visiting the site of every vendor and carryout a sequence of database operations.

A recent example is worm program developed at xerox PARC. It is designed in such a way to transfer detailed computer by using of idle computers.

The environment accepting a mobile agent should determine which local resources are allowed to access. Because there is a potential security threat to the resources by mobile agents.

* Thin Clients:-

A Window-based user interface on a computer local to the user while running applications on a remote computer is supported by a software layer referred to as thin client.



Similar to the Network Computer Scheme,

this architecture is, has the same low

management and hardware costs unlike

network computer scheme, it runs applications

on a Computer Server, instead of downloading

the code into the user's computers.

A Computer Server has the capacity to

run applications simultaneously. It is a

powerful computer that runs a multiprocessor

version of an operating system (UNIX)

The highly interactive graphical activities

such as CAD and image processing are

the limitations of the thin client architecture.

Since, the transfer of image and vector

information between the client and application

process include both network and processor

system potentials, the delay experienced by

users is increased.

* The limited processing and communication capacities of computers and networks:

The performance issues arising from

the limited processing and communication

capacities of computers and networks are

illustrated below.

1. Responsiveness:

The interactive application users always

expect fast and consistent response, but

frequently client programs require accessing

of Shared resources. In case of remote Service, the response speed can be obtained not only with the load and performance of the Server and network, but also by means of delays in every involved slow components such as client and server OS involved in the communication, middleware services, and the program code used for service implementation. This can be overcome if the systems contain relatively few software layers, and if small amount of data is transferred b/w the client & server.

In web browsing clients, the locally cached pages and images which are held by the client application can be accessed in faster response mode. However, in case of graphical images that contain large volume of data, interactive response time is very slow.

3. Throughput -

It can be defined as the rate of completing the computational task, it is the traditional performance measure for computer systems.

For distributed system, throughput determines the ability of performing work for all its users. The factors such as processing speeds and data transfer rates

Affect distributed Systems

Data present on remote server, if needed to be passed from a server program to a client program, then it has to pass through all software layers present in both computers thus the throughput of network and software layer is needed to be determined.

* distributed system must be developed such that it should enable applications and service program. This is done, for concurrent proceeding by making use of available computational resources instead of going ahead for the same resources. The loads is said to be eliminated from the web server causing service improvement.

* Quality of Service:-

In a distributed system, as soon as the users are provided with the functionality they are in need of services like file service. Adaptability and resource availability are the important characteristics of service quality. Quality of service gets affected by some non-functional properties of system such as reliability, security and performance. Most of the computer system's design basically depends on the issues of reliability and security. From the performance

performance point of view, quality of service was basically referred to as, the response time and computational throughput, according to the definition, a distributed system, from performance point of view is strongly related to the interaction model.

In certain applications, the transfer of data streams from one program to another is done at a fixed amount of time, this is called time-critical data. This signifies that a video clip should be retrieved from the server, and then it is presented on the screen of the user. If the program involves displaying of successive clip with in some specified time limits, then its, then its service quality is said to be satisfied.

Service quality can be achieved based on the availability of the necessary computing and resources of network at the particular instance-of-time. But the performance of today's network, even with quite good performance characteristics is getting deteriorated due to heavily loaded traffic.

Web-caching protocols

The responses to request by "client" to web servers are cached by both web browsers and proxy servers. Because of this, the request a client can be satisfied with the response cached by the browser or by the proxy server between it and the web server. But updation can be done or made to relax in order to allow performance, availability and disconnect operation. To validate the response cached by browser or proxy, the updation web or original web-server is checked and if it fails the test, then the state-response is replaced with the fresh response returned by the web server. If the cached response is found to be sufficiently fresh, then validation is not required. Even if the web server acknowledges with the updation, of a request then also it doesn't inform the browsers and proxies.

The web server can inform it if it maintains a record of interested browsers and proxies for each of its resources. Servers enable browsers and proxies orders to determine if the cached responses are stale; it does this by assigning approximate; it does this by assigning approximate expiry times time.

The browsers and proxies along with their cached response stores the expiry time and server perform composition of response age the expiry time in order to determine if a cached response is state. Here response and the sum of the time taken to the cached the response and the server time. In this concept, the calculation is independent of computer clocks.

* Dependability Issues -

The dependability for computer system can be defined as correctness, security and fault tolerance. In most applications domain dependability is the crucial requirement. It is not only important in command and control activities but also play an important role in various commercial applications.

These applications may include internet commerce in which financial safety and soundness of the participants is dependent on the systems dependability on which they operate.

The issues arising in dependability are considered under the following sub-headings:

Fault Tolerance

Dependable applications should have the ability to continue its functions even in the presence of faults that may arise in software, hardware, and networks. Reliability of an application can be obtained by means of redundancy. Where in redundancy can be achieved by employing several resources. With the help of these resources, the system and application software can reconfigure and proceed its function even after the occurrence of faults.

But, redundancy is cost effective limiting its application at times, there by limiting the degree of fault tolerance.

9. Security

The requirement for security stresses on the consideration of confidential data and data from other resources that can be safeguarded against all the attacks for an instance, consider the database of a bank containing customer records with sensitive components and widely available components. The sensitive components must be visible to only some specific bank authorities. Then it is inappropriate to delete a system that can load the whole record of the patient into the desktop.

Computer may not provide a secure environment

* System Model:-
System model is a model used to create accurate assumptions regarding the systems that are to be modeled. It is used to create generalizations are considered about the systems that are to be modeled based on the assumptions made. These generalizations are considered as efficient generalizations for the purpose, algorithm or the properties which are required and guaranteed.

Factors of fundamental models

The factors that are to be considered for the construction of fundamental models.

(a) Security:-

The distributed systems are vulnerable to attacks because of their openness and modular nature. The security in fundamental model defines attacks from system design with the threat analysis and threats that are resistant to

(b) Interaction:-

It is the communication and co-ordination between the processes with the exchange of message

The interaction model in distributed systems must reflect the following facts.

- * Communication involves delay of specific duration
- * The delays limit the accuracy of process coordination
- * Difficult to maintain same time content across the components

(C) Failure:-

The operation of a distributed system is vulnerable to failure with the occurrence of failure in the network or in the computer which is connected to it. fundamental models describes the faults and analyses their effects. It also designs fault tolerable systems capable of running even in the presence of faults.

* Interacting processes:-

Interacting processes are responsible for the execution of activities in the distributed systems. Each process has a state that cannot be accessed or updated by another process. A state consists of variables and a set of data that can be updated and accessed by its state.

Factors effect the interacting process:-

In distributed system, two factors effect the interacting process:

1. performance of communication channels

9. Variation in Time

Performance of Communication Channels:-

The performance of communication channel is restricted. The communication in models the place by implementation of streams or by the transmission of message over the computer network. There are three performance characteristics associated with the communication over a computer network.

(a) Latency - Latency is the delay that occurs when the message is transmitted from one process to another. The delay includes the limit in between the initiation of message transmission of its reception. Latency consists of the following three factors:

(*) The variation in time consumption by the operating system communication at the source and destination processes is on the basis of current load on the operating system.

(*) The time consumed for the transmission of the first stream of bits of string to reach the destination via network. This can be seen in radio signals that travel back and forth the satellite for the transmission of a message.

(*) The increase in delay to access the network is due to heavy traffic in it.

(b) Jitter:-

It is the difference in time during the transmission of message. This difference in time can be seen when a series of audio data samples with variant time intervals are played resulting into distortion in time.

(c) Bandwidth:-

It is the amount of data that can be transmitted on network in a given time period. If same network is being used by many communication channels, then they have to share the bandwidth.

2. Variation in Time:-

In distributed system, every computer has an internal clock that gives the current time to the local processes. Hence the processes running on distinct computers can relate the time stamps with their events. But this does not guarantee the supply of accurate time. This is due to the clock drifts whose rate varies from one process to another. A clock drift rate is the relative rate at which a computer clock drifts from an accurate clock that is to be referred. It is obvious that clock timings vary significantly unless some corrections have been applied on it.

The Correction of time on Computer clock can be done in many ways. Radio receivers can be used by the Computers to know the accurate time from the Global Positioning System (GPS). It gives the time with an accuracy of about 1 microsecond. The cost of GPS receiver is high and also they can't be operated inside the buildings. A Computer, that possesses accurate time giving services like GPS can transmit the messages containing accurate time to other Computers. However there may be delay in messages which ultimately affect the agreement made by the local Computers.

Variants of Interaction Model

There are two variants of interaction

- model:-
- (i) Synchronous distributed Systems
 - (ii) Asynchronous distributed Systems.

Synchronous distributed System is based on the assumptions of time, whereas asynchronous distributed System has no assumptions about time.

(i) Synchronous Distributed Systems
A Synchronous distributed System is defined to be system which has

Known upper and lower bounds associated with the execution of each and every phase of a process. The transmitted messages are received within a bounded time. The process has local clocks, the drift rate of which is associated with a known bound.

However, the values of upper and lower bounds chosen are not guaranteed.

The behaviour of a Synchronous System in a real distributed systems, the process that performs the tasks should be familiar with the resource requirements. These requirements must guarantee the processor cycles and network capacity sufficient enough to perform task process should also be supplied with clocks having bounded drift rate.

(ii) Asynchronous Distributed System :-

An asynchronous distributed system is a system in which there are no bounds on execution speed of a process, message transmission delays and drift rates between the local clocks.

Asynchronous models do not have any timing conventions. This criteria models the internet in which a server or a network load neither has intrinsic bounds on the server nor on their duration.

This can be seen in the transmission of file via FTP. However an email message can take too many days to deliver.

Some problems associated with the design can be dealt even with out the time assumptions. This can be seen in web browsers. When web is unable to provide a response within a specific time

A solution reasonable for asynchronous distributed system is also effective for distributed system.

In general, the distributed systems are asynchronous because it requires the sharing of processors and networks among the processes and communication channels respectively. This can be seen when too many processes of unknown character share a single processor, then the performance of any one process among them cannot be guaranteed.

* Event Handling in Distributed System

In distributed system, event ordering is essential to describe the execution of system on the basis of its events and their ordering.

Consider the events involving a group of users (user 1, user 2, user 3, user 4) who

exchange the email messages among themselves. Here, user 1 sends the email to which user 2 and user 3 reply back. In actual, user 1 has sent the message first, to which user 2 replies immediately then user 3 leads user 1's sent and user 2's replied messages and replies accordingly.

Due to delay on the delivery of message and replies accordingly that they users may view these messages in wrong order. For example, user 4 may view the messages in the following order:

- Order of message delivery seen by user 4
- 1st → user 3
 - 2nd → user 1
 - 3rd → user 2

If user 1, user 2, and user 3 have synchronized clocks on their computers, then the sent messages includes the time of the local computer's clock. This allows the reception of messages to its users according to their time order.

However, in a distributed system, no clock is perfectly synchronized. So, a model of logical time was proposed by Lamport to provide events ordering delivered messages can be inferred without

any alternative

following are some of the aspects of logical ordering that can be applied to the problem of event ordering. Suppose there are two users, user1, user2, send the messages msg1 and msg2, then it is obvious that:

* A message is received only after it was sent. user1 $\xrightarrow{\text{Sends msg1}}$ before user receives msg1
user2 $\xrightarrow{\text{receives msg2}}$ before user1 $\xrightarrow{\text{receives msg2}}$

* A message is replied only after it has received any message.
user2 $\xleftarrow{\text{receives msg1}}$ before sending msg2

These two aspects can be used further

When there are too many users and messages to be transmitted and received.

Then a number can be assigned to each event based on its logical ordering.

* Failure Model:

Failure model defines methods of failure occurrences to facilitate an understanding of failure effects.

The taxonomy that distinguishes between failures of processes and communication channels can be understood by the

following concepts:

1. Omission failures
2. Arbitrary failures
3. Timing failures

Omission failures

The faults classified as omission failure refer to cases when a process or communication channel fails to perform actions that supposed to do.

Process omission failures

The primary reason in a process for omission failure is its crash. That means process has been terminated completely. The design services that can resist the failure can be simplified on the assumption that the services crash cleanly. In a clean crash a process terminates or executes correctly. Due to crash, the process does not respond to invocation messages repeatedly. This leads other processes in identifying the crash. Timeouts allot a specific time period to process with which it should occur. In an asynchronous system a time out implies that a process is not responding. The reason for this may be its crash or slow performance or undelivery of messages.

A crash that is identified by other processes is called fail-stop. This can be introduced in a synchronous system as:

* if message delivery is correct

* if processes employ timeouts identify failure of other processes

Consider an example, if the reply to a process is not delivered within the specified time limit, then the process that is waiting for the reply concludes that sender process has been failed.

(b) Communication omission failures

Send and receive are two communication primitives. Consider the example where send operation is performed by process A. It inserts a message in its outgoing message buffer. Communication channel transmits this message to process B's incoming message buffer. Receive operation is performed by process B that takes the message from its incoming message buffer and delivers it.

The communication channel raises an omission failure if the message from A's outgoing message buffer is not transmitted to B's incoming message buffer. This is called the incoming message buffer.

dropping messages which is due to insufficient buffer space at the receiver site or at the intervening gateway or by error at network transmission. These consequences are identified by the checksum that is attached with the message content.

Following are the failures resulting due to loss of messages:

(i) Send-omission failures:
occur in between sending process and outgoing message buffer.

(ii) Channel-omission failures:
occur in between outgoing message buffer and incoming message buffer.

(iii) Receive-omission failures:
occur in between incoming message buffer and receiving process.

2. Arbitrary Failures

Arbitrary or Byzantine failure describes the semantics of a worst failure in which there is possibility of occurrence of any type of error.

In a process an arbitrary failure is caused when it skips any of its execution phase or adopts a new phase for execution. Hence, detection of arbitrary failures in processes cannot be known.

by checking the responses to invocations because there is a possibility for the process to skip the process reply. For example, a process may set incorrect values in its data items or it may return an incorrect value.

In communication channels, arbitrary failures can take the form of a corrupt message, delivery of a wrong message or delivery of a correct message twice or many times. However, these occur very few times. In communication channels, arbitrary failures can occur because the communication slows down. They can identify them and discard corrupt messages. For example, checksums determine corrupt messages and a message sequence identifies wrong or duplicated messages.

Omission & arbitrary failures

| class of failures | Affects | Description |
|-------------------|---------|---|
| fail-stop | process | process halts and remains halted. other processes may detect this state. |
| crash | process | process halts & remains halted. other processes may not be able to detect this state. |

| | | |
|---------------|-----------------|---|
| omission | channel | A message inserted in an outgoing message buffer never arrives at the other end |
| Send-omission | process | A process completes a send but the message is not put in its outgoing message buffer |
| Receive | process | A message is put in process incoming message buffer, omission that process does not receive it |
| Arbitrary | process/channel | process/channel exhibits arbitrary behaviour. It sends/transmits arbitrary messages at arbitrary times, commit omission a process may stop or take an incorrect state |

3: Timing failures:-

It is applicable in synchronous distributed system where as in asynchronous distributed system, it is not guaranteed

| class of failure | affects | Description |
|------------------|---------|---|
| clock | process | process's local clock exceeds the bounds on its rate drift from real time |
| Performance | process | process exceeds the bound the interval between two states |
| Performance | channel | A message transmission takes more than the stated bound |

The term reliable communication is defined in terms of validity and integrity as follows:

Validity: any message in the outgoing message buffer is eventually delivered to the incoming message buffer.

Integrity: The message received is identical to one sent, and no messages are delivered twice.

(a) Synchronous Distributed System: In this system, time limits are fixed on process execution time, message delivery time and clock drift rate. A process rate is effected if it exceeds its time bounds. A channel transmission exceeds the specified bounds if the time taken for the transmission exceeds the specified bounds.

A clock failure effects its process if its local time exceeds its drift rate bounds from real time. There three failures fail to respond within the specified bounds.

(b) Asynchronous Distributed System: In asynchronous distributed system, timing failure cannot be guaranteed, and there is a possibility of a heavily loaded server to respond late.

Threats to integrity is derived from two independent sources, protocols and intruders.

Protocols: It is the protocol that retransmits the messages without discarding the message that is received twice. Such protocols assign sequence number to the messages to identify those messages that are delivered twice.

Intruders: It may send fake messages or interface with old messages. In such cases, measures should be taken to preserve integrity.

* Security model:

A Security model is constructed on the basis of distributed system architecture where the processes, objects and interaction channel have to be protected against threats. This model provides the analysis of threats, different forms of attacks by threats, evaluation of risks and consequences for each at minimum cost.

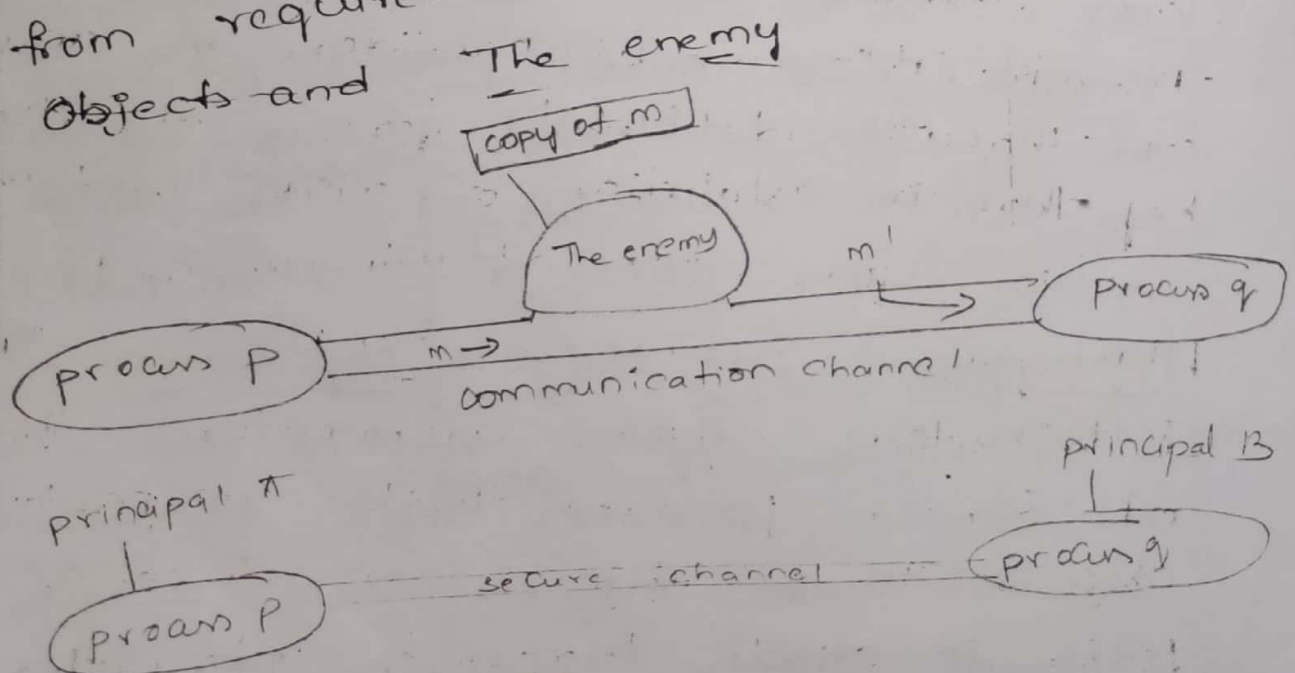
Management of objects by Server: The distributed system architecture consists of Server processes that encapsulate objects. These objects are accessed by different clients through interactions. The

user/client process invokes the server to perform operation on the objects. Server performs the operation for each invocations and sends results to client.

The whole process requires the association of authority to both the processes (client/server) such authorized processes are known as principal.

The server's authority is to verify the identity of client principal. It checks whether the client principal has sufficient access rights to perform operation on the particular invoked object or not.

On the other side, the client verifies the identity of server principal and confirms whether the result received is from required server or not.



Secure channels

→ Threats generated from potential enemies

The threats that can be generated from a potential enemy are,

1. Threats to processes
2. Threats to Communication channels
3. Denial of Service

1. Threats to processes:-

In a distributed system, a process reserved for handling incoming requests

This process might receive any message

from another process, it is not necessary

that it can identify the identity of

process, the sender. Although communication

protocol like IP of each message contains

the source computer address, there is

possibility for an enemy to generate a

message with a fake sender address.

Thus it leads to a threat for both client

and servers to function correctly.

(i) Servers:- A server receives invocation

from several clients which makes it

difficult to identify the identity of an

invocation.

Actually the server can identify the

identity of any invocation. But, there

is a possibility of generating false

identity by an enemy due to which

the server cannot make a decision of

performing the operation or rejecting it
(ii) Client: The client that receives the message from the server is unable to determine whether the message is from the desired server or from an enemy. This received message might not match the original in location.

2. Threats to Communication Channels
Across communication channels an enemy can attack over information by means of altering, copying or injecting. This leads to a threat for privacy and integrity of the information. An enemy can attempt to serve messages copy in order to reply them at later time.

3. Denial of Service
It is an attack where the enemy overloads the physical resources of authentic users. The enemy does this by making several unnecessary in locations on service or transmitting the messages in a network. The main reason behind this attack is to delay or prevent the actions of other users.

INTERPROCESS COMMUNICATION

* API for Internet protocols:-

→ characteristics of Interprocess Communication:-

Message passing between a pair of processes can be supported by two message communication operations: send and receive.

In order for one process to communicate with another, one will send a message to destination and other process at the destination

requires the message. This activity involves the communication of data from the sending process to the receiving process and may involve the synchronization of the two processes.

Synchronous and Asynchronous Communication:-

→ In Synchronous inter process communication both send and receive operations are blocking system calls. It means that sending process stops execution until it receives a reply from receiver and receiving process stops execution until it receives from sender which ensures both are synchronized.

→ In asynchronous inter process communication, in this sending process uses a non blocking send operation. AS a result

Sending process doesn't wait after
Sending message instead it continues
execution. However the receiving process can
use either blocking or non-blocking
receiving operation. If blocking receive
is used, process waits until it receives
message else process doesn't wait and
continues its execution.

Characteristics

Message Destination in a distributed
System, messages are specified using
internet address; port address pairs.

The port address refers to process on
destination system which is supposed
to receive a message. This mechanism
doesn't provide transparency following
approaches are used.

- a) program can refer services by
specific names and there are special
name into server location
- b) most operating system uses location
independent identifiers which are later
mapped to low level addresses.

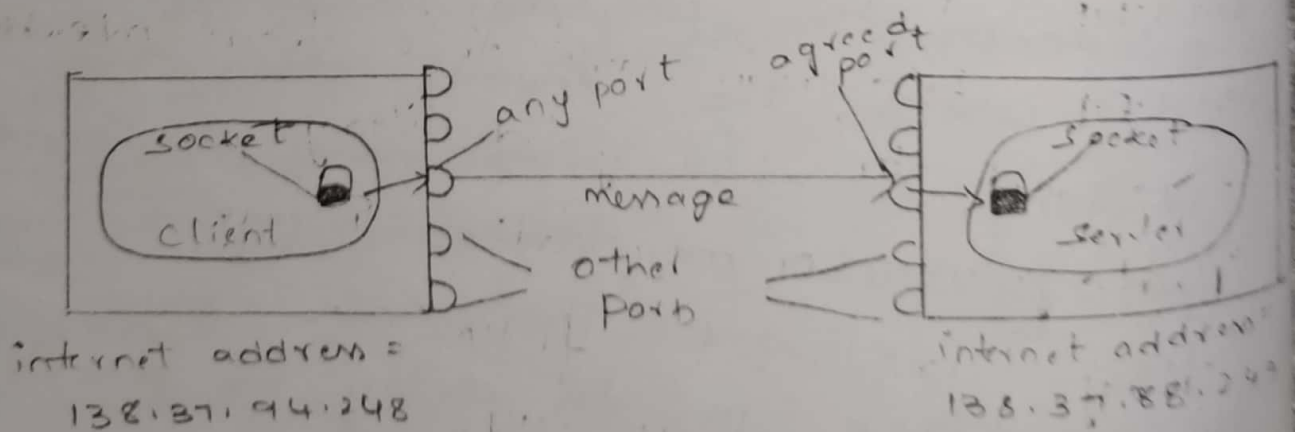
Reliability: Validity and integrity are
two important ingredients of reliable
communications. Reliability refers to

guarantee that messages are delivered to destination without being dropped and integrity refers to packets are received with damage and duplication.

Ordering: - messages are said to be in order if they arrive at receiver side in the same order that they were sent by sender. Some applications will not accept packets if they are not received in order.

→ Sockets:

A socket acts as an endpoint in a connection between a client and a server present in a network using a socket, a process running on one computer can communicate with other process running on distinct computer.



Sockets & ports

→ UDP Datagram Communication

A datagram sent by UDP is transmitted from a sending process to a receiving process without acknowledgment or retries. If failure occurs, the message may not arrive. A datagram transmitted between processes when one process sends it and another receives it. To send or receive messages a process must first create a socket bound to an Internet address of the local host and local port. A server will bind its socket to a server port.

→ uses of UDP

The uses of UDP are listed below.

- (1) UDP implements DNS that looks up the DNS names into the internet.
- (2) UDP allows VoIP to run over it.
- (3) UDP datagram are free from overhead released to message delivery.

→ Datagram Communication Issues

They are as follows :-

1. Size of message
2. Blocking
3. Receive message from any source

III. Size of Message: Size of packet allowed by IP protocol ranges upto 2¹⁶ bytes along with message packet holds header also. A program that receives the message header specifies array of bytes to hold message. If message exceeds the allowable range then it is truncated after message has been received.

Blocking: usually sockets provide non-blocking send operation and blocking receive operation for datagram communication.

(a) Send operation sends message to relevant UDP and IP protocols. These protocols transfer message to desired destination. The received message is placed in a queue for that socket which is intended for destination port and if program doesn't have socket message is discarded.

(b) The receive operation is blocked till a datagram received and time out has been set on socket. In case program that invokes receive operation is found to be busy, it must assign new thread to use.

Receive operation does not mention messages from any source. The

Source of message. Its invocation will address a message from any source. It enables recipient to check the source of message by returning the internet address and local port of the sender and also receive operation. enables connection from socket to remote port.

→ Java API for UDP Datagram!

Java API provides two types of classes, namely Datagram packet and Datagram socket.

Datagram packet:- Datagram packet class provides a constructor that creates an instance which can be transmitted from one process to another such that one process sends and other process receives it. The instances are created by array of bytes consisting of message length of message, internet address and local port number of destination socket.

Datagram socket:- Datagram socket class provides support to socket for sending and receiving UDP datagrams. It provides a constructor which accepts

port number as arguments. This type of constructor is used by programs to deal with the specific port.

→ Tcp Stream Communication:-

The API to the TCP protocol, which originates from BSD 4.x UNIX, provides an abstraction of a stream of bytes to which data may be written and from which data may be read. The following characteristics of the network are hidden by the stream abstraction:

Message Sizes:- The application can choose how much data it writes to a stream or reads from it. It may deal in very small or very large sets of data. The underlying implementation of a TCP stream decides how much data to collect before transmitting it as one or more IP packets. On arrival, the data is handed to the application as requested. Applications can, if necessary, force data to be sent immediately.

Lost Messages:- The TCP protocol uses an acknowledgment scheme. As an example of a simple scheme, the sending end keeps a record of each IP packet sent and the receiving end acknowledges all arrivals. If the sender does not receive an acknowledgment within a time out, it retransmits the message. The more

Sophisticated sliding window scheme cuts down on the number of acknowledgements sent managers required.

Flow Control: - The TCP protocol attempts to match the speeds of the processes that read from and write to a stream. If the writer is too fast for the reader, then it is blocked until the reader has consumed sufficient data.

Message duplication & ordering: - message identifiers are associated with each IP packet which enables the recipient to detect and reject duplicates, or to reorder messages that do not arrive in sender order.

Message destinations: - A pair of communicating processes establish a connection before they can communicate over a stream. Once a connection is established, the processes simply read from and write to the stream without needing to use internet addresses and ports. Establishing a connection involves a connect request from client to server followed by an accept request from server to client before any communication can take place. This could be considerable overhead for a single client-server

request and reply.

Uses of TCP:-

The Services that run across TCP connection with reserved port number are as follows:

HTTP :- Hyper Text transfer protocol:-

This service is utilized in order to establish communication between web browsers and web services.

FTP (File Transfer protocol):- This service is utilized in order to incorporate directories running on remote computer to be browsed on other computer.

Telnet:- This service is utilized by terminal session in order to gain access to remote computer.

SMTP:- It is used for sending mails between computer.

→ Issues Related to Stream Communication:-

The issues are:

(i) Matching of Data Items:- In TCP stream communication the two processes must agree on the contents of data transmission across stream. It is essential because the contents written by one process must be read in order by other process.

(2) Blocking

In this communication data is written to stream is stored in the form of queue across destination socket. When data is available at destination tries to read data from channel else if unavailable it blocks itself until data is available.

(3) Threads

In this communication when client establishes connection with server. It accepts request and forms separate thread during this it blocks other client threads with which it already established.

* External data Representation & marshalling

The information stored in running programs is represented as data structures - for example by sets of interconnected objects - whereas the information in messages consists of sequences of bytes. Irrespective of the form of communication used, the data structures must be flattened before transmission and rebuilt on arrival. The individual primitive data items transmitted in messages can be

data values of many different types, and not all computers store primitive values such as integers in the same order. The representation of floating-point numbers also differs between architectures. There are two variants for the ordering of integers. The so-called big-endian order, in which the most significant byte comes first, and little-endian order, in which it comes last. Another issue is the set of codes used to represent characters. For example, the majority of applications on systems such as UNIX use ASCII character coding, taking one byte per character, whereas the unicode standard allows for the representation of texts in many different languages and takes two bytes per character.

One of the following methods can be used to enable any two computers to exchange binary data values:

* The values are converted to an agreed external format before transmission and converted to the local form on receipt; if the two computers are known to be the same type, the conversion to external form can be omitted.

* The values are transmitted in the send format together with an indication of the format used, and the recipient converts the values if necessary.

Note, however that bytes themselves are never altered during transmission. To support RMI or RPC, any data type that can be passed as an argument or returned as a result must be able to be flattened and the individual primitive data values represented in an agreed format. An agreed standard for the representation of data structures and primitive values is called an external data representation.

Marshalling: It is the process of taking a collection of data items and assembling them into a form suitable for transmission in a message. unmarshalling is the process of disassembling them on arrival to produce an equivalent collection of data items at the destination. Thus marshalling consists of the translation of structured data items and primitive values into an external data representation.

* CORBA's Common data representation which is concerned with an external representation for the structured and primitive types that can be passed as arguments and results of RMI in CORBA. It can be used by a variety of programming languages.

* XML or extensible markup language, which defines a textual format for representation of structured data.

In first two cases, the marshalling and unmarshalling activities are intended to be carried out by a middleware layer without any involvement on the part of the application programmer. Even the case of xml, which is textual and therefore more accessible to hand encoding, software for marshalling and unmarshalling is available for all commonly used platforms and programming environment.

CORBA's Common Data Representation (CDR):

CORBA CDR is the external data representation defined with CORBA 2.0. CDR can represent all of data types that can be used as arguments and return values in remote invocations in CORBA. These consist of 15 primitive types, which include Short (16-bit), long (32-bit), unsigned short, unsigned long, float (32-bit), double (64-bit), char, boolean (TRUE-FALSE), octet (8-bit) and any together with a range of composite types. Each argument or result in remote invocation is represented by a sequence of bytes in the invocation or result message.

Primitive type :- CDR defines a representation for both big-endian and little-endian orderings. The values are transmitted in the sender's ordering, which is specified in each message. The recipient translates it requires different orderings. For example, a 16-bit short occupies two bytes in the message and for big-endian ordering, the most significant bits occupy the first byte and the least significant bits occupy the second bytes.

CORBA CDR for Constructed types

| Type | Representation |
|------------|---|
| Sequence | Specifies length followed by elements in order |
| String | Specifies length followed by character in order |
| Array | Array elements in order |
| Struct | order of declaration of components |
| Enumerated | Specifies unsigned long |

CORBA CDR message :-

Each primitive value that contains constructed type are added to the sequence of bytes in order.

Index in sequence
of bytes

4 bytes

0 - 3

7 → length of string

4 - 7

masa → masarat

8 - 11

rat

⋮

16 - 19

9 → length of string

"Hyde" → Hyderabad

20 - 23

"rabad"

24 - 27

"d"

32 - 35

⋮

Constructed types:- The primitive values that comprise each constructed type are added to a sequence of bytes in a particular order.

2. Java object serialization:- An object instantiated from Java class. The both objects and primitive data values existing in Java RMI is passed as arguments and outcomes of method invocations. For instance, Java class equivalent to struct in CORBA IDL is

```
public class Stud implements Serializable
```

{

```
    private String name;
```

```
    private String location;
```

```
    private int year;
```

```
    public Stud (String bname, String blocation, int byear)
```

```

    {
        name = bname;
        location = blocation;
        year = byear;
    }
}

```

The above stated class implements Serializable interface and does not include methods. As the class implements Serializable interface located in java.io package, it allows serialization of instances. It refers to an act of transferring an object or connected set of objects at equals levels in order to achieve serial formation.

3. XML:-

It is a simple markup language for describing structure of document. It has been derived from standard generalized markup language for large scale electronic publishing and for exchange of data in web. XML has a set of rules for creating other markup languages called as "meta markup" language.

XML Schema:-

It is for constructing xml documents

-nts.

1. XML Schema depends on XML Syntaxes
for their documentation:-

We can directly make use of XML editors as well as parsers for generating XML Schemas and parsing.

2. XML Schema documents support usage
of data types:-

We can impose restrictions on data and define few data patterns. Data belonging to one data type can be easily converted to other.

3. using XML Schema we can ensure
Secure data communication between
Various entities:-

Whenever any of two distinct entities resort to data communication, they both have to agree on same patterns to ensure error free communication.

4. As XML Schemas are written in
XML, they are extensible in nature:-

Basing on one standard type, we can derive numerous data types are required.

A given Schema can be implemented
or reuse other Schema etc.
5. xml Schemas can be taken as
Successors to DTD:-

xml Schemas extensively support
namespaces as well as data types.
They have address sets of provisions
which make xml Schemas more powerful
and richer than DTD.

creation of xml schema:-

```
<xs:element name="d.o.b", type="xsd:date">
```

```
</xs:element>
```

For attributes declaration.

```
<xs:attribute name="name of attribute"  
type="format of attribute">
```

Building blocks of xml:-

Building blocks of xml are as follows

1. Tags:- Tags usually mark up elements
embedded between them. There are
two tags namely opening and
closing tags.

eg: <heading> final notice </heading>

2. Elements:- There are major entities
to represent or to structure data
in these documents.

Syntax:-

<ELEMENT name-of-element (Content) >

<ELEMENT name-of-element (EMPTY) >

<ELEMENT name-of-element (#PCDATA) >

<ELEMENT name-of-element (CDATA) >

3. Attribute:-

Attributes are defined after elements. Attributes are followed by their value.

Syn!

<ATTNIST name-of-element, name-of-attribute, attribute-type, default-value >

4. PCDATA:-

It is a data which will be parsed by a given parser. They may carry certain tags which are nothing but mark up elements and finally the entities.

5. CDATA:- This is just opposite to PCDATA. There is no entity expansion. tags are not treated as mark up element and finally CDATA is not parsed at all.

6. Entities:-

We use few variables in xml instead of some frequently used text. Such variables are nothing but entities. There are two types of entities. Internal entities and external entities.

→ Client-Server Communication

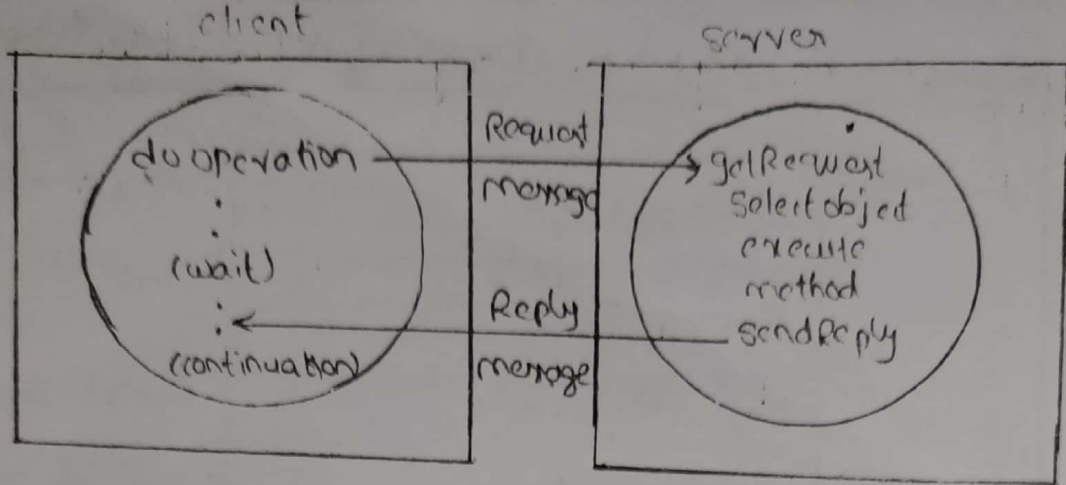
This form of communication is designed to support the roles and message exchanges in typical client-server interactions.

A protocol built over datagrams avoids unnecessary overheads associated with the Tcp Stream protocol.

- particular:
- acknowledgements are redundant, since requests are followed by replies.
 - establishing a connection involves two extra pairs of messages in addition to the pair required for a request and a reply.
 - flow control is redundant for the majority of invocations, which pass only small arguments and results.

Request-Reply protocol

The following protocol is based on trio of communication primitives: dooperation, getRequest and sendReply. RMI that it passes a remote object reference for the object whose method be invoked in the request message.



The do operation :- The method is used by clients to invoke remote operations. Its arguments specify the remote object and which method to invoke, together with additional information required by the method. The first argument of do operation, is an instance of the class RemoteObjectRef.

GetRequest :- It is used by a server program to acquire service requests. When the server has invoked the method in the specified object it then uses sendReply to send the reply message to the client. When the reply message is received by the client the original do operation is unblocked and execution of the client program continues.

The information to be transmitted in a request message or a reply message. The first field indicates whether the message is a Request or Reply message. The second field requested contains identifier

Request reply message Structure

| | |
|------------------|------------------------|
| Message Type | int (0 = req, 1 = rep) |
| Request Id | int |
| object Reference | Remote Object Ref |
| Method Id | int or method |
| arguments | ll array of bytes |

→ RPC exchange protocols - of handling
Three protocols, which produce different behaviours in the presence of communication failures, are used for implementing various types of RPC. They were originally identified by Spector [1982]:

- The request (R) protocol
- The request-reply (RR) protocol
- The request-reply-acknowledge-reply (RRA) protocol.

Request-protocol:-

In this protocol, client sends single request to server. In this client sends next request immediately after first request is sent as client does not need any confirmation that request has passed to server.

2. Request-Reply Protocol:-

Here client sends request to server and server sends reply message to client. It can be used when there is a value to be returned from remote method.

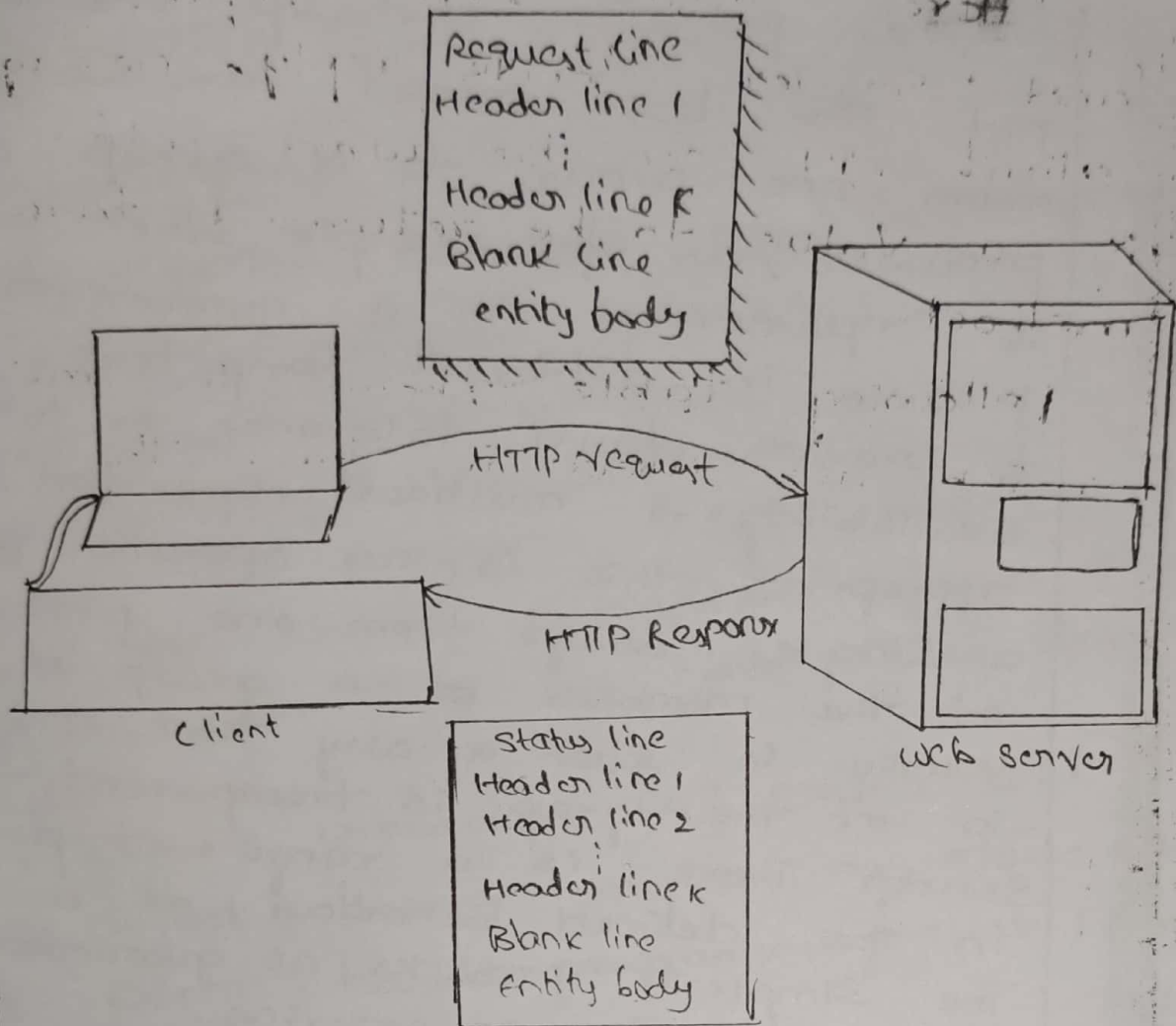
Following table shows the message which are used in these three protocols.

| Message Passed by | Name of the protocol | | |
|-------------------|----------------------|---------|----------------------|
| | R | RR | RRA |
| client | Request | Request | Request |
| server | | Reply | Reply |
| client | | | Acknowledge Reply |

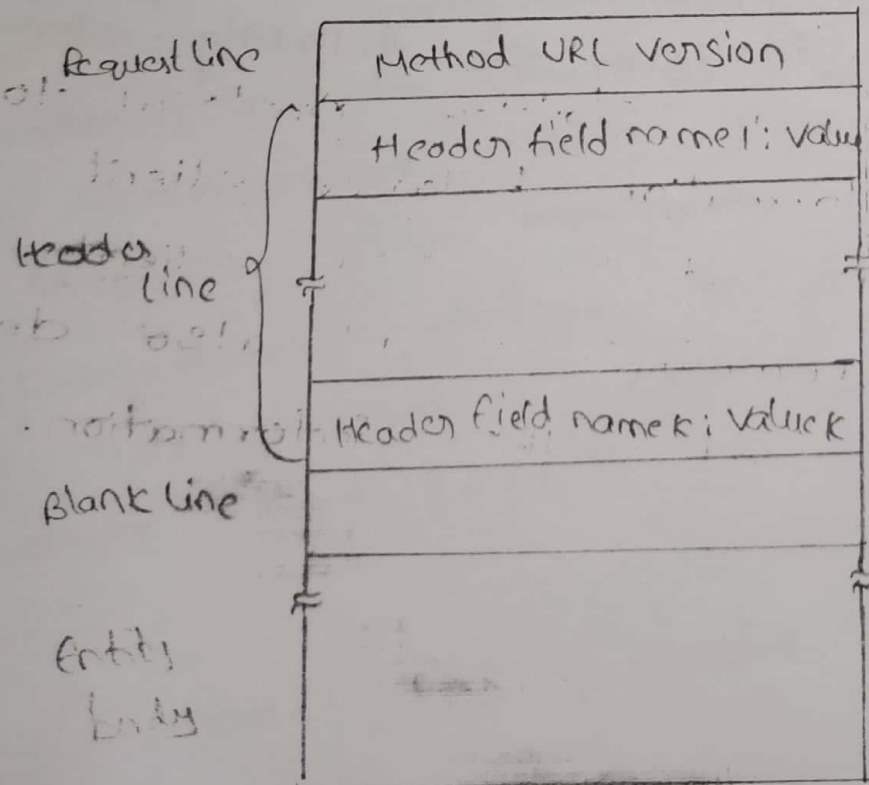
Http Protocol:-

It is a stateless protocol to transmit information between client and server. It defines how client and server communicate. It also defines structure of this information.

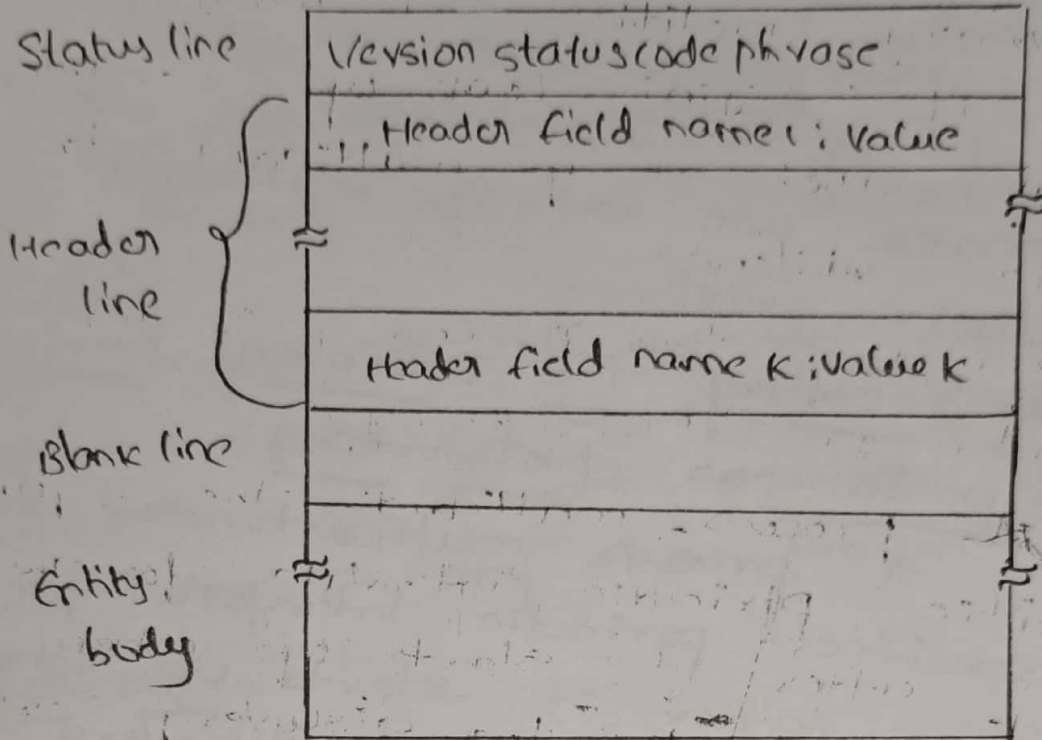
Basic HTTP Model



HTTP Request :-



HTTP Response:-



→ Reliability and ordering of multicast:-
effects of reliability and ordering
are as follows.

1. Fault tolerance depending upon replicate

Services:-

Consider a replicated service which consists of a group of servers. All the servers start at same initial state and perform the same operation in the same order in order to be consistent with each other.

2. Spreading the Event Notification

A specific application defines the necessary qualities of multicast. Example of it is Jini look up service which make use of IP multicast to notify their existence.

3. Locating the Discovery Servers in Spontaneous Networking

A process multicasts the requests at specific periodic intervals for the same time after it starts up when it wants to locate the discovery servers. Any request which is lost is not considered.

4. Enhanced performance through Replicated Data

Consider such a case where replicate data by itself distributed by means of multicast messages instead of operations on data. The loss of messages, and inconsistent sequence have an effect which is based on technique of replicating and importance of updating all replicas.

Group Communication:-

The pairwise exchange of messages is not the best model for communication from one process to a group of other processes, as for example when a service is implemented as a number of different processes in different computers, perhaps to provide fault tolerance or to enhance availability. A multicast operation is more appropriate - this is an operation that sends a single message from one process to each of the members of a group of processes usually in such a way that the membership of the group is transparent to the sender. There is a range of possibilities in the desired behaviour of a multicast. The simplest provides no guarantees about message delivery or ordering.

Multicast messages provide a useful infrastructure for constructing distributed systems with the following characteristics

Fault

DISTRIBUTED OBJECTS & REMOTE INVOCATION

Introduction:

Application consists of cooperative programs (distributed programs) running in different processes. These processes may generally related to different computers. The programs might have a need to invoke operations from other processes. In this we have three design models.

1. Remote procedure call model (RPC)
2. Object based programming model (OBP)
3. Event based programming model (EBP)

RPC :-

This model is followed to develop client/server programs in different processes, where the client program can communicate with server programs by invoking the methods of objects to procedures define in server programs. This process is called remote procedure call (RPC)

OBP:

This model is followed to develop

Objects in different process, where the objects can communicate one another by invoking the methods, of objects living in any other process. This process is called Remote method invocation (RMI)

EBP:-

This model is used to develop the programs where the objects can receive notification of events occurred at other objects living in any other process.

Middleware

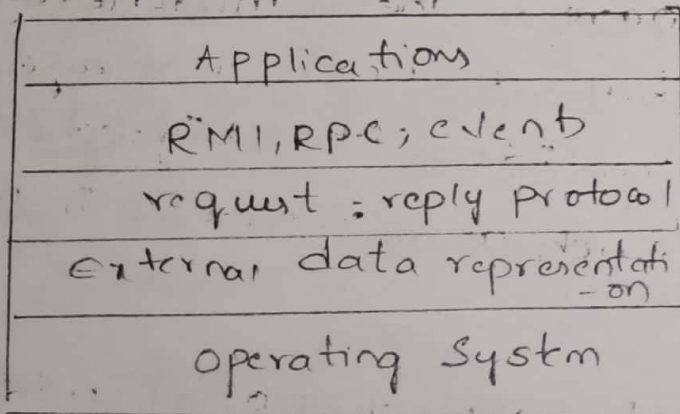
Software that provides a programming model above the basic building blocks of processes and message passing is called middleware. The middleware layer uses protocols based on message between processes to provide its higher-level abstractions such as remote invocations and events.

An important aspect of middleware is the provision of location transparency and independence from the details of communication protocols, operating systems and computer hardware. Some forms of middleware allow the separate components to be written in different programming

languages.

Location transparency: In RPC the client that calls a procedure cannot tell whether the procedure runs in the same process or in a different process, possibly on a different computer. Similarly in RMI the object making the invocation cannot tell whether the object it invokes is local or not and does not need to know its location. In distributed event-based programs the objects generating events and the objects that receive notifications of these events need not be aware of one another's locations.

Middleware layer



communication protocols

The protocols that support the middleware abstractions are independent of the underlying transport protocols.

for example, the request-reply protocol can be implemented over either UDP or TCP.

Computer hardware -

Three agreed standards for external data representation are described above. These are used when marshalling and unmarshalling messages. They hide the differences due to hardware architectures such as byte ordering.

Operating systems -

The higher-level abstractions provided by the middleware layer are independent of the underlying operating systems.

Use of several programming languages

Some middleware is designed to allow distributed applications to use more than one programming language.

In particular, CORBA allows client written in one language to invoke methods

in objects that live in server programs written in another language. This is

achieved by using an interface definition language or IDL to define interfaces.

Interfaces

The interface of a module specifies the procedures and the variables that can be accessed from other modules. Modules are implemented so as to hide all the information about them except that which is available through its

interface

Interface in distributed systems

The CORBA IDL interfaces can specify attributes, which seems to break this rule. However, the attributes are not accessed directly but by means of some getter and setter procedures added automatically to the interface.

Service interfaces

In the client-server model, each server provides a set of procedures that are available for use by clients. For example a file server would provide procedures for reading and writing files. The term service interface is used to refer to the specification of the procedures offered by a server, defining the types of the input and output arguments of each of the procedures.

Remote interface:-

In the distributed object model, a remote interface specifies the methods of an object that are available for invocation by objects in other processes, defining the type of input and output arguments of each of them.

Interface definition languages:-

IDL provides the notations for defining interfaces. These interfaces allow the objects written in different programming languages to communicate to another.

The following example uses COBRA as IDL

```
// In file employee.idl
```

```
struct employee {
```

```
    string ename;
```

```
    string eloc;
```

```
    long age;
```

```
};
```

```
interface EmployeeList {
```

```
    readonly attribute string listname;
```

```
    void addEmployee(in employee e);
```

```
    void getEmployee(in string ename, out Employee e);
```

```
    long number();
```

```
};
```

The above example shows an interface

"EmployeeList". It specifies the method implemented by remote objects. They are addEmployee(), getEmployee(), number().

"addEmployee" method has one argument which is specified as "in" it means that it is an input argument.

"getEmployee" method has two arguments. The first argument is specified as "in" and the second argument is specified as "out". The second argument shows that it is an output argument.

* Communication between distributed objects.

The object-based model for a distributed system extends the model supported by object-oriented programming languages to make it apply to distributed objects.

→ Object model: A brief description of an object-oriented program, for example in Java or C++, consists of a collection of interacting objects, each of which consists of a set of data and a set of methods. An object communicates with other objects by invoking their methods, generally passing arguments and receiving results. Objects can encapsulate their

data and the code of their methods

-ds. Some languages, for example Java and C++, allow programmers to define objects whose instances

can be accessed directly. But for use in a distributed object system, an object's data should be accessible only via its methods.

Object references:- objects can be accessed

via object references. For example

in Java, a variable that appears to

hold an object actually holds a reference

to that object. To invoke a method on

a object, the object reference and

method name are given, together with

any necessary arguments. The object

whose method is invoked is sometimes

called the 'target' and sometimes

the 'receiver'.

Interfaces:- An interface provides a

definition of the signatures of a set

of methods, without specifying their

implementation. An object will provide

a particular interface if its class

contains code that implements the

methods of that interface.

Actions:- The receiver executes the

appropriate method and then returns

control to the invoking object, sometimes supplying a result. An invocation of a method can have three effects:

1. The state of the receiver may be changed.
2. A new object may be instantiated for example, by using a constructor in Java or C++.
3. Further invocations on methods in other objects may take place.

Exceptions - Programs can encounter many sorts of errors and unexpected conditions of varying seriousness. During the execution of a method, many different problems may be discovered. A block of code may define to throw an exception whenever a particular unexpected condition or errors arise. This means that control passes to another block of code that catches the exception. Control does not return to the place where the exception was thrown.

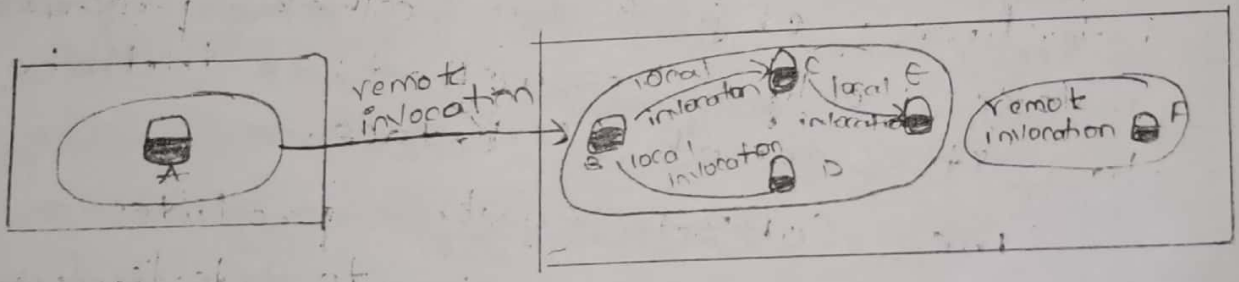
Garbage Collection - It is necessary to provide a means of freeing the space occupied by objects when they are no longer needed. A language, for example Java that can detect automatically when an object is no longer accessible recovers the space & makes it available

for allocation to other objects.
 This process is called garbage collection.

Distributed object model:-

Each process contains a collection of objects, some of which can receive both local and remote invocations, whereas the other objects can receive only local invocations. Method invocations between objects in different processes, where the same computer or not, are known as remote method invocations. Method invocations between objects in the same process are local method invocations.

Remote and local method invocations:-



All objects can receive local invocations, although they can receive them only for from other objects that hold references to them.

references:- The notion of object reference is extended to allow any object that can receive an RMI to have a remote object reference. A remote object reference is an identifier that can be used throughout a distributed system to refer to a particular unique remote object.

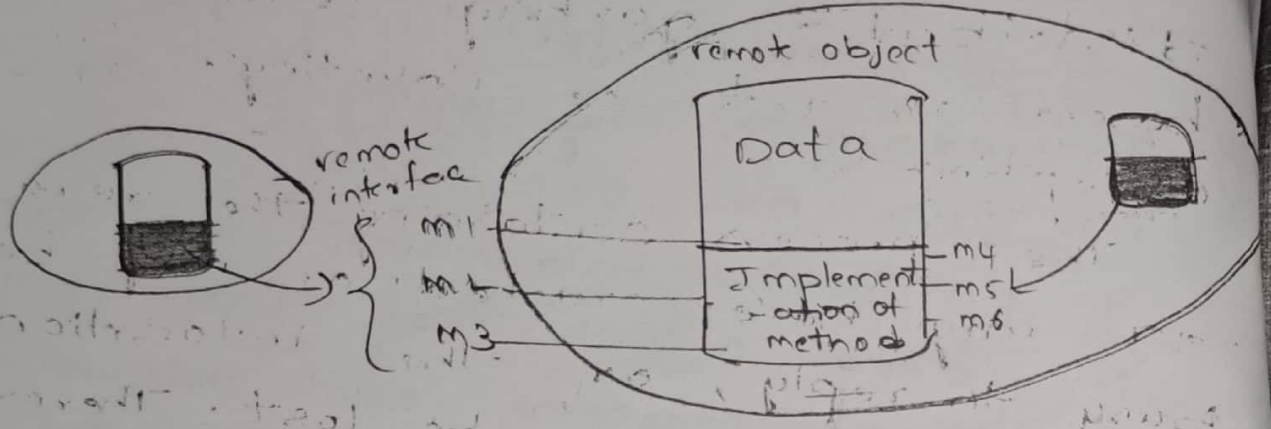
1. The remote object to receive a remote method invocation is specified by the invoker as a remote object reference.

2. remote object references may be passed as arguments and results of remote method invocations.

Interfaces:- The class of a remote object implements the method of its remote interface. Objects in other programs can invoke only the methods that belong to its remote interface.

The CORBA system provides an IDL, which is used to defining remote interfaces. For example of a remote interface defined in CORBA IDL. The classes of remote objects and the client programs may be implemented in any language such as C++, Java or Python for which an IDL compiler

is available. CORBA clients need not use the same language as the remote object in order to invoke its methods remotely. remote object and its remote interface.



Actions:-

As in the non-distributed case, an action is initiated by a method invocation, which may result in further invocations on methods in other objects. But in distributed case, the objects involved in a chain of related invocations may be located in different programs or different computers. When an invocation crosses the boundary of a program or computer, RMI is used and the remote reference of the object must be available to the invoking

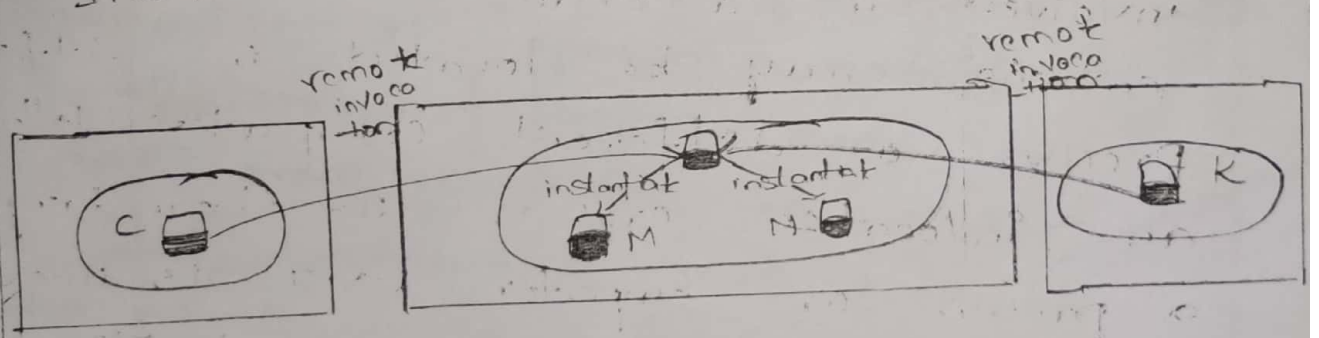
Garbage Collection:- In a language, for example Java, supports automatic garbage collection, then any associated RMI

System should allow garbage collection of remote objects. It is generally achieved by cooperation between the existing local garbage collector and an address module that carries out a form of distributed garbage collection, usually based on reference counting.

Exceptions :-

- The process containing the remote object may have crashed or may be too busy to reply, or the invocation or result message may be lost. Therefore, the rmi should be able to raise exceptions such as time outs that are due to distribution as well as those raised during the execution of the method invoked.

Instantiation of remote objects :-



Design Issues for RMI:

We have two design issues:-

1. The choice of invocation semantics. Although local invocations are executed exactly once, this cannot always be the case for remote method invocations.
2. The level of transparency that is desirable for RMI.

RMI Invocation Semantics:-

The main choices are:-

Retry request message:- Whether to retransmit the request message until either a reply is received or the server is assumed to have failed.

Duplicate filtering:- When retransmissions are used, whether to filter out duplicate requests at the server.

Retransmission of results:- Whether to keep a history of result messages to enable lost results to be retransmitted without re-executing the operations at the server.

→ local method invocations the semantics are exactly once, that means every method is executed exactly once. The choices of RMI invocation semantics are defined as follows:-

Maybe Invocation Semantics! - with maybe invocation semantics, the remote method may be executed once or not at all. may be semantics arises when none of the fault tolerance measure is applied.

This can suffer from the following types of failure:

- > Omission failures if the invocation or result message is lost;
- > Crash failures when the server containing the remote object fails

| Invocation Semantics - fault tolerance measures | | | Invocation Semantics |
|---|---------------------|----------------------|----------------------|
| retransmit request message | Duplicate filtering | Re-execute Procedure | |
| NO | Not applicable | Not applicable | maybe |
| Yes | NO | re-execute procedure | At-least-once |
| Yes | Yes | retransmit reply | At-most-once. |

If the invocation message was lost, then the method will not have been executed. on the other hand, the method may have been executed and the result message lost. -> crash

failure may occur, either before or after the method is executed. Moreover, in an asynchronous system, the result of executing the method may arrive after the timeout.

At-least-once invocation Semantics:-

With at-least-once invocation Semantics, the invoker receives either a result, in which case, the invoker knows that the method was executed at least once or an exception informing it that no result was received. At-least-once

invocation Semantics can be achieved by retransmission of request messages.

At-least-once invocation Semantics can suffer from the following types of failures

Crash failures when the server containing the remote object fails:

arbitrary failures. In cases when the invocation message is retransmitted, the remote object may receive it and execute the method more than once.

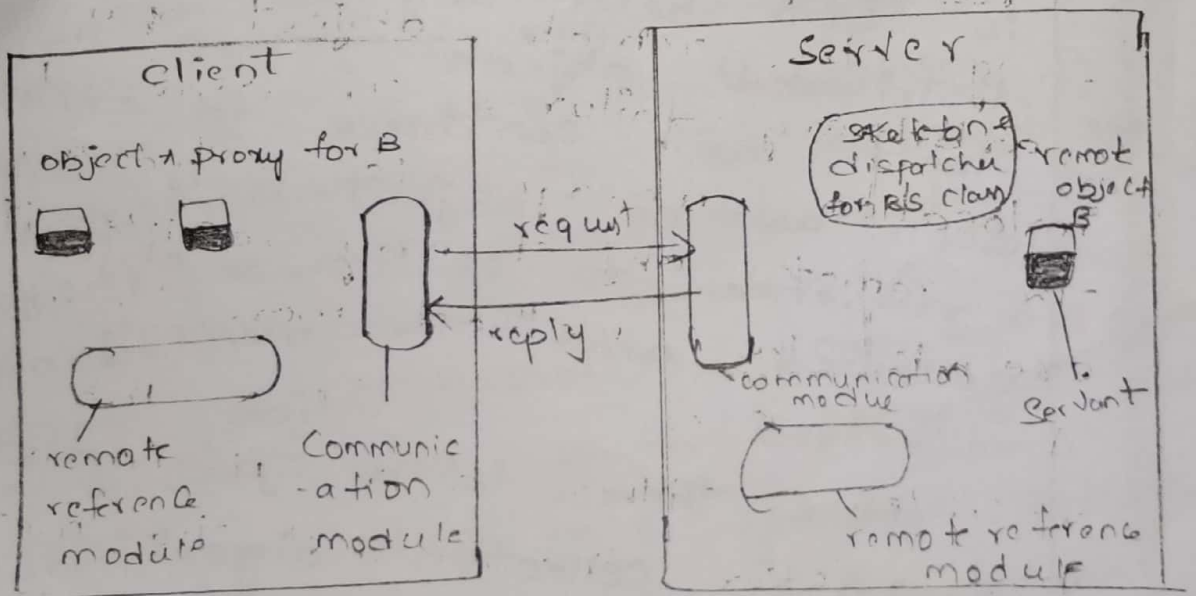
At-most-Once Semantics:- The invoker receives either a result in which case the invoker knows that the method was executed exactly once, or an exception informing it that no result was received, in which case the method will have been executed either

once or not at all.

Implementation of RMI

Several objects and modules are involved in achieving a remote method invocation in which an application-level object A invokes a method in remote application-level object B for which it holds a remote object reference. This section discusses the roles of each of the components, dealing first with the communication and remote reference modules and then with the RMI software that runs over them.

role of proxy and skeleton in RMI:



Communication module:- The two cooperating Communication modules, carry out the request reply protocol, which transmits request and reply messages between client and server. The communication module uses only the first three items which specify the message type, its requestId and the remote reference of the object to be invoked. The methodId and all the marshalling and unmarshalling is the concern of RMI software.

The communication module in the server is the dispatcher for the class of the object to be invoked, passing on its local reference, which it gets from the remote reference module in return for the remote object identifier in the request message.

Remote reference module:- It is responsible for translating between local and remote object references for creating remote object references. The remote reference module in each process has a remote object table. The table includes:-

- An entry for all remote objects held by the process
- An entry for each local proxy.

Servants:- A servant is an instance of a class which provides the body of

of remote object. It is the servant that eventually handles the remote request passed on by the corresponding skeleton.

→ The RMI Software:-

This consists of a layer of software between the application-level objects and the communication and remote reference modules. The roles of middle layer objects are as follows:-

Proxy:- The role of proxy is to make RMI transparent to clients by behaving like a local object to the invoker, but instead of executing an invocation forwards it in a message to a remote object. It hides the details of the remote object reference, the marshalling of arguments, unmarshalling of results and sending & receiving of messages from the client.

Dispatcher:- A server has one dispatcher and skeleton for each class representing a remote object. In our example the server has a dispatcher and skeleton for the class of remote object B.

The dispatcher receives the request message from the communication module. It uses the method id to select the appropriate method in the skeleton, passing on the request message. The

dispatcher and proxy use the allocation of method IDs to the same methods of the remote interface. The class of a remote object has a skeleton, which implements the methods in remote interface. A skeleton method unmarshals the arguments in the request message and invokes the corresponding method in the servant. It waits for the invocation to complete and then marshals the result together with any exceptions in a reply message to the sending proxy's method.

* Distributed Garbage Collection

The aim of a distributed garbage collector is to ensure that if a local or remote reference to an object is still held anywhere in a set of distributed objects, then the object itself will continue to exist, but as soon as no object any longer holds a reference to it, the object will be collected and the memory it uses recovered.

Here there is a Java distributed garbage collection algorithm, which is similar to the one described by Birrell et al. [1995]. It is based on reference counting. Whenever a remote

object reference enters a proxy, a proxy will be created and will stay there for as long as it is needed. The process when there is no longer a object lives should be informed of the new proxy at the client. Then later when there is no longer a proxy at the client, the server should be informed. The distributed garbage collector works in cooperation with the local garbage collector as follows:

* Each server process maintains a set of the names of the processes that hold remote object reference for each of its remote objects. For example, B holders is the set of client processes that have proxies for object B. This set can be held in an additional column in the remote object table.

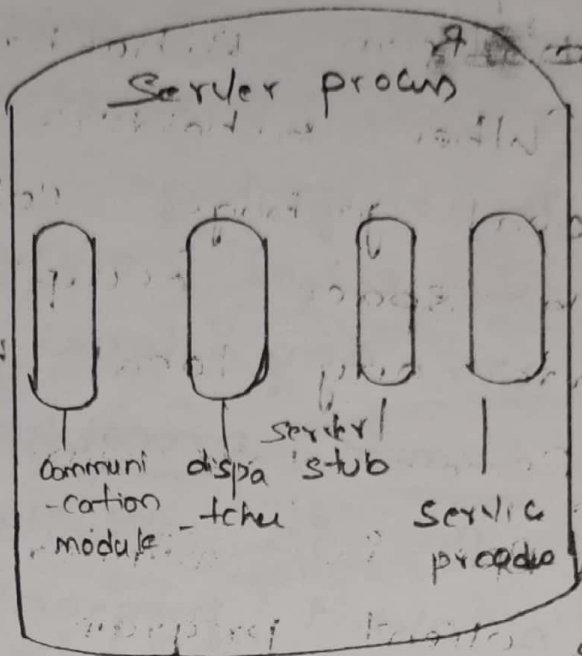
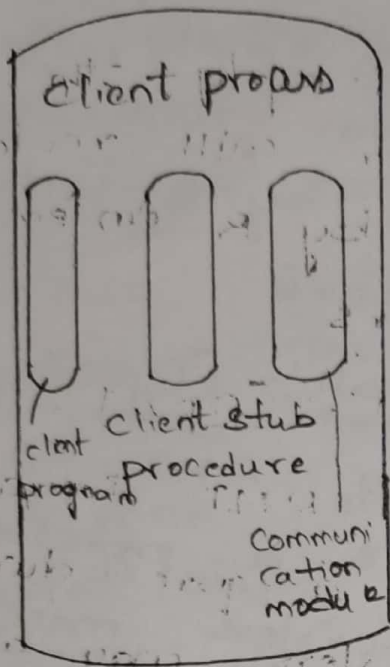
* When a client c first receives a remote reference to a particular remote object, B, it makes an address (B) invocation to the server of that remote object and then creates a proxy. The server adds c to B holders.

* When a client c's garbage collector notices that a proxy for remote object B is no longer reachable, it makes a remoteRef (B) invocation to the corresponding server and then

deletes the proxy, the server removes
c from B.holders.

* When B.holders is empty, the server's
local garbage collector will reclaim
the space occupied by B unless there
are any local holders.

* Remote procedure call :- [RPC] :-
Rpc is similar to RMI in which
a client program calls a procedure of
server program. The servers may act as
clients of other servers to allow chains
of rpc's. The procedures available for
calling remotely are defined by the server
program in its service interface. This type
of service is like a single remote
object containing state and methods. It
does not support remote object interferences
because it does not have the ability to
create new instances of objects.
like RMI, Rpc can also be implemented
to have one of the invocation
Semantics choices such as at-least-one
or at-most-once that are generally chosen.
It is implemented generally over a
request-reply protocol that the omission
of object references simplify from request
messages.



Sun RPC is an open networking Computer - ng Rfc, as described by the RFC 183 Srinivasan (1995). It was developed for client server communication in the Sun network File system. It is available with NFS installation and provides a part of the various Sun and other unix operating systems. There is liberty for implementors over using RPC over either UDP or TCP. When operating with UDP the length of request and reply messages is limited in length theoretically upto 64 kb, but practically to 8 or 9 kb. This Sun Rfc with UDP make use of at least once call Semant Interface Definition language! - XDR, the Sun interface language was developed basically for specifying representation of data and later it was extended to become an interface definition

language. The `sun.rpc` package defines services of `sun.rpc` by elucidating a set of procedure definitions together with supporting type definitions.

Input-output parameter: Only input-output parameter is permitted. Hence, procedures utilizing multiple parameters have to, then as part of single structure. Single result returns the output parameter.

Procedure Signature: It includes the result type, the name of the procedure and the input parameter type. The type of both input-output parameter may indicate either a single value or a structure consisting several values.

Interface Compiler: The `rp.gen` generates the following

1. Client stub procedure
2. The main server procedure, dispatcher and server stub procedure
3. For use by dispatcher and client and server stub procedures. XDR marshalling and unmarshalling procedures.

* Events & Notifications:-

The idea behind the use of events is that one object can react to a change occurring in another object. Notifications of events are essentially asynchronous and determined by their

receivers
events and notifications can be
used in a wide variety of different
applications, for example, to communicate
a shape added to a drawing, a
modification to a document, the fact
that a person has entered, or left a
room, or that a piece of equipment
or an electronically tagged book is at a
new location. Distributed event-based
systems have two main characteristics.

Heterogeneous: When event notifications
are used as a means of communication
between distributed objects, components
in a distributed system that
were not designed to interoperate, can
be made to work together. All that
is required is that event-generating
objects publish the types of events
they offer, and that other objects
subscribe to events and provide an
interface for receiving notifications.

Asynchronous: Notifications are sent
asynchronously by event-generating objects
to all the objects that have subscribed
to them to prevent publishers needing to
synchronize with subscribers - publishers
and subscribers need to be decoupled

Mushroom is a distributed, event-based system designed to support collaborative work, in which the user interface displays objects representing users and information objects such as document and notepad within shared workspaces called network places.

Participants in distributed event notification

The architecture is designed to decouple the publishers from the subscribers, allowing publishers to be developed independently of their subscribers, as far as possible limiting the work imposed on publishers by subscribers. The main component is an event service that maintains a database of published events and of subscribers interests.

The role of participating objects are as follows:-

The object of interest:- This is an object that experiences changes of state, as a result of its operation being invoked. Its changes of state might be of interest to other objects. This description allows for events such as a person wearing an active badge entering a room, in which case the room is the object of interest and operation consists of adding

information about the new person to its record of who is in the room. The object of interest is considered as part of event service if it transmits notifications.

Event: An event occurs at an object of interest as the result of the completion of a method execution.

Notification: It is an object that contains information about an event.

Subscriber: It is an object that has subscribed to some type of events in another object.

Publisher: It is an object that declares that it will generate notifications of particular type of event. A publisher may be an object of interest or an observer.

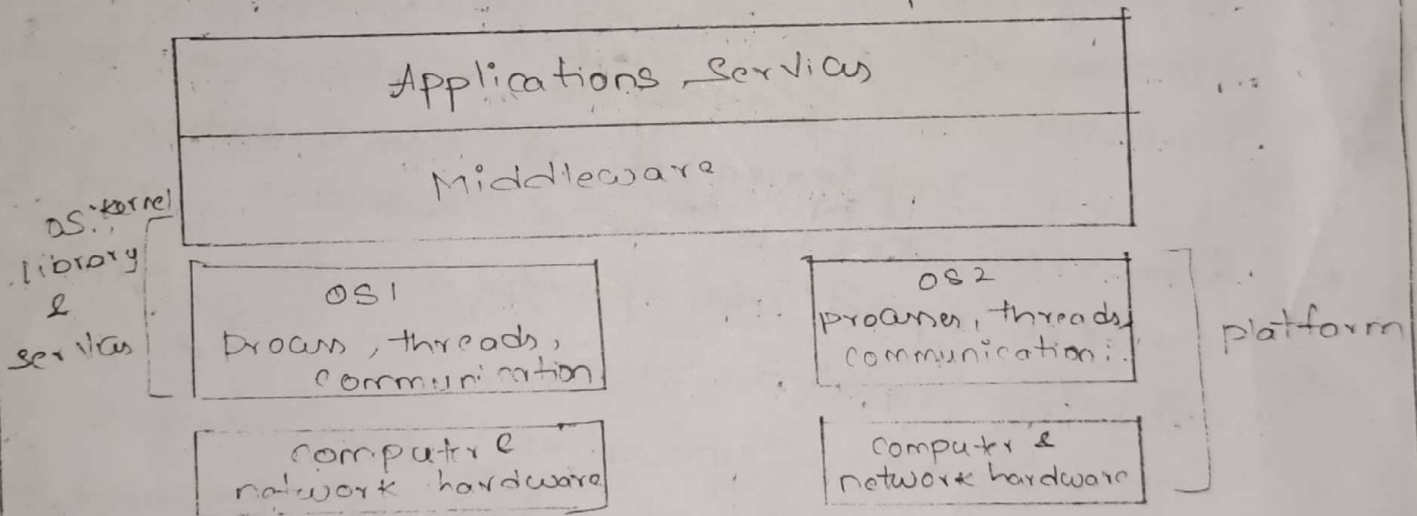
Observer objects: The main purpose of an observer is to decouple an object of interest from its subscribers. An object of interest can have many different subscribers with different interests.

Operating System Support

→ Operating System layer:-

Users will only be satisfied if their middleware - OS combination has good performance. Middleware runs on a variety of OS - hardware combinations at the nodes of a distributed system. The OS running at a node - a kernel and associated user-level services. Middleware utilizes a combination of these local resources, to implement its mechanism for remote interactions between objects or processes at the nodes.

System Layers



Encapsulation:- They should provide a useful service interface to their resources that is a set of operations that meet their clients needs. Details such as management of memory and devices used to

implement resources should be hidden from clients.

Protection:- Resources require protection from illegitimate access - for example, files are protected from being read by users without read permissions and device registers are protected from application processes.

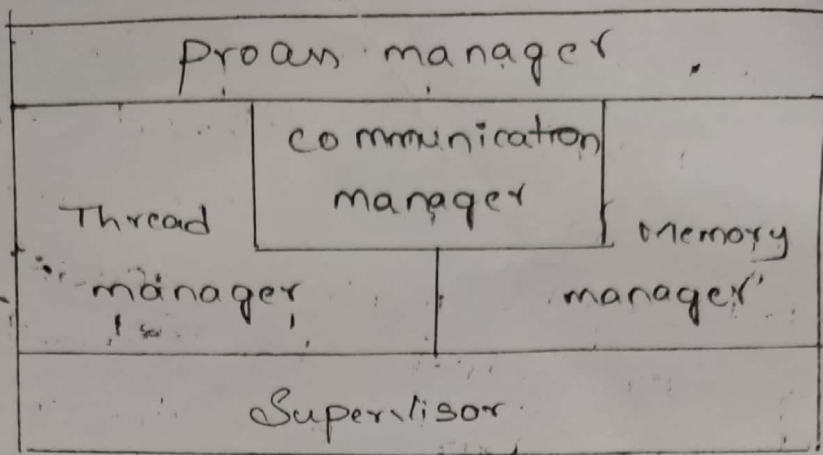
Concurrent processing:- clients may share resources and access them concurrently. Resource managers are responsible for achieving concurrent transparency.

clients access resources by making, for example rmi to a server object or system calls to a kernel. We call a means of accessing an encapsulated resource an invocation mechanism, however it is implemented. A combination of libraries, kernels and servers may be called upon to perform the following invocation-related tasks.

Communication:- operation parameters and results have to be passed to and from resource managers over a network or within a computer.

Scheduling:- When an operation is invoked, its processing must be scheduled within the kernel or server. ✓

Core OS functionality:-



The core OS components are as follows:

Process manager:- handles the creation of and operations upon processes. A process is a unit of resource management, including an address space and one or more threads.

Thread Manager:- Thread creation, synchronization and scheduling. threads are schedulable activities attached to processes.

Communication manager:- communication b/w threads attached to different processes on the same computer. Some kernels also support communication b/w threads in remote processes.

Memory manager:- management of physical and virtual memory.

Supervisor:- dispatching of interrupts, system call traps and other exceptions, control memory management unit and hardware caches.

Protection:

Resources require protection from illegitimate access. The threat to a system's integrity does not come only from malicious controlled code. Benign code that contains a bug or which has unanticipated behavior may cause part of the rest of the system to behave incorrectly.

We can also employ hardware support to protect modules from one another at the level of individual invocations, regardless of the language in which they are written. To operate the scheme on a general-purpose computer, we require a kernel.

Kernels and protection:

Kernel is a program that is distinguished by the facts that it always runs and its code is executed with complete access privileges for the physical resources on its host computer. It can control the memory management unit and set the processor registers so that no other code may access the machine's physical resources except in acceptable ways.

A kernel program executes with the processor in supervisor (privileged) mode. The kernel arranges that the other

Processes execute in user (unprivileged) mode. The kernel also sets up address spaces to protect itself and other processes from the actions of an aberrant process.

When a process executes application code it executes in a distinct user-level address space for that application. When the same process executes kernel code it executes in kernel's address space. The process can safely transfer from a user-level address space to the kernel's address space via an exception such as an interrupt or system-call trap. It is implemented by a machine-level TRAP instruction, which puts the processor into Supervisor mode and switches to the kernel address space.

* Processes and Threads:

Process consists of an execution environment together with one or more threads. A thread is the operating system abstraction of an activity. An execution environment is the unit of resource management: a collection of locally kernel-managed resources to which its threads have access. An execution environment primarily consists of:
→ an address space;

→ Thread Synchronization and Communication resources such as semaphores and communication interface

→ higher-level resources such as open files and windows

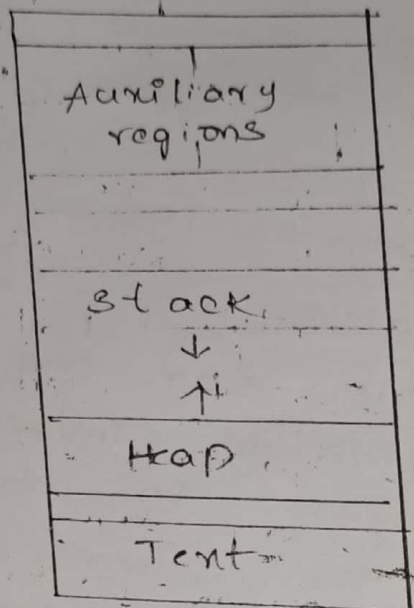
Execution environment are normally expensive to create and manage, but several threads can share them. Threads can be created and destroyed dynamically as needed.

An execution environment provides protection from threads outside it, but some kernels allow the controlled sharing of resources such as physical memory between execution environments residing at the same computer.

Address Spaces

It is a unit of management of a process's virtual memory. It is large (up to 32 bytes and sometimes up to 64 bytes) and consists of one or more regions separated by inaccessible areas of virtual memory. A region is an area of contiguous virtual memory that is accessible by the threads of the owning process. regions do not overlap. Each region is specified by the following properties :-

- its extent (lowest virtual address)
- read/write/execute permissions for program's threads
- whether it can be grown upwards/downwards



The model is page-oriented. There is need to support a separate stack for each thread. Allocating a separate stack region to each thread makes it possible to detect attempts to exceed the stack limits and to control each stack's growth. The unallocated virtual memory lies beyond each stack region and attempts to access this will cause an exception. The alternative is to allocate stacks for threads on the heap, but then it is difficult to detect when a thread has exceeded its stack limit.

Another motivation is to enable files in general - and not just the text

and data sections of binary files -
to be mapped into the address space.
A mapped file is one that is accessed
as an array of bytes in memory.
The virtual memory system ensures that
accesses made in memory are reflected in
the underlying file storage.

The need to share memory between
processes and the kernel is another factor
leading to extra regions in address
space. A shared memory region is
one that is backed by the same physi-
-cal memory as one or more regions
belonging to other address space.

Creation of a new process:-

The creation of a new process has
traditionally been an indivisible operation
provided by the operating system. The
design of the process creation mechanism
is to take account of the utilization
of multiple computers. Consequently, the process
support infrastructure is divided into
separate system services.

The creation of new process can be
separated into two independent aspects:

→ The choice of a target host. For
example, the host may be chosen from
among the nodes in a cluster of

Computers acting as a Computer Server.

→ The creation of an execution environment.

→ Choice of process host :-

The choice of node at which the new process will reside - the process allocation decision - is a matter of policy.

In general, process allocation policies range from always running new processes at their originator's workstation to sharing the processing load between a set of computers.

The transfer policy determines whether to situate a new process locally or remotely. This may depend.

The location policy determines which node should host a new process selected for transfer. This decision may depend on the relative loads of nodes, on their machine architectures.

Process location policies may be static or adaptive. There are two cases, in first case one load manager component and in the second there are several organized in tree structure. Load manager collect information about the nodes and use it to allocate new processes to nodes.

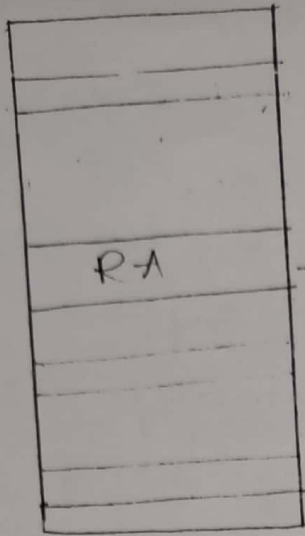
↳ Creation of a new execution environment
Once the host computer has been selected, a new program requires an execution environment consisting of an address space with initialized contents.

There are two approaches to defining and initializing the address space of newly created programs. The first approach is used where the address space is of statically defined format.

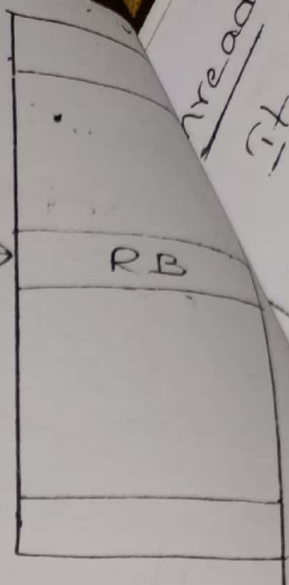
Alternatively, the address space can be defined with respect to an existing execution environment. For ex, the newly created child program physically shares the parent's text region and has heap and stack regions that are copies of the parent's in extent. This scheme has been generalized so that each region of parent program may be inherited by the child program. An inherited region may either be shared with or logically copied from the parent's region. When parent and child share a region, the page frames belonging to the parent's region are mapped simultaneously into the corresponding child region.

There are some techniques to determine this environment:-

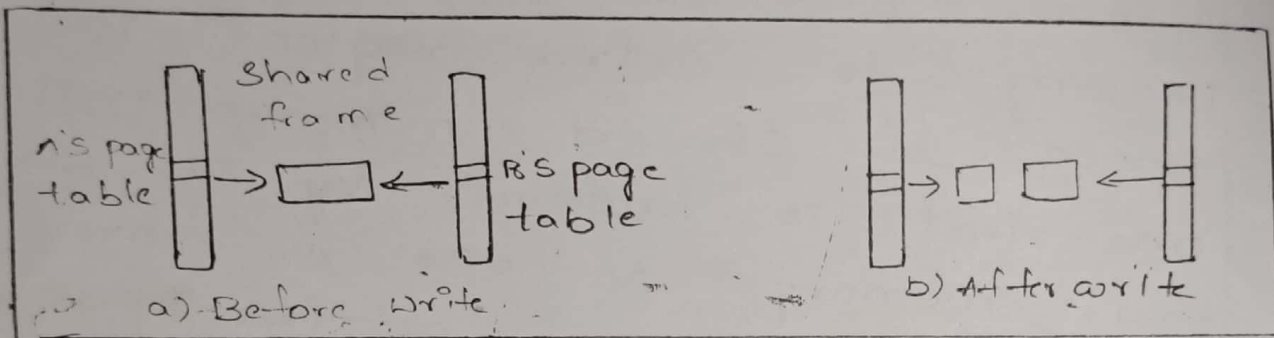
→ Copy-on-write.
content and



RB Copied from RA



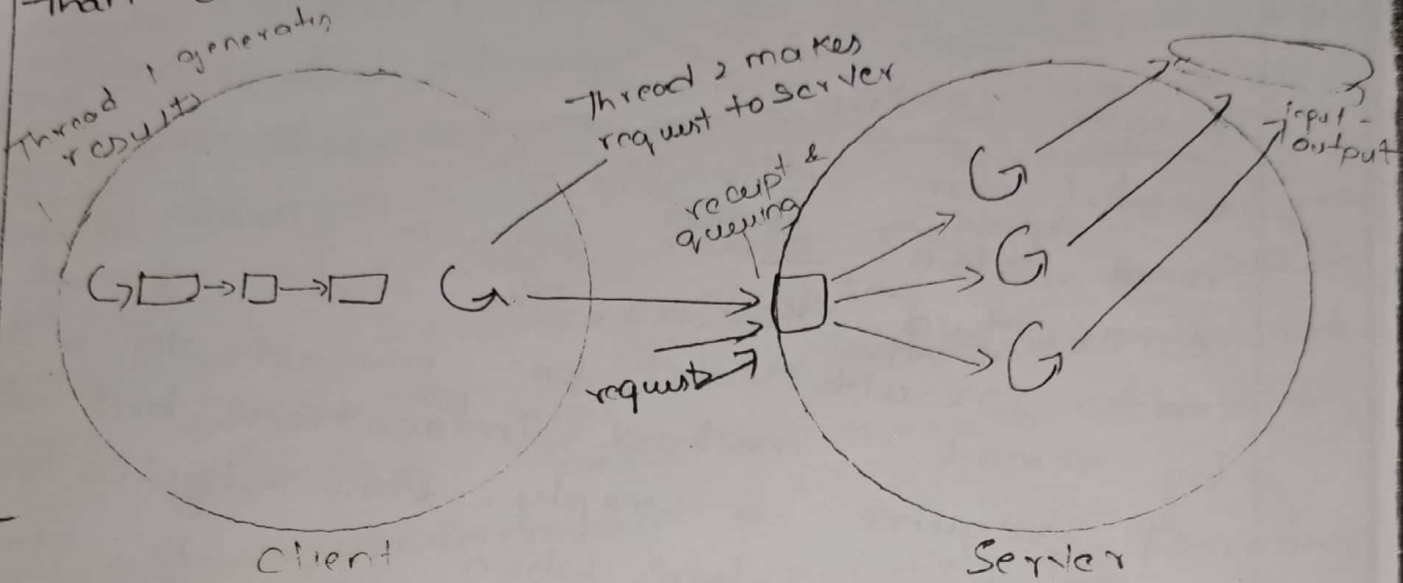
Kernel



Copy-on-write is a general technique for example, it is also used in copying large messages, so we take some time to explain its operation. For example of regions RA & RB whose memory is shared copy-on-write between two processes A and B.

Threads :-

It examines the advantage of enabling client and server processes to possess more than one thread.



The server has a pool of one or more threads, each of which repeatedly removes a request from a queue of received requests and processes it. We shall not concern ourselves for the moment with how the request and processes it. We do not concern for the moment with how the requests are received and queued up for the threads. Also, for the sake of simplicity, we assume that each thread applies the same procedure to process the requests.

consider the maximum server throughput measured in client requests handled per second, for different numbers of threads.

↓
If a single thread has to perform all processing, then the turnaround time for handling any request is on average $t + t = 10$ milliseconds. So, this server can handle 100 client request per second. Any new request messages that arrive while the server is handling a request are queued at the server port.

SET-3 Architectures for multi-threaded servers:-

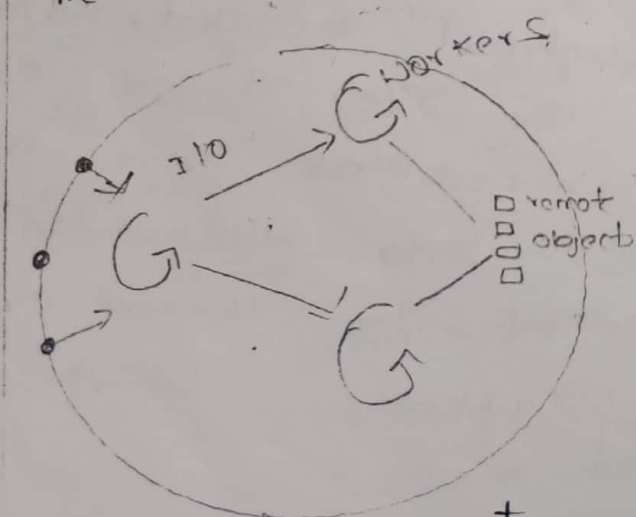
① To describe the various ways of mapping requests to threads within a server, we summarize Schmidt [1998], who describes the threads architecture of various implementations of CORBA object Request Broker (ORB).

Thread-per-request architecture-

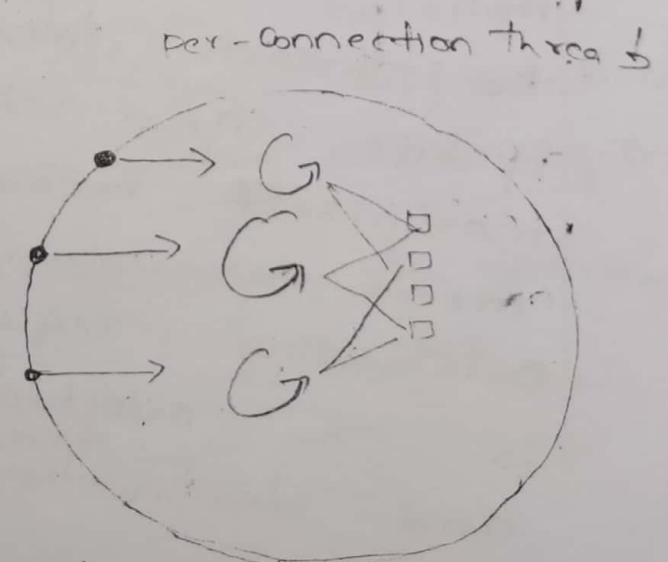
The I/O thread spawns a new worker thread for each request, and that worker destroys itself when it has processed the request against its designated remote object. This architecture has the advantage that the threads do not contend for a shared queue, and throughput is potentially maximized because the I/O thread can create as many workers as there are outstanding requests. Its disadvantage is the overhead of the thread creation and destruction operations.

Thread-per-Connection architecture It associates a thread with each connection. The server creates a new worker thread when a client makes a connection and destroys the thread when the client closes the connection. In between, the client may make many requests over the connection targeted at one or more remote objects. The thread-per-object architecture associates a thread with each remote object. An I/O thread receives request and queues them for the workers, but this time there is a per-object-queue.

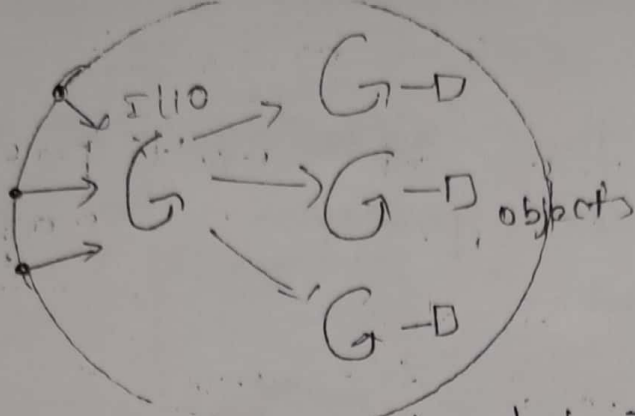
In each of these last two architectures the server benefits from lowered thread-management overheads compared with the thread-per-request architecture. Their disadvantage is that clients may be delayed while a worker thread has several outstanding requests but another thread has no work to perform. ..



Thread-per-request



Thread-per-connection



Threads within clients :- Threads can be useful for clients as well as servers and also ^{shows} server by a client process with two threads. The first thread generates results to be passed to a server by remote method invocation, but does not require a reply. RMI typically blocks the caller, even when there is strictly no need to wait. This client process can incorporate a second thread, which performs the RMI and blocks while the first thread is able to continue computing further results. The first thread places its results in buffers, which are emptied by the second thread. It is only blocked when all the buffers are full.

The case for multithreaded clients is also evident in the example of web browsers. Users experience substantial delays while pages are fetched. It is essential, therefore, for browsers to handle multiple concurrent requests for web pages.

Threads Versus Multiple Processes:-

In the case of multi processor, with other computation. The reader may have noted, however, that the same overlap could be achieved through the use of multiple single-threaded processes. Why we prefer multi-threaded processes is threads are cheaper to create and manage than processes and resource sharing can be achieved more efficiently between threads than between processes, because threads share an execution environment state associated with execution environments and threads:-

| Execution Environment | Thread |
|--|--|
| Address space tables Communication interfaces, open files Semaphores, other synchronization objects List of thread identifiers | Scaled processor registers priority and execution state (Such as Blocked) Software interrupt handling information execution environment identifier. |

The table shows that an execution environment and the threads belonging to it are both associated with pages belonging to the address space held in main memory and data and

1. Schedule another thread to run.

instructions held in hardware caches.

We can summarize a comparison of processes and threads as follows:

- * Creating a new thread within an existing process is cheaper than creating a process.
- * More importantly, switching to a different thread within the same process is cheaper than switching between threads belonging to different processes.
- * Threads within a process may share data and other resources conveniently and efficiently compared with separate processes.
- * But, by the same token, threads within a process are not protected from one another.

Threads programming:

Thread programming is concurrent programming, as traditionally studied in, for example, the field of operating systems. It refers to the concurrent programming concepts, which are explained by Bacon [2002]: race condition, critical section, monitor, condition, variable, semaphore. Much threads programming is done in a conventional language such as C, which has been augmented with a threads library.

Java Thread Constructor & management methods

Thread (ThreadGroup group, Runnable target, String name)
Creates a new thread in the SUSPENDED state, which will belong to group and be identified as name; the thread will execute the run() method of target.

setPriority(int newPriority), getPriority()
Set and return the thread's priority.

run()

A thread executes the run() method of its target object, if it has one, and otherwise its own run() method (Thread implements Runnable).

start()

change the state of the thread from SUSPENDED to RUNNABLE

sleep(int millisecs)

Cause the thread to enter the SUSPENDED state for the specified time

yield()

enter the READY state and invoke the scheduler

destroy()

Destroy the thread

Thread lifetimes:- A new thread on the same Java virtual machine starts its creator, in the SUSPENDED state. After it is made RUNNABLE with its start() method, it executes the run() method of an object designated by its constructor. The JVM and the threads on top of it all execute in a process on top of the underlying operating system. Threads can be assigned a priority, so that a Java implementation that supports priorities will run a portion with lower priority. A thread ends its life when it returns from the run() method or when its destroy() method is called.

Thread Synchronization:-

Java provides the synchronized keyword for programmers to designate the well-known monitor construct for thread coordination. Programmers designate either entire methods or arbitrary blocks of code as belonging to a monitor associated with an individual object. We could serialize the actions of I/O and worker threads by designating add() and removeFrom() methods in the queue class as synchronized methods.

machala Thread Synchronization calls.

thread, Join(int millisecs)
Blocks the calling thread for upto specified time until thread has terminated.

Thread, interrupt()
interrupts thread: Causes it to return from a blocking method call such as sleep()

Object, wait(long millisecs, int nanosecs)
Blocks the calling thread until a call made to notify() or notifyAll() on object wakes the thread, or the thread is interrupted or the specified time has elapsed.

Object, notify(), Object, notifyAll()
wakes, respectively, one or all of any threads that have called wait() on object.

Threads Scheduling:

It is between preemptive and non-preemptive scheduling of threads. In preemptive scheduling, a thread may be suspended at any point to make way for another thread. In non-preemptive scheduling, a thread runs until it makes a call to the threading system when the system may de-schedule.

it and schedule another thread to run.

The advantage of non-preemptible scheduling is that any section of code that does not contain a call to the threading system is automatically a critical section. On the other hand, non-preemptively scheduled threads cannot take advantage of a multiprocessor, since they run exclusively. Care must be taken over long running sections of code that do not contain calls to the threading system.

DISTRIBUTED FILE SYSTEMS

Introduction:

File Systems were originally developed for centralized computer systems and desktop computers as an OS facility providing a convenient programming interface to disk storage. They subsequently acquired features such as access control and file-locking mechanisms that made them useful for the sharing of data and programs. Distributed file systems support the sharing of information in the form of files and hardware resources in the form of persistent storage throughout an intranet. A well-designed file service provides access to files stored at a server with performance and reliability similar to, and in some cases better than, files stored on local disks. Their design is adapted to the performance and reliability similar to, and is characteristic of local networks and hence they are most effective in providing shared persistent storage for use in intranets.

File System modules:-

- Directory module relates file names to file IDs
- File module relates file IDs to particular files
- Access control module checks permission for operation requested.

file access module reads or writes file attributes.

Block module accesses and allocates blocks

Device module disk i/o and buffer

Characteristics of file systems:

File systems are responsible for the organization, storage, retrieval, naming, sharing and protection of files. They provide a programming interface that characterizes the file abstraction, freeing programmers from concern with the details of storage allocation and layout. Files are stored on disks or other non-volatile storage media.

Files contain both data and attributes. The data consist of sequence of data items (8-bit bytes), accessible by operations to read and write any portion of the sequence. The attributes are held as a single record contains information such as length of the file, timestamps, file type, owner's identify and access control lists. File systems are designed to store and manage large numbers of files, with facilities for creating, naming and deleting files.

file attribute record structure.

file length

creation timestamp

Read timestamp

write timestamp

Attribute timestamp

Reference count

owner

file type

Access Control list

26 File System operations!

The main operations on files that are available to applications in UNIX systems. These are the system calls implemented by the kernel; application programmers usually access them through library procedures.

The file system is responsible for applying access control for files. In local file systems such as UNIX, it does so when each file is opened checking the rights allowed for the user's identity in the access control list against the mode of access requested in the open system call.

UNIX file System operations:

$filedes = open(name, mode)$

$filedes = \overset{create}{open}(name, mode)$

$status = close(filedes)$

$count = read(filedes, buffer, n)$

$count = write(filedes, buffer, n)$

$pos = lseek(filedes, offset, whence)$

$status = unlink(name)$

$status = link(name1, name2)$

$status = stat(name, buffer)$

opens an existing file with given name

^{creates}
~~open~~ a new file

closes the open file

transfers n bytes from the file referenced by $filedes$ to $buffer$.

transfers n bytes to the file referenced by $filedes$ from $buffer$.

moves the read-write pointer to $offset$.

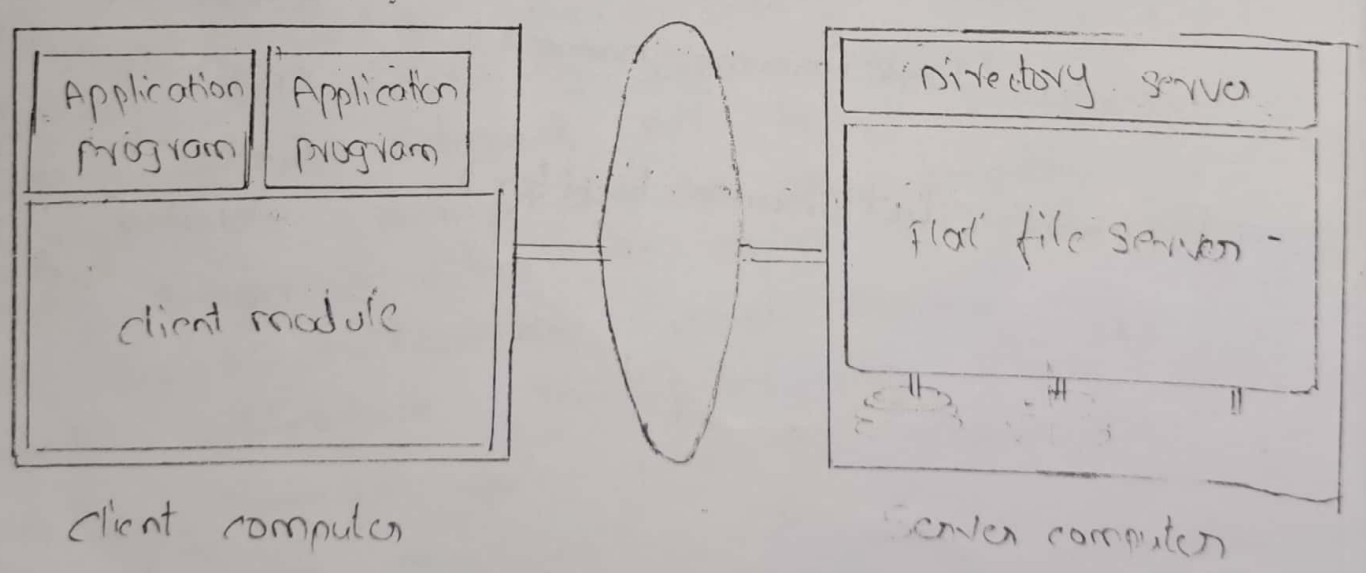
removes the file name from directory

adds a new name ($name2$) for file ($name1$)

gets the file attributes for file name into $buffer$.

File Service Architecture:-

An architecture that offers a clear separation of main concerns in providing access to files is obtained by structuring the file service as three components: a flat file service, a directory service, and a client module. A flat file service and directory service each export an interface for use by client programs and their RPC interfaces, taken together, provide a comprehensive set of operations for access to files. The client modules provide single programming interface with operations on files.



Flat file service: The flat file concerned with implementing operations on the contents of files. Unique file identifiers (UFIDs) are used to refer to files. All requests for flat file service operations. The division of responsibilities between the file service and directory service based upon the use of UFIDs. UFIDs are long sequences of bits chosen such that each file has UFID that is unique among all of these files in distributed system.

2b Flat file service operations

| | |
|---|---|
| $\text{Read}(\text{Field}, i, n) \rightarrow \text{Data}$ -throws BadPosition | $\text{If } 1 \leq i \leq \text{length}(\text{File})$: Read sequence of upto n items from a file starting at item i & returns it in Data |
| $\text{write}(\text{Field}, i, \text{Data})$ -throws BadPosition | $\text{If } 1 \leq i \leq \text{length}(\text{File}) + 1$: write sequence of Data to a file |
| $\text{create}() \rightarrow \text{Field}$ | creates a new file of length 0 & defines a UFID |
| $\text{Delete}(\text{Field})$ | removes file from file store |
| $\text{GetAttributes}(\text{Field}) \rightarrow \text{attr}$ | returns the file attributes for the file |
| $\text{SetAttributes}(\text{Field}, \text{attr})$ | sets the file attributes |

File
Per
File

Directory Service:- The directory service provides a mapping between text names for files and their UIDs. Clients may obtain the UID of a file by quoting its text name to the directory service. The directory service provides the functions needed to generate directories, to add new file names to directories and to obtain UIDs from directories. It is a client of the flat file service; its directory files are stored in files of flat file service. When a hierarchic file naming scheme is adopted, as in UNIX directories hold references to other-

Directory Service operations:-

lookup (Dir, Name) -> Field
- throws NotFound

locates the text name in directory & returns the relevant UID. if name not found it throws an exception.

AddName (Dir, Name, Field)
- throws NameDuplicate

If Name is not in the directory, add (Name, File) to the directory & update's the file's attribute record

UnName (Dir, Name)
- throws NotFound

if Name is in the directory: the entry containing Name removed from the directory.

GetName (Dir, pattern) -> NameSeq

if name returns all text names in directory

Client module :- A client module runs in each client computer, integrating and extending the operations of the flat file service and the directory service under a single application programming interface that is available to user-level programs in client computers. *

FSA :- It is an abstract architectural model that underpins both NFS and AFS. It is based upon a division of responsibilities between three modules. The architecture is designed to enable a stateless implementation of server module.

SUN NFS :-

SUN Microsystems's Network File System has been widely adopted in industry and in academic environment since its introduction in 1985. The design and development of NFS were undertaken by staff at SUN Microsystems in 1984. Although several distributed file services had already been developed and used in universities and research laboratories, NFS was the first file service that was designed as a product. The design and implementation of NFS have achieved success both technically and commercially.

To encourage its adoption as a standard, the definitions of the key

interfaces were placed in the public domain, enabling other vendors to produce implementations, and the source code for a reference implementation was made available to other computer vendors under license. It is now supported by many vendors, and the NFS protocol is an Internet Standard, defined in RFC 1813.

NFS provides transparent access to remote files for client programs running on UNIX and other systems. The client-server relationship is symmetrical, each computer in an NFS network can act as both client and server.

An important goal of NFS is to achieve a high level of support for hardware and OS heterogeneity. The design is OS independent. Implementation of NFS on high-performance multiprocessor hosts have been developed by several vendors and these are widely used to meet storage requirements in intranets with many concurrent users.

Andrew File System:-

Andrew is a distributed computing environment developed at Carnegie Mellon University for use as a campus computing and information system. The design of AFS reflects an intention

to support information sharing
large scale by minimizing client-
communication. This was achieved
transferring whole files between server
and client computers and eaching to
at clients until the server receives
a more up-to-date version. ∞.

* \$

* peer-to-peer Systems: introduction:-

Peer-to-peer systems aim to support useful distributed services and applications using data and computing resources available in the personal computers and workstations that are present on the Internet and other networks in ever-increasing numbers. This is increasingly attractive as the performance difference between desktop and server machines narrows, and broadband network connections proliferate. Peer-to-peer applications enhance their scalability, reliability, and security. These systems provide access to information resources located on computers throughout a network whether it be the Internet or a corporate network. Algorithms for the placement and subsequent retrieval of information objects are a key aspect of system design. Their design aims to deliver a service that is fully decentralized and self-organizing, dynamically balancing the storage.

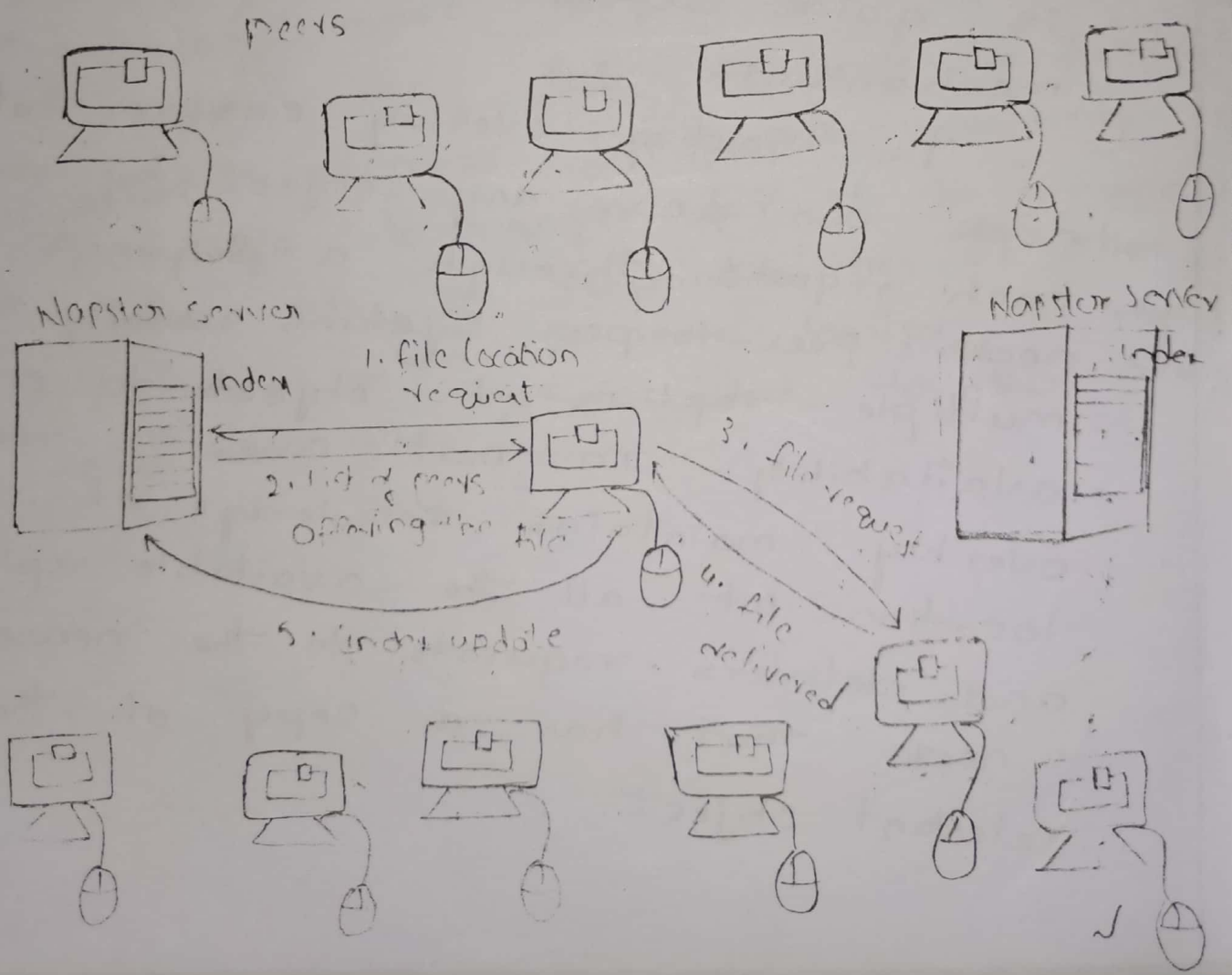
* Characteristics:-

- * Their design ensures that each other contributes resources to the system.
- * Although they may differ in the resources that they contribute, all the nodes in a peer-to-peer system have the same function.

→ Napster took advantage of the special characteristics of the application for which it was designed in other ways

→ Music files are never updated, avoiding any need to make all the replicas of files consistent after updates

→ No guarantees are required concerning the availability of individual files - if a music file is temporarily unavailable it can be downloaded later. This reduces the requirement for dependability of individual computers and their connections to the internet



* peer-to-peer middleware :-

peer-to-peer middleware systems are designed specifically to meet the need for the automatic placement and subsequent location of the distributed objects managed by peer-to-peer systems and applications. Functional requirements:- The function of peer-to-peer middleware is to simplify the construction of services that are implemented across many hosts in a widely distributed network. To achieve this it must enable clients to locate and communicate with any individual resource made available to a service, even though the resources are widely distributed amongst the hosts.

Non-functional requirements:- To perform effectively, peer-to-peer middleware must also address the following.

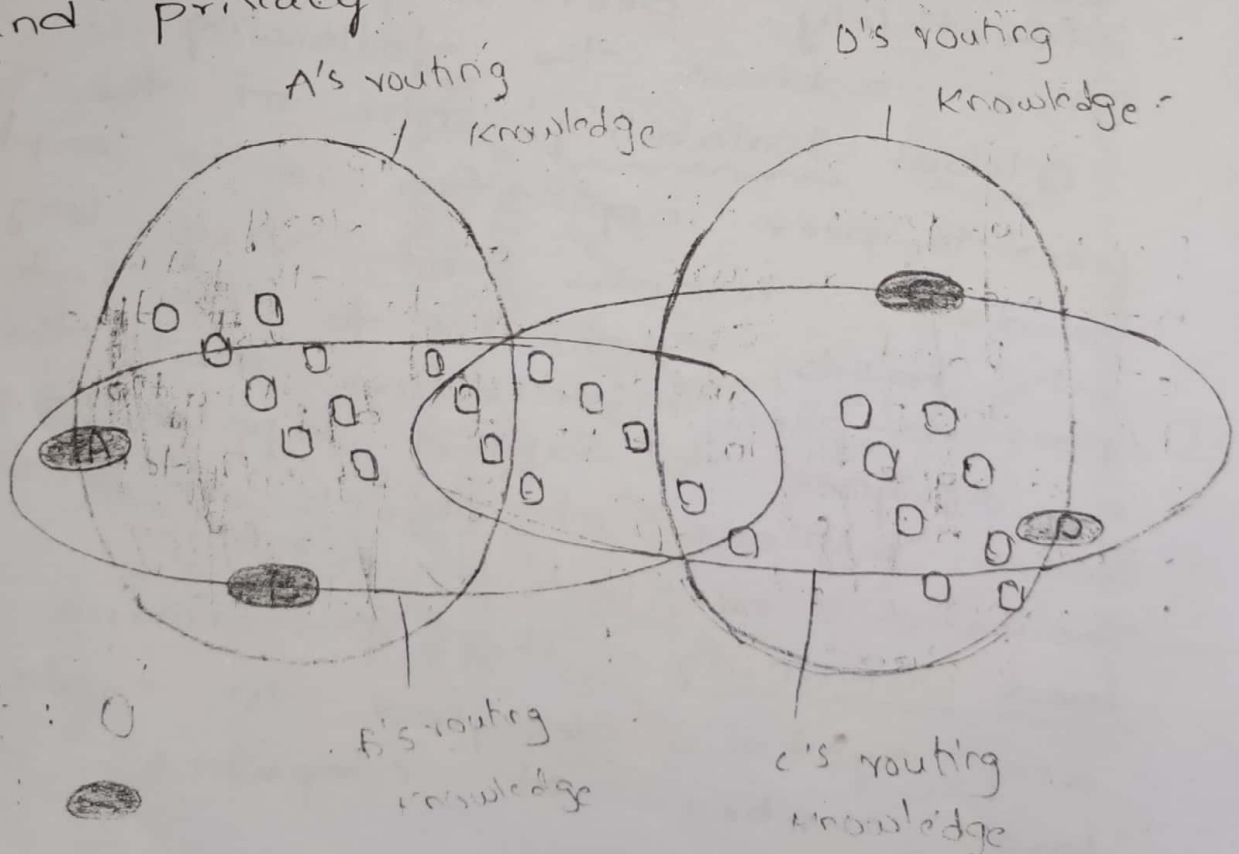
Global scalability:- One of the aims of peer-to-peer application is to exploit the hardware resources of very large number of hosts connected to the Internet. Peer-to-peer middleware must therefore be designed to support applications that access millions of objects on tens of thousands of hosts.

load balancing:- The performance of any system designed to exploit a large number of computers depends

upon the balanced distribution of workload across them. For the systems we are considering, this will be achieved by a random placement of resources together with the use of replicas of heavily-used resources.

Optimization for local interaction between neighbouring peers: The network distance between nodes that interact has a substantial impact on the latency of individual interactions, such as client requests for access to resources.

Security systems diverse with participating hosts, or ownership, trust must be built up by the use of authentication and encryption mechanisms to ensure the integrity and privacy of information.



* 2¹ routing overlays:-

A distributed algorithm of the type known as a routing overlay takes responsibility for locating nodes and objects. The name denotes the fact that middle-ware takes the form of a layer that is responsible for routing requests from any client to a host that holds the object to which the request is addressed. The objects of interest may be placed and subsequently relocated to any node in the network without client involvement. It is termed an overlay since it implements a routing mechanism in the application layer that is quite separate from any other routing mechanisms.

The routing overlay ensures that any node can access any object by routing each request through a sequence of nodes. Peer-to-peer systems usually store multiple replicas of objects to ensure availability. In that case, the routing overlay maintains knowledge of the location of all the available replicas and delivers requests to the nearest live node that has a copy of the relevant object.

3. The main task of a routing overlay is the following:-

1. A client wishing to invoke an operation on an object submits a request including the object's GUID to the routing overlay which routes the request to a node at which a replica of the object resides.

2. A node wishing to make a new object available to a peer-to-peer service computes a GUID for the object and announces it to the routing overlay, which then ensures that the object is reachable by all other clients.

3. When clients request the removal of objects from the service the routing overlay must make them unavailable.

Nodes may join & leave the service. When a node joins, the routing overlay

arranges for it to assume some of the responsibilities of other

nodes. When a node leaves, the responsibilities are distributed to other nodes.

COORDINATION & AGREEMENT

Introduction:

failure assumptions:

- A various failure assumptions include:
- The connection between each pair of processes is made over reliable channels because these channels mark the failures of network components.
 - The processes do not depend on each other for forwarding the messages.
 - Hardware redundancy exists at certain points in the synchronous system which results in delivering the message on time without failure.
 - It is also assumed that process can fail only by crashing.

Failure Detectors:

It is a fundamental abstraction in distributed computing. Failure detectors are used as building blocks to simplify the design of reliable distributed algorithms. In particular, we describe how failure detectors can factor out timing assumptions to detect failures in distributed algorithms.

Failure detectors can be defined as service which can be used to find whether the process has failed. For detecting failures locally, each computer

runs algorithms that detect failures of each local process which is called as local failure detector.

The concept of unreliable failure detectors characterize two properties completeness and accuracy. It results Suspected and unsuspected. The term Suspected suggests that process has failed because it has been long time that it has received any message in this method, suspicion may be wrong because the process may be running slowly because of some reason. The term unsuspected suggests that the process has not failed due to any reason and it is working correctly.

A reliable failure detector is accurate because it detects whether the process has failed or crashed, or it is getting executed correctly very accurately. If a process is failed it queried & if it does not receive any message for long time than it is declared as failed, or crashed.

unreliable failure detector can be implemented using algorithm. In this algorithm every process sends message. 'p is here' and maintain time record and if the process is taking more time to transmit

a message than usual and a
get message with $T+D$ seconds
the last message which it received
A report is sent to q by p that
p is suspicious or else (p is here)
then it reports to q that p is (fine)

* Distributed mutual exclusion:-

If a collection of processes shared
collection of resources then often
mutual exclusion is required to prevent
interference and ensure consistency when
accessing the resources. This is the critical
section problem, familiar in operating
system. In some cases shared resource
are managed by servers that also
provide mechanisms for mutual exclusion

Algorithms for mutual exclusion:-

The performance of mutual exclusion
algorithm can be calculated with
following measures:-

1. Bandwidth:- The no. of messages transmi-
tted in each entry and exit operati-
on is directly proportion to the band-
width utilized

2. client Delay:- The delay that occur
due to processes at each entry and
exit operation of critical section.

System throughput! - The number of critical section request fulfilled per unit time

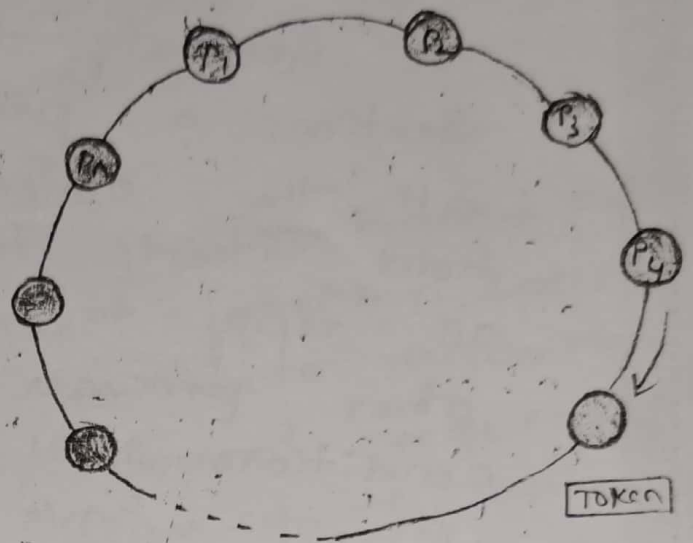
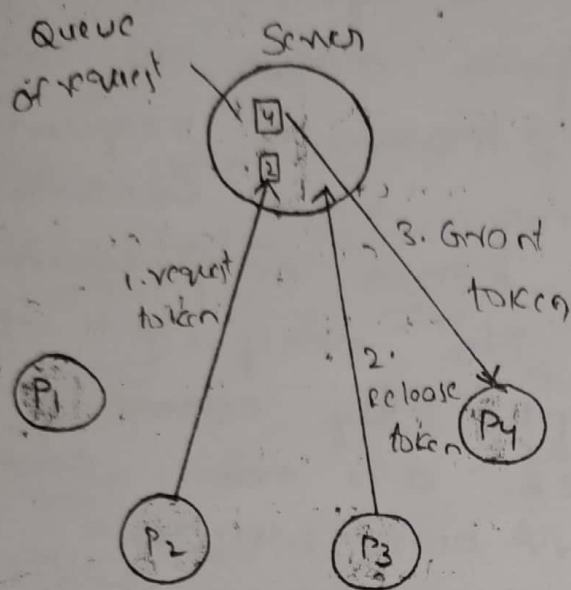
4. Synchronization delay! - The time elapsed between when some process leaves its critical section. System's throughput increases as their is decrease in Synchronization delay

Client algorithms

(1) Central server algorithm!

It can be seen as a token based algorithm. In this, the system maintains a token and makes sure that one process at a time gets that token. The process of granting a token is done by a server that grants the permission to enter the critical section on the request of client. If a process does not have a token then it cannot enter the critical section. In order to enter the critical section, a process requests the central server by sending a request message and if no other process is in critical section, then the central server immediately replies by granting the tokens. If other process is in critical section then the server adds the requests in the queue list and waits till

the process exists, as the process exists it asks the server to return its token. Then the server grants the token to the next process.



Server managing a mutual-exclusion token for a set of processes

A ring of processes transferring a mutual exclusion token

The process P3 requests for token which is added to the queue, that already consists of P2 requests on the other hand where P1 exists the critical section, the server grants the permission to P2, process P4 is constant as it doesn't need entry in the critical section. Although this algorithm is simple to implement but it can become a bottleneck. In order to overcome this problem, shared data can be distributed among several servers.

Ricart-Agrawala Algorithm:-

Ricart and Agrawala developed a distributed algorithm using multicast and logical clocks in the year 1981.

Generally to gain access to critical section, a process makes a request to all the other processes and considers to hold itself from access only when they all reply to say its okay. All the other processes reply by ensuring that conditions $M1 \leq M2$ are met. Each process $P1, P2$ have their own unique identifiers through which they send request. Use Lamport clocks and their IDs. The main steps involved in this algorithm is

On Initialization

state := RELEASED;

To enter the section

state := WANTED

multicast request to all processes;

T = request's timestamp;

wait until (number of replies received = $(N-1)$);

state := HELD;

on receipt of request $\langle T_i, P_i \rangle$ at P_j ($i \neq j$)

if (state = HELD or (state = WANTED and $(T, P) < (T_i, P_i)$))

then

queue request from P_i without replying;

else

reply immediately to P_i ;

end it

To exit critical section.

state = RELEASED;

reply to any queued requests.

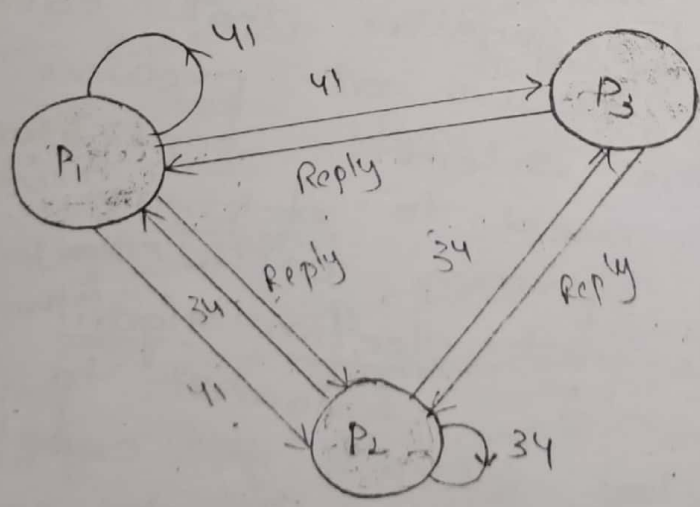
All the processes send request to enter critical section including the timestamp and identity P_i . The process record their status when they are out of critical section as RELEASED and when they want to enter critical section as WANTED, & the process which is in critical section will record its state as HELD.

Whenever a process request entry is the critical section and all the processes reply at once then the process can gain entry at that time only. If processes do not reply to the request because process is already in the critical section then the process who requested entry in the critical section has to wait till the other process exists the critical section.

When 2 or more processes request entry in critical section it will not accept from other processes requesting, until it has recorded its timestamp.

This multicast and logical locks algorithm meets the conditions of, ME1, ME2, ME3.

Consider an example of 3
 processes P_1, P_2, P_3 . Whereas P_1 & P_2 require
 to enter critical section simultaneously
 and P_3 is not willing to enter, P_1
 recorded its timestamp as 19 and P_2
 recorded its timestamp as 23. When P_3
 receives requests from P_1 and P_2 , it
 replies to them immediately whereas when
 P_1 receives P_2 's request it finds that
 its own recorded timestamp is lower
 than P_2 and it does not reply. When
 P_2 receives P_1 's request it finds that
 P_1 's timestamp is lower compared to
 its own and replies thus P_1 gains
 entry in critical section. When P_1
 exits it grants request to P_2 to
 enter.



P_1 & P_2 request entering the critical
 section simultaneously.

(1) This algorithm is expensive because it uses more bandwidth. Performance can be improved by deciding that the process which enters first in C/S must again send request to all the processes again.

Maekawa's Voting algorithm:-

It is the first quorum-based mutual exclusion algorithm. The algorithm specifies that for a process to have access to the critical section, does not require its peer to grant access to the critical section. Therefore, a single process can get access to the critical section. Here the permission is granted to the process to enter subsets of their respective peers and the permission is given until the subsets of any two processes intersect. The permission is use of voting concept, i.e. each process votes for one another for gaining access to the critical section. An adequate amount of votes must be collected by a candidate process so as to enter the critical section. When the processes under the intersection about two sets of votes occurs, a safety property met is established.

The algorithm defines the voting set $V S_i$ with each process P_i ($i = 1, 2, 3, 4, \dots, N$) & $V S_i \subseteq \{P_1, P_2, \dots, P_N\}$. The voting sets $V S_i$ are chosen such that $\forall i, j = 1, 2, 3, 4, \dots, N$ is,

1. $P_i \in V S_i$

2. There exist at most one member common to any two voting sets.

$$V S_i \cap V S_j \neq \emptyset$$

3. Each and every process must hold a voting set of equal size

$$|V S_i| = Q$$

4. Each and every process P_i , belongs to m of voting sets $V S_i$.

The main steps involved in algorithm is

Initialization

stat = RELEASED;

voted = FALSE;

process P_i request entering the C.S

stat = WANTED;

making a multicast request to all processes present in $V S_i$;

wait until (No. of replies received is = Q);

stat = HELD;

upon receiving the request from P_i at P_j

if (stat = HELD) OR voted = TRUE

Then queue the request from P_i

else send a reply to P_i

voted = TRUE;

end if.

in order for P_i to exit c.s

stat = RELEASED;

if (queue of requests is non-empty) then

Remove the head of queue - from P_k , say;

Send reply to P_k ;

voted = TRUE;

else voted = FALSE;

end if

* Elections:-

* ELECTIONS:-

Election algorithm:- It is algorithm used for solving the problem of choosing an election coordinator.

Goal:- It ensure that each time an election occurs, all processes must agree upon the selected coordinator. Moreover, it is used to determine the non-crashed process with the highest ID.

Features of election algorithm:- election algorithm is designed to choose a coordinator which is capable of choosing a process that can play the role of server. When the process which is playing the role of server retires

Another election is needed to select the new Server or new Coordinator for.

Any process calls the election to choose the new coordinator or server, a single process cannot call the election more than once at a time. Whereas, in principle the N processes can call N concurrent elections.

A prominent requirement is elected process should be unique. The requirements for any specific run of the algorithm are as follows,

- (i) e_1 (safety): All processes agree on the same elected process with the largest identifier. At the end of a run and the elected one should be a process that did not crash yet.
- (ii) e_2 (live): All processes participated or crash during the process. The performance of an election algorithm can be measured by the total bandwidth utilization and turn around time.

Bully algorithm

It is a method in distributed computing for dynamically selecting a coordinator by process ID number.

This algorithm was devised by Garcia-Molina in 1982. Some of the assumptions of this algorithm are,

i. Each process knows the IP and address of every other process.
e. The system is synchronous and users timeout for identifying process failure/crash.

There are three types of messages in this algorithm. They are as follows,

- (a) Election Message:- It is sent to announce the election.
- (b) Answer Message:- This responds to the election message
- (c) Coordinator message:- It is sent to announce the identity of the elected process

When a process notices that, the coordinator is no longer responding to request because of message timeout or failure of the coordinator, it initiates the election.

process P holds an election as follows,
(i) P sends an election message to all other processes with higher process IDs.

(ii) if no one responds with a higher process ID that it, wins the election and become coordinator.

(iii) if one of the process with a higher ID responds, P waits a certain amount of time for that process to broadcast itself as the leader. if it does not receive this message in time, it rebroadcasts the election message.

(iv) if it gets an election message from

→ Multicast Communication is of various types that are as follows:-

Reliable Multicast:-

It should satisfy the following properties:-

1. Integrity: A correct process P must deliver a message at most once, P should belong to group (m) and m was sent by a multicast operation by sender (m) .
2. Validity: If a correct process multicasts the message m then it will eventually transport m .
3. Agreement: If a correct process transports message m , then all other correct processes in group (m) will eventually transport m .

Implementation of Reliable Multicast:-

The two implementations are:-

1. B-multicast
2. IP multicast

B-multicast: It can be explained with the help of following algorithm, that contains the primitives R-multicast and R-deliver.

On initialization

Received = { };

for process P to R-multicast message m to group g B-multicast (g, m) ; // $P \in g$ is entered as destination on B-deliver(m)

at process q with $g = \text{group}(m)$ if C_m Then from SR

Then

Required: = Received $\forall \{m\}$;
if $(q \neq p)$ then B-multicast (g, m) ; end if
R-deliver m ;
end if.

IP multicast: IP multicast is a method of sending Internet Protocol (IP) datagrams to a group of interested receivers in a single transmission. It is a form of point-to-multipoint communication employed for streaming media & other applications.

FIFO Ordering:

When a true process issues multicast (g, m) and multicast (g, m') then each true process delivering m' will deliver m before m' .

Implementation of FIFO ordering: In FIFO implementation the multicast we achieve using sequence number is similar to what we achieve in one-to-one communication.

In this implementation consider only non-overlapping groups. The FIFO ordering is guaranteed by the reliable multicast protocol that is defined on top of IP multicast. A FIFO-ordered multicast can be constructed on top of any given basic multicast, for this, the

Variables such as S_{g^p} & R_{g^q} can be used. These variables are held at process p from the reliable multicast protocol. S_{g^p} refers to the count of messages sent by p to g , whereas R_{g^q} refers to the recorded sequence number of the latest message that has been delivered from g by p to group g .

Fault tolerance:-

It is an important feature that plays a key role in distributed systems. In ds, if one of the critical component of a system fail to function properly, it results in a partial failure. To overcome such partial failures, ds incorporate a design which tolerate faults and continues to operate even in the existence of faults. Such characteristics feature in ds is referred to as fault tolerance.

Due to some faults in h/w or sw programs generate inappropriate results and sometime may stop before the completion of the computation needed. So, the ability of a system to continue functioning in the event of partial system failure. There are two approaches:

1. Hardware Redundancy

2. Software Redundancy

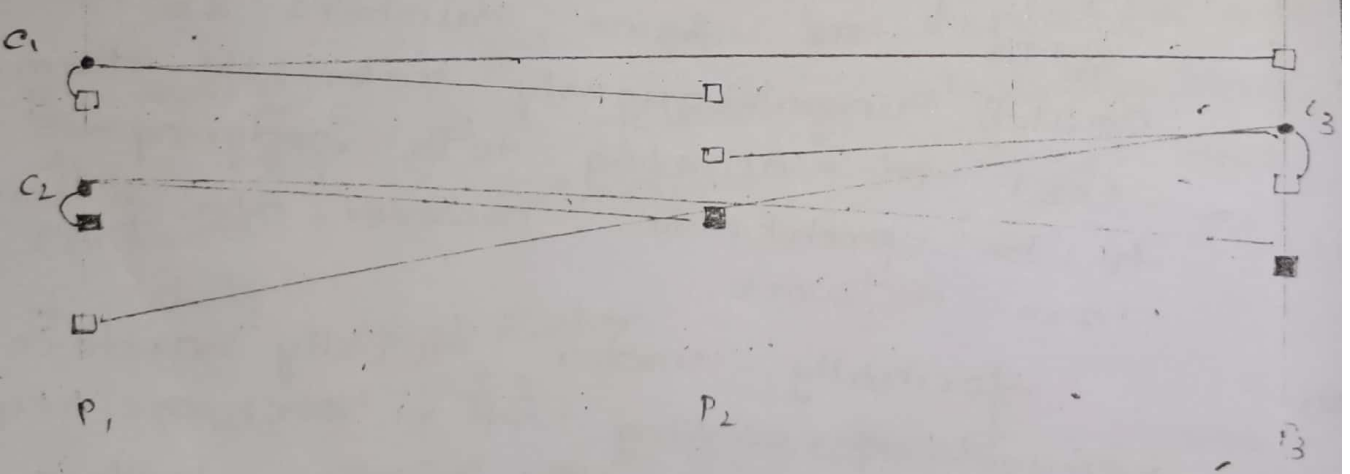
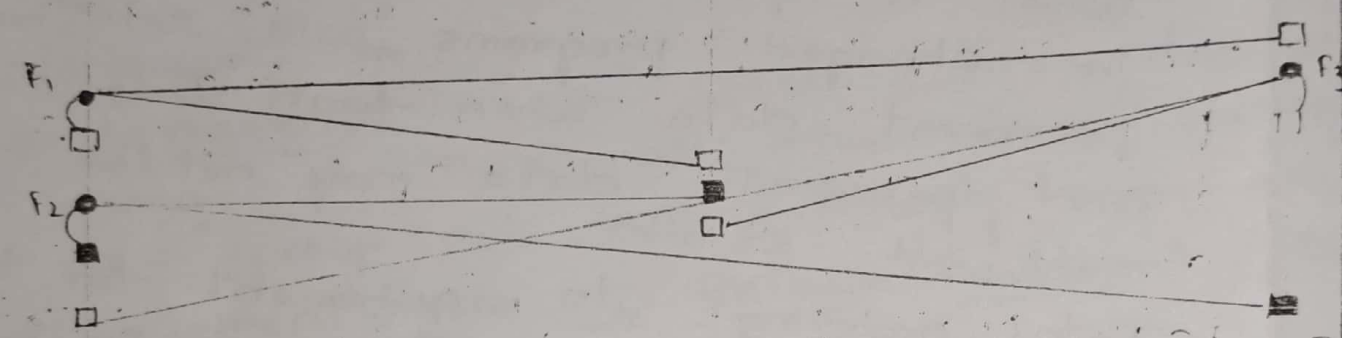
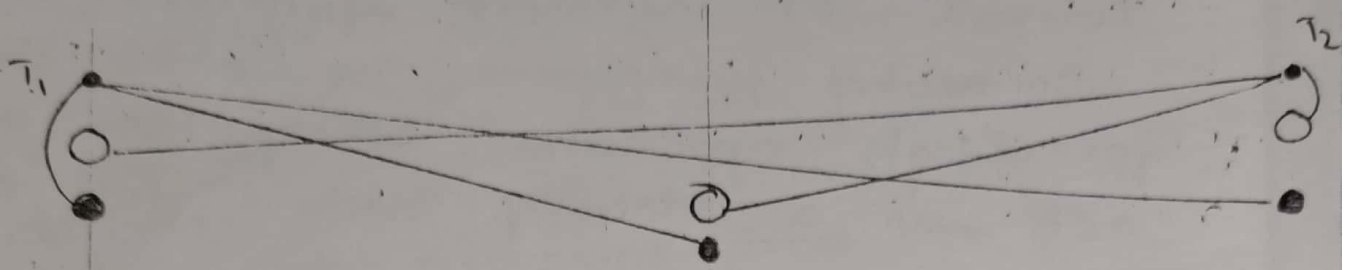
Hardware Redundancy:- It uses redundant components to produce systems that are tolerant of network failures, two inter-connected computers are often employed for single application where one of them acts as a standby machine for other i.e. ATM.

Software Redundancy:- It implies that the data can be rolled back or recovered when a fault is detected. In other words, the changed programs will set back to permanent data when fault occurs and permanent data may not be in consistent state.

Total ordering:- A requirement for total ordering is exemplified by the bulletin board system, where it would be convenient if replica managers could label items with the same numbers so that users could unambiguously refer to them. The cost of achieving total ordering is liable to be prohibitive, however, over a wide area network.

formally, under totally ordered request processing, if r_1 & r_2 are requests then either r_1 is processed before r_2 at all replica managers or r_2 is processed

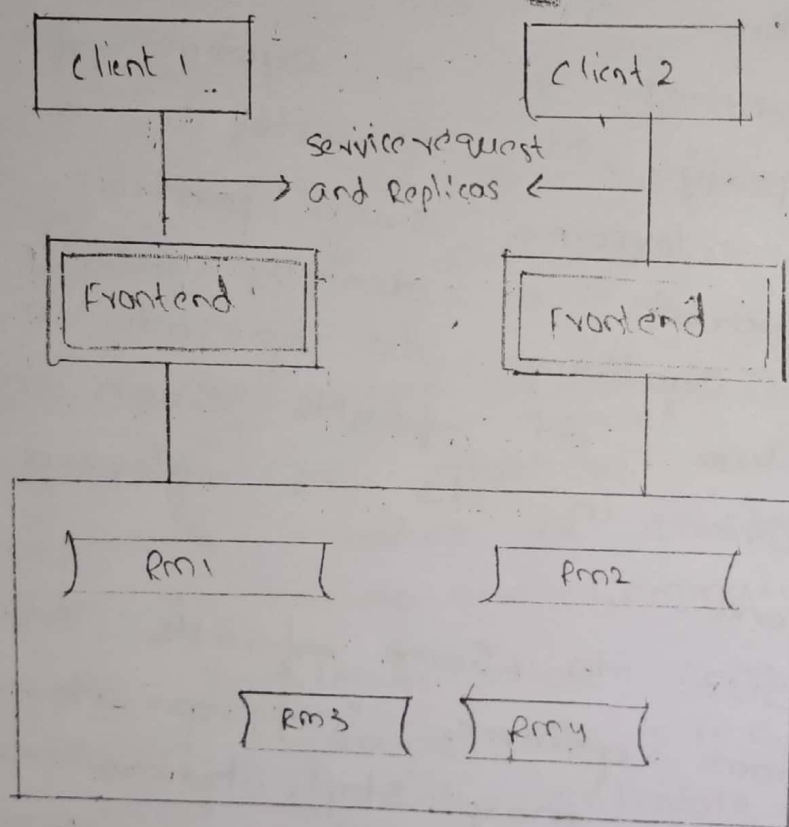
before r_1 at all replica managers, total ordering is a general relation over extent in distributed system.



6: TRANSACTIONS & REPLICATIONS.

System model & Group Communication

③ System model: In this model, a logical object is implemented using multiple physical copies called replicas. These replicas are physical copies stored in computers. They may not be consistent at all times because one replica can get information that other replicas did not get. The data is replicated and held by unique replica manager.



Different clients request for service through front end. The frontend replica managers apply operations and when a crash happens, they do not leave inconsistent

result:

Generally, every replica manager maintains replicas of every object. However different object's replica can also be maintained by different sets of replica managers. Moreover, the managers can be either dynamic or static. In dynamic system, new replica manager can appear and it can crash and leave the system whereas in static system, manager cannot appear and it cannot crash although it can leave the system for a time period.

There are five phases involved in the performance of single request upon the replicated objects.

Example: (A single request upon replicated object is - per) → Service that offers fault-tolerant services will differ from a service that offers disconnected operations.

Phase-1 Issue Request! In this phase, the front end request is a single replica manager and this manager passes the request to other replica manager. It also multicast request to replica manager.

Phase - 2 Coordination: In this phase, replica managers decide whether to apply the request or not. They also decide the order of a request in relation to other request in one of following orders...

FIFO order: In this order, at a front-end request r_1 and then request r_2 then replica manager will handle r_1 first and then handle r_2 .

Causal order: In this order, when a request r_1 happens before request r_2 , then replica manager will handle r_1 first and then r_2 .

Total order: In this phase, when a replica manager handle request r_1 and then r_2 .

Phase-3 Execution: In this phase, replica manager ^{executes} agrees on the request, effect.

Phase-4 Agreement: In this phase, replica manager agrees on the request effect.

Phase-5 Response: In this phase, front-end receives request from one or more replica managers.

SET-3
3 Group Communication:

It is also known as multicast communication. It is useful specially,

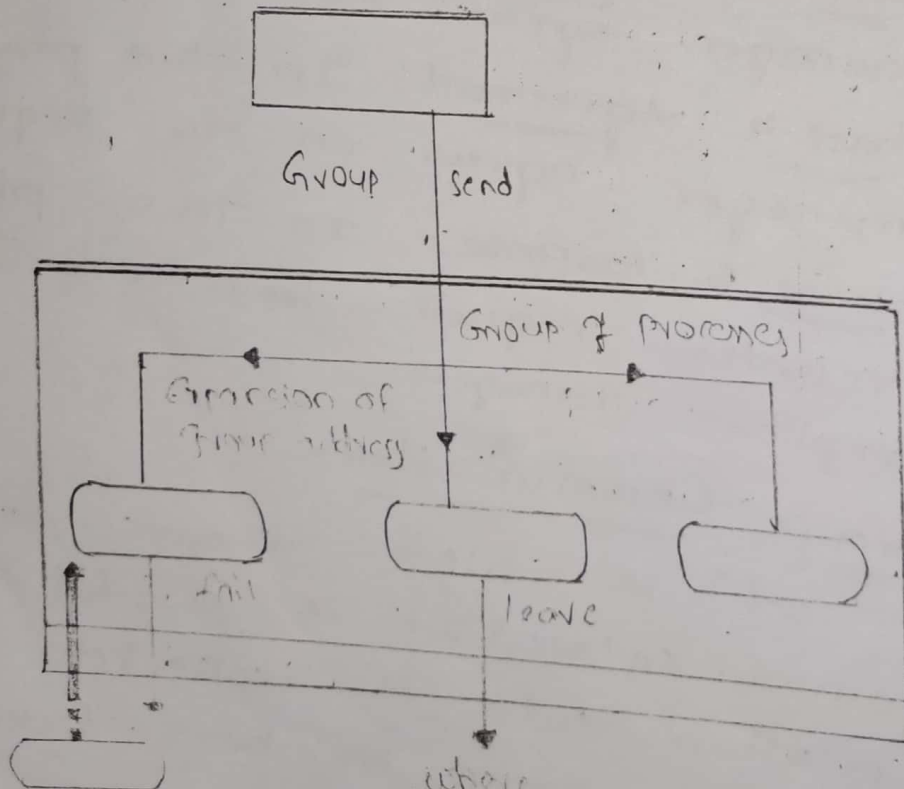
1. When data is replicated
2. When group members take common

message stream and

3. When the system process that cooperates to achieve common target through messages are managed by groups;

In general, a process may join or leave a group and the membership of a group can be either static or dynamic. Service. The dynamic membership is managed by group membership service. For example, in replicated data service - users may add or withdraw a replica manager or the manager may crash and need to be replaced.

Moreover group membership management multicast are connected to each other or depicted as below



where fail, leave, join are group membership management

→ four main tasks performed by group membership service

1. Interface to Add/Remove members: The membership service offers operations through which process groups are created and destroyed. This service also allows to add and remove a process from a group.

2. Incorporate a Failure Detector: The membership service includes a failure detector. This detector checks for failure when a crash occurs and also when a group member becomes unreachable due to connection failure.

3. Communicate Group member: When a process is added to a group or removed from a group, the service communicates it to other group members.

4. Expand Group Address: A group identifier is supplied in the group when a process multicasts a message. This identifier is expanded by member management service into current group members. The service can also coordinate delivery by controlling address expansion. In other words, the service is capable to decide about the message delivery even when the membership is modified during its delivery.

View delivery: The group management

service in a group provides any member process $p \in G$ set of views $view_0(G), view_1(G), \dots$ fundamental requirements

Order of view: when a process p delivers a $view(G)$ and again delivers $view'(G)$ then other process $\neq p$ cannot deliver $view'(G)$ before dealing $view(G)$.

This order of view is always change in same order for different process.

Integrity: when process p delivers a $view(G)$ then $p \in view(G)$. This is to perform sanity checking.

Non-triviality: when a process p_1 joins a group and if it becomes unreachable from process $\neq p_1$ then p_1 is always in view that process p_2 will deliver. Moreover,

when a group is partitioned, the view delivered in one partition will not contain any process in other partition of a group. It prevents against trivial solutions for process.

→ View - Synchronous group communication
This system ensures delivery of multicast messages view notification through

Agreement. At a given time same sequence of views are delivered by process along with same set of messages.

Integrity: This ensures that when a process p_1 delivers a message m_1 ,

then m_1 will not be delivered twice. It also ensures that the sender will be in same view that delivered m_1 . Validity In general, a correct process always detects messages, however if a system fails to detect a message - a process p_i , then the system informs the remaining processes through a new view delivery and this message will not be delivered to p_i .

* Concurrency Control in Distributed Transactions:-

locks in Distributed Transaction:- The locks in a distributed transaction on an object are maintained in the same server by a lock manager. On one hand, this lock manager decides to make the lock by requesting transaction to wait or grant a lock to it. On the other hand, it cannot unlock the object held by transaction until it gets information that transaction is either aborted or committed at all servers.

In locking mechanism, when an object is locked by a transaction it is inaccessible by other transaction. However, when a transaction is aborted its locks are unlocked after phase 1 of locking protocol.

Optimistic Concurrency Control:

every transaction is validated before being committed & distributed transaction is validated by many independent servers that are involved in a transaction. However each server validate its own objects and operators. This validation is performed in phase-1 of the two-phase protocol.

In general, transactions suffer from committed deadlock. Consider table which shows interleaving of transactions that suffer from a committed deadlock.

| Transaction P | | Transaction Q | |
|---------------|--------|---------------|--------|
| operation | server | operation | Server |
| Read(L) | at M | Read(K) | at N |
| write(L) | at M | write(K) | at N |
| Read(K) | at N | Read(L) | at M |
| write(K) | at N | write(L) | at M |

This table shows that transaction P access object before transaction Q on M server whereas transaction Q access object before transaction P on N server. However validation protocol states that only one transaction is allowed to validate and update at a time, the other need to wait. This situation is called committed deadlock. Therefore a new method called parallel validation

Should be used to overcome
Commit deadlock problem.

Parallel Validation Protocol

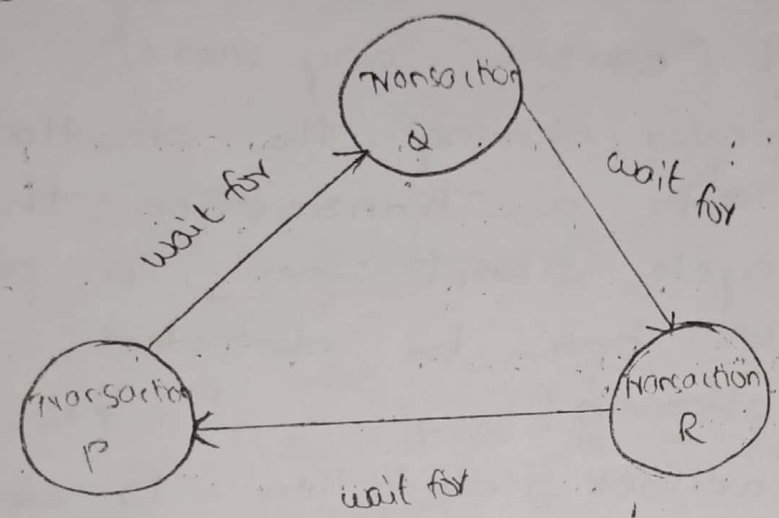
If it is used, transactions will
suffer from commitment deadlock.
However, if servers simply perform
independent validations, it is possible
that time different servers of a
distributed transaction may serialize
the same set of transactions in differ-
ent orders.

② Timestamp ordering Concurrency Control
every transaction coordinator must
issue a timestamp which should be
unique globally. When a transaction acquires
a coordinator, the client is issued a
globally unique transaction timestamp
by the first coordinator. This timestamp
is supplied to the transaction coordinator
when a transaction is performed in a
server.

Moreover, when a distributed
transaction is performed the servers
which perform are combinely responsible
in a sequential manner. This can be
achieved by timestamp maintenance. It
is made up of \langle local timestamp, server id
pair. Moreover, maintaining transactions
of same order is possible in all the
servers in spite of local clocks.

Distributed Deadlocks:-

Distributed transactions can lead to distributed deadlocks. A global wait for graph can be formed by local wait for graphs. This global graph can contain a cycle which does not belong to any single local wait for graph. This cycle is called as distributed deadlock.



or from above server M, transaction R wait for its lock. Similarly at server K, transaction L wait for transaction P. Finally, transaction M wait for transaction R at server L. These local wait for graphs are built by each server's lock manager. In order to solve this they should be detected by finding a cycle in global transaction wait for graph. This cycle can be found by communicating different servers that were

involved in transaction. Distributed deadlock can be solved by using centralized deadlock detection method. In this method, one server is assigned the responsibility of global detector.

Some advantages of centralized servers depend on single server is:

1. Lack of fault tolerance
2. Lack of scalability
3. Poor availability
4. Expensive communication among different servers
5. Deadlock detection may consume more time, if global graph is not collected frequently.

Preventing Deadlocks in DS:-

It can be done in two ways:-

Deadlock can be prevented if for a transaction all the required objects are locked before the transaction is started. This must be done in a single operation in order to avoid deadlock with other transactions.

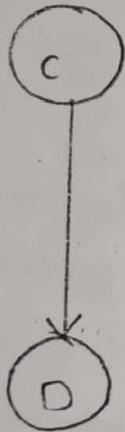
However this method of locking has some issues associated with it such as,

→ Access to shared resource is restricted if it is locked

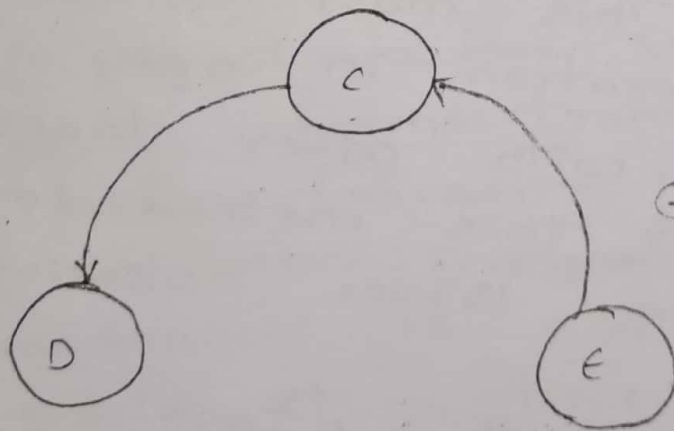
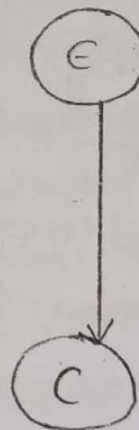
→ Sometime the objects required to perform transaction are not known before starting.

→ Phantom Deadlock:
 It is a type of deadlock that is detected but it is not a real deadlock. The mechanism of distributed deadlock detection involve transmission of wait for relation information between the servers. During this transmission, if there is a deadlock then it is detected. However this method consume sometime during which a lock can be released by a transaction and deadlock does not exit.

Server A



Server B



Global deadlock detector wait for graph

The above figure shows global and local wait for graphs. Suppose that

→ At Server A transaction C unlocks object and request for the object B server locked by transaction E and then

→ The global detector get B server's local graph before receiving server A's local graph.

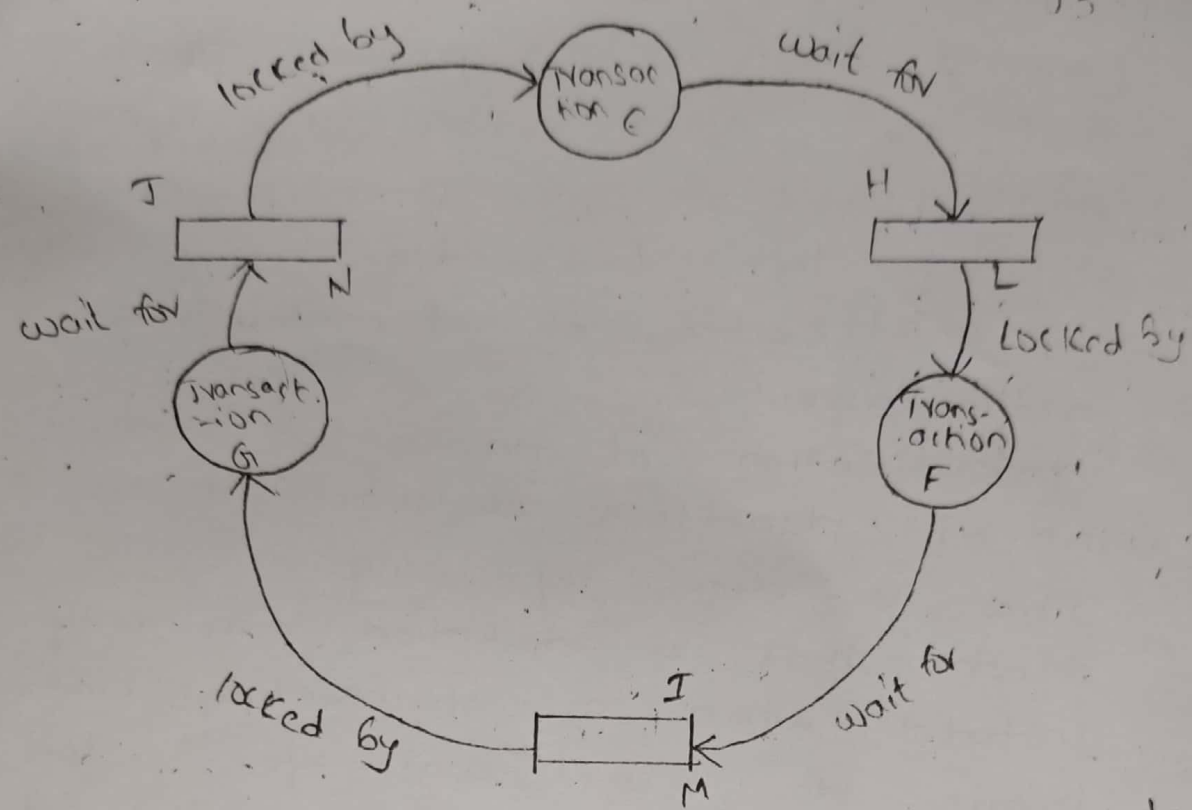
The global detector will detect the cycle $C \rightarrow D \rightarrow E \rightarrow C$ in which $C \rightarrow D$ cycle does not exist any more.

Moreover during the deadlock detection method, if a transaction that wait in a cycle abort then a phantom deadlock can be detected.

→ Edge chasing!

Deadlock detection is done by using a technique known as edge chasing. This technique is also known as path pushing. In this technique, the global wait for graph is not constructed instead every server involved in transaction has some information about its edges. These servers find cycles by using a message known as probes. These probes proceed the graph edges all the way in distributed systems. These message contain wait for relationship that request a path in the global wait-for graph. However time to send a probe is

Complicated task



1. Server N has added a new edges $G \rightarrow E$ in its local graph. At this point transaction E is not waiting for an object. Therefore no need to send a probe.
2. Server L has added a new $E \rightarrow F$ edge in its local graph at the same time transaction F wait for I object. This I object is locked at M. Server by transaction G. This forms a cycle $G \rightarrow T_1 \rightarrow T_2 \dots E \rightarrow F \rightarrow G$ involving many servers.

Edge chasing algorithm has 3 phases

1. Initiation
2. Detection
3. Resolution

Initiation:- In this phase, a server initiates a phase when it sees that, for ex transaction x wait for transaction y and y wait for z transaction. The server add an edge $x \rightarrow y$ in a probe and send to the k server. when transaction y is locked on object. In case, if lock is shared by other transaction then all the transactions that share this lock are also sent the probe.

Detection:- In detection phase, server k receive a probe $x \rightarrow B$. for example stating that transaction x is waiting for an object locked by transaction B . This server k also verifies if transaction B is waiting for transaction c then it is also added in the probe and form a cycle $x \rightarrow B \rightarrow c$. Similarly, if transaction c wait for some other transactions then other transactions are also added into the probe.

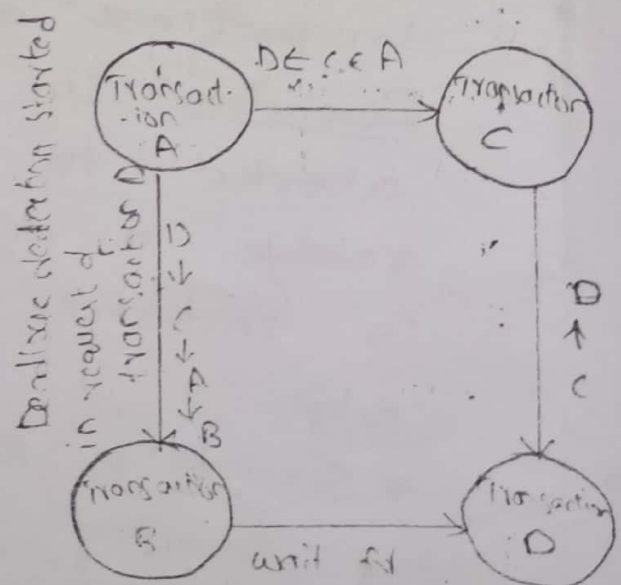
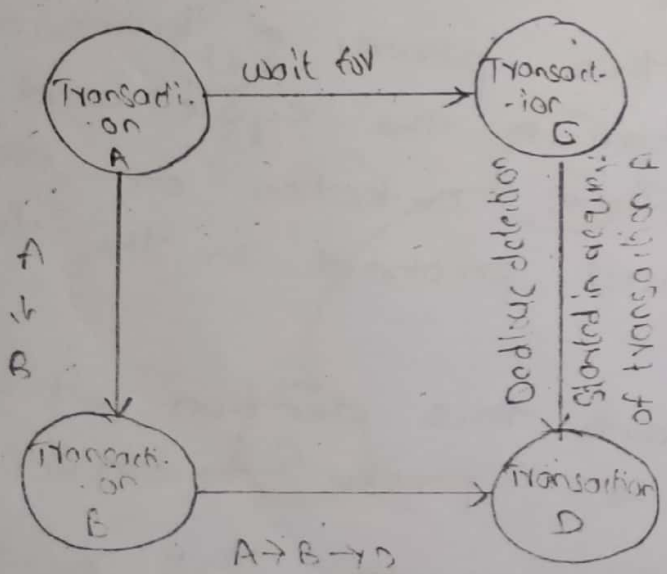
Resolution:- In this phase, a transaction is aborted to break the cycle and to resolve a deadlock. Detection of deadlock is initiated and probed in the following steps.

- * Server p starts this detection of deadlock by sending $L \rightarrow I$ probe to server q .
- * Once q server receives this probe it check that object N is locked by

transaction K. It further add to transact
 - on K and produce $L \rightarrow J \rightarrow K$ probe
 * finally, R server receive $L \rightarrow J \rightarrow K$ probe
 and it see that object O is locked
 by the transaction L and create
 $(L) \rightarrow J \rightarrow K \rightarrow L$ probe. This probe forma
 cycle and a deadlock is detected.

→ Prioritization of Transaction!

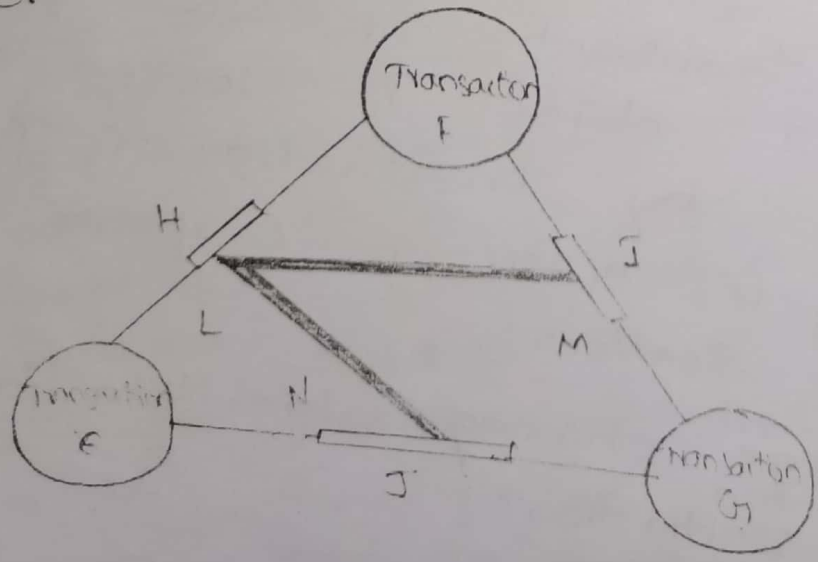
transactions are assigned priority
 to ensure that in a cycle only one
 transaction is aborted in case of a
 deadlock. If priority is not assigned
 to a transaction then there are chances
 that multiple transactions are aborted in
 a cycle. + transaction can be prioritize
 based on the timestamp and a transaction
 that has lowest priority in a cycle
 is aborted and a deadlock is detected.



from the above for it shows, for example as $Z > C$, where A has high priority than C. Suppose that transaction $A > B > C > D$ and when any cycle either $A \rightarrow B \rightarrow D \rightarrow C \rightarrow A$ or $D \rightarrow C \rightarrow A \rightarrow B \rightarrow D$ is detected then transaction D is aborted.

prioritization of transactions offer some advantages. Such as,

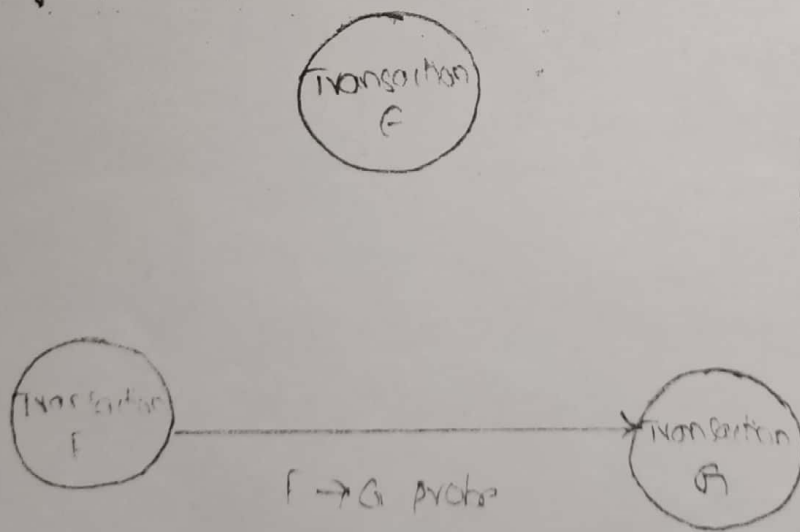
1. prioritization of transaction reduce the no. of deadlock detection to be initiated. This is done by using a rule which state that a detection can only be initiated when a transaction has higher priority to start, wait for a low priority transaction. This rule reduce the detection intention by approximately 50 percent.
2. prioritization of transactions also reduce the no. of probes to be forwarded. This is done by applying a rule which state that a probe should be sent from higher priority transaction to lower



transaction E start to wait for F where E is already waiting to G and return transaction G wait for E . In this circumstance, if transactions are not prioritized then a probe $E \rightarrow F$ can be sent and cycle can be detected. If transactions are prioritized then this probe cannot be sent as $E < F$ as a result deadlock cannot be detected.

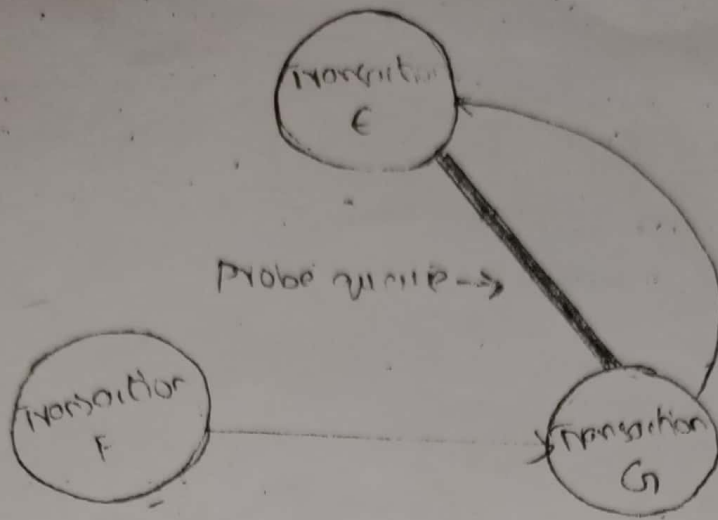
To overcome this problem, transaction coordinator showed save the probe copies sent by transaction in a queue known as probe queue and it should be forwarded to the server that wait for an object.

Example: where transaction E starts to wait for transaction G , the probe $F \rightarrow G$ is saved by transaction coordinator of G .



After this G transaction start to wait for transaction E , this forms a

queue $F \rightarrow G$ and $G \rightarrow E$



when transaction E start to wait for object P then $F \rightarrow G \rightarrow E$ probe queue is passed to P server. This P check probe $F \rightarrow E$ and add this is in for to the queue to form $F \rightarrow G \rightarrow E \rightarrow F$ & this way deadlock is detected.

* Transaction Recovery!-

transaction recovery def!-

It ensures that the server objects are durable which requires that objects must be saved in permanent storage. It also ensures that failure of transaction has atomic effect. This recovery is done by a mechanism called recovery manager.

The recovery manager performs the following tasks.

1. It Save objects in permanent storage for the transactions that are committed
2. It rearrange the recovery file in order to enhance the recovery performance
3. It restore the objects of server after a crash occur.
4. It reclaim storage space in recovery file.

* Replication :- It is a process through which multiple copies of data are maintained in different locations.

Motivations for Replication:- This technique is employed to

1. Improving performance
2. To keep system available at most of the time
3. To make system that can handle fault tolerance.

Improving performance: A server performance can be improved by using replication. For example, resources are cached in web browsers and proxy servers from web servers. Moreover workload is shared through binding all IP addresses of the server to a website DNS. This DNS lookup will results in better performance with minor cost. However, there is a limit of improving performance through this technique. For example, if

dynamic data is replicated it needs processing overheads to handle the changing data.

Enhanced Availability: In general, a server suffers from problems such as failure, network partition, unplanned disconnection and data locking. Due to these problems a server cannot be active and accessed 100% of time. These problems can be solved by replication of data.

In replication technique, data is automatically maintained in two or more locations on different servers that cannot be failed. These servers can be accessed when a default server fails.

$$1 - \text{prob} \text{ (when a server fails)} \\ = 1 - p_{\text{server fail}}^n$$

Where

n = number of servers

prob = probability of server failure.

Fault Tolerance: Replication of data guarantees that accuracy of data through fault tolerance service although few faults can occur. The accuracy is maintained on client's data. Sometimes accuracy is also related to timeliness of response from

Service:

n

IV B.Tech II Semester Regular/Supplementary Examinations, July - 2021

DISTRIBUTED SYSTEMS

(Common to Computer Science and Engineering and Information Technology)

Time: 3 hours**Max. Marks: 70***Question paper consists of Part-A and Part-B**Answer ALL sub questions from Part-A**Answer any FOUR questions from Part-B*

PART-A(14 Marks)

1. a) Define and Give examples of Distributed systems. [3]
- b) Write about IGMP. [3]
- c) Why distributed garbage collection is important? [2]
- d) What is the role of kernel in OS? [2]
- e) What are the goals of distributed file system? [2]
- f) List the advantages of Data replication. [2]

PART-B(4x14 = 56 Marks)

2. a) Discuss various issues and challenges involved in the implementation of Distributed Systems. [7]
- b) How the security model ensures security to the interacting processes in a Distributed System? Explain. [7]
3. a) Draw the structure of UDP datagram and explain about various structures available in JAVA API for UDP transmission. [7]
- b) What is meant by Multicast transmission in Distributed Systems? Explain some of the important applications of Multicast Transmission in Distributed systems. [7]
4. a) With a neat sketch, Explain the implementation of Remote Method Invocation. [7]
- b) Discuss the issues in design and implementation of RMI in Distributed Systems. [7]
5. a) What are the design issues of distributed operating system? [7]
- b) Explain any five advantages of creating Threads over multiple execution environments. [7]
6. a) Write the differences between Overlay networks and IP routing. [7]
- b) What are the requirements for mutual exclusion in Distributed systems? Explain about various metrics used for evaluating the performance of mutual exclusion algorithms in Distributed systems. [7]
7. a) Write about the Local and Global Wait-for graphs. [7]
- b) Explain the passive replication model for fault tolerance in distributed systems. [7]

Code No: R1642051

R16

Set No. 1

IV B.Tech II Semester Regular Examinations, September - 2020

DISTRIBUTED SYSTEMS

(Common to Computer Science and Engineering and Information Technology)

Time: 3 hours

Max. Marks: 70

Question paper consists of Part-A and Part-B

Answer ALL sub questions from Part-A

Answer any FOUR questions from Part-B

PART-A (14 Marks)

1. a) Discuss about client server resource sharing. [3]
- b) What is Multicast Transmission in Distributed systems? Discuss. [2]
- c) Discuss about Remote Procedure Calls. [2]
- d) What is meant by Address space? Discuss. [2]
- e) What is election process? Discuss about its goal? [2]
- f) What is replication? Differentiate between Active and passive replication. [3]

PART-B (4x14 = 56 Marks)

2. Explain the architectural and fundamental models of distributed systems? [14]
3. a) Explain the client server communication model. Also Discuss about marshaling in detail. [7]
- b) Discuss the issues relating to datagram communication. [7]
4. a) Explain the features of distributed object model [7]
- b) Explain the design issues of RMI. [7]
5. a) Briefly explain architecture for multi threaded servers. [7]
- b) What is the need for protection? Explain various protection mechanisms supported by operating systems. [7]
6. a) Discuss the mounting issues of remote file systems on NFS client. [7]
- b) Explain about overlay routing? Explain how it useful in peer communication. [7]
7. a) Describe various deadlock handling techniques. [7]
- b) Explain about concurrency control in distributed transactions. [7]

Code No: R1642051

R16

Set No. 2

IV B.Tech II Semester Regular Examinations, September - 2020

DISTRIBUTED SYSTEMS

(Common to Computer Science and Engineering and Information Technology)

Time: 3 hours

Max. Marks: 70

Question paper consists of Part-A and Part-B

Answer ALL sub questions from Part-A

Answer any FOUR questions from Part-B

PART-A (14 Marks)

1. a) What is meant by distributed system? Give any two examples. [2]
- b) Discuss about any three applications of Multicast Transmission in Distributed systems. [3]
- c) What is an event and notifications? [2]
- d) What is meant by multi threaded model. Discuss. [2]
- e) Define overlay routing? What is its importance? [2]
- f) What is dead lock? How deadlock can be handled. [3]

PART-B (4x14 = 56 Marks)

2. a) What do you mean by Scalability of a distributed system? Explain the principles for designing scalable distributed systems. [10]
- b) Explain the security challenges of distributed systems. [4]
3. a) Explain the different methods for inter-process communication. [7]
- b) Discuss the issues relating to datagram communication. [7]
4. Discuss the design and implementation issues in Remote Method Invocation. [14]
5. a) Explain the general architecture of operating systems for Distributed Systems [7]
- b) What is thread? Explain the life cycle of the thread, with neat state diagram. [7]
6. a) Explain how mutual exclusion is handled in distributed system. [7]
- b) Discuss the Napster and its legacy with respect to distributed file systems. [7]
7. a) Explain the basic architectural model for the management of Replicated data. [7]
- b) What is transaction? Briefly explain about flat and nested distributed transactions. [7]

Code No: R1642051

R16

Set No. 3

IV B.Tech II Semester Regular Examinations, September - 2020

DISTRIBUTED SYSTEMS

(Common to Computer Science and Engineering and Information Technology)

Time: 3 hours

Max. Marks: 70

Question paper consists of Part-A and Part-B

Answer ALL sub questions from Part-A

Answer any FOUR questions from Part-B

PART-A (14 Marks)

1. a) What is meant by resource sharing? Discuss with an example. [2]
- b) Discuss about the characteristics of the IPC. [3]
- c) Differentiate static and dynamic invocation methods. [2]
- d) Differentiate between process and threads. [3]
- e) What is mutual exclusion? List its requirements. [2]
- f) Define replication? What is the importance of it? [2]

PART-B (4x14 = 56 Marks)

2. What is distributed systems? Explain its key characteristics of distributed system in detail. [14]
3. a) What is marshaling? Explain marshaling operations in detail. [7]
- b) Explain Multicast transmission in Distributed Systems? Discuss about important applications of Multicast Transmission in Distributed systems. [7]
4. a) What is the importance of distributed garbage collection? Explain the Distributed garbage collector algorithm. [7]
- b) Discuss about various Remote Procedure Calls. [7]
5. a) What is an Execution environment? Explain in detail about the process execution environment. [7]
- b) Describe the architecture for multi-threaded servers. [7]
6. a) What is distributed file system? Briefly explain the file service architecture. [7]
- b) What is the goal of an election algorithm? Explain it detail. [7]
7. a) What is concurrency? Write the importance of concurrency control in distributed systems. [7]
- b) What is distributed deadlock? Explain with example. [7]

Code No: R1642051

R16

Set No. 4

IV B.Tech II Semester Regular Examinations, September - 2020

DISTRIBUTED SYSTEMS

(Common to Computer Science and Engineering and Information Technology)

Time: 3 hours

Max. Marks: 70

Question paper consists of Part-A and Part-B

Answer ALL sub questions from Part-A

Answer any FOUR questions from Part-B

PART-A (14 Marks)

1. a) What is meant by failure handling? [2]
- b) Discuss about marshaling. [3]
- c) What is the importance of distributed garbage collection? [2]
- d) Discuss about protection mechanisms supported by operating systems. [3]
- e) What are the goals of election algorithm? [2]
- f) What is transaction? List the different types of transactions. [2]

PART-B (4x14 = 56 Marks)

2. a) Explain the design requirements and challenges for distributed systems. [10]
- b) Explain the client server resource sharing system. [4]
3. a) List and Explain the various socket primitives used in TCP stream communication. [7]
- b) Describe IP Multicast communication. [7]
4. a) With a neat sketch, Explain the implementation of Remote Method Invocation. [7]
- b) Explain communication between distributed objects, With a neat diagram. [7]
5. What is thread? Explain the issues related to thread programming, thread lifecycle, and thread synchronization. [14]
6. a) Explain the techniques to achieve high performance in distributed file systems. [7]
- b) Explain the main tasks of Routing Overlays. [7]
7. a) What is replication? Explain about Active and Passive replications [7]
- b) Compare and contrast the various methods of concurrency control. [7]

