# Distributed Systems
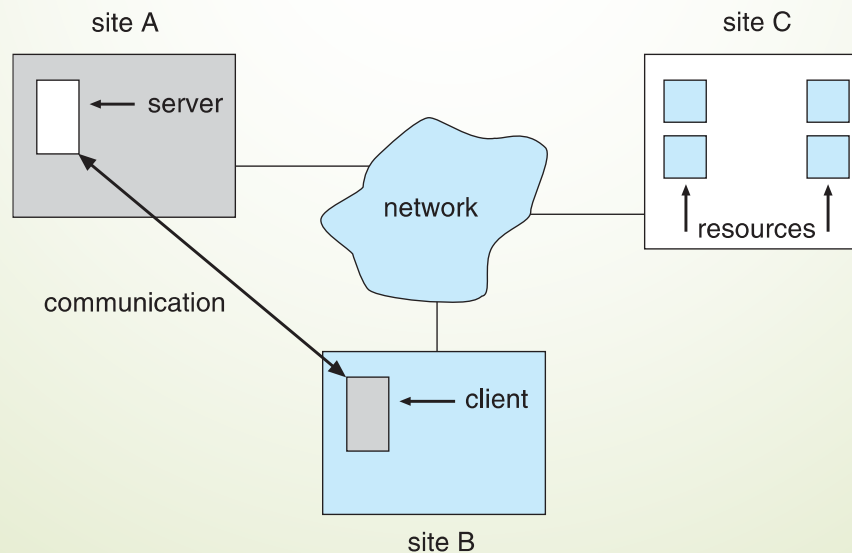
# Overview

n **Distributed system** is collection of loosely coupled processors interconnected by a communications network

n Processors variously called *nodes, computers, machines, hosts*

  l *Site* is location of the processor

  l Generally a *server* has a resource a *client* node at a different site wants to use

site A          site C

server

network

resources

communication

client

site B

# Reasons for Distributed Systems

- Reasons for distributed systems
  - **Resource sharing**
    - Sharing and printing files at remote sites
    - Processing information in a distributed database
    - Using remote specialized hardware devices
  - **Computation speedup** – **load sharing** or **job migration**
  - Reliability – detect and recover from site failure, function transfer, reintegrate failed site
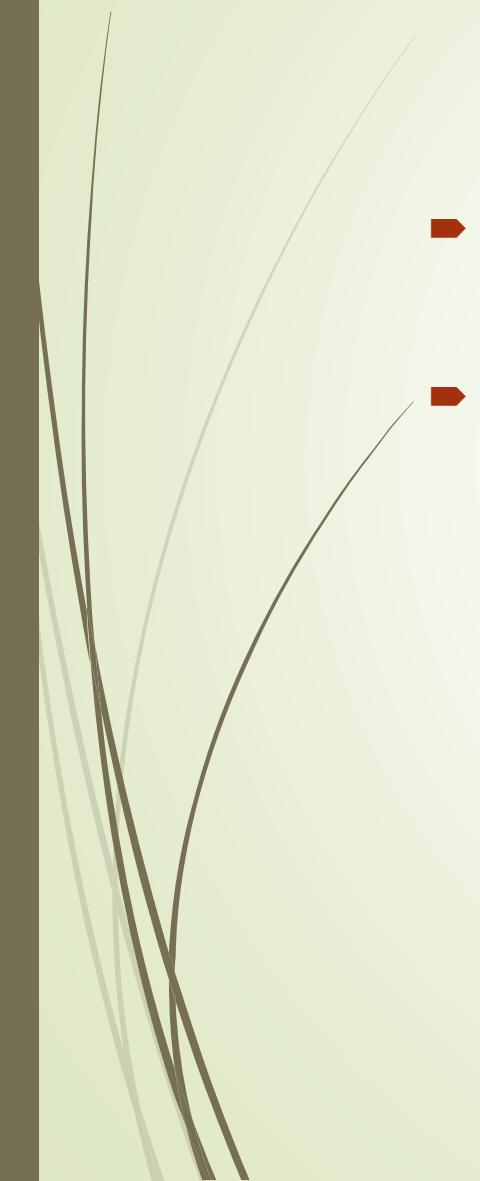  - Communication – **message** passing
    - All higher-level functions of a standalone system can be expanded to encompass a distributed system
  - Computers can be downsized, more flexibility, better user interfaces and easier maintenance by moving from large system to multiple smaller systems performing distributed computing

# Types of Distributed Operating Systems

- Network Operating Systems

- Distributed Operating Systems

# Network-Operating Systems

- Users are aware of multiplicity of machines

- Access to resources of various machines is done explicitly by:

  - Remote logging into the appropriate remote machine (telnet, ssh)

  - Remote Desktop (Microsoft Windows)

  - Transferring data from remote machines to local machines, via the File Transfer Protocol (FTP) mechanism

- Users must change paradigms – establish a **session**, give network-based commands

  - More difficult for users

# Distributed-Operating Systems

- Users not aware of multiplicity of machines
  - Access to remote resources similar to access to local resources
- **Data Migration** – transfer data by transferring entire file, or transferring only those portions of the file necessary for the immediate task
- **Computation Migration** – transfer the computation, rather than the data, across the system
  - Via remote procedure calls (RPCs)
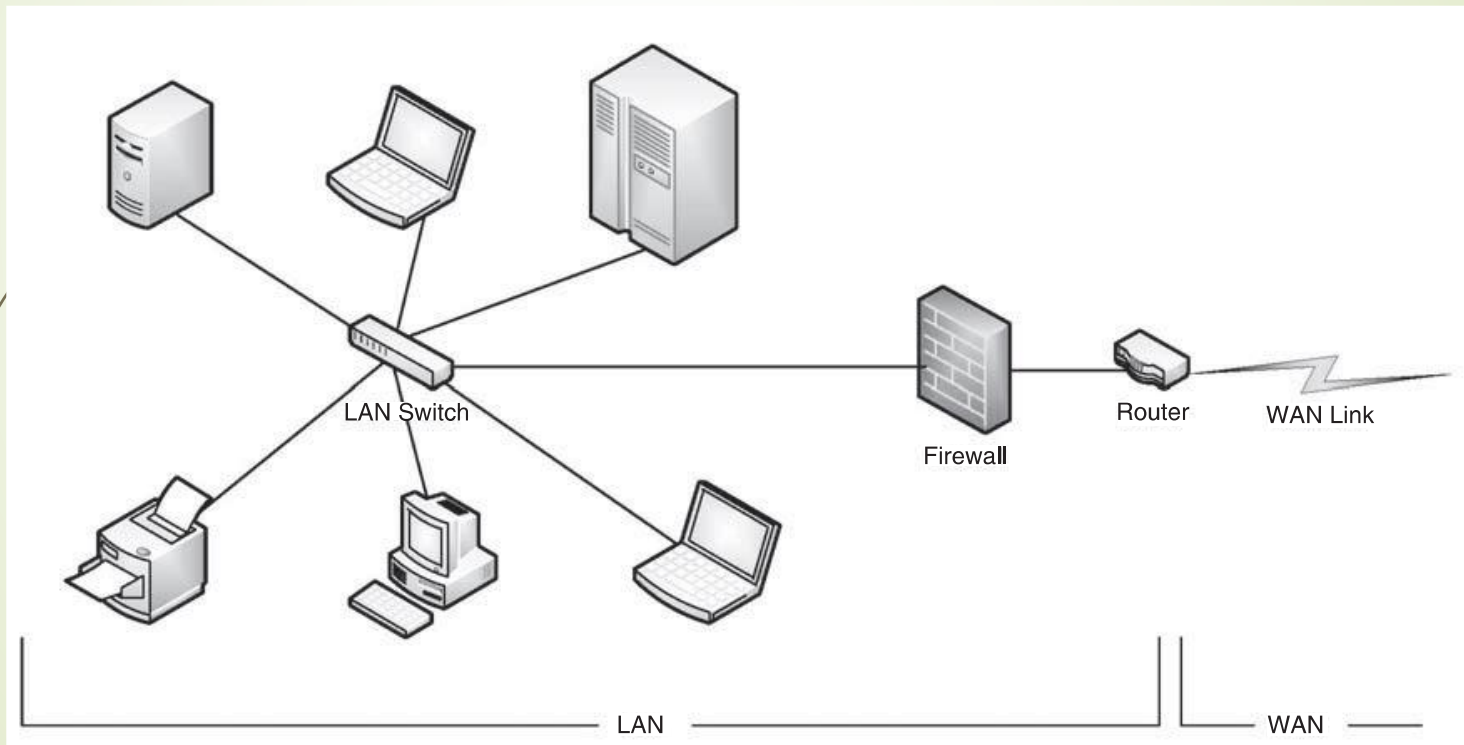  - or via messaging system

# Distributed-Operating Systems (Cont.)

- **Process Migration** – execute an entire process, or parts of it, at different sites

  - **Load balancing** – distribute processes across network to even the workload

  - **Computation speedup** – subprocesses can run concurrently on different sites

  - **Hardware preference** – process execution may require specialized processor

  - **Software preference** – required software may be available at only a particular site

  - **Data access** – run process remotely, rather than transfer all data locally

- Consider the World Wide Web

# Network Structure

- **Local-Area Network** (**LAN**) – designed to cover small geographical area
  - Multiple topologies like star or ring
  - Speeds from 1Mb per second (Appletalk, bluetooth) to 40 Gbps for fastest Ethernet over twisted pair copper or optical fibre
  - Consists of multiple computers (mainframes through mobile devices), peripherals (printers, storage arrays), routers (specialized network communication processors) providing access to other networks
  - Ethernet most common way to construct LANs
    - Multiaccess bus-based
    - Defined by standard IEEE 802.3
  - Wireless spectrum (**WiFi**) increasingly used for networking
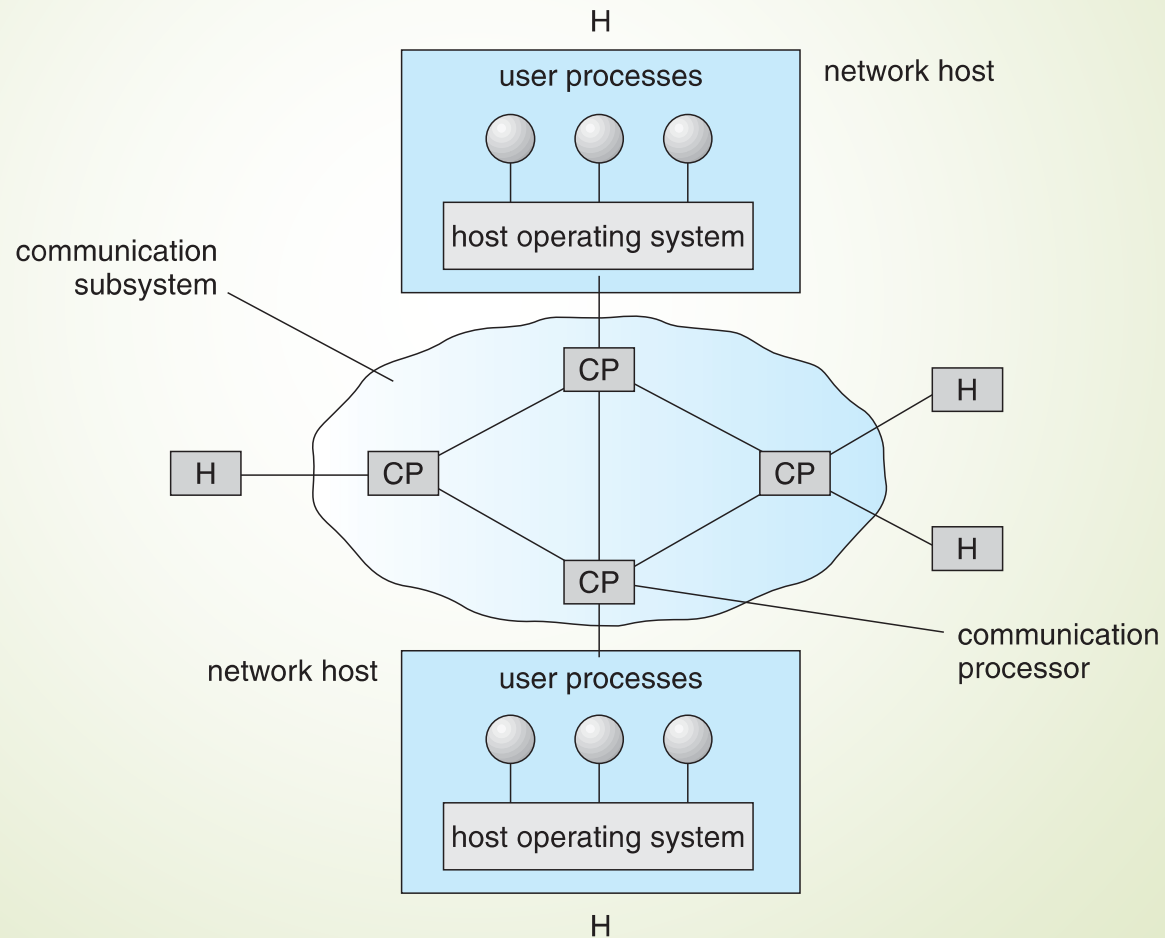    - I.e. IEEE 802.11g standard implemented at 54 Mbps

# Local-area Network

# Network Types (Cont.)

- **Wide-Area Network** (**WAN**) – links geographically separated sites
  - Point-to-point connections over long-haul lines (often leased from a phone company)
    - Implemented via **connection processors** known as **routers**
  - Internet WAN enables hosts world wide to communicate
    - Hosts differ in all dimensions but WAN allows communications
  - Speeds
    - T1 link is 1.544 Megabits per second
    - T3 is 28 x T1s = 45 Mbps
    - OC-12 is 622 Mbps
  - WANs and LANs interconnect, similar to cell phone network:
    - Cell phones use radio waves to cell towers
    - Towers connect to other towers and hubs

# Communication Processors in a Wide-Area Network

# Communication Structure

The design of a communication network must address four basic issues:

- **Naming and name resolution** - How do two processes locate each other to communicate?

- **Routing strategies** - How are messages sent through the network?

- **Connection strategies** - How do two processes send a sequence of messages?

- **Contention -** The network is a shared resource, so how do we resolve conflicting demands for its use?

# Routing Strategies

- **Fixed routing** - A path from *A* to *B* is specified in advance; path changes only if a hardware failure disables it

  - Since the shortest path is usually chosen, communication costs are minimized

  - Fixed routing cannot adapt to load changes

  - Ensures that messages will be delivered in the order in which they were sent

- **Virtual routing**-  A path from *A* to *B* is fixed for the duration of one session.  Different sessions involving messages from *A* to *B* may have different paths

  - Partial remedy to adapting to load changes

  - Ensures that messages will be delivered in the order in which they were sent

# Routing Strategies (Cont.)

- **Dynamic routing** - The path used to send a message form site *A* to site *B* is chosen only when a message is sent
  - Usually a site sends a message to another site on the link least used at that particular time
  - Adapts to load changes by avoiding routing messages on heavily used path
  - Messages may arrive out of order
    - This problem can be remedied by appending a sequence number to each message
  - Most complex to set up
- Tradeoffs mean all methods are used
  - UNIX provides ability to mix fixed and dynamic
  - Hosts may have fixed routes and **gateways** connecting networks together may have dynamic routes

# Routing Strategies (Cont.)

- **Router** is communications processor responsible for routing messages

- Must have at least 2 network connections

- Maybe special purpose or just function running on host

- Checks its tables to determine where destination host is, where to send messages

    - Static routing – table only changed manually

    - Dynamic routing – table changed via **routing protocol**

# Routing Strategies (Cont.)

- More recently, routing managed by intelligent software more intelligently than routing protocols
  - **OpenFlow** is device-independent, allowing developers to introduce network efficiencies by decoupling data-routing decisions from underlying network devices
- Messages vary in length – simplified design breaks them into **packets** (or **frames**, or **datagrams**)
- **Connectionless message** is just one packet
  - Otherwise need a connection to get a multi-packet message from source to destination
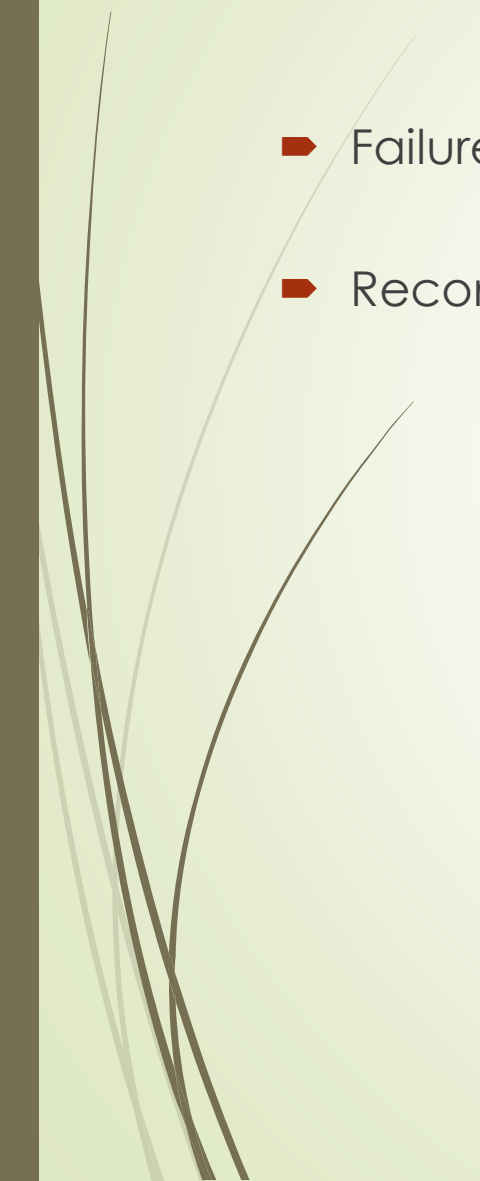
# Connection Strategies

- **Circuit switching** - A permanent physical link is established for the duration of the communication (i.e., telephone system)

- **Message switching** - A temporary link is established for the duration of one message transfer (i.e., post-office mailing system)

- **Packet switching** - Messages of variable length are divided into fixed-length packets which are sent to the destination

  - Each packet may take a different path through the network

  - The packets must be reassembled into messages as they arrive

- Circuit switching requires setup time, but incurs less overhead for shipping each message, and may waste network bandwidth

  - Message and packet switching require less setup time, but incur more overhead per message

# An Ethernet Packet

| bytes | | |
|---|---|---|
| 7 | preamble—start of packet | each byte pattern 10101010 |
| 1 | start of frame delimiter | pattern 10101011 |
| 2 or 6 | destination address | Ethernet address or broadcast |
| 2 or 6 | source address | Ethernet address |
| 2 | length of data section | length in bytes |
| 0–1500 | data | message data |
| 0–46 | pad (optional) | message must be > 63 bytes long |
| 4 | frame checksum | for error detection |

# Robustness

- Failure detection

- Reconfiguration

# Failure Detection

- Detecting hardware failure is difficult

- To detect a link failure, a **heartbeat** protocol can be used

- Assume Site A and Site B have established a link

  - At fixed intervals, each site will exchange an *I-am-up* message indicating that they are up and running

- If Site A does not receive a message within the fixed interval, it assumes either (a) the other site is not up or (b) the message was lost

- Site A can now send an *Are-you-up?* message to Site B

- If Site A does not receive a reply, it can repeat the message or try an alternate route to Site B

# Failure Detection (Cont.)

- If Site A does not ultimately receive a reply from Site B, it concludes some type of failure has occurred

- Types of failures:
  - Site B is down
  - The direct link between A and B is down
  - The alternate link from A to B is down
  - The message has been lost

- However, Site A cannot determine exactly **why** the failure has occurred

# Reconfiguration

- When Site A determines a failure has occurred, it must reconfigure the system:

    1. If the link from A to B has failed, this must be broadcast to every site in the system

    2. If a site has failed, every other site must also be notified indicating that the services offered by the failed site are no longer available

- When the link or the site becomes available again, this information must again be broadcast to all other sites