



MODULE 1: APPLICATION LAYER

1.1 Principles of Network Applications

- Network-applications are the driving forces for the explosive development of the internet.
- Examples of network-applications:
 - 1) Web
 - 2) File transfers
 - 3) E-mail
 - 4) P2P file sharing
 - 5) Social networking (Facebook, Twitter)
 - 6) Video distribution (YouTube)
 - 7) Real-time video conferencing (Skype)
 - 8) On-line games (World of Warcraft)
- In network-applications, program usually needs to
 - run on the different end-systems and
 - communicate with one another over the network.
- For ex: In the Web application, there are 2 different programs:
 - 1) The browser program running in the user's host (Laptop or Smartphone).
 - 2) The Web-server program running in the Web-server host.

1.1.1 Network Application Architectures

- Two approaches for developing an application:
 - 1) Client-Server architecture
 - 2) P2P (Peer to Peer) architecture

1.1.1.1 Client-Server Architecture

- In this architecture, there is a server and many clients distributed over the network (Figure 1.1a).
- The server is always-on while a client can be randomly run.
- The server is listening on the network and a client initializes the communication.
- Upon the requests from a client, the server provides certain services to the client.
- Usually, there is no communication between two clients.
- The server has a fixed IP address.
- A client contacts the server by sending a packet to the server's IP address.
- A server is able to communicate with many clients.
- The applications such as FTP, telnet, Web, e-mail etc use the client-server architecture.

1.1.1.1.1 Data Center

- Earlier, client-server architecture had a single-server host.
- But now, a single-server host is unable to keep up with all the requests from large no. of clients.
- For this reason, data-center is used.
- A data-center contains a large number of hosts.
- A data-center is used to create a powerful virtual server.
- In data center, hundreds of servers must be powered and maintained.
- For example:
 - Google has around 50 data-centers distributed around the world.
 - These 50 data-centers handle search, YouTube, Gmail, and other services.



COMPUTER NETWORKS

1.1.1.2 P2P Architecture

- There is no dedicated server (Figure 1.1b).
- Pairs of hosts are called peers.
- The peers communicate directly with each other.
- The peers are not owned by the service-provider. Rather, the peers are laptops controlled by users.
- Many of today's most popular and traffic-intensive applications are based on P2P architecture.
- Examples include file sharing (BitTorrent), Internet telephone (Skype) etc.
- Main feature of P2P architectures: self-scalability.
- For ex: In a P2P file-sharing system,
 - Each peer generates workload by requesting files.
 - Each peer also adds service-capacity to the system by distributing files to other peers.
- Advantage: Cost effective \'. Normally, server-infrastructure & server bandwidth are not required.
- Three challenges of the P2P applications:
 - 1) ISP Friendly**
 - Most residential ISPs have been designed for asymmetrical bandwidth usage.
 - Asymmetrical bandwidth means there is more downstream-traffic than upstream-traffic.
 - But P2P applications shift upstream-traffic from servers to residential ISPs, which stress on the ISPs.
 - 2) Security**
 - Since the highly distribution and openness, P2P applications can be a challenge to security.
 - 3) Incentive**
 - Success of P2P depends on convincing users to volunteer bandwidth & resources to the applications.

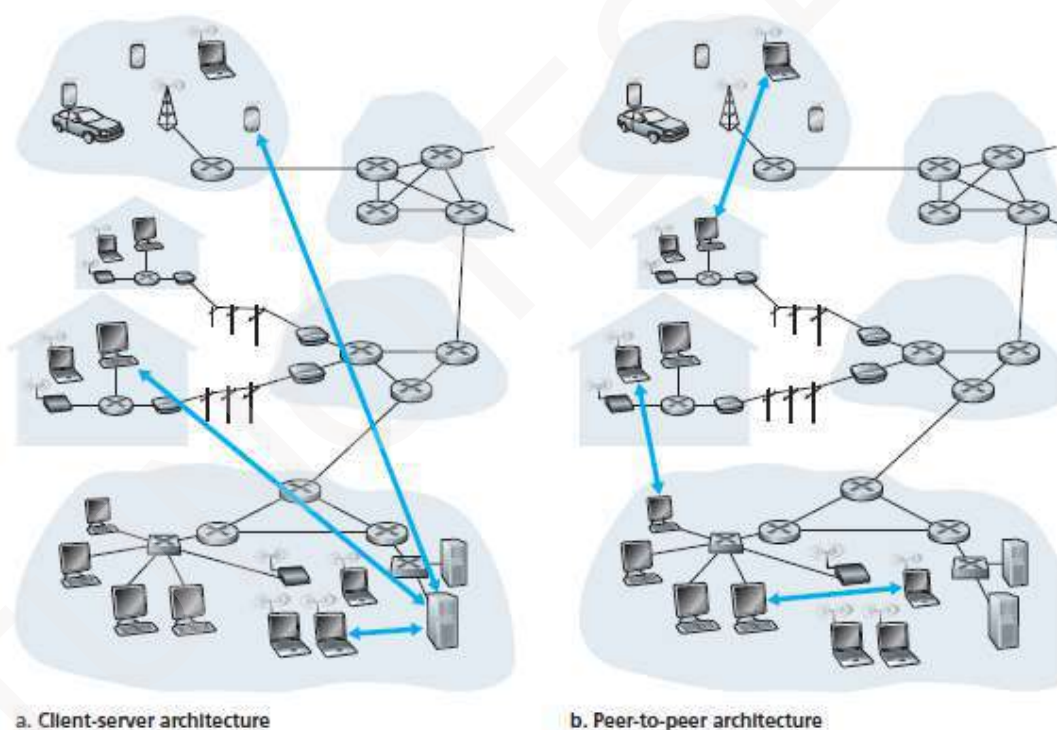


Figure 1.1: (a) Client-server architecture; (b) P2P architecture



COMPUTER NETWORKS

1.1.2 Processes Communicating

1.1.2.1 Process

- A process is an instance of a program running in a computer.(IPC → inter-process communication).
- The processes may run on the 1) same system or 2) different systems.
 - 1) The processes running on the same end-system can communicate with each other using IPC.
 - 2) The processes running on the different end-systems can communicate by exchanging messages.
 - i) A sending-process creates and sends messages into the network.
 - ii) A receiving-process receives the messages and responds by sending messages back.

1.1.2.1.1 Client & Server Processes

- A network-application consists of pairs of processes:
 - 1) The process that initiates the communication is labeled as the client.
 - 2) The process that waits to be contacted to begin the session is labeled as the server.
- For example:
 - 1) In Web application, a client-browser process communicates with a Web-server-process.
 - 2) In P2P file system, a file is transferred from a process in one peer to a process in another peer.

1.1.2.1.2 Interface between the Process and the Computer Network Socket

- Any message sent from one process to another must go through the underlying-network.
- A process sends/receives message through a software-interface of underlying-network called socket.
- Socket is an API between the application-layer and the transport layer within a host (Figure 1.2).
- The application-developer has complete control at the application-layer side of the socket.
- But, the application-developer has little control of the transport-layer side of the socket.

For ex: The application-developer can control:

- 1) The choice of transport-protocol: TCP or UDP. (API → Application Programming Interface)
- 2) The ability to fix parameters such as maximum-buffer & maximum-segment-sizes.

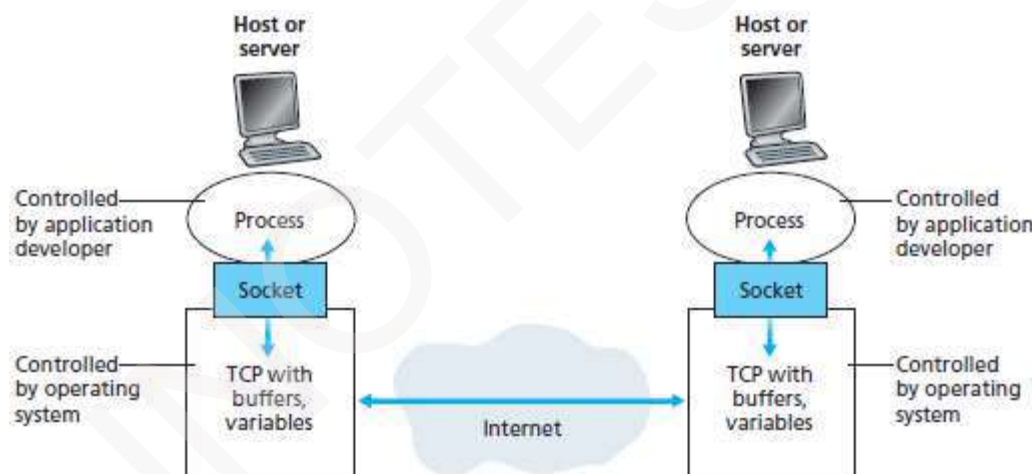


Figure 1.2: Application processes, sockets, and transport-protocol

1.1.2.1.3 Addressing Processes

- To identify the receiving-process, two pieces of information need to be specified:
 - 1) IP address of the destination-host.
 - 2) Port-number that specifies the receiving-process in the destination-host.
- In the Internet, the host is identified by IP address.
- An IP address is a 32-bit that uniquely identify the host.
- Sending-process needs to identify receiving-process '∴' a host may run several network-applications.
- For this purpose, a destination port-number is used.
- For example,
 - A Web-server is identified by port-number 80.
 - A mail-server is identified by port-number 25.



COMPUTER NETWORKS

1.1.3 Transport Services Available to Applications

- Networks usually provide more than one transport-layer protocols for different applications.
- An application-developer should choose certain protocol according to the type of applications.
- Different protocols may provide different services.

1.1.3.1 Reliable Data Transfer

- Reliable means guaranteeing the data from the sender to the receiver is delivered correctly.
For ex: TCP provides reliable service to an application.
- Unreliable means the data from the sender to the receiver may never arrive.
For ex: UDP provides unreliable service to an application.
- Unreliability may be acceptable for loss-tolerant applications, such as multimedia applications.
- In multimedia applications, the lost data might result in a small glitch in the audio/video.

1.1.3.2 Throughput

- Throughput is the rate at which the sending-process can deliver bits to the receiving-process.
- Since other hosts are using the network, the throughput can fluctuate with time.
- Two types of applications:
 - 1) Bandwidth Sensitive Applications**
 - These applications need a guaranteed throughput.
For ex: Multimedia applications
 - Some transport-protocol provides guaranteed throughput at some specified rate (r bits/sec).
 - 2) Elastic Applications**
 - These applications may not need a guaranteed throughput.
For ex: Electronic mail, File transfer & Web transfers.

1.1.3.3 Timing

- A transport-layer protocol can provide timing-guarantees.
- For ex: guaranteeing every bit arrives at the receiver in less than 100 msec.
- Timing constraints are useful for real-time applications such as
 - Internet telephony
 - Virtual environments
 - Teleconferencing and
 - Multiplayer games

1.1.3.4 Security

- A transport-protocol can provide one or more security services.
- For example,
 - 1) In the sending host, a transport-protocol can encrypt all the transmitted-data.
 - 2) In the receiving host, the transport-protocol can decrypt the received-data.



COMPUTER NETWORKS

1.1.4 Transport Services Provided by the Internet

- The Internet makes two transport-protocols available to applications, UDP and TCP.
- An application-developer who creates a new network-application must use either: UDP or TCP.
- Both UDP & TCP offers a different set of services to the invoking applications.
- Table 1.1 shows the service requirements for some selected applications.

Table 1.1: Requirements of selected network-applications

Application	Data Loss	Throughput Time	Sensitive
File transfer/download	No loss	Elastic	No
E-mail	No loss	Elastic	No
Web documents	No loss	Elastic (few kbps)	No
Internet-telephony/ Video-conferencing	Loss-tolerant	Audio: few kbps-1 Mbps Video: 10 kbps-5 Mbps	Yes: 100s of ms
Streaming stored audio/video	Loss-tolerant	Same as above	Yes: few seconds
Interactive games	Loss-tolerant	Few kbps-10 kbps	Yes: 100s of ms
Instant messaging	No loss	Elastic	Yes and no

1.1.4.1 TCP Services

- An application using transport-protocol TCP, receives following 2 services.
 - 1) Connection-Oriented Service**
 - Before the start of communication, client & server need to exchange control-information.
 - This phase is called handshaking phase.
 - Then, the two processes can send messages to each other over the connection.
 - After the end of communication, the applications must tear down the connection.
 - 2) Reliable Data Transfer Service**
 - The communicating processes must deliver all data sent without error & in the proper order.
- TCP also includes a congestion-control.
- The congestion-control throttles a sending-process when the network is congested.

1.1.4.2 UDP Services

- UDP is a lightweight transport-protocol, providing minimal services.
- UDP is connectionless, so there is no handshaking before the 2 processes start to communicate.
- UDP provides an unreliable data transfer service.
- Unreliable means providing no guarantee that the message will reach the receiving-process.
- Furthermore, messages that do arrive at the receiving-process may arrive out-of-order.
- UDP does not include a congestion-control.
- UDP can pump data into the network-layer at any rate.



COMPUTER NETWORKS

1.2 The Web & HTTP

- The appearance of Web dramatically changed the Internet.
- Web has many advantages for a lot of applications.
- It operates on demand so that the users receive what they want when they want it.
- It provides an easy way for everyone make information available over the world.
- Hyperlinks and search engines help us navigate through an ocean of Web-sites.
- Forms, JavaScript, Java applets, and many other devices enable us to interact with pages and sites.
- The Web serves as a platform for many killer applications including YouTube, Gmail, and Facebook.

1.2.1 Overview of HTTP

1.2.1.1 Web

- A web-page consists of objects (HTML → Hyper Text Markup Language).
- An object is a file such as an HTML file, a JPEG image, a Java applet, a video clip.
- The object is addressable by a single URL (URL → Uniform Resource Locator).
- Most Web-pages consist of a base HTML file & several referenced objects.
- For example:
 - If a Web-page contains HTML text and five JPEG images; then the Web-page has six objects:
 - 1) Base HTML file and
 - 2) Five images.
- The base HTML file references the other objects in the page with the object's URLs.
- URL has 2 components:
 - 1) The hostname of the server that houses the object and
 - 2) The object's path name.
- For example:
 - "http://www.someSchool.edu/someDepartment/picture.gif"
 - In above URL,
 - 1) Hostname = "www.someSchool.edu "
 - 2) Path name = "/someDepartment/picture.gif".
- The web browsers implement the client-side of HTTP. For ex: Google Chrome, Internet Explorer
- The web-servers implement the server-side of HTTP. For ex: Apache

1.2.1.2 HTTP

- HTTP is Web's application-layer protocol (Figure 1.3) (HTTP → HyperText Transfer Protocol).
- HTTP defines
 - how clients request Web-pages from servers and
 - how servers transfer Web-pages to clients.

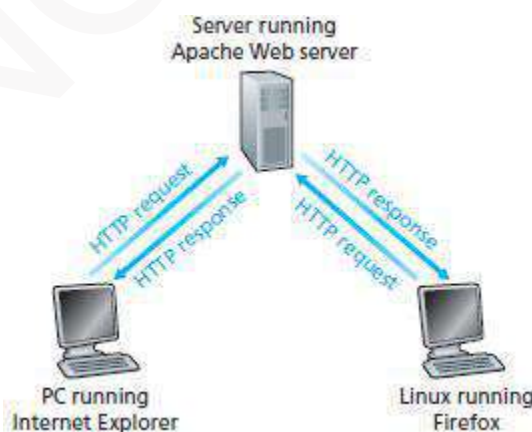


Figure 1.3: HTTP request-response behavior

- When a user requests a Web-page, the browser sends HTTP request to the server.
- Then, the server responds with HTTP response that contains the requested-objects.
- HTTP uses TCP as its underlying transport-protocol.
- The HTTP client first initiates a TCP connection with the server.
- After connection setup, the browser and the server-processes access TCP through their sockets.



COMPUTER NETWORKS

- HTTP is a stateless protocol.
- Stateless means the server sends requested-object to client w/o storing state-info about the client.
- HTTP uses the client-server architecture:
 - 1) Client**
 - Browser that requests receive and displays Web objects.
 - 2) Server**
 - Web-server sends objects in response to requests.

1.2.2 Non-Persistent & Persistent Connections

- In many internet applications, the client and server communicate for an extended period of time.
- When this client-server interaction takes place over TCP, a decision should be made:
 - 1) Should each request/response pair be sent over a separate TCP connection or
 - 2) Should all requests and their corresponding responses be sent over same TCP connection?
- These different connections are called non-persistent connections (1) or persistent connections (2).
- Default mode: HTTP uses persistent connections.



COMPUTER NETWORKS

1.2.2.1 HTTP with Non-Persistent Connections

- A non-persistent connection is closed after the server sends the requested-object to the client.
- In other words, the connection is used exactly for one request and one response.
- For downloading multiple objects, multiple connections must be used.
- Suppose user enters URL:
"http://www.someSchool.edu/someDepartment/home.index"
- Assume above link contains text and references to 10 jpeg images.

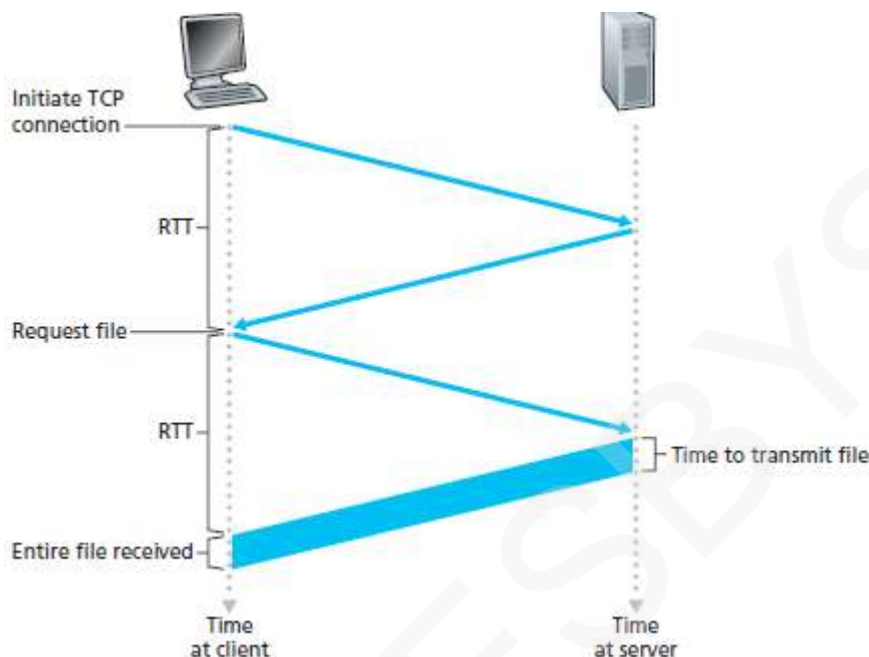


Figure 1.4: Back-of-the-envelope calculation for the time needed to request and receive an HTML file

- Here is how it works:

- 1a. HTTP client initiates TCP connection to HTTP server (process) at www.someSchool.edu on port 80
- 1b. HTTP server at host www.someSchool.edu waiting for TCP connection at port 80. "accepts" connection, notifying client
2. HTTP client sends HTTP request message (containing URL) into TCP connection socket. Message indicates that client wants object someDepartment/home.index
3. HTTP server receives request message, forms response message containing requested object, and sends message into its socket
4. HTTP server closes TCP connection.
5. HTTP client receives response message containing html file, displays html. Parsing html file, finds 10 referenced jpeg objects
6. Steps 1-5 repeated for each of 10 jpeg objects

- RTT is the time taken for a packet to travel from client to server and then back to the client.
- The total response time is sum of following (Figure 1.4):
 - i) One RTT to initiate TCP connection (RTT → Round Trip Time).
 - ii) One RTT for HTTP request and first few bytes of HTTP response to return.
 - iii) File transmission time.
 i.e. Total response time = (i) + (ii) + (iii) = 1 RTT+ 1 RTT+ File transmission time
= 2(RTT) + File transmission time

"To hell with circumstances; I create opportunities." -Bruce Lee



COMPUTER NETWORKS

1.2.2.2 HTTP with Persistent Connections

- Problem with Non-Persistent Connections:
 - 1) A new connection must be established and maintained for each requested-object.
 - Hence, buffers must be allocated and state info must be kept in both the client and server.
 - This results in a significant burden on the server.
 - 2) Each object suffers a delivery delay of two RTTs:
 - i) One RTT to establish the TCP connection and
 - ii) One RTT to request and receive an object.
- Solution: Use persistent connections.
- With persistent connections, the server leaves the TCP connection open after sending responses.
- Hence, subsequent requests & responses b/w same client & server can be sent over same connection
- The server closes the connection only when the connection is not used for a certain amount of time.
- Default mode of HTTP: Persistent connections with pipelining.
- Advantages:
 - 1) This method requires only one RTT for all the referenced-objects.
 - 2) The performance is improved by 20%.



COMPUTER NETWORKS

1.2.3 HTTP Message Format

- Two types of HTTP messages: 1) Request-message and 2) Response-message.

1.2.3.1 HTTP Request Message

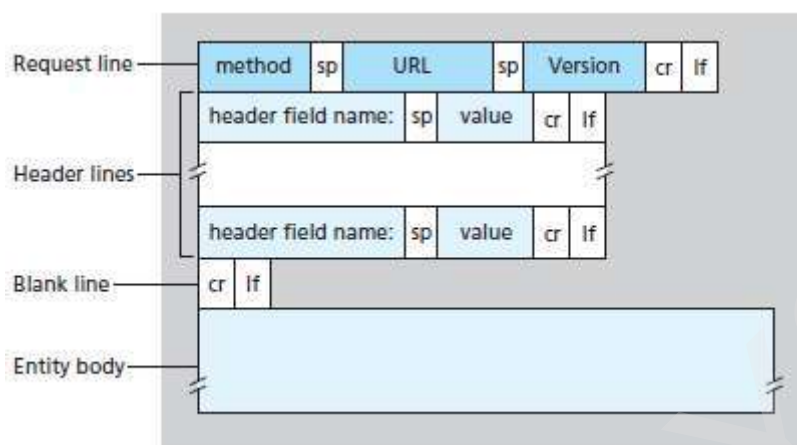


Figure 1.5: General format of an HTTP request-message

- An example of request-message is as follows:

```
GET /somedir/page.html HTTP/1.1
Host: www.someschool.edu
Connection: close
User-agent: Mozilla/5.0
Accept-language: eng
```

- The request-message contains 3 sections (Figure 1.5):
 - 1) Request-line
 - 2) Header-line and
 - 3) Carriage return.
- The first line of message is called the request-line. The subsequent lines are called the header-lines.
- The request-line contains 3 fields. The meaning of the fields is as follows:
 - 1) Method
 - "GET": This method is used when the browser requests an object from the server.
 - 2) URL
 - "/somedir/page.html": This is the object requested by the browser.
 - 3) Version
 - "HTTP/1.1": This is version used by the browser.
- The request-message contains 4 header-lines. The meaning of the header-lines is as follows:
 - 1) "Host: www.someschool.edu" specifies the host on which the object resides.
 - 2) "Connection: close" means requesting a non-persistent connection.
 - 3) "User-agent: Mozilla/5.0" means the browser used is the Firefox.
 - 4) "Accept-language: eng" means English is the preferred language.
- The method field can take following values: GET, POST, HEAD, PUT and DELETE.
 - 1) GET is used when the browser requests an object from the server.
 - 2) POST is used when the user fills out a form & sends to the server.
 - 3) HEAD is identical to GET except the server must not return a message-body in the response.
 - 4) PUT is used to upload objects to servers.
 - 5) DELETE allows an application to delete an object on a server.



COMPUTER NETWORKS

1.2.3.2 HTTP Response Message

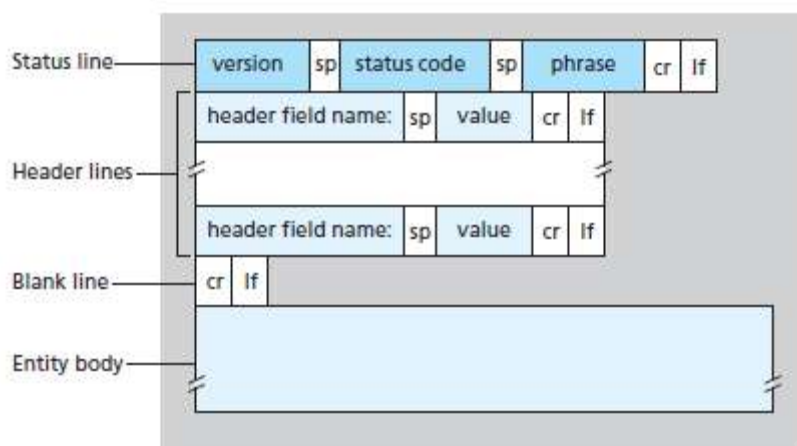


Figure 1.6: General format of an HTTP response-message

- An example of response-message is as follows:

```
HTTP/1.1 200 OK
Connection: close
Date: Tue, 09 Aug 2011 15:44:04 GMT
Server: Apache/2.2.3 (CentOS)
Last-Modified: Tue, 09 Aug 2011 15:11:03 GMT
Content-Length: 6821
Content-Type: text/html
(data data data data data ...)
```

- The response-message contains 3 sections (Figure 1.6):
 - 1) Status line
 - 2) Header-lines and
 - 3) Data (Entity body).
- The status line contains 3 fields:
 - 1) Protocol version
 - 2) Status-code and
 - 3) Status message.
- Some common status-codes and associated messages include:
 - 1) 200 OK: Standard response for successful HTTP requests.
 - 2) 400 Bad Request: The server cannot process the request due to a client error.
 - 3) 404 Not Found: The requested resource cannot be found.
- The meaning of the Status line is as follows:

"HTTP/1.1 200 OK": This line indicates the server is using HTTP/1.1 & that everything is OK.
- The response-message contains 6 header-lines. The meaning of the header-lines is as follows:
 - 1) Connection: This line indicates browser requesting a non-persistent connection.
 - 2) Date: This line indicates the time & date when the response was sent by the server.
 - 3) Server: This line indicates that the message was generated by an Apache Web-server.
 - 4) Last-Modified: This line indicates the time & date when the object was last modified.
 - 5) Content-Length: This line indicates the number of bytes in the sent-object.
 - 6) Content-Type: This line indicates that the object in the entity body is HTML text.



COMPUTER NETWORKS

1.2.4 User-Server Interaction: Cookies

- Cookies refer to a small text file created by a Web-site that is stored in the user's computer.
- Cookies are stored either temporarily for that session only or permanently on the hard disk.
- Cookies allow Web-sites to keep track of users.
- Cookie technology has four components:
 - 1) A cookie header-line in the HTTP response-message.
 - 2) A cookie header-line in the HTTP request-message.
 - 3) A cookie file kept on the user's end-system and managed by the user's browser.
 - 4) A back-end database at the Web-site.

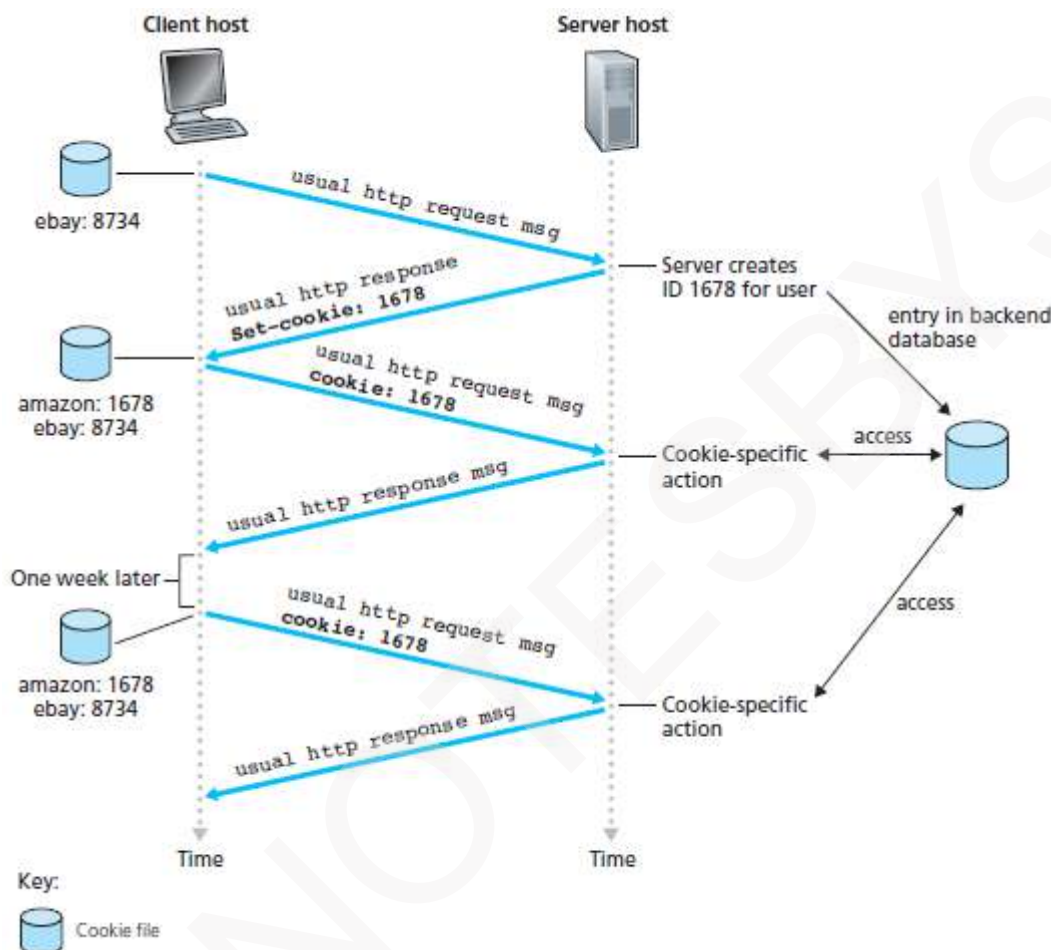


Figure 1.7: Keeping user state with cookies

- Here is how it works (Figure 1.7):
 - 1) When a user first time visits a site, the server
 - creates a unique identification number (1678) and
 - creates an entry in its back-end database by the identification number.
 - 2) The server then responds to user's browser.
 - HTTP response includes `Set-cookie:` header which contains the identification number (1678)
 - 3) The browser then stores the identification number into the cookie-file.
 - 4) Each time the user requests a Web-page, the browser
 - extracts the identification number from the cookie file, and
 - puts the identification number in the HTTP request.
 - 5) In this manner, the server is able to track user's activity at the web-site.



COMPUTER NETWORKS

1.2.5 Web Caching

- A Web-cache is a network entity that satisfies HTTP requests on the behalf of an original Web-server.
- The Web-cache has disk-storage.
- The disk-storage contains copies of recently requested-objects.

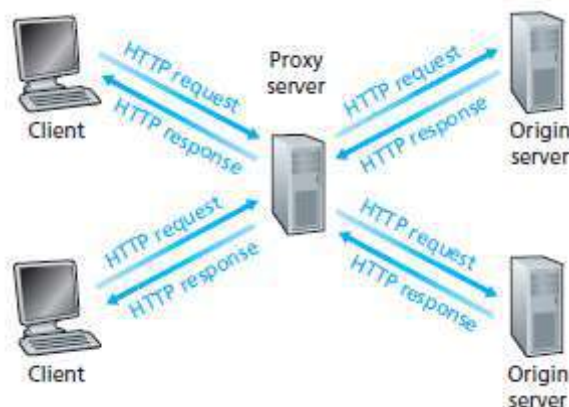


Figure 1.8: Clients requesting objects through a Web-cache (or Proxy Server)

- Here is how it works (Figure 1.8):
 - 1) The user's HTTP requests are first directed to the web-cache.
 - 2) If the cache has the object requested, the cache returns the requested-object to the client.
 - 3) If the cache does not have the requested-object, then the cache
 - connects to the original server and
 - asks for the object.
 - 4) When the cache receives the object, the cache
 - stores a copy of the object in local-storage and
 - sends a copy of the object to the client.
- A cache acts as both a server and a client at the same time.
 - 1) The cache acts as a server when the cache
 - receives requests from a browser and
 - sends responses to the browser.
 - 2) The cache acts as a client when the cache
 - requests to an original server and
 - receives responses from the origin server.
- Advantages of caching:
 - 1) To reduce response-time for client-request.
 - 2) To reduce traffic on an institution's access-link to the Internet.
 - 3) To reduce Web-traffic in the Internet.

1.2.6 The Conditional GET

- Conditional GET refers a mechanism that allows a cache to verify that the objects are up to date.
- An HTTP request-message is called conditional GET if
 - 1) Request-message uses the GET method and
 - 2) Request-message includes an If-Modified-Since: header-line.
- The following is an example of using conditional GET:

```

GET /fruit/kiwi.fig HTTP1.1
Host: www.exoriguecuisine.com
If-modified-since: Wed, 7 Sep 2011 09:23:24
  
```

- The response is:

```

HTTP/1.1 304 Not Modified
Date: Sat, 15 Oct 2011 15:39:29
  
```



COMPUTER NETWORKS

1.3 File Transfer: FTP

- FTP is used by the local host to transfer files to or from a remote-host over the network.
- FTP uses client-server architecture (Figure 1.9).
- FTP uses 2 parallel TCP connections (Figure 1.10):

1) Control Connection

- The control-connection is used for sending control-information b/w local and remote-hosts.
- The control-information includes:
 - user identification
 - password
 - commands to change directory and
 - commands to put & get files.

2) Data Connection

- The data-connection is used to transfer files.

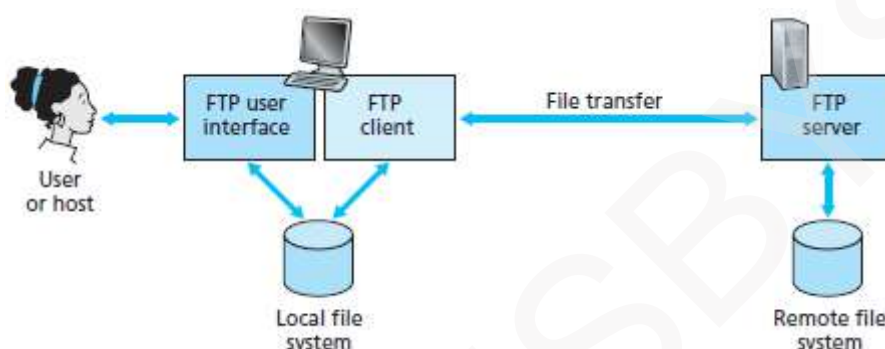


Figure 1.9: FTP moves files between local and remote file systems



Figure 1.10: Control and data-connections

- Here is how it works:
 - 1) When session starts, the client initiates a control-connection with the server on port 21.
 - 2) The client sends user-identity and password over the control-connection.
 - 3) Then, the server initiates data-connection to the client on port 20.
 - 4) FTP sends exactly one file over the data-connection and then closes the data-connection.
 - 5) Usually, the control-connection remains open throughout the duration of the user-session.
 - 6) But, a new data-connection is created for each file transferred within a session.
- During a session, the server must maintain the state-information about the user.
- For example:
 - The server must keep track of the user's current directory.
- Disadvantage:
 - Keeping track of state-info limits the no. of sessions maintained simultaneously by a server.



COMPUTER NETWORKS

1.3.1 FTP Commands & Replies

- The commands are sent from client to server.
- The replies are sent from server to client.
- The commands and replies are sent across the control-connection in 7-bit ASCII format.
- Each command consists of 4-uppercase ASCII characters followed by optional arguments.
- For example:
 - 1) USER username
 - Used to send the user identification to the server.
 - 2) PASS password
 - Used to send the user password to the server.
 - 3) LIST
 - Used to ask the server to send back a list of all the files in the current remote directory.
 - 4) RETR filename
 - Used to retrieve a file from the current directory of the remote-host.
 - 5) STOR filename
 - Used to store a file into the current directory of the remote-host.
- Each reply consists of 3-digit numbers followed by optional message.
- For example:
 - 1) 331 Username OK, password required
 - 2) 125 Data-connection already open; transfer starting
 - 3) 425 Can't open data-connection
 - 4) 452 Error writing file



COMPUTER NETWORKS

1.4 Electronic Mail in the Internet

- e-mail is an asynchronous communication medium in which people send and read messages.
- e-mail is fast, easy to distribute, and inexpensive.
- e-mail has features such as
 - messages with attachments
 - hyperlinks
 - HTML-formatted text and
 - embedded photos.
- Three major components of an e-mail system (Figure 1.11):
 - 1) User Agents**
 - User-agents allow users to read, reply to, forward, save and compose messages.
 - For example: Microsoft Outlook and Apple Mail
 - 2) Mail Servers**
 - Mail-servers contain mailboxes for users.
 - A message is first sent to the sender's mail-server.
 - Then, the sender's mail-server sends the message to the receiver's mail-server.
 - If the sender's server cannot deliver mail to receiver's server, the sender's server
 - holds the message in a message queue and
 - attempts to transfer the message later.
 - 3) SMTP (Simple Mail Transfer Protocol)**
 - SMTP is an application-layer protocol used for email.
 - SMTP uses TCP to transfer mail from the sender's mail-server to the recipient's mail-server.
 - SMTP has two sides:
 - 1) A client-side, which executes on the sender's mail-server.
 - 2) A server-side, which executes on the recipient's mail-server.
 - Both the client and server-sides of SMTP run on every mail-server.
 - When a mail-server receives mail from other mail-servers, the mail-server acts as a server. When a mail-server sends mail to other mail-servers, the mail-server acts as a client.

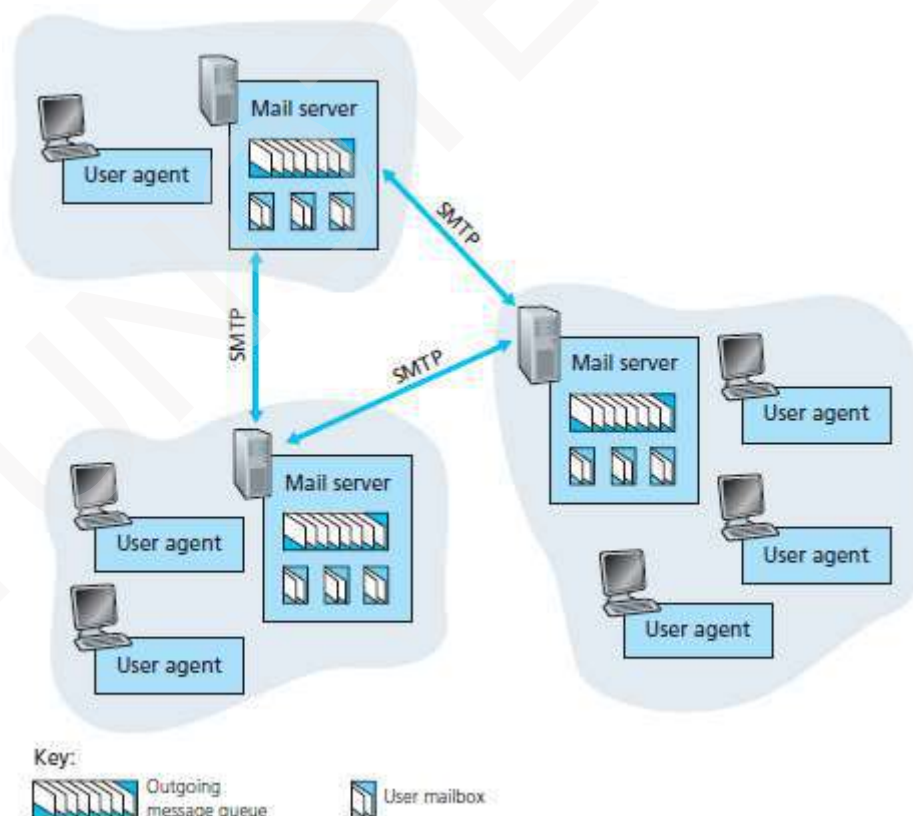


Figure 1.11: A high-level view of the Internet e-mail system



COMPUTER NETWORKS

1.4.1 SMTP

- SMTP is the most important protocol of the email system.
- Three characteristics of SMTP (that differs from other applications):
 - 1) Message body uses 7-bit ASCII code only.
 - 2) Normally, no intermediate mail-servers used for sending mail.
 - 3) Mail transmissions across multiple networks through mail relaying.
- Here is how it works:
 - 1) Usually, mail-servers are listening at port 25.
 - 2) The sending server initiates a TCP connection to the receiving mail-server.
 - 3) If the receiver's server is down, the sending server will try later.
 - 4) If connection is established, the client & the server perform application-layer handshaking.
 - 5) Then, the client indicates the e-mail address of the sender and the recipient.
 - 6) Finally, the client sends the message to the server over the same TCP connection.

1.4.2 Comparison of SMTP with HTTP

- 1) HTTP is mainly a pull protocol. This is because
 - someone loads information on a web-server and
 - users use HTTP to pull the information from the server.
- On the other hand, SMTP is primarily a push protocol. This is because
 - the sending mail-server pushes the file to receiving mail-server.
- 2) SMTP requires each message to be in seven-bit ASCII format.
 - If message contains binary-data, the message has to be encoded into 7-bit ASCII format.
 - HTTP does not have this restriction.
- 3) HTTP encapsulates each object of message in its own response-message.
 - SMTP places all of the message's objects into one message.

1.4.3 Mail Access Protocols

- It is not realistic to run the mail-servers on PC & laptop. This is because
 - mail-servers must be always-on and
 - mail-servers must have fixed IP addresses
- Problem: How a person can access the email using PC or laptop?
- Solution: Use mail access protocols.
- Three mail access protocols:
 - 1) Post Office Protocol (POP)
 - 2) Internet Mail Access Protocol (IMAP) and
 - 3) HTTP.



COMPUTER NETWORKS

1.4.3.1 POP

- POP is an extremely simple mail access protocol.
- POP server will listen at port 110.
- Here is how it works:
 - The user-agent at client's computer opens a TCP connection to the main server.
 - POP then progresses through three phases:
 - 1) Authentication**
 - The user-agent sends a user name and password to authenticate the user.
 - 2) Transaction**
 - The user-agent retrieves messages.
 - Also, the user-agent can
 - mark messages for deletion
 - remove deletion marks &
 - obtain mail statistics.
 - The user-agent issues commands, and the server responds to each command with a reply.
 - There are two responses:
 - i) +OK: used by the server to indicate that the previous command was fine.
 - ii) -ERR: used by the server to indicate that something is wrong.
 - 3) Update**
 - After user issues a quit command, the mail-server removes all messages marked for deletion.
- Disadvantage:
 - The user cannot manage the mails at remote mail-server. For ex: user cannot delete messages.

1.4.3.2 IMAP

- IMAP is another mail access protocol, which has more features than POP.
- An IMAP server will associate each message with a folder.
- When a message first arrives at server, the message is associated with recipient's INBOX folder
- Then, the recipient can
 - move the message into a new, user-created folder
 - read the message
 - delete the message and
 - search remote folders for messages matching specific criteria.
- An IMAP server maintains user state-information across IMAP sessions.
- IMAP permits a user-agent to obtain components of messages.
 - For example, a user-agent can obtain just the message header of a message.

1.4.3.3 Web-Based E-Mail

- HTTPs are now used for Web-based email accessing.
- The user-agent is an ordinary Web browser.
- The user communicates with its remote-server via HTTP.
- Now, Web-based emails are provided by many companies including Google, Yahoo etc.



COMPUTER NETWORKS

1.5 DNS – The Internet's Directory Service

- DNS is an internet service that translates domain-names into IP addresses.
For ex: the domain-name "www.google.com" might translate to IP address "198.105.232.4".
- Because domain-names are alphabetic, they are easier to remember for human being.
- But, the Internet is really based on IP addresses (DNS → Domain Name System).

1.5.1 Services Provided by DNS

- The DNS is
 - 1) A distributed database implemented in a hierarchy of DNS servers.
 - 2) An application-layer protocol that allows hosts to query the distributed database.
- DNS servers are often UNIX machines running the BIND software.
- The DNS protocol runs over UDP and uses port 53. (BIND → Berkeley Internet Name Domain)
- DNS is used by application-layer protocols such as HTTP, SMTP, and FTP.
- Assume a browser requests the URL www.someschool.edu/index.html.
- Next, the user's host must first obtain the IP address of www.someschool.edu
- This is done as follows:
 - 1) The same user machine runs the client-side of the DNS application.
 - 2) The browser
 - extracts the hostname "www.someschool.edu" from the URL and
 - passes the hostname to the client-side of the DNS application.
 - 3) The client sends a query containing the hostname to a DNS server.
 - 4) The client eventually receives a reply, which includes the IP address for the hostname.
 - 5) After receiving the IP address, the browser can initiate a TCP connection to the HTTP server.
- DNS also provides following services:
 - 1) Host Aliasing**
 - A host with a complicated hostname can have one or more alias names.
 - 2) Mail Server Aliasing**
 - For obvious reasons, it is highly desirable that e-mail addresses be mnemonic.
 - 3) Load Distribution**
 - DNS is also used to perform load distribution among replicated servers.
 - Busy sites are replicated over multiple servers & each server runs on a different system.

1.5.2 Overview of How DNS Works

- Distributed database design is more preferred over centralized design because:
 - 1) A Single Point of Failure**
 - If the DNS server crashes then the entire Internet will not stop.
 - 2) Traffic Volume**
 - A Single DNS Server cannot handle the huge global DNS traffic.
 - But with distributed system, the traffic is distributed and reduces overload on server.
 - 3) Distant Centralized Database**
 - A single DNS server cannot be "close to" all the querying clients.
 - If we put the single DNS server in Mysore,
then all queries from USA must travel to the other side of the globe.
 - This can lead to significant delays.
 - 4) Maintenance**
 - The single DNS server would have to keep records for all Internet hosts.
 - This centralized database has to be updated frequently to account for every new host.



COMPUTER NETWORKS

1.5.2.1 A Distributed, Hierarchical Database

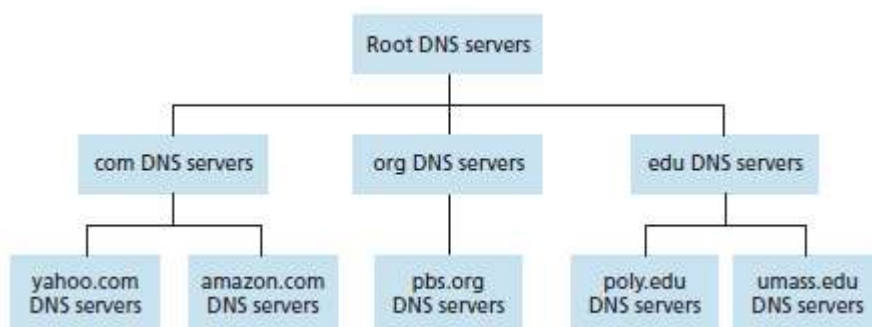


Figure 1.12: Portion of the hierarchy of DNS servers

- Suppose a client wants to determine IP address for hostname "www.amazon.com" (Figure 1.12):
 - 1) The client first contacts one of the root servers, which returns IP addresses for TLD servers.
 - 2) Then, the client contacts one of these TLD servers.
 - The TLD server returns the IP address of an authoritative-server for "amazon.com".
 - 3) Finally, the client contacts one of the authoritative-servers for amazon.com.
 - The authoritative-server returns the IP address for the hostname "www.amazon.com".

1.5.2.1.1 Recursive Queries & Iterative Queries

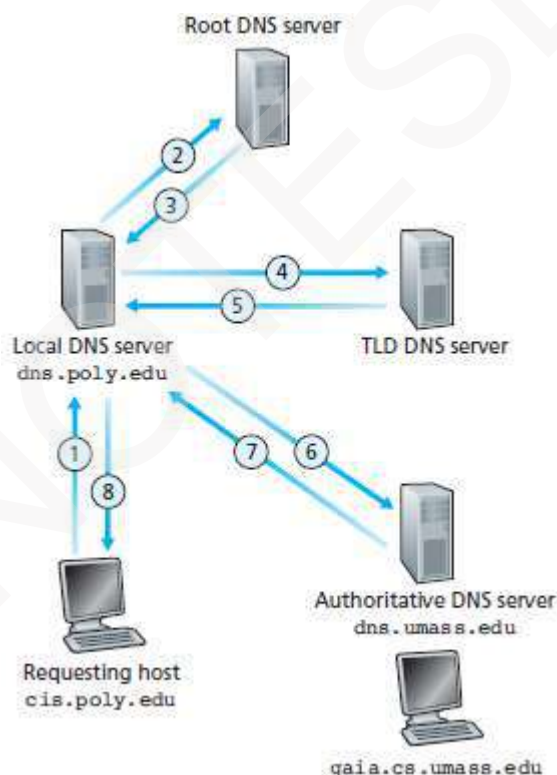


Figure 1.13: Interaction of the various DNS servers

- The example shown in Figure 1.13 makes use of both recursive queries and iterative queries.
- The query 1 sent from cis.poly.edu to dns.poly.edu is a recursive query. This is because
 - the query asks dns.poly.edu to obtain the mapping on its behalf.
- But the subsequent three queries 2, 4 and 6 are iterative. This is because
 - all replies are directly returned to dns.poly.edu.



COMPUTER NETWORKS

1.5.3 DNS Records & Messages

- The DNS server stores resource-records (RRs).
- RRs provide hostname-to-IP address mappings.
- Each DNS reply message carries one or more resource-records.
- A resource-record is a 4-tuple that contains the following fields: (Name, Value, Type, TTL)
- TTL (time to live) determines when a resource should be removed from a cache.
- The meaning of Name and Value depend on Type:
 - 1) If Type=A, then Name is a hostname and Value is the IP address for the hostname.
 - Thus, a Type A record provides the standard hostname-to-IP address mapping.
For ex: (relay1.bar.foo.com, 145.37.93.126, A)
 - 2) If Type=NS, then
 - i) Name is a domain (such as foo.com) and
 - ii) Value is the hostname of an authoritative DNS server.
 - This record is used to route DNS queries further along in the query chain.
For ex: (foo.com, dns.foo.com, NS) is a Type NS record.
 - 3) If Type=CNAME, then Value is a canonical hostname for the alias hostname Name.
 - This record can provide querying hosts the canonical name for a hostname.
For ex: (foo.com, relay1.bar.foo.com, CNAME) is a CNAME record.
 - 4) If Type=MX, Value is the canonical name of a mail-server that has an alias hostname Name.
 - MX records allow the hostnames of mail-servers to have simple aliases.
For ex: (foo.com, mail.bar.foo.com, MX) is an MX record.



COMPUTER NETWORKS

1.5.3.1 DNS Messages

- Two types of DNS messages: 1) query and 2) reply.
- Both query and reply messages have the same format.

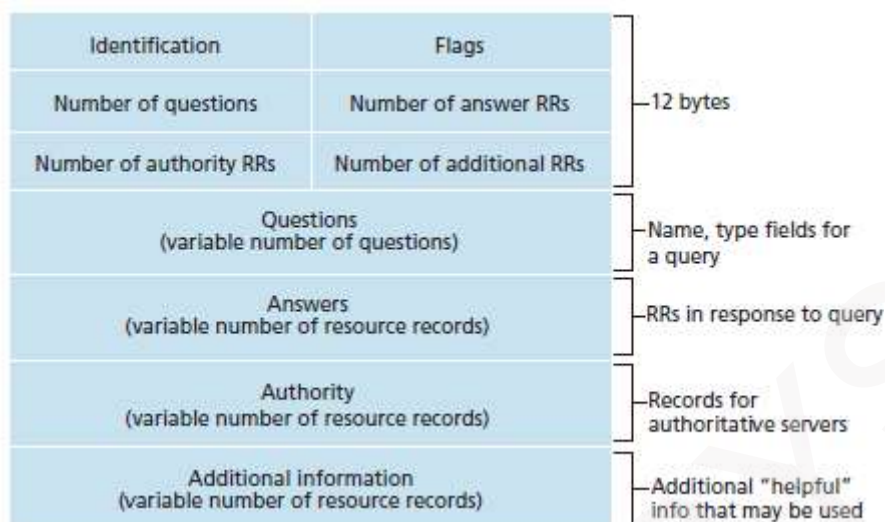


Figure 1.14: DNS message format

- The various fields in a DNS message are as follows (Figure 1.14):

1) Header Section

- The first 12 bytes is the header-section.
- This section has following fields:

i) Identification

- This field identifies the query.
- This identifier is copied into the reply message to a query.
- This identifier allows the client to match received replies with sent queries.

ii) Flag

- This field has following 3 flag-bits:

a) Query/Reply

- ✗ This flag-bit indicates whether the message is a query (0) or a reply (1).

b) Authoritative

- ✗ This flag-bit is set in a reply message when a DNS server is an authoritative-server.

c) Recursion Desired

- ✗ This flag-bit is set when a client desires that the DNS server perform recursion.

iii) Four Number-of-Fields

- These fields indicate the no. of occurrences of 4 types of data sections that follow the header.

2) Question Section

- This section contains information about the query that is being made.
- This section has following fields:

i) Name

- This field contains the domain-name that is being queried.

ii) Type

- This field indicates the type of question being asked about the domain-name.

3) Answer Section

- This section contains a reply from a DNS server.
- This section contains the resource-records for the name that was originally queried.
- A reply can return multiple RRs in the answer, since a hostname can have multiple IP addresses.

4) Authority Section

- This section contains records of other authoritative-servers.

5) Additional Section

- This section contains other helpful records.



COMPUTER NETWORKS

1.6 Peer-to-Peer Applications

- Peer-to-peer architecture is different from client-server architecture.
- In P2P, each node (called peers) acts as a client and server at the same time.
- The peers are not owned by a service-provider.
- The peers not supposed to be always listening on the Internet.
- The peers are dynamic, i.e., some peers will join some peers will leave from time to time.

1.6.1 P2P File Distribution

- One popular P2P file distribution protocol is BitTorrent
- Consider the following scenarios:

Suppose a server has a large file and 'N' computers want to download the file (Figure 1.15).

- 1) In client-server architecture, each of the N computers will
 - connect to the server &
 - download a copy of the file to local-host.
- 2) In P2P architecture, a peer need not necessarily download a copy from the server.
 - Rather, the peer may download from other peers.

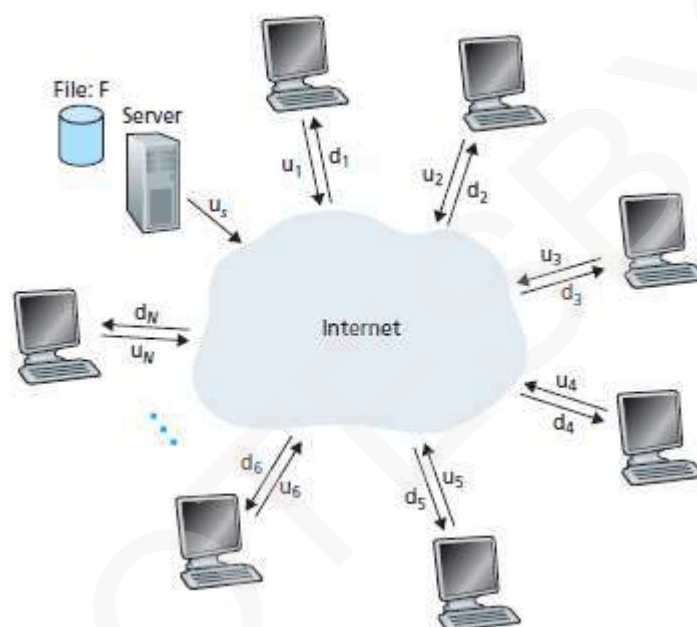


Figure 1.15: An illustrative file distribution problem

- Let us define the following:
 - Length of the file = F
 - Upload rate for the server = u_s
 - Upload rate for i th computer = u_i
 - Download-rate for i th computer = d_i
 - Distribution-time for the client-server architecture = D_{CS}
 - Server needs transmit = N_F bits.
 - Lowest download-rate of peer = d_{min}
 - Distribution-time for the P2P architecture = D_{P2P}

Case 1: Client-Server Architecture

- We have 2 observations:
 - 1) The server must transmit one copy of the file to each of the N peers.
 - Since the server's upload rate is u_s , the distribution-time is at least N_F/u_s .
 - 2) The peer with the lowest download-rate cannot obtain all F bits of the file in less than F/d_{min} .
 - Thus, the minimum distribution-time is at least F/d_{min} .
- Putting above 2 observations together, we have

$$D_{CS} = \max \left\{ \frac{NF}{u_s}, \frac{F}{d_{min}} \right\}$$



COMPUTER NETWORKS

Case 2: P2P Architecture

- We have 2 observations:
 - 1) At the beginning of the distribution, only the server has the file.
 - So, the minimum distribution-time is at least F/u_s .
 - 2) The peer with the lowest download-rate cannot obtain all F bits of the file in less than F/d_{\min} .
 - Thus, the minimum distribution-time is at least F/d_{\min} .
 - 3) The total upload capacity of the system as a whole is $u_{\text{total}} = u_s + u_1 + u_2 \dots + u_N$.
 - The system must deliver F bits to each of the N peers.
 - Thus, the minimum distribution-time is at least $NF/(u_s + u_1 + u_2 \dots + u_N)$.
- Putting above 3 observations together, we have

$$D_{\text{P2P}} \geq \max \left\{ \frac{F}{u_s}, \frac{F}{d_{\min}}, \frac{NF}{u_s + \sum_{i=1}^N u_i} \right\}$$

- Figure 1.16 compares the minimum distribution-time for the client-server and P2P architectures.

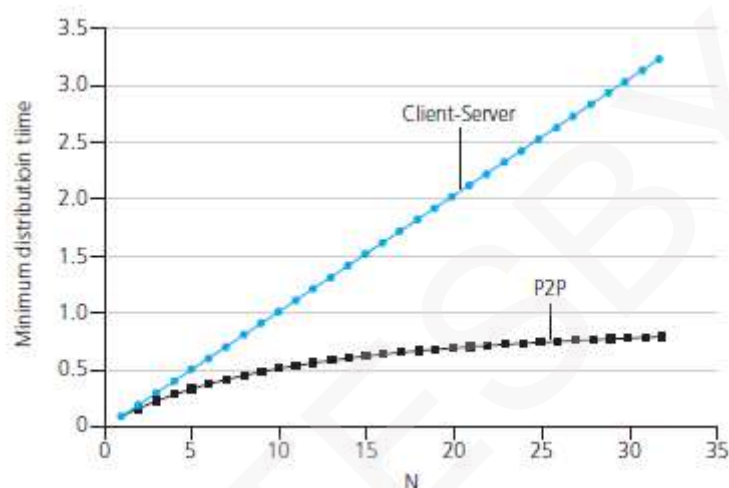


Figure 1.16: Distribution-time for P2P and client-server architectures

- The above comparison shows that when N is large,
 - i) P2P architecture is consuming less distribution-time.
 - ii) P2P architecture is self-scaling.



COMPUTER NETWORKS

1.6.1.1 BitTorrent

- The collection of all peers participating in the distribution of a particular file is called a torrent.
- Peers download equal-size chunks of the file from one another. Chunk size = 256 KBytes.
- The peer also uploads chunks to other peers.
- Once a peer has acquired the entire file, the peer may leave the torrent or remain in the torrent.
- Each torrent has an infrastructure node called tracker.
- Here is how it works (Figure 1.17):
 - 1) When a peer joins a torrent, the peer
 - registers itself with the tracker and
 - periodically informs the tracker that it is in the torrent.
 - 2) When a new peer joins the torrent, the tracker
 - randomly selects a subset of peers from the set of participating peers and
 - sends the IP addresses of these peers to the new peer.
 - 3) Then, the new peer tries to establish concurrent TCP connections with all peers on this list.
 - All peers on the list are called neighboring-peers.
 - 4) Periodically, the new peer will ask each of the neighboring-peers for the set of chunks.
- To choose the chunks to download, the peer uses a technique called rarest-first.
- Main idea of rarest-first:
 - Determine the chunks that are the rarest among the neighbors and
 - Request then those rarest chunks first.

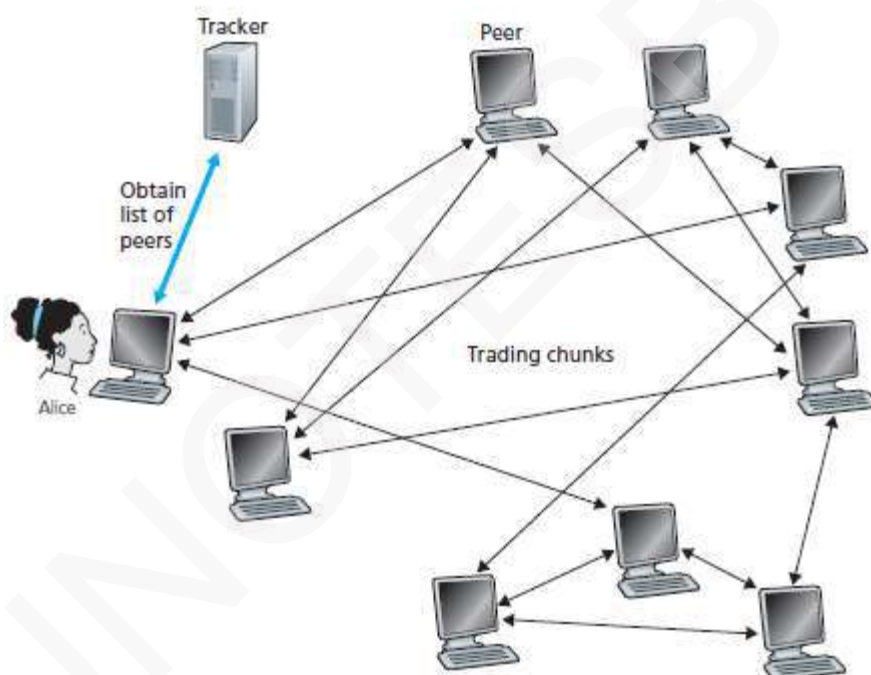


Figure 1.17: File distribution with BitTorrent



COMPUTER NETWORKS

1.6.2 Distributed Hash Table

- We can use P2P architecture to form a distributed database.
- We consider a simple database, which contains (key, value) pairs.
- Each peer will only hold a small subset of the data.
- Any peer can query the distributed database with a particular key.
- Then, the database will
 - locate the peers that have the corresponding (key, value) pairs and
 - return the key-value to the querying peer.
- Any peer will also be allowed to insert new key-value pairs into the database.
- Such a distributed database is referred to as a distributed hash table (DHT).
- To construct the database, we need to use some hash-functions.
- The input of a hash-function can be a large number.
 - But the output of the hash-function is of fixed-size bit-string.
- The outline of building a DHT is as follows:
 - 1) Assign an identifier to each peer, where the identifier is an n-bit string.
 - So, we can view the identifier as an integer at the range from 0 to $2^n - 1$.
 - 2) For a data pair (key, value), the hash value of the key is computed.
 - Then the data is stored in the peer whose identifier is closest to the key.
 - 3) To insert or retrieve data, first we need to find the appropriate peer.
 - Problem: It is not realistic to let the peer to store all of the other peer's identifiers.
 - Solution: Use a circular arrangement.

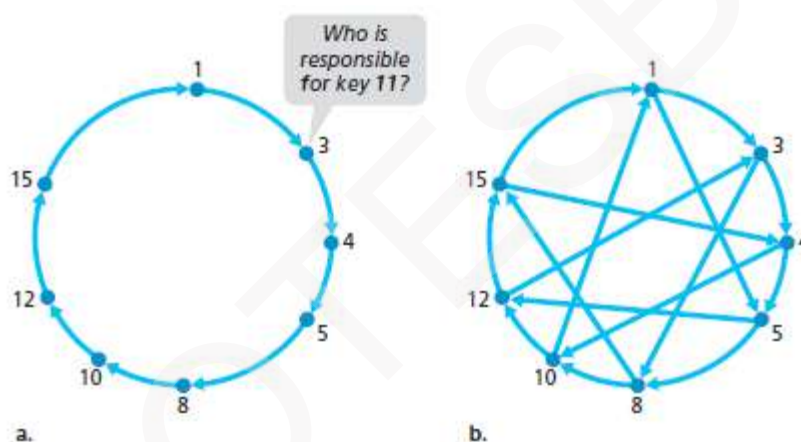


Figure 1.18: (a) A circular DHT. Peer 3 wants to determine who is responsible for key 11. (b) A circular DHT with shortcuts

1.6.2.1 Circular Arrangement

- In this example, the identifiers are range $[0, 15]$ (4-bit strings) and there are eight peers.
- Each peer is only aware of its immediate successor and predecessor (Figure 1.18a).
- So, each peer just needs to track two neighbors.
- When a peer asks for a key, the peer sends the message clockwise around the circle.
- For example:
 - The peer 4 sends a message saying "Who is responsible for key 11?"
 - The message will forward through peers 5, 8, 10 and reach peer 12.
 - The peer 12 determines that it is the closest peer to key 11.
 - At this point, peer 12 sends a message back to the querying peer 4.
- But when the circle is too large, a message may go through a large number of peers to get answer.
- Problem: There is trade off between
 - i) Number of neighbors each peer has to track and
 - ii) Number of message to be sent to resolve a single query.
- Solution: To reduce the query time, add "shortcuts" to the circular arrangement (Figure 1.18b).



COMPUTER NETWORKS

1.7 Socket Programming: Creating Network Applications

- Two types of network-applications:
 - 1) First type is an implementation whose operation is specified in a protocol standard (RFC)
 - Such an application is referred to as "open".
 - The client & server programs must conform to the rules dictated by the RFC.
 - 2) Second type is a proprietary network-application.
 - The client & server programs use an application-layer protocol not openly published in a RFC.
 - A single developer creates both the client and server programs.
 - The developer has complete control over the code.
- During development phase, the developer must decide whether the application uses TCP or UDP.



COMPUTER NETWORKS

1.7.1 Socket Programming with UDP

- Consider client-server application in which following events occurs:
 - The client
 - reads a line of characters (data) from the keyboard and
 - sends the data to the server.
 - The server receives the data and converts the characters to uppercase.
 - The server sends the modified data to the client.
 - The client receives the modified data and displays the line on its screen.

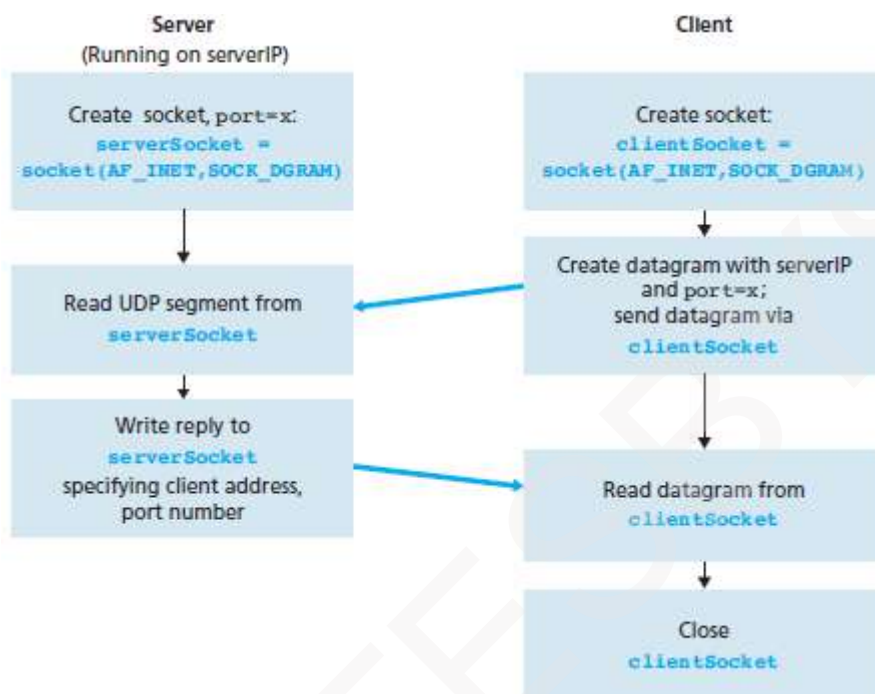


Figure 1.19: The client-server application using UDP

- The client-side of the application is as follows (Figure 1.19):

```

from socket import * //This line declares socket within the program.
serverName = 'hostname' // This line sets the server name to "hostname".
serverPort = 12000 // This line sets the server port# to "12000".
clientSocket = socket(socket.AF_INET, socket.SOCK_DGRAM) //This line creates the client's socket
message = raw_input('Input lowercase sentence:') //This line reads the data at the client
clientSocket.sendto(message,(serverName, serverPort)) //This line sends data into the socket
modifiedMessage, serverAddress = clientSocket.recvfrom(2048) // This line receives data from socket
print modifiedMessage //This line displays received-data on the screen
clientSocket.close() //This line closes the socket. The process then terminates.
  
```

Here, AF_INET indicates address family

SOCK_DGRAM indicates UDP as socket type contains server's IP address & port#

- The server-side of the application is as follows:

```

from socket import *
serverPort = 12000
serverSocket = socket(AF_INET, SOCK_DGRAM)
serverSocket.bind(('', serverPort)) // This line assigns the port# 12000 to the server's socket.
print "The server is ready to receive"
while 1:
    message, clientAddress = serverSocket.recvfrom(2048)
    modifiedMessage = message.upper() // This line converts data to upper case
    serverSocket.sendto(modifiedMessage, clientAddress)
  
```



COMPUTER NETWORKS

1.7.2 Socket Programming with TCP

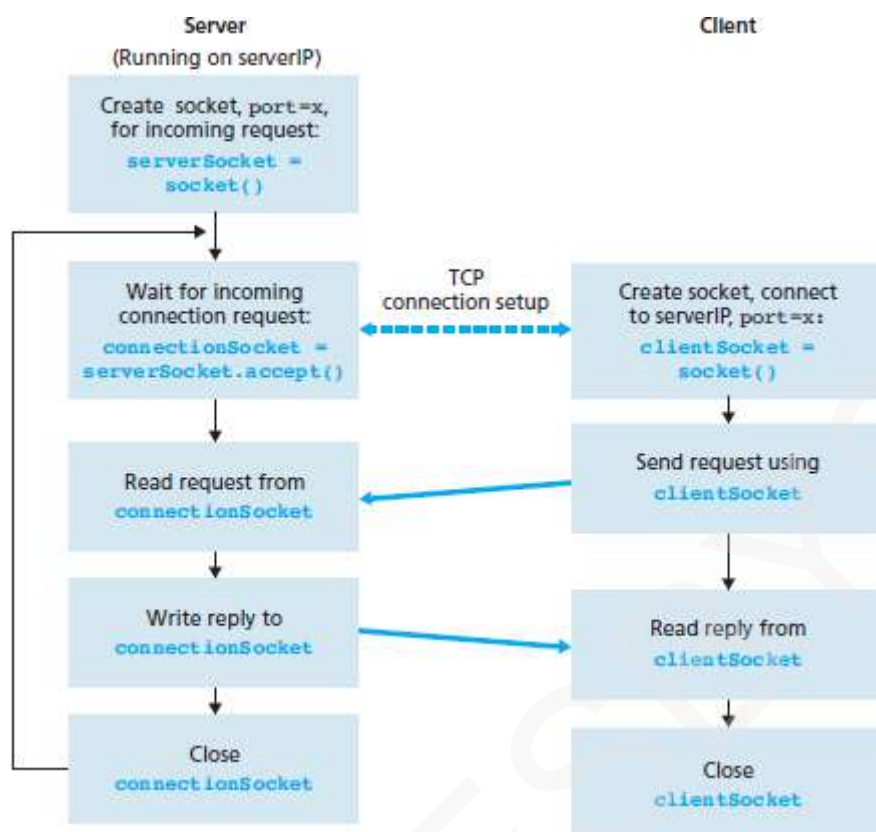


Figure 1.20: The client-server application using TCP

- The client-side of the application is as follows (Figure 1.20):

```

from socket import *
serverName = 'servername'
serverPort = 12000
clientSocket = socket(AF_INET, SOCK_STREAM)
clientSocket.connect((serverName, serverPort)) // This line initiates TCP connection b/w client & server
sentence = raw_input('Input lowercase sentence:')
clientSocket.send(sentence)
modifiedSentence = clientSocket.recv(1024)
print 'From Server:', modifiedSentence
clientSocket.close()

```

- The server-side of the application is as follows:

```

from socket import *
serverPort = 12000
serverSocket = socket(AF_INET, SOCK_STREAM)
serverSocket.bind(('', serverPort))
serverSocket.listen(1) // This line specifies no. of connection-requests from the client to server
print 'The server is ready to receive'
while 1:
    connectionSocket, addr = serverSocket.accept() //allows server to accept connection request from client
    sentence = connectionSocket.recv(1024)
    capitalizedSentence = sentence.upper()
    connectionSocket.send(capitalizedSentence)
    connectionSocket.close()

```

**MODULE-WISE QUESTIONS****PART 1**

- 1) Explain client-server & P2P architecture. (8*)
- 2) With block diagram, explain how application processes communicate through a socket. (8*)
- 3) Explain 4 transport services available to applications. (4)
- 4) Briefly explain 2 transport layer protocols. (4)
- 5) With block diagram, explain the working of Web & HTTP. (8*)
- 6) Explain HTTP non-persistent & persistent connections. (8*)
- 7) With general format, explain HTTP request- & response-messages. (8*)
- 8) With a diagram, explain how cookies are used in user-server interaction. (6*)
- 9) With a diagram, explain the working of web caching. (6*)
- 10) With a diagram, explain the working of FTP. (6*)

PART 2

- 11) With a diagram, explain the working of e-mail system. (6*)
- 12) Briefly explain 3 mail access protocols. (6*)
- 13) Briefly explain the working of DNS. (8*)
- 14) With general format, explain DNS messages. (6*)
- 15) Explain P2P File Distribution. (6)
- 16) With a diagram, explain the working of BitTorrent. (6*)
- 17) With a diagram, explain the working of Distributed Hash Table. (6)
- 18) Draw flow diagram for the client-server application using TCP. Also, write code for the client- & server-sides of the application. (8*)
- 19) Draw flow diagram for the client-server application using UDP. Also, write code for the client- & server-sides of the application. (8)



MODULE 2: TRANSPORT LAYER

- 2.1 Introduction and Transport Layer Services
 - 2.1.1 Relationship between Transport and Network Layers
 - 2.1.2 Overview of the Transport Layer in the Internet
- 2.2 Multiplexing and Demultiplexing
 - 2.2.1 Endpoint Identification
 - 2.2.2 Connectionless Multiplexing and Demultiplexing
 - 2.2.3 Connection Oriented Multiplexing and Demultiplexing
 - 2.2.4 Web Servers and TCP
- 2.3 Connectionless Transport: UDP
 - 2.3.1 UDP Segment Structure
 - 2.3.2 UDP Checksum
- 2.4 Principles of Reliable Data Transfer
 - 2.4.1 Building a Reliable Data Transfer Protocol
 - 2.4.1.1 Reliable Data Transfer over a Perfectly Reliable Channel: rdt1.0
 - 2.4.1.2 Reliable Data Transfer over a Channel with Bit Errors: rdt2.0
 - 2.4.1.2.1 Sender Handles Garbled ACK/NAKs: rdt2.1
 - 2.4.1.2.2 Sender uses ACK/NAKs: rdt2.2
 - 2.4.1.3 Reliable Data Transfer over a Lossy Channel with Bit Errors: rdt3.0
 - 2.4.2 Pipelined Reliable Data Transfer Protocols
 - 2.4.3 Go-Back-N (GBN)
 - 2.4.3.1 GBN Sender
 - 2.4.3.2 GBN Receiver
 - 2.4.3.3 Operation of the GBN Protocol
 - 2.4.4 Selective Repeat (SR)
 - 2.4.4.1 SR Sender
 - 2.4.4.2 SR Receiver
 - 2.4.5 Summary of Reliable Data Transfer Mechanisms and their Use
- 2.5 Connection-Oriented Transport: TCP
 - 2.5.1 The TCP Connection
 - 2.5.2 TCP Segment Structure
 - 2.5.2.1 Sequence Numbers and Acknowledgment Numbers
 - 2.5.2.2 Telnet: A Case Study for Sequence and Acknowledgment Numbers
 - 2.5.3 Round Trip Time Estimation and Timeout
 - 2.5.3.1 Estimating the Round Trip Time
 - 2.5.3.2 Setting and Managing the Retransmission Timeout Interval
 - 2.5.4 Reliable Data Transfer
 - 2.5.4.1 A Few Interesting Scenarios
 - 2.5.4.1.1 First Scenario
 - 2.5.4.1.2 Second Scenario
 - 2.5.4.1.3 Third Scenario
 - 2.5.4.2 Fast Retransmit
 - 2.5.5 Flow Control
 - 2.5.6 TCP Connection Management
 - 2.5.6.1 Connection Setup & Data Transfer
 - 2.5.6.2 Connection Release
- 2.6 Principles of Congestion Control
 - 2.6.1 The Causes and the Costs of Congestion
 - 2.6.1.1 Scenario 1: Two Senders, a Router with Infinite Buffers
 - 2.6.1.2 Scenario 2: Two Senders and a Router with Finite Buffers
 - 2.6.1.3 Scenario 3: Four Senders, Routers with Finite Buffers, and Multihop Paths



COMPUTER NETWORKS

- 2.6.2 Approaches to Congestion Control
- 2.6.3 Network Assisted Congestion Control Example: ATM ABR Congestion Control
 - 2.6.3.1 Three Methods to indicate Congestion
- 2.7 TCP Congestion Control
 - 2.7.1 TCP Congestion Control
 - 2.7.1.1 Slow Start
 - 2.7.1.2 Congestion Avoidance
 - 2.7.1.3 Fast Recovery
 - 2.7.1.4 TCP Congestion Control: Retrospective
 - 2.7.2 Fairness
 - 2.7.2.1 Fairness and UDP
 - 2.7.2.2 Fairness and Parallel TCP Connections



MODULE 2: TRANSPORT LAYER

2.1 Introduction and Transport Layer Services

- A transport-layer protocol provides logical-communication b/w application-processes running on different hosts.
- Transport-layer protocols are implemented in the end-systems but not in network-routers.
- On the sender, the transport-layer
 - receives messages from an application-process
 - converts the messages into the segments and
 - passes the segment to the network-layer.
- On the receiver, the transport-layer
 - receives the segment from the network-layer
 - converts the segments into the messages and
 - passes the messages to the application-process.
- The Internet has 2 transport-layer protocols: TCP and UDP

2.1.1 Relationship between Transport and Network Layers

- A transport-layer protocol provides logical-communication b/w processes running on different hosts. Whereas, a network-layer protocol provides logical-communication between hosts.
- Transport-layer protocols are implemented in the end-systems but not in network-routers.
- Within an end-system, a transport protocol
 - moves messages from application-processes to the network-layer and vice versa.
 - but doesn't say anything about how the messages are moved within the network-core.
- The routers do not recognize any info. which is appended to the messages by the transport-layer.

2.1.2 Overview of the Transport Layer in the Internet

- When designing a network-application, we must choose either TCP or UDP as transport protocol.
 - 1) UDP (User Datagram Protocol)**
 - UDP provides a connectionless service to the invoking application.
 - The UDP provides following 2 services:
 - i) Process-to-process data delivery and
 - ii) Error checking.
 - UDP is an unreliable service i.e. it doesn't guarantee data will arrive to destination-process.
 - 2) TCP (Transmission Control Protocol)**
 - TCP provides a connection-oriented service to the invoking application.
 - The TCP provides following 3 services:
 - 1) Reliable data transfer i.e. guarantees data will arrive to destination-process correctly.
 - 2) Congestion control and
 - 3) Error checking.



COMPUTER NETWORKS

2.2 Multiplexing and Demultiplexing

- A process can have one or more sockets.
- The sockets are used to pass data from the network to the process and vice versa.

1) Multiplexing

- At the sender, the transport-layer
 - gathers data-chunks at the source-host from different sockets
 - encapsulates data-chunk with header to create segments and
 - passes the segments to the network-layer.
- The job of combining the data-chunks from different sockets to create a segment is called multiplexing.

2) Demultiplexing

- At the receiver, the transport-layer
 - examines the fields in the segments to identify the receiving-socket and
 - directs the segment to the receiving-socket.
 - The job of delivering the data in a segment to the correct socket is called demultiplexing.
- In Figure 2.1,
 - In the middle host, the transport-layer must demultiplex segments arriving from the network-layer to either process P1 or P2.
 - The arriving segment's data is directed to the corresponding process's socket.

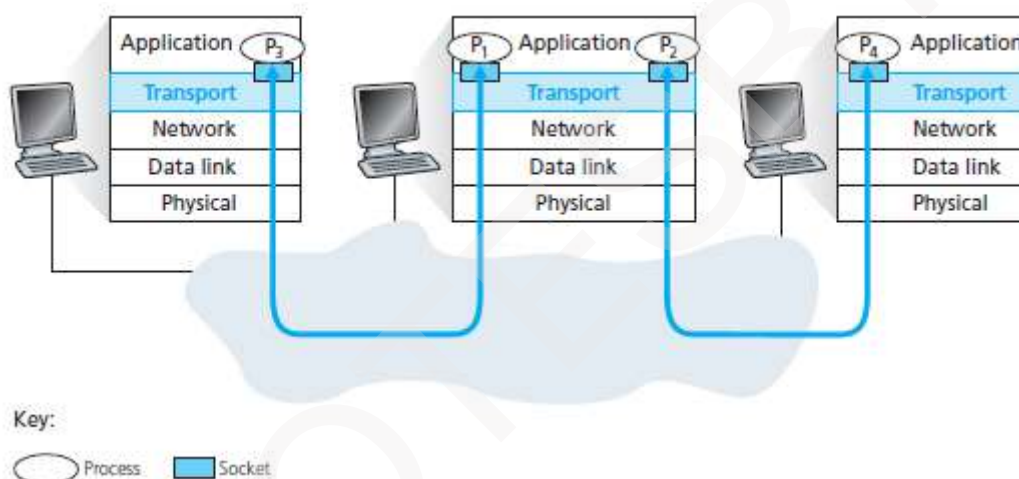


Figure 2.1: Transport-layer multiplexing and demultiplexing



COMPUTER NETWORKS

2.2.1 Endpoint Identification

- Each socket must have a unique identifier.
- Each segment must include 2 header-fields to identify the socket (Figure 2.2):
 - 1) Source-port-number field and
 - 2) Destination-port-number field.
- Each port-number is a 16-bit number: 0 to 65535.
- The port-numbers ranging from 0 to 1023 are called well-known port-numbers and are restricted.
For example: HTTP uses port-no 80
FTP uses port-no 21
- When we develop a new application, we must assign the application a port-number, which are known as ephemeral ports (49152–65535).

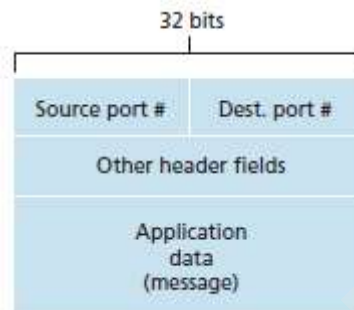


Figure 2.2: Source and destination-port-no fields in a transport-layer segment

- How the transport-layer implements the demultiplexing service?
- Answer:
 - Each socket in the host will be assigned a port-number.
 - When a segment arrives at the host, the transport-layer
 - examines the destination-port-no in the segment
 - directs the segment to the corresponding socket and
 - passes then the segment to the attached process.



COMPUTER NETWORKS

2.2.2 Connectionless Multiplexing and Demultiplexing

- At client side of the application, the transport-layer automatically assigns the port-number. Whereas, at the server side, the application assigns a specific port-number.
- Suppose process on Host-A (port 19157) wants to send data to process on Host-B (port 46428).

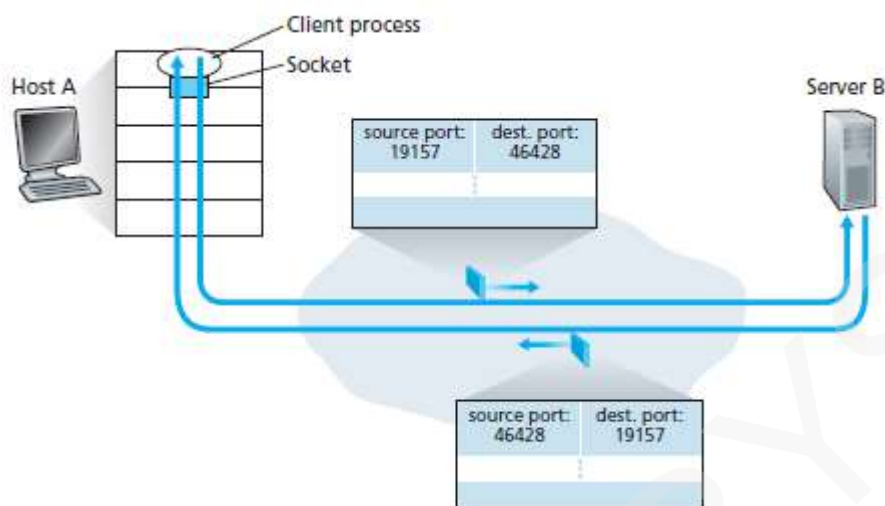


Figure 2.3: The inversion of source and destination-port-nos

- At the sender A, the transport-layer
 - creates a segment containing source-port 19157, destination-port 46428 & data and
 - passes then the resulting segment to the network-layer.
- At the receiver B, the transport-layer
 - examines the destination-port field in the segment and
 - delivers the segment to the socket identified by port 46428.
- A UDP socket is identified by a two-tuple:
 - 1) Destination IP address &
 - 2) Destination-port-no.
- As shown in Figure 2.3, Source-port-no from Host-A is used at Host-B as "return address" i.e. when B wants to send a segment back to A.



COMPUTER NETWORKS

2.2.3 Connection Oriented Multiplexing and Demultiplexing

- Each TCP connection has exactly 2 end-points. (Figure 2.4).
- Thus, 2 arriving TCP segments with different source-port-nos will be directed to 2 different sockets, even if they have the same destination-port-no.
- A TCP socket is identified by a four-tuple:
 - 1) Source IP address
 - 2) Source-port-no
 - 3) Destination IP address &
 - 4) Destination-port-no.

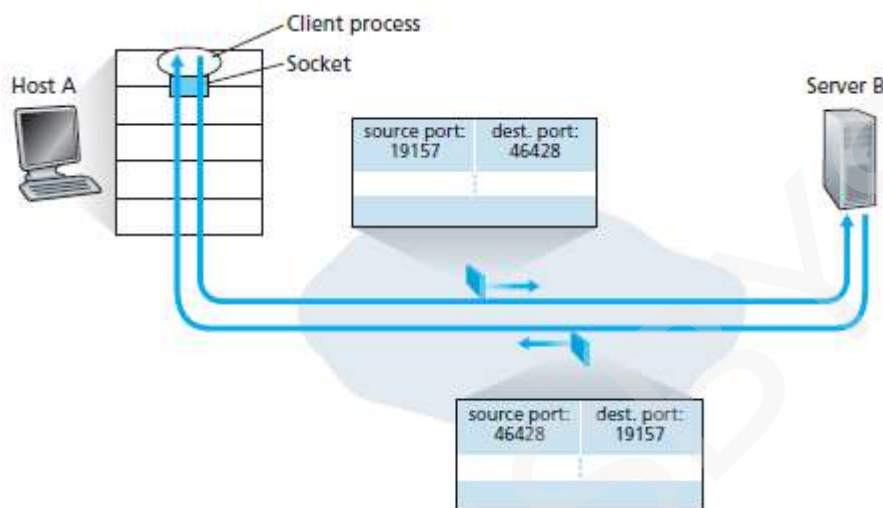


Figure 2.4: The inversion of source and destination-port-nos

- The server-host may support many simultaneous connection-sockets.
- Each socket will be
 - attached to a process.
 - identified by its own four tuple.
- When a segment arrives at the host, all 4 fields are used to direct the segment to the appropriate socket. (i.e. Demultiplexing).



COMPUTER NETWORKS

2.2.4 Web Servers and TCP

- Consider a host running a Web-server (ex: Apache) on port 80.
- When clients (ex: browsers) send segments to the server, all segments will have destination-port 80.
- The server distinguishes the segments from the different clients using two-tuple:
 - 1) Source IP addresses &
 - 2) Source-port-nos.

• Figure 2.5 shows a Web-server that creates a new process for each connection.

- The server can use either i) persistent HTTP or ii) non-persistent HTTP

i) Persistent HTTP

➤ Throughout the duration of the persistent connection the client and server exchange HTTP messages via the same server socket.

ii) Non-persistent HTTP

- A new TCP connection is created and closed for every request/response.
- Hence, a new socket is created and closed for every request/response.
- This can severely impact the performance of a busy Web-server.

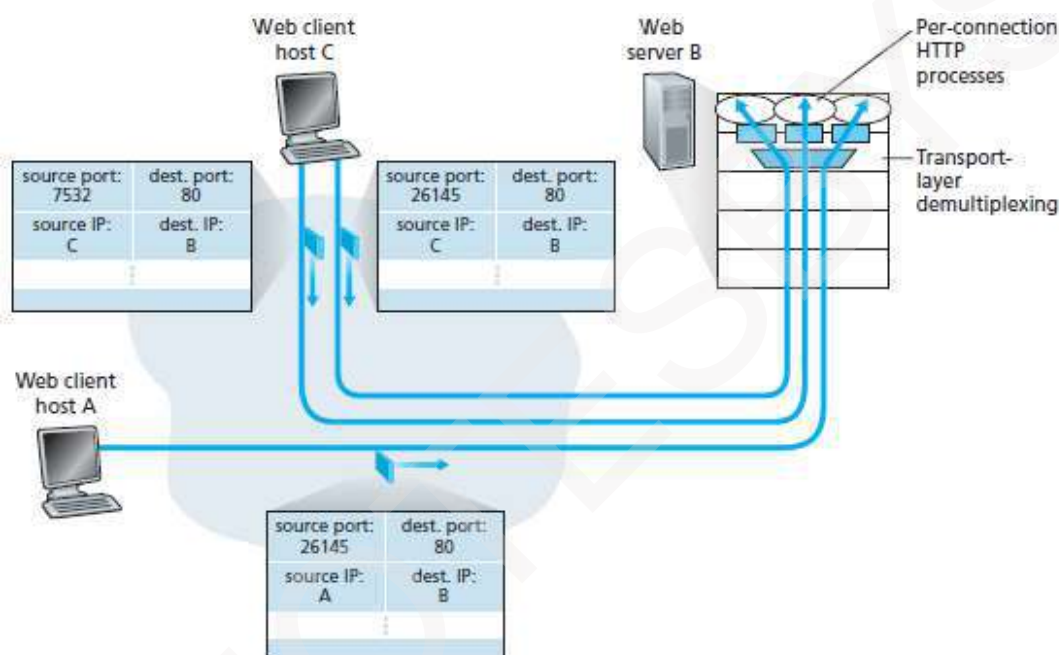


Figure 2.5: Two clients, using the same destination-port-no (80) to communicate with the same Web-server application



COMPUTER NETWORKS

2.3 Connectionless Transport: UDP

- UDP is an unreliable, connectionless protocol.
 - Unreliable service means UDP doesn't guarantee data will arrive to destination-process.
 - Connectionless means there is no handshaking b/w sender & receiver before sending data.
- It provides following 2 services:
 - i) Process-to-process data delivery and
 - ii) Error checking.
- It does not provide flow, error, or congestion control.
- At the sender, UDP
 - takes messages from the application-process
 - attaches source- & destination-port-nos and
 - passes the resulting segment to the network-layer.
- At the receiver, UDP
 - examines the destination-port-no in the segment and
 - delivers the segment to the correct application-process.
- It is suitable for application program that
 - needs to send short messages &
 - cannot afford the retransmission.
- UDP is suitable for many applications for the following reasons:
 - 1) Finer Application Level Control over what Data is Sent, and when.**
 - When an application-process passes data to UDP, the UDP
 - packs the data inside a segment and
 - passes immediately the segment to the network-layer.
 - On the other hand,
 - In TCP, a congestion-control mechanism throttles the sender when the n/w is congested
 - 2) No Connection Establishment.**
 - TCP uses a three-way handshake before it starts to transfer data.
 - UDP just immediately passes the data without any formal preliminaries.
 - Thus, UDP does not introduce any delay to establish a connection.
 - That's why, DNS runs over UDP rather than TCP.
 - 3) No Connection State.**
 - TCP maintains connection-state in the end-systems.
 - This connection-state includes
 - receive and send buffers
 - congestion-control parameters and
 - sequence- and acknowledgment-number parameters.
 - On the other hand,
 - In UDP, no connection-state is maintained.
 - 4) Small Packet Header Overhead.**
 - The TCP segment has 20 bytes of header overhead in every segment.
 - On the other hand, UDP has only 8 bytes of overhead.

Table 2.1: Popular Internet applications and their underlying transport protocols

Application	Application-Layer Protocol	Underlying Transport Protocol
Electronic mail	SMTP	TCP
Remote terminal access	Telnet	TCP
Web	HTTP	TCP
File transfer	FTP	TCP
Remote file server	NFS	Typically UDP
Streaming multimedia	typically proprietary	UDP or TCP
Internet telephony	typically proprietary	UDP or TCP
Network management	SNMP	Typically UDP
Routing protocol	RIP	Typically UDP
Name translation	DNS	Typically UDP



COMPUTER NETWORKS

2.3.1 UDP Segment Structure

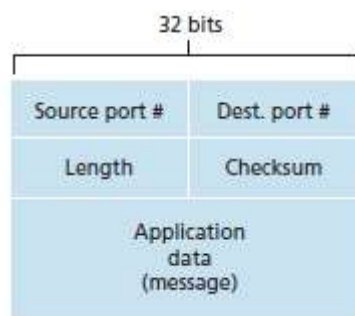


Figure 2.6: UDP segment structure

- UDP Segment contains following fields (Figure 2.6):
 - 1) Application Data:** This field occupies the data-field of the segment.
 - 2) Destination Port No:** This field is used to deliver the data to correct process running on the destination-host. (i.e. demultiplexing function).
 - 3) Length:** This field specifies the number of bytes in the segment (header plus data).
 - 4) Checksum:** This field is used for error-detection.

2.3.2 UDP Checksum

- The checksum is used for error-detection.
- The checksum is used to determine whether bits within the segment have been altered.
- How to calculate checksum on the sender:
 - 1) All the 16-bit words in the segment are added to get a sum.
 - 2) Then, the 1's complement of the sum is obtained to get a result.
 - 3) Finally, the result is added to the checksum-field inside the segment.
- How to check for error on the receiver:
 - 1) All the 16-bit words in the segment (including the checksum) are added to get a sum.
 - i) For no errors: In the sum, all the bits are 1. (Ex: 1111111)
 - ii) For any error: In the sum, at least one of the bits is a 0. (Ex: 1011111)

Example:

- On the sender:
 - Suppose that we have the following three 16-bit words:


```
0110011001100000
0101010101010101    → three 16 bits words
1000111100001100
```
 - The sum of first two 16-bit words is:


```
0110011001100000
0101010101010101
1011101110110101
```
 - Adding the third word to the above sum gives:


```
1011101110110101    → sum of 1st two 16 bit words
1000111100001100    → third 16 bit word
0100101011000010    → sum of all three 16 bit words
```
 - Taking 1's complement for the final sum:


```
0100101011000010    → sum of all three 16 bit words
10111010100111101    → 1's complement for the final sum
```
- The 1's complement value is called as checksum which is added inside the segment.
- On the receiver
 - All four 16-bit words are added, including the checksum.
 - i) If no errors are introduced into the packet, then clearly the sum will be 1111111111111111.
 - ii) If one of the bits is a 0, then errors have been introduced into the packet.

COMPUTER NETWORKS

2.4 Principles of Reliable Data Transfer

- Figure 2.7 illustrates the framework of reliable data transfer protocol.

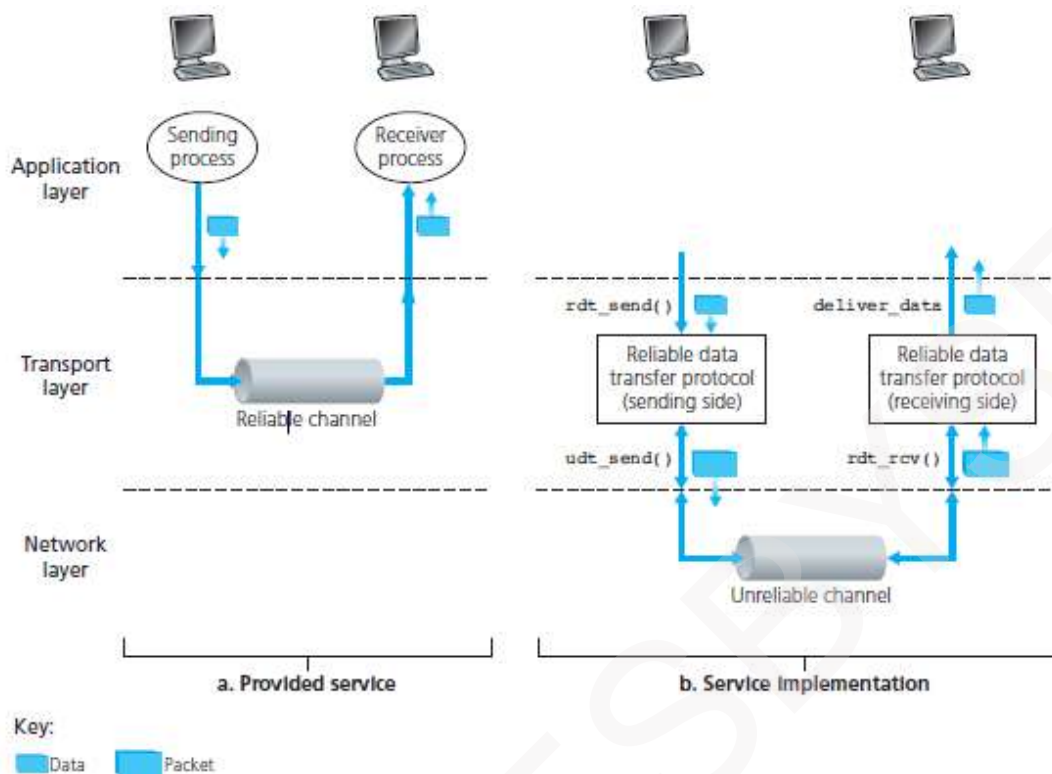


Figure 2.7: Reliable data transfer: Service model and service implementation

- On the sender, `rdt_send()` will be called when a packet has to be sent on the channel.
- On the receiver,
 - i) `rdt_rcv()` will be called when a packet has to be received on the channel.
 - ii) `deliver_data()` will be called when the data has to be delivered to the upper layer



COMPUTER NETWORKS

2.4.1 Building a Reliable Data Transfer Protocol

2.4.1.1 Reliable Data Transfer over a Perfectly Reliable Channel: rdt1.0

- Consider data transfer over a perfectly reliable channel.
- We call this protocol as rdt1.0.

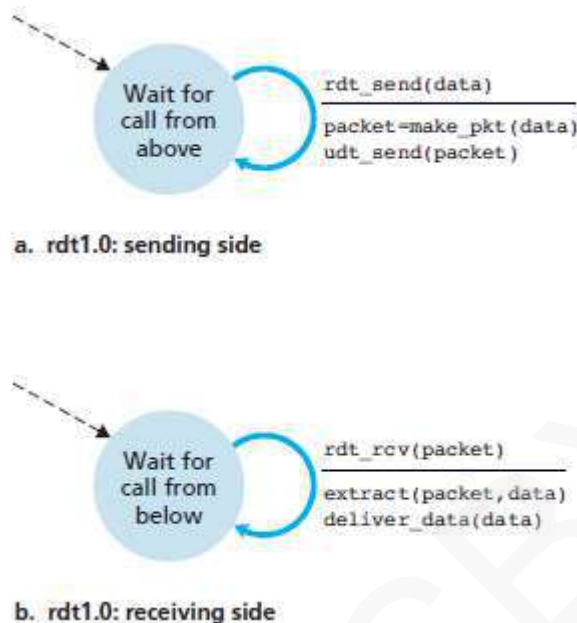


Figure 2.8: rdt1.0 – A protocol for a completely reliable channel

- The finite-state machine (FSM) definitions for the rdt1.0 sender and receiver are shown in Figure 2.8.
- The sender and receiver FSMs have only one state.
- In FSM, following notations are used:
 - i) The arrows indicate the transition of the protocol from one state to another.
 - ii) The event causing the transition is shown above the horizontal line labelling the transition.
 - iii) The action taken when the event occurs is shown below the horizontal line.
 - iv) The dashed arrow indicates the initial state.
- On the sender, rdt
 - accepts data from the upper layer via the `rdt_send(data)` event
 - creates a packet containing the data (via the action `make_pkt(data)`) and
 - sends the packet into the channel.
- On the receiver, rdt
 - receives a packet from the underlying channel via the `rdt_rcv(packet)` event
 - removes the data from the packet (via the action `extract(packet, data)`) and
 - passes the data up to the upper layer (via the action `deliver_data(data)`).



COMPUTER NETWORKS

2.4.1.2 Reliable Data Transfer over a Channel with Bit Errors: rdt2.0

- Consider data transfer over an unreliable channel in which bits in a packet may be corrupted.
- We call this protocol as rdt2.0.
- The message dictation protocol uses both
 - positive acknowledgements (ACK) and
 - negative acknowledgements (NAK).
- The receiver uses these control messages to inform the sender about
 - what has been received correctly and
 - what has been received in error and thus requires retransmission.
- Reliable data transfer protocols based on the retransmission are known as ARQ protocols.
- Three additional protocol capabilities are required in ARQ protocols:
 - 1) Error Detection**
 - A mechanism is needed to allow the receiver to detect when bit-errors have occurred.
 - UDP uses the checksum field for error-detection.
 - Error-correction techniques allow the receiver to detect and correct packet bit-errors.
 - 2) Receiver Feedback**
 - Since the sender and receiver are typically executing on different end-systems.
 - The only way for the sender to learn about status of the receiver is by the receiver providing explicit feedback to the sender.
 - For example: ACK & NAK
 - 3) Retransmission**
 - A packet that is received in error at the receiver will be retransmitted by the sender.
- Figure 2.9 shows the FSM representation of rdt2.0.

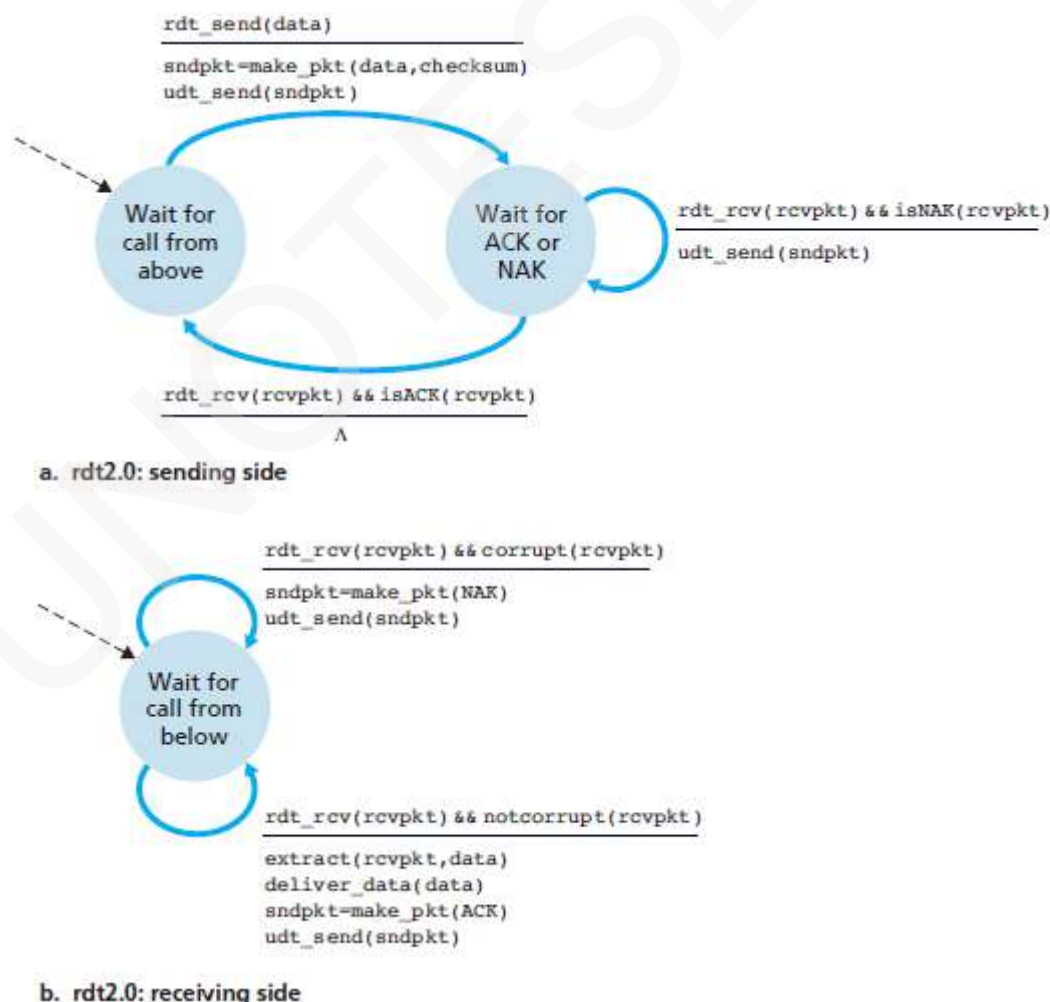


Figure 2.9: rdt2.0-A protocol for a channel with bit-errors



COMPUTER NETWORKS

Sender FSM

- The sender of rdt2.0 has 2 states:
 - 1) In one state, the protocol is waiting for data to be passed down from the upper layer.
 - 2) In other state, the protocol is waiting for an ACK or a NAK from the receiver.
 - i) If an ACK is received, the protocol
 - knows that the most recently transmitted packet has been received correctly
 - returns to the state of waiting for data from the upper layer.
 - ii) If a NAK is received, the protocol
 - retransmits the last packet and
 - waits for an ACK or NAK to be returned by the receiver.
- The sender will not send a new data until it is sure that the receiver has correctly received the current packet.
- Because of this behaviour, protocol rdt2.0 is known as stop-and-wait protocols.

Receiver FSM

- The receiver of rdt2.0 has a single state.
- On packet arrival, the receiver replies with either an ACK or a NAK, depending on the received packet is corrupted or not.

**COMPUTER NETWORKS****2.4.1.2.1 Sender Handles Garbled ACK/NAKs: rdt2.1**

- Problem with rdt2.0:
 - If an ACK or NAK is corrupted, the sender cannot know whether the receiver has correctly received the data or not.
- Solution: The sender resends the current data packet when it receives garbled ACK or NAK packet.
 - Problem: This approach introduces duplicate packets into the channel.
 - Solution: Add sequence-number field to the data packet.
 - The receiver has to only check the sequence-number to determine whether the received packet is a retransmission or not.
- For a stop-and-wait protocol, a 1-bit sequence-number will be sufficient.
- A 1-bit sequence-number allows the receiver to know whether the sender is sending
 - previously transmitted packet (0) or
 - new packet (1).
- We call this protocol as rdt2.1.
- Figure 2.10 and 2.11 shows the FSM description for rdt2.1.

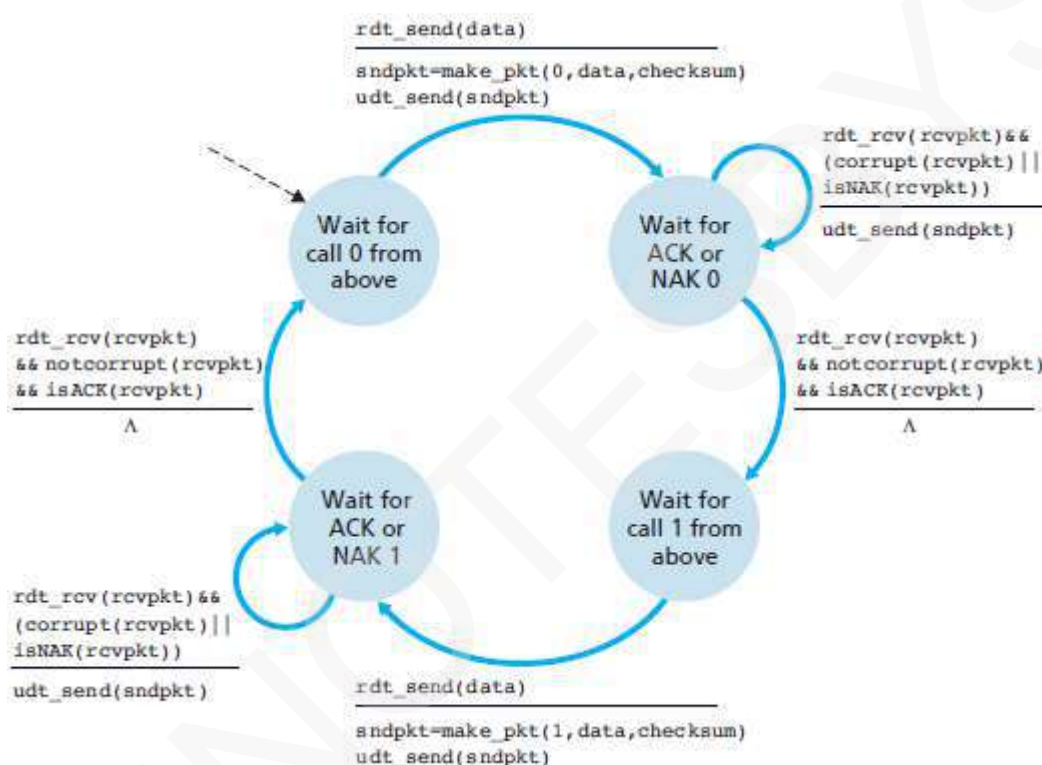


Figure 2.10: rdt2.1 sender

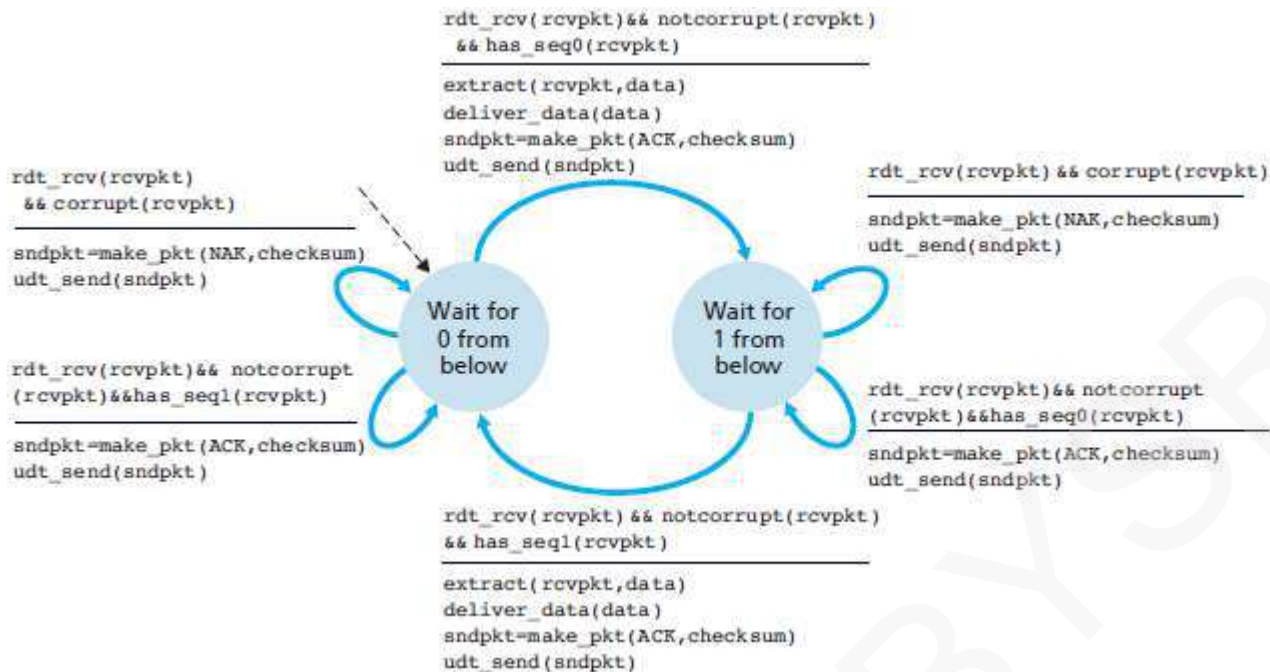


Figure 2.11: rdt2.1 receiver



COMPUTER NETWORKS

2.4.1.2.2 Sender uses ACK/NAKs: rdt2.2

- Protocol rdt2.2 uses both positive and negative acknowledgments from the receiver to the sender.
 - i) When out-of-order packet is received, the receiver sends a positive acknowledgment (ACK).
 - ii) When a corrupted packet is received, the receiver sends a negative acknowledgment (NAK).
- We call this protocol as rdt2.2. (Figure 2.12 and 2.13).

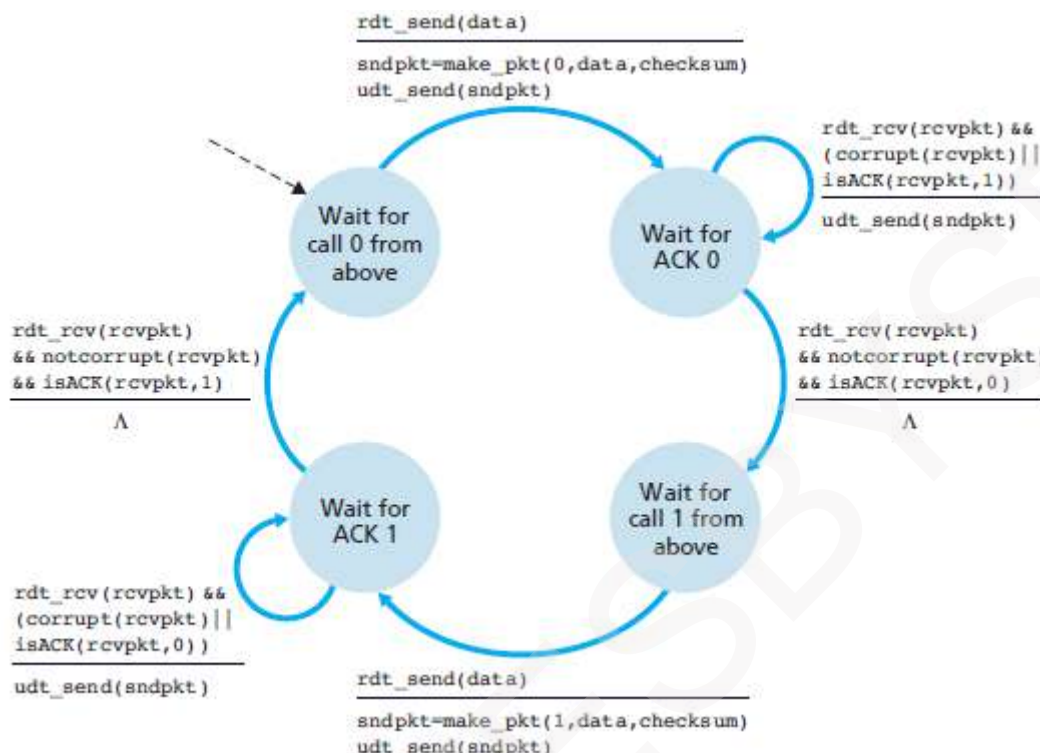


Figure 2.12: rdt2.2 sender

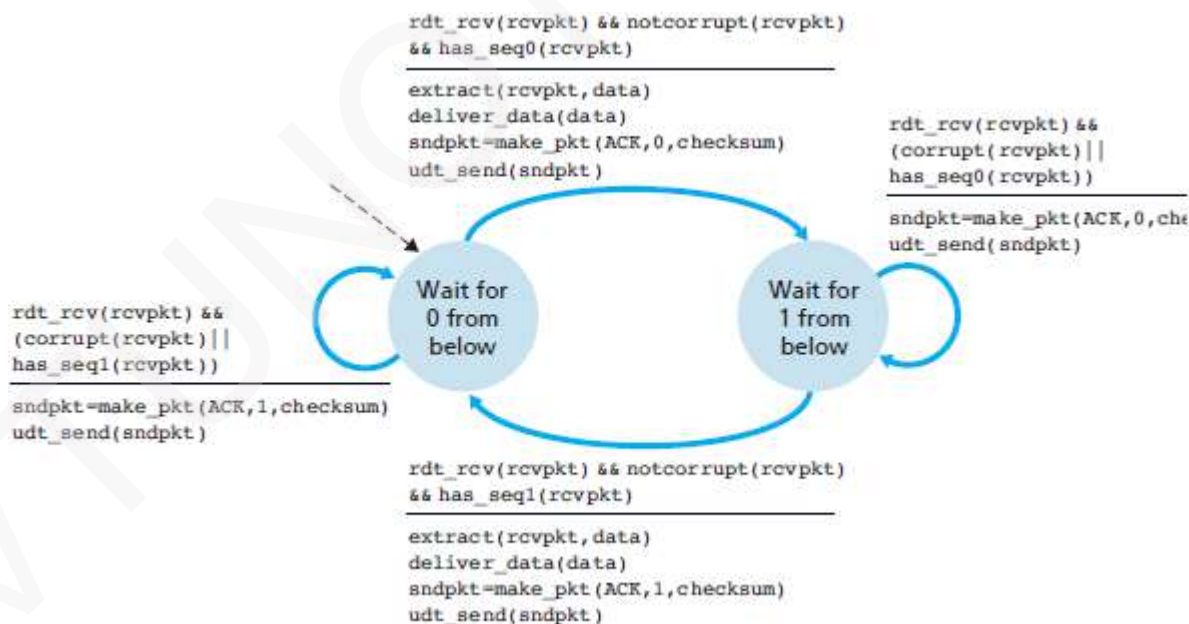


Figure 2.13: rdt2.2 receiver



COMPUTER NETWORKS

2.4.1.3 Reliable Data Transfer over a Lossy Channel with Bit Errors: rdt3.0

- Consider data transfer over an unreliable channel in which packet loss may occur.
- We call this protocol as rdt3.0.
- Two problems must be solved by the rdt3.0:
 - 1) How to detect packet loss?
 - 2) What to do when packet loss occurs?
- Solution:
 - The sender
 - sends one packet & starts a timer and
 - waits for ACK from the receiver (okay to go ahead).
 - If the timer expires before ACK arrives, the sender retransmits the packet and restarts the timer.
- The sender must wait at least as long as
 - 1) A round-trip delay between the sender and receiver plus
 - 2) Amount of time needed to process a packet at the receiver.
- Implementing a time-based retransmission mechanism requires a countdown timer.
- The timer must interrupt the sender after a given amount of time has expired.
- Figure 2.14 shows the sender FSM for rdt3.0, a protocol that reliably transfers data over a channel that can corrupt or lose packets;
- Figure 2.15 shows how the protocol operates with no lost or delayed packets and how it handles lost data packets.
- Because sequence-numbers alternate b/w 0 & 1, protocol rdt3.0 is known as alternating-bit protocol.

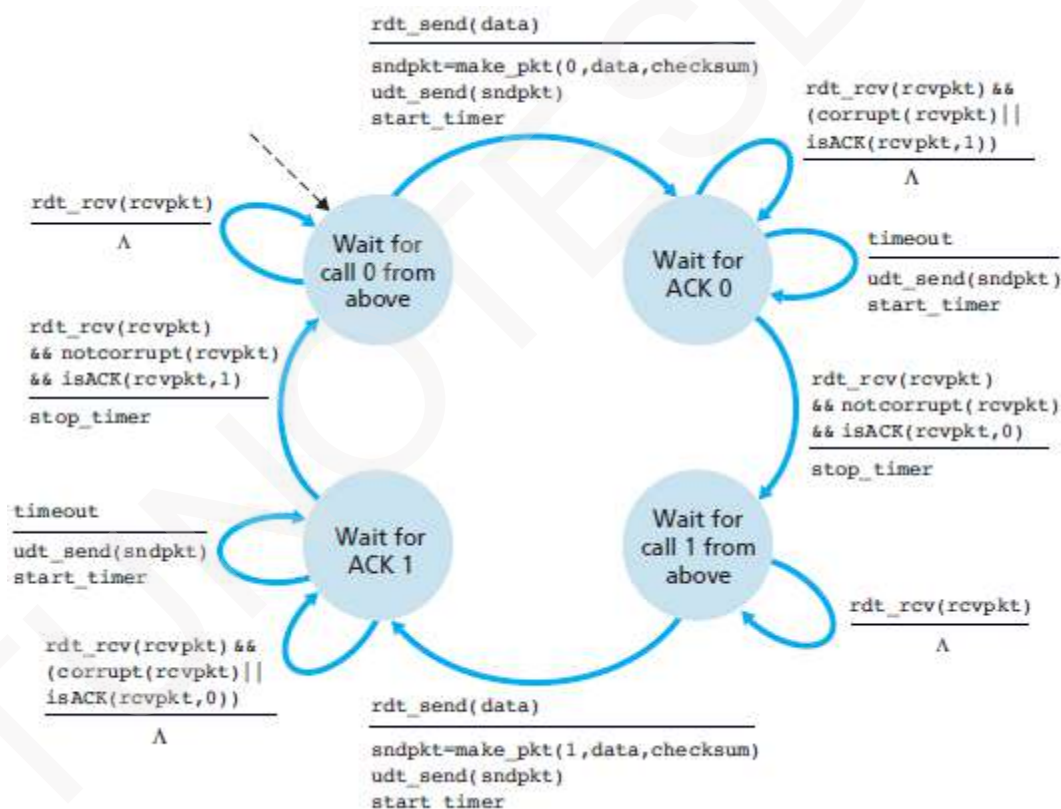


Figure 2.14: rdt3.0 sender



COMPUTER NETWORKS

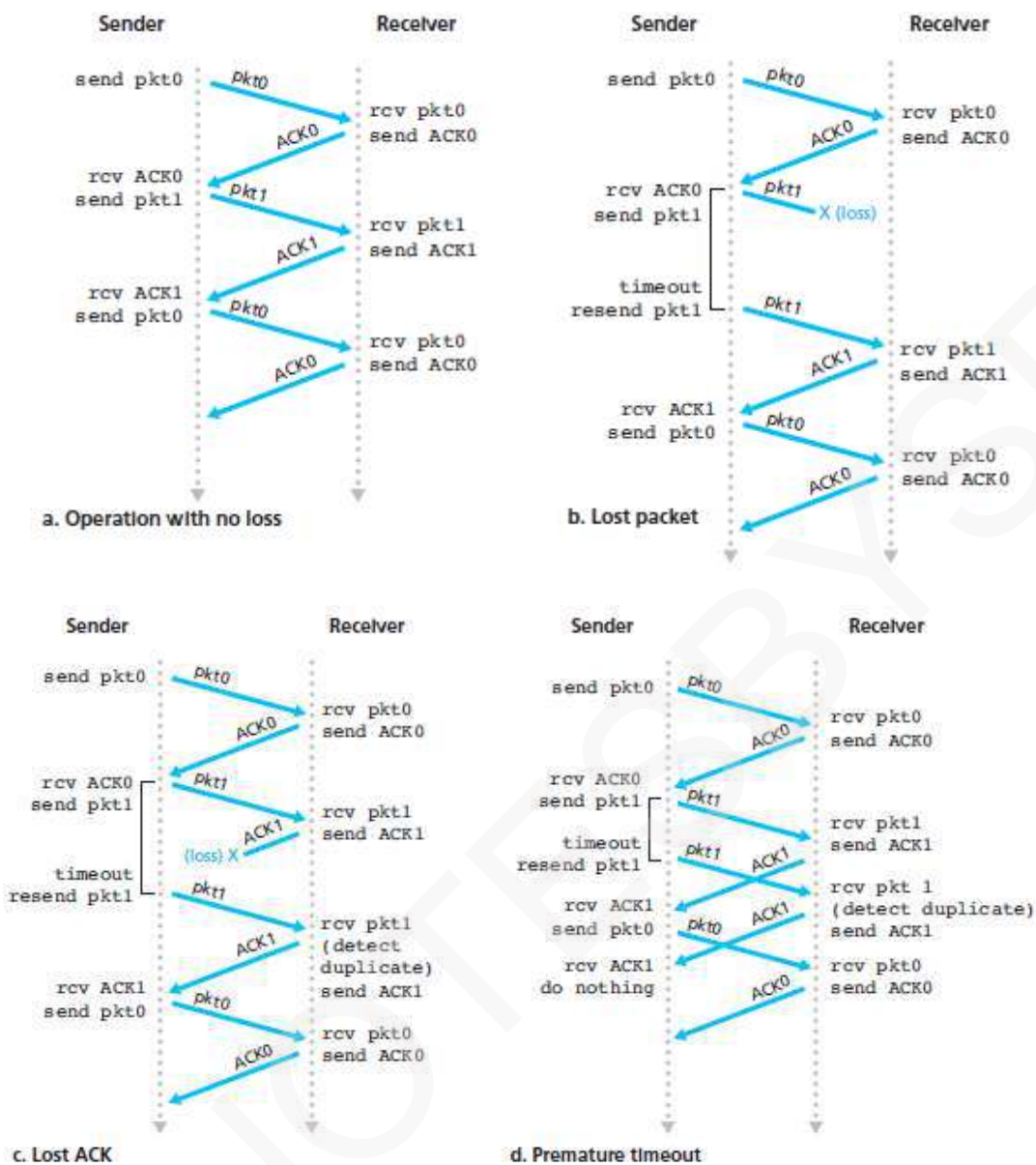


Figure 2.15: Operation of rdt3.0, the alternating-bit protocol



COMPUTER NETWORKS

2.4.2 Pipelined Reliable Data Transfer Protocols

- The sender is allowed to send multiple packets without waiting for acknowledgments.
- This is illustrated in Figure 2.16 (b).
- Pipelining has the following consequences:
 - 1) The range of sequence-numbers must be increased.
 - 2) The sender and receiver may have to buffer more than one packet.
- Two basic approaches toward pipelined error recovery can be identified:
 - 1) Go-Back-N and 2) Selective repeat.

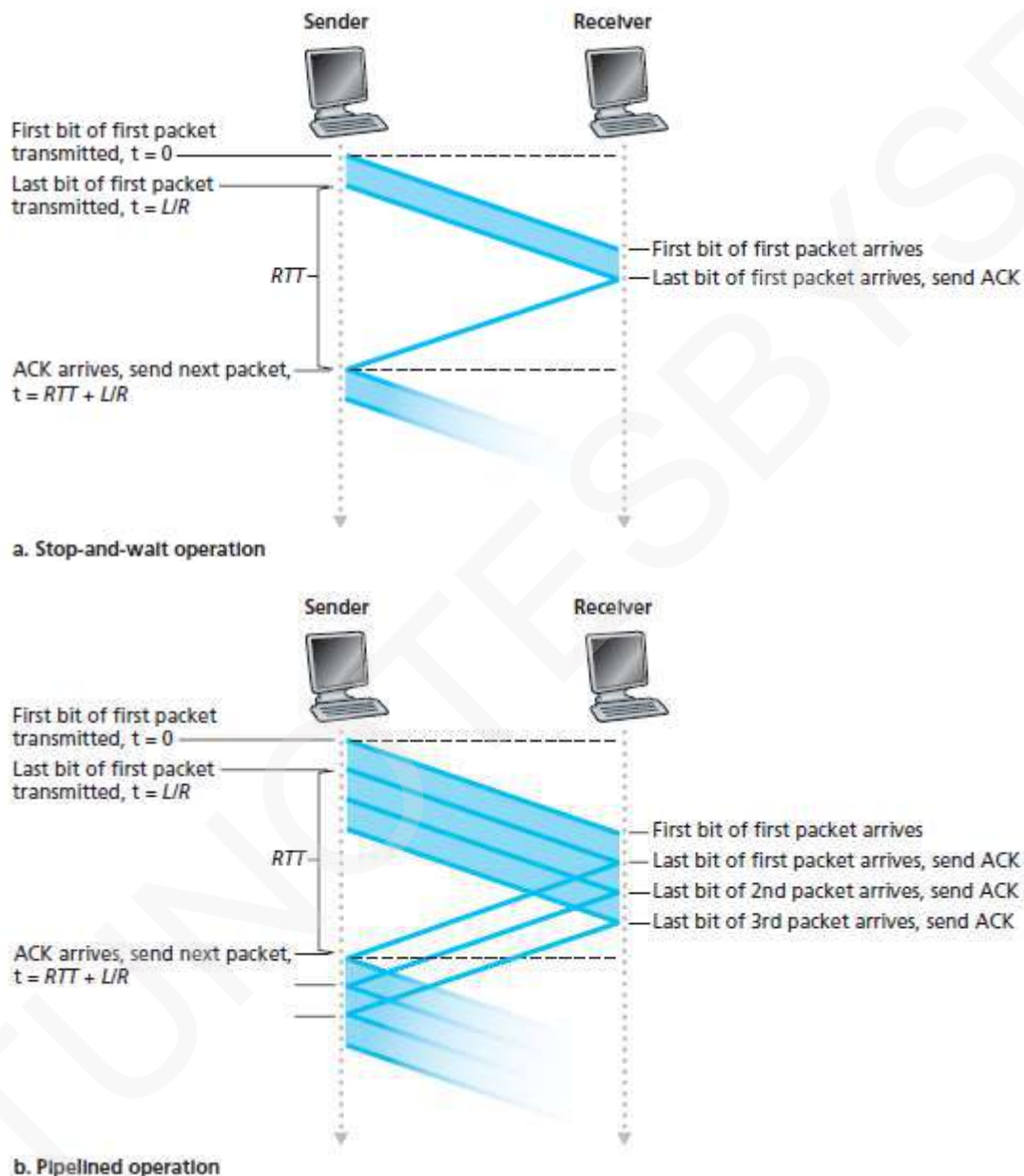


Figure 2.16: Stop-and-wait and pipelined sending



COMPUTER NETWORKS

2.4.3 Go-Back-N (GBN)

- The sender is allowed to transmit multiple packets without waiting for an acknowledgment.
- But, the sender is constrained to have at most N unacknowledged packets in the pipeline.
 - Where N = window-size which refers maximum no. of unacknowledged packets in the pipeline
- GBN protocol is called a sliding-window protocol.
- Figure 2.17 shows the sender's view of the range of sequence-numbers.

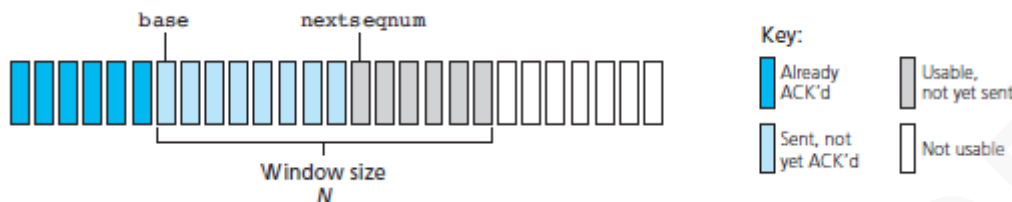


Figure 2.17: Sender's view of sequence-numbers in Go-Back-N

- Figure 2.18 and 2.19 give a FSM description of the sender and receivers of a GBN protocol.

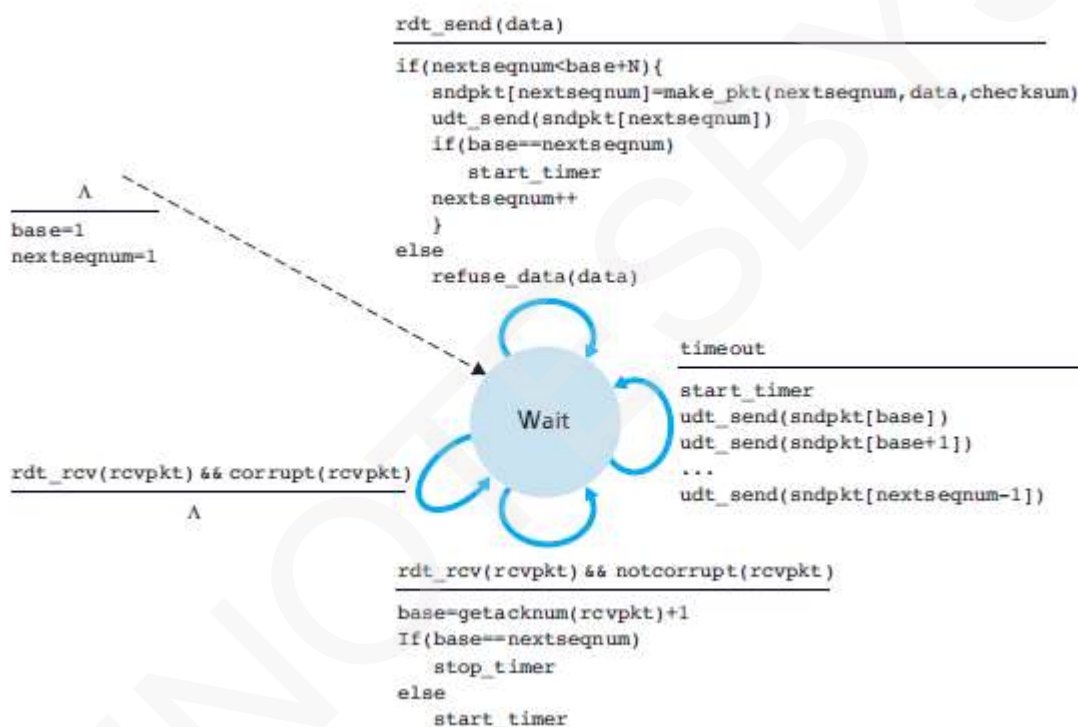


Figure 2.18: Extended FSM description of GBN sender

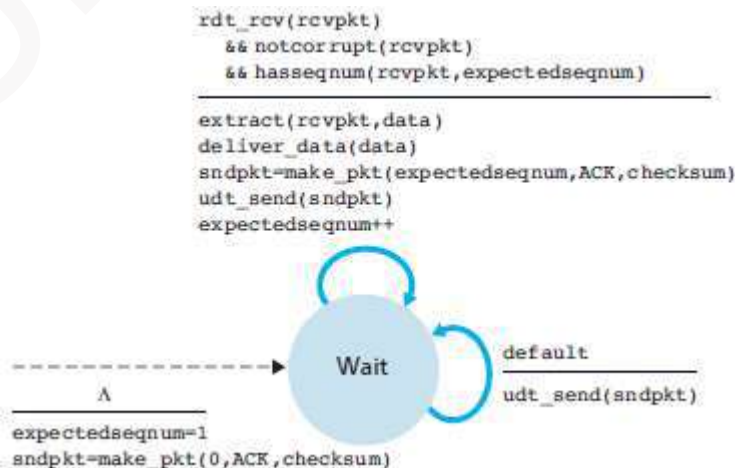


Figure 2.19: Extended FSM description of GBN receiver

“Try not to become a man of success, but rather try to become a man of value.” —Albert Einstein



COMPUTER NETWORKS

2.4.3.1 GBN Sender

- The sender must respond to 3 types of events:

1) Invocation from above.

- When `rdt_send()` is called from above, the sender first checks to see if the window is full i.e. whether there are N outstanding, unacknowledged packets.
 - i) If the window is not full, the sender creates and sends a packet.
 - ii) If the window is full, the sender simply returns the data back to the upper layer. This is an implicit indication that the window is full.

2) Receipt of an ACK.

- An acknowledgment for a packet with sequence-number n will be taken to be a cumulative acknowledgment.
- All packets with a sequence-number up to n have been correctly received at the receiver.

3) A Timeout Event.

- A timer will be used to recover from lost data or acknowledgment packets.
 - i) If a timeout occurs, the sender resends all packets that have been previously sent but that have not yet been acknowledged.
 - ii) If an ACK is received but there are still additional transmitted but not yet acknowledged packets, the timer is restarted.
 - iii) If there are no outstanding unacknowledged packets, the timer is stopped.

2.4.3.2 GBN Receiver

- If a packet with sequence-number n is received correctly and is in order, the receiver
 - sends an ACK for packet n and
 - delivers the packet to the upper layer.
- In all other cases, the receiver
 - discards the packet and
 - resends an ACK for the most recently received in-order packet.



COMPUTER NETWORKS

2.4.3.3 Operation of the GBN Protocol

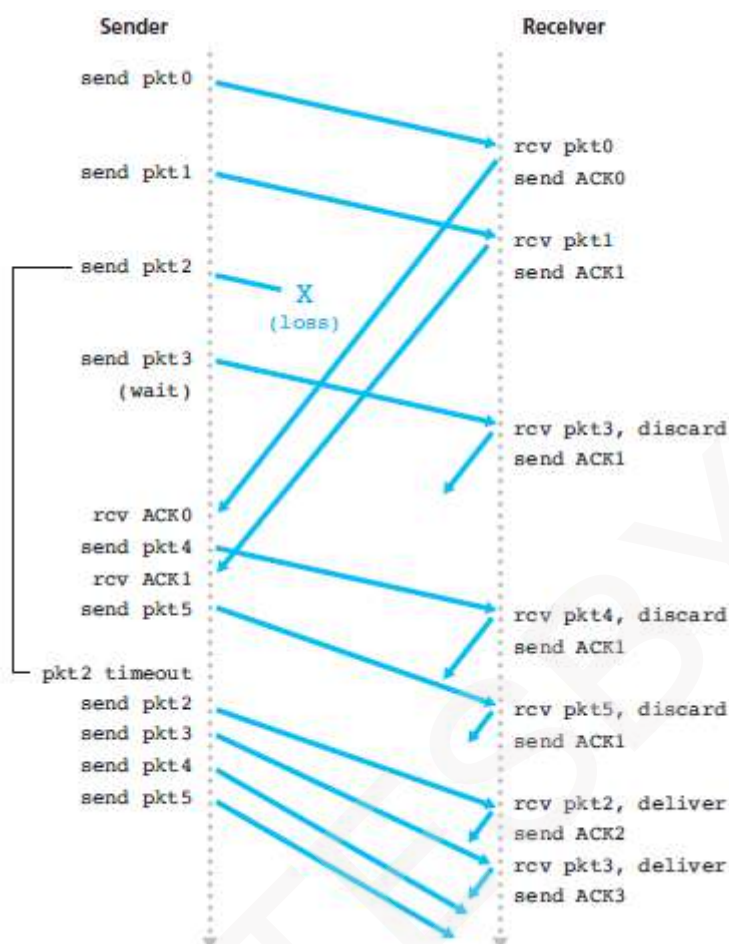


Figure 2.20: Go-Back-N in operation

- Figure 2.20 shows the operation of the GBN protocol for the case of a window-size of four packets.
- The sender sends packets 0 through 3.
- The sender then must wait for one or more of these packets to be acknowledged before proceeding.
- As each successive ACK (for ex, ACK0 and ACK1) is received, the window slides forward and the sender transmits one new packet (pkt4 and pkt5, respectively).
- On the receiver, packet 2 is lost and thus packets 3, 4, and 5 are found to be out of order and are discarded.



COMPUTER NETWORKS

2.4.4 Selective Repeat (SR)

- Problem with GBN:
 - GBN suffers from performance problems.
 - When the window-size and bandwidth-delay product are both large, many packets can be in the pipeline.
 - Thus, a single packet error results in retransmission of a large number of packets.
- Solution: Use Selective Repeat (SR).

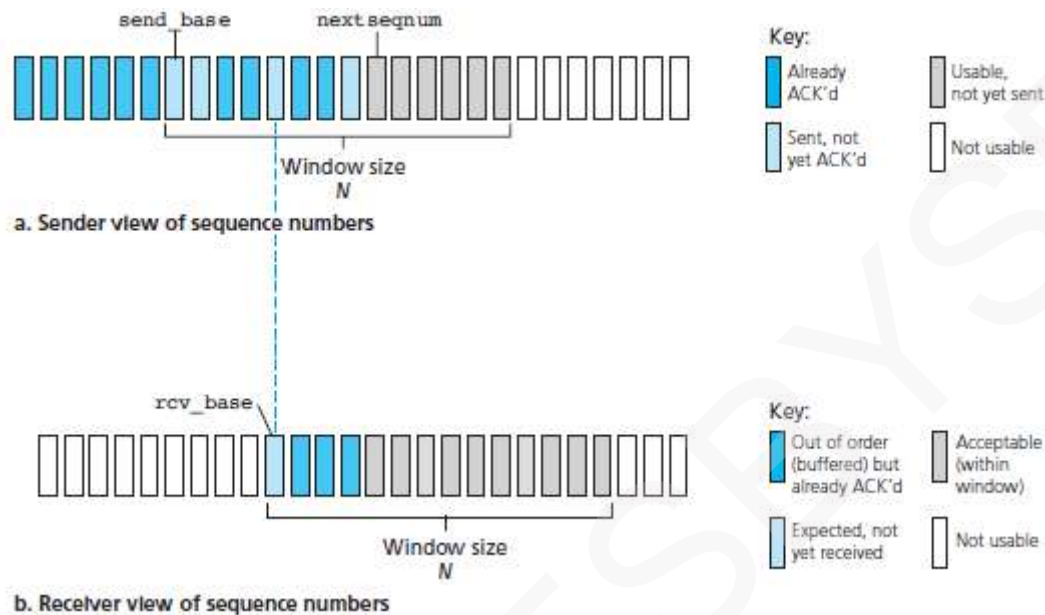


Figure 2.21: Selective-repeat (SR) sender and receiver views of sequence-number space

- The sender retransmits only those packets that it suspects were erroneous.
- Thus, avoids unnecessary retransmissions. Hence, the name "selective-repeat".
- The receiver individually acknowledge correctly received packets.
- A window-size N is used to limit the no. of outstanding, unacknowledged packets in the pipeline.
- Figure 2.21 shows the SR sender's view of the sequence-number space.

2.4.4.1 SR Sender

- The various actions taken by the SR sender are as follows:

1) Data Received from above.

- When data is received from above, the sender checks the next available sequence-number for the packet.
- If the sequence-number is within the sender's window;
 - Then, the data is packetized and sent;
 - Otherwise, the data is buffered for later transmission.

2) Timeout.

- Timers are used to protect against lost packets.
- Each packet must have its own logical timer. This is because
 - only a single packet will be transmitted on timeout.

3) ACK Received.

- If an ACK is received, the sender marks that packet as having been received.
- If the packet's sequence-number is equal to `send_base`, the window base is increased by the smallest sequence-number.
- If there are untransmitted packets with sequence-numbers that fall within the window, these packets are transmitted.



COMPUTER NETWORKS

2.4.4.2 SR Receiver

- The various actions taken by the SR receiver are as follows:

1) Packet with sequence-number in $[rcv_base, rcv_base+N-1]$ is correctly received.

- In this case,
 - received packet falls within the receiver's window and
 - selective ACK packet is returned to the sender.
- If the packet was not previously received, it is buffered.
- If this packet has a sequence-number equal to rcv_base , then this packet, and any previously buffered and consecutively numbered packets are delivered to the upper layer.
- The receive-window is then moved forward by the no. of packets delivered to the upper layer.
- For example: consider Figure 2.22.
 - ✕ When a packet with a sequence-number of $rcv_base=2$ is received, it and packets 3, 4, and 5 can be delivered to the upper layer.

2) Packet with sequence-number in $[rcv_base-N, rcv_base-1]$ is correctly received.

- In this case, an ACK must be generated, even though this is a packet that the receiver has previously acknowledged.

3) Otherwise.

- Ignore the packet.

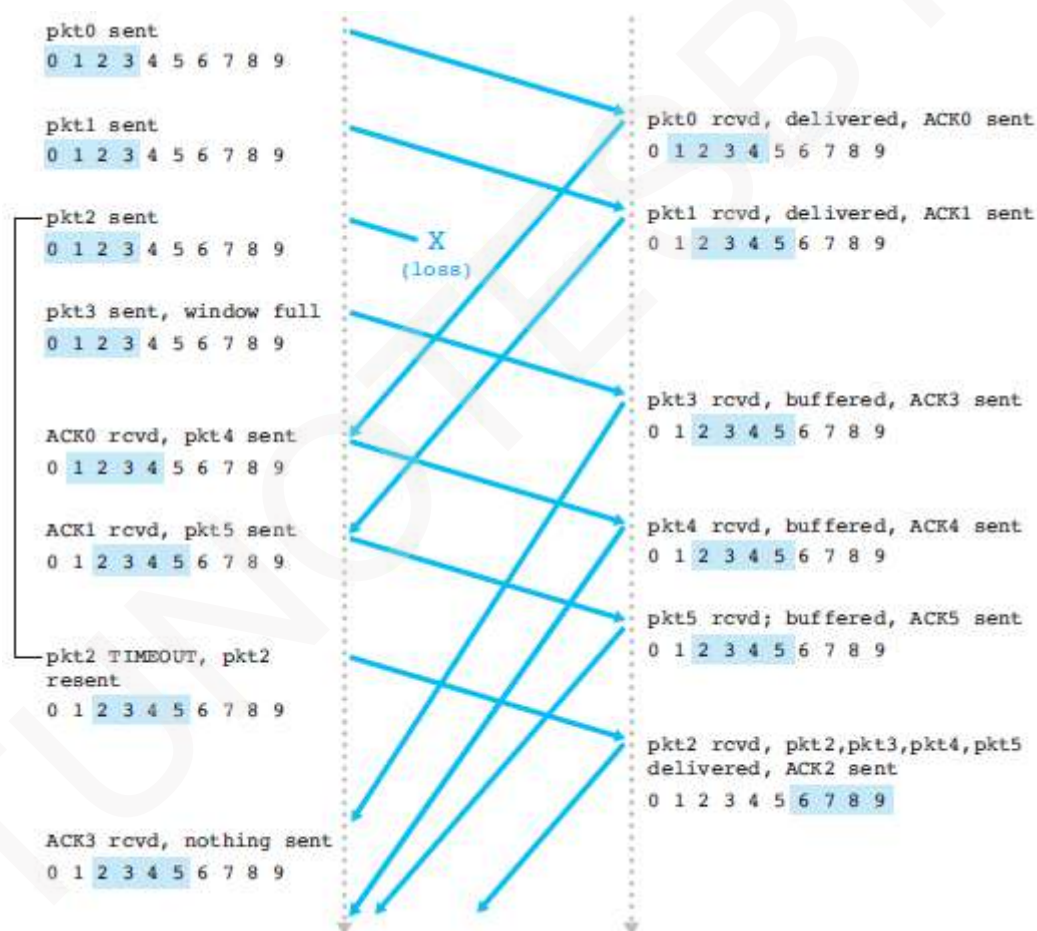


Figure 2.22: SR operation

**2.4.5 Summary of Reliable Data Transfer Mechanisms and their Use**

Table 2.2: Summary of reliable data transfer mechanisms and their use

Mechanism	Use, Comments
Checksum	Used to detect bit errors in a transmitted packet.
Timer	Used to timeout/retransmit a packet because the packet (or its ACK) was lost. Because timeouts can occur when a packet is delayed but not lost, duplicate copies of a packet may be received by a receiver.
Sequence-number	Used for sequential numbering of packets of data flowing from sender to receiver. Gaps in the sequence-numbers of received packets allow the receiver to detect a lost packet. Packets with duplicate sequence-numbers allow the receiver to detect duplicate copies of a packet.
Acknowledgment	Used by the receiver to tell the sender that a packet or set of packets has been received correctly. Acknowledgments will typically carry the sequence-number of the packet or packets being acknowledged. Acknowledgments may be individual or cumulative, depending on the protocol.
Negative acknowledgment	Used by the receiver to tell the sender that a packet has not been received correctly. Negative acknowledgments will typically carry the sequence-number of the packet that was not received correctly.
Window, pipelining	The sender may be restricted to sending only packets with sequence-numbers that fall within a given range. By allowing multiple packets to be transmitted but not yet acknowledged, sender utilization can be increased over a stop-and-wait mode of operation.



COMPUTER NETWORKS

2.5 Connection-Oriented Transport: TCP

- TCP is a reliable connection-oriented protocol.
 - Connection-oriented means a connection is established b/w sender & receiver before sending the data.
 - Reliable service means TCP guarantees that the data will arrive to destination-process correctly.
- TCP provides flow-control, error-control and congestion-control.

2.5.1 The TCP Connection

- The features of TCP are as follows:

1) Connection Oriented

- TCP is said to be connection-oriented. This is because The 2 application-processes must first establish connection with each other before they begin communication.
- Both application-processes will initialize many state-variables associated with the connection.

2) Runs in the End Systems

- TCP runs only in the end-systems but not in the intermediate routers.
- The routers do not maintain any state-variables associated with the connection.

3) Full Duplex Service

- TCP connection provides a full-duplex service.
- Both application-processes can transmit and receive the data at the same time.

4) Point-to-Point

- A TCP connection is point-to-point i.e. only 2 devices are connected by a dedicated-link
- So, multicasting is not possible.

5) Three-way Handshake

- Connection-establishment process is referred to as a three-way handshake. This is because 3 segments are sent between the two hosts:
 - i) The client sends a first-segment.
 - ii) The server responds with a second-segment and
 - iii) Finally, the client responds again with a third segment containing payload (or data).

6) Maximum Segment Size (MSS)

- MSS limits the maximum amount of data that can be placed in a segment.
For example: MSS = 1,500 bytes for Ethernet

7) Send & Receive Buffers

- As shown in Figure 2.23, consider sending data from the client-process to the server-process.

At Sender

- i) The client-process passes a stream-of-data through the socket.
- ii) Then, TCP forwards the data to the send-buffer.
- iii) Each chunk-of-data is appended with a header to form a segment.
- iv) The segments are sent into the network.

At Receiver

- i) The segment's data is placed in the receive-buffer.
- ii) The application reads the stream-of-data from the receive-buffer.

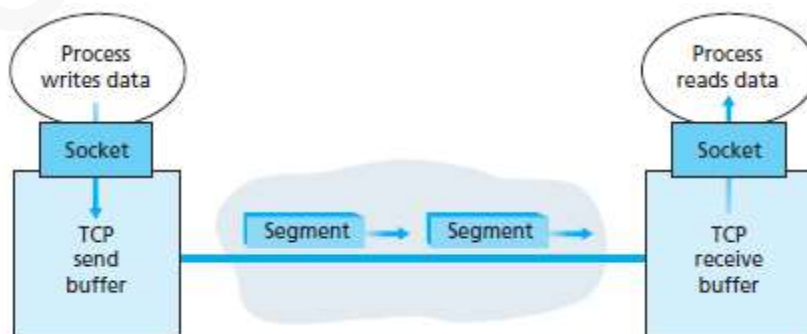


Figure 2.23: TCP send and receive-buffers



COMPUTER NETWORKS

2.5.2 TCP Segment Structure

- The segment consists of header-fields and a data-field.
- The data-field contains a chunk-of-data.
- When TCP sends a large file, it breaks the file into chunks of size MSS.
- Figure 2.24 shows the structure of the TCP segment.

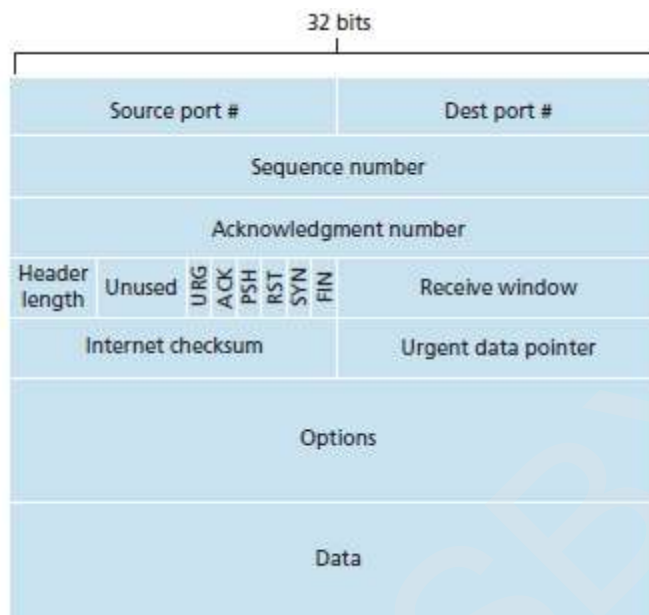


Figure 2.24: TCP segment structure

- The fields of TCP segment are as follows:

1) Source and Destination Port Numbers

- These fields are used for multiplexing/demultiplexing data from/to upper-layer applications.

2) Sequence Number & Acknowledgment Number

- These fields are used by sender & receiver in implementing a reliable data-transfer-service.

3) Header Length

- This field specifies the length of the TCP header.

4) Flag

- This field contains 6 bits.

i) ACK

- ✗ This bit indicates that value of acknowledgment field is valid.

ii) RST, SYN & FIN

- ✗ These bits are used for connection setup and teardown.

iii) PSH

- ✗ This bit indicates the sender has invoked the push operation.

iv) URG

- ✗ This bit indicates the segment contains urgent-data.

5) Receive Window

- This field defines receiver's window size
- This field is used for flow control.

6) Checksum

- This field is used for error-detection.

7) Urgent Data Pointer

- This field indicates the location of the last byte of the urgent data.

8) Options

- This field is used when a sender & receiver negotiate the MSS for use in high-speed networks.



COMPUTER NETWORKS

2.5.2.1 Sequence Numbers and Acknowledgment Numbers

Sequence Numbers

- The sequence-number is used for sequential numbering of packets of data flowing from sender to receiver.

- Applications:

- Gaps in the sequence-numbers of received packets allow the receiver to detect a lost packet.
- Packets with duplicate sequence-numbers allow the receiver to detect duplicate copies of a packet.

Acknowledgment Numbers

- The acknowledgment-number is used by the receiver to tell the sender that a packet has been received correctly.

- Acknowledgments will typically carry the sequence-number of the packet being acknowledged.

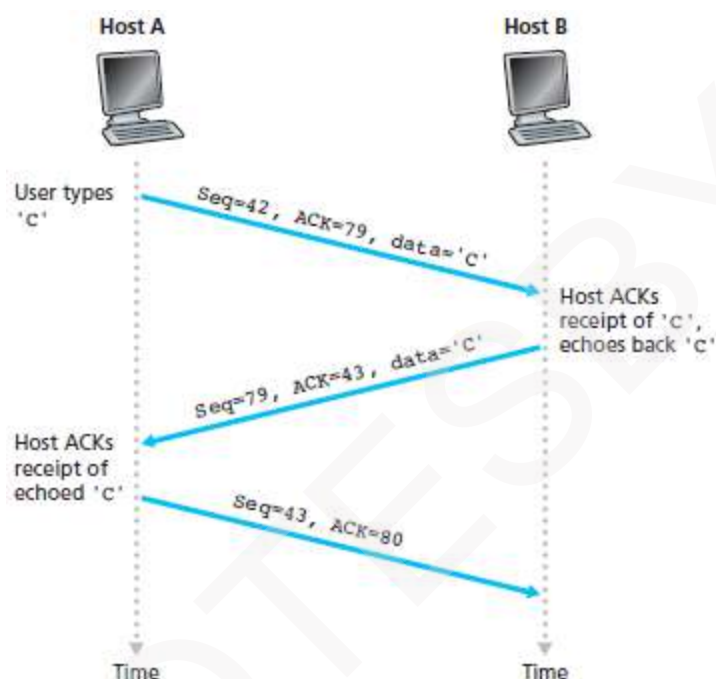


Figure 2.25: Sequence and acknowledgment-numbers for a simple Telnet application over TCP

- Consider an example (Figure 2.25):

- A process in Host-A wants to send a stream-of-data to a process in Host-B.
- In Host-A, each byte in the data-stream is numbered as shown in Figure 2.26.

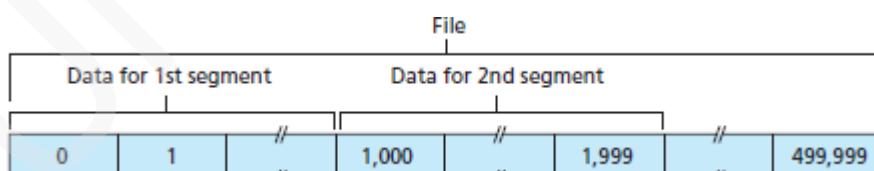


Figure 2.26: Dividing file data into TCP segments

- The first segment from A to B has a sequence-number 42 i.e. Seq=42.
- The second segment from B to A has a sequence-number 79 i.e. Seq=79.
- The second segment from B to A has acknowledgment-number 43, which is the sequence-number of the next byte, Host-B is expecting from Host-A. (i.e. ACK=43).
- What does a host do when it receives out-of-order bytes?

Answer: There are two choices:

- The receiver immediately discards out-of-order bytes.
- The receiver
 - keeps the out-of-order bytes and
 - waits for the missing bytes to fill in the gaps.

“You will not be punished for your anger; you will be punished by your anger.” —Buddha



COMPUTER NETWORKS

2.5.2.2 Telnet: A Case Study for Sequence and Acknowledgment Numbers

- Telnet is a popular application-layer protocol used for remote-login.
- Telnet runs over TCP.
- Telnet is designed to work between any pair of hosts.
- As shown in Figure 2.27, suppose client initiates a Telnet session with server.
- Now suppose the user types a single letter, 'C'.
- Three segments are sent between client & server:

1) First Segment

- The first-segment is sent from the client to the server.
- The segment contains
 - letter 'C'
 - sequence-number 42
 - acknowledgment-number 79

2) Second Segment

- The second-segment is sent from the server to the client.
- Two purpose of the segment:
 - i) It provides an acknowledgment of the data the server has received.
 - ii) It is used to echo back the letter 'C'.
- The acknowledgment for client-to-server data is carried in a segment carrying server-to-client data.
- This acknowledgment is said to be piggybacked on the server-to-client data-segment.

3) Third Segment

- The third segment is sent from the client to the server.
- One purpose of the segment:
 - i) It acknowledges the data it has received from the server.

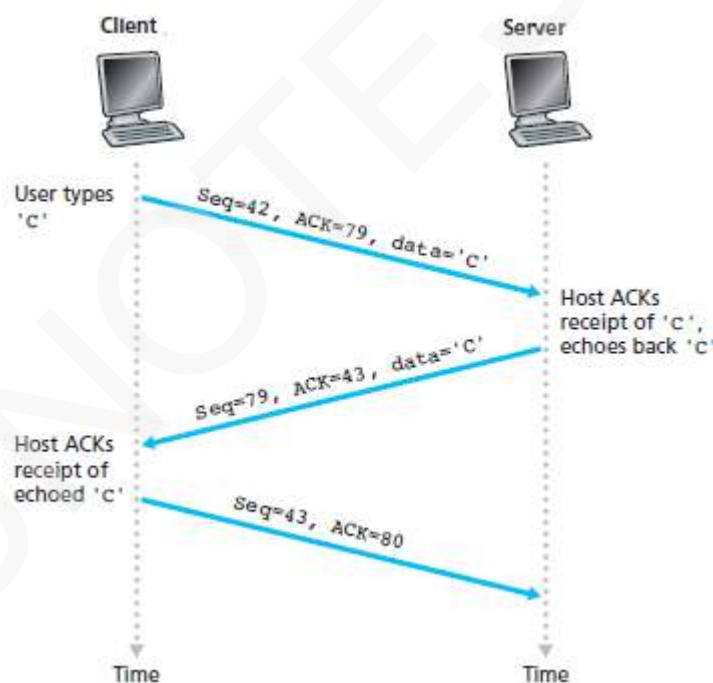


Figure 2.27: Sequence and acknowledgment-numbers for a simple Telnet application over TCP



COMPUTER NETWORKS

2.5.3 Round Trip Time Estimation and Timeout

- TCP uses a timeout/retransmit mechanism to recover from lost segments.
- Clearly, the timeout should be larger than the round-trip-time (RTT) of the connection.

2.5.3.1 Estimating the Round Trip Time

- SampleRTT is defined as
"The amount of time b/w when the segment is sent and when an acknowledgment is received."
- Obviously, the SampleRTT values will fluctuate from segment to segment due to congestion.
- TCP maintains an average of the SampleRTT values, which is referred to as EstimatedRTT.

$$\text{EstimatedRTT} = (1 - \alpha) \cdot \text{EstimatedRTT} + \alpha \cdot \text{SampleRTT}$$

- DevRTT is defined as
"An estimate of how much SampleRTT typically deviates from EstimatedRTT."

$$\text{DevRTT} = (1 - \beta) \cdot \text{DevRTT} + \beta \cdot |\text{SampleRTT} - \text{EstimatedRTT}|$$

- If the SampleRTT values have little fluctuation, then DevRTT will be small.
If the SampleRTT values have huge fluctuation, then DevRTT will be large.

2.5.3.2 Setting and Managing the Retransmission Timeout Interval

- What value should be used for timeout interval?
- Clearly, the interval should be greater than or equal to EstimatedRTT.
- Timeout interval is given by:

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 \cdot \text{DevRTT}$$



COMPUTER NETWORKS

2.5.4 Reliable Data Transfer

- IP is unreliable i.e. IP does not guarantee data delivery.
 - IP does not guarantee in-order delivery of data.
 - IP does not guarantee the integrity of the data.
- TCP creates a reliable data-transfer-service on top of IP's unreliable-service.
- At the receiver, reliable-service means
 - data-stream is uncorrupted
 - data-stream is without duplication and
 - data-stream is in sequence.

2.5.4.1 A Few Interesting Scenarios

2.5.4.1.1 First Scenario

- As shown in Figure 2.28, Host-A sends one segment to Host-B.
- Assume the acknowledgment from B to A gets lost.
- In this case, the timeout event occurs, and Host-A retransmits the same segment.
- When Host-B receives retransmission, it observes that the sequence-no has already been received.
- Thus, Host-B will discard the retransmitted-segment.

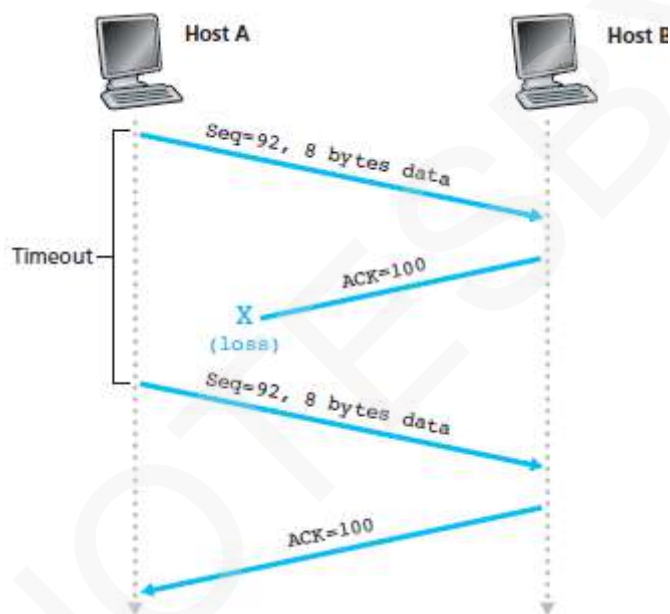


Figure 2.28: Retransmission due to a lost acknowledgment



COMPUTER NETWORKS

2.5.4.1.2 Second Scenario

- As shown in Figure 2.29, Host-A sends two segments back-to-back.
- Host-B sends two separate acknowledgments.
- Suppose neither of the acknowledgments arrives at Host-A before the timeout.
- When the timeout event occurs, Host-A resends the first-segment and restarts the timer.
- The second-segment will not be retransmitted until ACK for the second-segment arrives before the new timeout.

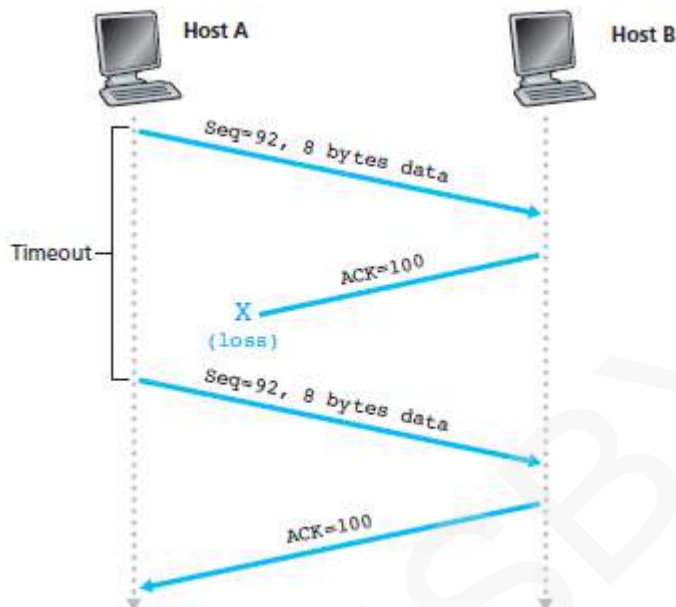


Figure 2.29: Segment 100 not retransmitted

2.5.4.1.3 Third Scenario

- As shown in Figure 2.30, Host-A sends the two segments.
- The acknowledgment of the first-segment is lost.
- But just before the timeout event, Host-A receives an acknowledgment-no 120.
- Therefore, Host-A knows that Host-B has received all the bytes up to 119.
- So, Host-A does not resend either of the two segments.

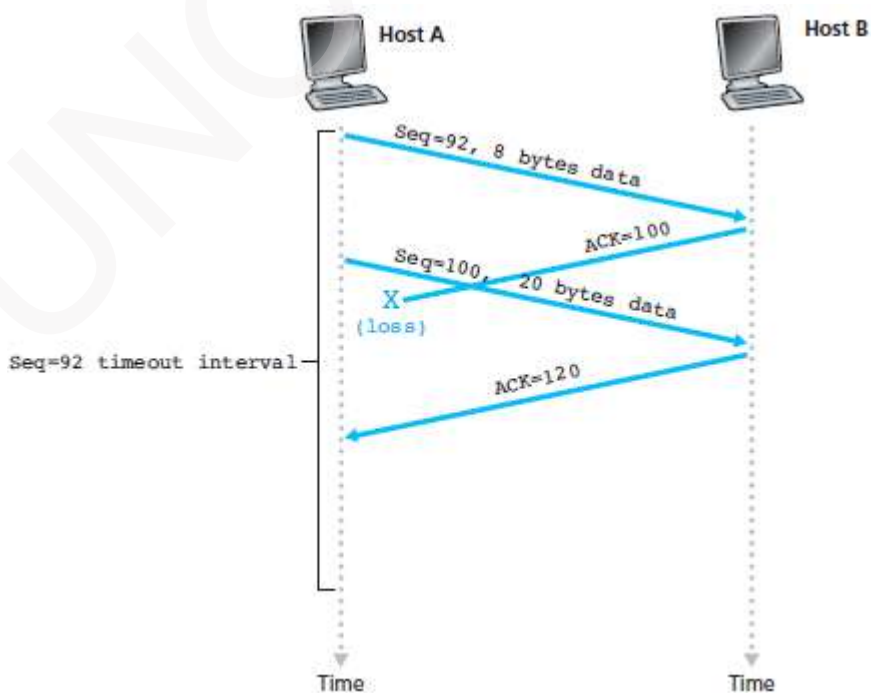


Figure 2.30: A cumulative acknowledgment avoids retransmission of the first-segment

“Faith is the bird that feels the light when the dawn is still dark.” —Rabindranath Tagore



COMPUTER NETWORKS

2.5.4.2 Fast Retransmit

- The timeout period can be relatively long.
- The sender can often detect packet-loss well before the timeout occurs by noting duplicate ACKs.
- A duplicate ACK refers to ACK the sender receives for the second time. (Figure 2.31).

Table 2.3: TCP ACK Generation Recommendation

Event	TCP Receiver Action
Arrival of in-order segment with expected sequence-number. All up to expected sequence-number already acknowledged.	Delayed ACK. Wait up to 500 msec for arrival of another in-order segment. If next in-order segment does not arrive in this interval, send an ACK.
Arrival of in-order segment with expected sequence-number. One other in-order segment waiting for ACK transmission.	Immediately send single cumulative ACK, ACKing both in-order segments.
Arrival of out-of-order segment with higher-than-expected sequence-number. Gap detected.	Immediately send duplicate ACK, indicating sequence-number of next expected-byte.
Arrival of segment that partially or completely fills in gap in received-data.	Immediately send ACK.

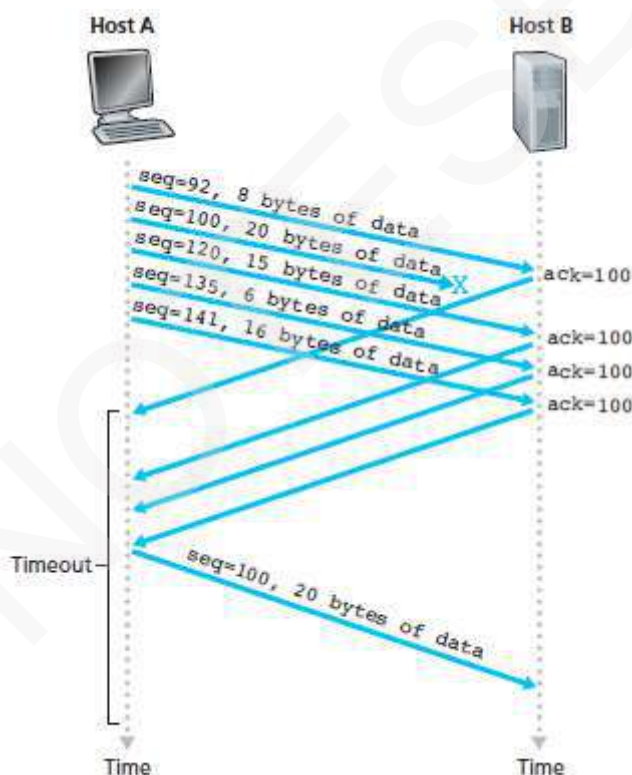


Figure 2.31: Fast retransmit: retransmitting the missing segment before the segment's timer expires



COMPUTER NETWORKS

2.5.5 Flow Control

- TCP provides a flow-control service to its applications.
- A flow-control service eliminates the possibility of the sender overflowing the receiver-buffer.

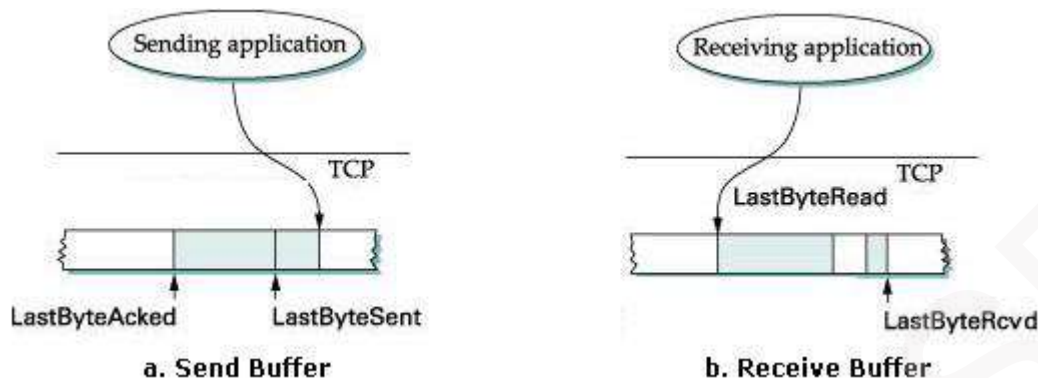


Figure 2.32: a) Send buffer and b) Receive Buffer

- As shown in Figure 2.32, we define the following variables:
 - 1) MaxSendBuffer: A send-buffer allocated to the sender.
 - 2) MaxRcvBuffer: A receive-buffer allocated to the receiver.
 - 3) LastByteSent: The no. of the last bytes sent to the send-buffer at the sender.
 - 4) LastByteAcked: The no. of the last bytes acknowledged in the send-buffer at the sender.
 - 5) LastByteRead: The no. of the last bytes read from the receive-buffer at the receiver.
 - 6) LastByteRcvd: The no. of the last bytes arrived & placed in receive-buffer at the receiver.

Send Buffer

- Sender maintains a send buffer, divided into 3 segments namely
 - 1) Acknowledged data
 - 2) Unacknowledged data and
 - 3) Data to be transmitted
- Send buffer maintains 2 pointers: LastByteAcked and LastByteSent. The relation b/w these two is:

$$\text{LastByteAcked} \leq \text{LastByteSent}$$

Receive Buffer

- Receiver maintains receive buffer to hold data even if it arrives out-of-order.
- Receive buffer maintains 2 pointers: LastByteRead and LastByteRcvd. The relation b/w these two is:

$$\text{LastByteRead} \leq \text{LastByteRcvd} + 1$$

Flow Control Operation

- Sender prevents overflowing of send buffer by maintaining

$$\text{LastByteWritten} - \text{LastByteAcked} \leq \text{MaxSendBuffer}$$
- Receiver avoids overflowing receive buffer by maintaining

$$\text{LastByteRcvd} - \text{LastByteRead} \leq \text{MaxRcvBuffer}$$
- Receiver throttles the sender by advertising a window that is smaller than the amount of free space that it can buffer as:

$$\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$



COMPUTER NETWORKS

2.5.6 TCP Connection Management

2.5.6.1 Connection Setup & Data Transfer

- To setup the connection, three segments are sent between the two hosts. Therefore, this process is referred to as a three-way handshake.
- Suppose a client-process wants to initiate a connection with a server-process.
- Figure 2.33 illustrates the steps involved:

Step 1: Client sends a connection-request segment to the Server

- The client first sends a connection-request segment to the server.
- The connection-request segment contains:
 - 1) SYN bit is set to 1.
 - 2) Initial sequence-number (client_isn).
- The SYN segment is encapsulated within an IP datagram and sent to the server.

Step 2: Server sends a connection-granted segment to the Client

- Then, the server
 - extracts the SYN segment from the datagram
 - allocates the buffers and variables to the connection and
 - sends a connection-granted segment to the client.
- The connection-granted segment contains:
 - 1) SYN bit is set to 1.
 - 2) Acknowledgment field is set to client_isn+1.
 - 3) Initial sequence-number (server_isn).

Step 3: Client sends an ACK segment to the Server

- Finally, the client
 - allocates buffers and variables to the connection and
 - sends an ACK segment to the server
- The ACK segment acknowledges the server.
- SYN bit is set to zero, since the connection is established.

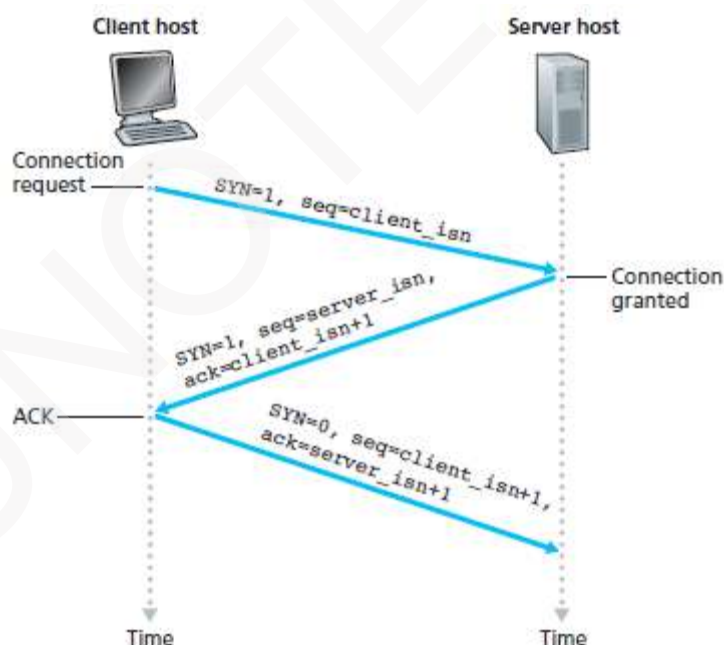


Figure 2.33: TCP three-way handshake: segment exchange



COMPUTER NETWORKS

2.5.6.2 Connection Release

- Either of the two processes in a connection can end the connection.
- When a connection ends, the “resources” in the hosts are de-allocated.
- Suppose the client decides to close the connection.
- Figure 2.34 illustrates the steps involved:
 - 1) The client-process issues a close command.
 - × Then, the client sends a shutdown-segment to the server.
 - × This segment has a FIN bit set to 1.
 - 2) The server responds with an acknowledgment to the client.
 - 3) The server then sends its own shutdown-segment.
 - × This segment has a FIN bit set to 1.
 - 4) Finally, the client acknowledges the server’s shutdown-segment.

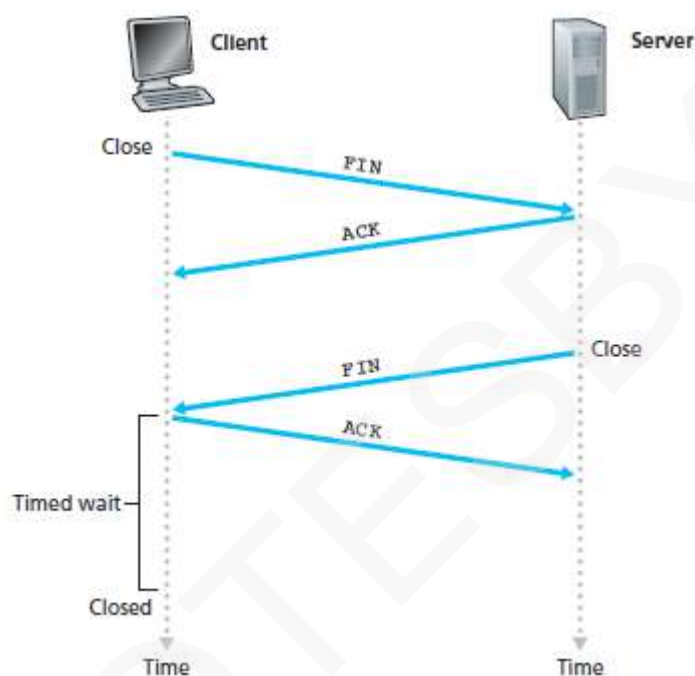


Figure 2.34: Closing a TCP connection



COMPUTER NETWORKS

2.6 Principles of Congestion Control

2.6.1 The Causes and the Costs of Congestion

2.6.1.1 Scenario 1: Two Senders, a Router with Infinite Buffers

- Two hosts (A & B) have a connection that shares a single-hop b/w source & destination.
- This is illustrated in Figure 2.35.

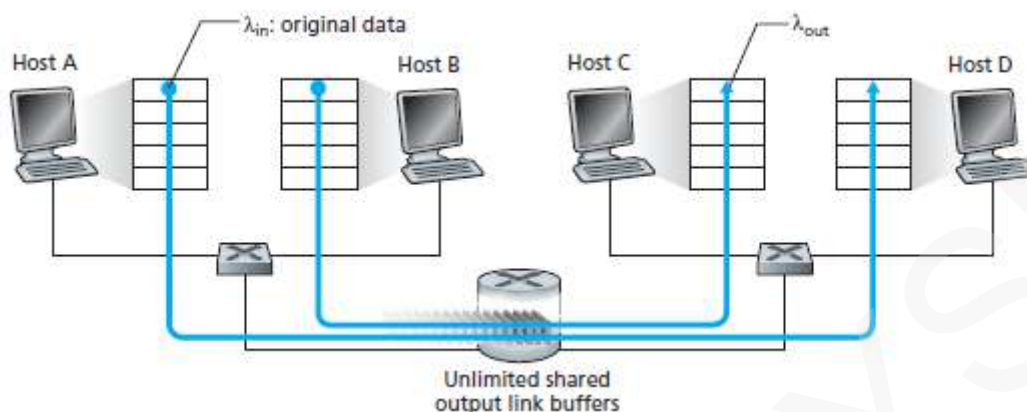


Figure 2.35: Congestion scenario 1: Two connections sharing a single hop with infinite buffers

- Let
 - Sending-rate of Host-A = λ_{in} bytes/sec
 - Outgoing Link's capacity = R
- Packets from Hosts A and B pass through a router and over a shared outgoing link.
- The router has buffers.
- The buffers stores incoming packets when packet-arrival rate exceeds the outgoing link's capacity.

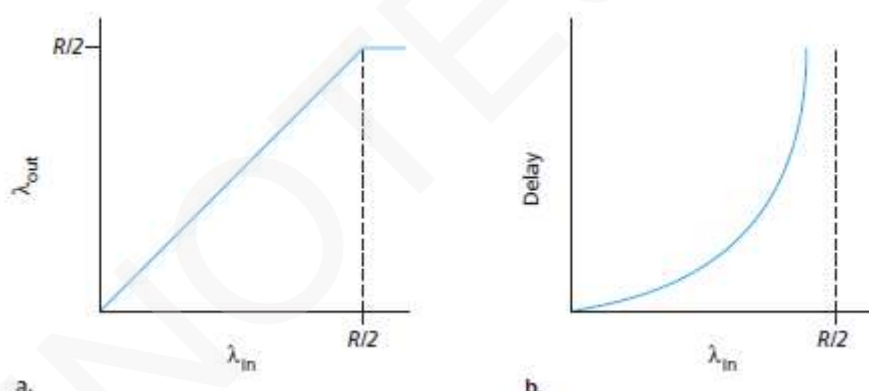


Figure 2.36: Congestion scenario 1: Throughput and delay as a function of host sending-rate

- Figure 2.36 plots the performance of Host-A's connection.

Left Hand Graph

- The left graph plots the per-connection throughput as a function of the connection-sending-rate.
- For a sending-rate b/w 0 and $R/2$, the throughput at the receiver equals the sender's sending-rate. However, for a sending-rate above $R/2$, the throughput at the receiver is only $R/2$. (Figure 2.36a)
- Conclusion: The link cannot deliver packets to a receiver at a steady-state rate that exceeds $R/2$.

Right Hand Graph

- The right graph plots the average delay as a function of the connection-sending-rate (Figure 2.36b).
- As the sending-rate approaches $R/2$, the average delay becomes larger and larger. However, for a sending-rate above $R/2$, the average delay becomes infinite.
- Conclusion: Large queuing delays are experienced as the packet arrival rate nears the link capacity.

**COMPUTER NETWORKS****2.6.1.2 Scenario 2: Two Senders and a Router with Finite Buffers**

- Here, we have 2 assumptions (Figure 2.37):

- The amount of router buffering is finite.
 - Packets will be dropped when arriving to an already full buffer.
- Each connection is reliable.
 - If a packet is dropped at the router, the sender will eventually retransmit it.

- Let

Application's sending-rate of Host-A = λ_{in} bytes/sec

Transport-layer's sending-rate of Host-A = λ'_{in} bytes/sec (also called offered-load to network)

Outgoing Link's capacity = R

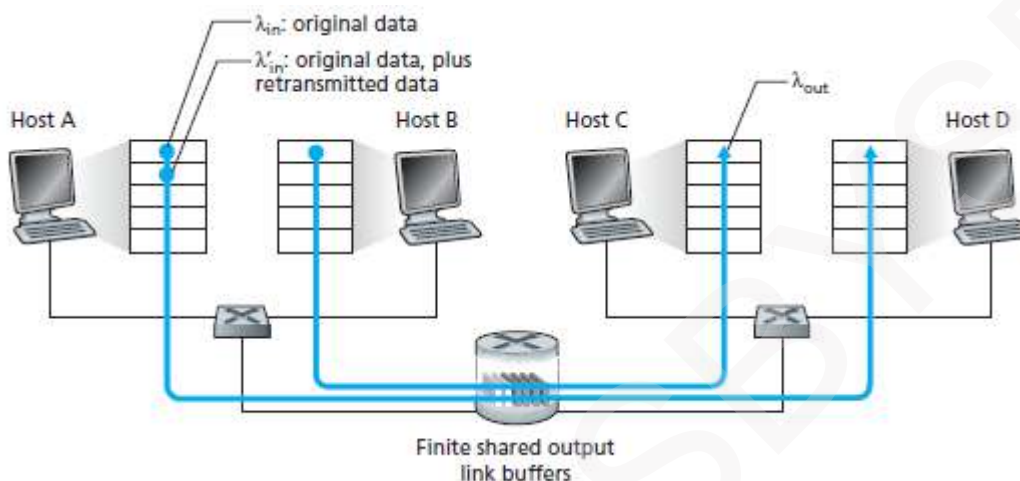


Figure 2.37: Scenario 2: Two hosts (with retransmissions) and a router with finite buffers

Case 1 (Figure 2.38(a)):

- Host-A sends a packet only when a buffer is free.
- In this case,
 - no loss occurs
 - λ_{in} will be equal to λ'_{in} , and
 - throughput of the connection will be equal to λ_{in} .
- The sender retransmits only when a packet is lost.
- Consider the offered-load $\lambda'_{in} = R/2$.
- The rate at which data are delivered to the receiver application is $R/3$.
- The sender must perform retransmissions to compensate for lost packets due to buffer overflow.

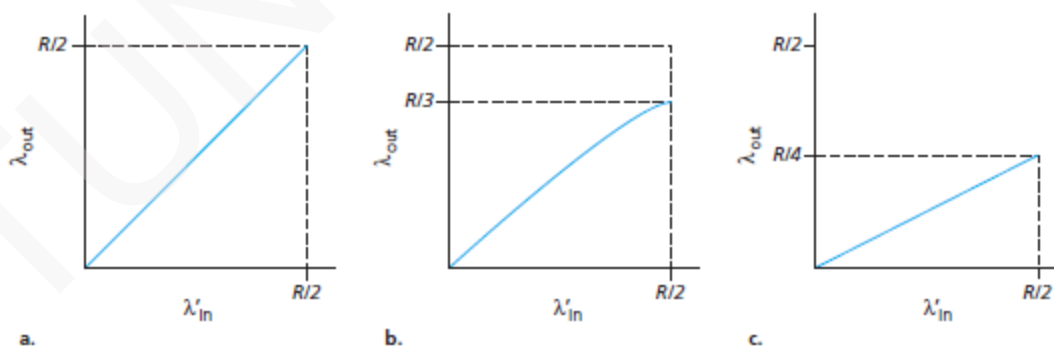


Figure 2.38: Scenario 2 performance with finite buffers

Case 3 (Figure 2.38(c)):

- The sender may time out & retransmit a packet that has been delayed in the queue but not yet lost.
- Both the original data packet and the retransmission may reach the receiver.
- The receiver needs one copy of this packet and will discard the retransmission.
- The work done by the router in forwarding the retransmitted copy of the original packet was wasted.

**COMPUTER NETWORKS****2.6.1.3 Scenario 3: Four Senders, Routers with Finite Buffers, and Multihop Paths**

- Four hosts transmit packets, each over overlapping two-hop paths.
- This is illustrated in Figure 2.39.

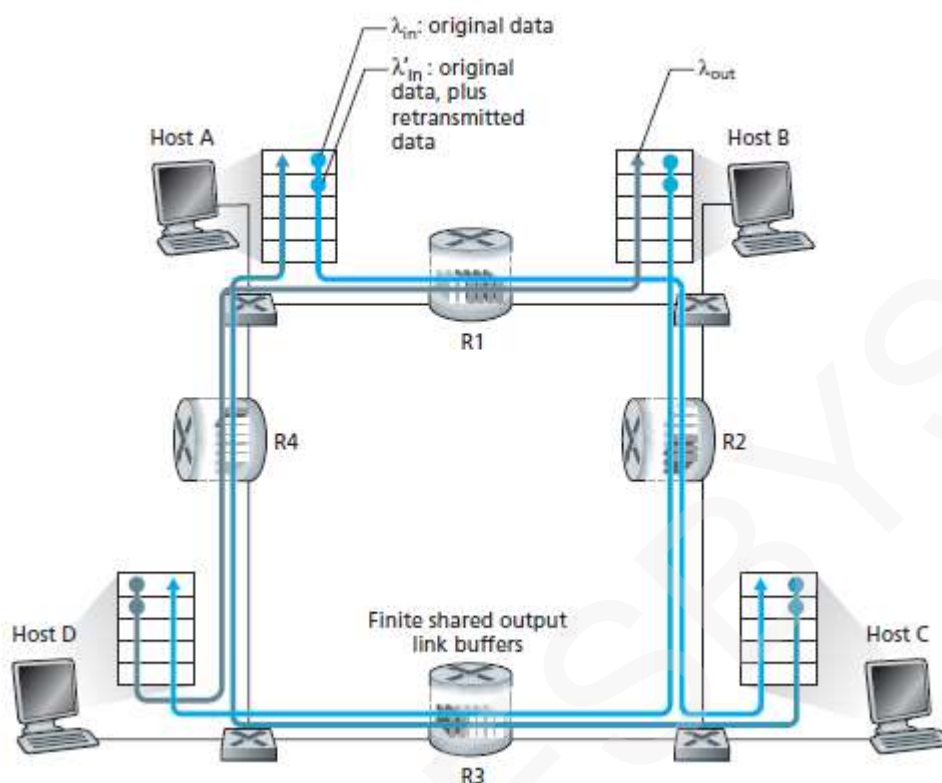


Figure 2.39: Four senders, routers with finite buffers, and multihop paths

- Consider the connection from Host-A to Host C, passing through routers R1 and R2.
- The A-C connection
 - shares router R1 with the D-B connection and
 - shares router R2 with the B-D connection.
- Case-1: For extremely small values of λ_{in} ,
 - buffer overflows are rare (as in congestion scenarios 1 and 2) and
 - the throughput approximately equals the offered-load.
- Case-2: For slightly larger values of λ_{in} , the corresponding throughput is also larger. This is because
 - more original data is transmitted into the network
 - data is delivered to the destination and
 - overflows are still rare.
- Case-3: For extremely larger values of λ_{in} .
 - Consider router R2.
 - The A-C traffic arriving to router R2 can have an arrival rate of at most R regardless of the value of λ_{in} .
 - where $R =$ the capacity of the link from R1 to R2,.
 - If λ_{in}' is extremely large for all connections, then the arrival rate of B-D traffic at R2 can be much larger than that of the A-C traffic.
 - The A-C and B-D traffic must compete at router R2 for the limited amount of buffer-space.
 - Thus, the amount of A-C traffic that successfully gets through R2 becomes smaller and smaller as the offered-load from B-D gets larger and larger.
 - In the limit, as the offered-load approaches infinity, an empty buffer at R2 is immediately filled by a B-D packet, and the throughput of the A-C connection at R2 goes to zero.
 - When a packet is dropped along a path, the transmission capacity ends up having been wasted.



COMPUTER NETWORKS

2.6.2 Approaches to Congestion Control

• Congestion-control approaches can be classified based on whether the network-layer provides any explicit assistance to the transport-layer:

1) End-to-end Congestion Control

- The network-layer provides no explicit support to the transport-layer for congestion-control.
- Even the presence of congestion must be inferred by the end-systems based only on observed network-behavior.
- Segment loss is taken as an indication of network-congestion and the window-size is decreased accordingly.

2) Network Assisted congestion Control

- Network-layer components provide explicit feedback to the sender regarding congestion.
- This feedback may be a single bit indicating congestion at a link.
- Congestion information is fed back from the network to the sender in one of two ways:
 - i) Direct feedback may be sent from a network-router to the sender (Figure 2.40).
 - ✗ This form of notification typically takes the form of a choke packet.
 - ii) A router marks a field in a packet flowing from sender to receiver to indicate congestion.
 - ✗ Upon receipt of a marked packet, the receiver then notifies the sender of the congestion indication.
 - ✗ This form of notification takes at least a full round-trip time.

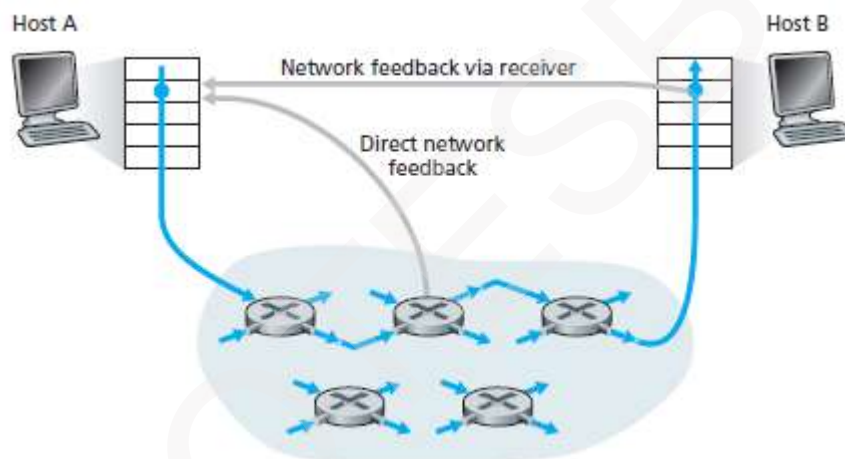


Figure 2.40: Two feedback pathways for network-induced congestion information



COMPUTER NETWORKS

2.6.3 Network Assisted Congestion Control Example: ATM ABR Congestion Control

- ATM (Asynchronous Transfer Mode) protocol uses network-assisted approach for congestion-control.
- ABR (Available Bit Rate) has been designed as an elastic data-transfer-service.
 - i) When the network is underloaded, ABR has to take advantage of the spare available bandwidth.
 - ii) When the network is congested, ABR should reduce its transmission-rate.

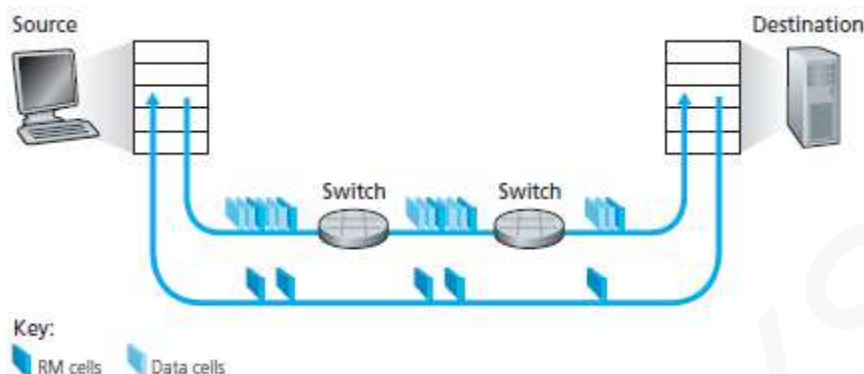


Figure 2.41: Congestion-control framework for ATM ABR service

- Figure 2.41 shows the framework for ATM ABR congestion-control.
- Data-cells are transmitted from a source to a destination through a series of intermediate switches.
- RM-cells are placed between the data-cells. (RM → Resource Management).
- The RM-cells are used to send congestion-related information to the hosts & switches.
- When an RM-cell arrives at a destination, the cell will be sent back to the sender
- Thus, RM-cells can be used to provide both
 - direct network feedback and
 - network feedback via the receiver.

2.6.3.1 Three Methods to indicate Congestion

- ATM ABR congestion-control is a rate-based approach.
- ABR provides 3 mechanisms for indicating congestion-related information:
 - 1) EFCI Bit**
 - Each data-cell contains an EFCI bit. (EFCI → Explicit forward congestion indication)
 - A congested-switch sets the EFCI bit to 1 to signal congestion to the destination.
 - The destination must check the EFCI bit in all received data-cells.
 - If the most recently received data-cell has the EFCI bit set to 1, then the destination
 - sets the CI bit to 1 in the RM-cell (CI → congestion indication)
 - sends the RM-cell back to the sender.
 - Thus, a sender can be notified about congestion at a network switch.
 - 2) CI and NI Bits**
 - The rate of RM-cell interspersion is a tunable parameter.
 - The default value is one RM-cell every 32 data-cells. (NI → No Increase)
 - The RM-cells have a CI bit and a NI bit that can be set by a congested-switch.
 - A switch
 - sets the NI bit to 1 in a RM-cell under mild congestion and
 - sets the CI bit to 1 under severe congestion conditions.
 - 3) ER Setting**
 - Each RM-cell also contains an ER field. (ER → explicit rate)
 - A congested-switch may lower the value contained in the ER field in a passing RM-cell.
 - In this manner, ER field will be set to minimum supportable rate of all switches on the path.



COMPUTER NETWORKS

2.7 TCP Congestion Control

2.7.1 TCP Congestion Control

- TCP has congestion-control mechanism.
- TCP uses end-to-end congestion-control rather than network-assisted congestion-control
- Here is how it works:
 - Each sender limits the rate at which it sends traffic into its connection as a function of perceived congestion.
 - i) If sender perceives that there is little congestion, then sender increases its data-rate.
 - ii) If sender perceives that there is congestion, then sender reduces its data-rate.
- This approach raises three questions:
 - 1) How does a sender limit the rate at which it sends traffic into its connection?
 - 2) How does a sender perceive that there is congestion on the path?
 - 3) What algorithm should the sender use to change its data-rate?
- The sender keeps track of an additional variable called the congestion-window (cwnd).
- The congestion-window imposes a constraint on the data-rate of a sender.
- The amount of unacknowledged-data at a sender will not exceed minimum of (cwnd & rwnd), that is:
$$\text{LastByteSent} - \text{LastByteAcked} \leq \min\{\text{cwnd}, \text{rwnd}\}$$
- The sender's data-rate is roughly cwnd/RTT bytes/sec.
- Explanation of Loss event:
 - A "loss event" at a sender is defined as the occurrence of either
 - timeout or
 - receipt of 3 duplicate ACKs from the receiver.
 - Due to excessive congestion, the router-buffer along the path overflows. This causes a datagram to be dropped.
 - The dropped datagram, in turn, results in a loss event at the sender.
 - The sender considers the loss event as an indication of congestion on the path.
- How congestion is detected?
 - Consider the network is congestion-free.
 - Acknowledgments for previously unacknowledged segments will be received at the sender.
 - TCP
 - will take the arrival of these acknowledgments as an indication that all is well and
 - will use acknowledgments to increase the window-size (& hence data-rate).
 - TCP is said to be self-clocking because
 - acknowledgments are used to trigger the increase in window-size
 - Congestion-control algorithm has 3 major components:
 - 1) Slow start
 - 2) Congestion avoidance and
 - 3) Fast recovery.



COMPUTER NETWORKS

2.7.1.1 Slow Start

- When a TCP connection begins, the value of $cwnd$ is initialized to 1 MSS.
- TCP doubles the number of packets sent every RTT on successful transmission.
- Here is how it works:
 - As shown in Figure 2.42, the TCP
 - sends the first-segment into the network and
 - waits for an acknowledgment.
 - When an acknowledgment arrives, the sender
 - increases the congestion-window by one MSS and
 - sends out 2 segments.
 - When two acknowledgments arrive, the sender
 - increases the congestion-window by one MSS and
 - sends out 4 segments.
 - This process results in a doubling of the sending-rate every RTT.
- Thus, the TCP data-rate starts slow but grows exponentially during the slow start phase.

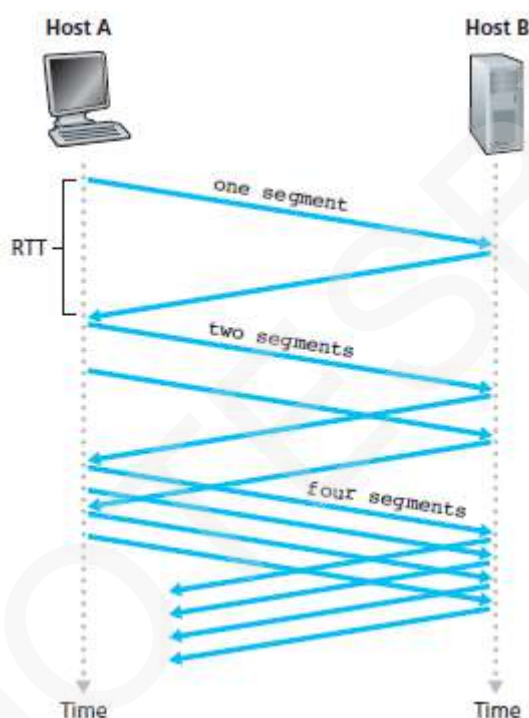


Figure 2.42: TCP slow start

- When should the exponential growth end?
 - Slow start provides several answers to this question.
 - 1) If there is a loss event, the sender
 - sets the value of $cwnd$ to 1 and
 - begins the slow start process again. ($ssthresh \rightarrow$ "slow start threshold")
 - sets the value of $ssthresh$ to $cwnd/2$.
 - 2) When the value of $cwnd$ equals $ssthresh$, TCP enters the congestion avoidance state.
 - 3) When three duplicate ACKs are detected, TCP
 - performs a fast retransmit and
 - enters the fast recovery state.
- TCP's behavior in slow start is summarized in FSM description in Figure 2.43.

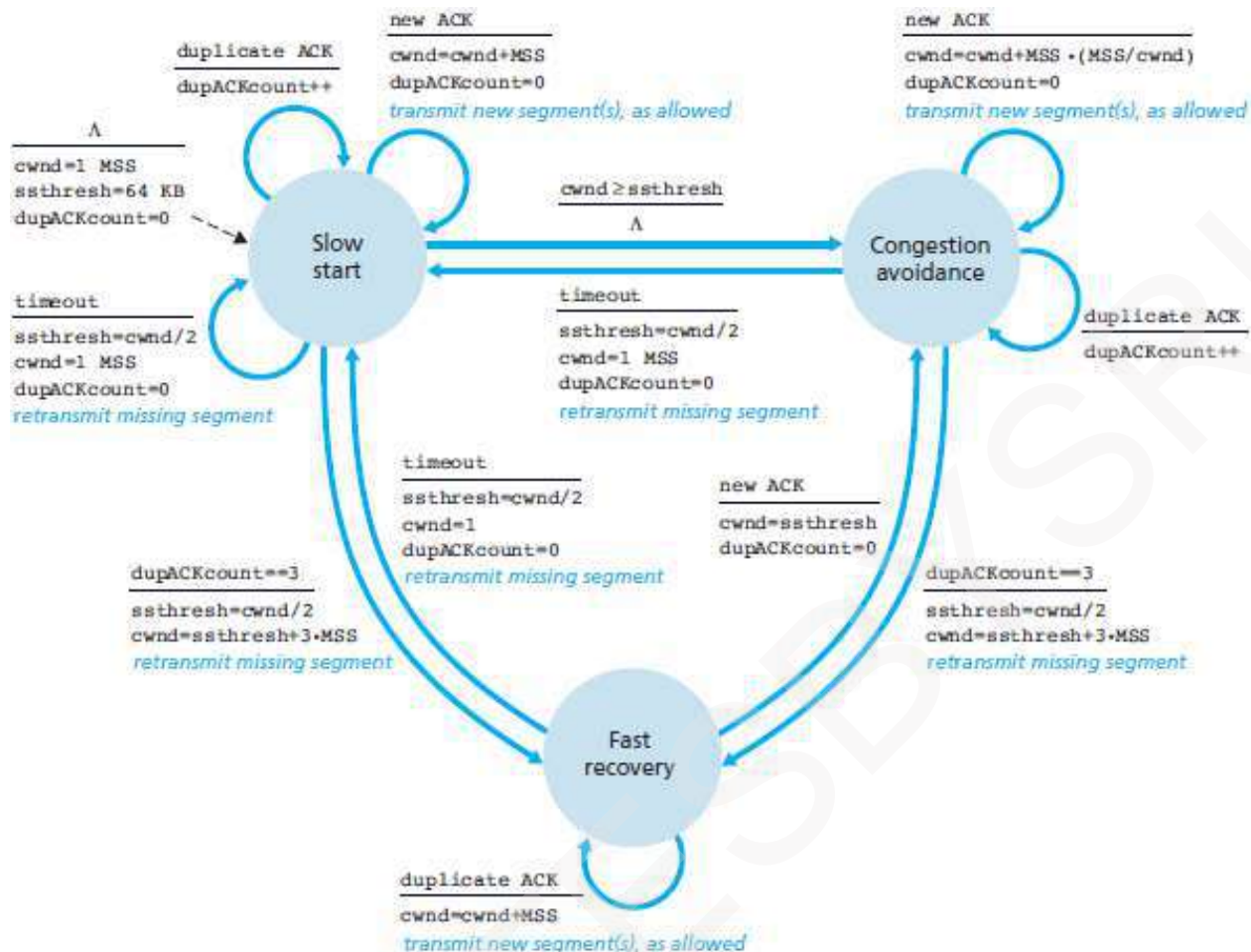


Figure 2.43: FSM description of TCP congestion-control



COMPUTER NETWORKS

2.7.1.2 Congestion Avoidance

- On entry to congestion-avoidance state, the value of cwnd is approximately half its previous value.
- Thus, the value of cwnd is increased by a single MSS every RTT.
- The sender must increase cwnd by MSS bytes ($MSS/cwnd$) whenever a new acknowledgment arrives
- When should linear increase (of 1 MSS per RTT) end?

1) When a timeout occurs.

- When the loss event occurred,
 - value of cwnd is set to 1 MSS and
 - value of ssthresh is set to half the value of cwnd.

2) When triple duplicate ACK occurs.

- When the triple duplicate ACKs were received,
 - value of cwnd is halved.
 - value of ssthresh is set to half the value of cwnd.

2.7.1.3 Fast Recovery

- The value of cwnd is increased by 1 MSS for every duplicate ACK received.
- When an ACK arrives for the missing segment, the congestion-avoidance state is entered.
- If a timeout event occurs, fast recovery transitions to the slow-start state.
- When the loss event occurred
 - value of cwnd is set to 1 MSS, and
 - value of ssthresh is set to half the value of cwnd.

- There are 2 versions of TCP:

1) TCP Tahoe

- An early version of TCP was known as TCP Tahoe.
- TCP Tahoe
 - cut the congestion-window to 1 MSS and
 - entered the slow-start phase after either
 - i) timeout-indicated or
 - ii) triple-duplicate-ACK-indicated loss event.

2) TCP Reno

- The newer version of TCP is known as TCP Reno.
- TCP Reno incorporated fast recovery.
- Figure 2.44 illustrates the evolution of TCP's congestion-window for both Reno and Tahoe.

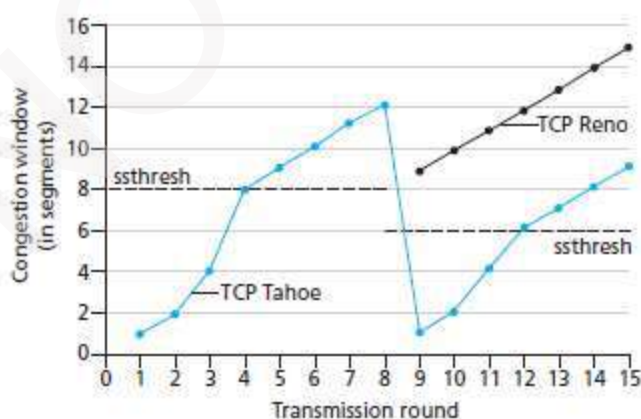


Figure 2.44: Evolution of TCP's congestion-window (Tahoe and Reno)



COMPUTER NETWORKS

2.7.1.4 TCP Congestion Control: Retrospective

- TCP's congestion-control consists of (AIMD → additive increase, multiplicative decrease)
 - Increasing linearly (additive) value of cwnd by 1 MSS per RTT and
 - Halving (multiplicative decrease) value of cwnd on a triple duplicate-ACK event.
- For this reason, TCP congestion-control is often referred to as an AIMD.
- AIMD congestion-control gives rise to the "saw tooth" behavior shown in Figure 2.45.
- TCP
 - increases linearly the congestion-window-size until a triple duplicate-ACK event occurs and
 - decreases then the congestion-window-size by a factor of 2



Figure 2.45: Additive-increase, multiplicative-decrease congestion-control



COMPUTER NETWORKS

2.7.2 Fairness

- Congestion-control mechanism is fair if each connection gets equal share of the link-bandwidth.
- As shown in Figure 2.46, consider 2 TCP connections sharing a single link with transmission-rate R .
- Assume the two connections have the same MSS and RTT.

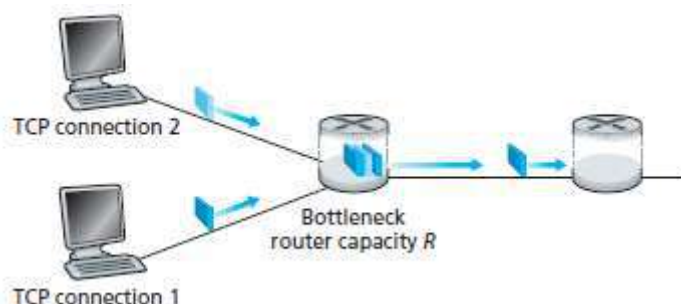


Figure 2.46: Two TCP connections sharing a single bottleneck link

- Figure 2.47 plots the throughput realized by the two TCP connections.
 - If TCP shares the link-bandwidth equally b/w the 2 connections, then the throughput falls along the 45-degree arrow starting from the origin.

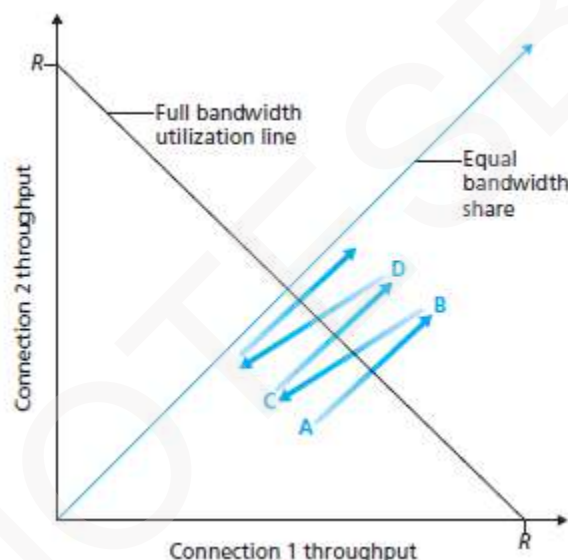


Figure 2.47: Throughput realized by TCP connections 1 and 2

2.7.2.1 Fairness and UDP

- Many multimedia-applications (such as Internet phone) often do not run over TCP.
- Instead, these applications prefer to run over UDP. This is because
 - applications can pump their audio into the network at a constant rate and
 - occasionally lose packets.

2.7.2.2 Fairness and Parallel TCP Connections

- Web browsers use multiple parallel-connections to transfer the multiple objects within a Web page.
- Thus, the application gets a larger fraction of the bandwidth in a congested link.
- ∴ Web-traffic is so pervasive in the Internet; multiple parallel-connections are common nowadays.

**MODULE-WISE QUESTIONS****PART 1**

- 1) With a diagram, explain multiplexing and demultiplexing. (6*)
- 2) Explain the significance of source and destination-port-no in a segment. (4*)
- 3) With a diagram, explain connectionless multiplexing and demultiplexing. (4)
- 4) With a diagram, explain connection oriented multiplexing and demultiplexing. (4)
- 5) Briefly explain UDP & its services. (6*)
- 6) With general format, explain various field of UDP segment. Explain how checksum is calculated (8*)
- 7) With a diagram, explain the working of rdt1.0. (6)
- 8) With a diagram, explain the working of rdt2.0. (6*)
- 9) With a diagram, explain the working of rdt2.1. (6)
- 10) With a diagram, explain the working of rdt3.0. (6*)
- 11) With a diagram, explain the working of Go-Back-N. (6*)
- 12) With a diagram, explain the working of selective repeat. (6*)
- 13) Explain the following terms: (8)
 - i) Sequence-number
 - ii) Acknowledgment
 - iii) Negative acknowledgment
 - iv) Window, pipelining

PART 2

- 14) Briefly explain TCP & its services. (6*)
- 15) With general format, explain various field of TCP segment. (6*)
- 16) With a diagram, explain the significance of sequence and acknowledgment numbers. (4*)
- 17) With a diagram, explain the reliable data transfer with few interesting scenarios. (8)
- 18) With a diagram, explain fast retransmit in TCP. (6*)
- 19) With a diagram, explain flow Control in TCP. (6)
- 20) With a diagram, explain connection management in TCP. (8*)
- 21) With a diagram, explain the causes of congestion with few scenarios. (8)
- 22) Briefly explain approaches to congestion control. (6*)
- 23) With a diagram, explain ATM ABR congestion control. (8)
- 24) With a diagram, explain slow start in TCP. (6*)
- 25) With a diagram, explain fast recovery in TCP. (6*)



MODULE 3: THE NETWORK LAYER

- 3.1 Introduction
 - 3.1.1 Forwarding & Routing
 - 3.1.2 Network Service Models
- 3.2 Virtual Circuit & Datagram Networks
 - 3.2.1 Virtual Circuit Networks
 - 3.2.2 Datagram Networks
 - 3.2.3 Comparison of Virtual Circuit & Datagram Networks
- 3.3 What's Inside a Router?
 - 3.3.1 Switching
 - 3.3.1.1 Switching via Memory
 - 3.3.1.2 Switching via a Bus
 - 3.3.1.3 Switching via an Interconnection Network
 - 3.3.2 Output Processing
 - 3.3.3 Where Does Queueing Occur?
- 3.4 IP: Forwarding & Addressing in the Internet
 - 3.4.1 IPv4 Datagram Format
 - 3.4.2 Fragmentation
 - 3.4.2.1 Maximum Transfer Unit
 - 3.4.2.2 Fields Related to Fragmentation & Reassembly
 - 3.4.3 IPv4 Addressing
 - 3.4.3.1 IPv4 Classful Addressing
 - 3.4.3.1.1 Subnet Addressing
 - 3.4.3.2 CIDR
 - 3.4.3.3 Obtaining a Block of Addresses
 - 3.4.3.4 Obtaining a Host Address: DHCP
 - 3.4.3.4.1 DHCP Protocol
 - 3.4.3.5 NAT
 - 3.4.4 ICMP
 - 3.4.5 IPv6
 - 3.4.5.1 Changes from IPv4 to IPv6
 - 3.4.5.2 IPv6 Datagram Format
 - 3.4.5.3 IPv4 Fields not present in IPv6
 - 3.4.5.4 Difference between IPv4 & IPv6
 - 3.4.5.5 Transitioning from IPv4 to IPv6
 - 3.4.5.5.1 Dual Stack Approach
 - 3.4.5.5.2 Tunneling
 - 3.4.6 A Brief Foray into IP Security
- 3.5 Routing Algorithms
 - 3.5.1 Routing Algorithm Classification
 - 3.5.1.1 Global or Decentralized
 - 3.5.1.2 Static or Dynamic
 - 3.5.1.3 Load Sensitive or Load Insensitive
 - 3.5.2 LS Routing Algorithm
 - 3.5.2.1 Dijkstra's Algorithm
 - 3.5.3 DV Routing Algorithm
 - 3.5.3.1 Bellman Ford Algorithm
 - 3.5.4 A Comparison of LS and DV Routing-algorithms
 - 3.5.5 Hierarchical Routing
 - 3.5.5.1 Intra-AS Routing Protocol
 - 3.5.5.2 Intra-AS Routing Protocol



COMPUTER NETWORKS

3.6 Routing in the Internet

3.6.1 Intra-AS Routing in the Internet: RIP

3.6.1.1 RIP Protocol

3.6.2 Intra-AS Routing in the Internet: OSPF

3.6.3 Inter-AS Routing: BGP

3.6.3.1 Basics

3.6.3.2 Path Attributes & Routes

3.6.3.3 Route Selection

3.6.3.4 Routing Policy

3.7 Broadcast & Multicast Routing

3.7.1 Broadcast Routing Algorithms

3.7.1.1 N-way Unicast

3.7.1.2 Uncontrolled Flooding

3.7.1.3 Controlled Flooding

3.7.1.4 Spanning - Tree Broadcast

3.7.1.4.1 Center Based Approach

3.7.2 Multicast

3.7.2.1 IGMP

3.7.2.2 Multicast Routing Algorithms

3.7.2.3 Multicast Routing in the Internet



MODULE 3: THE NETWORK LAYER

3.1 Introduction

3.1.1 Forwarding & Routing

- The role of the network-layer is to move packets from a sending-host to a receiving-host.
- Two important functions of network-layer:
 - 1) Forwarding**
 - Forwarding refers to transferring a packet from incoming-link to outgoing-link within a router.
 - Forwarding is a router-local action.
 - 2) Routing**
 - Routing means determining the path taken by packets from a sender to a receiver.
 - Routing is a network-wide process.

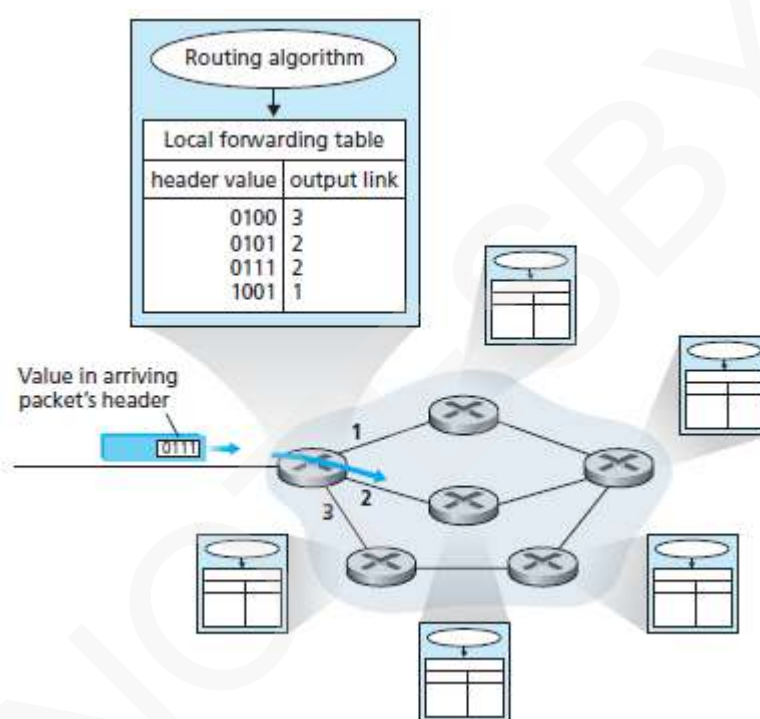


Figure 3.1: Routing-algorithms determine values in forwarding-tables

- The algorithms that determine the paths are referred to as routing-algorithms.
- Each router has a forwarding-table.
- As shown in Figure 3.1, a forwarding-table contains 2 columns:
 - 1) Header value and
 - 2) Output link.
- How forwarding is done?
 - 1) Firstly, a router examines the header-value of an arriving packet.
 - 2) Then, the router uses the header-value to index into the forwarding-table.
 - 3) Finally, the router forwards the packet.



COMPUTER NETWORKS

3.1.2 Network Service Models

- This defines the characteristics of end-to-end transport of packets b/w sending & receiving systems.
- The network-layer provides following services:
 - 1) Guaranteed Delivery**
 - This service guarantees that the packet will eventually arrive at its destination.
 - 2) Guaranteed Delivery with Bounded Delay**
 - This service guarantees delivery of the packet within a specified host-to-host delay bound.
 - 3) In-order Packet Delivery**
 - This service guarantees packets arrive at the destination in the order that they were sent.
 - 4) Guaranteed Minimal Bandwidth**
 - This service imitates the behavior of a link of a specified bit rate b/w sender & receiver.
 - 5) Guaranteed Maximum Jitter**
 - This service guarantees
The amount of time b/w the transmissions of 2 successive packets at the sender is equal to the amount of time b/w their receipts at the destination.
 - 6) Security Services**
 - The network-layer provides security-services such as
 - confidentiality
 - data-integrity and
 - source-authentication.
- How confidentiality is provided?
 - 1) Using a secret key, the sender encrypts the data being sent to the receiver.
 - 2) Then, the receiver decrypts the data using the same secret key.



COMPUTER NETWORKS

3.2 Virtual Circuit & Datagram Networks

- A network-layer provides 2 types of services:
 - 1) Connectionless service and
 - 2) Connection-oriented service.
- Two categories of computer-networks:
 - 1) Virtual Circuit (VC) Networks**
 - This provides only a connection-oriented service at the network-layer.
 - 2) Datagram Networks**
 - This provides only a connectionless service at the network-layer.
 - For example: The Internet.



COMPUTER NETWORKS

3.2.1 Virtual Circuit Networks

- A VC consists of
 - 1) A path between the source and destination.
 - 2) VC number: This is one number for each link along the path.
 - 3) Entries in the forwarding-table in each router.
- A packet belonging to a virtual-circuit will carry a VC number in its header.
- At intervening router, the VC number of traversing packet is replaced with a new VC number.
- The new VC number is obtained from the forwarding-table.

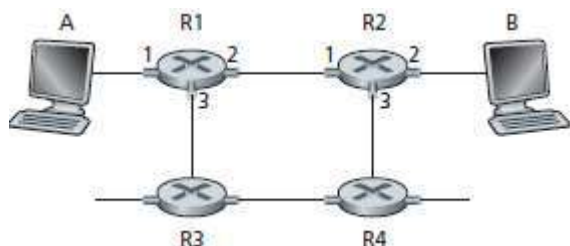


Figure 3.2: A simple virtual-circuit network

Incoming Interface	Incoming VC #	Outgoing Interface	Outgoing VC #
1	12	2	22
2	63	1	18
3	7	2	17
1	97	3	87

Table 3.1: Forwarding-table in R1

- Q: How does router determine the replacement VC number for a packet traversing the router?
Answer: Each router's forwarding-table includes VC number translation.
- The forwarding-table in R1 is shown in Table 3.1.

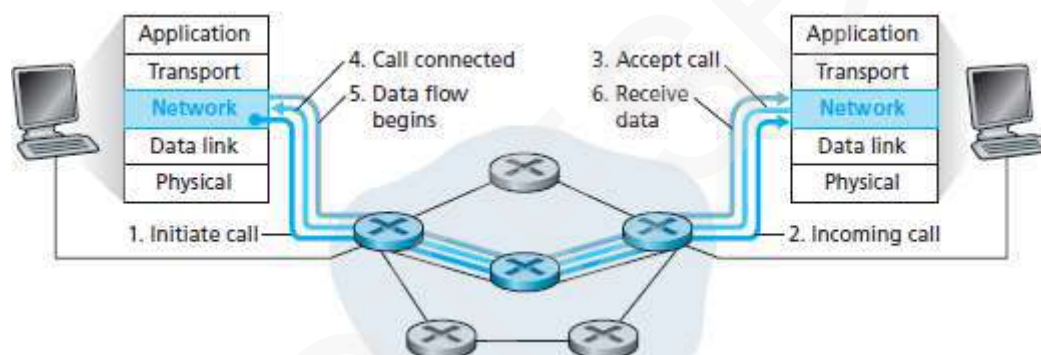


Figure 3.3: Virtual-circuit setup

- Why a packet does not use the same VC number on each link along the path?
Answer: 1) Replacing the number from link to link reduces length of the VC field in the packet-header.
2) VC setup is simplified by permitting a different VC number at each link along the path.

- Disadvantage:
The routers must maintain connection state information for the ongoing connections.
- Three phases in a virtual-circuit (Figure 3.3):

1) VC Setup

- During the setup phase, the sending transport-layer
 - contacts the network-layer
 - specifies the receiver's address and
 - waits for the network to set-up the VC.
- The network-layer determines the path between sender and receiver.
- The network-layer also determines the VC number for each link along the path.
- Finally, the network-layer adds an entry in the forwarding-table in each router.
- During VC setup, the network-layer may also reserve resources.

2) Data Transfer

- Once the VC has been established, packets can begin to flow along the VC.

3) VC Teardown

- This is initiated when the sender/receiver wants to terminate the VC.
- The network-layer
 - informs the other end-system of the call termination and
 - removes the appropriate entries in the forwarding-table in each router.

“Do what you have always done and you will get what you have always got.” —Sue Knight



COMPUTER NETWORKS

3.2.2 Datagram Networks

- The source attaches the packet with the address of the destination.
- The packets are injected into the network.
- The packets are routed independent of each other.
- No advance circuit setup is needed. So, routers do not maintain any connection state information.
- As a packet is transmitted from source to destination, it passes through a series of routers.
- Each router uses the packet's destination-address to forward the packet.

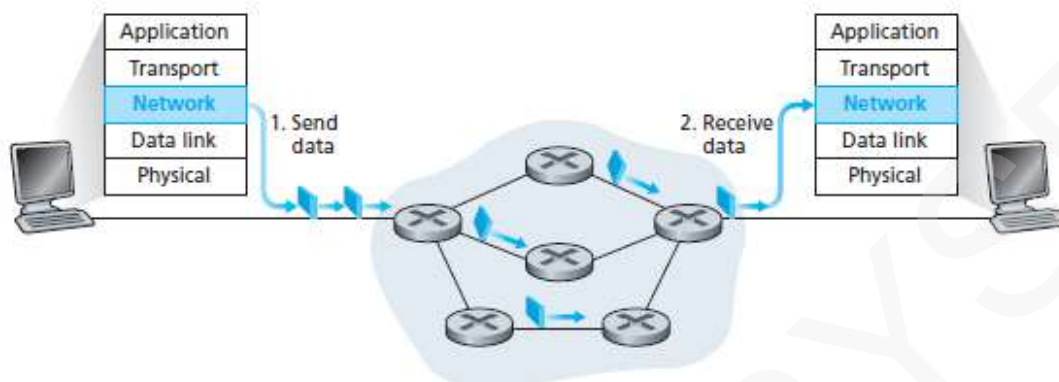


Figure 3.4: Datagram network

- Suppose the router R1 has four links, numbered 0 through 3 (Figure 3.4).
- Forwarding-table of R1 is as follows (Table 3.2):

Prefix Match	Link Interface
11001000 00010111 00010	0
11001000 00010111 00011000	1
11001000 00010111 00011	2
otherwise	3

Table 3.2: Forwarding-table of R1

- The router matches a prefix of the packet's destination-address with the entries in the table;
 - 1) If both are equal, the router forwards the packet to an associated link. (0, 1 or 2)
 - 2) If both are unequal, the router forwards the packet to a default link (otherwise 3).
- When there are multiple matches, the router uses the longest prefix matching rule.

3.2.3 Comparison of Virtual Circuit & Datagram

Issue	Datagram	Virtual Circuit
Connection Setup	None	Required
Addressing	Packet contains full source and destination-address	Packet contains short virtual-circuit number identifier
State Information	None other than router table containing destination-network	Each virtual-circuit number entered to table on setup, used for routing
Routing	Packets routed independently	Route established at setup, all packets follow same route
Effect of Router Failure	Only on packets lost during crash	All virtual circuits passing through failed router terminated



COMPUTER NETWORKS

3.3 What's Inside a Router?

- The router is used for transferring packets from an incoming-links to the appropriate outgoing-links.

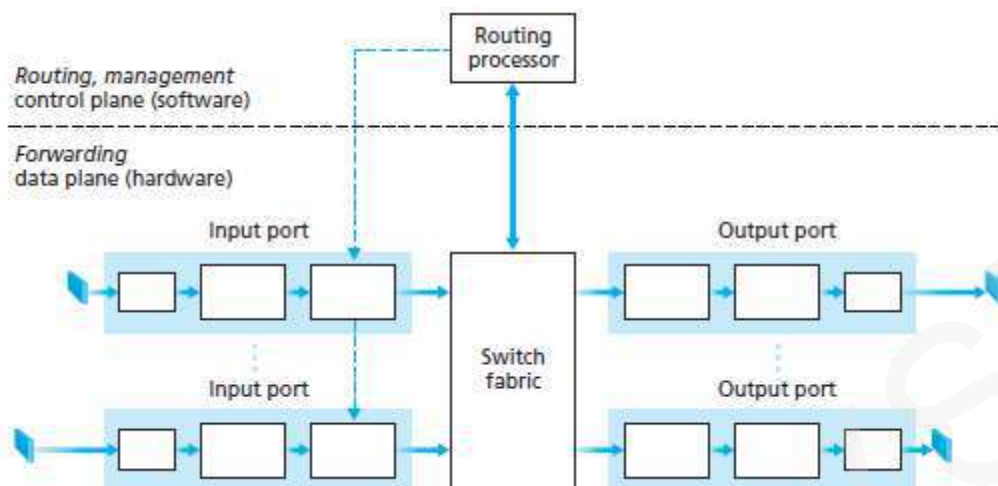


Figure 3.5: Router architecture

- Four components of router (Figure 3.5):

1) Input Ports

- An input-port is used for terminating an incoming physical link at a router (Figure 3.6).
- It is used for interoperating with the link layer at the other side of the incoming-link.
- It is used for lookup function i.e. searching through forwarding-table looking for longest prefix match.
- It contains forwarding-table.
- Forwarding-table is consulted to determine output-port to which arriving packet will be forwarded.
- Control packets are forwarded from an input-port to the routing-processor.
- Many other actions must be taken:
 - Packet's version number, checksum and time-to-live field must be checked.
 - Counters used for network management must be updated.

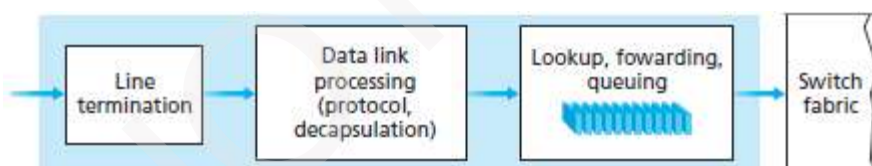


Figure 3.6: Input port processing

2) Switching Fabric

- The switching fabric connects the router's input-ports to its output-ports.
- In fabric, the packets are switched (or forwarded) from an input-port to an output-port.
- In fact, fabric is a network inside of a router.
- A packet may be temporarily blocked if packets from other input-ports are currently using the fabric.
- A blocked packet will be queued at the input-port & then scheduled to send at a later point in time.

3) Output Ports

- An output-port
 - stores packets received from the switching fabric and
 - transmits the packets on the outgoing-link.
- For a bidirectional link, an output-port will typically be paired with the input-port.

4) Routing Processor

- The routing-processor
 - executes the routing protocols
 - maintains routing-tables & attached link state information and
 - computes the forwarding-table.
- It also performs the network management functions.



COMPUTER NETWORKS

3.3.1 Switching

- Three types of switching fabrics (Figure 3.7):
 - 1) Switching via memory
 - 2) Switching via a bus and
 - 3) Switching via an interconnection network.

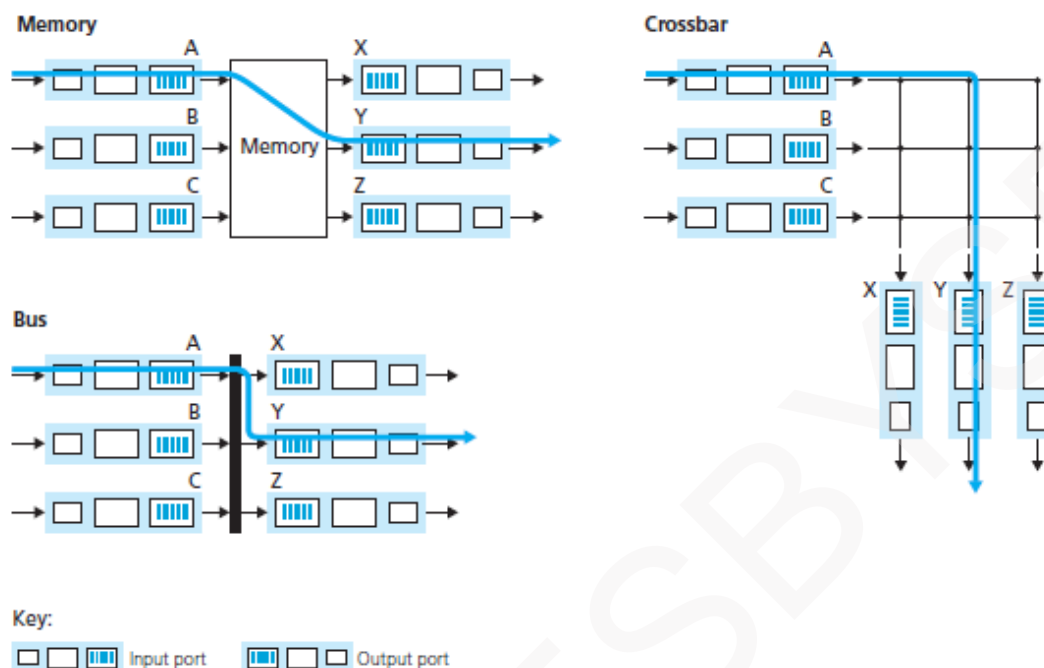


Figure 3.7: Three switching techniques

3.3.1.1 Switching via Memory

- Switching b/w input-ports & output-ports is done under direct control of CPU i.e. routing-processor.
- Input and output-ports work like a traditional I/O devices in a computer.
- Here is how it works (Figure 3.7a):
 - i) On arrival of a packet, the input-port notifies the routing-processor via an interrupt.
 - ii) Then, the packet is copied from the input-port to processor-memory.
 - iii) Finally, the routing-processor
 - extracts the destination-address from the header
 - looks up the appropriate output-port in the forwarding-table and
 - copies the packet into the output-port's buffers.
- Let memory-bandwidth = B packets per second.
Thus, the overall forwarding throughput must be less than B/2.
- Disadvantage:
 - Multiple packets cannot be forwarded at the same time. This is because
 - only one memory read/write over the shared system bus can be done at a time.

3.3.1.2 Switching via a Bus

- Switching b/w input-ports & output-ports is done without intervention by the routing-processor.
- Here is how it works (Figure 3.7b):
 - i) The input-port appends a switch-internal label (header) to the packet.
 - The label indicates the local output-port to which the packet must be transferred.
 - ii) Then, the packet is received by all output-ports.
 - But, only the port that matches the label will keep the packet.
 - iii) Finally, the label is removed at the output-port.
- Disadvantages:
 - i) Multiple packets cannot be forwarded at the same time. This is because
 - only one packet can cross the bus at a time.
 - ii) The switching speed of the router is limited to the bus-speed.



COMPUTER NETWORKS

3.3.1.3 Switching via an Interconnection Network

- A crossbar switch is an interconnection network.
- The network consists of $2N$ buses that connect N input-ports to N output-ports.
- Each vertical bus intersects each horizontal bus at a crosspoint.
- The crosspoint can be opened or closed at any time by the switch-controller.
- Here is how it works (Figure 3.7c):
 - 1) To move a packet from port A to port Y, the switch-controller closes the crosspoint at the intersection of buses A and Y.
 - 2) Then, port A sends the packet onto its bus, which is picked up by bus Y.
- Advantage:
 - Crossbar networks are capable of forwarding multiple packets in parallel.
 - For ex: A packet from port B can be forwarded to port X at the same time. This is because
 - A-to-Y and B-to-X packets use different input and output buses.
- Disadvantage:
 - If 2 packets have to use same output-port, then one packet has to wait. This is because
 - only one packet can be sent over any given bus at a time.

3.3.2 Output Processing

- Output-port processing
 - takes the packets stored in the output-port's memory and
 - transmits the packets over the output link (Figure 3.8).
- This includes
 - selecting and dequeuing packets for transmission and
 - performing the linklayer and physical-layer transmission functions.

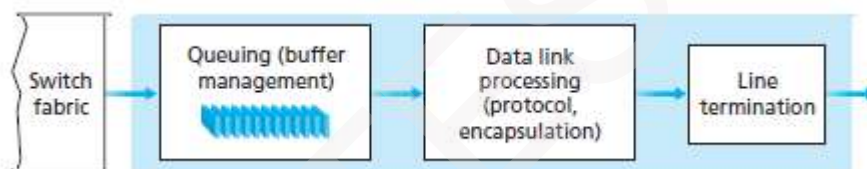


Figure 3.8: Output port processing

COMPUTER NETWORKS

3.3.3 Where Does Queueing Occur?

- Packet queues may form at both the input-ports & the output-ports (Figure 3.9).
- As the queues grow large, the router's memory can be exhausted and packet loss will occur.
- The location and extent of queueing will depend on
 - 1) The traffic load
 - 2) The relative speed of the switching fabric and
 - 3) The line speed
- Switching fabric transfer rate R_{switch} is defined as
"The rate at which packets can be moved from input-port to output-port".
- If R_{switch} is N times faster than R_{line} , then only negligible queueing will occur at the input-ports.

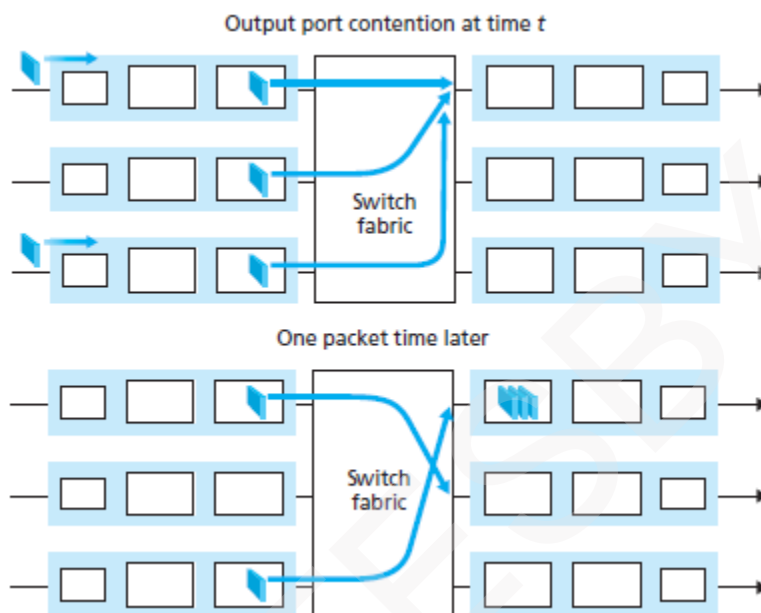


Figure 3.9: Output port queueing

- At output-port, packet-scheduler is used to choose one packet among those queued for transmission.
- The packet-scheduling can be done using
 - first-come-first-served (FCFS) or
 - weighted fair queueing (WFQ).
- Packet scheduling plays a crucial role in providing QoS guarantees.
- If there is less memory to hold an incoming-packet, a decision must be made to either
 - 1) Drop the arriving packet (a policy known as drop-tail) or
 - 2) Remove one or more already-queued packets to make room for the newly arrived packet.



COMPUTER NETWORKS

3.4 IP: Forwarding & Addressing in the Internet

- IP(Internet Protocol) is main protocol responsible for packetizing, forwarding & delivery of a packet at network-layer.
- It is a connection-less & unreliable protocol.
 - i) Connection-less means there is no connection setup b/w the sender and the receiver.
 - ii) Unreliable protocol means
 - IP does not make any guarantee about delivery of the data.
 - Packets may get dropped during transmission.
- It provides a best-effort delivery service.
- Best effort means IP does its best to get the packet to its destination, but with no guarantees.
- If reliability is important, IP must be paired with a TCP which is reliable transport-layer protocol.
- IP does not provide following services
 - flow control
 - error control
 - congestion control services.

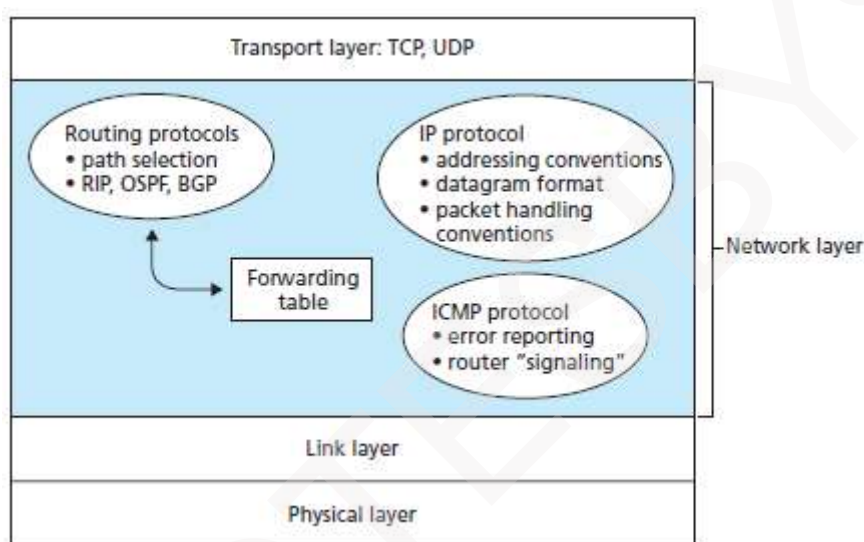


Figure 3.10: A look inside the Internet's network-layer

- Two important components of IP:
 - 1) Internet addressing and
 - 2) Forwarding
- There are two versions of IP in use today.
 - 1) IP version 4 (IPv4) and
 - 2) IP version 6 (IPv6)
- As shown in Figure 3.10, the network-layer has three major components:
 - 1) IP protocol
 - 2) Routing component determines the path a data follows from source to destination
 - 3) Network-layer is a facility to report errors in datagrams



COMPUTER NETWORKS

3.4.1 IPv4 Datagram Format

- IP uses the packets called datagrams.
- A datagram consist of 2 parts: 1) Payload (or Data) 2) Header.

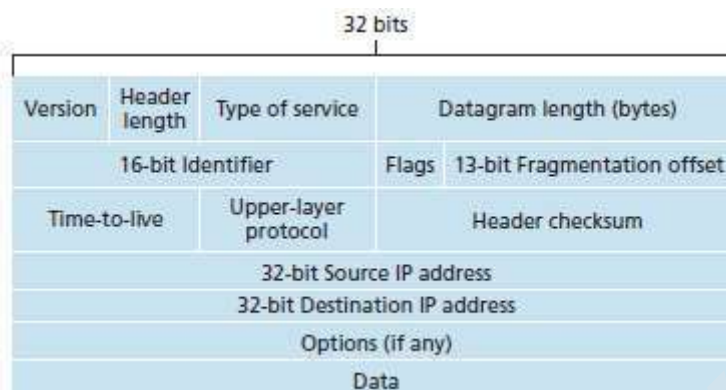


Figure 3.11: IPv4 datagram format

1) Payload (or Data)

- This field contains the data to be delivered to the destination.

2) Header

- Header contains information essential to routing and delivery.
- IP header contains following fields (Figure 3.11):

1) Version

- This field specifies version of the IPv4 datagram, i.e. 4.

2) Header Length

- This field specifies length of header.
- Without options field, header-length = 5 bytes.

3) Type of Service (TOS)

- This field specifies priority of packet based on parameters such as delay, throughput, reliability & cost.

4) Datagram Length

- This field specifies the total length of the datagram (header + data).
- Maximum length = 65535 bytes.

5) Identifier, Flags, Fragmentation Offset

- These fields are used for fragmentation and reassembly.
- Fragmentation occurs when the size of the datagram is larger than the MTU of the network.

i) **Identifier:** This field uniquely identifies a datagram packet.

ii) **Flags:** It is a 3-bit field. The first bit is not used.
The second bit D is called the do not fragment bit.
The third bit M is called the more fragment bit.

iii) **Fragmentation Offset:** This field identifies location of a fragment in a datagram.

6) Time-To-Live (TTL)

- This defines lifetime of the datagram (default value 64) in hops.
- Each router decrements TTL by 1 before forwarding. If TTL is zero, the datagram is discarded.

7) Protocol

- This field specifies upper-layer protocol used to receive the datagram at the destination-host.
- For example, TCP=6 and UDP=17.

8) Header Checksum

- This field is used to verify integrity of header only.
- If the verification process fails, the packet is discarded.

9) Source IP Address & Destination IP Address

- These fields contain the addresses of source and destination respectively.

10) Options

- This field allows the packet to request special features such as
 - security level
 - route to be taken by packet at each router.



COMPUTER NETWORKS

3.4.2 Fragmentation

3.4.2.1 Maximum Transfer Unit

- Each network imposes a restriction on maximum size of packet that can be carried. This is called the MTU (maximum transmission unit).
- For example:
 - MTU Ethernet = 1500 bytes
 - MTU FDDI = 4464 bytes
- Fragmentation means
 - “The datagram is divided into smaller fragments when size of a datagram is larger than MTU”
- Each fragment is routed independently (Figure 3.12).
- A fragmented datagram may be further fragmented, if it encounters a network with a smaller MTU.
- Source/router is responsible for fragmentation of original datagram into the fragments.
 - Only destination is responsible for reassembling the fragments into the original datagram.

3.4.2.2 Fields Related to Fragmentation & Reassembly

- Three fields in the IP header are used to manage fragmentation and reassembly:
 - 1) Identification
 - 2) Flags
 - 3) Fragmentation offset.

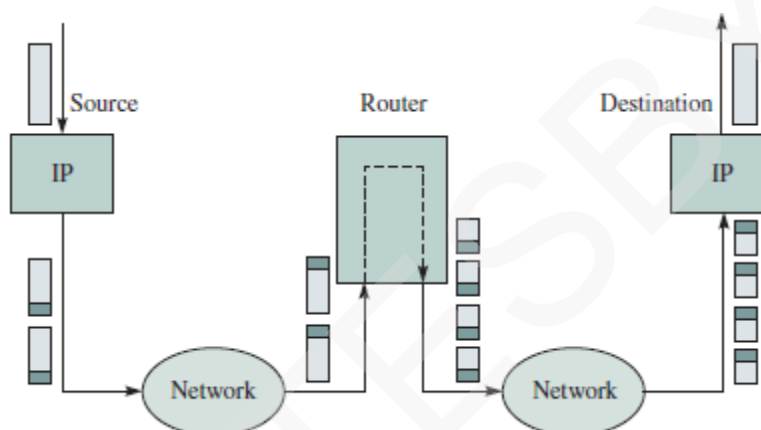


Figure 3.12: IP fragmentation and reassembly

1) Identification

- This field is used to identify to which datagram a particular fragment belongs to (so that fragments for different packets do not get mixed up).
- When a datagram is created, the source attaches the datagram with an identification-number.
- When a datagram is fragmented, the value in the identification-field is copied into all fragments.
- The identification-number helps the destination in reassembling the datagram.

2) Flags

- This field has 3 bits.
 - i) The first bit is not used.
 - ii) DF bit (Don't Fragment):
 - a) If DF=1, the router should not fragment the datagram. Then, the router discards the datagram.
 - b) If DF=0, the router can fragment the datagram.
 - iii) MF bit (More Fragment):
 - a) If MF=1, there are some more fragments to come.
 - b) If MF=0, this is last fragment.

3) Fragmentation Offset

- This field identifies location of a fragment in a datagram.
- This field is the offset of the data in the original datagram.



COMPUTER NETWORKS

3.4.3 IPv4 Addressing

- IP address is a numeric identifier assigned to each machine on the internet.
- IP address consists of two parts: network ID(NID) and host ID(HID).
 - 1) NID identifies the network to which the host is connected. All the hosts connected to the same network have the same NID.
 - 2) HID is used to uniquely identify a host on that network.
- HID is assigned by the network-administrator at the local site. NID for an organization may be assigned by the ISP (Internet Service Provider).
- IPv4 uses 32-bit addresses, i.e., approximately 4 billion addresses (2^{32}).
- IP addresses are usually written in dotted-decimal notation. The address is broken into four bytes. For example, an IP address of 10000000 10000111 01000100 00000101 is written as 128.135.68.5
- IP address can be classified as
 - 1) Classful IP addressing &
 - 2) Classless IP addressing (CIDR → Classless Inter Domain Routing)

3.4.3.1 IPv4 Classful Addressing

- In classful addressing, the address space is divided into five classes: A, B, C, D and E.
- IP address class is identified by MSBs in binary.
- Classes A, B and C are used for unicast addressing. (Figure 3.13).
- Class D was designed for multicasting and class E is reserved.

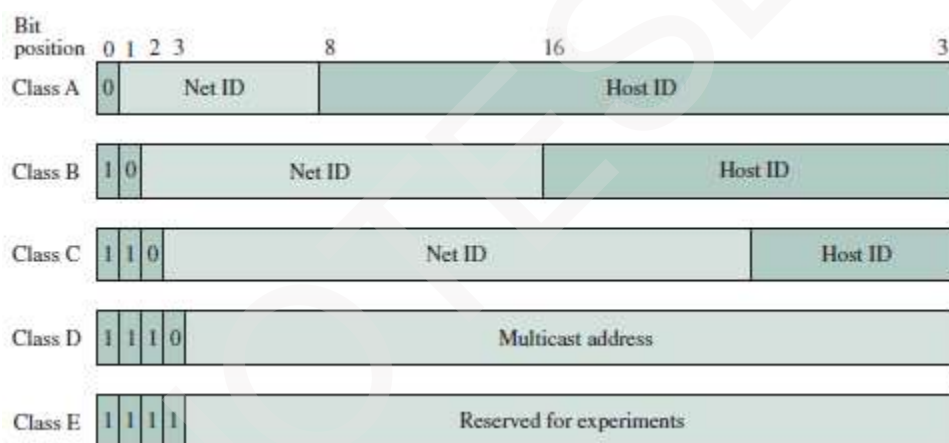


Figure 3.13: The five classes of IP addresses

Class	No. of networks	Max. No. of hosts per network	Designed for
A	126	$2^{24} - 2$	WAN
B	16,382	65,534	Campus networks
C	2^{21}	254	LAN

Table 3.3: Classful Addressing

- Analysis:
 - In classful addressing, a large part of the available addresses were wasted, since Class A and B were too large for most organizations (Table 3.3).
 - Class C is suited only for small organization and reserved addresses were sparingly used.



COMPUTER NETWORKS

3.4.3.1.1 Subnet Addressing

- Problem with classful addressing:
 - Consider an organization has a Class B address which can support about 64,000 hosts.
 - It will be a huge task for the network-administrator to manage all 64,000 hosts.
- Solution: Use subnet addressing.
- Subnetting reduces the total number of network-numbers by assigning a single network-number to many adjacent physical networks.
- Each adjacent physical network is referred to as subnet. (Figure 3.14).
- All nodes on a subnet are configured with a subnet mask. For example: 255.255.255.0.
- The 1's in the subnet-mask represent the positions that refer to the network or subnet-numbers.
- The 0's represent the positions that refer to the host part of the address.
- The bitwise AND of IP address and its subnet mask gives the subnet number.
- Advantage:
 - The subnet-addressing scheme is oblivious to the network outside the organization.
 - Inside the organization the network-administrator is free to choose any combination of lengths for the subnet & host ID fields.

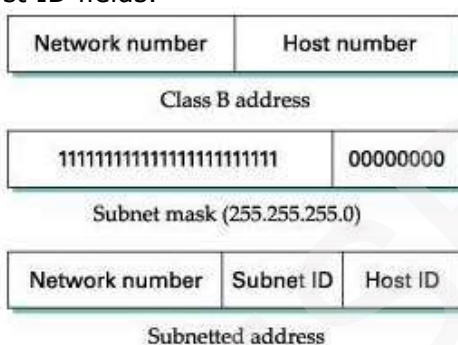


Figure 3.14: Subnet addressing

Question: If a packet with a destination IP address of 150.100.12.176 arrives at site from the outside network, which subnet should a router forward this packet to? Assume subnet mask is 255.255.255.128 (Figure 3.15).

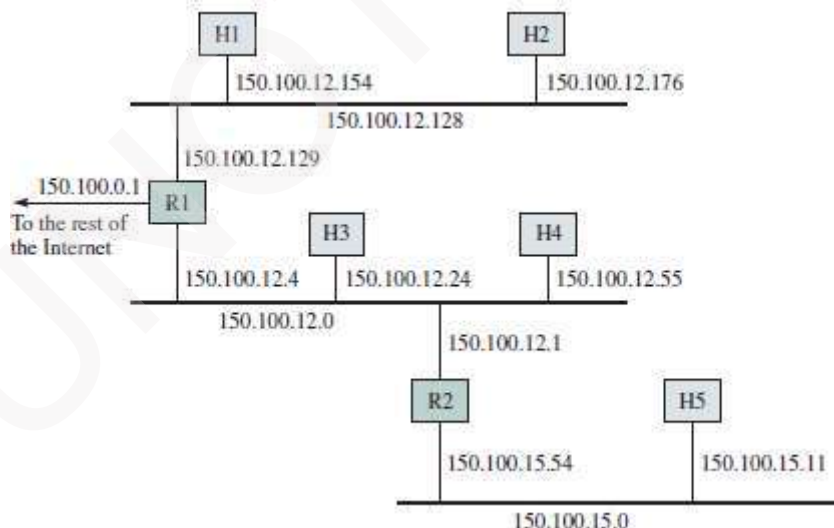


Figure 3.15: Example of address assignment with subnetting

Solution: The router can determine the subnet number by performing a binary AND between the subnet mask and the IP address.

```

IP address:      10010110 01100100 00001100 10110000(150.100.12.176)
Subnet mask:    11111111 11111111 11111111 10000000(255.255.255.128)
Subnet number:  10010110 01100100 00001100 10000000(150.100.12.128)

```

This number (150.100.12.128) is used to forward the packet to the correct subnet work inside the organization.



COMPUTER NETWORKS

3.4.3.2 CIDR

- Problem with classful IP addressing:
 - Consider an organization needs about 500 hosts.
 - Obviously, the organization will get a Class B license, even though it has far fewer than 64,000 hosts.
 - At most, over 64,000 addresses can go unused.
 - This results in inefficient usage of the available address-space.
- Solution: Use CIDR (Classless Inter Domain Routing).
 - A single IP address can be used to designate many unique IP addresses. This is called supernetting.
 - A CIDR IP address looks like a normal IP address except that
 - the address ends with a slash followed by a number, called the IP network prefix.
 - For ex: 205.100.0.0/22
 - CIDR addresses
 - reduce the size of routing-tables and
 - make more IP addresses available within organizations.

3.4.3.3 Obtaining a Block of Addresses

- To obtain a block of IP addresses for use within an organization's subnet, a network-administrator contacts the ISP.
- IP addresses are managed under the authority of the ICANN.
- The responsibility of the ICANN (Internet Corporation for Assigned Names and Numbers):
 - to allocate IP addresses,
 - to manage the DNS root servers.
 - to assign domain names and resolve domain name disputes.
 - to allocate addresses to regional Internet registries.

3.4.3.4 Obtaining a Host Address: DHCP

- Two ways to assign an IP address to a host:
 - 1) Manual Configuration**
 - Operating systems allow system-administrator to manually configure IP address.
 - 2) Dynamic Host Configuration Protocol (DHCP)**
 - DHCP enables auto-configuration of IP address to host.



COMPUTER NETWORKS

3.4.3.4.1 DHCP Protocol

- DHCP enables auto-configuration of IP address to host.
- DHCP assigns dynamic IP addresses to devices on a network.
- Dynamic address allocation is required
 - when a host moves from one network to another or
 - when a host is connected to a network for the first time.

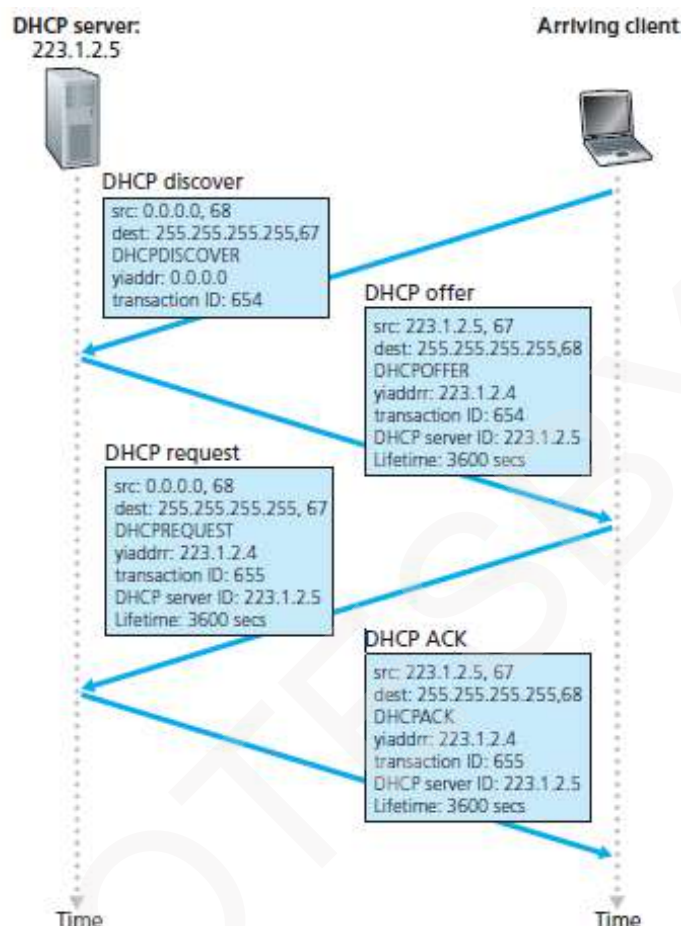


Figure 3.16: DHCP client-server interaction

- Four steps in DHCP protocol (Figure 3.16):
 - 1) DHCP Server Discovery**
 - DHCP server contains a range of unassigned addresses to be assigned to hosts on-demand.
 - To contact DHCP server, a client broadcasts a DHCPDISCOVER message with destination IP address 255.255.255.255.
 - 2) DHCP Server Offer**
 - DHCP server broadcasts DHCPOFFER message containing
 - client's IP address
 - network mask and
 - IP address lease time (i.e. the amount of time for which the IP address will be valid).
 - 3) DHCP Request**
 - The client sends a DHCPREQUEST message, requesting the offered address.
 - 4) DHCP ACK**
 - The DHCP server acknowledges with a DHCPACK message containing the requested configuration.



COMPUTER NETWORKS

3.4.3.5 NAT

- Network Address Translation (NAT) enables hosts to use Internet without the need to have globally unique addresses.
- NAT enables organization to have a large set of addresses internally and one address externally.
- The organization must have single connection to the Internet through a NAT-enabled router.
- NAT allows a single device (such as a router) to act as an agent b/w
 - 1) Internet (or "public network") and
 - 2) Local (or "private") network.
- This means only a single, unique IP address is required to represent an entire group of computers.
- Figure 3.17 shows the operation of a NAT-enabled router.

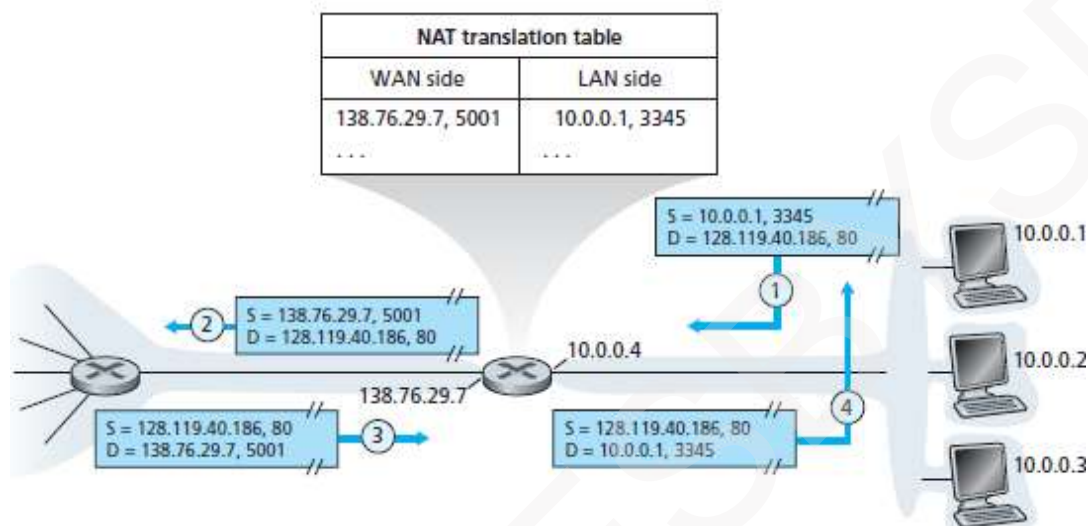


Figure 3.17: Network address translation

- The private addresses only have meaning to devices within a given network.
- The NAT-enabled router does not look like a router to the outside world.
- Instead, the NAT-enabled router behaves to the outside world as a single device with a single IP address.
- In Figure 3.17,
 - 1) All traffic leaving the home-router for the Internet has a source-address of 138.76.29.7.
 - 2) All traffic entering the home-router must have a destination-address of 138.76.29.7.
- The NAT-enabled router is hiding the details of the home-network from the outside world.
- At the NAT router, NAT translation-table includes
 - 1) Port numbers and
 - 2) IP addresses.
- IETF community is against the use of NAT. This is because of following reasons:
 - 1) They argue, port numbers are to be used for addressing processes, not for addressing hosts.
 - 2) They argue routers are supposed to process packets only up to layer 3.
 - 3) They argue the NAT protocol violates the end-to-end argument.
 - 4) They argue, we should use IPv6 to solve the shortage of IP addresses.
 - 5) NAT interferes with P2P applications. If Peer B is behind NAT, Peer B cannot act as a server.



COMPUTER NETWORKS

3.4.4 ICMP

- ICMP is a network-layer protocol. (ICMP → Internet Control Message Protocol).
- This is used to handle error and other control messages.
- Main responsibility of ICMP: To report errors that occurs during the processing of the datagram.
- ICMP does not correct errors; ICMP simply reports the errors to the source.
- 12 types of ICMP messages are defined as shown in Table 3.4.
- Each ICMP message type is encapsulated in an IP packet.

ICMP Type	Code	Description
0	0	echo reply (to ping)
3	0	destination network unreachable
3	1	destination host unreachable
3	2	destination protocol unreachable
3	3	destination port unreachable
3	6	destination network unknown
3	7	destination host unknown
4	0	source quench (congestion control)
8	0	echo request
9	0	router advertisement
10	0	router discovery
11	0	TTL expired
12	0	IP header bad

Table 3.4: ICMP message types

1) Destination Unreachable (Type=3)

- This message is related to problem reaching the destinations.
- This message uses different codes (0 to 15) to define type of error-message.
- Possible values for code field:

Code 0 = network unreachable	Code 1 = host unreachable
Code 2 = protocol unreachable	Code 3 = port unreachable

2) Source Quench (Type=4)

- The main purpose is to perform congestion control.
- This message
 - informs the sender that network has encountered congestion & datagram has been dropped.
 - informs the sender to reduce its transmission-rate.

3) Echo Request & Echo Reply (Type=8 & Type=0)

- These messages are used to determine whether a remote-host is alive.
- A source sends an echo request-message to destination;
 - If the destination is alive, the destination responds with an echo reply message.
- Type=8 is used for echo request;
 - Type=0 is used for echo reply.
- These messages can be used in two debugging tools: ping and traceroute.
 - i) **Ping**
 - The ping program can be used to find if a host is alive and responding.
 - The source-host sends ICMP echo-request-messages.
 - The destination, if alive, responds with ICMP echo-reply messages.
 - ii) **Traceroute**
 - The traceroute program can be used to trace the path of a packet from source to destination.
 - It can find the IP addresses of all the routers that are visited along the path.
 - The program is usually set to check for the maximum of 30 hops (routers) to be visited.



COMPUTER NETWORKS

3.4.5 IPv6

- CIDR, subnetting and NAT could not solve address-space exhaustion faced by IPv4.
- IPv6 was evolved to solve this problem.

3.4.5.1 Changes from IPv4 to IPv6 (Advantages of IPv6)

1) Expanded Addressing Capabilities

- IPv6 increases the size of the IP address from 32 to 128 bits (Supports upto 3.4×10^{38} nodes).
- In addition to unicast & multicast addresses, IPv6 has an anycast address.
- Anycast address allows a datagram to be delivered to only one member of the group.

2) A Streamlined 40-byte Header

- A number of IPv4 fields have been dropped or made optional.
- The resulting 40-byte fixed-length header allows for faster processing of the IP datagram.
- A new encoding of options field allows for more flexible options processing.

3) Flow Labeling & Priority

- A flow can be defined as
"Labeling of packets belonging to particular flows for which the sender requests special handling".
- For example:
Audio and video transmission may be treated as a flow.



COMPUTER NETWORKS

3.4.5.2 IPv6 Datagram Format

The format of the IPv6 datagram is shown in Figure 3.18.

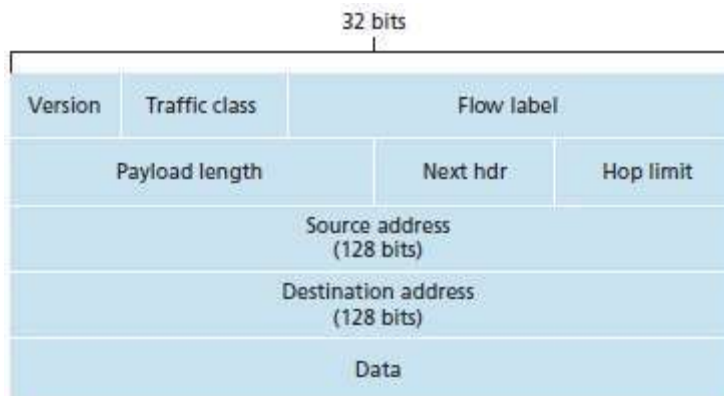


Figure 3.18: IPv6 datagram format

The following fields are defined in IPv6:

1) Version

This field specifies the IP version, i.e., 6.

2) Traffic Class

This field is similar to the TOS field in IPv4.

This field indicates the priority of the packet.

3) Flow Label

This field is used to provide special handling for a particular flow of data.

4) Payload Length

This field shows the length of the IPv6 payload.

5) Next Header

This field is similar to the options field in IPv4 (Figure 3.19).

This field identifies type of extension header that follows the basic header.

6) Hop Limit

This field is similar to TTL field in IPv4.

This field shows the maximum number of routers the packet can travel.

The contents of this field are decremented by 1 by each router that forwards the datagram.

If the hop limit count reaches 0, the datagram is discarded.

7) Source & Destination Addresses

These fields show the addresses of the source & destination of the packet.

8) Data

This field is the payload portion of the datagram.

When the datagram reaches the destination, the payload will be

→ removed from the IP datagram and

→ passed on to the upper layer protocol (TCP or UDP).

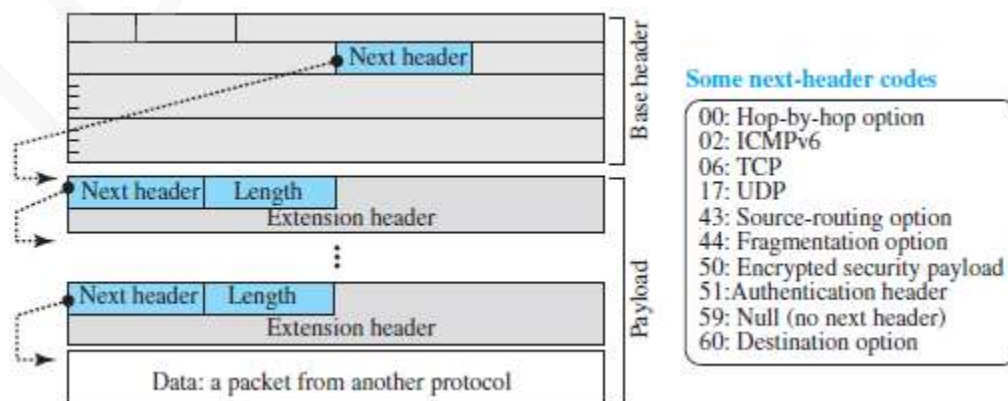


Figure 3.19: Payload in IPv6 datagram



COMPUTER NETWORKS

3.4.5.3 IPv4 Fields not present in IPv6

1) Fragmentation/Reassembly

- Fragmentation of the packet is done only by the source, but not by the routers.
The reassembling is done by the destination.
- Fragmentation & reassembly is a time-consuming operation.
- At routers, the fragmentation is not allowed to speed up the processing in the router.
- If packet-size is greater than the MTU of the network, the router
 - drops the packet.
 - sends an error message to inform the source.

2) Header Checksum

- In the Internet layers, the transport-layer and link-layer protocols perform check summing.
- This functionality was redundant in the network-layer.
- So, this functionality was removed to speed up the processing in the router.

3) Options

- In, IPv6, next-header field is similar to the options field in IPv4.
- This field identifies type of extension header that follows the basic header.
- To support extra functionalities, extension headers can be placed b/w base header and payload.

3.4.5.4 Difference between IPv4 & IPv6

	IPv4	IPv6
1	IPv4 addresses are 32 bit length	IPv6 addresses are 128 bit length
2	Fragmentation is done by sender and forwarding routers	Fragmentation is done only by sender
3	Does not identify packet flow for QoS handling	Contains Flow Label field that specifies packet flow for QoS handling
4	Includes Options up to 40 bytes	Extension headers used for optional data
5	Includes a checksum	Does not includes a checksum
6	Address Resolution Protocol (ARP) is available to map IPv4 addresses to MAC addresses	Address Resolution Protocol (ARP) is replaced with Neighbor Discovery Protocol (NDP)
7	Broadcast messages are available	Broadcast messages are not available
8	Manual configuration (Static) of IP addresses or DHCP (Dynamic configuration) is required to configure IP addresses	Auto-configuration of addresses is available
9	IPSec is optional, external	IPSec is required



COMPUTER NETWORKS

3.4.5.5 Transitioning from IPv4 to IPv6

- IPv4-capable systems are not capable of handling IPv6 datagrams.
- Two strategies have been devised for transition from IPv4 to IPv6:
 - 1) Dual stack and
 - 2) Tunneling.

3.4.5.5.1 Dual Stack Approach

- IPv6-capable nodes also have a complete IPv4 implementation. Such nodes are referred to as IPv6/IPv4 nodes.
- IPv6/IPv4 node has the ability to send and receive both IPv4 and IPv6 datagrams.
- When interoperating with an IPv4 node, an IPv6/IPv4 node can use IPv4 datagrams.
When interoperating with an IPv6 node, an IPv6/IPv4 node can use IPv6 datagrams.
- IPv6/IPv4 nodes must have both IPv6 and IPv4 addresses.
- IPv6/IPv4 nodes must be able to determine whether another node is IPv6-capable or IPv4-only.
- This problem can be solved using the DNS.
If the node name is resolved to IPv6-capable, then the DNS returns an IPv6 address
Otherwise, the DNS return an IPv4 address.
- If either the sender or the receiver is only IPv4-capable, an IPv4 datagram must be used.
- Two IPv6-capable nodes can send IPv4 datagrams to each other.

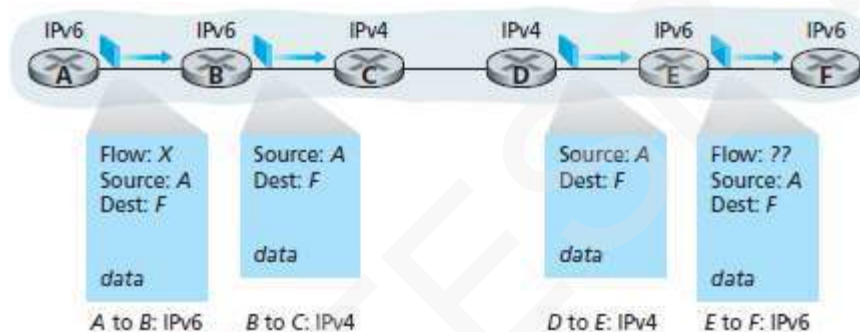


Figure 3.20: A dual-stack approach

- Dual stack is illustrated in Figure 3.20.
- Here is how it works:
 - 1) Suppose IPv6-capable Node-A wants to send a datagram to IPv6-capable Node-F.
 - 2) IPv6-capable Node-B creates an IPv4 datagram to send to IPv4-capable Node-C.
 - 3) At IPv6-capable Node-B, the IPv6 datagram is copied into the data field of the IPv4 datagram and appropriate address mapping can be done.
 - 4) At IPv6-capable Node-E, the IPv6 datagram is extracted from the data field of the IPv4 datagram.
 - 5) Finally, IPv6-capable Node-E forwards an IPv6 datagram to IPv6-capable Node-F.
- Disadvantage: During transition from IPv6 to IPv4, few IPv6-specific fields will be lost.



COMPUTER NETWORKS

3.4.5.5.2 Tunneling

- Tunneling is illustrated in Figure 3.21.
- Suppose two IPv6-nodes B and E
 - want to interoperate using IPv6 datagrams and
 - are connected by intervening IPv4 routers.
- The intervening-set of IPv4 routers between two IPv6 routers are referred as a tunnel.
- Here is how it works:
 - On the sending side of the tunnel:
 - IPv6-node B takes & puts the IPv6 datagram in the data field of an IPv4 datagram.
 - The IPv4 datagram is addressed to the IPv6-node E.
 - On the receiving side of the tunnel: The IPv6-node E
 - receives the IPv4 datagram
 - extracts the IPv6 datagram from the data field of the IPv4 datagram and
 - routes the IPv6 datagram to IPv6-node F

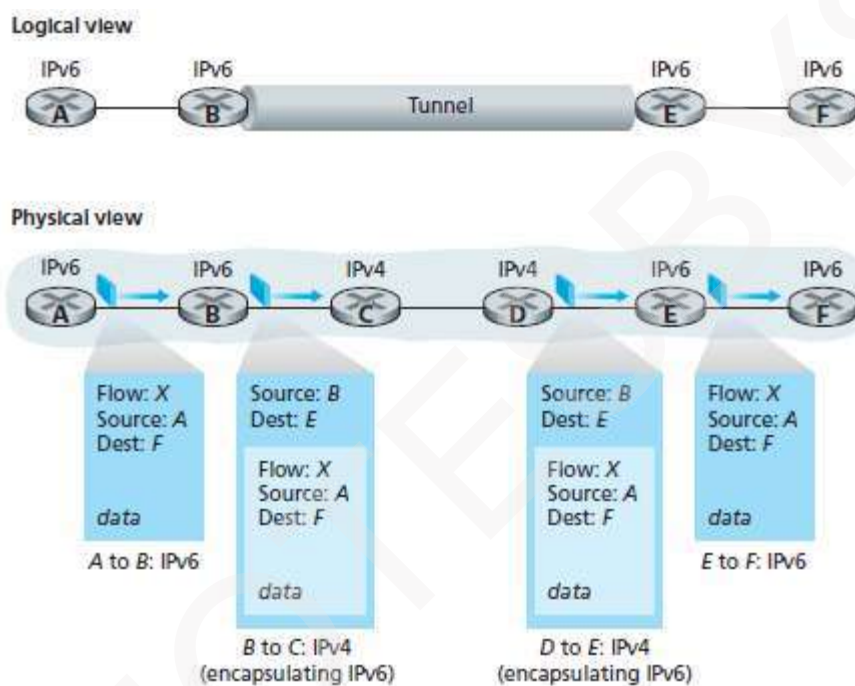


Figure 3.21: Tunneling



COMPUTER NETWORKS

3.4.6 A Brief Foray into IP Security

- IPsec is a popular secure network-layer protocol.
- It is widely deployed in Virtual Private Networks (VPNs).
- It has been designed to be backward compatible with IPv4 and IPv6.
- It can be used to create a connection-oriented service between 2 entities.
- In transport mode, 2 hosts first establish an IPsec session between themselves.
- All TCP and UDP segments sent between the two hosts enjoy the security services provided by IPsec.
- On the source-side,
 - 1) The transport-layer passes a segment to IPsec.
 - 2) Then, IPsec
 - encrypts the segment
 - appends additional security fields to the segment and
 - encapsulates the resulting payload in a IP datagram.
 - 3) Finally, the sending-host sends the datagram into the Internet.
 - The Internet then transports the datagram to the destination-host.
- On the destination-side,
 - 1) The destination receives the datagram from the Internet.
 - 2) Then, IPsec
 - decrypts the segment and
 - passes the unencrypted segment to the transport-layer.
- Three services provided by an IPsec:
 - 1) Cryptographic Agreement**
 - This mechanism allows 2 communicating hosts to agree on cryptographic algorithms & keys.
 - 2) Encryption of IP Datagram Payloads**
 - When the sender receives a segment from the transport-layer, IPsec encrypts the payload.
 - The payload can only be decrypted by IPsec in the receiver.
 - 3) Data Integrity**
 - IPsec allows the receiver to verify that the datagram's header fields.
 - The encrypted payload is not modified after transmission of the datagram into the n/w.
 - 4) Origin Authentication**
 - The receiver is assured that the source-address in datagram is the actual source of datagram.



COMPUTER NETWORKS

3.5 Routing Algorithms

- A routing-algorithm is used to find a “good” path from source to destination.
- Typically, a good path is one that has the least cost.
- The least-cost problem: Find a path between the source and destination that has least cost.

3.5.1 Routing Algorithm Classification

- A routing-algorithm can be classified as follows:
 - 1) Global or decentralized
 - 2) Static or dynamic
 - 3) Load-sensitive or Load-insensitive

3.5.1.1 Global or Decentralized

Global Routing Algorithm

- The calculation of the least-cost path is carried out at one centralized site.
- This algorithm has complete, global knowledge about the network.
- Algorithms with global state information are referred to as link-state (LS) algorithms.

Decentralized Routing Algorithm

- The calculation of the least-cost path is carried out in an iterative, distributed manner.
- No node has complete information about the costs of all network links.
- Each node has only the knowledge of the costs of its own directly attached links.
- Each node performs calculation by exchanging information with its neighboring nodes.

3.5.1.2 Static or Dynamic

Static Routing Algorithms

- Routes change very slowly over time, as a result of human intervention.
- For example: a human manually editing a router’s forwarding-table.

Dynamic Routing Algorithms

- The routing paths change, as the network-topology or traffic-loads change.
- The algorithm can be run either
 - periodically or
 - in response to topology or link cost changes.
- Advantage: More responsive to network changes.
- Disadvantage: More susceptible to routing loop problem.

3.5.1.3 Load Sensitive or Load Insensitive

Load Sensitive Algorithm

- Link costs vary dynamically to reflect the current level of congestion in the underlying link.
- If high cost is associated with congested-link, the algorithm chooses routes around congested-link.

Load Insensitive Algorithm

- Link costs do not explicitly reflect the current level of congestion in the underlying link.
- Today’s Internet routing-algorithms are load-insensitive. For example: RIP, OSPF, and BGP



COMPUTER NETWORKS

3.5.2 LS Routing Algorithm

3.5.2.1 Dijkstra's Algorithm

- Dijkstra's algorithm computes the least-cost path from one node to all other nodes in the network.
- Let us define the following notation:

- 1) u : source-node
- 2) $D(v)$: cost of the least-cost path from the source u to destination v .
- 3) $p(v)$: previous node (neighbor of v) along the current least-cost path from the source to v .
- 4) N' : subset of nodes; v is in N' if the least-cost path from the source to v is known.

Link-State (LS) Algorithm for Source Node u

```

1 Initialization:
2    $N' = \{u\}$ 
3   for all nodes  $v$ 
4     if  $v$  is a neighbor of  $u$ 
5       then  $D(v) = c(u,v)$ 
6     else  $D(v) = \infty$ 
7
8 Loop
9   find  $w$  not in  $N'$  such that  $D(w)$  is a minimum
10  add  $w$  to  $N'$ 
11  update  $D(v)$  for each neighbor  $v$  of  $w$  and not in  $N'$ :
12     $D(v) = \min( D(v), D(w) + c(w,v) )$ 
13  /* new cost to  $v$  is either old cost to  $v$  or known
14   least path cost to  $w$  plus cost from  $w$  to  $v$  */
15 until  $N' = N$ 

```

- Example: Consider the network in Figure 3.22 and compute the least-cost paths from u to all possible destinations.

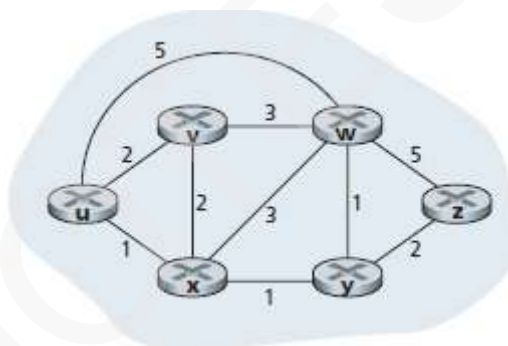


Figure 3.22: Abstract graph model of a computer network

Solution:

- Let's consider the few first steps in detail.
 - 1) In the initialization step, the currently known least-cost paths from u to its directly attached neighbors, v , x , and w , are initialized to 2, 1, and 5, respectively.
 - 2) In the first iteration, we
 - look among those nodes not yet added to the set N' and
 - find that node with the least cost as of the end of the previous iteration.
 - 3) In the second iteration,
 - nodes v and y are found to have the least-cost paths (2) and
 - we break the tie arbitrarily and
 - add y to the set N' so that N' now contains u , x , and y .
 - 4) And so on. . . .
 - 5) When the LS algorithm terminates,
 - We have, for each node, its predecessor along the least-cost path from the source.
- A tabular summary of the algorithm's computation is shown in Table 3.5.



step	N'	$D(v),p(v)$	$D(w),p(w)$	$D(x),p(x)$	$D(y),p(y)$	$D(z),p(z)$
0	u	2,u	5,u	1,u	∞	∞
1	ux	2,u	4,x		2,x	∞
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					

Table 3.5: Running the link-state algorithm on the network in Figure 3.20

- Figure 3.23 shows the resulting least-cost paths for u for the network in Figure 3.22.

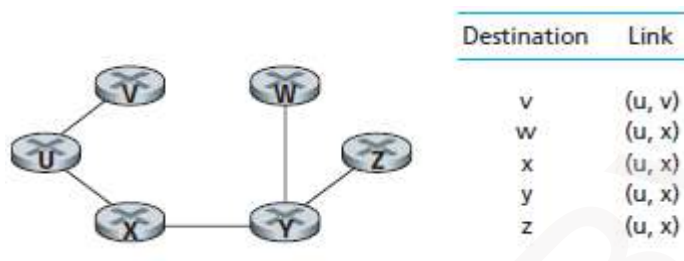


Figure 3.23: Least cost path and forwarding-table for node u



COMPUTER NETWORKS

3.5.3 DV Routing Algorithm

3.5.3.1 Bellman Ford Algorithm

- Distance vector (DV) algorithm is 1) iterative, 2) asynchronous, and 3) distributed.
 - It is distributed. This is because each node
 - receives some information from one or more of its directly attached neighbors
 - performs the calculation and
 - distributes then the results of the calculation back to the neighbors.
 - It is iterative. This is because
 - the process continues on until no more info is exchanged b/w neighbors.
 - It is asynchronous. This is because
 - the process does not require all of the nodes to operate in lockstep with each other.
- The basic idea is as follows:
 - Let us define the following notation:
 - $D_x(y)$ = cost of the least-cost path from node x to node y , for all nodes in N .
 - $D_x = [D_x(y): y \text{ in } N]$ be node x 's distance vector of cost estimates from x to all other nodes y in N .
 - Each node x maintains the following routing information:
 - For each neighbor v , the cost $c(x,v)$ from node x to directly attached neighbor v
 - Node x 's distance vector, that is, $D_x = [D_x(y): y \text{ in } N]$, containing x 's estimate of its cost to all destinations y in N .
 - The distance vectors of each of its neighbors, that is, $D_v = [D_v(y): y \text{ in } N]$ for each neighbor v of x .
 - From time to time, each node sends a copy of its distance vector to each of its neighbors.
 - The least costs are computed by the Bellman-Ford equation:

$$D_x(y) = \min_v \{c(x,v) + D_v(y)\} \quad \text{for each node } y \text{ in } N$$
 - If node x 's distance vector has changed as a result of this update step, node x will then send its updated distance vector to each of its neighbors.

Distance-Vector (DV) Algorithm

At each node, x :

```

1  Initialization:
2  for all destinations  $y$  in  $N$ :
3   $D_x(y) = c(x,y)$  /* if  $y$  is not a neighbor then  $c(x,y) = \infty$  */
4  for each neighbor  $w$ 
5   $D_w(y) = ?$  for all destinations  $y$  in  $N$ 
6  for each neighbor  $w$ 
7  send distance vector  $D_x = [D_x(y): y \text{ in } N]$  to  $w$ 
8
9  loop
10 wait (until I see a link cost change to some neighbor  $w$  or
11 until I receive a distance vector from some neighbor  $w$ )
12
13 for each  $y$  in  $N$ :
14  $D_x(y) = \min_v \{c(x,v) + D_v(y)\}$ 
15
16 if  $D_x(y)$  changed for any destination  $y$ 
17 send distance vector  $D_x = [D_x(y): y \text{ in } N]$  to all neighbors
18
19 forever

```

**COMPUTER NETWORKS**

- Figure 3.24 illustrates the operation of the DV algorithm for the simple three node network.

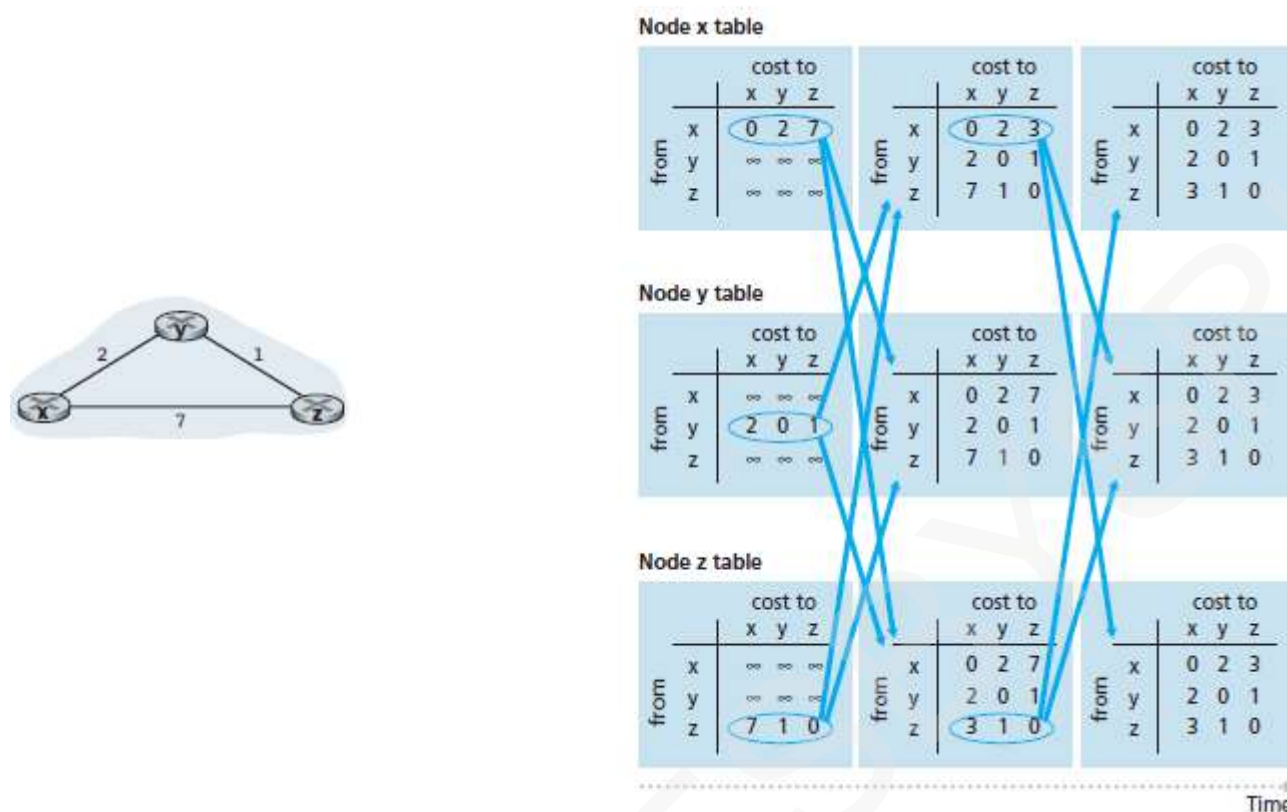


Figure 3.24: Distance-vector (DV) algorithm

- The operation of the algorithm is illustrated in a synchronous manner. Here, all nodes simultaneously
 - receive distance vectors from their neighbours
 - compute their new distance vectors, and
 - inform their neighbours if their distance vectors have changed.
- The table in the upper-left corner is node x's initial routing-table.
- In this routing-table, each row is a distance vector.
- The first row in node x's routing-table is $D_x = [D_x(x), D_x(y), D_x(z)] = [0, 2, 7]$.
- After initialization, each node sends its distance vector to each of its two neighbours.
- This is illustrated in Figure 3.24 by the arrows from the first column of tables to the second column of tables.
- For example, node x sends its distance vector $D_x = [0, 2, 7]$ to both nodes y and z. After receiving the updates, each node recomputes its own distance vector.
- For example, node x computes

$$D_x(x) = 0$$

$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\} = \min\{2 + 0, 7 + 1\} = 2$$

$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\} = \min\{2 + 1, 7 + 0\} = 3$$
- The second column therefore displays, for each node, the node's new distance vector along with distance vectors just received from its neighbours.
- Note, that node x's estimate for the least cost to node z, $D_x(z)$, has changed from 7 to 3.
- The process of receiving updated distance vectors from neighbours, recomputing routing-table entries, and informing neighbours of changed costs of the least-cost path to a destination continues until no update messages are sent.
- The algorithm remains in the quiescent state until a link cost changes.

**COMPUTER NETWORKS****3.5.4 A Comparison of LS and DV Routing-algorithms**

Distance Vector Protocol	Link State Protocol
Entire routing-table is sent as an update	Updates are incremental & entire routing-table is not sent as update
Distance vector protocol send periodic update at every 30 or 90 second	Updates are triggered not periodic
Updates are broadcasted	Updates are multicasted
Updates are sent to directly connected neighbour only	Update are sent to entire network & to just directly connected neighbour
Routers don't have end to end visibility of entire network.	Routers have visibility of entire network of that area only.
Prone to routing loops	No routing loops
Each node talks to only its directly connected neighbors	Each node talks with all other nodes (via broadcast)



COMPUTER NETWORKS

3.5.5 Hierarchical Routing

- Two problems of a simple routing-algorithm:
 - 1) Scalability**
 - As no. of routers increases, overhead involved in computing & storing routing info increases.
 - 2) Administrative Autonomy**
 - An organization should be able to run and administer its network.
 - At the same time, the organization should be able to connect its network to internet.
- Both of these 2 problems can be solved by organizing routers into autonomous-system (AS).
- An autonomous system (AS) is a group of routers under the authority of a single administration.
For example: same ISP or same company network.
- Two types of routing-protocol:
 - 1) Intra-AS routing protocol: refers to routing inside an autonomous system.
 - 2) Inter-AS routing protocol: refers to routing between autonomous systems.

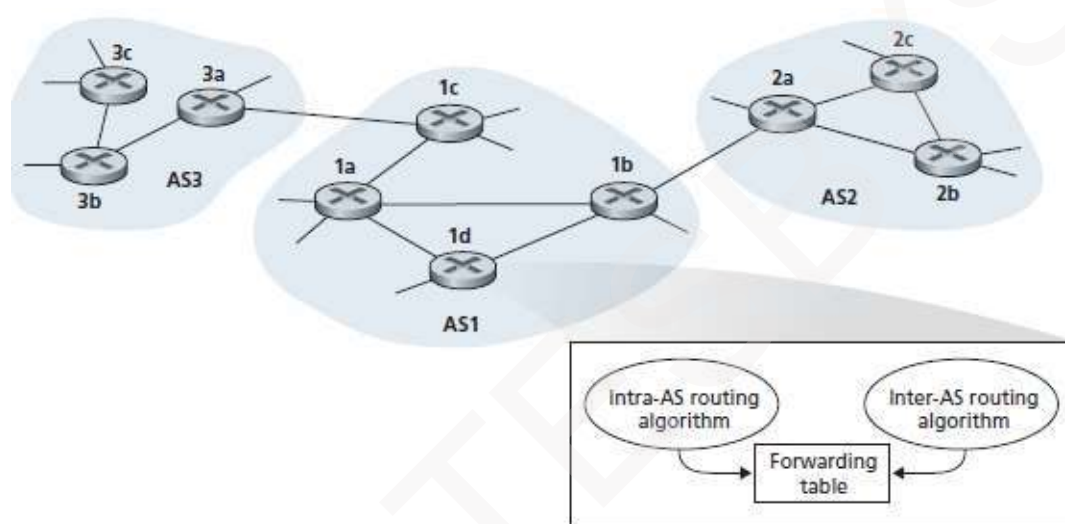


Figure 3.25: An example of interconnected autonomous-systems

3.5.5.1 Intra-AS Routing Protocol

- The routing-algorithm running within an autonomous-system is called intra-AS routing protocol.
- All routers within the same AS must run the same intra-AS routing protocol. For ex: RIP and OSPF
- Figure 3.25 provides a simple example with three ASs: AS1, AS2, and AS3.
- AS1 has four routers: 1a, 1b, 1c, and 1d. These four routers run the intra-AS routing protocol.
- Each router knows how to forward packets along the optimal path to any destination within AS1.

3.5.5.2 Intra-AS Routing Protocol

- The routing-algorithm running between 2 autonomous-systems is called inter-AS routing protocol.
- Gateway-routers are used to connect ASs to each other.
- Gateway-routers are responsible for forwarding packets to destinations outside the AS.
- Two main tasks of inter-AS routing protocol:
 - 1) Obtaining reachability information from neighboring Ass.
 - 2) Propagating the reachability information to all routers internal to the AS.
- The 2 communicating ASs must run the same inter-AS routing protocol. For ex: BGP.
- Figure 3.26 summarizes the steps in adding an outside-AS destination in a router's forwarding-table.

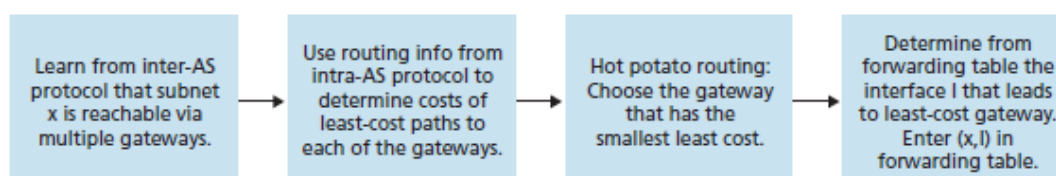


Figure 3.26: Steps in adding an outside-AS destination in a router's forwarding-table

“It is more important to know where you are going than to get there quickly.” —Mabel Newcomber



COMPUTER NETWORKS

3.6 Routing in the Internet

- Purpose of Routing protocols:
To determine the path taken by a datagram between source and destination.
- An autonomous-system (AS) is a collection of routers under the same administrative control.
- In AS, all routers run the same routing protocol among themselves.

3.6.1 Intra-AS Routing in the Internet: RIP

- Intra-AS routing protocols are also known as interior gateway protocols.
- An intra-AS routing protocol is used to determine how routing is performed within an AS.
- Most common intra-AS routing protocols:
1) Routing-information Protocol (RIP) and 2) Open Shortest Path First (OSPF)
- OSPF deployed in upper-tier ISPs whereas RIP is deployed in lower-tier ISPs & enterprise-networks.

3.6.1.1 RIP Protocol

- RIP is widely used for intra-AS routing in the Internet.
- RIP is a distance-vector protocol.
- RIP uses hop count as a cost metric. Each link has a cost of 1.
- Hop count refers to the no. of subnets traversed along the shortest path from source to destination.
- The maximum cost of a path is limited to 15.
- The distance vector is the current estimate of shortest path distances from router to subnets in AS.
- Consider an AS shown in Figure 3.27.

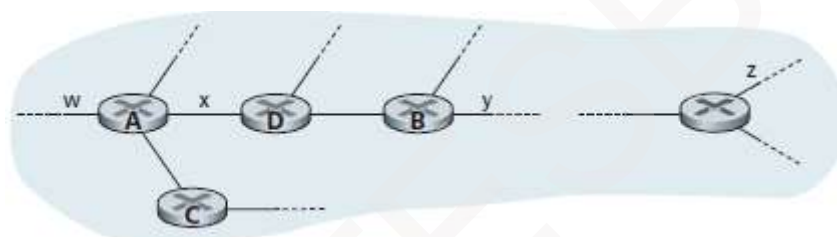


Figure 3.27: A portion of an autonomous-system

- Each router maintains a RIP table known as a routing-table.
- Figure 3.28 shows the routing-table for router D.

Destination Subnet	Next Router	Number of Hops to Destination
w	A	2
y	B	2
z	B	7
x	—	1

Figure 3.28: Routing-table in router D before receiving advertisement from router A

- Routers can send types of messages: 1) Response-message & 2) Request-message
 - 1) Response Message**
 - Using this message, the routers exchange routing updates with their neighbors every 30 secs.
 - If a router doesn't hear from its neighbor every 180 secs, then that neighbor is not reachable.
 - When this happens, RIP
 - modifies the local routing-table and
 - propagates then this information by sending advertisements to its neighbors.
 - The response-message contains
 - list of up to 25 destination subnets within the AS and
 - sender's distance to each of those subnets.
 - Response-messages are also known as advertisements.
 - 2) Request Message**
 - Using this message, router requests info about its neighbor's cost to a given destination.
- Both types of messages are sent over UDP using port# 520.
- The UDP segment is carried between routers in an IP datagram.



COMPUTER NETWORKS

3.6.2 Intra-AS Routing in the Internet: OSPF

- OSPF is widely used for intra-AS routing in the Internet.
- OSPF is a link-state protocol that uses
 - flooding of link-state information and
 - Dijkstra least-cost path algorithm.
- Here is how it works:
 - 1) A router constructs a complete topological map (a graph) of the entire autonomous-system.
 - 2) Then, the router runs Dijkstra's algorithm to determine a shortest-path tree to all subnets.
 - 3) Finally, the router broadcasts link state info to all other routers in the autonomous-system.
 - Specifically, the router broadcasts link state information
 - periodically at least once every 30 minutes and
 - whenever there is a change in a link's state. For ex: a change in up/down status.
- Individual link costs are configured by the network-administrator.
- OSPF advertisements are contained in OSPF messages that are carried directly by IP.
- HELLO message can be used to check whether the links are operational.
- The router can also obtain a neighboring router's database of network-wide link state.
- Some of the advanced features include:
 - 1) Security**
 - Exchanges between OSPF routers can be authenticated.
 - With authentication, only trusted routers can participate within an AS.
 - By default, OSPF packets between routers are not authenticated.
 - Two types of authentication can be configured: 1) Simple and 2) MD5.
 - i) Simple Authentication**
 - ✗ The same password is configured on each router.
 - ✗ Clearly, simple authentication is not very secure.
 - ii) MD5 Authentication**
 - ✗ This is based on shared secret keys that are configured in all the routers.
 - ✗ Here is how it works:
 - 1) The sending router
 - computes a MD5 hash on the content of packet
 - includes the resulting hash value in the packet and
 - sends the packet
 - 2) The receiving router
 - computes an MD5 hash of the packet
 - compares computed-hash value with the hash value carried in packet and
 - verifies the packet's authenticity
 - 2) Multiple Same Cost Paths**
 - When multiple paths to a destination have same cost, OSPF allows multiple paths to be used.
 - 3) Integrated Support for Unicast & Multicast Routing**
 - Multicast OSPF (MOSPF) provides simple extensions to OSPF to provide for multicast-routing.
 - MOSPF
 - uses the existing OSPF link database and
 - adds a new type of link-state advertisement to the existing broadcast mechanism.
 - 4) Support for Hierarchy within a Single Routing Domain**
 - An autonomous-system can be configured hierarchically into areas.
 - In area, an area-border-router is responsible for routing packets outside the area.
 - Exactly one OSPF area in the AS is configured to be the backbone-area.
 - The primary role of the backbone-area is to route traffic between the other areas in the AS.



COMPUTER NETWORKS

3.6.3 Inter-AS Routing: BGP

- BGP is widely used for inter-AS routing in the Internet.
- Using BGP, each AS can
 - 1) Obtain subnet reachability-information from neighboring ASs.
 - 2) Propagate the reachability-information to all routers internal to the AS.
 - 3) Determine good routes to subnets based on i) reachability-information and ii) AS policy.
- Using BGP, each subnet can advertise its existence to the rest of the Internet.

3.6.3.1 Basics

- Pairs of routers exchange routing-information over semi-permanent TCP connections using port-179.
- One TCP connection is used to connect 2 routers in 2 different autonomous-systems.
Semipermanent TCP connection is used to connect among routers within an autonomous-system.
- Two routers at the end of each connection are called peers.
The messages sent over the connection is called a session.
- Two types of session:
 - 1) External BGP (eBGP) session
 - This refers to a session that spans 2 autonomous-systems.
 - 2) Internal BGP (iBGP) session
 - This refers to a session between routers in the same AS.

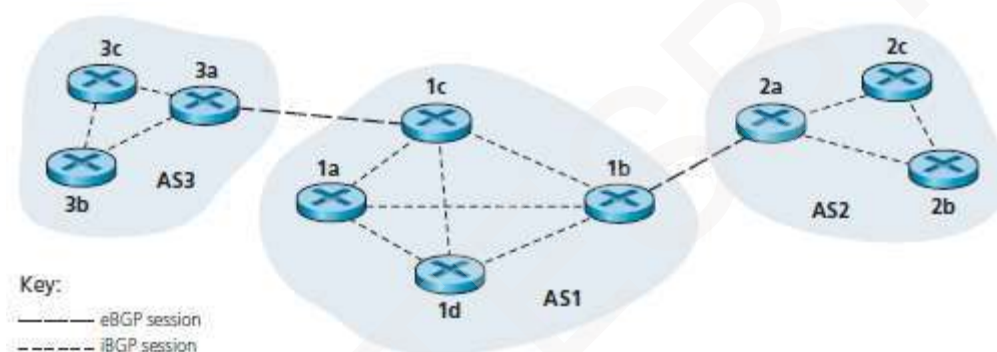


Figure 3.29: eBGP and iBGP sessions

- BGP operation is shown in Figure 3.29.
- The destinations are not hosts but instead are CIDRized prefixes.
- Each prefix represents a subnet or a collection of subnets.

3.6.3.2 Path Attributes & Routes

- An autonomous-system is identified by its globally unique ASN (Autonomous-System Number).
- A router advertises a prefix across a session.
- The router includes a number of attributes with the prefix.
- Two important attributes: 1) AS-PATH and 2) NEXT-HOP
 - 1) **AS-PATH**
 - This attribute contains the ASs through which the advertisement for the prefix has passed.
 - When a prefix is passed into an AS, the AS adds its ASN to the ASPATH attribute.
 - Routers use the AS-PATH attribute to detect and prevent looping advertisements.
 - Routers also use the AS-PATH attribute in choosing among multiple paths to the same prefix.
 - 2) **NEXT-HOP**
 - This attribute provides the critical link between the inter-AS and intra-AS routing protocols.
 - This attribute is the router-interface that begins the AS-PATH.
- BGP also includes
 - attributes which allow routers to assign preference-metrics to the routes.
 - attributes which indicate how the prefix was inserted into BGP at the origin AS.
- When a gateway-router receives a route-advertisement, the gateway-router decides
 - whether to accept or filter the route and
 - whether to set certain attributes such as the router preference metrics.



COMPUTER NETWORKS

3.6.3.3 Route Selection

- For 2 or more routes to the same prefix, the following elimination-rules are invoked sequentially:
 - 1) Routes are assigned a local preference value as one of their attributes.
 - 2) The local preference of a route
 - will be set by the router or
 - will be learned by another router in the same AS.
 - 3) From the remaining routes, the route with the shortest AS-PATH is selected.
 - 4) From the remaining routes, the route with the closest NEXT-HOP router is selected.
 - 5) If more than one route still remains, the router uses BGP identifiers to select the route.

3.6.3.4 Routing Policy

- Routing policy is illustrated as shown in Figure 3.30.
- Let A, B, C, W, X & Y = six interconnected autonomous-systems.
 - W, X & Y = three stub-networks.
 - A, B & C = three backbone provider networks.
- All traffic entering a stub-network must be destined for that network.
 - All traffic leaving a stub-network must have originated in that network.
- Clearly, W and Y are stub-networks.
- X is a multihomed stub-network, since X is connected to the rest of the n/w via 2 different providers
- X itself must be the source/destination of all traffic leaving/entering X.
- X will function as a stub-network if X has no paths to other destinations except itself.
- There are currently no official standards that govern how backbone ISPs route among themselves.

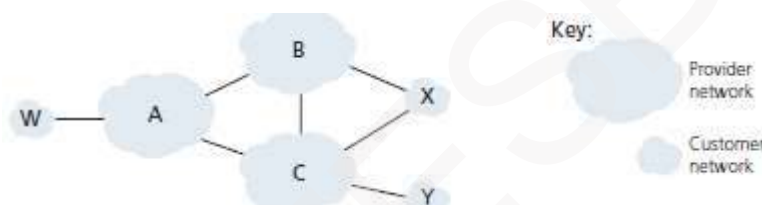


Figure 3.30: A simple BGP scenario



COMPUTER NETWORKS

3.7 Broadcast & Multicast Routing

3.7.1 Broadcast Routing Algorithms

- Broadcast-routing means delivering a packet from a source-node to all other nodes in the network.

3.7.1.1 N-way Unicast

- Given N destination-nodes, the source-node
 - makes N copies of the packet and
 - transmits then the N copies to the N destinations using unicast routing (Figure 3.31).
- Disadvantages:
 - 1) Inefficiency**
 - If source is connected to the n/w via single link, then N copies of packet will traverse this link.
 - 2) More Overhead & Complexity**
 - An implicit assumption is that the sender knows broadcast recipients and their addresses.
 - Obtaining this information adds more overhead and additional complexity to a protocol.
 - 3) Not suitable for Unicast Routing**
 - It is not good idea to depend on the unicast routing infrastructure to achieve broadcast.

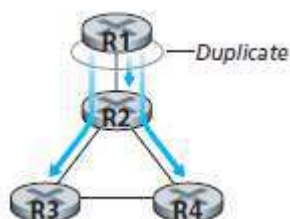


Figure 3.31: Duplicate creation/transmission

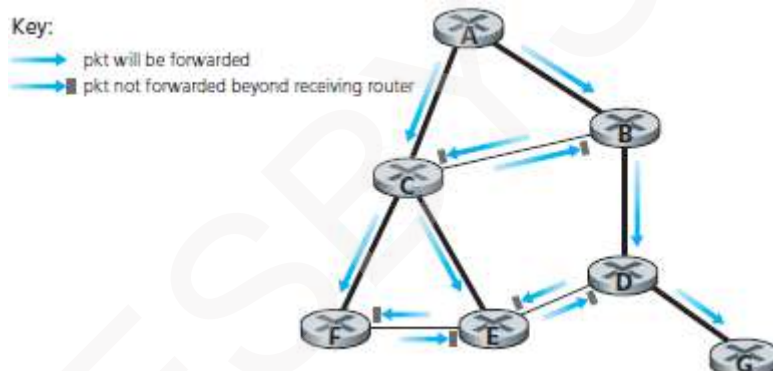


Figure 3.32: Reverse path forwarding

3.7.1.2 Uncontrolled Flooding

- The source-node sends a copy of the packet to all the neighbors.
- When a node receives a broadcast-packet, the node duplicates & forwards packet to all neighbors.
- In connected-graph, a copy of the broadcast-packet is delivered to all nodes in the graph.
- Disadvantages:
 - 1) If the graph has cycles, then copies of each broadcast-packet will cycle indefinitely.
 - 2) When a node is connected to 2 other nodes, the node creates & forwards multiple copies of packet
- Broadcast-storm refers to

“The endless multiplication of broadcast-packets which will eventually make the network useless.”

3.7.1.3 Controlled Flooding

- A node can avoid a broadcast-storm by judiciously choosing
 - when to flood a packet and when not to flood a packet.
- Two methods for controlled flooding:
 - 1) Sequence Number Controlled Flooding**
 - A source-node
 - puts its address as well as a broadcast sequence-number into a broadcast-packet
 - sends then the packet to all neighbors.
 - Each node maintains a list of the source-address & sequence# of each broadcast-packet.
 - When a node receives a broadcast-packet, the node checks whether the packet is in this list.
 - If so, the packet is dropped; if not, the packet is duplicated and forwarded to all neighbors.
 - 2) Reverse Path Forwarding (RPF)**
 - If a packet arrived on the link that has a path back to the source;
 - Then the router transmits the packet on all outgoing-links.
 - Otherwise, the router discards the incoming-packet.
 - Such a packet will be dropped. This is because
 - the router has already received a copy of this packet (Figure 3.32).



COMPUTER NETWORKS

3.7.1.4 Spanning - Tree Broadcast

- This is another approach to providing broadcast. (MST → Minimum Spanning Tree).
- Spanning-tree is a tree that contains each and every node in a graph.
- A spanning-tree whose cost is the minimum of all of the graph's spanning-trees is called a MST.
- Here is how it works (Figure 3.33):
 - 1) Firstly, the nodes construct a spanning-tree.
 - 2) The node sends broadcast-packet out on all incident links that belong to the spanning-tree.
 - 3) The receiving-node forwards the broadcast-packet to all neighbors in the spanning-tree.

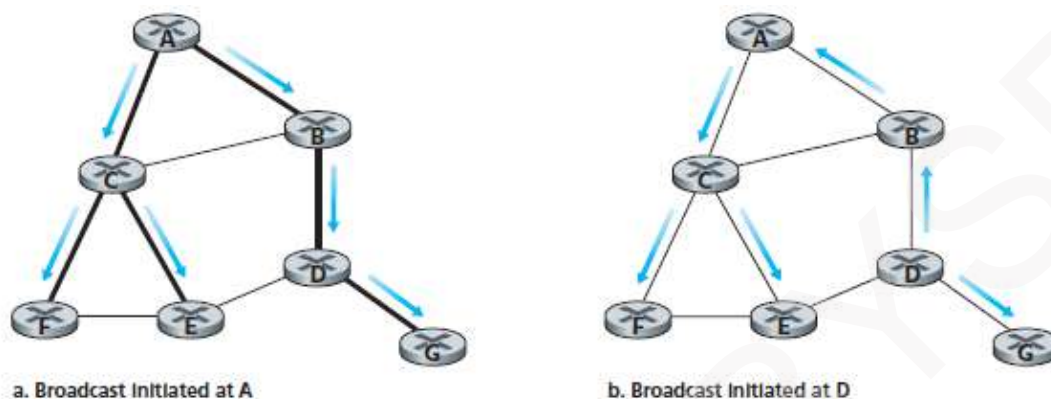


Figure 3.33: Broadcast along a spanning-tree

- Disadvantage:
Complex: The main complexity is the creation and maintenance of the spanning-tree.

3.7.1.4.1 Center Based Approach

- This is a method used for building a spanning-tree.
- Here is how it works:
 - 1) A center-node (rendezvous point or a core) is defined.
 - 2) Then, the nodes send unicast tree-join messages to the center-node.
 - 3) Finally, a tree-join message is forwarded toward the center until the message either
 - arrives at a node that already belongs to the spanning-tree or
 - arrives at the center.
- Figure 3.34 illustrates the construction of a center-based spanning-tree.

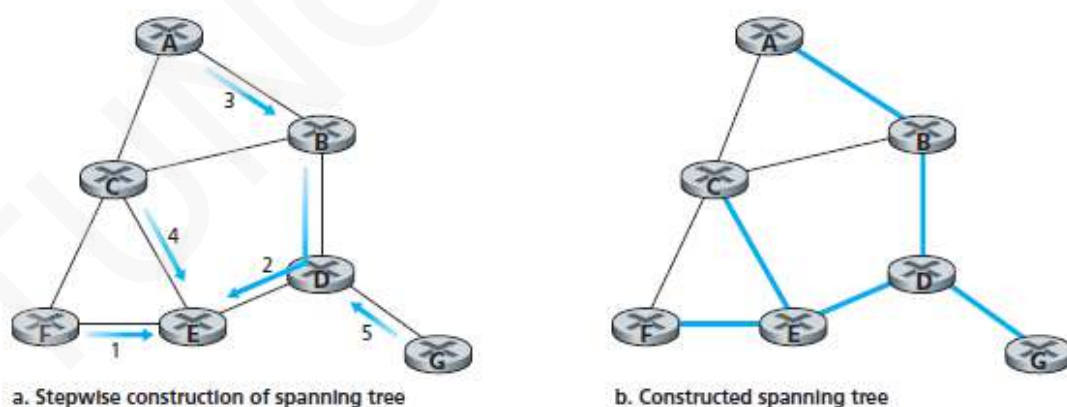


Figure 3.34: Center-based construction of a spanning-tree



COMPUTER NETWORKS

3.7.2 Multicast

- Multicasting means a multicast-packet is delivered to only a subset of network-nodes.
- A number of emerging network applications requires multicasting. These applications include
 - 1) Bulk data transfer (for ex: the transfer of a software upgrade)
 - 2) Streaming continuous media (for ex: the transfer of the audio/video)
 - 3) Shared data applications (for ex: a teleconferencing application)
 - 4) Data feeds (for ex: stock quotes)
 - 5) Web cache updating and
 - 6) Interactive gaming (for ex: multiplayer games).
- Two problems in multicast communication:
 - 1) How to identify the receivers of a multicast-packet.
 - 2) How to address a packet sent to these receivers.
- A multicast-packet is addressed using address indirection.
- A single identifier is used for the group of receivers.
- Using this single identifier, a copy of the packet is delivered to all multicast receivers.
- In the Internet, class-D IP address is the single identifier used to represent a group of receivers.
- The multicast-group abstraction is illustrated in Figure 3.35.

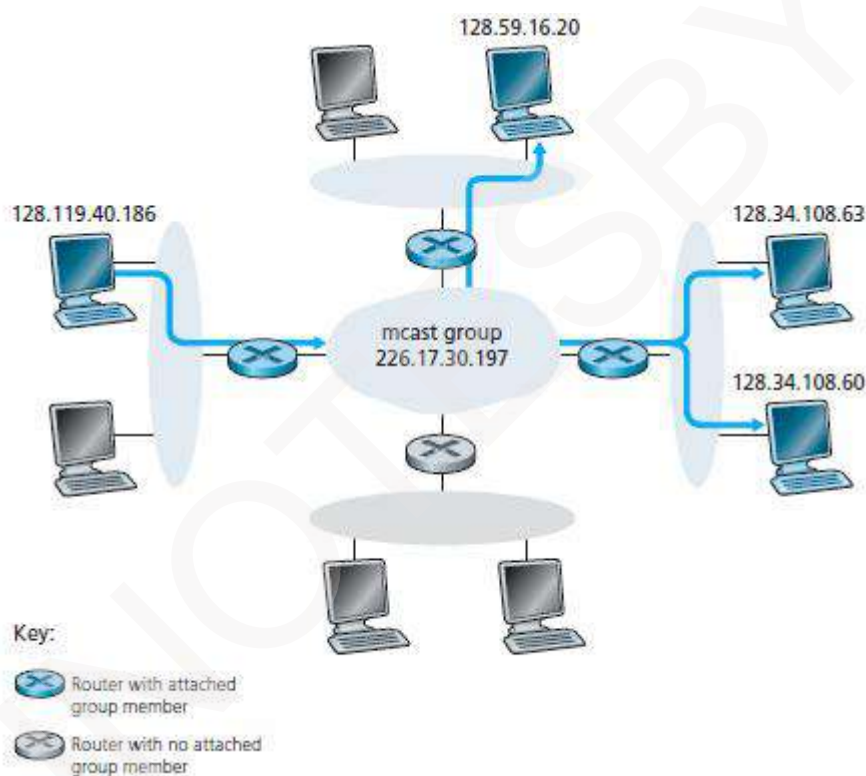


Figure 3.35: The multicast group: A datagram addressed to the group is delivered to all members of the multicast group



COMPUTER NETWORKS

3.7.2.1 IGMP

- In the Internet, the multicast consists of 2 components:
 - 1) IGMP (Internet Group Management Protocol)**
 - IGMP is a protocol that manages group membership.
 - It provides multicast-routers info about the membership-status of hosts connected to the n/w
 - The operations are i) Joining/Leaving a group and ii) monitoring membership
 - 2) Multicast Routing Protocols**
 - These protocols are used to coordinate the multicast-routers throughout the Internet.
 - A host places a multicast address in the destination address field to send packets to a set of hosts belonging to a group.
- The IGMP protocol operates between a host and its attached-router.
- Figure 3.36 shows three first-hop multicast-routers.

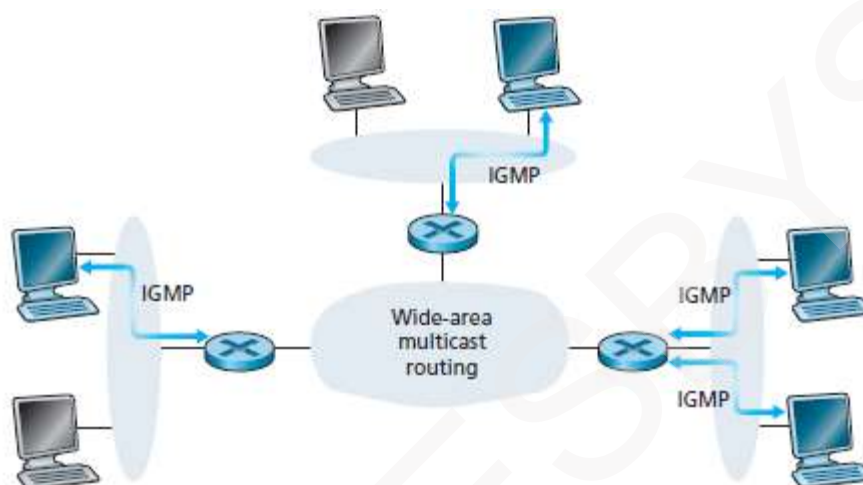


Figure 3.36: The two components of network-layer multicast in the Internet: IGMP and multicast-routing protocols

- IGMP messages are encapsulated within an IP datagram.
- Three types of message: 1) membership_query 2) membership_report 3) leave_group
 - 1) membership_query**
 - A host sends a membership-query message to find active group-members in the network.
 - 2) membership_report**
 - A host sends membership_report message when an application first joins a multicast-group.
 - The host sends this message w/o waiting for a membership_query message from the router.
 - 3) leave_group**
 - This message is optional.
 - The host sends this message to leave the multicast-group.

- How does a router detect when a host leaves the multicast-group?

Answer: The router infers that a host is no longer in the multicast-group if it no longer responds to a membership_query message. This is called soft state.



COMPUTER NETWORKS

3.7.2.2 Multicast Routing Algorithms

- The multicast-routing problem is illustrated in Figure 3.37.
- Two methods used for building a multicast-routing tree:
 - 1) Single group-shared tree.
 - 2) Source-specific routing tree.

1) Multicast Routing using a Group Shared Tree

- A single group-shared tree is used to distribute the traffic for all senders in the group.
- This is based on
 - Building a tree that includes all edge-routers & attached-hosts belonging to the multicast-group.
- In practice, a center-based approach is used to construct the multicast-routing tree.
- Edge-routers send join messages addressed to the center-node.
- Here is how it works:
 - 1) A center-node (rendezvous point or a core) is defined.
 - 2) Then, the edge-routers send unicast tree-join messages to the center-node.
 - 3) Finally, a tree-join message is forwarded toward the center until it either
 - arrives at a node that already belongs to the multicast tree or
 - arrives at the center.

2) Multicast Routing using a Source Based Tree

- A source-specific routing tree is constructed for each individual sender in the group.
- In practice, an RPF algorithm is used to construct a multicast forwarding tree.
- The solution to the problem of receiving unwanted multicast-packets under RPF is known as pruning.
- A multicast-router that has no attached-hosts will send a prune message to its upstream router.

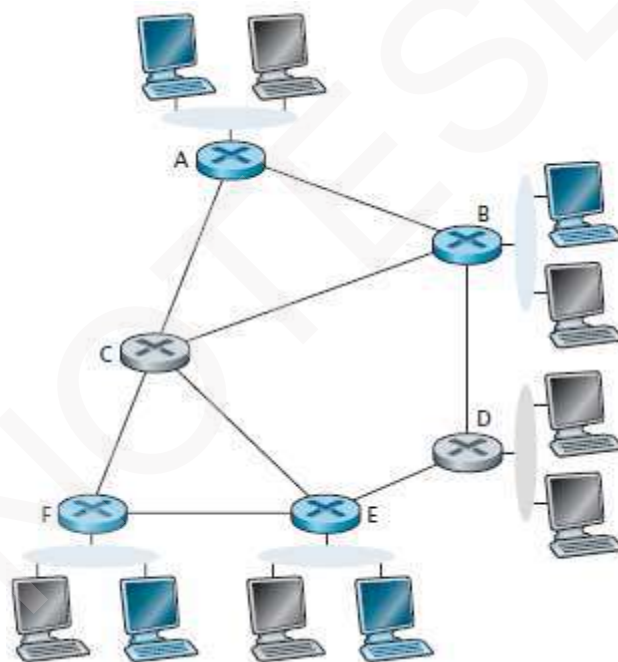


Figure 3.37: Multicast hosts, their attached routers, and other routers



COMPUTER NETWORKS

3.7.2.3 Multicast Routing in the Internet

- Three multicast routing protocols are:
 - 1) Distance Vector Multicast Routing Protocol (DVMRP)
 - 2) Protocol Independent Multicast (PIM) and
 - 3) Source Specific Multicast (SSM)

1) DVMRP

- DVMRP was the first multicast-routing protocol used in the Internet.
- DVMRP uses an RPF algorithm with pruning. (Reverse Path Forwarding).

2) PIM

- PIM is the most widely used multicast-routing protocol in the Internet.
- PIM divides multicast routing into sparse and dense mode.
 - i) Dense Mode**
 - Group-members are densely located.
 - Most of the routers in the area need to be involved in routing the data.
 - PIM dense mode is a flood-and-prune reverse path forwarding technique.
 - i) Sparse Mode**
 - The no. of routers with attached group-members is small with respect to total no. of routers.
 - Group-members are widely dispersed.
 - This uses rendezvous points to set up the multicast distribution tree.

3) SSM

- Only a single sender is allowed to send traffic into the multicast tree. This simplifies tree construction & maintenance.

**MODULE-WISE QUESTIONS****PART 1**

- 1) List out network-layer services. (2)
- 2) With a diagram, explain virtual circuit network. (6*)
- 3) With a diagram, explain datagram network. (6*)
- 4) Compare datagram & virtual-circuit networks. (4*)
- 5) With a diagram, explain 4 components of router. (6*)
- 6) With a diagram for each, explain 3 types of switching fabrics. (8*)
- 7) Briefly explain IP & its services. (4*)
- 8) With general format, explain various field of IPv4. (8*)
- 9) Explain fragmentation. Explain 3 fields related to fragmentation (6*)
- 10) Briefly explain IPv4 addressing. (8)
- 11) With a diagram, explain subnet addressing. (6*)
- 12) With a diagram, explain the working of DHCP. (6*)
- 13) With a diagram, explain the working of NAT. (6*)
- 14) Briefly explain ICMP. (8)
- 15) Explain changes from IPv4 to IPv6. (4*)
- 16) With general format, explain various field of IPv6. (8*)
- 17) Briefly explain IPv4 fields not present in IPv6. (4)
- 18) Compare IPv4 & IPv6. (4*)
- 19) Explain 2 strategies to make transition from IPv4 to IPv6. (8*)
- 20) Briefly explain IPSec & its services. (6)

PART 2

- 21) Define routing algorithm. Explain 3 classification of routing algorithm. (6*)
- 22) Explain Dijkstra's algorithm with example. (8*)
- 23) Explain Bellman-Ford algorithm with example. (8*)
- 24) Compare distance vector & link state protocols. (4*)
- 25) Explain the intra-AS and inter-AS routing protocols. (4*)
- 26) With a diagram, explain the working of RIP. (6*)
- 27) With a diagram, explain the working of OSPF. (6*)
- 28) With a diagram, explain the working of BGP. (8*)
- 29) What are disadvantages of N-way Unicast and Uncontrolled Flooding? (4)
- 30) Explain the following broadcast routing algorithms: (8*)
 - i) Controlled flooding
 - ii) Spanning - tree broadcast
- 31) Briefly explain multicasting & its applications. (4)
- 32) Briefly explain the working of IGMP. (4)
- 33) Briefly explain two methods used for building a multicast-routing tree. (6*)
- 34) Briefly explain three multicast routing protocols. (4)



MODULE 4: WIRELESS AND MOBILE NETWORKS

- 4.1 Cellular Internet Access
 - 4.1.1 An Overview of Cellular Network Architecture
 - 4.1.1.1 Cellular Network Architecture, 2G: Voice Connections to the Telephone Network
 - 4.1.2 3G Cellular Data Networks: Extending the Internet to Cellular Subscribers
 - 4.1.2.1 3G Core Network
 - 4.1.2.2 3G Radio Access Network: The Wireless Edge
 - 4.1.3 On to 4G: LTE
- 4.2 Mobility Management: Principles
 - 4.2.1 Addressing
 - 4.2.2 Routing to a Mobile Node
 - 4.2.2.1 Indirect Routing to a Mobile Node
 - 4.2.2.2 Direct Routing to a Mobile Node
 - 4.2.2.2.1 Challenges in Direct Routing
- 4.3 Mobile IP
 - 4.3.1 Agent Discovery
 - 4.3.1.1 Agent Advertisement
 - 4.3.1.2 Agent Solicitation
 - 4.3.2 Registration with the Home Agent
- 4.4 Managing Mobility in Cellular Networks
 - 4.4.1 Routing Calls to a Mobile User
 - 4.4.2 Handoffs in GSM
 - 4.4.2.1 Problem of Handoffs
- 4.5 Wireless and Mobility: Impact on Higher-Layer Protocols
 - 4.5.1 Commonalities between Mobile IP & GSM
 - 4.5.2 TCP Congestion Control Approaches



MODULE 4: WIRELESS AND MOBILE NETWORKS

4.1 Cellular Internet Access

4.1.1 An Overview of Cellular Network Architecture

- Cellular technology can be classified into following generations:
 - 1) First Generation (1G)**
 - 1G systems were analog FDMA systems designed exclusively for voice-only communication.
 - 2) Second Generation (2G)**
 - 2G systems were also designed for voice (GSM → Global System for Mobile communication).
 - Later, the 2G was extended to support data (i.e., Internet) service.
 - GSM was the major system evolved in the second generation:
 - 3) Third Generation (3G)**
 - 3G systems were also designed for voice and data.
 - More emphasis was given on data capabilities and higher-speed radio access links.



COMPUTER NETWORKS

4.1.1.1 Cellular Network Architecture, 2G: Voice Connections to the Telephone Network

- The region covered by cellular-network is divided into no. of geographic coverage-areas called cells.
- Each cell contains a BTS (Base Transceiver Station) (Figure 4.1).
- BTS is responsible for delivering the signals to/from the mobile-stations in the cell.
- The coverage-area of a cell depends on following factors:
 - 1) The transmitting power of the BTS.
 - 2) The transmitting power of the user devices.
 - 3) Obstructing buildings in the cell.
 - 4) The height of base-station antennas.
- The 2G systems use combined FDM/TDM for the air-interface.
- In combined FDM/TDM systems,
 - 1) The channel is divided into a number of frequency sub-bands.
 - 2) Within each sub-band, time is partitioned into frames and slots.

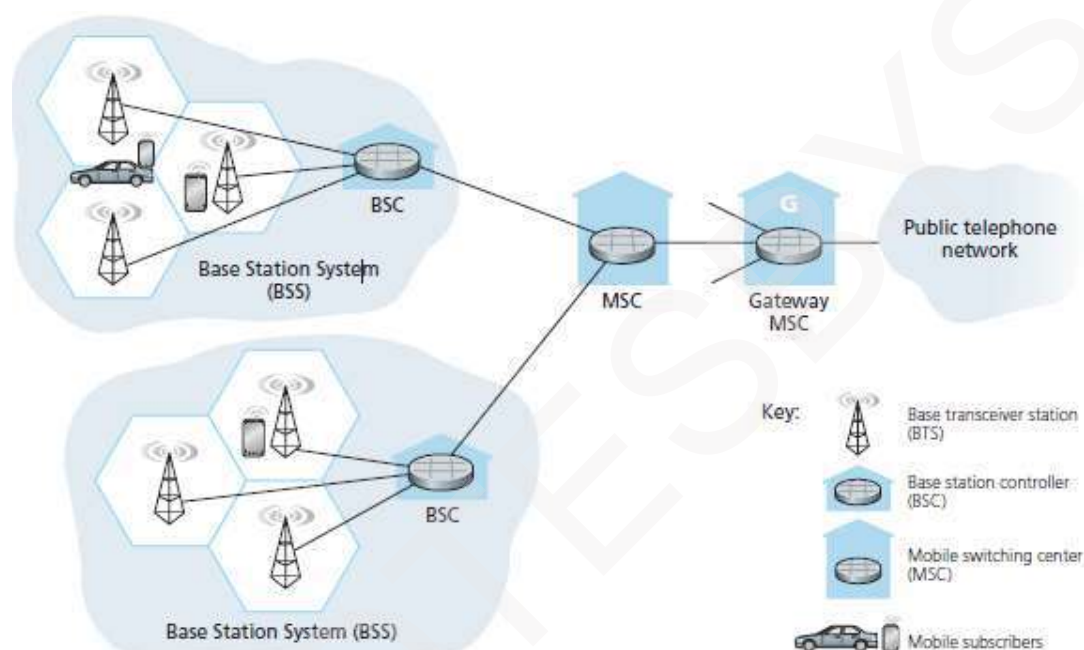


Figure 4.1: Components of the GSM 2G cellular network architecture

- The GSM network contains many BSCs (Base Station Controllers).
- Main responsibilities of the BSC:
 - 1) Providing service to many BTSs.
 - 2) Allocating radio-channels to mobile-users.
 - 3) Performing paging.
 - 4) Performing handoff of mobile-users.
- BSS (Base Station System) contains the BSC and its controlled BTSs.
- A MSC (Mobile Switching Center) contains upto 5 BSCs. This results in approx 200K subscribers/MSC.
- Main responsibilities of the MSC:
 - 1) User authorization & accounting
 - 2) Call establishment & teardown and
 - 3) Handoff.
- A cellular-provider's network will have a number of special MSCs known as gateway MSCs.
- Gateway MSCs are used to connect the provider's cellular-network to the public telephone-network.

COMPUTER NETWORKS

4.1.2 3G Cellular Data Networks: Extending the Internet to Cellular Subscribers

4.1.2.1 3G Core Network

- 3G system architecture is shown in Figure 4.2.
- Main responsibilities of the core-network:
 - 1) Connects radio access-networks (RANs) to the public Internet.
 - 2) Interoperates with components of the existing voice-network.

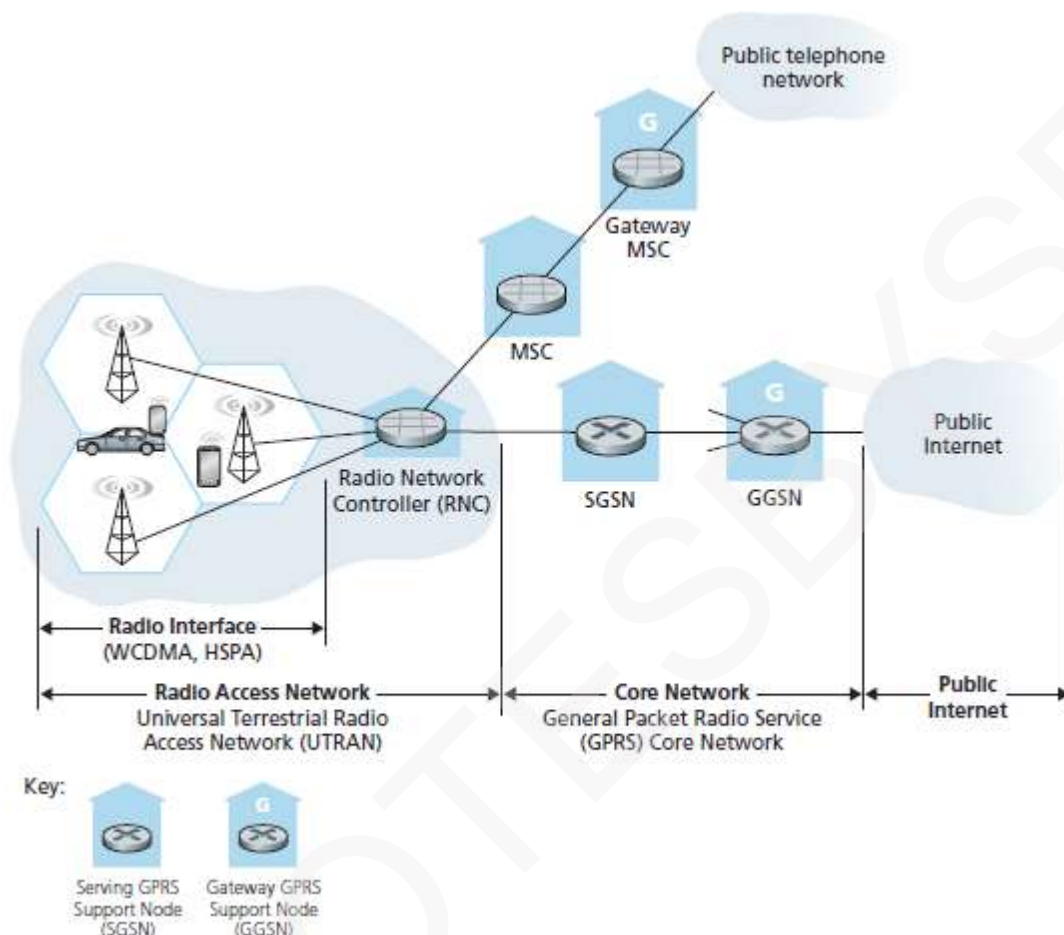


Figure 4.2: 3G system architecture

- The idea of 3G designers:
 - “Leave the existing voice-network untouched;
 - Add additional data functionality in parallel to the existing voice-network.”
 - Two types of nodes in the core-network:
 - 1) Serving GPRS Support Node (SGSN) and
 - 2) Gateway GPRS Support Node (GGSN).
- 1) SGSN**
- An SGSN is responsible for delivering data to/from the mobile-nodes in the RAN.
 - Main responsibilities of the SGSN:
 - 1) Interacting with the MSC of voice-network.
 - 2) Providing user authorization and handoff.
 - 3) Maintaining location information about active mobile-nodes.
 - 4) Performing data forwarding between a GGSN & mobile-nodes in the RAN.
- 2) GGSN**
- A GGSN acts as a gateway.
 - The GGSN is used to connect multiple SGSNs into the larger Internet.
 - To the outside world, the GGSN looks like any other router.
 - The mobility of the nodes within the GGSN's network is hidden from the outside world.



COMPUTER NETWORKS

4.1.2.2 3G Radio Access Network: The Wireless Edge

- The RAN is the wireless first-hop network that the 3G user sees.
- The RNC (Radio Network Controller) typically controls several cell BTSs
- Each cell's wireless-link operates between the mobile-nodes and a BTS.
- The RNC connects to both the circuit-switched voice-network and the packet-switched Internet.
- UMTS (Universal Mobile Telecommunications Service) is a widely deployed 3G technology.
- UMTS uses CDMA technique known as DS-WCDMA within TDMA slots.
- TDMA slots, in turn, are available on multiple frequencies.
- The data-service associated with the WCDMA specification is known as HSP.
(HSP → High Speed Packet access DS-WCDMA → Direct Sequence Wideband CDMA)

4.1.3 On to 4G: LTE

- The 4G systems have 2 main improvements over 3G systems:

- 1) Evolved Packet Core (EPC) and
- 2) LTE Radio Access-network.

1) EPC

- It is an "all-IP" network i.e. both voice & data will be carried in IP datagrams.
- Main responsibilities of the EPC:
 - 1) Combines the circuit-switched voice-network and the packet-switched data-network.
 - 2) Manages network-resources to provide high QoS.
 - 3) Allows multiple types of access-networks (such as 2G/3G) to attach to the core-n/w.

2) LTE Radio Access Network

- LTE uses a combination of FDM & TDM on the downstream-channel known as OFDM (Orthogonal Frequency Division Multiplexing).
- Each mobile-node is allocated one or more time-slots in one or more channel-frequencies.
- Figure 4.3 shows an allocation of 8 time-slots over 4 frequencies.
- By allocating more time-slots, a mobile-node is able to achieve higher data-rates.
- Multiple-input, multiple output (MIMO) antennas can be used to increase the data-rate.
- For example:

In the downstream direction, maximum data-rate = 100 Mbps.

In the upstream direction, maximum data-rate = 50 Mbps.

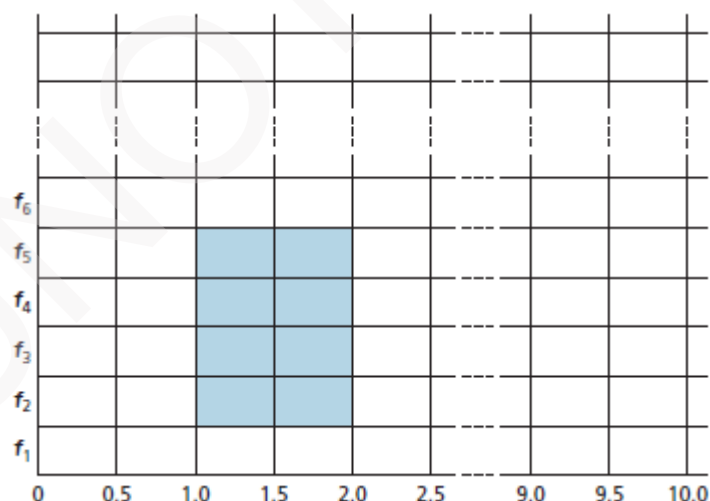


Figure 4.3: Twenty 0.5 ms slots organized into 10 ms frames at each frequency. An eight-slot allocation is shown shaded.



COMPUTER NETWORKS

4.2 Mobility Management: Principles

- A mobile-node is one that changes its point of attachment into the network over time.

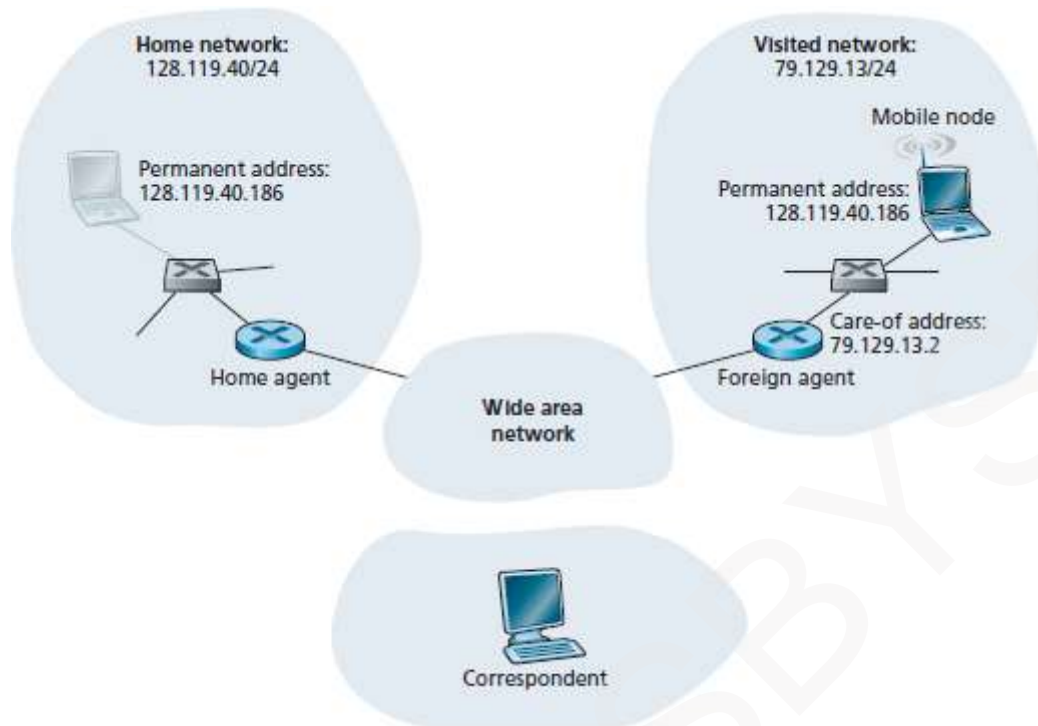


Figure 4.4: Initial elements of a mobile network architecture

- Five elements of mobile-network architecture:

1) Home Network

- Home network is a network that is the permanent home of the mobile-node (ex: smartphone)

2) Home Agent

- The home-agent is a router within the home-network. (COA → care-of-address)
- The home-agent performs the mobility management functions on behalf of the mobile-node.
- The home-agent interacts with a foreign-agent to track the mobile-node's COA.

3) Foreign Network

- Foreign network is a network to which the mobile-node moves.
- The foreign-network is also known as visited-network.

4) Foreign Agent

- The foreign-agent is a router within the foreign-network.
- The foreign-agent performs the mobility management functions on behalf of the mobile-node.
- Two roles of foreign-agent:
 - 1) Create a care-of-address (COA) for the mobile-node.
 - ✘ The network portion of the COA must match with the foreign-network.
 - ✘ Two addresses are associated with a mobile-node:
 - i) Permanent-address and
 - ii) COA (known as a foreign address).
 - 2) Inform the home-agent that the mobile-node is resident in the foreign-agent's n/w.
 - ✘ The foreign-agent has the given COA.

5) Correspondent

- The entity wishing to communicate with the mobile-node is known as a correspondent.
- Figure 4.4 illustrates these concepts.



COMPUTER NETWORKS

4.2.1 Addressing

- When a mobile-node moves from one network to another, the mobile-node must keep its address.
- Thus, user-mobility will be transparent to network-applications.
- When a mobile-node is in a foreign-network, the mobile-node's traffic is routed to foreign-network.
- The foreign-network advertises to the neighbors that it knows a route to mobile-node's permanent-address.
- Then, these neighbors propagate the routing-information throughout the network.
- When the mobile-node moves from one foreign-network to another, the new foreign-network advertises a new route to the mobile-node.
- Disadvantage:
 - Scalability: If mobility management is the responsibility of routers, the routers have to maintain forwarding-table entries for potentially millions of mobile-nodes.



COMPUTER NETWORKS

4.2.2 Routing to a Mobile Node

- Two approaches are 1) indirect routing and 2) direct routing.

4.2.2.1 Indirect Routing to a Mobile Node

- Four steps are involved. Figure 4.5 illustrates the 4 steps.

Step 1

- The correspondent
 - addresses the datagram to the mobile-node's permanent-address and
 - routes the datagram to the mobile-node's home-network.

Step 2

- Home-agent encapsulates the correspondent's original datagram within a larger datagram.
- This larger datagram is addressed & delivered to the mobile-node's COA.

Step 3

- The foreign-agent receives and decapsulates the datagram.
- The foreign-agent forwards the original datagram to the mobile-node.

Step 4

- The mobile-node directly routes the datagram to the correspondent.
- There is no need to route the datagram back through the home-agent.

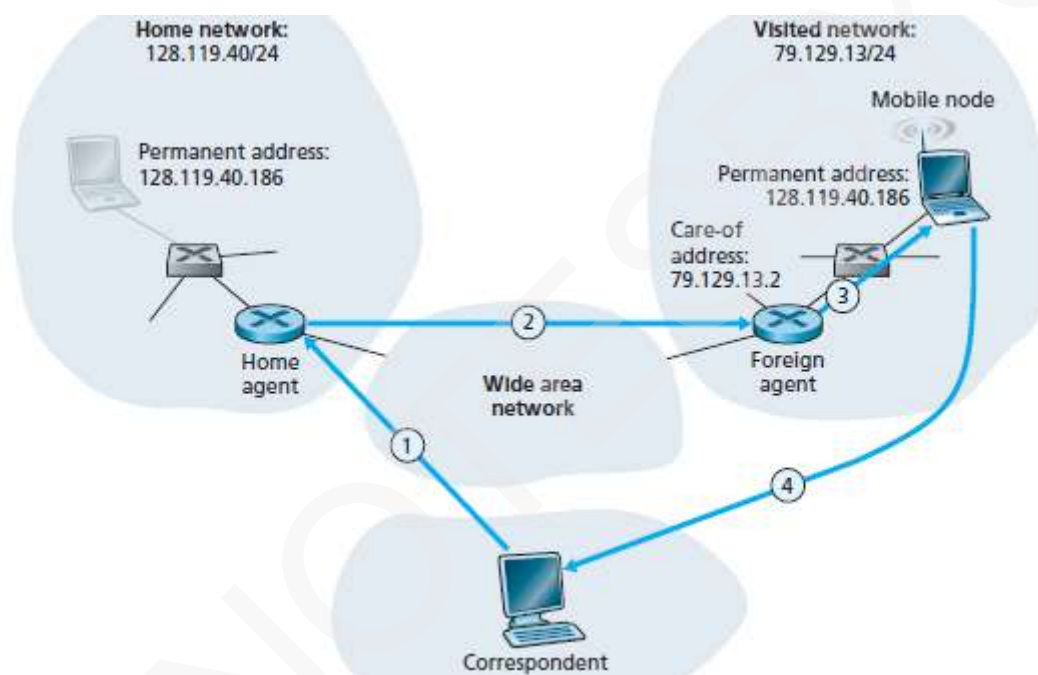


Figure 4.5: Indirect routing to a mobile node

- New functionality required to support mobility:

1) A Mobile-node-to-Foreign-agent Protocol

- The mobile-node will register with the foreign-agent when attaching to the foreign-network.
- Similarly, a mobile-node will deregister with the foreign-agent leaving the foreign-network.

2) A Foreign-agent-to-Home-agent Registration Protocol

- The foreign-agent will register the mobile-node's COA with the home-agent.
- A foreign-agent need not explicitly deregister a COA when a mobile-node leaves its network.

3) A Home-agent Datagram Encapsulation Protocol

- Encapsulation of correspondent's original datagram within a datagram addressed to the COA.

4) A Foreign-agent Decapsulation Protocol

- Extraction of the correspondent's original datagram from the encapsulated-datagram.
- Then, the forwarding of the original datagram to the mobile-node.

- Disadvantage of Indirect Routing: Suffers from triangle routing problem: The datagrams addressed to the mobile-node must be routed first to the home-agent and then to the foreign-network, even when an efficient route exists b/w the correspondent and the mobile-node.

- Solution: Use direct routing.

COMPUTER NETWORKS

4.2.2.2 Direct Routing to a Mobile Node

- Four steps are involved. Figure 4.6 illustrates the 4 steps.

Steps 1 & 2

- A correspondent-agent in the correspondent's n/w first learns the COA of the mobile-node.
- This can be done by having the correspondent-agent query the home-agent.

Steps 3 & 4

- Then, the correspondent-agent forwards datagrams directly to the mobile-node's COA.

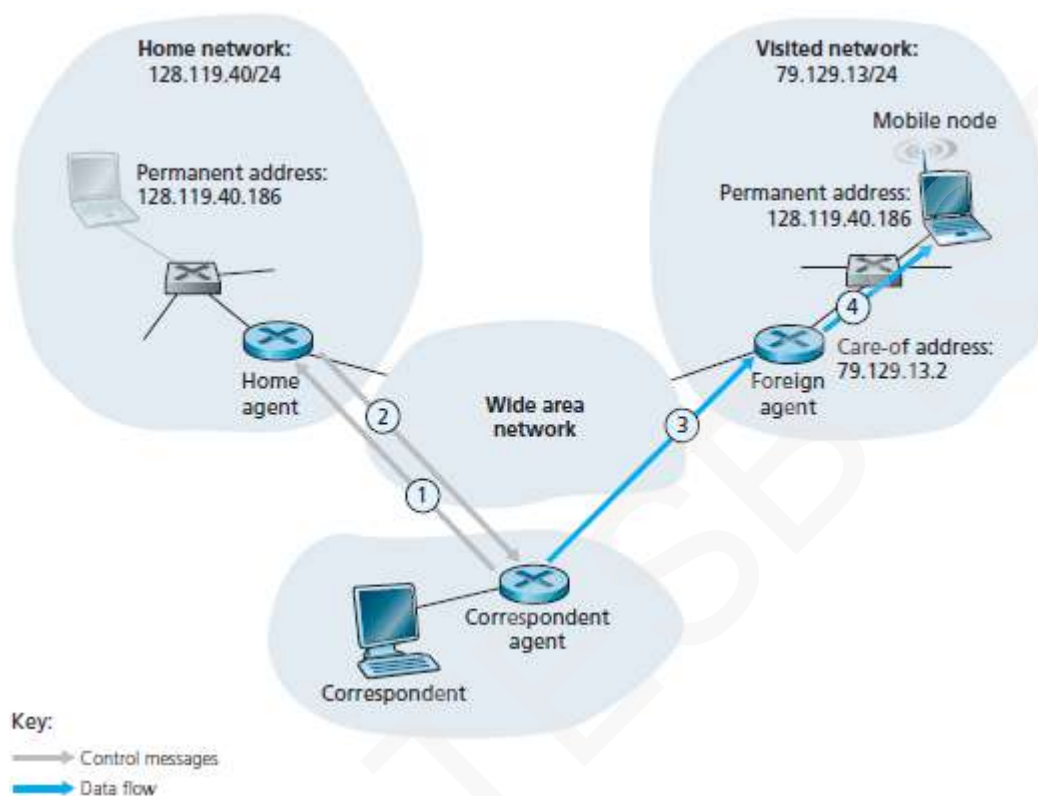


Figure 4.6: Direct routing to a mobile user



COMPUTER NETWORKS

4.2.2.2.1 Challenges in Direct Routing

• Two additional challenges:

1) The correspondent-agent needs a mobile location protocol to query the home-agent to obtain the mobile-node's COA (steps 1 & 2 in Figure 4.6).

2) Problem: When the mobile-node moves from one foreign-network to another, how will data be forwarded to the new foreign-network?

Solution: Use anchor foreign-agent.

- An anchor foreign-agent refers to a foreign-agent in the foreign-network where the mobile-node was first found. (step 1 in Figure 4.7).
- When the mobile-node moves to a new foreign-network (step 2), the mobile-node registers with the new foreign-agent (step 3).
- The new foreign-agent provides the anchor foreign-agent with the mobile-node's new COA (step 4).
- When the anchor foreign-agent receives an encapsulated-datagram, the anchor re-encapsulates and forwards the datagram to the mobile-node (step 5).

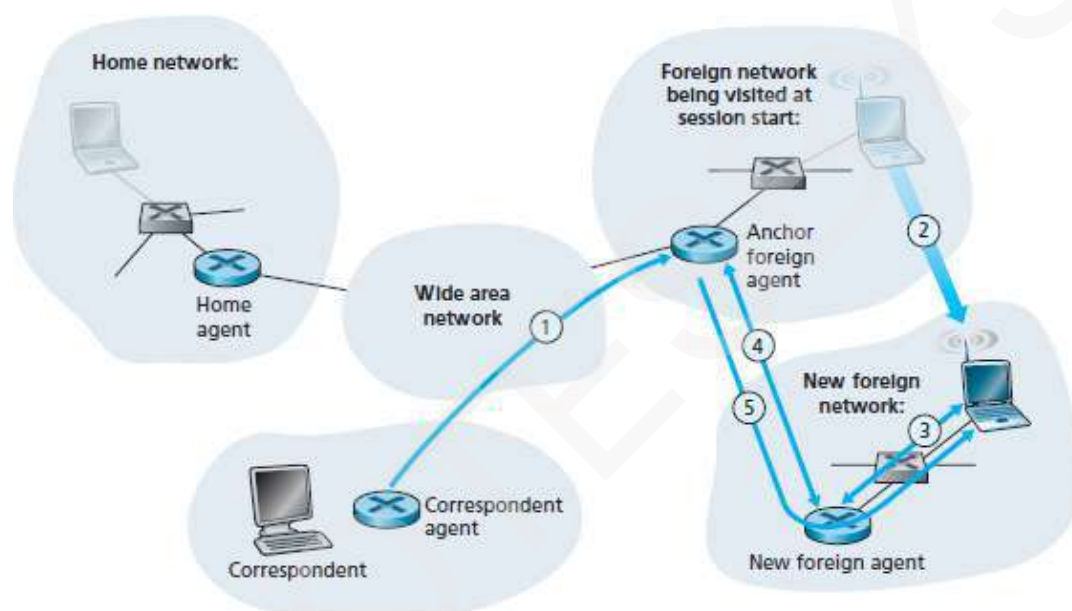


Figure 4.7: Mobile transfer between networks with direct routing



COMPUTER NETWORKS

4.3 Mobile IP

- Mobile IP is the extension of IP protocol.
- Mobile IP allows laptops (or smartphones) to be connected to the Internet.
- Services of Mobile IP:
 - 1) Support for many different modes of operation.
 - 2) Multiple ways for agents and mobile-nodes to discover each other.
 - 3) Use of single or multiple COAs.
 - 4) Multiple forms of encapsulation.
- Three main parts of mobile IP:
 - 1) Agent Discovery**
 - Mobile IP defines the protocols used by a home or foreign-agent to advertise its services to mobile-nodes.
 - It also defines the protocols for mobile-nodes to solicit the services of a foreign or home-agent.
 - 2) Registration with the Home Agent**
 - Mobile IP defines the protocols used by the mobile-node to register COAs with the home-agent.
 - 3) Indirect Routing of Datagrams**
 - Mobile IP defines the manner in which datagrams are forwarded to mobile-nodes by a home-agent.
 - It also defines
 - rules for forwarding datagrams
 - rules for handling error conditions and
 - several forms of encapsulation



COMPUTER NETWORKS

4.3.1 Agent Discovery

- A mobile-node arriving to a new network must learn the identity of the corresponding foreign or home-agent. This process is known as agent discovery.
- Two methods to perform agent discovery:
 - 1) Via agent advertisement and
 - 2) Via agent solicitation.

4.3.1.1 Agent Advertisement

- A foreign or home-agent advertises its services using a router discovery protocol.
- The agent periodically broadcasts a router discovery message on all links.
- The router discovery message contains
 - 1) IP address of the agent and
 - 2) A mobility agent advertisement extension.
- Five main fields in the extension:
 - 1) Home Agent (H)**
 - This bit indicates that the agent is a home-agent for the network in which it resides.
 - 2) Foreign Agent (F)**
 - This bit indicates that the agent is a foreign-agent for the network in which it resides.
 - 3) Registration required (R)**
 - This bit indicates that a mobile-user in this network must register with a foreign-agent.
 - 4) M, G Encapsulation**
 - These bits indicate whether an encapsulation other than IP-in-IP encapsulation will be used.
 - 5) Care-of-address (COA) Fields**
 - This field indicates a list of one or more care-of-addresses provided by the foreign-agent.
 - Figure 4.8 illustrates some of the key fields in the agent advertisement message.

4.3.1.2 Agent Solicitation

- A mobile-node wanting to learn about agents can broadcast an agent solicitation message.
- An agent receiving the solicitation will unicast an agent advertisement directly to the mobile-node.

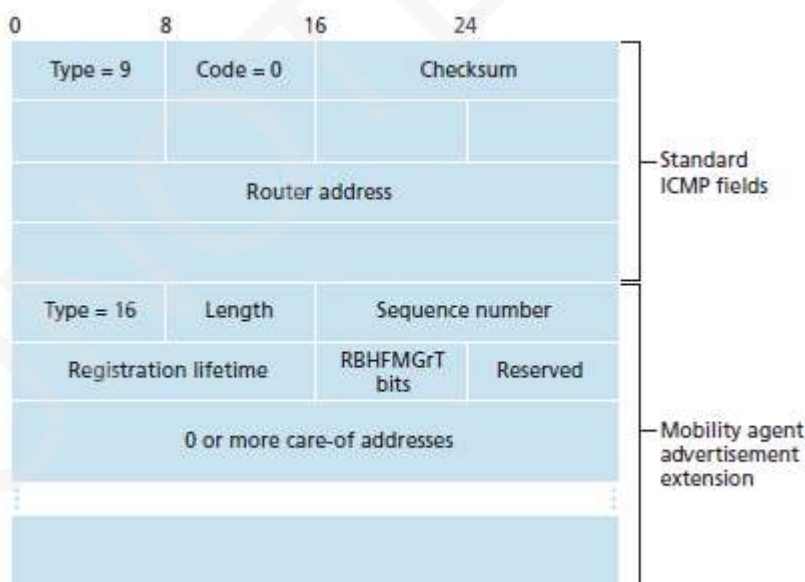


Figure 4.8: ICMP router discovery message with mobility agent advertisement extension



COMPUTER NETWORKS

4.3.2 Registration with the Home Agent

- Address must be registered with the home-agent. This can be done in 2 ways:
 - Via the foreign-agent who then registers the COA with the home-agent.
 - By the mobile IP node itself.

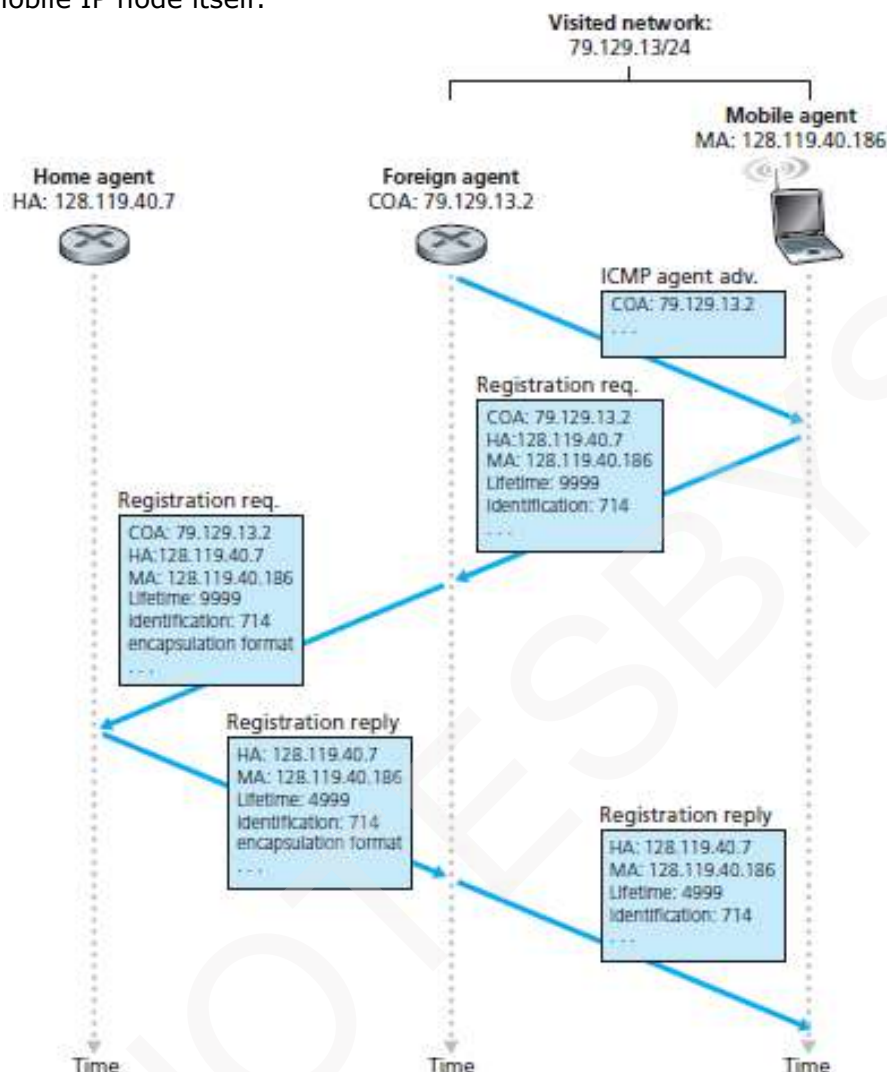


Figure 4.9: Agent advertisement and mobile IP registration

- Four steps are involved. Figure 4.9 illustrates the 4 steps.
 - When a mobile receives a foreign-agent advertisement, the mobile sends a registration-request to the foreign-agent.
 - The registration-request contains
 - COA advertised by the foreign-agent
 - address of the home-agent (HA)
 - permanent-address of the mobile (MA)
 - registration identification and
 - requested lifetime of the registration.
 - The requested registration lifetime indicates number of seconds the registration is valid.
 - If registration is not renewed within the specified lifetime, the registration will become invalid.
 - When the foreign-agent receives the registration-request, the foreign-agent records the mobile's permanent IP address.
 - The foreign-agent then sends a registration-request to the home-agent.
 - When home-agent receives the registration-request, the home-agent checks for correctness.
 - The home-agent binds the mobile's permanent IP address with the COA.
 - The home-agent sends a registration-reply.
 - The foreign-agent receives and forwards the registration-reply to the mobile-node.



COMPUTER NETWORKS

4.4 Managing Mobility in Cellular Networks

- GSM adopts an indirect routing approach.
- The mobile-user's home-network is referred to as home public land mobile-network (home PLMN).
- The home-network is the cellular-provider with which the mobile-user has a subscription.
- Mobile-user's visited-network is referred to as the visited public land mobile-network (visited PLMN).
- The visited PLMN is the network in which the mobile-user is currently residing.
- The responsibilities of the home and visited-networks are quite different:
 - 1) The home-network maintains a database known as the HLR (home location register).
 - The HLR contains
 - permanent phone-number
 - subscriber profile information.
 - information about the current locations of the subscribers.
 - In home-network, a home MSC is contacted by correspondent when a call is placed to mobile.
 - 2) The visited-network maintains a database known as the VLR (visitor location register).
 - The VLR contains an entry for each mobile-user that is currently in the network.
 - Thus, VLR entries come and go as mobile-users enter and leave the network.
 - A VLR is co-located with MSC that coordinates the setup of a call to & from the visited-n/w.

4.4.1 Routing Calls to a Mobile User

- Three steps are involved. Figure 4.10 illustrates the 3 steps.
 - 1) The correspondent dials the mobile-user's phone-number.
 - The call is routed from the correspondent via PSTN to the home MSC in the home-network.
 - 2) The home MSC
 - receives the call and
 - interrogates the HLR to determine the location of the mobile-user.
 - The HLR returns the roaming number.
 - If HLR doesn't have the roaming number, it returns the address of VLR in the visited-network.
 - Then, the home MSC queries the VLR to obtain the roaming number of the mobile-node.
 - 3) The home MSC sets up the call through the network to the MSC in the visited-network.
 - The call is completed.

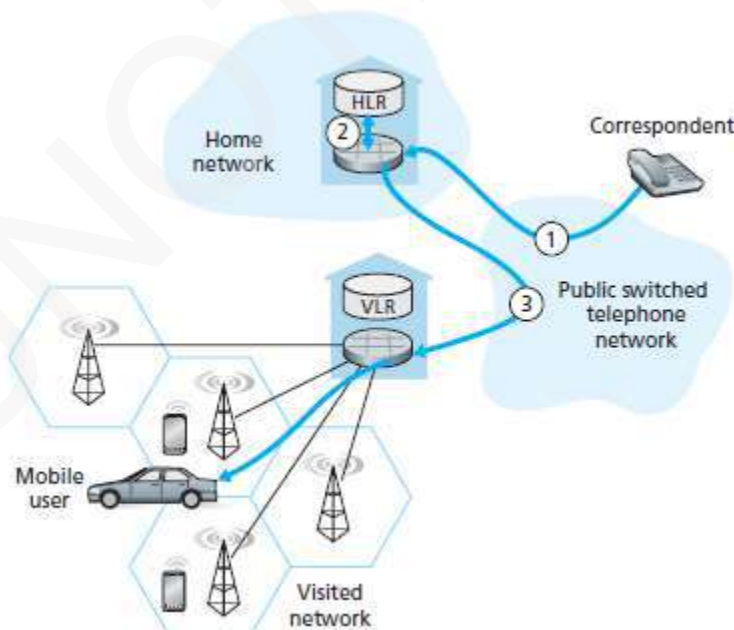


Figure 4.10: Placing a call to a mobile user: indirect routing



COMPUTER NETWORKS

4.4.2 Handoffs in GSM

- A handoff occurs when a mobile-station moves from one base-station to another during a call.
- As shown in Figure 4.11,
 - 1) Before handoff, a call is initially routed to the mobile through old base-station.
 - 2) After handoff, the call is routed to the mobile through another new base-station.
- Two reasons for handoff:
 - 1) The Call may be Dropped**
 - Because the signal between the current base-station and the mobile may have weakened.
 - 2) To reduce Congestion**
 - Because a cell may be overloaded because of handling a large number of calls.
 - This congestion may be reduced by handing off mobiles to less congested cells.

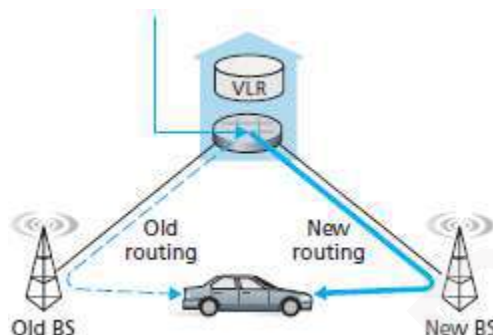


Figure 4.11: Handoff scenario between base stations with a common MSC

- Eight steps are involved. Figure 4.12 illustrates the steps involved when a hand off occurs.
 - 1) Old base-station (BS) informs both visited MSC & new BS that a handoff is about to happen.
 - 2) The visited MSC performs following tasks:
 - i) Initiates path setup to the new BS.
 - ii) Allocates the resources needed to carry the rerouted call.
 - iii) Signals the new BS that a handoff is about to occur.
 - 3) The new BS allocates and activates a radio-channel for the mobile.
 - 4) The new BS informs both visited MSC and old BS that the new path is set up.
 - 5) The mobile is informed to perform a handoff.
 - 6) The mobile & new BS exchange signaling messages to fully activate the new channel.
 - 7) The mobile sends a handoff complete message to the new BS.
 - ✕ This message is then forwarded to the visited MSC.
 - ✕ The visited MSC then reroutes the ongoing-call to the mobile via the new BS.
 - 8) The resources allocated along the path to the old BS are released.

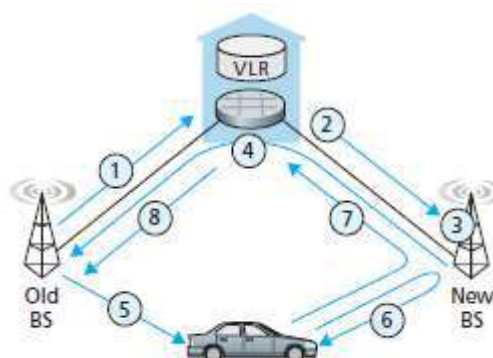


Figure 4.12: A handoff between base stations with a common MSC

COMPUTER NETWORKS

4.4.2.1 Problem of Handoffs

- Problem: How inter-MSC handoff occurs?
 - Solution: Use an anchor MSC.
 - This operation is shown in Figure 4.13.
 - The anchor MSC is the MSC visited by the mobile when a call first begins.
 - When a mobile moves from one MSC to another, the ongoing-call is rerouted from the anchor MSC to the new visited MSC.

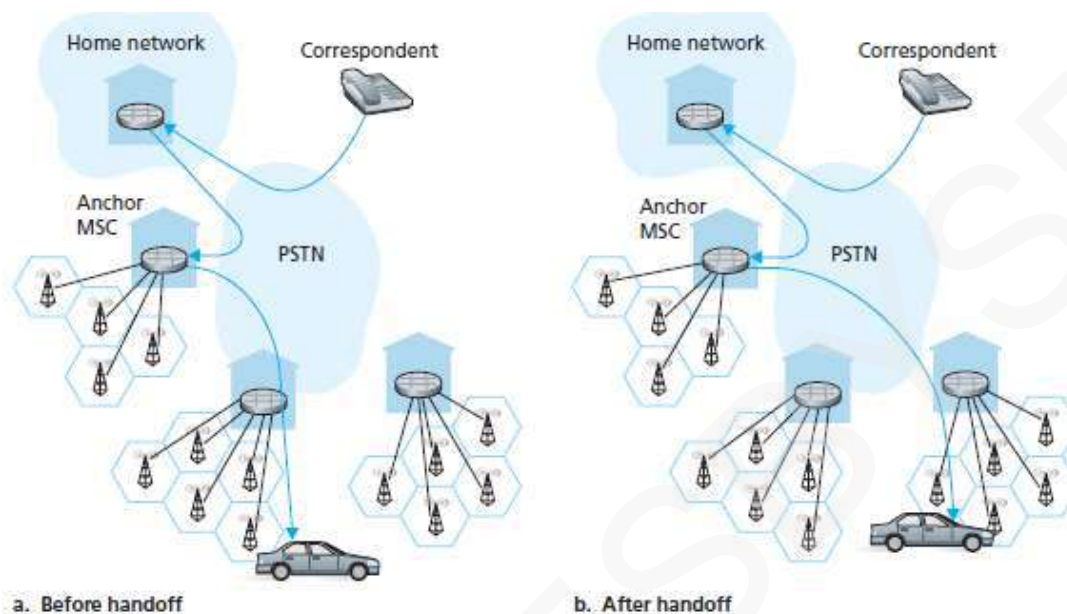


Figure 4.13: Rerouting via the anchor MSC

**COMPUTER NETWORKS****4.5 Wireless and Mobility: Impact on Higher-Layer Protocols****4.5.1 Commonalities between Mobile IP & GSM**

Table 4.1: Commonalities between mobile IP and GSM mobility

GSM element	Comment on GSM element	Mobile IP element
Home system	Network to which the mobile user's permanent phone number belongs.	Home network
Gateway mobile switching center or simply home MSC, Home location register (HLR)	Home MSC: point of contact to obtain routable address of mobile user. HLR: database in home system containing permanent phone number, profile information, current location of mobile user, subscription information.	Home agent
Visited system	Network other than home system where mobile user is currently residing.	Visited network.
Visited mobile services switching center, Visitor location register (VLR)	Visited MSC: responsible for setting up calls to/from mobile nodes in cells associated with MSC. VLR: temporary database entry in visited system, containing subscription information for each visiting mobile user.	Foreign agent
Mobile station roaming number (MSRN) or simply roaming number	Routable address for telephone call segment between home MSC and visited MSC, visible to neither the mobile nor the correspondent.	Care-of-address

4.5.2 TCP Congestion Control Approaches

- Three approaches for dealing with TCP's congestion-control: 1) Local recovery 2) TCP sender awareness of wireless-links and 3) Split-connection approaches

1) Local Recovery

- These protocols recover from bit-errors. For ex: ARQ protocol.

2) TCP Sender Awareness of Wireless Links

- The sender and receiver must be aware of the existence of a wireless-link.
- The sender and receiver must be to distinguish between
 - congestive losses occurring at the wired-network and
 - congestive losses occurring at the wireless-network.
- Sender & receiver invoke congestion-control only in response to congestive wired-n/w losses.

3) Split Connection

- The end-to-end connection b/w mobile-user & other end-point are broken into 2 connections:
 - 1) One connection from the mobile-host to the wireless access-point.
 - 2) Another connection from the wireless access-point to the other end-point.
- Thus, the end-to-end connection is formed by concatenation of a wireless-part & a wired-part.

**MODULE-WISE QUESTIONS****PART 1**

- 1) With a diagram, explain various components of 2G system architecture. (6*)
- 2) With a diagram, explain various components of 3G system architecture. (8*)
- 3) With a diagram, explain various components of mobile network architecture. (6*)
- 4) With a diagram for each, explain indirect & direct routing. (8*)
- 5) With a diagram, explain the problem & its solution in direct routing. (4)

PART 2

- 6) Briefly explain mobile IP & its services. (4*)
- 7) With a diagram, explain the following with respect to mobile IP: (8*)
 - i) Agent Discovery
 - ii) Registration with the home agent
- 8) With a diagram, explain how a call is placed to a mobile user. (4*)
- 9) With a diagram, explain the concept of handoffs. (6*)
- 10) Compare mobile IP and GSM mobility. (4)
- 11) Briefly explain TCP congestion control approaches. (4)



MODULE 5: MULTIMEDIA NETWORKING

- 5.1 Multimedia Networking Applications
 - 5.1.1 Properties of Video
 - 5.1.2 Properties of Audio
 - 5.1.2.1 PCM Encoder
 - 5.1.2.2 PCM Decoder
 - 5.1.3 Types of Multimedia Network Applications
 - 5.1.3.1 Streaming Stored Audio & Video
 - 5.1.3.2 Conversational Voice- and Video-over-IP
 - 5.1.3.3 Streaming Live Audio & Video
- 5.2 Streaming Stored Video
 - 5.2.1 UDP Streaming
 - 5.2.2 HTTP Streaming
 - 5.2.2.1 Client Application Buffer & TCP Buffers
 - 5.2.2.2 Early Termination & Repositioning the Video
 - 5.2.3 Adaptive Streaming & DASH
 - 5.2.3.1 DASH
 - 5.2.4 CDN
 - 5.2.4.1 Motivation for CDN
 - 5.2.4.2 CDN Types
 - 5.2.4.3 CDN Operation
 - 5.2.4.4 Cluster Selection Strategies
- 5.3 Voice-over-IP
 - 5.3.1 Limitations of the Best Effort IP Service
 - 5.3.1.1 Packet Loss
 - 5.3.1.2 End-to-End Delay
 - 5.3.1.3 Packet Jitter
 - 5.3.2 Removing Jitter at the Receiver for Audio
 - 5.3.3 Recovering from Packet Loss
 - 5.3.3.1 FEC
 - 5.3.3.2 Interleaving
 - 5.3.3.3 Error Concealment
- 5.4 Protocols for Real-Time Conversational Applications
 - 5.4.1 RTP
 - 5.4.1.1 RTP Basics
 - 5.4.1.2 RTP Packet Header Fields
 - 5.4.2 SIP
 - 5.4.2.1 Setting up a Call to a Known IP Address
 - 5.4.2.2 SIP Messages
 - 5.4.2.3 Name Translation & User Location
- 5.5 Network Support for Multimedia
 - 5.5.1 Dimensioning Best-Effort Networks
 - 5.5.2 Providing Multiple Classes of Service
 - 5.5.2.1 Motivating Scenarios
 - 5.5.2.1.1 Approaches for Providing Isolation among Traffic Classes
 - 5.5.2.1.1.1 Using Traffic Policing
 - 5.5.2.1.1.2 Using Packet Scheduling
 - 5.5.2.2 Scheduling Mechanisms
 - 5.5.2.2.1 FIFO
 - 5.5.2.2.2 Priority Queuing
 - 5.5.2.2.3 RRQ
 - 5.5.2.2.4 WFQ



COMPUTER NETWORKS

5.5.2.3 Policing: The Leaky Bucket

5.5.2.3.1 Leaky Bucket Operation

5.5.3 DiffServ

5.5.4 Per-Connection QoS Guarantees: Resource Reservation & Call Admission

5.5.4.1 Mechanisms for Guaranteed QoS

VTUNOTESBYSRI



MODULE 5: MULTIMEDIA NETWORKING

5.1 Multimedia Networking Applications

- A multimedia network application can be defined as any application that employs audio or video.

5.1.1 Properties of Video

1) High Bit Rate

- Video distributed over the Internet use
 - 100 kbps for low-quality video conferencing.
 - 3 Mbps for streaming high-definition (HD) movies.
- The higher the bit-rate,
 - better the image quality and
 - better the overall user viewing experience.

2) Video Compression

- A video can be compressed, thereby trading off video-quality with bit-rate.
- A video is a sequence of images, displayed at a constant rate.
For example: 24 or 30 images per second.
- An uncompressed digital image consists of an array of pixels.
- Each pixel is encoded into a number of bits to represent luminance and color.
- There are two types of redundancy in video:
 - 1) Spatial Redundancy**
 - An image that consists of mostly white space has a high degree of redundancy.
 - These images can be efficiently compressed without sacrificing image quality.
 - 2) Temporal Redundancy**
 - Temporal redundancy reflects repetition from image to subsequent image.
 - For example:
 - If image & subsequent image are same, re-encoding of subsequent image can be avoided.



COMPUTER NETWORKS

5.1.2 Properties of Audio

- PCM (Pulse Code Modulation) is a technique used to change an analog signal to digital data (digitization).
- PCM consists of 1) Encoder at the sender and 2) Decoder at the receiver.

5.1.2.1 PCM Encoder

- Digital audio has lower bandwidth requirements than video.
- Consider how analog audio is converted to a digital-signal:
- The analog audio-signal is sampled at some fixed rate. This operation is referred to as sampling.
- For example: 8000 samples per second.
- The value of each sample is an arbitrary real number.
- Each sample is then rounded to one of a finite number of values. This process is called quantization.
- The number of such finite values is called as quantization-values.
- The number of quantization-values is typically a power of 2. For ex: $256(2^8)$ quantization-values.
- Each of the quantization-values is represented by a fixed number of bits.
- For example:
 - If there are $256(2^8)$ quantization-values, then each value is represented by 8 bits.
- Bit representations of all values are then concatenated to form digital representation of the signal. This process is called encoding.
- For example:
 - If an analog-signal is sampled at 8000 samples per second & each sample is represented by 8 bits, then the digital-signal will have a rate of 64000 bits per second ($8000*8=64000$).

5.1.2.2 PCM Decoder

- For playback through audio speakers, the digital-signal can be converted back to an analog-signal. This process is called decoding.
- However, the decoded analog-signal is only an approximation of the original signal.
- The sound quality may be noticeably degraded.
- The decoded signal can better approximate the original analog-signal by increasing
 - i) sampling rate and
 - ii) number of quantization-values,
- Thus, there is a trade-off between
 - quality of the decoded signal and
 - bit-rate & storage requirements of the digital-signal.



COMPUTER NETWORKS

5.1.3 Types of Multimedia Network Applications

- Three broad categories of multimedia applications:
 - 1) Streaming stored audio/video
 - 2) Conversational voice/video-over-IP and
 - 3) Streaming live audio/video.

5.1.3.1 Streaming Stored Audio & Video

- The underlying medium is prerecorded video. For example: a movie.
- These prerecorded videos are placed on servers.
- The users send requests to the servers to view the videos on-demand.
- Nowadays, many Internet companies provide streaming video. For example: YouTube.
- Three key distinguishing features of streaming stored video:

1) Streaming

- The client begins video playout within few seconds after it begins receiving the video from the server.
- At the same time,
 - i) The client will be playing out from one location in the video.
 - ii) The client will be receiving later parts of the video from the server.
- This technique avoids having to download the entire video-file before playout begins.

2) Interactivity

- The media is prerecorded, so the user may pause, reposition or fast-forward through video-content.
- The response time should be less than a few seconds.

3) Continuous Playout

- Once playout of the video begins, it should proceed according to the original timing of the recording.
- The data must be received from the server in time for its playout at the client.
Otherwise, users experience video-frame skipping (or freezing).

5.1.3.2 Conversational Voice- and Video-over-IP

- Real-time conversational voice over the Internet is often referred to as Internet telephony.
- It is also commonly called Voice-over-IP (VoIP).
- Conversational video includes the video of the participants as well as their voices.
- Most of today's voice applications allow users to create conferences with three or more participants.
- Nowadays, many Internet companies provide voice application. For example: Skype & Google Talk.
- Two parameters are particularly important for voice applications:
 - 1) Timing considerations and
 - 2) Tolerance of data loss
- Timing considerations are important because voice applications are highly delay-sensitive.
- Loss-tolerant means
Occasional loss only causes occasional glitches in audio playback & these losses can be partially/fully hidden.

5.1.3.3 Streaming Live Audio & Video

- These applications are similar to broadcast radio, except that transmission takes place over Internet.
- These applications allow a user to receive a live radio transmitted from any corner of the world.
- For example: live cricket commentary.
- Today, thousands of radio stations around the world are broadcasting content over the Internet.
- Live broadcast applications have many users who receive the same audio program at the same time.
- The network must provide an average throughput that is larger than the video consumption rate.



COMPUTER NETWORKS

5.2 Streaming Stored Video

- Prerecorded videos are placed on servers.
- Users send requests to these servers to view the videos on-demand.
- The media is prerecorded, so the user may pause, reposition or fast-forward through video-content.
- Three categories of applications:
 - 1) UDP streaming
 - 2) HTTP streaming and
 - 3) Adaptive HTTP streaming.
- A main characteristic of video-streaming is the extensive use of client-side buffering.
- Two advantages of client-side buffering:
 - 1) Client-side buffering can mitigate effects of varying end-to-end delays
 - 2) This can mitigate effects of varying amounts of available bandwidth b/w server & client.

5.2.1 UDP Streaming

- The server transmits video at a rate that matches the client's video consumption rate.
- The server transmits the video-chunks over UDP at a steady rate.
- UDP does not employ a congestion-control mechanism.
- Therefore, the server can push packets into the network at the video consumption rate.
- Typically, UDP streaming uses a small client-side buffer. (RTP → Real-Time Transport Protocol).
- Using RTP, the server encapsulates the video-chunks within transport packets.
- The client & server also maintain a control-connection over which the client sends commands (such as pause, resume and reposition).
- The RTSP (Real-Time Streaming Protocol) is a popular open protocol for a control-connection.
- Disadvantages:
 - 1) Unreliability**
 - UDP streaming can fail to provide continuous playout of varying amt of available bandwidth
 - 2) Costly & Complex**
 - A media control server (RTSP) is required
 - to process client-to-server interactivity requests and
 - to track client-state for each ongoing client-session.
 - This increases the overall cost and complexity of deploying a large-scale application.
 - 3) Firewall Problem**
 - Many firewalls are configured to block UDP traffic.
 - This prevents the users behind the firewalls from receiving the video.



COMPUTER NETWORKS

5.2.2 HTTP Streaming

- The video is stored in an HTTP server as an ordinary file with a specific URL.
- Here is how it works:
 - 1) When a user wants to see the video, the client
 - establishes a TCP connection with the server and
 - issues an HTTP GET request for that URL.
 - 2) Then, the server responds with the video file, within an HTTP response message.
 - 3) On client side, the bytes are collected in a client application buffer.
 - 4) Once no. of bytes in this buffer exceeds a specific threshold, the client begins playback.
- Advantages:
 - 1) Not Costly & Complex**
 - Streaming over HTTP avoids the need for a media control server (RTSP).
 - This reduces the cost of deploying a large-scale application.
 - 2) No Firewall Problem**
 - The use of HTTP over TCP also allows the video to traverse firewalls and NATs more easily.
 - 3) Prefetching Video**
 - The client downloads the video at a rate higher than the consumption rate.
 - Thus, prefetching video-frames that are to be consumed in the future.
 - This prefetched video is stored in the client application buffer
- Nowadays, most video-streaming applications use HTTP streaming. For example: YouTube

5.2.2.1 Client Application Buffer & TCP Buffers

- Figure 5.1 illustrates the interaction between client and server for HTTP streaming.
- On the server side,
 - 1) The bytes of the video file are sent into the server's socket.
 - 2) Then, the bytes are placed in the TCP send buffer before.
 - 3) Finally, the bytes are transmitted into the Internet.
- On the client side,
 - 1) The application (media-player) reads bytes from the TCP receive-buffer (thro client-socket)
 - 2) Then, the application places the bytes into the client-buffer.
 - 3) At the same time, the application periodically
 - grabs video-frames from the client-buffer
 - decompresses the frames and
 - displays the frames on the user's screen.

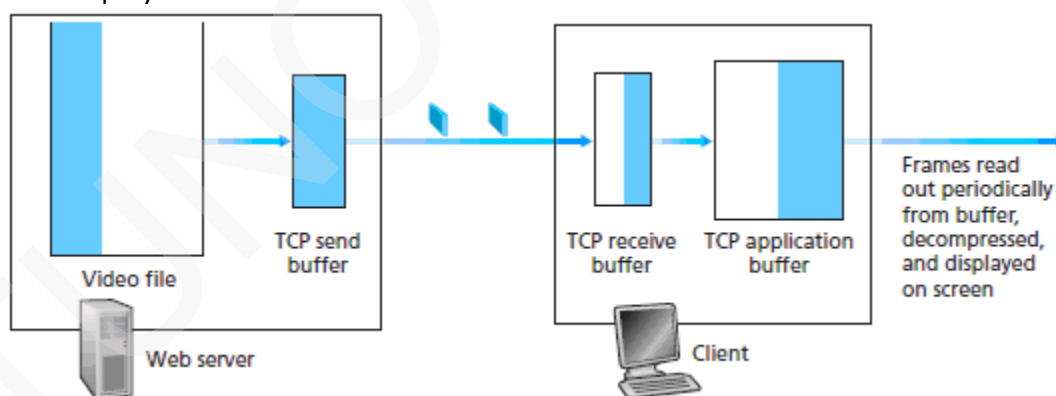


Figure 5.1: Streaming stored video over HTTP/TCP

5.2.2.2 Early Termination & Repositioning the Video

- HTTP streaming systems make use of the byte-range header in the HTTP GET request message.
- Byte-range header specifies the range of bytes the client currently wants to retrieve from the video.
- This is particularly useful when the user wants to reposition to a future point in the video.
- When the user repositions to a new position, the client sends a new HTTP request.
- When server receives new HTTP request, the server sends the requested-bytes.
- Disadvantage:
 - When a user repositions to a future point in the video, some prefetched-but-not-yet-viewed data will go unwatched. This results in a waste of bandwidth and server-resources.

“Too many of us are not living our dreams because we are living our fears.” —Les Brown



COMPUTER NETWORKS

5.2.3 Adaptive Streaming & DASH

- Problem with HTTP streaming:
All clients receive the same encoding of video, despite the large variations in the amount of bandwidth available to different clients.
Solution: Use DASH (Dynamic Adaptive Streaming over HTTP).

5.2.3.1 DASH

- The video is encoded into several different versions.
- Each version has a different bit-rate and a different quality level.
- Two main tasks:
 - 1) The client dynamically requests video-chunks from the different versions: low & high.
 - i) When the available bandwidth is high, the client selects chunks from a high-rate version.
For ex: Fiber connections can receive a high-quality version.
 - ii) When the available bandwidth is low, the client naturally selects from a low-rate version.
For ex: 3G connections can receive a low-quality version.
 - 2) The client adapts to the available bandwidth if end-to-end bandwidth changes during session.
 - This feature is particularly important for mobile-users.
 - The mobile-users see their bandwidth fluctuate as they move with respect to base-stations.
- HTTP server stores following files:
 - 1) Each video version with a different URL.
 - 2) Manifest file provides a URL for each version along with its bit-rate.
- Here is how it works:
 - 1) First, the client requests the manifest file and learns about the various versions.
 - 2) Then, the client selects one chunk at a time by specifying
 - URL and
 - byte range in an HTTP GET request message.
 - 3) While downloading chunks, the client
 - measures the received bandwidth and
 - runs a rate determination-algorithm.
 - i) If measured-bandwidth is high, client will choose chunk from high-rate version.
 - ii) If measured-bandwidth is low, client will choose chunk from low-rate version
 - 4) Therefore, DASH allows the client to freely switch among different quality-levels.
- Advantages:
 - 1) DASH can achieve continuous playout at the best possible quality level w/o frame freezing.
 - 2) Server-side scalability is improved: Because
 - the client maintains the intelligence to determine which chunk to send next.
 - 3) Client can use HTTP byte-range request to precisely control the amount of prefetched video.



COMPUTER NETWORKS

5.2.4 CDN

5.2.4.1 Motivation for CDN

- The streaming video service can be provided as follows:
 - 1) Build a single massive data-center.
 - 2) Store all videos in the data-center and
 - 3) Stream the videos directly from the data-center to clients worldwide.
- Three major problems with the above approach:
 - 1) More Delay**
 - If links provides a throughput lesser than consumption-rate, the end-to-end throughput will also be below the consumption-rate.
 - This results in freezing delays for the user.
 - 2) Network Bandwidth is wasted**
 - A popular video may be sent many times over the same links.
 - 3) Single Point of Failure:**
 - If the data-center goes down, it cannot distribute any video streams.

Problem Solution: Use CDN (Content Distribution Network).

5.2.4.2 CDN Types

- A CDN
 - manages servers in multiple geographically distributed locations
 - stores copies of the videos in its servers, and
 - attempts to direct each user-request to a CDN that provides the best user experience.
- The CDN may be a private CDN or a third-party CDN.
 - 1) Private CDN**
 - A private CDN is owned by the content provider itself.
 - For example:
 - Google's CDN distributes YouTube videos
 - 2) Third Party CDN**
 - A third-party CDN distributes content on behalf of multiple content providers CDNs.
 - Two approaches for server placement:
 - i) Enter Deep**
 - ✗ The first approach is to enter deep into the access networks of ISPs.
 - ✗ Server-clusters are deployed in access networks of ISPs all over the world.
 - ✗ The goal is to get close to end users.
 - ✗ This improves delay/throughput by decreasing no. of links b/w end user & CDN cluster
 - ii) Bring Home**
 - ✗ The second approach is to bring the ISPs home.
 - ✗ Large clusters are built at a smaller number of key locations.
 - ✗ These clusters are connected using a private high-speed network.
 - ✗ Typically, clusters are placed at a location that is near the PoPs of many tier-1 ISPs.
 - For example: within a few miles of both Airtel and BSNL PoPs in a major city.
 - ✗ Advantage:
 - Lower maintenance and management overhead.
 - ✗ Disadvantage:
 - Higher delay and lower throughput to end users.



COMPUTER NETWORKS

5.2.4.3 CDN Operation

- When a browser wants to retrieve a specific video, the CDN intercepts the request.
- Then, the CDN
 - 1) determines a suitable server-cluster for the client and
 - 2) redirects the client's request to the desired server.
- Most CDNs take advantage of DNS to intercept and redirect requests.
- CDN operation is illustrated in Figure 5.2.

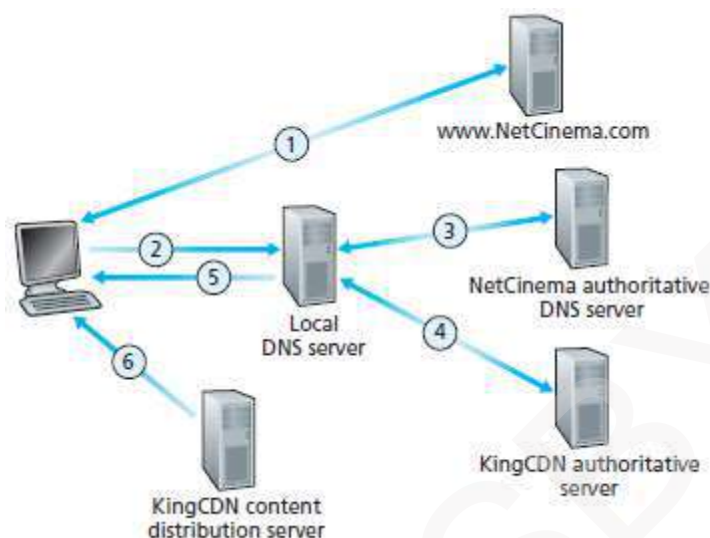


Figure 5.2: DNS redirects a user's request to a CDN server

- Suppose a content provider "NetCinema" employs the CDN company "KingCDN" to distribute videos.
- Let URL = <http://video.netcinema.com/6Y7B23V>
- Six events occur as shown in Figure 5.2:
 - 1) The user visits the Web page at NetCinema.
 - 2) The user clicks on the following link:
<http://video.netcinema.com/6Y7B23V>,
 - Then, the user's host sends a DNS query for "video.netcinema.com".
 - 3) The user's local-DNS-server (LDNS) forwards the DNS-query to an authoritative-DNS-server "NetCinema".
 - The server "NetCinema" returns to the LDNS a hostname in the KingCDN's domain.
 - For example: "a1105.kingcdn.com".
 - 4) The user's LDNS then sends a second query, now for "a1105.kingcdn.com".
 - Eventually, KingCDN's DNS system returns the IP addresses of a "KingCDN" server to LDNS.
 - 5) The LDNS forwards the IP address of the "KingCDN" server to the user's host.
 - 6) Finally, the client
 - establishes a TCP connection with the server
 - issues an HTTP GET request for the video.



COMPUTER NETWORKS

5.2.4.4 Cluster Selection Strategies

- Cluster-selection strategy is used for dynamically directing clients to a server-cluster within the CDN.
- The CDN learns the IP address of the client's LDNS server via the client's DNS lookup.
- After learning this IP address, the CDN selects an appropriate cluster based on this IP address.
- Three approaches for cluster-selection:

1) Geographically Closest

- The client is assigned to the cluster that is geographically closest.
- Using geo-location databases, each LDNS IP address is mapped to a geographic location.
- When a DNS request is received from LDNS, the CDN chooses geographically closest-cluster.
- Advantage:
 - This solution can work reasonably well for a large fraction of the clients.
- Disadvantages: The solution may perform poorly. This is because
 - 1) Geographically closest-cluster may not be the closest-cluster along the path.
 - 2) The LDNs location may be far from the client's location.

2) Based on Current Traffic Conditions

- The best cluster can be determined for a client based on the current traffic-conditions.
- CDNs perform real-time measurements of delay/loss performance b/w their clusters & clients.
- In a CDN, each cluster periodically sends probes to all of the LDNSs around the world.
- Disadvantage:
 - Many LDNSs are configured to not respond to the probes.

3) IP Anycast

- The idea behind IP anycast:
 - In Internet, the routers must route the packets to the closest-cluster, as determined by BGP.
- IP anycast is illustrated in Figure 5.3.

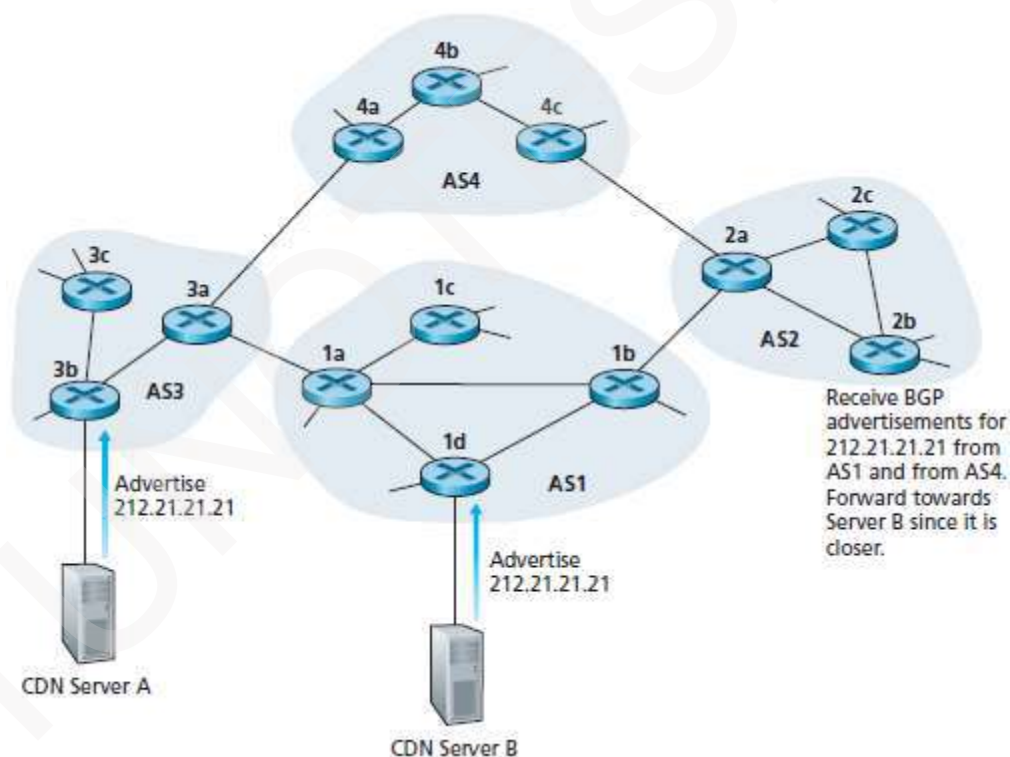


Figure 5.3: Using IP anycast to route clients to closest CDN cluster

- Here is how it works:
 - 1) During the IP-anycast configuration stage, the CDN company
 - assigns the same IP address to each clusters and
 - uses BGP to advertise the IP address from different cluster locations.
 - 2) The BGP router treats the multiple route advertisements as different paths to the same physical location.
 - 3) Then, the BGP router picks the "best" route to the IP address.



COMPUTER NETWORKS

5.3 Voice-over-IP

- Real-time voice over the Internet is often referred to as Internet telephony.
- It is also commonly called Voice-over-IP (VoIP).

5.3.1 Limitations of the Best-Effort IP Service

- The Internet's network-layer protocol IP provides best-effort service.
- The IP makes best effort to move each datagram from source to destination.
- But IP does not guarantee deliver of the packet to the destination.
- Three main challenges to the design of real-time applications:
 - 1) Packet-loss
 - 2) Packet delay and
 - 3) Packet jitter.

5.3.1.1 Packet Loss

- By default, most existing VoIP applications run over UDP.
- The UDP segment is encapsulated in an IP datagram.
- The datagram passes through router buffers in the path from sender to receiver
- Problem:
 - There is possibility that one or more buffers are full.
 - In this case, the arriving IP datagram may be discarded.
- Possible solution:
 - Loss can be eliminated by sending the packets over TCP rather than over UDP.
 - However, retransmissions are unacceptable for real-time applications \.' they increase delay.
 - Packet-loss results in a reduction of sender's transmission-rate, leading to buffer starvation.

5.3.1.2 End-to-End Delay

- End-to-end delay is the sum of following delays:
 - 1) Transmission, processing, and queuing delays in routers.
 - 2) Propagation delays in links and
 - 3) Processing delays in end-systems.
- For VoIP application,
 - delays smaller than 150 msec are not perceived by a human listener.
 - delays between 150 and 400 msec can be acceptable but are not ideal and
 - delays exceeding 400 msec can seriously hinder the interactivity in voice conversations.
- Typically, the receiving-side will discard any packets that are delayed more than a certain threshold.
- For example: more than 400 msec.

5.3.1.3 Packet Jitter

- Jitter refers to varying queuing delays that a packet experiences in the network's routers.
- If the receiver
 - ignores the presence of jitter and
 - plays out audio-chunks,then the resulting audio-quality can easily become unintelligible.
- Jitter can often be removed by using sequence numbers, timestamps, and a playout delay



COMPUTER NETWORKS

5.3.2 Removing Jitter at the Receiver for Audio

- For VoIP application, receiver must provide periodic playout of voice-chunks in presence of random jitter
- This is typically done by combining the following 2 mechanisms:
 - 1) Prepending each Chunk with a Timestamp**
 - The sender attaches each chunk with the time at which the chunk was generated.
 - 2) Delaying Playout of Chunks at the Receiver**
 - The playout delay of the received chunks must be long.
 - So, the most of the packets are received before their scheduled playout times.
 - This playout delay can either be
 - fixed throughout the duration of the session or
 - vary adaptively during the session-lifetime.



COMPUTER NETWORKS

5.3.3 Recovering from Packet Loss

- Loss recovery schemes attempt to preserve acceptable audio-quality in the presence of packet-loss.
- Here, packet-loss is defined in a 2 broad sense:
 - i) A packet is lost if the packet never arrives at the receiver or
 - ii) A packet is lost if the packet arrives after its scheduled playout time.
- VoIP applications often use loss anticipation schemes.
- Here, we consider 2 types of loss anticipation schemes:
 - 1) Forward error correction (FEC) and
 - 2) Interleaving.
- We also consider an error concealment scheme.

5.3.3.1 FEC

- The basic idea of FEC: Redundant information is added to the original packet stream.
- The redundant information can be used to reconstruct approximations of some of the lost-packets.
- Two FEC mechanisms:

1) Block Coding

- A redundant encoded chunk is sent after every n chunks.
- The redundant chunk is obtained by exclusive OR-ing the n original chunks.
- If any one packet in a group is lost, the receiver can fully reconstruct the lost-packet.
- Disadvantages:
 - 1) If 2 or more packets in a group are lost, receiver cannot reconstruct the lost-packets.
 - 2) Increases the playout delay. This is because
 - receiver must wait to receive entire group of packets before it can begin playout.

2) Lower Resolution Redundant Information

- A lower-resolution audio-stream is sent as the redundant information.
- For example: The sender creates
 - nominal audio-stream and
 - corresponding low-resolution, low-bit-rate audio-stream.
- The low-bit-rate stream is referred to as the redundant-stream.
- As shown in Figure 5.4, the sender constructs the n th packet by
 - taking the n th chunk from the nominal stream and
 - appending the n th chunk to the $(n-1)$ st chunk from the redundant-stream
- Advantage:
 - Whenever there is packet-loss, receiver can hide the loss by playing out low-bit-rate chunk.

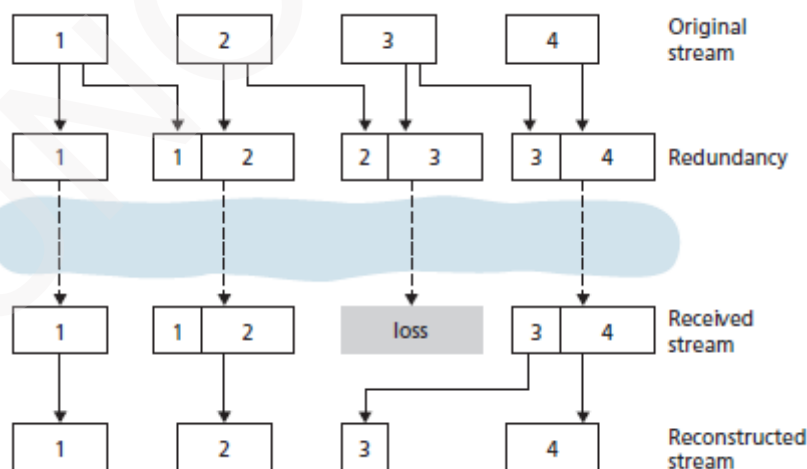


Figure 5.4: Piggybacking lower-quality redundant information



COMPUTER NETWORKS

5.3.3.2 Interleaving

- A VoIP application can send interleaved audio.
- The sender resequences units of audio-data before transmission.
- Thus, originally adjacent units are separated by a certain distance in the transmitted-stream.
- Interleaving can mitigate the effect of packet-losses.
- Interleaving is illustrated in Figure 5.5.

- For example:

If units are 5 msec in length and chunks are 20 msec (that is, four units per chunk), then

→ the first chunk contains the units 1, 5, 9, and 13

→ the second chunk contains the units 2, 6, 10 & 14 and so on.

- Advantages:

- 1) Improves the perceived quality of an audio-stream.
- 2) Low overhead.
- 3) Does not increase the bandwidth requirements of a stream.

- Disadvantage:

- 1) Increases latency. This limits use for VoIP applications.

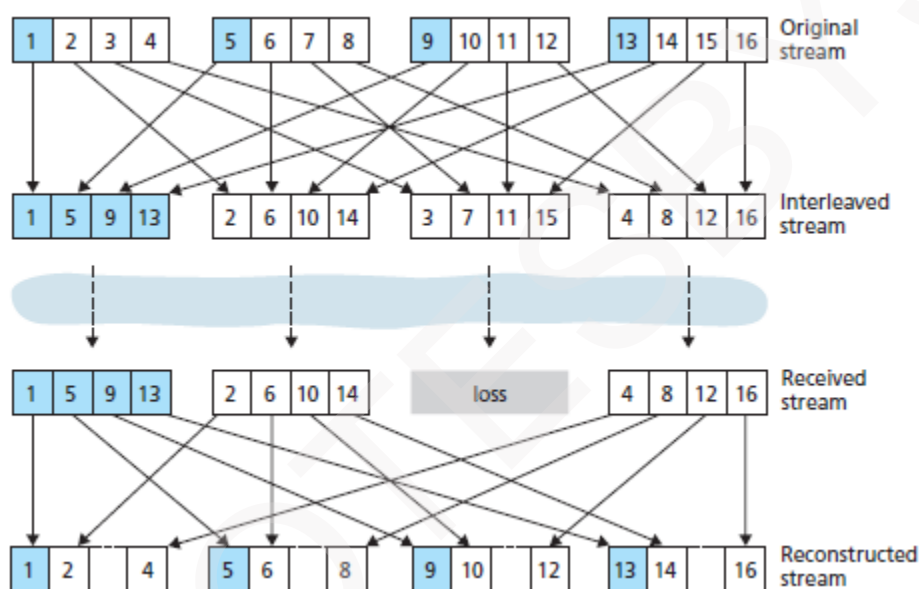


Figure 5.5: Sending interleaved audio

5.3.3.3 Error Concealment

- This scheme attempts to produce a replacement for a lost-packet that is similar to the original.
- This is possible since audio-signals exhibit large amounts of short-term self-similarity.
- Two receiver-based recovery techniques:

1) Packet Repetition

- This replaces lost-packets with copies of packets that arrived immediately before the loss.

- Advantage:

- 1) Low computational complexity.

2) Interpolation

- This uses audio before and after the loss to interpolate a suitable packet to cover the loss.

- Advantage:

- 1) Performs better than packet repetition.



COMPUTER NETWORKS

5.4 Protocols for Real-Time Conversational Applications

- Real-time applications are very popular. For ex: VoIP and video conferencing.
- Two standards bodies are working for real-time applications: 1) IETF and 2) ITU
- Both standards (IETF & ITU) are enjoying widespread implementation in industry products.

5.4.1 RTP

- RTP can be used for transporting common formats such as
 - MP3 for sound and
 - MPEG for video
- It can also be used for transporting proprietary sound and video formats.
- Today, RTP enjoys widespread implementation in many products and research prototypes.
- It is also complementary to other important real-time interactive protocols, such as SIP.

5.4.1.1 RTP Basics

- RTP runs on top of UDP.
- The RTP packet is composed of i) RTP header & ii) audio chunk
- The header includes
 - i) Type of audio encoding
 - ii) Sequence number and
 - iii) Timestamp.
- The application appends each chunk of the audio-data with an RTP header.
- Here is how it works:
 - 1) At sender-side:
 - i) A media chunk is encapsulated within an RTP packet.
 - ii) Then, the packet is encapsulated within a UDP segment.
 - iii) Finally, the UDP segment is handed over to IP.
 - 2) At receiving-side:
 - i) The RTP packet is extracted from the UDP segment.
 - ii) Then, the media chunk is extracted from the RTP packet.
 - iii) Finally, the media chunk is passed to the media-player for decoding and rendering
- If an application uses RTP then the application easily interoperates with other multimedia applications
- For example:
 - If 2 different companies use RTP in their VoIP product, then users will be able to communicate.
- What RTP does not provide?
 - i) It does not provide any mechanism to ensure timely delivery of data.
 - ii) It does not provide quality-of-service (QoS) guarantees.
 - iii) It does not guarantee delivery of packets.
 - iv) It does not prevent out-of-order delivery of packets.
- RTP encapsulation is seen only at the end systems.
- Routers do not distinguish between
 - i) IP datagrams that carry RTP packets and
 - ii) IP datagrams that don't carry RTP packets.
- RTP allows each source to be assigned its own independent RTP stream of packets.
- For example:
 - 1) For a video conference between two participants, four RTP streams will be opened
 - i) Two streams for transmitting the audio (one in each direction) and
 - ii) Two streams for transmitting the video (again, one in each direction).
 - 2) Encoding technique MPEG bundles audio & video into a single stream.
 - In this case, only one RTP stream is generated in each direction.
- RTP packets can also be sent over one-to-many and many-to-many multicast trees.



COMPUTER NETWORKS

5.4.1.2 RTP Packet Header Fields

- Four header fields of RTP Packet (Figure 5.6):
 - 1) Payload type
 - 2) Sequence number
 - 3) Timestamp and
 - 4) Source identifier.
- Header fields are illustrated in Figure 5.6.

Payload type	Sequence number	Timestamp	Synchronization source identifier	Miscellaneous fields
--------------	-----------------	-----------	-----------------------------------	----------------------

Figure 5.6: RTP header fields

Table 5.1: Audio payload types supported by RTP

Payload-Type Number	Audio Format	Sampling Rate	Rate
0	PCM μ -law	8 kHz	64 kbps
1	1016	8 kHz	4.8 kbps
3	GSM	8 kHz	13 kbps
7	LPC	8 kHz	2.4 kbps
9	G.722	16 kHz	48–64 kbps
14	MPEG Audio	90 kHz	—
15	G.728	8 kHz	16 kbps

Table 5.2: Some video payload types supported by RTP

Payload-Type Number	Video Format
26	Motion JPEG
31	H.261
32	MPEG 1 video
33	MPEG 2 video

1) Payload Type

- For an audio-stream, this field is used to indicate type of audio encoding that is being used.
 - For example: PCM, delta modulation.
 - Table 5.1 lists some of the audio payload types currently supported by RTP.
- For a video stream, this field is used to indicate the type of video encoding.
 - For example: motion JPEG, MPEG.
 - Table 5.2 lists some of the video payload types currently supported by RTP.

2) Sequence Number

- This field increments by one for each RTP packet sent.
- This field may be used by the receiver to detect packet loss and to restore packet sequence.

3) Timestamp

- This field reflects the sampling instant of the first byte in the RTP data packet.
- The receiver can use timestamps
 - to remove packet jitter in the network and
 - to provide synchronous playout at the receiver.
- The timestamp is derived from a sampling clock at the sender.

4) Source Identifier (SRC)

- This field identifies the source of the RTP stream.
- Typically, each stream in an RTP session has a distinct SRC.



COMPUTER NETWORKS

5.4.2 SIP

- SIP (Session Initiation Protocol) is an open and lightweight protocol.
- Main functions of SIP:
 - 1) It provides mechanisms for establishing calls b/w a caller and a callee over an IP network.
 - 2) It allows the caller to notify the callee that it wants to start a call.
 - 3) It allows the participants to agree on media encodings.
 - 4) It also allows participants to end calls.
 - 5) It provides mechanisms for the caller to determine the current IP address of the callee.
 - 6) It provides mechanisms for call management, such as
 - adding new media streams during the call
 - changing the encoding during the call
 - inviting new participants during the call,
 - call transfer and
 - call holding.



COMPUTER NETWORKS

5.4.2.1 Setting up a Call to a Known IP Address

- SIP call-establishment process is illustrated in Figure 5.7.

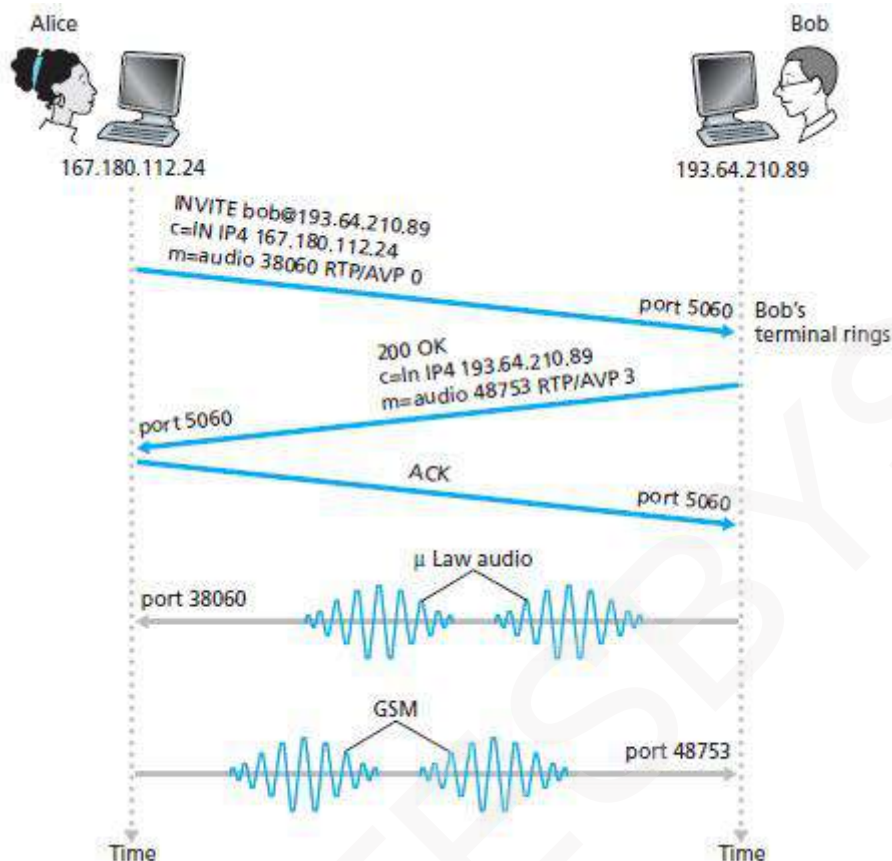


Figure 5.7: SIP call establishment when Alice knows Bob's IP address

- Consider an example: Alice wants to call Bob.
- Alice's & Bob's PCs are both equipped with SIP-based software for making and receiving phone calls.
- The following events occur:
 - 1) An SIP session begins when Alice sends Bob an INVITE message.
 - This INVITE message is sent over UDP to the well-known port 5060 for SIP.
 - The INVITE message includes
 - i) An identifier for Bob (bob@193.64.210.89)
 - ii) An indication of Alice's current IP address
 - iii) An indication that Alice desires to receive audio, which is encoded in format AVP 0.
 - 2) Then, Bob sends an SIP response message (which resembles an HTTP response message).
 - The response message is sent over UDP to the well-known port 5060 for SIP.
 - The response message includes
 - i) 200 OK
 - ii) An indication of Bob's current IP address
 - iii) An indication that Bob desires to receive audio, which is encoded in format AVP 3.
 - 3) Then, Alice sends Bob an SIP acknowledgment message.
 - 4) Finally, Bob and Alice can talk.
- Three key characteristics of SIP:
 - 1) SIP is an out-of-band protocol
 - The SIP message & the media-data use different sockets for sending and receiving.
 - 2) The SIP messages are ASCII-readable and resemble HTTP messages.
 - 3) SIP requires all messages to be acknowledged, so it can run over UDP or TCP.



COMPUTER NETWORKS

5.4.2.2 SIP Messages

- Suppose that Alice wants to initiate a VoIP call to Bob.
- Then, her message will look something like this:

```
L1) INVITE sip:bob@domain.com SIP/2.0
L2) Via: SIP/2.0/UDP 167.180.112.24
L3) From: sip:alice@hereway.com
L4) To: sip:bob@domain.com
L5) Call-ID: a2e3a@pigeon.hereway.com
L6) Content-Type: application/sdp
L7) Content-Length: 885
L8) c=IN IP4 167.180.112.24
L9) m=audio 38060 RTP/AVP 0
```

- Line by line explanation is as follows:
 - L1) The INVITE line includes the SIP version.
 - L2) Via header indicates the IP address of the SIP device.
 - L3) Similar to an e-mail message, the SIP message includes a From header line.
 - L4) Similar to an e-mail message, the SIP message includes a To header line.
 - L5) Call-ID uniquely identifies the call.
 - L6) Content-Type defines the format used to describe the message-content.
 - L7) Content-Length provides the length in bytes of the message-content.
 - L8) A carriage return followed by line feed.
 - L9) The message contains the content.



COMPUTER NETWORKS

5.4.2.3 Name Translation & User Location

- IP addresses are often dynamically assigned with DHCP.
- Suppose that Alice knows only Bob's e-mail address, bob@domain.com
- In this case, Alice needs to obtain the IP address of the device associated with the bob@domain.com.
- How to find IP address?
 - 1) Alice creates & sends an INVITE message to an SIP proxy.
 - 2) The proxy responds with an SIP reply.
 - The reply includes the IP address of the device associated with the bob@domain.com.
- How can the proxy server determine the current IP address for bob@domain.com?

Answer: Every SIP user has an associated registrar.

- First, a user launches an SIP application on a device.
- Then, the application sends an SIP register message to the registrar.
- Finally, IP address of the device will be registered in the registrar.
- The user's registrar keeps track of his current IP address.
- When the user switches to a new device, IP address of new device will be registered in the registrar.
- The registrar is analogous to a DNS authoritative name server:
 - 1) The DNS server translates fixed host names to fixed IP addresses;
 - 2) The SIP registrar translates human identifiers (ex: bob@domain.com) to dynamic IP address.

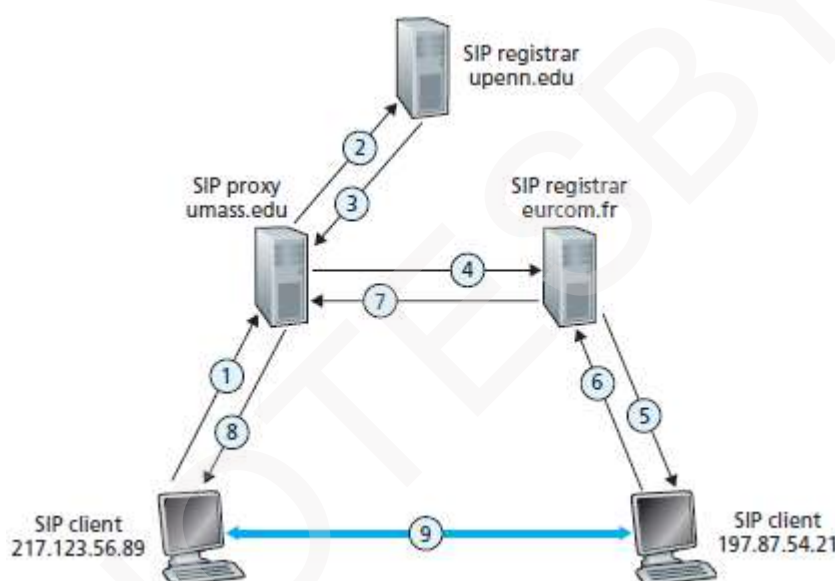


Figure 5.8: Session initiation, involving SIP proxies and registrars

- Session initiation process is illustrated in Figure 5.8.
- jim@umass.edu (217.123.56.89) wants to initiate VoIP session with keith@upenn.edu (197.87.54.21)
- The following steps are taken:
 - 1) Jim sends an INVITE message to the umass SIP proxy.
 - 2) The proxy
 - performs DNS lookup on the SIP registrar upenn.edu and
 - forwards then the message to the registrar server upenn.
 - 3) keith@upenn.edu is not registered at the upenn registrar.
 - Therefore, the upenn registrar sends a redirect response to umass proxy.
 - 4) The umass proxy sends an INVITE message to the eurecom SIP registrar.
 - 5) The eurecom registrar
 - knows the IP address of keith@eurecom.fr and
 - forwards INVITE message to the host 197.87.54.21 which is running Keith's SIP client.
 - 6-8) An SIP response is sent back through registrars/proxies to SIP client on 217.123.56.89.
 - 9) Media is sent directly between the two clients.



COMPUTER NETWORKS

5.5 Network Support for Multimedia

- Table 5.3 summarizes 3 approaches to provide network-level support for multimedia applications.

1) Making the Best of Best-effort Service

- The application-level mechanisms can be successfully used in a n/w where packet-loss rarely occur.
- When demand increases are forecasted, ISPs can deploy additional bandwidth & switching capacity.
- This ensures satisfactory delay and packet-loss performance.

2) Differentiated Service

- In this, one traffic-type can be given priority over another one when both are queued at a router
- For ex:

Packets belonging to a real-time application can be given priority over non-real-time application.

3) Per-connection QoS Guarantees

- In this, each instance of an application explicitly reserves end-to-end bandwidth.
- Thus, end-to-end performance is guaranteed.
 - i) A hard guarantee means the application will receive its requested QoS with certainty.
 - ii) A soft guarantee means the application will receive its requested QoS with high probability.

Table 5.3: Three network-level approaches to supporting multimedia applications

Approach	Granularity	Guarantee	Mechanisms	Complexity	Deployment to date
Making the best of best-effort service.	all traffic treated equally	none, or soft	application-layer support, CDNs, overlays, network-level resource provisioning	minimal	everywhere
Differentiated service	different classes of traffic treated differently	none, or soft	packet marking, policing, scheduling	medium	some
Per-connection Quality-of-Service (QoS) Guarantees	each source-destination flows treated differently	soft or hard, once flow is admitted	packet marking, policing, scheduling; call admission and signaling	light	little



COMPUTER NETWORKS

5.5.1 Dimensioning Best Effort Networks

- Bandwidth-provisioning is defined as
 - “The bandwidth capacity required at network links to achieve a given performance.”
- Network dimensioning is defined as
 - “The design of a network topology to achieve a given performance.”
- Three issues must be addressed to achieve a given performance:
 - 1) Models of traffic demand between network end points.
- Models have to be specified at both the call-level and at the packet-level.
 - i) Example for call level: Users arriving to the network and starting up end-to-end applications.
 - ii) Example for packet level: Packets being generated by ongoing applications.
- 2) Well-defined performance requirements.
- For example
 - Consider delay-sensitive traffic, such as a conversational multimedia application.
 - Here, a performance requirement can be:
 - probability the end-to-end delay of the packet is greater than a maximum tolerable delay
- 3) Models to predict end-to-end performance for a given workload model.
- Techniques to find a least cost bandwidth allocation that results in all user requirements being met.



COMPUTER NETWORKS

5.5.2 Providing Multiple Classes of Service

- The traffic can be divided into different classes.
- Different priority can be provided to the different traffic-classes.
- For example:
 - ISP provides a higher priority to delay-sensitive VoIP traffic than to elastic traffic email/HTTP.

5.5.2.1 Motivating Scenarios

- Figure 5.9 shows a simple network scenario.
- Here, two flows
 - originate on Hosts H1 & H2 on one LAN and
 - are destined for Hosts H3 and H4 on another LAN.
- The routers on the two LANs are connected by a 1.5 Mbps link.

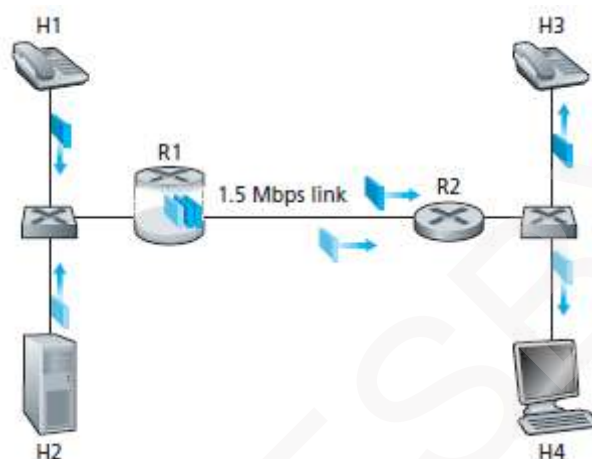


Figure 5.9: Competing audio and HTTP applications

- Consider the best-effort Internet.
- At output-queue of R1, the audio and HTTP packets are mixed & transmitted in a FIFO order.
- Problem: HTTP packet from Web-server fills up the queue. This causes delay/loss of audio packets.
- Solution:
 - Use a priority scheduling discipline.
 - Here, an audio packet will always be transmitted before HTTP packets.

Insight 1:

- Packet marking allows a router to distinguish among packets belonging to different classes of traffic.
- If audio application starts sending packets at 1.5 Mbps or higher, then the HTTP packets will starve.

Insight 2:

- It is desirable to provide a degree of traffic isolation among classes.
- Thus, one class is not adversely affected by another class of traffic that misbehaves.



COMPUTER NETWORKS

5.5.2.1.1 Approaches for Providing Isolation among Traffic Classes

- Two approaches:
 - 1) Traffic policing can be performed.
 - 2) Packet-scheduling mechanism explicitly allocates a fixed amount of bandwidth to each class.

5.5.2.1.1.1 Using Traffic Policing

- Traffic policing can be performed. This scenario is shown in Figure 5.10.
- If a traffic class meets certain criteria, then a policing mechanism ensures that these criteria are observed. (For example: the audio flow should not exceed a peak-rate of 1 Mbps).
- If the policed application misbehaves, the policing mechanism will take some action. (For example: drop or delay packets that are in violation of the criteria)
- The leaky bucket mechanism is the most widely used policing mechanism.
- This scenario is shown in Figure 5.10.

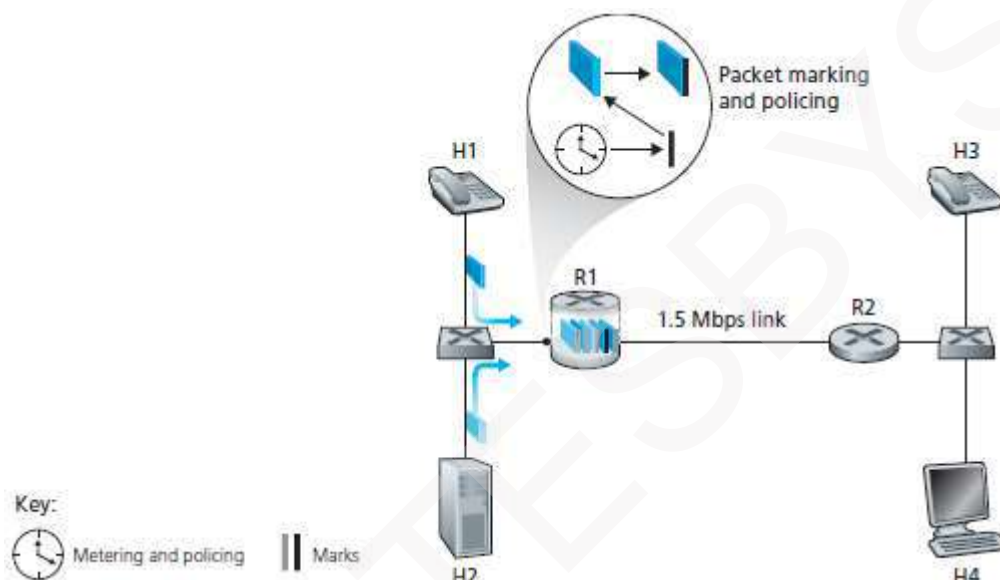


Figure 5.10: Policing (and marking) the audio and HTTP traffic classes

5.5.2.1.1.2 Using Packet Scheduling

- Packet-scheduling mechanism explicitly allocates a fixed amount of bandwidth to each class.
- For example:
 - the audio class will be allocated 1 Mbps at R1
 - the HTTP class will be allocated 0.5 Mbps.
- This scenario is shown in Figure 5.11.

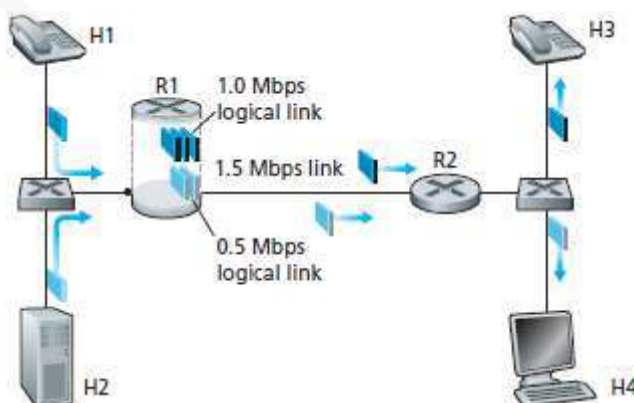


Figure 5.11: Logical isolation of audio and HTTP traffic classes

Insight 3:

- While providing isolation among classes, it is desirable to use resources as efficiently as possible.

"Life is a learning process; learning is a lifelong process. You cant separate them." -Vivienne Forrester



COMPUTER NETWORKS

5.5.2.2 Scheduling Mechanisms

- Queue-scheduling refers to
"The manner in which queued packets are selected for transmission on the link."

5.5.2.2.1 FIFO

- FIFO (First-In-First-Out) is illustrated in Figure 5.12 & Figure 5.13.
 - Packets are transmitted in order of their arrival at the queue.
 - Packets arriving at the queue wait for transmission if the link is currently busy.
 - Packets are discarded when they arrive at a full buffer.
 - When a packet is completely transmitted over outgoing-link, the packet is removed from the queue.
 - Disadvantages:
 - 1) This is not possible to provide different information flows with different QoS.
 - 2) Hogging occurs when a user
 - sends packets at a high rate and
 - fills the buffers in the system
- Thus, depriving other users of access to the buffer.

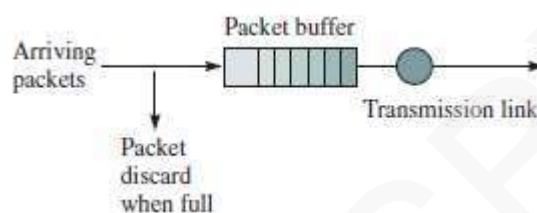


Figure 5.12: FIFO queuing

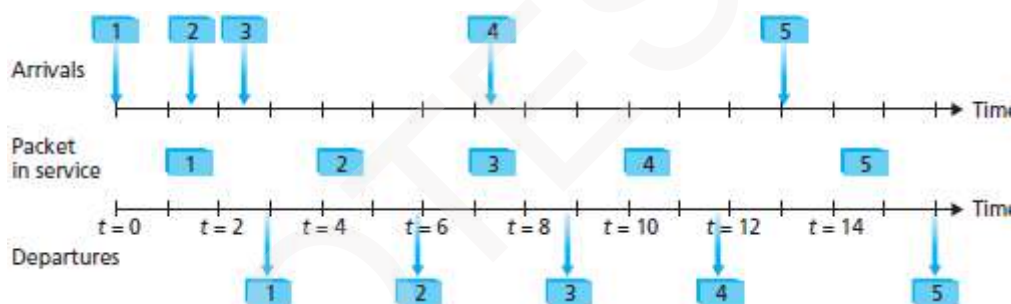


Figure 5.13: Operation of the FIFO queue



COMPUTER NETWORKS

5.5.2.2.2 Priority Queuing

- PQ (Priority Queuing) is illustrated in Figure 5.14 & Figure 5.15.
- Number of priority classes is defined.
- A separate buffer is maintained for each priority class.
- Each time the transmission link becomes available, the next packet for transmission is selected from the highest priority queue.
- Typically, the choice among packets in the same priority-class is done in a FIFO manner.
- In a non-preemptive PQ, the transmission of a packet is not interrupted once it has begun.
- Disadvantages:
 - 1) This does not discriminate among users of the same priority.
 - 2) This does not allow for providing some degree of guaranteed access to transmission bandwidth to the lower priority classes.
 - 3) Hogging occurs.

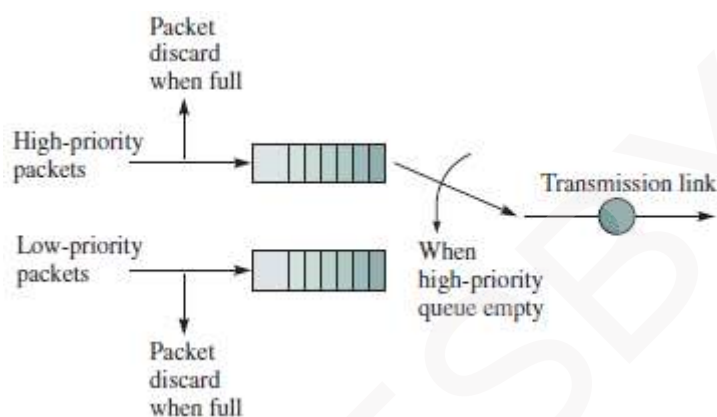


Figure 5.14 : Priority Queueing

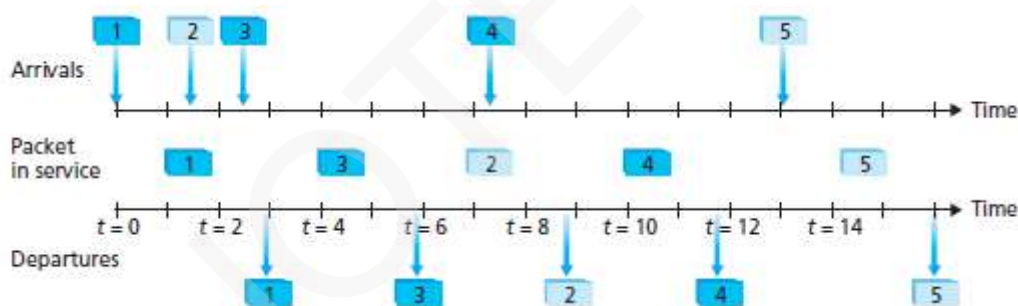


Figure 5.15: Operation of the priority queue

COMPUTER NETWORKS

5.5.2.2.3 RRQ

- RRQ (Round Robin Queuing) is illustrated in Figure 5.16 & Figure 5.17.
- The transmission bandwidth is divided equally among the buffers.
- Each user flow has its own logical buffer.
- Round-robin scheduling is used to service each non-empty buffer one bit at a time.
- In the simplest form, a class 1 packet is transmitted, followed by a class 2 packet, followed by a class 1 packet, followed by a class 2 packet, and so on.
- RRQ is a work-conserving queuing discipline.
- Thus, RRQ will immediately move on to the next class when it finds an empty queue.
- Disadvantage: Extensive processing at the destination.

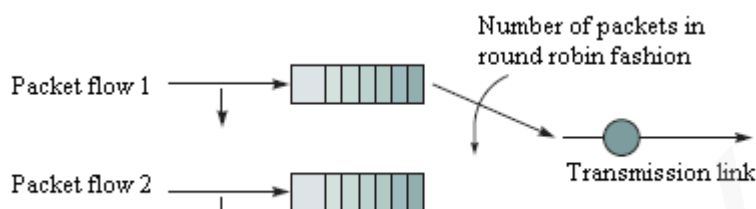


Figure 5.16: Round-robin queuing

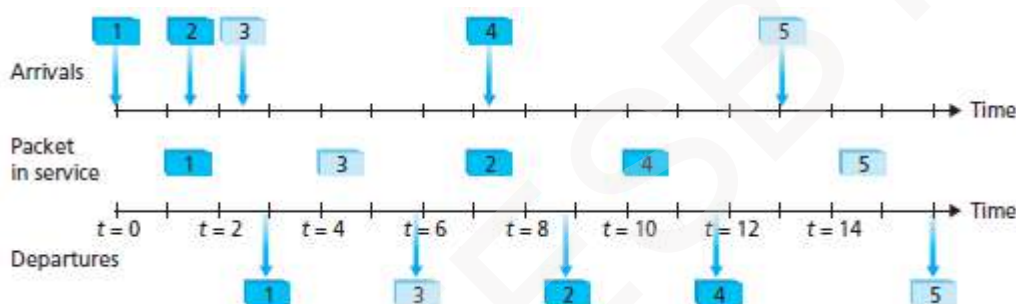


Figure 5.17: Operation of the two-class round robin queue

5.5.2.2.4 WFQ

- WFQ (Weighted Fair Queuing) is illustrated in Figure 5.18.
- Each user flow has its own buffer, but each user flow also has a weight that determines its relative share of the bandwidth.
- If buffer 1 has weight 1 and buffer 2 has weight 3, then buffer 1 will receive 1/4 of the bandwidth and buffer 2 will receive 3/4 of the bandwidth.
- In each round, each non-empty buffer would transmit a number of packets proportional to its weight.
- WFQ systems are means for providing QoS guarantees.
- WFQ is also a work-conserving queuing discipline.
- Thus, WFQ will immediately move on to the next class when it finds an empty queue.

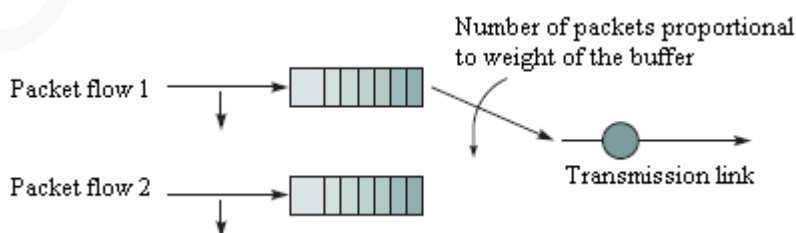


Figure 5.18: Weighted fair queuing



COMPUTER NETWORKS

5.5.2.3 Policing: The Leaky Bucket

- Policing is an important QoS mechanism
- Policing means the regulation of the rate at which a flow is allowed to inject packets into the network.
- Three important policing criteria:
 - 1) **Average Rate**
 - This constraint limits amount of traffic that can be sent into n/w over a long period of time.
 - 2) **Peak Rate**
 - This constraint limits maximum no. of packets that can be sent over a short period of time
 - 3) **Burst Size**
 - This constraint limits the maximum no. of packets that can be sent into n/w over a very short period of time.

5.5.2.3.1 Leaky Bucket Operation

- Policing-device can be implemented based on the concept of a leaky bucket.
- Tokens are generated periodically at a constant rate.
- Tokens are stored in a bucket.
- A packet from the buffer can be taken out only if a token in the bucket can be drawn.
- If the bucket is full of tokens, additional tokens are discarded.
- If the bucket is empty, arriving packets have to wait in the buffer until a sufficient no. of tokens is generated.

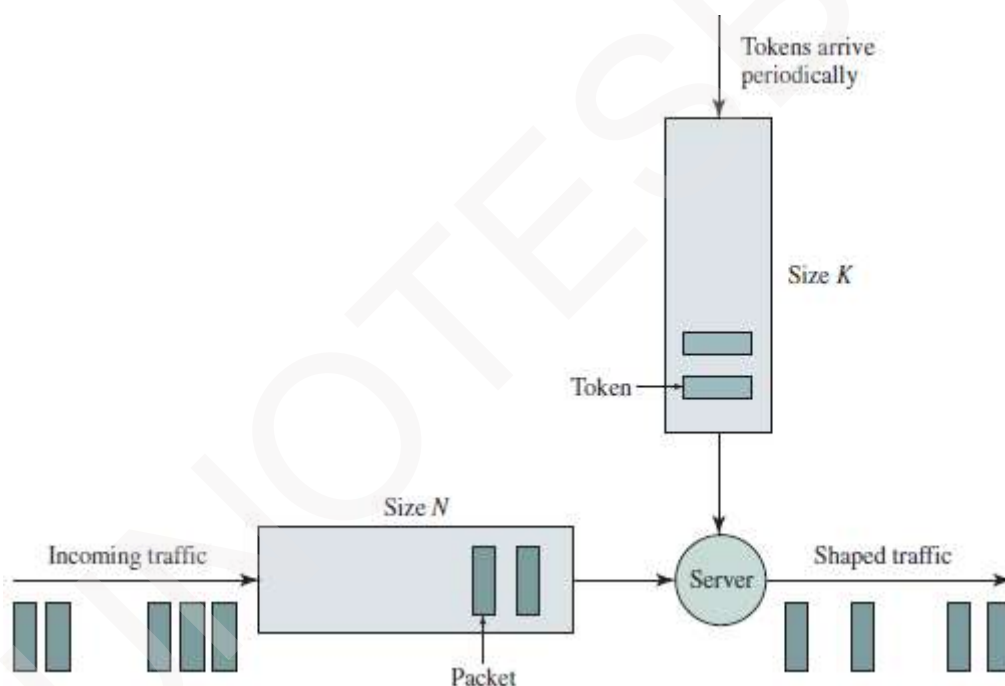


Figure 5.19: The leaky bucket policer



COMPUTER NETWORKS

5.5.3 DiffServ

- This provides QoS support to a broad class of applications.
- This provides service differentiation.
- Differentiation is defined as
“The ability to handle different classes of traffic in different ways within the Internet”.

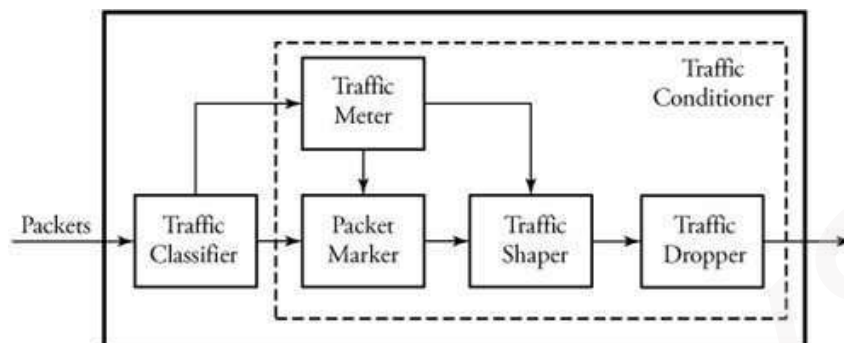


Figure 5.21: Overview of DiffServ operation

- The Diffserv architecture consists of 2 functional elements:

1) Packet Classification & Traffic Conditioning

- The traffic-classifier routes packets to specific outputs, based on the values of one or more header-fields.
- The traffic-profile contains a limit on the peak-rate of the flow.
- The traffic-conditioner detects and responds if any packet has violated the negotiated traffic-profile.
- The traffic-conditioner has 4 major components:

i) Meter

- The meter measures the traffic to make sure that packets do not exceed their traffic profiles

ii) Marker

- The marker marks or unmarks packets in order to keep track of their situations in the Diffserv node.

iii) Shaper

- The shaper delays any packet that is not complaint with the traffic-profile

iv) Dropper

- The dropper discards any packet that violates its traffic-profile

2) Core Function: Forwarding

- The per-hop behavior (PHB) is performed by Diffserv-capable routers.
- A router forwards marked-packet onto its next hop according to the PHB (per-hop behavior).
- PHB influences how network-resources are shared among the competing classes of traffic.
- Two types of PHB are: i) expedited forwarding and ii) assured forwarding.

i) Expedited Forwarding (EF) PHB

- This specifies that the departure rate of a class of traffic from a router must equal or exceed a configured rate.

ii) Assured Forwarding (AF) PHB

- This divides traffic into 3 classes: good, average and poor.
- Here, each class is guaranteed to be provided with some minimum amount of bandwidth and buffering.

- PHB is defined as

“A description of the externally observable forwarding behavior of a Diffserv node applied to a particular Diffserv behavior aggregate”

- From above definition, we have 3 considerations:

- i) A PHB can result in different classes of traffic receiving different performance.
- ii) A PHB does not dictate any particular mechanism for differentiating performance (behavior) among classes.
- iii) Differences in performance must be observable and hence measurable.



COMPUTER NETWORKS

5.5.4 Per-Connection QoS Guarantees: Resource Reservation & Call Admission

- Consider two 1 Mbps audio applications transmitting their packets over 1.5 Mbps link (Figure 5.22).
- The combined data-rate of the two flows (2 Mbps) exceeds the link capacity.
- There is lesser bandwidth to accommodate the needs of both applications at the same time.
- Problem: What should the network do?

Answer:

- One of the applications should be allowed to use the full 1 Mbps.
- While the other application flows should be blocked.
- For example:
 - In telephone n/w, if the required resources cannot be allocated to the call, the call is blocked.

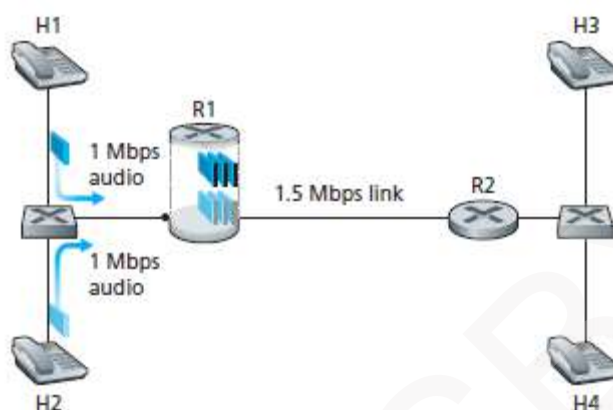


Figure 5.22: Two competing audio applications overloading R1-to-R2 link

- Admission process is defined as
The process of declaring flow's QoS requirement, & then deciding either to accept or block the flow.

Insight 4:

- If sufficient resources are not available, and QoS is to be guaranteed, a call admission process is needed to decide whether to accept or block the flow.

COMPUTER NETWORKS

5.5.4.1 Mechanisms for Guaranteed QoS

- Three mechanisms to provide guaranteed QoS:

1) Resource Reservation

- The resources are explicitly allocated to the call to meet the desired QoS.
- After reservation, the call has on-demand access to the resources throughout its duration.

2) Call Admission

- The network must have a mechanism for calls to request & reserve resources.
- If the requested resources are not available, the call will be blocked.
- Such a call admission is performed by the telephone network.
- For ex: In telephone network, we request resources when we dial a number.
 - i) If the required resources are allocated, the call is completed.
 - ii) If the required resources cannot be allocated, the call is blocked.

3) Call Setup Signaling

- A signaling protocol can be used to coordinate following activities.
 - 1) The per-hop allocation of local resources (Figure 5.23)
 - 2) The overall end-to-end decision of whether or not the call has been able to reserve sufficient resources at each and every router on the end-to-end path.
- The RSVP protocol is a call setup protocol for providing QoS guarantees.

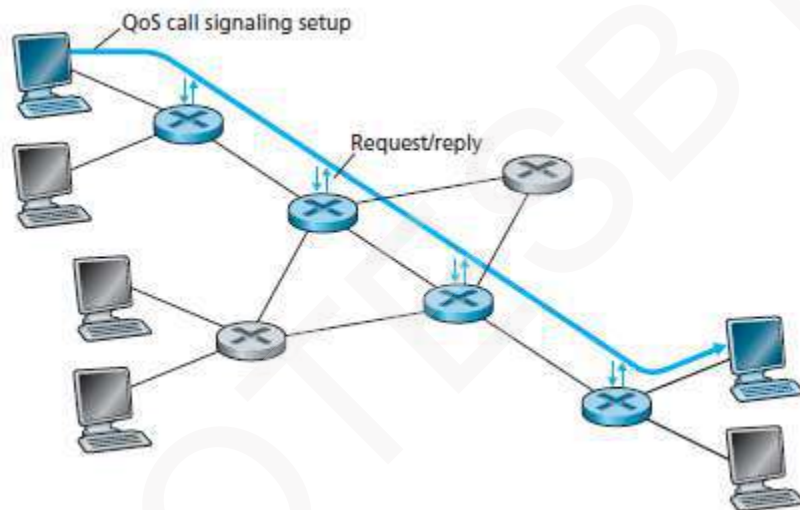


Figure 5.23: The call setup process

**MODULE-WISE QUESTIONS****PART 1**

- 1) Briefly explain properties of Video. (4*)
- 2) Briefly explain properties of Audio. (4)
- 3) Briefly explain three broad categories of multimedia network applications. (6*)
- 4) Briefly explain UDP streaming stored video application. (4*)
- 5) With a diagram, explain streaming stored video over HTTP/TCP. (6*)
- 6) Briefly explain adaptive streaming & DASH. (6*)
- 7) Briefly explain CDN, its motivation & its type. (6)
- 8) With a diagram, explain CDN Operation. (6*)
- 9) Briefly explain three approaches for cluster-selection. (6)
- 10) List out and explain three challenges to the design of real-time applications. (6)
- 11) Briefly explain two FEC mechanisms. (6*)
- 12) With a diagram, explain interleaving. (6*)

PART 2

- 13) Briefly explain RTP & its services. (4*)
- 14) With general format, explain various field of RTP. (6*)
- 15) Briefly explain SIP & its services. (4)
- 16) With a diagram, explain SIP call establishment. (6*)
- 17) With a diagram, explain SIP name translation & user location. (6)
- 18) Briefly explain three approaches to provide network-level support for multimedia applications. (6*)
- 19) Briefly explain dimensioning best effort networks. (6)
- 20) Briefly explain 2 approaches for providing isolation among traffic classes. (6)
- 21) With a diagram for each, explain FIFO and priority queue scheduling. (8*)
- 22) With a diagram for each, explain RRQ and WFQ scheduling. (8*)
- 23) With a diagram, explain the working of leaky bucket. (6*)
- 24) With a diagram, explain various components of DiffServ. (6*)
- 25) Briefly explain three mechanisms for guaranteed QoS. (6*)