# **LECTURENOTES**

# **COMPUTER NETWORK**

B.Tech,4<sup>TH</sup>Semester,CSE

Preparedby: MISS ALINA KUMARI SWAIN Assistant Professor in Computer Science & Engineering



# VikashInstitute ofTechnology, Bargarh

(ApprovedbyAICTE,NewDelhi&AffiliatedtoBPUT,Odisha) BarahagudaCanalChowk,Bargarh,Odisha-768040 www.vitbargarh.ac.in

# **DISCLAIMER**

- This document does not claim any originality and cannot be used as a substitute for prescribed textbooks.
- The information presented here is merely a collection by Miss ALINA KUMARI SWAIN with the inputs of students for their respective teaching assignments as an additional tool for the teaching- learning process.
- Various sources as mentioned at the reference of the document as well as freely available materials from internet were consulted for preparing this document.
- Further, this document is not intended to be used for commercial purpose and the authors are not accountable forany issues, legal or otherwise, arising out of use of this document.
- The author makes no representations or warranties with respect to the accuracy or completeness of the contents of this document and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose.

\*\*\*\*\*\*

# **COURSECONTENT COMPUTER NETWORK**

B.Tech,4<sup>TH</sup>Semester,CSE

#### Introduction to Networks $\succ$

Network hardware, Network software, OSI, TCP/IP Reference models, Example Networks: ARPANET, Internet.

Physical Layer: Data and signals: analog and digital, periodic analog signals, digital signals. transmission impairments, data rate limit, Guided transmission media twisted pairs, coaxial cable, fiber optics, Wireless transmission, unguided transmission media.

#### Datalink Layer

Design issues, framing, Error detection and correction, CRC codes

Elementary data link protocols: simplex protocol, A simplex stop and wait protocol for an error-free channel, A simplex stop and wait protocol for noisy channel.

Sliding Window protocols: A one-bit sliding window protocol, A protocol using Go-Back-N, A protocol using Selective Repeat, Example data link protocols.

Medium Access sub layer: The channel allocation problem, Multiple access protocols: ALOHA Carrier sense multiple access protocols, collision free protocols. Wireless LANs, Data link layer switching.

#### Connecting Devices

#### Learning bridges, spanning tree bridges, repeaters, hubs, bridges, switches, routers and gateways, definition of multiplexing and types.

Network Layer: Design issues, Routing algorithms: shortest path routing, Flooding, Hierarchical routing, Broadcast, Multicast, distance vector routing, link state protocols, path vector routing, Congestion Control Algorithms, Quality of Service

#### Internetworking

logical addressing, internet protocols, IP address, CIDR, IPv4 addressing, IPv6 Protocol addressing, addresses mapping, ICMP, IGMP, ARP, RARP, DHCP.

Transport Protocols: process to process delivery, UDP, TCP, TCP Service Model, TCP Sliding Window, TCP Congestion Control, congestion control and quality of service.

#### Application Layer

Introduction, providing services, Client server model, Standard client-server application-HTTP, FTP, electronic mail, TELNET, DNS

#### **{Page No. 96}**

#### **{Page No.172}**

**{Page No.266}** 

{**PageNo. 68**}

{PageNo. 1}

## REFERENCES COMPUTER NETWORK B.Tech,4<sup>TH</sup>Semester,CSE

#### **Books:**

- [1] James F. Kurose, K. W. Ross, "Computer Networking: A Top-Down Approach Featuring the Internet".
  Pearson EducationIntroduction to the Design and Analysis of Algorithms, Anany Levitin: 2nd Edition, 2009. Pearson.
- [2] An Engineering Approach to Computer Networks-S. Keshav, 2nd Edition, Pearson EducationA.V. Aho.
  - J. E. Hopcroft and J.D. Ullman, "The Design and Analysis of Computer Algorithms", Pearson Education.

#### **DigitalLearningResources:**

https://nptel.ac.in/courses/106105183

https://nptel.ac.in/courses/106105081

#### UNIT- I

#### **NETWORKS**

A network is a set of devices (often referred to as*nodes*)connected by communication links. A node can be a computer, printer, or any other device capable of sending and/or receiving data generated by other nodes on the network.

"Computer network" to mean a collection of autonomous computers interconnected by a single technology. Two computers are said to be interconnected if they are able to exchange information.

The connection need not be via a copper wire; fiber optics, microwaves, infrared, and communication satellites can also be used.

Networkscomeinmanysizes, shapes and forms, as we will see later.

Theyareusuallyconnectedtogethertomakelargernetworks, with the Internet being the most well-known example of a network of networks.

There is considerable confusion in the literature between acomputer **network** and a **distributed system**. The key distinction is that in a distributed system, a collection of independent computers appears to its users as a single coherent system. Usually, it has a single model or paradigm that it presents to the users. Often a layer of software on top of the operating system, called **middleware**, is responsible for implementing this model. Awell-known example of a distributed system is the **World Wide Web**. It runs on top of the Internet and presents a model in which everything looks like a document (Web page).

#### USESOFCOMPUTERNETWORKS

#### 1. BusinessApplications

- todistributeinformationthroughoutthecompany(**resource sharing**). sharingphysicalresourcessuchasprinters, and tapebackupsystems, is sharing information
- **client-servermodel**.Itiswidelyusedandformsthebasisofmuchnetwork usage.
- communication medium among employees.email (electronic mail), which employees generally use for a great deal of daily communication.
- Telephonecallsbetweenemployeesmaybecarriedbythecomputernetwork insteadofbythephonecompany. Thistechnologyiscalled IPtelephonyor VoiceoverIP(VoIP) when Internettechnologyisused.
- **Desktopsharing**letsremoteworkersseeandinteractwithagraphical computer screen
- doingbusinesselectronically, especially with customers and suppliers. This new model is called **e-commerce** (**electronic commerce**) and it has grown rapidly in recent years.

#### **2HomeApplications**

- **peer-to-peer**communication
- person-to-personcommunication

- electroniccommerce
- entertainment.(gameplaying,)

#### 3MobileUsers

- Textmessagingortexting
- Smartphones,
- GPS(GlobalPositioningSystem)
- m-commerce
- NFC(NearFieldCommunication)

#### **4SocialIssues**

Withthegoodcomesthebad, as this new-found freedom brings with it many unsolved social, political, and ethical issues.

Socialnetworks, messageboards, contents having sites, and a host of other applications allow people to share their views with like-minded individuals. As long as the subjects are restricted to technical topics or hobbies like gardening, not too many problems will arise.

Thetroublecomeswithtopicsthatpeopleactuallycareabout, likepolitics, religion, or sex. Views that are publicly posted may be deeply offensive to some people. Worse yet, they may not be politically correct. Furthermore, opinions need not be limited to text; high-resolution color photographs and video clipsare easily shared over computer networks. Some people take a live-and-let-live view, but others feel that posting certain material (e.g., verbal attacks on particular countries or religions, pornography, etc.) is simply unacceptable and that such content must be censored. Different countries have different and conflicting laws in this area. Thus, the debate rages.

Computer networks make it very easy to communicate. They also make it easy for the people who run the network to snoop on the traffic. This sets up conflictsoverissuessuchasemployeerightsversusemployerrights.Many people read and write email at work. Many employers have claimed the right to read and possibly messages, including censor employee messages sent from a home computer outside working hours. Not all employees agree with this, especially the latter part.

Anotherconflictiscenteredaroundgovernmentversuscitizen'srights.

A new twist with mobile devices is location privacy. As part of the process of providing service to your mobile device the network operators learn where you are at different times of day. This allows them to track your movements. They may know which nightclub you frequent and which medical center you visit.

**Phishing ATTACK:** *Phishing*is a type of social engineering *attack*often used tostealuserdata,includinglogincredentialsandcreditcardnumbers. Itoccurs when an attacker, masquerading as a trusted entity, dupes a victim into openinganemail, instantmessage, ortextmessage.

**BOTNET ATTACK:** Botnets can be used to perform <u>distributed denial-of-</u> <u>serviceattack</u>(DDoS attack), steal data, send spam, and allows the attacker to access the device and its connection.

The effectiveness of a data communications system depends on fourfundamental characteristics: delivery, accuracy, timeliness, and jitter.

I. **Delivery.** The system must deliver data to the correct destination. Data must be received by the intended device or user and only by that device or user.

2Accuracy. Thesystemmust deliver the data accurately. Data that have been

altered in transmission and left uncorrected are unusable.

3. **Timeliness**. Thesystemmust deliver data in a timely manner. Data delivered late are useless. In the case of video and audio, timely delivery means delivering data as they are produced, in the same order that they are produced, and without significant delay. This kind of delivery is called *real-time* transmission.

4. **Jitter**. Jitterreferstothevariationinthepacketarrivaltime. Itistheuneven delay in the delivery of audio or video packets. For example, let us assume that video packets are sent every 30 ms. If some of the packets arrive with 30-ms delay and others with 40-ms delay, an uneven quality in the video is the result.

Adatacommunicationssystemhasfivecomponents

I.**Message**.Themessageistheinformation(data)tobecommunicated. Popular forms of information include text, numbers, pictures, audio, and

video.2Sender.Thesenderisthedevicethatsendsthedatamessage.Itcanbea

computer, workstation, telephone handset, video camera, and so on.

3. **Receiver**.Thereceiveristhedevicethatreceivesthemessage.Itcanbea computer, workstation, telephone handset, television, and so on.

4. **Transmission medium**. The transmission medium is the physical path by which a message travels from sender to receiver. Some examples of transmissionmediaincludetwisted-pairwire, coaxialcable, fiber-opticcable, and radio waves.

5. **Protocol.** A protocol is a set of rules that govern data communications. It represents an agreement between the communicating devices. Without a protocol, two devices may be connected but not communicating, just as aperson speaking French cannot be understood by a person who speaks only Japanese.



#### DataRepresentation

Text
Numbers
Images
Audio
Video

#### DataFlow

Communication between two devices can be simplex, half-duplex, or full-duplex as shown in Figure.



*Simplex*Insimplexmode, the communication is unidirectional, as on a only street. Only one of the two devices on a link can transmit; the other can only receive (Figure a). Keyboards and traditional monitors are examples of simplex devices. *Half-Duplex* 

In half-duplex mode, each station can both transmit and receive, but not at the same time. When one device is sending, the other can only receive, and vice versa (Figure b).Walkie-talkies and CB (citizens band) radios are both half- duplex systems.

#### Full-Duplex

In full-duplex, both stations can transmit and receive simultaneously (Figure c). Onecommon exampleof full-duplexcommunication is thetelephonenetwork. Whentwopeoplearecommunicatingbyatelephoneline,bothcantalkand listenatthesametime.Thefull-duplexmodeisusedwhencommunicationin both directions is required all the time.

#### NetworkCriteria

Anetworkmustbeabletomeetacertainnumberofcriteria. Themost important of these are performance, reliability, and security.

#### Performance

Performancecanbemeasuredinmanyways, including transittime and response time. Transit time is the amount of time required for a message to travel from one device to another. Response time is the elapsed time between an inquiry and a response. The performance of a network depends on a number of factors, including the number of users, the type of transmission medium, the capabilities of the connected hardware, and the efficiency of the software.

Performance is often evaluated by two networking metrics: throughput anddelay. We often need more throughput and less delay. However, these two criteria are often contradictory. If we try to send more data to the network, we mayincreasethroughputbutweincreasethedelaybecauseoftraffic congestion in the network.

**Reliability:** In addition to accuracy of delivery, network reliability is measured bythefrequencyoffailure,thetimeittakesalinktorecoverfromafailure,and the network's robustness in a catastrophe.

*Security*:Network security issues include protecting data from unauthorized access, protecting data from damage and development, and implementing policies and procedures for recovery from breaches and data losses.

#### PhysicalStructures

Before discussing networks, we need to define some network attributes.

#### **TypeofConnection**

Anetworkistwoormoredevicesconnected throughlinks.Alinkisa communications pathway that transfers data from one device to another. There are two possible types of connections: point-to-point and multipoint.

Point-to-PointApoint-to-pointconnectionprovidesadedicatedlinkbetween

two devices. The entire capacity of the link is reserved for transmission

between those two devices. Most point-to-point connection suse an actual

lengthofwireorcabletoconnectthetwoends, but other options, such as microwave or satellite links, are also possible

When you change television channels by infrared remote control, you are establishing a point-to-point connection between the remote control and the television's control system.

**Multipoint**Amultipoint(alsocalledmulti-drop)connectionisoneinwhich more than two specific devices share a single link

In a multipoint environment, the capacity of the channel is shared, either spatially or temporally. If several devices can use the link simultaneously, it is a *spatially shared* connection. If users must take turns, it is a*timeshared* connection.



#### PhysicalTopology

The term*physical topology* refers to the way in which a network is laid out physically.

Two or more devices connect to a link; two or more links form a topology. The topology of a network is the geometric representation of the relationship of all the links and linking devices (usually called nodes) to one another.

Therearefourbasictopologiespossible:mesh,star,bus,andring



#### MESH:

A mesh topology is the one where every node is connected to every other nodein the network.



# A mesh topology can be a**full mesh topology**or a**partially connected mesh topology**.

Ina *fullmeshtopology*, every computer in the network has a connection to each of the other computers in that network. The number of connections in this

network can be calculated using the following formula (n is the number of computers in the network): n(n-1)/2

In a *partially connected mesh topology*, at least two of the computers in the network have connections to multiple other computers in that network. It is an inexpensive way toimplement redundancy in a network. In theevent that one of the primary computers or connections in the network fails, the rest of the network continues to operate normally.

#### Advantagesofameshtopology

- Canhandlehighamountsoftraffic,becausemultipledevicescantransmit data simultaneously.
- A failure of one device does not cause a break in the network or transmission of data.
- Addingadditionaldevicesdoesnotdisrupt datatransmissionbetweenother devices.

#### Disadvantagesofameshtopology

- Thecosttoimplementishigherthanothernetworktopologies, makingita less desirable option.
- Buildingandmaintaining thetopology is difficultand time consuming.
- Thechanceofredundantconnectionsishigh, which adds to the high costs and potential for reduced efficiency.

#### STAR:



Astarnetwork, startopology is one of the most common network setups. In thisconfiguration, every nodeconnects to acentral network device, like a hub, switch, or computer. The central network device acts as а serverand the peripheraldevicesactas clients. Dependingon the type of network cardused in each computer of the star topology, a coaxial cableor aRJ-45network cable is used to connect computers together.

Advantagesofstartopology

- Centralized management of the network, through the use of the central computer, hub, or switch.
- Easytoaddanothercomputertothenetwork.
- Ifonecomputeronthenetworkfails, the rest of the network continues to function normally.
- Thestartopologyisusedinlocal-areanetworks(LANs), High-speedLANs often use a star topology with a central hub.

Disadvantagesofstartopology

- Canhaveahighercosttoimplement, especially when using as witch or router the as central network device.
- Thecentralnetworkdevicedeterminestheperformanceandnumberof nodes the • network can handle.
- If the central computer, hub, or switch fails, the entire network goes down all and computers are disconnected from the network

#### **BUS:**



alinetopology, abustopology is a network setup in which each computer and network device are connected to a single cable or backbone. Advantagesofbustopology

- Itworkswellwhenyouhaveasmallnetwork.
- It's the easiest network topology for connecting computers or peripheralsin a linear fashion.
- It requires lesscablelengththanastartopology.

#### Disadvantages of bus topology

- It can be difficult to identify the problems if the whole network goes down.
- Itcanbehard totroubleshootindividualdevice issues.
- Bustopologyisnotgreatforlargenetworks.
- Terminatorsare required for bothends of the main cable.
- Additionaldevicesslowthe networkdown.
- If amaincableisdamaged, the networkfailsorsplits intotwo.



#### **RING:**

Aringtopology is a <u>network</u> configuration in which device connections create a circular <u>datapath</u>. In aring network, <u>packets</u> of data travel from one device to the next until they reach their destination. Most ring topologies allow packets to travelonly in one direction, called a **unidirectional** ring network. Others permit data to move in either direction, called **bidirectional**.

The major disadvantage of a ring topology is that if any individual connection in the ring is broken, the entire network is affected.

Ring topologies may be used in either local area networks (<u>LANs</u>) or wide area networks (<u>WANs</u>).

#### Advantagesofringtopology

- Alldataflowsinonedirection, reducing the chance of packet collisions.
- Anetworkserverisnotneededtocontrolnetworkconnectivitybetween each workstation.
- Datacantransferbetweenworkstationsathighspeeds.
- Additionalworkstationscanbeaddedwithoutimpactingperformanceof the network.

#### Disadvantagesofringtopology

- Alldatabeingtransferredoverthenetworkmustpassthrougheach workstation on the network, which can make it slower than a <u>star topology</u>.
- Theentirenetworkwillbeimpactedifoneworkstationshutsdown.
- The hardware needed to connect each workstation to the network is more expensive than Ethernet cards and hubs/switches.

**HybridTopology** A network can be hybrid. For example, we can have a mainstar topology with each branch connecting several stations in a bus topology as shown in Figure



#### TypesofNetworkbasedonsize

The types of network are classified based upon the size, the area it covers and its physical architecture. The three primary network categories are LAN, WAN and MAN. Each network differs in their characteristics such as distance, transmission speed, cables and cost.

Basictypes

#### LAN(LocalAreaNetwork)

Group of interconnected computers within a small area. (room, building, campus)

Twoormorepc'scanfromaLANtosharefiles, folders, printers, applications and other devices.

Coaxial or CAT 5 cables are normally used for connections.

Due to short distances, errors and noise are minimum.

Data transfer rate is 10 to 100 mbps.

Example: Acomputer labinaschool. MAN

(Metropolitan Area Network)Design to

extend over a large area.

ConnectingnumberofLAN'stoformlargernetwork, so that resources can be shared.

Networks can be up to 5 to 50 km.Owned

by organization or individual.Datatransfer

rateislow compare to LAN.

 $\label{eq:complete} Example: Organization with different branches located in the city.$ 

#### WAN (WideAreaNetwork)

Arecountryandworldwidenetwork.

Contains multiple LAN's and MAN's.

Distinguished in terms of geographical range.

Uses satellites and microwave relays.

Data transfer rate depends upon the ISP provider and varies over the location. Best example is the internet.

#### Othertypes

#### WLAN(WirelessLAN)

A LAN that uses high frequency radio waves for communication.

Provides short range connectivity with high speed data transmission.

#### PAN(PersonalAreaNetwork)

Networkorganizedbytheindividualuserforitspersonaluse.

#### SAN(StorageAreaNetwork)

Connects servers to data storage devices via fiber-optic cables.

E.g.: Used for daily backup of organization or a mirror copy

Atransmissionmediumcanbebroadlydefinedasanythingthatcancarry information from a source to a destination.



<u>**GuidedMedia</u>**:Guidedmedia, which are those that provide a medium from one device to another, include twisted-pair cable, coaxial cable, and fiber-optic cable.</u>

<u>Twisted-PairCable</u>: A twisted pair consists of two conductors (normallycopper), each with its own plastic insulation, twisted together. One of the wires is used to carry signals to the receiver, and the other is used only as a ground reference.



#### UnshieldedVersusShieldedTwisted-PairCable

The most common twisted-pair cable used in communications is referred to as unshielded twisted-pair (UTP). STP cable has a metal foil or braided mesh covering that encases each pair of insulated conductors. Although metal casing improves the quality of cable by preventing the penetration of noise orcrosstalk, it is bulkier and more expensive.





ThemostcommonUTP connector is RJ45(RJ stands forregistered jack)

#### **Applications**

Twisted-paircablesareusedintelephonelinestoprovidevoiceanddata channels. Local-areanetworks, suchasl0Base-Tandl00Base-T, alsousetwisted-pair cables.

#### **CoaxialCable**

Coaxial cable (or *coax*) carries signals of higher frequency ranges than those in twisted pair cable. coax has a central core conductor of solid or stranded wire (usuallycopper)enclosedinaninsulatingsheath, which is, inturn, encased in an outer conductor of metal foil, braid, or a combination of the two. The outer metallic wrapping serves both as a shield against noise and as the second conductor, which completes the circuit. This outer conductor is also enclosed in an insulating sheath, and the whole cable is protected by a plastic cover.



The most common type of connector used today is the Bayone-Neill-Concelman(BNe), connector.

#### Applications

Coaxialcablewaswidelyusedinanalogtelephonenetworks, digitaltelephone networks CableTVnetworksalsousecoaxialcables.

AnothercommonapplicationofcoaxialcableisintraditionalEthernetLANs

#### Fiber-OpticCable

Afiber-optic cableis madeof glassor plastic and transmitssignalsin theform of light. Light travels in a straight line as long as it is moving through a single uniform substance.

If a ray of light traveling through one substance suddenly enters another substance(ofadifferentdensity), the ray changes direction.

Bendingoflightray



Multimode is so named because multiple beams from a light source movethrough the core in different paths. How these beams move within the cable depends on the structure of the core, as shown in Figure.



In **multimode step-index fiber**, the density of the core remains constant from thecenter to the edges. A beam of lightmoves through this constant density ina straight line until it reaches the interface of the core and the cladding. The term *step index* refers to the suddenness of this change, which contributes to the distortion of the signal as it passes through the fiber.

A second type of fiber, called **multimode graded-index fiber**, decreases this distortion of the signal through the cable. The word*index*here refers to the index of refraction.

**Single-Mode:**Single-mode uses step-index fiber and a highly focused source oflightthatlimitsbeamstoasmallrangeofangles, all closet othehorizontal.



Thesubscriberchannel(SC)connector, Thestraight-tip(ST)connector, MT-RJ(mechanical transfer registered jack) is a connector

#### **Applications**

Fiber-optic cable is often found in backbone networks because its wide bandwidth is cost-effective..

SomecableTVcompaniesuseacombinationofopticalfiberandcoaxialcable, thus creating a hybrid network.

Local-area networks such as 100Base-FX network (Fast Ethernet) and 1000Base-X also use fiber-optic cable

#### <u>AdvantagesandDisadvantagesofOpticalFiber</u>

AdvantagesFiber-

opticcablehasseveraladvantagesovermetalliccable(twistedpairorcoaxial).

1 Higherbandwidth.

2 Less signal attenuation. Fiber-optic transmission distance is significantly greaterthan that of other guided media. A signal can run for 50 km without requiring regeneration. We need repeaters every 5 km for coaxial or twisted- pair cable.

3 Immunity to electromagnetic interference. Electromagnetic noise cannot affect fiber-optic cables.

4 Resistance to corrosive materials. Glass is more resistant to corrosive materials than copper.

5 Lightweight.Fiber-opticcablesaremuchlighterthancoppercables.

6 Greater immunity to tapping. Fiber-optic cables are more immune to tapping than copper cables. Copper cables create antenna effects that can easily be tapped.

**Disadvantages**Therearesomedisadvantagesintheuseofopticalfiber. 1Installation and maintenance

2 Unidirectionallightpropagation.Propagation of light isunidirectional. If we need bidirectional communication, two fibers are needed.

3 Cost. The cable and the interfaces are relatively more expensive than those of other guided media. If the demand for bandwidth is not high, often the use of optical fiber cannot be justified.

Unguided media transport electromagnetic waves without using a physical conductor. This type of communication is often referred to as wireless communication. RadioWaves Microwaves Infrared



Unguided signals can travel from the source to destination in several ways: ground propagation, sky propagation, and line-of-sight propagation, as shown in Figure



#### RadioWaves

Electromagnetic waves ranging in frequencies between 3 kHz and 1 GHz are normally called radio waves. Radio waves are omni directional. When anantennatransmitsradiowaves, they are propagated in all directions. This means that the sending and receiving antennas do not have to be aligned. A sending antenna sends waves that can be received by any receiving antenna. The omni directional property has a disadvantage, too. The radio waves transmitted by one antenna are susceptible to interference by another antenna that may send signals using the same frequency or band.

#### *OmnidirectionalAntenna*

Radiowavesuseomnidirectionalantennasthatsendoutsignalsinall directions. Based on the wavelength, strength, and the purpose of transmission, we can have several types of antennas. Figure shows an omnidirectionalantenna.



#### Applications

The Omni directional characteristics of radio waves make them useful for multicasting, in which there is one sender but many receivers. AM and FM radio, television, maritime radio, cordless phones, and paging are examples of multicasting.

#### Microwaves

Electromagnetic waves having frequencies between 1 and 300 GHz are called microwaves. Microwaves are unidirectional. The sending and receiving antennas needtobealigned. The unidirectional property has an obvious advantage. A pair of antennas can be aligned without interfering with another pair of aligned antennas *UnidirectionalAntenna* 

Microwaves need unidirectional antennas that send out signals in one direction. Two types of antennas are used for microwave communications: the parabolic dish and the horn





a. Dish antenna Applications:

b. Horn antenna

Microwavesareusedforunicastcommunicationsuchascellulartelephones, satellite networks, and wireless LANs

#### Infrared

Infrared waves, with frequencies from 300 GHz to 400 THz (wavelengths from 1 mm to 770 nm), can be used for short-range communication. Infrared waves, havinghighfrequencies, cannot penetrate walls. This advantageous

characteristic prevents interference between one system and another; a shortrangecommunicationsysteminoneroomcannotbeaffectedbyanother system in the next room.

When we use our infrared remote control, we do not interfere with the use of the remote by our neighbors. Infrared signals useless for long-range communication. In addition, we cannot use infrared waves outside a building because the sun's rays contain infrared waves that can interfere with the communication.

#### **Applications:**

Infrared signals can be used for short-range communication in a closed area using line-of-sight propagation.

#### Switching

A network is a set of connected devices. Whenever we have multiple devices, we have the problem of how to connect them to make one-to-one communication possible. One solution is to make a point-to-point connection betweeneachpairofdevices(ameshtopology)orbetweenacentraldevice and every other device (a star topology). These methods, however, are impractical and wasteful when applied to very large networks.

The number and length of the links require too much infrastructure to be cost-efficient, and the majority of those links would be idle most of the time.

A better solution is switching. A switched network consists of a series ofinterlinked nodes, called switches. Switches are devices capable of creatingtemporary connections between two or more devices linked to the switch. In a switched network, some of these nodes are connected to the end systems (computers or telephones, for example). Others are used only for routing. Figure shows a switched network.



We can then divide today's networks into three broad categories: circuit- switched networks, packet-switched networks, and message-switched. Packet- switched networks can further be divided into two subcategories-virtual-circuit networks and datagram networks as shown in Figure.



#### **CIRCUIT-SWITCHEDNETWORKS**

Acircuit-switchednetworkconsistsofasetofswitchesconnectedby physicallinks.Aconnectionbetweentwostationsisadedicatedpathmadeof oneormorelinks.However,eachconnectionusesonlyonededicatedchannel oneachlink.EachlinkisnormallydividedintonchannelsbyusingFDMorTDM.

In circuit switching, the resources need to be reserved during the setup phase;

the resources remain dedicated for the entire duration of data transfer until the teardown phase



#### **ThreePhases**

The actual communication in a circuit-switched network requires three phases: connection setup, data transfer, and connection teardown.

#### **SetupPhase**

Before the two parties (or multiple parties in a conference call) can communicate, a dedicated circuit (combination of channels in links) needs to be established. Connection setup means creating dedicated channels between the switches. For example, in Figure, when system A needs to connect to system M, it sends a setup request that includes the address of system M, to switch I. Switch I finds a channel between itself and switch IV that can be dedicated for thispurpose.SwitchIthensendstherequesttoswitchIV,whichfindsa dedicated channel between itself and switch III. Switch III informs system M of system A's intention at this time.

In the next step to making a connection, an acknowledgment from system M needs to be sent in the opposite direction to system A. Only after system A receives this acknowledgment is the connection established.

#### **DataTransferPhase**

After the establishment of the dedicated circuit (channels), the two parties can transfer data.

#### **TeardownPhase**

When one of the parties needs to disconnect, asignal issent to each switch to release the resources.

#### Efficiency

It can be argued that circuit-switched networks are not as efficient as the other two types of networks because resources are allocated during the entireduration of the connection. These resources are unavailable to other connections.

#### Delay

Although a circuit-switched network normally has low efficiency, the delay in thistypeofnetworkisminimal.Duringdatatransferthedataarenotdelayed at each switch; the resources are allocated for the duration of the connection.

The total delay is due to the time needed to create the connection, transfer data, and disconnect the circuit.

# Switchingatthephysicallayerinthetraditionaltelephonenetwork usesthecircuit-switchingapproach.

#### DATAGRAM NETWORKS

Inapacket-switchednetwork, there is no reservation; resources are allocated on demand. The allocation is done on a first come, first-served basis. When a switch receives a packet, no matter what is the source or destination, the packet must wait if there are other packets being processed. This lack of reservation may create delay. For example, if we do not have a reservation at a restaurant, we might have to wait.

In a datagram network, each packet is treated independently of all others. Packets in this approach are referred to as datagrams. Datagram switching is normally done at the network layer.

Figure shows how the datagram approach is used to deliver four packets from station A to station X. The switches in a datagram network are traditionally referred to as routers.

The datagram networks are sometimes referred to as connectionless networks. The term *connectionless* heremeans that the switch (packet switch) does not keep information about the connection state. There are no setup or teardown phases. Each packet is treated the same by a switch regardless of its source or destination.

A switch in a datagram network uses a routing table that is based on the destination address. The destination address in the header of a packet in a datagram network remains the same during the entire journey of the packet.



#### Efficiency

Theefficiencyofadatagramnetworkisbetterthanthatofacircuit-switched network; resources are allocated only when there are packets to be transferred. **Delay** There may be greater delay in a datagram network than in a virtual-circuit network. Although there are no setup and teardown phases, each packet may experience a wait at a switch before it is forwarded. In addition, since not all packets in a message necessarily travel through the same switches, the delay is not uniform for the packets of a message.

SwitchingintheInternetisdonebyusingthedatagramapproachto packet switching at the network layer.

#### VIRTUAL-CIRCUITNETWORKS

A virtual-circuit network is a cross between a circuit-switched networkandadatagramnetwork. It has some characteristics of both.



1. As in a circuit-switched network, there are setup and teardown phases in addition to the data transfer phase.

2. Resourcescanbeallocatedduringthesetupphase, as inacircuit-switched network, or ondemand, as in a datagram network.

3. As in a datagram network, data are packetized and each packet carries an address in the header. However, the address in the header has local jurisdiction (it defines what should be the next switch and the channel on which the packet is being carried), not end-to-end jurisdiction.

4. As in a circuit-switched network, all packets follow the same path established during the connection.

5. A virtual-circuit network is normally implemented in the data link layer, while a circuit-switched network is implemented in the physical layer and a datagram network in the network layer.

#### **Addressing**

Inavirtual-circuitnetwork, two types of addressing are involved: global and local (virtual-circuit identifier).

#### **GlobalAddressing**

A source or a destination needs to have a global address-an address that can be unique in the scope of the network.

#### <u>Virtual-CircuitIdentifier</u>

The identifier that is actually used for data transfer is called the virtual-circuit identifier (VCI). A VCI, unlike a global address, is a small number that has only switch scope; it is used by a frame between two switches. When a frame arrives at a switch, it has a VCI; when it leaves, it has a different VCl.

Figure shows how the VCI in a data frame changes from one switch to another. NotethataVCIdoesnotneedtobealargenumbersinceeachswitchcanuse its own unique set of VCls.



#### ThreePhases

Three phases in a virtual-circuit network: setup, data transfer, and teardown. We first discuss the data transfer phase, which is more straightforward; we then talk about the setup and teardown phases.

#### DataTransferPhase



To transfer a frame from a source to its destination, all switches need to haveatableentryforthisvirtualcircuit. The table, inits simplest form, has four columns.

We show later how the switches make their table entries, but for the moment we assume that each switch has a table with entries for all activevirtual circuits. Figure shows such a switch and its corresponding table.

Figure shows a frame arriving at port 1 with a VCI of 14. When the frame arrives, the switch looks in its table to find port 1 and a VCI of 14. When it is found, the switch knows to change the VCI to 22 and send out the frame from port 3.

Figure shows how a frame from source A reaches destination B and how its VCI changes during the trip.



EachswitchchangestheVClandroutestheframe.

The data transfer phase is active until the source sends all its frames tothe destination. The procedure at the switch is the same for each frame of a message. The process creates a virtual circuit, not a real circuit, between the source and destination.

#### SetupPhase

In the setup phase, a switch creates an entry for a virtual circuit. For example, suppose source A needs to create a virtual circuit to B. Two steps are required: the setup request and the acknowledgment.

<u>Setup Request</u> A setup request frame is sent from the source to the destination. Figure shows the process.



a. SourceAsendsasetupframetoswitch1.

b. Switch1 receives these tuprequest frame. It knows that a framegoing from A to B goes out through port 3. For the moment, assume that it knows theoutputport. Theswitchcreates an entry inits table for this virtual circuit, but it is only able to fill three of the four columns. The switch assigns the incoming port (1) and chooses an available incoming VCI (14) and the outgoing port (3). It does not yet know the outgoing VCI, which will be found during the acknowledgment step. The switch then forwards the frame through port 3 to switch 2.

c. Switch 2 receives the setup request frame. The same events happen here asat switch 1; three columns of the table are completed: in this case, incoming port (l), incoming VCI (66), and outgoing port (2).

d. Switch 3 receives the setup request frame. Again, three columns are completed: incoming port (2), incoming VCI (22), and outgoing port (3).

e. Destination B receives the setup frame, and if it is ready to receive frames fromA, itassignsaVCI to the incoming frames that come from A, in this case

77. This VCI lets the destination know that the frames come from A, and not other sources.

Acknowledgment A special frame, called the acknowledgment frame, completes the entries in the switching tables.



a. The destination sends an acknowledgment to switch 3. The acknowledgment carries the global source and destination addresses so the switch knows which entry in the table is to be completed. The frame also carries VCI 77, chosen by the destination as the incoming VCI for frames from A. Switch 3 uses this VCI to complete the outgoing VCI column for this entry. Note that 77 is the incoming VCI for destination B, but the outgoing VCI for switch 3.

b. Switch 3 sends an acknowledgment to switch 2 that contains its incoming VCI in the table, chosen in the previous step. Switch 2 uses this as the outgoing VCI in the table.

c. Switch 2 sends an acknowledgment to switch 1 that contains its incoming VCI in the table, chosen in the previous step. Switch 1 uses this as the outgoing VCI in the table.

d. Finally switch 1 sends an acknowledgment to source A that contains its incoming VCI in the table, chosen in the previous step.

e. The source uses this as the outgoing VCI for the data frames to be sent to destination B.

#### TeardownPhase

In this phase, source A, after sending all frames to B, sends a special frame calleda*teardownrequest*.DestinationBrespondswithateardown confirmation frame. All switches delete the corresponding entry from their tables.

#### Efficiency

In virtual-circuit switching, all packets belonging to the same source and destination travel the same path; but the packets may arrive at the destination with different delays if resource allocation is on demand.

#### Delay

In a virtual-circuit network, there is a one-time delay for setup and a one-time delay for teardown. If resources are allocated during the setup phase, there isno wait time for individual packets. Figure shows the delay for a packettraveling through two switches in a virtual-circuit network



Switching at the data link layer in a switched WAN is normally implemented byusing virtual-circuittechniques.

#### <u>Comparison</u>

Issue	Datagram network	Virtual-circuit network
Circuit setup	Not needed	Required
Addressing	Each packet contains the full source and destination address	Each packet contains a short VC number
State information	Routers do not hold state information about connections	Each VC requires router table space per connection
Routing	Each packet is routed independently	Route chosen when VC is set up; all packets follow it
Effect of router failures	None, except for packets lost during the crash	All VCs that passed through the failed router are terminated
Quality of service	Difficult	Easy if enough resources can be allocated in advance for each VC
Congestion control	Difficult	Easy if enough resources can be allocated in advance for each VC

 ${\it Diagrams from Tanenbaum Textbook}$ 



Figure 2-42. (a) Circuit switching. (b) Packet switching.



#### <u>OSI</u>

- OSIstandsforOpenSystemsInterconnection
- CreatedbyInternationalStandardsOrganization(ISO)
- Wascreatedasaframeworkandreferencemodeltoexplainhowdifferent networking technologies work together and interact
- Itisnotastandardthatnetworkingprotocolsmustfollow
- Eachlayerhasspecificfunctionsitisresponsiblefor
- Alllayers worktogetherinthecorrect orderto movedataarounda network



How Data Is Referred to in the OSI Model		
Data	• Application, Presentation, and Session layers	
Segment	• Transport layer	
Packet	• Networking layer	
Frame	• Data Link layer	
Bits	• Physical layer	

#### <u>PhysicalLayer</u>

- Dealswithallaspectsofphysicallymovingdatafromonecomputertothenext
- Convertsdatafromtheupperlayersinto1sand0sfortransmissionovermedia
- Defineshowdataisencodedontothemediatotransmitthedata
- Definedonthislayer:Cablestandards,wirelessstandards,andfiberoptic standards. Copper wiring, fiber optic cable, radio frequencies, anything that can be used to transmit data is defined on the Physical layer of the OSI Model
- Deviceexample:Hub
- Usedtotransmitdata

#### <u>DataLinkLayer</u>

- Isresponsibleformovingframesfromnodetonodeorcomputertocomputer
- Can move frames from one adjacent computer to another, cannot move frames across routers
- Encapsulation =frame
- RequiresMACaddressorphysicaladdress
- ProtocolsdefinedincludeEthernetProtocolandPoint-to-PointProtocol (PPP)
- Deviceexample:Switch
- Two sublayers: LogicalLinkControl(LLC)andthe MediaAccessControl(MAC)
- LogicalLinkControl(LLC)
- -DataLinklayeraddressing,flowcontrol,addressnotification,errorcontrol
- MediaAccessControl(MAC)
- -Determineswhichcomputerhasaccesstothenetworkmediaatanygiven time
- -Determineswhereoneframeendsandthenextonestarts, calledframe synchronization

#### <u>NetworkLayer</u>

- Responsibleformovingpackets(data)fromoneendofthenetworktothe other, called *end-to-end communications*
- Requires logical addresses such as IP addresses
- Deviceexample:Router
- -Routing is the ability of various network devices and their related software to move data packets from source to destination

#### <u>TransportLayer</u>

- Takes data from higher levels of OSI Model and breaks it into segments that can be sent to lower-level layers for data transmission
- Conversely, reassembles datasegments into data that higher level protocols and applications can use
- Alsoputssegmentsincorrectorder(calledsequencing)sotheycanbe reassembled in correct order at destination
- Concernedwiththereliabilityofthetransportofsentdata
- Mayusea*connection-orientedprotocol*suchasTCPtoensuredestination received segments
- Mayusea*connectionlessprotocol*suchasUDPtosendsegmentswithout assurance of delivery
- Usesportaddressing

### <u>SessionLayer</u>

- Responsibleformanagingthedialogbetweennetworkeddevices
- Establishes, manages, and terminates connections
- Providesduplex, half-duplex, or simplex communications between devices
- Providesproceduresforestablishingcheckpoints, adjournment, termination, and restart or recovery procedures

#### <u>PresentationLayer</u>

- Concernedwithhowdatais presentedtothenetwork
- Handlesthreeprimarytasks:-Translation,-Compression,-Encryption



### **ApplicationLayer**

- Containsallservicesorprotocolsneededbyapplicationsoftwareoroperating system to communicate on the network
- Examples
- -FirefoxwebbrowserusesHTTP(Hyper-TextTransportProtocol)
- -E-mailprogrammayusePOP3(PostOfficeProtocolversion3)toreade-mails and SMTP (Simple Mail Transport Protocol) to send e-mails

#### TheinteractionbetweenlayersintheOSImodel



Physical communication

AnexchangeusingtheOSImodel



SUMMARY:



#### TCP/IPModel(TransmissionControlProtocol/InternetProtocol)

-Aprotocolsuiteisalargenumberofrelatedprotocolsthatworktogetherto allow networked computers to communicate



Relationship of layers and addresses in TCP/IP

#### **ApplicationLayer**

- Application layer protocols define the rules when implementing specific network applications
- Rely on the underlyinglayerstoprovide accurate and efficient data delivery
- Typicalprotocols:
- FTP-FileTransferProtocol
- Forfiletransfer
- Telnet-Remoteterminalprotocol
- Forremoteloginonanyothercomputeron thenetwork
- SMTP-SimpleMailTransferProtocol
- Formailtransfer
- HTTP-HypertextTransferProtocol
- ForWebbrowsing
- Encompasses same functions as these OSI Model layers Application Presentation Session

## <u>TransportLayer</u>

#### TCP&UDP

- TCPisaconnection-orientedprotocol
- Doesnotmeanithasaphysicalconnectionbetweensenderandreceiver
- TCPprovidesthefunctiontoallowaconnectionvirtuallyexists-alsocalled virtual circuit
- UDPprovidesthefunctions:
- o Dividingachunkofdataintosegments
- Reassemblysegmentsintotheoriginalchunk
- $\circ \ \ \, \text{Provide} further the functions such as reordering and data resend$
- Offeringareliablebyte-streamdeliveryservice
- Functions thesame as the Transportlayerin OSI
- Synchronize source and destination computers to set up the session between the respective computers

#### InternetLayer

• The network layer, also called the internet layer, deals with packets and connects independent networks to transport the packets across network boundaries. The network layer protocols are the IP and the Internet Control Message Protocol (<u>ICMP</u>), which is used for error reporting.

#### Host-to-networklayer

The Host-to-networklayer is the lowest layer of the TCP/IP reference model. It combines the link layer and the physical layer of the ISO/OSI model. At this layer, data is transferred between adjacent network nodes in a WAN or between nodes on the same LAN.



OSI MODEL	TCP/IP MODEL
Contains 7 Layers	Contains 4 Layers
Uses Strict Layering resulting in vertical layers.	Uses Loose Layering resulting in horizontal layers.
Supports both connectionless & connection-oriented communication in the Network layer, but only connection-oriented communication in Transport Layer	Supports only connectionless communication in the Network layer, but both connectionless & connection- oriented communication in Transport Layer
It distinguishes between Service, Interface and Protocol.	Does not clearly distinguish between Service, Interface and Protocol.
Protocols are better hidden and can be replaced relatively easily as technology changes (No transparency)	Protocols are not hidden and thus cannot be replaced easily. (Transparency) Replacing IP by a substantially different protocol would be virtually impossible
OSI reference model was devised before the corresponding protocols were designed.	The protocols came first and the model was a description of the existing protocols

#### THEINTERNET

The Internet has revolutionized many aspects of our daily lives. It has affected the way we do business as well as the way we spend our leisure time. Count thewaysyou'veusedtheInternetrecently.Perhapsyou'vesentelectronic mail (email) to a business associate, paid a utility bill, read a newspaper from a distant city, or looked up a local movie schedule-all by using the Internet. Or maybe you researched a medical topic, booked a hotel reservation, chatted with a fellow Trekkie, or comparison-shopped for a car. The Internet is a communication system that has brought a wealth of information to our fingertips and organized it for our use.

#### ABriefHistory

A network is a group of connected communicating devices such as computers and printers. An internet (note the lowercase letter i) is two or more networks that can communicate with each other. The most notable internetiscalledtheInternet(uppercaseletterI),acollaborationofmorethan hundreds of thousands of interconnected networks. Private individuals as well as various organizations such as government agencies, schools, research facilities, corporations, and libraries in more than 100 countries use the Internet. Millions of people are users. Yet this extraordinary communication system only came into being in 1969.

In the mid-1960s, mainframe computers in research organizations were standalone devices. Computers from different manufacturers were unable to communicatewithoneanother.TheAdvancedResearchProjectsAgency
(ARPA) in the Department of Defense (DoD) was interested in finding a way to connect computers so that the researchers they funded could share their findings, thereby reducing costs and eliminating duplication of effort.

In1967, atan Association for Computing Machinery (ACM) meeting, ARPA presented its ideas for ARPANET, a small network of connected computers. The idea was that each host computer (not necessarily from the same manufacturer) would be attached to a specialized computer, called an *inteiface message processor* (IMP). The IMPs, in turn, would be connected to one another. Each IMP had to be able to communicate with other IMPs as well as with its own attached host. By 1969, ARPANET was a reality. Four nodes, at the University of California at Los Angeles (UCLA), the University of California at Santa Barbara (UCSB), Stanford Research Institute (SRI), and the University of Utah, were connected via the IMPs to form a network. Software called the *Network Control Protocol* (NCP) provided communication between the hosts.

In 1972, Vint Cerf and Bob Kahn, both of whom were part of the core ARPANET group, collaborated on what they called the*Internetting Projec1*. Cerf and Kahn's landmark 1973 paper outlined the protocols to achieve end-to-end delivery of packets. This paper on Transmission Control Protocol (TCP) includedconceptssuchasencapsulation, thedatagram, and the functions of agateway. Shortly thereafter, authorities made adecision to split TCP into two protocols: Transmission Control Protocol (TCP) and Internetworking Protocol (IP). IP would handle datagram routing while TCP would be responsible for higher-level functions such as segmentation, reassembly, and error detection. The internetworking protocol became known as TCPIIP. TheInternetToday

The Internet has come a long way since the 1960s. The Internet today is not a simple hierarchical structure. It is made up of many wide- and local-area networks joined by connecting devices and switching stations. It is difficult to give an accurate representation of the Internet because it is continually changing-new networks are being added, existing networks are adding addresses, and networks of defunct companies are being removed. Today most end users who want Internet connection use the services of Internet service providers (ISPs). There are international service providers, national service providers, regional service providers, and local service providers. The Internet today is run by private companies, not the government. Figure 1.13 shows a conceptual (not geographic) view of the Internet.



b. Interconnection of national ISPs

#### InternationalInternetServiceProviders:

At the top of the hierarchy are the international service providers that connect nations together.

### NationalInternetServiceProviders:

The national Internet service providers are backbone networks created and maintained by specialized companies. There are many national ISPs operating in North America; some of the most well known are SprintLink, PSINet, UUNet Technology, AGIS, and internet Mel. To provide connectivity between the end users, these backbone networks are connected by complex switching stations (normally run by a third party) called network access points (NAPs). Some national ISP networks are also connected to one another by private switching stations called *peering points*. These normally operate at a high data rate (up to 600 Mbps).

### RegionalInternetServiceProviders:

Regional internet service providers or regional ISPs are smaller ISPs that are connected to one or more national ISPs. They are at the third level of the hierarchywithasmallerdatarate.*LocalInternetServiceProviders*:

Local Internet service providers provide direct service to the end users. The local ISPs can be connected to regional ISPs or directly to national ISPs. Most end users are connected to the local ISPs. Note that in this sense, a local ISPcanbeacompanythatjustprovidesInternetservices, acorporation with network а services that supplies to its own employees, nonprofit or а organization, such as a college or a university, that runsits own network. Each of these local ISPs can be connected to a regional or national service provider.

# UNIT-II

# DATALINKLAYERFUNCTIONS(SERVICES)

### 1. Providingservices to thenetwork layer:

1<u>Unacknowledgedconnectionlessservice</u>.

Appropriateforlowerrorrateandreal-timetraffic.Ex:Ethernet

2. <u>Acknowledgedconnectionlessservice</u>.

Usefulinunreliablechannels,WiFi.Ack/Timer/Resend

3. <u>Acknowledgedconnection-orientedservice</u>.

Guarantee frames are received exactly once and in the right order. Appropriate over long, unreliable links such as a satellite channel or a long- distance telephone circuit

- Framing: Framesare thest reams of bits received from the network layer intomanage abled at a units. This division of stream of bits is done by Data Link Layer.
- 3. **Physical Addressing:** The Data Link layer adds a header to the frame in order to define physical address of the sender or receiver of the frame, if the frames are to be distributed to different systems on the network.
- 4. FlowControl:Areceivingnodecanreceivetheframesatafasterrate than it can process the frame. Without flow control, the receiver's buffer can overflow, and frames can get lost. To overcome this problem, the data link layer uses the flow control to prevent the sending node on one side of the link from overwhelming the receiving node on another side of the link. This prevents traffic jam at the receiver side.
- 5. ErrorControl:Error control isachieved by addinga trailer attheendof the frame. Duplication of frames are also prevented by using this mechanism. Data Link Layers adds mechanism to prevent duplication of frames.

**Error detection:** Errors can be introduced by signal attenuation and noise. Data Link Layer protocol provides a mechanism to detect one or more errors. This is achieved by adding error detection bits in the frame and then receiving node can perform an error check.

**Error correction:** Error correction is similar to the Error detection, except that receiving node not only detects the errors but also determine where the errors have occurred in the frame.

- 6. Access Control: Protocols of this layer determine which of the devices has control over the link at any given time, when two or more devices are connected to the same link.
- 7. **Reliabledelivery**:DataLinkLayerprovidesareliabledeliveryservice, i.e., transmits the network layer datagram without any error. A reliable deliveryserviceisaccomplishedwithtransmissionsand acknowledgements.Adatalinklayermainlyprovidesthereliabledelivery

service over the links as they have higher error rates and they can be corrected locally, link at which an error occurs rather than forcing to retransmit the data.

8. Half-Duplex& Full-Duplex: In a Full-Duplex mode, both the nodes can transmit the data at the same time. In a Half-Duplex mode, only one node can transmit the data at the same time.

### FRAMING:

To provide service to the network layer, the data link layer must use the service provided to it by the physical layer. What the physical layer does is accept a raw bit stream and attempt to deliver it to the destination. This bit stream is not guaranteed to be error free. The number of bits received may be less than, equal to, or more than the number of bits transmitted, and they may have different values. It is up to the data link layer to**detect and, if necessary,** correct errors. The usual approach is for the data link layer to break the bit discrete frames compute stream up into and the checksum for eachframe(framing). When a frame arrives at the destination, the checksum is recomputed. If the newly computed checksum is different from the one contained in the frame, the data link layer knows that an error has occurred and takes steps to deal with it (e.g., discarding the bad frame and possibly also sending back an error report). We will look at four framing methods:

- 1. Charactercount.
- 2. Flagbyteswithbytestuffing.
- 3. Starting and ending flags, with bits tuffing.
- 4. Physicallayercodingviolations.

**Character count** method uses a field in the header to specify the number of characters in the frame. When the data link layer at the destination sees the character count, it knows how many characters follow and hence where the end of the frame is. This technique is shown in Fig. (a) For four frames of sizes 5, 5, 8, and 8 characters, respectively.



The trouble with this algorithm is that the count can be garbled by a transmission error. For example, if the character count of 5 in the second frame of Fig. (b) becomes a 7, the destination will get out of synchronization and will be unable to locate the start of the next frame. Even if the checksum is incorrect so the destination knows that the frame is bad, it still has no way of telling where the frame starts. Sending а frame back to the next source askingforaretransmissiondoesnothelpeither, since the destination does not know how many characters to skip over to get to the start of the retransmission. For this reason, the character count method is rarely used anymore.

**Flag bytes with byte stuffing**method gets around the problem of resynchronization after an error by having each frame start and end with special bytes. In the past, the starting and ending bytes were different, but in recent years most protocols have used the same byte, called a flag byte, as both the starting and ending delimiter, as shown in Fig. (a) as FLAG.In this way, if the receiver ever loses synchronization, it can just search for the flag byte to find the end of the current frame. Two consecutive flag bytes indicate the end of one frame and start of the next one.





It may easily happen that the flag byte's bit pattern occurs in the data. This situation will usually interfere with the framing. One way to solve this problem is to have the sender's data link layer insert a special escape byte (ESC)justbeforeeach"accidental"flagbytein thedata.Thedatalinklayer onthereceivingendremovestheescapebytebeforethedataaregivento thenetworklayer.Thistechniqueiscalledbytestuffingorcharacterstuffing.

Thus, a framing flag byte can be distinguished from one in the data by the absenceorpresenceofanescapebytebeforeit.

What happens if an escape byte occurs in the middle of the data? The answer is that, it too is stuffed with an escape byte. Thus, any single escape byte is part of an escape sequence, whereas a doubled one indicates that a singleescapeoccurrednaturallyinthedata.Someexamplesareshownin Fig. (b). In all cases, the byte sequence delivered after de stuffing is exactly the same as the original byte sequence.

A major disadvantage of using this framing method is that it is closelytied to the use of 8-bit characters. Not all character codes use 8-bit characters. For example UNICODE uses 16-bit characters, so a new technique had to be developed to allow arbitrary sized characters

**Startingandendingflags, withbitstuffing**allowsdataframesto contain an arbitrary number of bits and allows character codes with an arbitrary numberofbitspercharacter.Itworkslikethis.Eachframe beginsandends withaspecialbitpattern,01111110(infact,aflagbyte).Wheneverthe sender's data link layer encounters five consecutive 1s in the data, it automatically stuffs a 0 bit into the outgoing bit stream. This bit stuffing is analogous to byte stuffing, in which an escape byte is stuffed into the outgoing character stream before a flag byte in the data.

Whenthereceiverseesfiveconsecutiveincoming1bits,followedbya0bit, it automatically de- stuffs (i.e., deletes) the 0 bit. Just as byte stuffing is completely transparent to the network layer in both computers, so is bit stuffing. If the user data contain the flag pattern, 01111110, this flag is transmittedas011111010butstoredinthereceiver'smemoryas0111110.

(a) 01101111111111111111110010

(b) 011011111011111011111010010 Stuffed bits

(c) 0110111111111111111110010

Fig:Bitstuffing.(a)Theoriginaldata.(b)Thedataastheyappearontheline.

(c) The data as they are stored in the receiver's memory after destuffing.

With bit stuffing, the boundary between two frames can be unambiguously recognized by the flag pattern. Thus, if the receiver loses track of where it is, all it has to do is scan the input for flag sequences, since they can only occurat frame boundaries and never within the data.

**Physical layer coding violations** method of framing is only applicable to networks in which the encoding on the physical medium contains some redundancy. For example, some LANs encode 1 bit of data by using 2 physical bits. Normally, a 1 bit is a high-low pair and a 0 bit is a low-high pair. The scheme means that everydatabithasa transitioninthemiddle, makingit easyforthereceivertolocatethebitboundaries. The combinationshigh-

highandlow-lowarenotusedfordatabutareusedfordelimitingframesin some protocols.



As a final note on framing, many data link protocols use combination of a character count with one of the other methods for extra safety. When a frame arrives, the count field is used to locate the end of the frame. Only if the appropriate delimiter is present at that position and the checksum is correct is the frame accepted as valid. Otherwise, the input stream is scanned for the next delimiter



# ELEMENTARYDATALINKPROTOCOLS

# SimplestProtocol



It is very simple. The sender sends a sequence of frames without even thinking about the receiver. Data are transmitted in one direction only. Both sender& receiver always ready. Processing time can be ignored. Infinite buffer space is available. And best of all, the communication channel between the data link layers never damages or loses frames. This thoroughly unrealistic protocol, which we will nickname "Utopia," .The utopia protocol is unrealistic becauseit doesnothandleeitherflowcontrolorerrorcorrection Stop-and-waitProtocol



It is still very simple. The sender sends one frame and waits for feedback from the receiver. When the ACK arrives, the sender sends the next frame

It is Stop-and-Wait Protocol because the sender sends one frame, stops until it receivesconfirmationfromthereceiver(okaytogoahead),andthensends the next frame. We still have unidirectional communication for data frames,but auxiliary ACK frames (simple tokens of acknowledgment) travel from the other direction. We add flow control to our previous protocol.

### NOISYCHANNELS

Although the Stop-and-Wait Protocol gives us an idea of how to add flowcontrol to its predecessor, noiseless channels are nonexistent. We can ignore the error (as we sometimes do), or we need to add error control to our protocols. We discuss three protocols in this section that use error control. SlidingWindowProtocols: 1 Stop-and-Wait Automatic Repeat Request

2 Go-Back-N Automatic Repeat Request

### 3 SelectiveRepeat AutomaticRepeatRequest

# 1 Stop-and-WaitAutomaticRepeatRequest

To detect and correct corrupted frames, we need to add redundancy bits to our data frame. When the frame arrives at the receiver site, it is checked and if it is corrupted, it is silently discarded. The detection of errors in this protocol is manifested by the silence of the receiver.

Lost frames are more difficult to handle than corrupted ones. In our previous protocols, there was no way to identify a frame. <u>The received framecouldbethecorrectone,oraduplicate,oraframeoutoforder</u>. Thesolutionis to number the frames. When the receiver receives a data frame that is out of order, this means that frames were either lost or duplicated

The lost frames need to be resent in this protocol. If the receiver does not respond when there is an error, how can the sender know which frame to resend? To remedy this problem, the sender keeps a copy of the sent frame. At the sametime, itstartsatimer. If the timer expires and there is no ACK for the sent frame, the frame is resent, the copy is held, and the timer is restarted. Since the protocol uses the stop-and-wait mechanism, there is only one specific frame that needs an ACK

Error correction in Stop-and-Wait ARQ is done by keeping a copy of the sent frame and retransmitting of the frame when the timer expires

## In Stop-and-Wait ARQ, we use sequence numbers to number the frames. The sequence numbers are based on modulo-2 arithmetic.

In Stop-and-Wait ARQ, the acknowledgment number always announces in modulo-2 arithmetic the sequence number of the next frame expected.



# BandwidthDelayProduct:

Assume that, in a Stop-and-Wait ARQ system, the bandwidth of the line is 1 Mbps, and 1 bit takes 20 ms to make a round trip. What is the bandwidth-delay product? If the system data frames are 1000 bits in length, what is the utilization percentage of the link?

$$(1 \times 10^6) \times (20 \times 10^{-3}) = 20,000$$
 bits

The link utilization is only 1000/20,000, or 5 percent. For this reason, for a link with a high bandwidth or long delay, the use of Stop-and-Wait ARQ wastes the capacity of the link.

### 2.Go-Back-NAutomaticRepeatRequest

To improve the efficiency of transmission (filling the pipe), multiple framesmustbeintransitionwhilewaitingforacknowledgment.Inother words,weneedtoletmorethanoneframebeoutstandingtokeepthe channel busy while the sender is waiting for acknowledgment. The first is called Go-Back-N Automatic Repeat. In this protocol we can sendseveral frames before receiving acknowledgments; we keep a copy of these frames until the acknowledgments arrive.

IntheGo-Back-NProtocol, these quence numbers are modulo  $2^m$ , where m is the size of the sequence number field in bits. The sequence numbers range from 0 to 2 power m- 1. For example, if m is 4, the only sequence numbers are 0 through 15 inclusive.



b. Send window after sliding

The **sender window** at any time divides the possible sequence numbers into four regions.

The first region, from the far left to the left wall of the window, defines thesequencenumbersbelongingtoframesthatarealreadyacknowledged. The sender does not worry about these frames and keeps no copies of them.

Thesecondregion, colored in Figure (a), defines the range of sequence numbers belonging to the frames that are sent and have an unknown status. The sender needs to wait to find out if these frames have been received or werelost. We call these outstanding frames.

The third range, white in the figure, defines the range of sequence numbersforframesthatcanbesent; however, the corresponding data packets have not yet been received from the network layer.

Finally, the fourth region defines sequence numbers that cannot be used until the window slides

The send windowis anabstract concept defining animaginaryboxof size 2<sup>*m*</sup>– 1 with three variables: S<sub>f</sub>, S<sub>n</sub>, and S<sub>size</sub>. The variable S*f* defines thesequencenumberofthefirst(oldest)outstandingframe.ThevariableS*n* holdsthesequencenumberthatwillbeassignedtothenextframetobesent. Finally,thevariableSsizedefinesthesizeofthewindow.

Figure (b) shows how a send window can slide one or more slots to the right when an acknowledgment arrives from the other end. The acknowledgments in this protocol are cumulative, meaning that more than one framecanbeacknowledgedbyanACKframe.InFigure,frames0,I,and2are

acknowledged, so the window has slide to the right three slots. Note that the value of *Sf* is 3 because frame 3 is now the first outstanding frame. **The send window can slide one or more slots when a valid acknowledgmentarrives.** 

### **<u>Receiverwindow</u>**:variable*Rn*(receivewindow,nextframeexpected).

The sequence numbers to the left of the window belong to the frames already received and acknowledged; the sequence numbers to the right of this window define the frames that cannot be received. Any received frame with а sequence these two regions is discarded. number in Only frame with а а sequencenumbermatchingthevalueof*Rn*isacceptedandacknowledged. receive The window also slides, but only one slot at a time. When a correct frame is received (and a frame is received only one at a time), the window slides. (see below figure for receiving window)

The receive window is an abstract concept defining an imaginary box of size 1 with one single variable Rn. The window slides when a correct frame has arrived; sliding occurs one slot at a time



b. Window after sliding

Fig:Receiverwindow(beforesliding(a),Aftersliding(b))

## <u>Timers</u>

Although there can be a timer for each frame that is sent, in our protocol we use only one. The reason is that the timer for the first outstanding framealways expires first; we send all outstanding frames when this timer expires.

### <u>Acknowledgment</u>

Thereceiversendsapositiveacknowledgmentifaframehasarrivedsafeand sound and in order. If a frame is damaged or is received out of order, the receiver is silent and will discard all subsequent frames until it receives the one it is of causes expecting. The silence the receiver the timer of the unacknowledgedframeatthesendersidetoexpire.This, inturn, causes the sender to go back and resend all frames, beginning with the one with the expired timer. The receiver does not have to acknowledge frame each received. It can send on ecumulative acknowledgment for several frames.

### <u>ResendingaFrame</u>

When the timer expires, the sender resends all outstanding frames. For example, suppose these nderhas already sent frame 6, but the timer for frame 3 expires. This means that frame 3 has not been acknowledged; the sendergoes back and sends frames 3, 4, 5, and 6 again. That is why the protocol is called *Go-Back-N* ARQ.

#### Belowfigureisanexample(ifacklost) of a case where the forward channel is

reliable, but the reverse is not. No data frames are lost, but some ACKs are delayed and one is lost. The example also shows how cumulative acknowledgmentscanhelpifacknowledgmentsaredelayedorlost



#### Belowfigureisanexample(ifframelost)



Stop-and-Wait ARQ is a special case of Go-Back-N ARQ in which the size of the send window is 1.

# 3SelectiveRepeat AutomaticRepeat Request

In Go-Back-N ARQ, The receiver keeps track of only one variable, and there isno need to buffer out-of- order frames; they are simply discarded. However, this protocol is very inefficient for a noisy link.

In a noisy link a frame has a higher probability of damage, which means the resending of multiple frames. This resending uses up the bandwidth and slows down the transmission.

For noisy links, there is another mechanism that does not resend *N* frameswhen just one frame is damaged; only the damaged frame is resent. This mechanism is called Selective Repeat ARQ.

It is more efficient for noisy links, but the processing at the receiver is more complex.

<u>Sender Window</u>(explain go-back N sender window concept (before& after sliding.) The only difference in sender window between Go-back N and Selective Repeat is Window size)



#### **Receiverwindow**

The receiver window in Selective Repeat is totally different from the one in Go Back-N.First, thesize of the receive window is thesameas thesize of the send window  $(2^{m-1})$ .

TheSelectiveRepeatProtocolallowsasmanyframesasthesizeofthe receiver window to arrive out of order and be kept until there is a set of in- order frames to be delivered to the network layer. Because the sizes of thesend window and receive window are the same, all the frames in the sendframe can arrive out of order and be stored until they can be delivered. However the receiver never delivers packets out of order to the network layer. Above Figure shows the receive window. Those slots inside the window that are colored define frames that have arrived out of order and are waiting for their neighbors to arrive before delivery to the network layer.

In Selective Repeat ARQ, the size of the sender and receiver window must be at most one-half of  $2^{\rm m}$ 

### DeliveryofDatain SelectiveRepeatARQ:



a. Before delivery FlowDiagram

	R <sub>r</sub>	ı		ackNo sent: 3					
0 1 2	3	4	5	6	7	0	1	2	3

b. After delivery



### DifferencesbetweenGo-BackN&SelectiveRepeat

One main difference is the number of timers. Here, each frame sent or resent needs a timer, which means that the timers need to be numbered (0, 1, 2, and 3). The timer for frame 0 starts at the first request, but stops when the ACK for this frame arrives.

There are two conditions for the delivery of frames to the network layer: First, a set of consecutive frames must have arrived. Second, the set startsfrom the beginning of the window. After the first arrival, there was only one frame and it started from the beginning of the window. After the last arrival, there are three frames and the first one starts from the beginning of the window.

Another importantpointisthataNAKissent. ThenextpointisabouttheACKs.NoticethatonlytwoACKsaresenthere. Thefirstoneacknowledgesonlythefirstframe;thesecondoneacknowledges three frames. In Selective Repeat, ACKs are sent when data are delivered to thenetworklayer.Ifthedatabelongingto*n*framesaredeliveredinoneshot, onlyoneACKissentforallofthem.

## Piggybacking

A technique called**piggybacking** is used to improve the efficiency of the bidirectional protocols. When a frame is carrying data from A to B, it can also carrycontrol informationabout arrived(orlost) frames fromB; whena frame is carrying data from B to A, it can also carry control information about the arrived (or lost) frames from A.

# RANDOMACCESSPROTOCOLS

We can consider the data link layer as two sub layers. The upper sub layer is responsible for data link control, and the lower sub layer is responsible for resolving access to the shared media

Data link layer



The upper sub layer that is responsible for flow and error control is called the logical link control (LLC) layer; the lower sub layer that is mostly responsible for multiple access resolution is called the media access control (MAC) layer. When nodes or stations are connected and use a common link, called a multipoint or broadcast link, we need a multiple-access protocol to coordinate access to the link.



# RANDOM ACCESS

In random access or contention methods, no station is superior to another station and none is assigned the control over another.

Twofeaturesgivethismethoditsname.First,thereisnoscheduledtime for a station to transmit. Transmission is random among the stations. That is whythesemethodsarecalled *randomaccess*.Second,norulesspecifywhich station should send next. Stations compete with one another to access the medium.Thatiswhythesemethodsarealsocalled*contention*methods.

# ALOHA

# <u>1PureALOHA</u>

The original ALOHA protocol is called pure ALOHA. This is a simple, but elegant protocol. The idea is that each station sends a frame whenever it has a frameto send. However, since there is only one channel to share, there is the possibility of collision between frames from different stations. Below Figure shows an example of frame collisions in pure ALOHA.



 $\label{eq:InpureALOHA} In {\tt pureALOHA}, the station stransmit frames whenever they have data to send.$ 

- When two or more stations transmit simultaneously, there is collision and the frames are destroyed.
- InpureALOHA, whenever any station transmits a frame, it expects the acknowledgement from the receiver.
- If acknowledgement is not received within specified time, the station assumes that the frame (or acknowledgement) has been destroyed.
- If the frame is destroyed because of collision the station waits for a random amount of time and sends it again. This waiting time must be random otherwise same frames will collide again and again.
- ThereforepureALOHAdictatesthatwhentime-outperiodpasses, eachstation must wait for a random amount of time before resending its frame. This randomness will help avoid more collisions.

<u>Vulnerabletime</u>Letusfindthelengthoftime, the vulnerabletime, in which there is a possibility of collision. We assume that the stations send fixed- length frames with each frame taking TfrS to send. Below Figure shows the vulnerable time for station A.



Station A sends a frame at timet. Now imagine station B has already sent a frame between t - Tfr and t. This leads to a collision between the frames from station A and station B. The end of B's frame collides with the beginning of A's frame. On the other hand, suppose that station C sends a frame between t and t + Tfr . Here, there is a collision between frames from station A and station C. The beginning of C's frame collides with the end of A's frame

LookingatFigure, we see that the vulnerable time, during which a collision may occur in pure ALOHA, is 2 times the frame transmission time. <u>PureALOHA</u> <u>vulnerable time = 2 x Tfr</u>



# ProcedureforpureALOHAprotocol

### Example

ApureALOHAnetworktransmits200-bitframesonasharedchannelof200kbps. What is the requirement to make this frame collision-free?

### Solution

Average frame transmission time Tfr is 200 bits/200 kbps or 1 ms. The vulnerable time is 2 x 1 ms = 2 ms. This means no station should send later than 1 ms before this station starts transmission and no station should start sending during the onel-msperiod that this station is sending.

The throughput for pure ALOHA is  $S = G \times e -2G$ . The maximum throughput Smax = 0.184 when G= (1/2).

#### PROBLEM

A pure ALOHA network transmits 200-bit frames on a shared channel of 200 kbps. What is the throughput if the system (all stations together) produces a. 1000framespersecondb.500framespersecondc.250framespersecond.

Theframetransmissiontimeis200/200kbpsor1ms.

a. If the system creates 1000 frames persecond, this is 1 frameper

millisecond. The load is 1. In this case  $S=G \times e^{-2G}$  or S=0.135 (13.5 percent). This means that the throughput is  $1000 \times 0.135 = 135$  frames. Only 135 frames out of 1000 will probably survive.

- b. If the system creates 500 frames persecond, this is (1/2) frameper millisecond. The load is (1/2). In this case S=G×e<sup>-2G</sup> or S=0.184 (18.4 percent). This means that the throughput is 500×
  - 0.184 = 92andthatonly92framesoutof500willprobablysurvive.Note that this is the maximum throughput case, percentage wise.
- c. If the system creates 250 frames per second, this is (1/4) frame per millisecond. The load is (1/4). In this case S = G × e  $-^{2G}$  or S = 0.152 (15.2 percent). This means that the throughput is 250 × 0.152 = 38. Only 38 frames out of 250 will probably survive.

## 2SlottedALOHA

PureALOHAhasavulnerabletimeof2x*Tfr*. Thisissobecause there is no rule that defines when the station can send. A station may send soon after another station has started or soon before another station has finished. Slotted ALOHA was invented to improve the efficiency of pure ALOHA.

In slotted ALOHA we divide the time into slots of Tfr s and force the station to sendonlyatthebeginningofthetimeslot.Figure3showsanexampleof frame collisions in slotted ALOHA



Because a station is allowed to send only at the beginning of the synchronized time slot, if a station misses this moment, it must wait until the beginning of the next time slot. This means that the station which started at the beginning of this slot has already finished sending its frame. Of course, there is still the possibility of collision if two stations try to send at the beginning of the same time slot. However, the vulnerable time is now reduced to one-half, equal to *Tfr* Figure 4 shows the situation

Below fig shows that the vulnerable time for slotted ALOHA isone-half that of pure ALOHA. Slotted ALOHA vulnerable time = Tfr



The throughput for slotted ALOHA is  $S = G \times e-G$ . The maximum throughput Smax = 0.368 when G = 1.

AslottedALOHAnetworktransmits200-bitframesusingasharedchannelwith a 200-Kbps bandwidth. Find the throughput if the system (all stations together) produces

a.1000framespersecondb.500framespersecondc.250 frames per second

## Solution

This situation is similar to the previous exercise except that the network is usingslottedALOHAinsteadofpureALOHA.Theframetransmissiontimeis 200/200 kbps or 1 ms.

- a. In this case G is 1. So S =G xe-G or S =0.368 (36.8 percent). This means that the throughput is 1000 x 0.0368 =368 frames. Only 368 out of 1000 frameswillprobablysurvive.Notethatthisisthemaximumthroughputcase, percentagewise.
- b. Here G is 1/2 In this case S =G xe-Gor S =0.303 (30.3 percent). Thismeansthatthethroughputis500x0.0303=151.Only151framesoutof 500 will probably survive.
- c. Now G is 1/4. In this case S =G x *e-G* or S =0.195 (19.5 percent). This meansthatthethroughputis250x0.195=49.Only49framesoutof250 will probably survive

ComparisonbetweenPureAloha&SlottedAloha



### CarrierSenseMultipleAccess(CSMA)

To minimize the chance of collision and, therefore, increase the performance, the CSMA method was developed. The chance of collision can be reduced if a station senses the medium before trying to use it. <u>Carrier sensemultiple access (CSMA)requires that each station first listen to the medium (orcheck the state of themedium) before sending. In other words, CSMA is basedon the principle "sensebefore transmit" or "listen before talk."</u>

CSMA can reduce the possibility of collision, but it cannot eliminate it. The reason for this is shown in below Figure. Stations are connected to a shared channel (usually a dedicated medium).

The possibility of collision still exists because of propagation delay; station may sense the medium and find it idle, only because the first bit sent by another station has not yet been received.

At time t/station B senses the medium and finds it idle, so it sends a frame. At time t2 (t2 > t1)'station C senses the medium and finds it idlebecause, at this time, the first bits from station B

have not reached station C. Station C also sends a frame. The two signalscollide and both frames are destroyed.



Space/timemodelofthecollisioninCSMA

## VulnerableTime

<u>The vulnerable time for CSMA is the propagation time Tp</u>. This is the time needed for a signal to propagate from one end of the medium to the other. When a station sends a frame, and any other station tries to send a frameduring this time, a collision will result. But if the first bit of the frame reaches the end of the medium, every station will already have heard the bit and will refrain from sending



# PersistenceMethods

What should a station do if the channel is busy? What should a station do if the channel is idle? Three methods have been devised to answer these questions: the 1-persistent method, the non-persistent method, and the p-persistent method



**1-Persistent:**In this method, after the station finds the line idle, it sends its frame immediately (with probability 1). This method has the highest chance of collision because two or more stations may find the line idle and send their frames immediately.

**Non-persistent:** astationthathasaframetosendsensestheline. If the line is idle, it sends immediately. If the line is not idle, it waits a random amount of time and then senses the line again. This approach reduces the chance of collision because it is unlikely that two or more stations will wait the same amount of time and retry to send simultaneously. However, this methodreduces the efficiency of the network because the medium remains idle when there may be stations with frames to send.

**p-Persistent:** This is used if the channel has time slots with a slot duration equal to or greater than the maximum propagation time. The p-persistent approach combines the advantages of the other two strategies. It reduces the chance of collision and improves efficiency.

Inthismethod, after the station finds the line idle it follows these steps:

- 1. Withprobabilityp, the stations ends its frame.
- 2. With probabilityq=1 -p, the station waits for the beginning of the nexttime slot and checks the line again.
  - a. If the line is idle, it goest ost ep 1.
  - b. If the line is busy, it acts as though a collision has occurred and uses the backoff procedure.



## CarrierSenseMultipleAccess withCollisionDetection(CSMA/CD)

The CSMA method does not specify the procedure following a collision. Carrier sense multiple access with collision detection (CSMA/CD) augments the algorithm to handle the collision.

In this method, a station monitors the medium after it sends a frame to see if the transmission was successful. If so, the station is finished. If, however, thereisacollision,theframeissentagain.

To better understand CSMA/CD, let us look at the first bits transmitted by the two stations involved in the collision. Although each station continues to send bits in the frame until it detects the collision, we show what happens as the first bits collide. In below Figure, stations A and C are involved in the collision.



# CollisionofthefirstbitinCSMA/CD

Attime t1, station A has executed its persistence procedure and starts sending the bits of its frame. At time t2, station C has not yet sensed the first bit sent by A. Station C executes its persistence procedure and starts sending the bits in its frame, which propagate both to the left and to the right. The collision occurs sometime after time t2. Station C detects a collision at time t3 when it receives the first bit of A's frame. Station C immediately (or after a short time, but we assume immediately) aborts transmission.

Station A detects collision at time *t4* when it receives the first bit of C's frame; italsoimmediatelyabortstransmission.Lookingatthefigure,weseethatA

transmitsfortheduration t4-tl; Ctransmitsfortheduration t3-t2.

# MinimumFrameSize

For CSMAICD to work, we need a restriction on the frame size. Before sending the last bit of the frame, the sending station must detect a collision, if any, and abort the transmission. This is so because the station, once the entire frame is sent, does not keep a copy of the frame and does not monitor the line for collision detection. Therefore, the frame transmission time Tfr must be at leasttwo times themaximum propagation time Tp. To understand the reason, let us think about the worst-case scenario. If the two stations involved in a collisionare the maximum distanceapart, the signal from the first takes time Tp to reach the first. So the requirement is that the first station must still be transmitting after 2Tp.



# CollisionandabortioninCSMA/CD



#### PROBLEM

A network using CSMA/CD has a bandwidth of 10 Mbps. If the maximum propagation time (including the delays in the devices and ignoring the time needed to send a jamming signal, as we see later) is 25.6  $\mu$ s, what is the minimum size of the frame?

SOL

The frame transmission time is  $Tfr=2\times Tp=51.2$  µs. This means, in the worstcase, a station needs to transmit for a period of 51.2 µs to detect the collision.Theminimumsizeoftheframeis10Mbps×51.2µs=512bitsor64

 $bytes. This is actually the minimum size of the frame for {\tt Standard Ethernet} \ .$ 

#### DIFFERENCESBETWEENALOHA&CSMA/CD

Thefirst difference is the addition of the persistence process. We need to sense the channel before we start sending the frame by using one of the persistence processes

The second difference is the frame transmission. In ALOHA, we first transmit the entire frame and then wait for an acknowledgment. In *CSMA/CD*, transmissionandcollisiondetectionisacontinuousprocess.Wedonotsend the entire frame and then look for a collision. The station transmits and receives continuouslyand simultaneously The third difference is the sending of a short jamming signal that enforces the collision in case other stations have not yet sensed the collision.

### CarrierSenseMultiple Access with CollisionAvoidance(CSMA/CA)

We need to avoid collisions on wireless networks because they cannot be detected.Carriersensemultipleaccesswithcollisionavoidance(CSMAICA) was invented for wirelesss network. Collisions are avoided through the use of CSMA/CA'sthreestrategies:the<u>interframespace,thecontentionwindow,and</u>



#### TiminginCSMA/CA

#### InterframeSpace(IFS)

First, collisionsareavoidedbydeferringtransmissionevenifthechannel is found idle. When an idle channel is found, the station does not sendimmediately. It waitsfor a period of time called the inter frame space or IFS.

Even though the channel may appear idle when it is sensed, a distant station may have already started transmitting. The distant station's signal has not yet reached this station. The IFS time allows the front of the transmitted signal by the distant station to reach this station. If after the IFS time the channel is still idle, the station can send, but it still needs to wait a time equal to the contention time. The IFS variable can also be used to prioritize stations or frame types. For example, a station that is assigned shorter IFS has a higher priority.

In CSMA/CA, the IFS can also be used to define the priority of a station or aframe.

#### ContentionWindow

The contention window is an amount of time divided into slots. A station that is ready to send chooses a random number of slots as its wait time. The number of slots in the window changes according to the binary exponential backoff strategy. This means that it is set to one slot the first time and then doubles each time the station cannot detect an idle channel after the IFS time. This is very similar to the p-persistent method except that a random outcome defines the number of slots taken by the waiting station.

One interesting point about the contention window is that the station needs to sense the channel after each time slot. However, if the station finds the channel busy, it does not restart the process; it just stops the timer and restarts it when the channel is sensed as idle. This gives priority to the station with the longest waiting time. In CSMA/CA, if the station finds the channel busy, it does not restart the timer of the contention window; it stops the timer and restarts it when the channel becomes idle.

### Acknowledgment

With all these precautions, there still may be a collision resulting in destroyed data. In addition, the data may be corrupted during the transmission. Thepositive acknowledgment and the time-out timer can help guarantee that the receiver has received the frame.



ThisistheCSMAprotocolwithcollision avoidance.

- The station ready to transmit, senses the line by using one of the persistent strategies.
- Assoonasitfindstheline tobeidle,thestationwaitsforanIFS(Interframe space)amountoftime.
- If then waits for some random time and sends the frame.
- After sending the frame, it sets a timer and waits for the acknowledgementfrom the receiver.
- If the acknowledgement is received before expiry of the timer, then the transmission is successful.
- But if the transmitting station does not receive the expected acknowledgementbeforethetimerexpirythenitincrementsthebackoff

# **Controlled AccessProtocols**

In controlled access, the stations seek information from one another to find which station has the right to send. It allows only one node to send at a time, to avoidcollisionofmessagesonsharedmedium. The three controlled-access methods are:

1Reservation2Polling3TokenPassing

# Reservation

- In the reservation method, a station needs to make a reservation before sending data.
- Thetime linehas two kindsof periods:
  - 1. Reservationintervaloffixedtimelength
  - 2. Datatransmissionperiodofvariableframes.
- If there are M stations, the reservation interval is divided into M slots, and each station has one slot.
- Suppose if station 1 has a frame to send, ittransmits 1 bit during the slot
  - 1. No otherstation isallowed totransmitduringthisslot.
- In general, i <sup>th</sup> station may announce that it has a frame to send by insertinga1bitinto i <sup>th</sup>slot.Afterall Nslotshavebeenchecked,eachstation knows which stations wish to transmit.
- The stations which have reserved their slots transfer their frames in that order.
- Afterdatatransmissionperiod, nextreservation interval begins.
- Sinceeveryoneagreesonwhogoesnext,therewillneverbeany collisions.

The following figure shows a situation with five stations and a five slot reservation frame. In the first interval, only stations 1, 3, and 4 have made reservations. In the second interval, only station 1 has made a reservation.



# Polling

- Polling process is similar to the roll-call performed in class. Just like the teacher, a controller sends a message to each node in turn.
- In this, one acts as a primary station(controller) and the others are secondary stations. All data exchanges must be made through the controller.
- Themessagesentbythecontrollercontainstheaddressofthenode being selected for granting access.
- Althoughallnodesreceivethemessagebuttheaddressedoneresponds to it and sends data, if any. If there is no data, usually a "poll reject" (NAK) message is sent back.
- Problems include high overhead of the polling messages and high dependence on the reliability of the controller.



# TokenPassing

- Intokenpassingscheme, the stations are connected logically to each other in form of ring and access of stations is governed by tokens.
- Atokenisaspecialbitpatternorasmallmessage,whichcirculatefrom one station to the next in the some predefined order.

- In Token ring, token is passed from one station to another adjacent station intheringwhereasincaseofTokenbus, each station uses the bus to send the token to the next station in some predefined order.
- Inbothcases, token represents permission to send. If a station has a frame queued for transmission when it receives the token, it can send that frame before it passes the token to the next station. If it has no queued frame, it passes the token simply.
- After sending a frame, each station must wait for all N stations (including itself)tosendthetokentotheirneighborsandtheotherN-1stationstosend a frame, if they have one.
- There exists problems like duplication of token or token is lost or insertion of new station, removal of a station, which need be tackled for correct and reliable operation of this scheme.



# ErrorDetection

### Error

A condition when the receiver's information does not matches with the sender's information. During transmission, digital signals suffer from noise that can introduce errors in the binary bits travelling from sender to receiver. That meansa0bitmaychangeto1ora1bitmaychangeto0.

ErrorDetectingCodes(ImplementedeitheratDatalinklayeror Trans						port						
			Layer			of		OSI	OSI		Model)	
Whe	Whenever a message is transmitted, it may get scrambled by noise or data may							ay				
get	corrupted.	То	avoid	this,	we	use	error-o	detecting	codes	which	are	
addi	additionaldataaddedtoagivendigitalmessagetohelpusdetectifanyerror											
hasoccurredduringtransmissionofthemessage.												
Basi	capproachus	edfor	errorde	etectio	onisth	euse	ofredun	dancybits,	where			

additional bits are added to facilitate detection of errors. Some popular techniques for error detection are:

- 1. SimpleParitycheck
- 2. Two-dimensionalParitycheck
- 3. Checksum
- 4. Cyclicredundancycheck

## SimpleParitycheck

Blocks of data from the source are subjected to a check bit or parity bit generator form, where a parity of :1 is added to the block if it contains odd number of 1's, and

0isaddedifitcontainsevennumberof1's

This scheme makes the total number of 1's even, that is why it is called even parity checking.



### Two-dimensionalParitycheck

Parity check bits are calculated for each row, which is equivalent to a simple parity check bit. Parity check bits are also calculated for all columns, then both are sent along with the data. At the receiving end these are compared with the parity bits calculated on the received data.

10011001	11100010	0010	0100	10000100
			Row p	oarities
	10011	001	0	
	11100	010	0	
	00100	100	0	
	10000	100	0	
Column	11011	011	0	



### Checksum

- Inchecksumerrordetectionscheme, the data is divided intoksegments each of m bits.
- In the sender's end the segments are added using 1's complement arithmetictogetthesum. Thesumiscomplementedtogetthechecksum.
- Thechecksumsegmentissentalongwiththedatasegments.
- At the receiver's end, all received segments are added using 1's complement arithmetic to get the sum. The sum is complemented.
- If the result is zero, the received data is accepted; otherwise discarded.

	Original L	Data		
	10011001	11100010	00100100	10000100
	1	2	З	4
	k=4, m=8			Reciever
	Sender		1	10011001
1	1001100	1	2	11100010
2	1110001	0	1	01111011
	1)0111101	1		01111100
	0111110	0	3	00100100
3	0010010	0		10100000
	1010000	0	4	10000100
4	1000010	0	(1)	00100100
(	1)0010010	0	9	1
C	4	1		00100101
Sum:	0010010	1	1.00	11011010
CheckSun	.1101101	0	Sum:	11111111
CHECKJUH		Co	mplement:	00000000
		Co	onclusion: A	ccept Data

Cyclicredundancycheck(CRC)



- Unlikechecksumscheme, which is based on addition, CRC is based on binary division.
- InCRC, as equence of redundant bits, called cyclic redundancycheck bits, are appended to the end of data unit so that the resulting data unit becomes exactly divisible by a second, predetermined binary number.
- Atthedestination, the incoming data unit is divided by the same number. If at this step there is no remainder, the data unit is assumed to be correct and is therefore accepted.
- Aremainderindicatesthatthedataunithasbeendamagedintransitand thereforemustberejected.

# ErrorCorrection

ErrorCorrectioncodesareusedtodetectandcorrecttheerrorswhendatais transmitted from the sender to the receiver.

ErrorCorrectioncanbehandledintwoways:

Backwarderrorcorrection:Oncetheerrorisdiscovered, thereceiver requests the sender to retransmit the entire data unit.

Forwarderrorcorrection: Inthiscase, thereceiveruses the error-correcting codewhich automatically corrects the errors.

Asingleadditionalbitcandetecttheerror, butcannot correctit.

For correcting the errors, one has to know the exact position of the error. For example, If we want to calculate a single-biterror, the error correction code will determine which one of seven bits is in error. To achieve this, we have to add some additional redundant bits.

Suppose r is the number of redundant bits and d is the total number of the data bits. The number of redundant bits r can be calculated by using the formula:  $2^{r} >= d+r+1$ 

The value of r is calculated by using the above formula. For example, if thevalue of d is 4, then the possible smallest value that satisfies the above relation would be 3.

Todeterminethepositionofthebitwhichisinerror, atechniquedevelopedby R.W Hamming is Hamming code which can be applied to any length of the dataunit and uses the relationship between data units and redundant units.

# HammingCode

Paritybits: The bit which is appended to the original data of binary bits so that the total number of 1 sise venorodd.

Even parity: Tocheckforevenparity, if the total number of 1s is even, then the value of the parity bit is 0. If the total number of 1s occurrences is odd, then the value of the parity bit is 1.

OddParity:Tocheckforoddparity,ifthetotalnumberof1siseven,thenthe value of parity bit is 1. If the total number of 1s is odd, then the value of parity bit is 0.

## AlgorithmofHammingcode:

Aninformationof'd'bitsareaddedtotheredundantbits'r'toformd+r.

The location of each of the (d+r) digits is assigned a decimal value.

The'r'bitsareplacedinthepositions1,2, ......2<sup>k-1</sup>

At the receiving end, the parity bits are recalculated. The decimal value of the parity bits determines the position of an error.

Relationshipb/wErrorposition&binarynumber.

Error Position	Binary Number
0	000
1	001
2	010
3	011
4	100
5	101
6	110
7	111

 $\label{eq:label} Let's understand the concept of Hamming code through an example:$ 

Suppose theoriginal data is 1010 which is to be sent.

Totalnumberofdatabits'd'=4 Numberofredundantbitsr:2<sup>r</sup>>=d+r+1
# $2^{r} > = 4 + r + 1$

Therefore, the value of ris3 that satisfies the above relation. Total number of bits = d+r=4+3=7;

#### Determiningthepositionoftheredundantbits

The number of redundant bits is 3. The three bits are represented by r1, r2, r4. The position of the redundant bits is calculated with corresponds to the raised power of 2. Therefore, their corresponding positions are 1,  $2^1$ ,  $2^2$ . The position of r1=1. The position of r2=2. The position of r4=4

RepresentationofDataontheadditionofparitybits:



# DeterminingtheParity bits

Determining the r1 bit: The r1 bit is calculated by performing a parity check on the bit positions whose binary representation includes 1 in the first position.



We observe from the above figure that the bit position that includes 1 in thefirst position are 1, 3, 5, 7. Now, we perform the even-parity check at these bit positions. The total number of 1 at these bit positions corresponding to r1 is even, therefore, the value of the r1 bit is 0.

Determining r2 bit: The r2 bit is calculated by performing a parity check on the bit positions whose binary representation includes 1 in the second position



We observe from the above figure that the bit positions that includes 1 in the secondpositionare2,3,6,7.Now,weperformtheeven-paritycheckatthese

bit positions. The total number of 1 at these bit positions corresponding to r2 is odd, therefore, the value of the r2 bit is 1.

Determining r4 bit: The r4 bit is calculated by performing a parity check on the bit positions whose binary representation includes 1 in the third position.



We observe from the above figure that the bit positions that includes 1 in the third position are 4, 5, 6, 7. Now, we perform the even-parity check at these bit positions. The total number of 1 at these bit positions corresponding to r4 is even, therefore, the value of the r4 bit is 0.

#### Datatransferredisgivenbelow:



Suppose the 4th bit is changed from 0 to 1 at the receiving end, then parity bits are recalculated.

#### R1bit



We observe from the above figure that the binary representation of r1 is 1100. Now, we perform the even-parity check, the total number of 1s appearing in the r1 bit is an even number. Therefore, the value of r1 is 0.

#### R2bit

Thebitpositionsofr2bitare2,3,6,7.



We observe from the above figure that the binary representation of r2 is 1001. Now, we perform the even-parity check, the total number of 1s appearing in the r2 bit is an even number. Therefore, the value of r2 is 0.

#### R4bit

Thebitpositionsofr4bitare4,5,6,7.



We observe from the above figure that the binary representation of r4 is 1011. Now, we perform the even-parity check, the total number of 1s appearing in the r4 bit is an odd number. Therefore, the value of r4 is 1.

The binary representation of redundant bits, i.e., r4r2r1 is 100, and its corresponding decimal value is 4. Therefore, the error occurs in a 4th bit position. The bit value must be changed from 1 to 0 to correct the error.

#### WiredLANs:Ethernet

In 1985, the Computer Society of the IEEE started a project, called Project802, to set standards to enable intercommunication among equipment from a variety of manufacturers. Project 802 is a way of specifying functions of the physical layer and the data link layer of major LAN protocols.

The relationship of the 802 Standard to the traditional OSI model is shown in below Figure. The IEEE has subdivided the data link layer into two sub layers: logical link control (LLC) and media access control).

IEEE has also created several physical layer standards for different LANprotocols



IEEEstandard forLANs

# STANDARDETHERNET

TheoriginalEthernetwascreatedin1976atXerox'sPaloAltoResearch Center (PARC). Since then, it has gone through four generations.

StandardEthernet(lOMbps),FastEthernet(100Mbps),GigabitEthernet(l Gbps), and Ten-Gigabit Ethernet (lO Gbps),

 $We briefly discuss the {\tt Standard} (or traditional) {\tt Ethernet} in this section$ 



Ethernetevolutionthroughfourgenerations

# **MACSublayer**

In Standard Ethernet, the MAC sublayer governs the operation of the access method. It also frames data received from the upper layer and passes them to the physical layer.

# <u>FrameFormat</u>

The Ethernet frame contains seven fields: preamble, SFD, DA, SA, length or typeofprotocoldataunit(PDU),upper-layerdata,andtheCRC.Ethernetdoes not provide any mechanism for acknowledging received frames, making it what is known as an unreliable medium. Acknowledgments must be implemented at the higher layers. The format of the MAC frame is shown in below figure



#### 802.3MACframe

<u>Preamble.</u> The first field of the 802.3 frame contains 7 bytes (56 bits) of alternating 0s and 1s that alerts the receiving system to the coming frame and enables it to synchronize its input timing. The pattern provides only an alertand a timing pulse. The 56-bit pattern allows the stations to miss some bits at the beginning of the frame. The preamble is actually added at the physicallayer and is not (formally) part of the frame.

Start frame delimiter (SFD). The second field (l byte: 10101011) signals the beginning of the frame. The SFD warns the station or stations that this is the lastchancefor synchronization. The last2bits is 11 and alerts there ever that the next field is the destination address.

<u>Destination address (DA)</u>. The DA field is 6 bytes and contains the physical address of the destination station or stations to receive the packet.

<u>Source address (SA).</u> The SA field is also 6 bytes and contains the physical address of the sender of the packet.

Length or type. This field is defined as a type field or length field. The original Ethernet used this field as the type field to define the upper-layer protocol usingtheMACframe.TheIEEEstandarduseditasthelengthfieldtodefine the number of bytes in the data field. Both uses are common today.

<u>Data</u>. This fieldcarriesdataencapsulated from the upper-layerprotocols. It isa minimum of 46 and a maximum of 1500 bytes.

<u>CRC</u>. The last field contains error detection information, in this case a CRC - 32

# <u>FrameLength</u>

Ethernet has imposed restrictions on both the minimum and maximum lengthsof a frame, as shown in below Figure



#### Minimumandmaximumlengths

An Ethernet frame needs to have a minimum length of 512 bits or 64 bytes. Part of this length is the header and the trailer. If we count 18 bytes of header and trailer (6 bytes of source address, 6 bytes of destination address, 2 bytes of length or type, and 4 bytes of CRC), then the minimum length of data fromtheupperlayeris64-18 =46bytes.Iftheupper-layerpacketisless than 46 bytes, padding is added to make up the difference

The standard defines the maximum length of a frame (without preamble and SFD field) as 1518 bytes. If we subtract the 18 bytes of header and trailer,

themaximumlengthofthepayloadis1500bytes.

The maximum length restriction has two historical reasons.

First, memory was very expensive when Ethernet was designed: a maximum length restriction helped to reduce the size of the buffer.

Second, the maximum length restriction prevents one station frommonopolizing the shared medium, blocking other stations that have data to send.

# <u>Addressing</u>

TheEthernetaddressis6bytes(48bits),normallywritteninhexadecimal notation, with a colon between the bytes.

 ${\it Example of an Ethern et address in hexadecimal notation}$ 

06:01:02:01:2C:4B 6 bytes = 12 hex digits = 48 bits

<u>Unicast, Multicast, and Broadcast Addresses</u> A source address is always a unicast address-the frame comes from only one station. The destination address, however, can be**unicast, multicast, or broadcast**. Below Figure showshowtodistinguishaunicastaddressfromamulticast address.

If the least significant bit of the first by teinadestination address is 0, the address is unicast; otherwise, it is multicast.



# Unicastand multicast addresses

A unicast destination address defines only one recipient; the relationship between the sender and the receiver is one-to-one.

A multicast destination address defines a group of addresses; the relationship between the sender and the receivers is one-to-many.

The broadcast address is a special case of the multicast address; the recipients are all the stations on the LAN. A broadcast destination address is forty-eight1s. <u>AccessMethod</u>:CSMA/CD

Standard Ethernet uses I-persistent CSMA/CD

Slot Time In an Ethernet network.

Slottime=round-triptime+timerequiredtosendthejamsequence

 $The slottime in {\tt Ethernet} is defined in bits. \\ {\tt It is the time required for a station}$ 

tosend512bits. This means that the actual slot time depends on the data rate; for traditional 10-Mbps Ethernet it is 51.2 micro sec.

<u>Slot Time and Maximum Network Length</u> There is a relationship between theslot time and the maximum length of the network (collision domain). It is dependent on the propagation speed of the signal in the particular medium.

Inmosttransmissionmedia, the signal propagates at  $2 \times 10^8$  m/s (two-thirds of the rate for propagation in air).

FortraditionalEthernet,wecalculate MaxLength =PropagationSpeedx (SlotTime/2) MaxLength=  $(2 \times 10^8) \times (51.2 \times 10^{-6})/2 = 5120$ m

Ofcourse, we need to consider the delay times in repeaters and interfaces, and the time required to send the jam sequence. These reduce the maximum- length of a traditional Ethernet network to 2500 m, just 48 percent of the theoretical calculation. MaxLength=2500 m

# <u>PhysicalLayer</u>

The Standard Ethernet defines several physical layer implementations; four of the most common, are shown in Figure



# **EncodingandDecoding**

Allstandardimplementationsusedigitalsignaling(baseband)at10Mbps.At the sender, data are converted to a digital signal using the Manchester scheme; at the receiver, the received signal is interpreted as Manchester and decoded into data. Manchester encoding is self-synchronous, providing a transition ateach bit interval. Figure shows the encoding scheme for Standard Ethernet



In Manchester encoding, the transition at the middle of the bit is used for synchronization



# IOBase5:ThickEthernet

The first implementation is called**10Base5**, thick Ethernet, or Thicknet. lOBase5 was the first Ethernet specification to use a bus topology with an externaltransceiver(transmitter/receiver) connected via a tap to a thick coaxial cable. Figure shows a schematic diagram of a lOBase5 implementation



# 10Base2:ThinEthernet

The second implementation is called 10 Base2, **thin** Ethernet, or Cheapernet. 10Base2 also uses a bus topology, but the cable is much thinner and more flexible. Figure shows the schematic diagram of a 10Base2 implementation.



thincoaxialcableislessexpensivethanthickcoaxial. Installationissimplerbecausethethincoaxialcableisveryflexible. However,thelengthofeachsegmentcannotexceed185m(closeto200m) due to the high level of attenuation in thin coaxial cable.

# 10Base-T: Twisted-PairEthernet

The third implementation is called 10Base-T or twisted-pair Ethernet. It uses a physical star topology. The stations are connected to a hub via two pairs of twisted cable, as shown in Figure

The maximum length of the twisted cable here is defined as 100 m, to minimize the effect of attenuation in the twisted cable



10Base-Timplementation

Although there are several types of optical fiber 10-Mbps Ethernet, the most commoniscalled10Base-F.10Base-Fusesastartopologytoconnectstations to a hub. The stations are connected to the hub using two fiber-optic cables, as shown in Figure



# UNIT-III

# NetworkLayerDesignIssues

- 1. Store-and-forwardpacketswitching
- 2. Servicesprovidedtotransportlayer
- 3. Implementationofconnectionlessservice
- 4. Implementationofconnection-orientedservice
- 5. Comparisonofvirtual-circuitanddatagramnetworks

# 1 Store-and-forwardpacketswitching



A host witha packetto send transmits it to the nearestrouter, either onits ownLANor over a pointto-point link to the ISP. The packet is stored there until it has fully arrived and the link has finished its processing by verifying the checksum. Then it is forwarded to the next router along the path until it reaches the destination host, where it is delivered. This mechanism is store-andforward packet switching.

# 2 Servicesprovidedtotransportlayer

The network layer providesservicestothetransportlayeratthenetwork layer/transport layer interface. The services need to be carefully designed with the following goals in mind:

- 1. Services independent of router technology.
- 2. Transportlayershieldedfromnumber,type,topologyofrouters.
- 3. Networkaddressesavailabletotransportlayeruseuniformnumberingplan
  - evenacrossLANsandWANs

# 3 Implementationofconnectionlessservice

If connection less service is offered, packets are injected into the network individually and routed independently of each other. No advances setup is needed. In this context, the packets

arefrequentlycalled **datagrams** (in analogy with telegrams) and the network is called a **datagramnetwork**.

	Router		1	
Process P1	Packet	B 2 C 2 2 2 2 2 2 2 2 2 2 2 2 2		P2 H2
A'stable(initially)A's	table(later)C'sTa	ble	E'sTable	
A ⊠ B B C C D B E C F C	A⊠BBCCDBEDFD	AABAC⊠DEEEFE	ACBDCCDDE⊠FF	

Let us assume for this example that the message is four times longer than the maximum packetsize, so the networklayer hasto breakitintofour packets, 1, 2, 3, and 4, and sendeach of them in turn to router *A*.

Every router has an internal table telling it where to send packets for each of the possible destinations. Each table entry is a pair(destination and the outgoing line). Only directly connected lines can be used.

 $\label{eq:alpha} A's initial routing table is shown in the figure under the label ``initially.''$ 

At *A*, packets 1, 2, and 3 are stored briefly, having arrived on the incoming link. Then each packet is forwarded according to *A*'s table, onto the outgoing link to *C* within a new frame. Packet 1 is then forwarded to *E* and then to *F*.

However, something different happens to packet 4. When it gets to A it is sent to router B, even though it is also destined for F. For some reason (traffic jam along ACE path), A decided to send packet 4 via a different route than that of the first three packets. Router A updated its routing table, as shown under the label "later."

The algorithm that manages the tables and makes the routing decisions is called the **routing algorithm**.

#### 4 Implementationofconnection-orientedservice



If connection-oriented service is used, a path from the source router all the way to the destination router mustbeestablished beforeany data packetscanbesent. Thisconnectionis called a **VC** (virtual circuit), and the network is called a virtual-circuit network

When a connection is established, a route from the source machine to the destinationmachine is chosenas partof the connectionsetupandstored intables inside the routers. That route is used for all traffic flowing over the connection, exactly the same way that the telephone system works. When the connection is released, the virtual circuit is also terminated. With connectionoriented service, each packet carries an identifier telling which virtual circuit it belongs to.

As an example, consider the situation shown in Figure. Here, host H1 has established connection 1 with host H2. This connection is remembered as the first entry in each of the routing tables. The first line of A's table says that if a packet bearing connection identifier 1 comes in from H1, it is to be sent to router C and given connection identifier 1. Similarly, the first entry at C routes the packet to E, also with connection identifier 1.

Now let us consider what happens if *H3* also wants to establish a connection to *H2*. It chooses connection identifier 1 (because it is initiating the connection and this is its only connection) and tells the network to establish the virtual circuit.

This leads to thesecond row in thetables. Note that we haveaconflict here becausealthough *A* can easily distinguish connection 1 packets from *H1* from connection 1 packets from *H3, C* cannot do this. For this reason, *A* assigns a different connection identifier to the outgoing traffic for the second connection. Avoiding conflicts of this kind is why routers need the ability to replace connection identifiers in outgoing packets.

In some contexts, this process is called **label switching**. An example of a connection-oriented network service is **MPLS** (**Multi Protocol Label Switching**).

Issue	Datagram network	Virtual-circuit network
Circuit setup	Not needed	Required
Addressing	Each packet contains the full source and destination address	Each packet contains a short VC number
State information	Routers do not hold state information about connections	Each VC requires router table space per connection
Routing	Each packet is routed independently	Route chosen when VC is set up; all packets follow it
Effect of router failures	None, except for packets lost during the crash	All VCs that passed through the failed router are terminated
Quality of service	Difficult	Easy if enough resources can be allocated in advance for each VC
Congestion control	Difficult	Easy if enough resources can be allocated in advance for each VC

# 5 Comparisonofvirtual-circuitanddatagramnetworks

# **RoutingAlgorithms**

ThemainfunctionofNL(NetworkLayer)isroutingpacketsfromthesourcemachinetothe destination machine.

Therearetwoprocessesinsiderouter:

- a) One of them handles each packet as it arrives, looking up the outgoing line to use for it inthe routing table. This process is forwarding.
- b) The other process is responsible for filling inand updating the routing tables. That is where the routing algorithm comes into play. This process is routing.

Regardless of whether routes are chosen independently for each packet or only when new connections are established, certain properties are desirable in a routing algorithm correctness, simplicity, robustness, stability, fairness, optimality

Routingalgorithmscanbegroupedintotwomajorclasses:

- 1) nonadaptive(StaticRouting)
- 2) adaptive.(DynamicRouting)

<u>Nonadaptive algorithm</u> do not base their routing decisions on measurements or estimates of the current traffic and topology. Instead, the choice of the route to use to get from I to J is computed in advance, off line, and downloaded to the routers when the network is booted. This procedure is sometimes called <u>static routing</u>.

<u>Adaptivealgorithm</u>, incontrast, change their routing decisions to reflect changes in the topology, and usually the traffic as well.

Adaptivealgorithmsdifferin

- 1) Where they get their information (e.g., locally, from adjacentrouters, or from all routers),
- 2) Whentheychangetheroutes(e.g.,every∆Tsec,whentheloadchangesorwhenthe topologychanges),and

3) What metric is used for optimization (e.g., distance, number of hops, or estimated transittime).

Thisprocedure is called dynamic routing

# DifferentRoutingAlgorithms

- Optimalityprinciple
- Shortestpathalgorithm
- Flooding
- Distancevectorrouting
- Linkstaterouting
- HierarchicalRouting

# **TheOptimalityPrinciple**

One can make a general statement about optimal routes without regard to network topology or traffic. This statement is known as the <u>optimality principle</u>.

It states that if router J is on the optimal path from router I to router K, then the optimal path from J to K also falls along the same

As a direct consequence of the optimality principle, we can see that the set of optimal routes from all sources to a given destination form a tree rooted at the destination. Such a tree is called a **sink tree**. The goal of all routing algorithms is to discover and use the sink trees for all routers



# ShortestPathRouting(Dijkstra's)

The idea is to build a graph of the subnet, with each node of the graph representing a routerand each arc of the graph representing a communication line or link.

To choose a route between a given pair of routers, the algorithm just finds the shortest pathbetween them on the graph

- 1. Start with the local node (router) as the root of the tree. Assigna cost of 0 to this node and make it the first permanent node.
- 2. Examineeachneighborofthenodethatwasthelastpermanentnode.
- 3. Assignacumulativecosttoeachnodeandmakeittentative
- 4. Amongthelistoftentativenodes
  - a. FindthenodewiththesmallestcostandmakeitPermanent
  - b. If a nodecanbereached frommore thanoneroutethen select theroutewiththe shortest cumulative cost.
- 5. Repeatsteps2to4untileverynodebecomespermanent



# Execution of Dijkstra's algorithm



Iteration	Permanent	tentative	D <sub>2</sub>	D <sub>3</sub>	$D_4$	D <sub>5</sub>	D <sub>6</sub>
Initial	{1}	{2,3,4}	3	2 🗸	5	×	×
1	{1,3}	{2,4,6}	3 🗸	2	4	×	3
2	{1,2,3}	{4,6,5}	3	2	4 `	7	3 🗸
3	{1,2,3,6}	{4,5}	3	2	4 🗸	5	3
4	{1,2,3,4,6}	{5}	3	2	4	5 🗸	3
5	{1,2,3,4,5,6}	{}	3	2	4	5	3

# **Flooding**

- Another static algorithm is flooding, in which every incoming packet is sent out on every outgoing line except the one it arrived on.
- Floodingobviously generates vast numbers of duplicate packets, infact, an infinite number unless some measures are taken to damp the process.
- One such measure is to have a <u>hop counter</u> contained in the header of each packet, which is decremented at each hop, with the packet being discarded when the counter reaches zero. Ideally, the hop counter should be initialized to the length of the path from source to destination.
- A variation of flooding that is slightly more practical is <u>selective flooding</u>. In this algorithm
  the routers do not send every incoming packet out on every line, only on those lines that are
  going approximately in the right direction.
- Floodingisnotpracticalinmostapplications.

# Intra-andInterdomainRouting

Anautonomous system(AS) is a group of networks and routers under the authority of a single administration.

Routinginsideanautonomoussystemisreferredtoasintradomainrouting.(DISTANCEVECTOR, LINK STATE)

Routing between autonomous systems is referred to as inter domain routing. (PATH VECTOR) Each autonomous system can choose oneor moreintradomain routingprotocols to handle routinginsidetheautonomoussystem. However, only one interdomain routing protocol handles routing between autonomous systems.



#### **DistanceVectorRouting**

In distance vector routing, the least-cost route between any two nodes is the <u>route</u> <u>withminimum distance</u>. In this protocol, as the name implies, each node maintains a vector(table) of minimum distances to every node.

Mainly3thingsinthis

Initialization Sharing Updating

#### Initialization

Each node can know only the distance between itself and its immediate neighbors, those directlyconnected to it. So forthe moment, we assume thateach node can senda message to the immediate neighbors and find the distance between itself and these neighbors. Below fig shows the initial tables for each node. The distance for any entry that is not a neighbor is marked as infinite (unreachable).

#### Initializationoftablesindistancevectorrouting



# <u>Sharing</u>

The whole idea of distance vector routing is the sharing of information between neighbors. Although node A does not know about node E, node C does. So if node C shares its routing table with A, node A can also know how to reach node E. On the other hand, node C does not know how to reach node D, but node A does. If node A shares its routing table with node C, node C also knows how to reach node D. In other words, nodes A and C, as immediate neighbors, can improve their routing tables if they help each other.

NOTE: In distance vector routing, each node shares its routing table with its immediate neighbors periodically and when there is a change

# Updating

Whenanodereceivesatwo-columntablefromaneighbor, it needs to update its routing table. Updating takes three steps:

1. Thereceiving nodeneeds to add the cost between itself and the sending node to each value in the second column. (x+y)

2. If the receiving node uses information from any row. The sending node is the next node in the route.

3. The receiving node needs to compare each row of its old table with the corresponding row of the modified version of the received table.

- a. If the next-node entry is different, there ceiving node chooses the row with the smaller cost. If there is a tie, the old one is kept.
- b. If the next-node entry is the same, the receiving node chooses the new row.

For example, suppose node C has previously advertised a route to node X with distance 3. Suppose that now there is no path between C and X; node C now advertises this route with a distance of infinity. Node A must notignore this value even though its oldentry is smaller. The old route does not exist anymore. The new route has a distance of infinity.



# **FinalDiagram**



# Whento Share

The questionnowis, Whendoesa nodesendits partial routing table (only two columns) to all its immediateneighbors? The table is sent both <u>periodically and when there is a change</u> in the table. <u>PeriodicUpdate</u> A nodesends its routing table, normally every 30 s, in a periodic update. The period depends on the protocol that is using distance vector routing.

<u>TriggeredUpdate</u>Anodesendsitstwo-columnroutingtabletoitsneighborsanytimethereis a change in its routing table. This is called a triggered update. The change can result from the following.

1. Anodereceivesatable from a neighbor, resulting inchanges inits own table after updating.

2. A node detects some failure in the neighboring links which results in a distance change toinfinity.



# Two-nodeinstability

# **SOLUTIONSFORINSTABILITY**

1. **Defining Infinity:** redefine infinity to a smaller number, such as 100. For our previous scenario, the system will be stable in less than 20 updates. As a matter of fact, most implementations of the distance vector protocold fine the distance between each node to

be 1 and define 16 as infinity. However, this means that the distance vector routing cannot beusedinlargesystems. Thesize of the network, in each direction, cannot exceed 15 hops.

- 2. **Split Horizon:**In this strategy, instead of flooding the table through each interface, each node sends **only part of its table** through each interface. If, according to its table, node B thinks thatthe optimumroute to reachXis via A,it does not need to advertise this piece of information to A; the information has come from A (A already knows). Taking information from node A, modifying it, and sending it back to node A creates the confusion. In our scenario, node B eliminates the last line of its routing table before it sends it to A. In this case, node A keeps the value of infinity as the distance to X. Later when node A sends its routing table to B, node B also corrects its routing table. The system becomes stable after the first update: both node A and B know that X is not reachable.
- 3. **Split Horizon and Poison Reverse** Using the split horizon strategy has one drawback. Normally, the distance vector protocol uses a timer, and if there is no news about a route, the node deletes the route from itstable. When nodeB in the previousscenario eliminates the route to X from its advertisement to A, node A cannot guess that this is due to the split horizon strategy (the source of information was A) or because B has not received any news about X recently. The split horizon strategy can be combined with the poison reverse strategy. Node B can still advertise the value for X, but if the source of information is A, it can replace the distance with infinity as a warning: "Do not use this value; what I know about this route comes from you."

#### TheCount-to-InfinityProblem

Α	В	С	D	E		Α	В	С	D	E	
•	•	•	•	•	Initially	•	1	2	3	4	Initially
	1	•	•	•	After 1 exchange		3	2	3	4	After 1 exchange
	1	2	•	•	After 2 exchanges		3	4	3	4	After 2 exchanges
	1	2	3	•	After 3 exchanges		5	4	5	4	After 3 exchanges
	1	2	3	4	After 4 exchanges		5	6	5	6	After 4 exchanges
					-		7	6	7	6	After 5 exchanges
							7	8	7	8	After 6 exchanges
								:			-
							•	•	•	•	
			(a)					(b)			

# **LinkStateRouting**

Link state routing is based on the assumption that, although the global knowledge about the topology is not clear, each node has partial knowledge: it knows the state (type, condition, and cost)ofitslinks. **Inotherwords,thewholetopologycanbecompiledfromthe** partialknowledge of each node



# **BuildingRoutingTables**

- 1. Creationofthestatesofthelinksbyeachnode,calledthelinkstatepacket(LSP).
- 2. DisseminationofLSPstoeveryotherrouter, called **flooding, inanefficientand** reliableway.
- 3. Formationofashortestpathtreeforeachnode.
- 4. Calculationofaroutingtablebasedontheshortestpath tree
- 1. **Creation of Link State Packet (LSP)** A link state packet can carry a large amount of information. For the moment, we assume that it carries a minimum amount of data: the node identity,thelistof links,asequence number,andage.Thefirsttwo,nodeidentityand the list of links, are needed to make the topology. The third, sequence number, facilitates flooding and distinguishes new LSPs from oldones. The fourth, age, prevents old LSPs from remaining in the domain for a long time.

LSPsaregeneratedontwooccasions:

1. When there is a change in the topology of the domain

2. on a periodic basis: The period in this case is much longer compared to distance vector. The timer set for periodic dissemination is normally in the range of **60 min or 2 h** based on the implementation. A longer periodensures thatflooding does not create too much traffic on the network.

- II. **FloodingofLSPs:**AfteranodehaspreparedanLSP,itmustbedisseminatedtoallother nodes, not only to its neighbors. The process is called flooding and based on the following
  - 1. ThecreatingnodesendsacopyoftheLSPoutofeachinterface

2. A node that receives an LSP compares it with the copy it may already have. If the newlyarrivedLSP isolder thantheone it has(foundbycheckingthesequencenumber), it discards the LSP. If it is newer, the node does the following:

a. ItdiscardstheoldLSPandkeepsthenewone.

b. It sends a copy of it out of each interface except the one from which the packet arrived. This guarantees that flooding stops somewhere in the domain (where a node has only one interface).

# III. FormationofShortestPathTree:DijkstraAlgorithm

A shortest path tree is a tree in which the path between the root and every other node is the shortest.

The Dijkstra algorithm creates a shortest path tree from a graph. The algorithm divides the nodes into two sets: **tentativeandpermanent.** It finds the neighbors of a current node, makes them tentative, examines them, and if they pass the criteria, makes them permanent.





#### IV. Calculationofaroutingtable

#### routingtablefornodeA

Node	Cost	Next Router
А	0	
В	5	
С	2	
D	3	
E	6	С

#### **PathVectorRouting**

Distance vector and link state routing are both intra domain routing protocols. They can be used inside an autonomous system, but not between autonomous systems. These two protocols are notsuitable for inter domain routing mostly because of scalability. Both of these routing protocols become intractable when the domain of operation becomes large. **Distance vector routing subjectto instability** in thedomain of operation. **Link state routing needs a** 

**huge amount of resources** to calculate routing tables. It also creates heavy traffic because of flooding. There is a need for a third routing protocol which we call path vector routing.

Path vector routing proved to be useful for inter domain routing. The principle of path vector routing is similar to that of distance vector routing. In path vector routing, we assume that there is one node (there can be more, but one is enough for our conceptual discussion) in each AS that acts on behalf of the entire AS. Let us call it the speaker node. The speaker node in an AS creates a routing table and advertises it to speaker nodes in the neighboring ASs. The idea is the same as for distance vector routing except that only speaker nodes in each AS can communicate with each other. However, what is advertised is different. A speaker node advertises the path, not the metric of the nodes, in its autonomous system or other autonomous systems

#### **Initialization**



#### **Sharing**

Just as in distance vector routing, in path vector routing, a speaker in an autonomous system sharesitstablewithimmediateneighbors.InFigure,nodeA1sharesitstablewithnodesB1

and C1. Node C1 shares its table with nodes D1, B1, and A1. Node B1 shares its table with C1 and A1. Node D1 shares its table with C1.

Dest.	Path	Dest.	Path	Dest.	Path	[	Dest.	Path
A1	AS1	A1	AS2-AS1	A1	AS3-AS1		A1	AS4-AS3-AS1
A5	AS1	A5	AS2-AS1	A5	AS3-AS1		A5	AS4-AS3-AS1
B1  B4	AS1-AS2  AS1-AS2	B1  B4	AS2  AS2	B1  B4	AS3-AS2  AS3-AS2		B1  B4	AS4-AS3-AS2  AS4-AS3-AS2
C1  C3	AS1-AS3  AS1-AS3	C1  C3	AS2-AS3  AS2-AS3	C1  C3	AS3  AS3		C1  C3	AS4-AS3  AS4-AS3
D1  D4	AS1-AS2-AS4  AS1-AS2-AS4	D1  D4	AS2-AS3-AS4  AS2-AS3-AS4	D1  D4	AS3-AS4  AS3-AS4		D1  D4	AS4  AS4
	A1 Table		B1 Table		C1 Table			D1 Table

<u>Updating</u>When a speaker node receives a two-column table from a neighbor, it updates its owntable byaddingthe nodes thatare notinits routingtableandaddingitsownautonomous systemandthe autonomous systemthat sent the table. Aftera while each speaker has a table and knows how to reach each node in other Ass

- a) **Loop prevention**. The instability of distance vector routing andthe creation of loops can be avoided in path vector routing. When a router receives a message, it checks to see if its AS is in the path list to the destination. If it is, looping is involved and the message is ignored.
- b) **Policy routing**. Policy routing can be easily implemented through path vector routing. When a router receives a message, it can check the path. If one of the AS listed in the path is against its policy, it can ignore that path and that destination. It does not update its routing table with this path, and it does not send this message to its neighbors.
- c) **Optimum path.** What is the optimum path in path vector routing? We are looking for apathto a destination thatis the bestfor the organizationthat runs the AS. One systemmay useRIP, which defines hopcountas the metric; another mayuse OSPF with minimum delay defined as the metric. In our previous figure, each AS may have more than one path to a destination. For example, a pathfrom AS4 to ASI can be AS4-AS3-AS2-AS1, or it can be AS4-AS3-AS1. For the tables, **we chose the one that had the smaller number of ASs**, but this is not always the case. Other criteria, such as security, safety, and reliability, can also be applied

# HierarchicalRouting

As networks grow in size, the router routing tables grow proportionally. Not only is router memory consumed by ever-increasing tables, but more CPU time is needed to scan them and more bandwidth is needed to send status reports about them.

At a certain point, the network may grow to the point where it is no longer feasible for every router to have an entry for every other router, so the routing will have to be done hierarchically, as it is in the telephone network.

When hierarchical routing is used, the routers are divided into what we will call regions. Each router knows all the details about how to route packets to destinations within its own region but knows nothing about the internal structure of other regions.

Forhugenetworks, a two-level hierarchy may beinsufficient; itmaybe necessaryto groupthe regionsinto clusters, the clusters intozones, the zones intogroups, and soon, until werunout of names for aggregations



 $\label{eq:constraint} When a single network becomes very large, an interesting question is ``how many levels should the hierarchy have?''$ 

For example, consider a network with 720 routers. If there is no hierarchy, each router needs 720 routing table entries.

If the network is partitioned into 24 regions of 30 routers each, each routerneeds 30 local entries plus 23 remote entries for a total of 53 entries.

If a three-level hierarchy is chosen, with 8 clusters each containing 9 regions of 10 routers, each router needs 10 entries for local routers, 8 entries for routing to other regions within its own cluster, and 7 entries for distant clusters, for a total of 25 entries

Kamoun and Kleinrock (1979) discovered that the optimal number of levels for an *N* router network is In *N*, requiring a total of *e* In *N* entries per router

# CONGESTIONCONTROLALGORITHMS

Too many packets present in (a part of) the network causes packet delay and loss that degrades performance. This situation is called **congestion**.

The network and transport layers share the responsibility for handling congestion. Since congestion occurs within the network, it is the network layer that directly experiences it and must ultimately determine what to do with the excess packets.

However, the most effective way to control congestionis to reduce the load that the transport layer is placing on the network. This requires the network and transport layers to work together. In this chapter we will look at the network aspects of congestion.



Offered Load (packet/sec)

When too much traffic is offered, congestion sets in and performance degrades sharply

Above Figuredepicts the onset of congestion. When the number of packets hosts send into the network is well within its carrying capacity, the number delivered is proportional to the number sent. If twice as many are sent, twice as many are delivered. However, as the offered load approaches the carrying capacity, bursts of traffic occasionally fill up the buffers inside routers and some packets are lost. These lost packets consume some of the capacity, so the numberofdeliveredpacketsfallsbelowtheidealcurve. Thenetwork isnowcongested.Unless the network is well designed, it may experience a **congestion collapse** 

# differencebetweencongestioncontrolandflowcontrol.

Congestion control has to do with making sure the network is able to carry the offered traffic. It is a global issue, involving the behavior of all the hosts and routers.

Flow control, in contrast, relates to the traffic between a particular sender and a particular receiver. Its job is to make sure that a fast sender cannot continually transmit data faster than the receiver is able to absorb it.

To see the difference between these two concepts, consider a network made up of 100-Gbps fiber optic links on which a supercomputer is trying to force feed a large file to a personal computer that is capable of handling only 1 Gbps. Although there is no congestion (the network itself is not in trouble), flow control is needed to force the supercomputer to stop frequently to give the personal computer a chance to breathe.

Atthe otherextreme, considera network with 1-Mbps lines and 1000 large computers, half of whichare trying to transfer files at 100 kbps to the other half. Here, the problemis not thatof fast senders overpowering slow receivers, but that the total offered traffic exceeds what the network can handle.

The reason congestion control and flow control are often confused is that the best way to handle both problems is to get the host to slow down. Thus, a host can get a "slow down" message either because the receiver cannot handle the load or because the network cannot handle it.

Severaltechniquescanbeemployed. These include:

- 1. Warningbit
- 2. Chokepackets
- 3. Loadshedding
- 4. Randomearlydiscard
- 5. Trafficshaping

The first 3 deal with congestion detection and recovery. The last 2 deal with congestion avoidance

# WarningBit

- 1. A special bit in the packet header is set by the router to warn the source when congestionis detected.
- 2. Thebitiscopiedand piggy-backedontheACKandsentto thesender.
- 3. ThesendermonitorsthenumberofACKpacketsitreceiveswiththewarningbitsetand adjusts its transmission rate accordingly.

# **ChokePackets**

- 1. Amoredirectwayoftellingthesourcetoslowdown.
- 2. Achokepacketisacontrolpacketgeneratedatacongestednodeandtransmittedto restrict traffic flow.
- 3. Thesource, onreceiving the chokepacket must reduce its transmission rate by a certain percentage.
- 4. AnexampleofachokepacketistheICMPSourceQuenchPacket. <u>Hop-</u> <u>by-Hop Choke Packets</u>
- 1. Overlongdistancesorathighspeedschokepacketsarenotveryeffective.
- 2. Amoreefficientmethodistosendtochokepacketshop-by-hop.
- 3. This requireseachhoptoreduceits transmissionevenbefore thechokepacketarriveat the source

# LoadShedding

- 1. Whenbuffersbecomefull, routerssimply discard packets.
- 2. Whichpacketischosentobe thevictimdependsontheapplicationandontheerror strategy used in the data link layer.
- 3. Forafile transfer, for, e.g. cannot discardolderpacketssincethis will cause a gap in the received data.
- 4. Forreal-timevoiceorvideo itisprobablybetterto throwawayolddata andkeepnew packets.
- 5. Gettheapplicationtomarkpacketswithdiscardpriority.

# RandomEarlyDiscard(RED)

- 1. Thisisaproactive approach in which therouter discards one or morepackets *before* the buffer becomes completely full.
- 2. Eachtimeapacketarrives, the RED algorithm computes the average queuelength, *avg*.
- 3. If *avg* is lower than some lower threshold, congestion is assumed to be minimal or non-existent and the packet is queued.
- 4. If *avg* is greater than some upper threshold, congestion is assumed to be serious and thepacket is discarded.
- 5. If *avg* is between the two thresholds, this might indicate the onset of congestion. The probability of congestion is then calculated.

# TrafficShaping

- 1. Anothermethodofcongestioncontrolisto"shape"thetrafficbeforeitentersthe network.
- 2. Trafficshapingcontrolsthe*rate*atwhichpacketsaresent(notjusthowmany).Usedin ATM and Integrated Services networks.
- 3. Atconnectionset-uptime, these nder and carrier negotiate atraffic pattern (shape).

Twotrafficshapingalgorithmsare: Leaky Bucket TokenBucket

The **Leaky Bucket Algorithm** used to control rate in a network. It is implemented as a singleserver queue with constant service time. If the bucket (buffer) overflows then packets are discarded.



- 1. The leaky bucket enforces a constant output rate (average rate) regardless of theburstiness of the input. Does nothing when input is idle.
- 2. The host injects one packet per clock tick onto the network. This results in a uniform flowof packets, smoothing out bursts and reducing congestion.
- 3. When packets are the same size (as in ATM cells), the one packet per tick is okay. For variable length packets though, it is better to allow a fixed number of bytes per tick. E.g. 1024 bytes per tick will allow one 1024-byte packet or two 512-byte packets or four 256- byte packets on 1 tick

# TokenBucketAlgorithm

- 1. Incontrastto theLB, the TokenBucketAlgorithm, allows the outputrate tovary, depending on the size of the burst.
- 2. IntheTBalgorithm,thebucketholdstokens.Totransmitapacket,thehostmustcapture and destroy one token.
- 3. Tokensaregeneratedbyaclockatthe rateofonetokenevery∆tsec.
- 4. Idle hostscancaptureandsaveuptokens(uptothe max.sizeof the bucket)in orderto send larger bursts later.



#### LeakyBucketvs.TokenBucket

- 1. LBdiscardspackets;TBdoesnot.TBdiscardstokens.
- 2. With TB, a packet can only be transmitted if there are enough tokens to cover its length in bytes.
- 3. LBsends packetsatanaverage rate. TBallowsforlarge bursts to be sentfaster byspeeding up the output.
- 4. TBallowssavinguptokens(permissions)tosendlargebursts.LBdoesnotallowsaving.

# **TRANSPORTLAYER**

The network layer provides end-to-end packet delivery using datagrams or virtual circuits.

The transport layer builds on the network layer to provide <u>datatransport</u> <u>from a process on a source machine to a process on adestination</u> <u>machine</u>with a desired level of reliability that is independent of the physical networks currently in use.

# <u>ServicesProvidedtotheUpperLayers</u>

The ultimate goal of the transport layer is to provide <u>efficient</u>, <u>reliable</u>, and <u>cost-effective</u> data transmission service to its users, normally processes in the application layer.

To achieve this, the transport layer makes use of the services provided by the network layer. The software and/or hardware within the transport layer that does the work is called the **transport entity**.



Thenetwork, transport, and application layers.

The connection-oriented transport service. connections have three phases: establishment, data transfer, and release.

Addressingandflowcontrol

The connection less transports ervice.




*IPaddressesversusportnumbers* 





### Figure23.6Multiplexing and demultiplexing



### Figure23.7Errorcontrol



### Figure23.8PositionofUDP,TCP,andSCTPinTCP/IPsuite



# **TransportServicePrimitives**

Primitive	Packet sent	Meaning	
LISTEN	(none)	Block until some process tries to connect	
CONNECT	CONNECTION REQ.	Actively attempt to establish a connection	
SEND	DATA	Send information	
RECEIVE	(none)	Block until a DATA packet arrives	
DISCONNECT	DISCONNECTION REQ.	This side wants to release the connection	

# Theprimitives for a simpletransport service.

Tostart with, the server executes a LISTENprimitive, typically by calling a library procedure that makes a system call that blocks the server until a client turns up.

When a client wants to talk to the server, it executes a CONNECT primitive. The transport entity carries out this primitive by blocking the caller and sending a packet to the server. The client's CONNECT call causes a CONNECTION REQUEST segment to be sent to the server. When it arrives, the transport entity checks to see that the server is blocked on a LISTEN (i.e., is interested in handling requests). If so, it then unblocks the server and sends a CONNECTION ACCEPTED segment back to the client. When this segment arrives, the client is unblocked and the connection is established.

Data can now be exchanged using the SEND and RECEIVE primitives. In the simplest form, either party can do a (blocking) RECEIVE to wait for the other party to do a SEND. When the segment arrives, the receiver is unblocked. It can then process the segment and send a reply. As long as both sides can keep track of whose turn it is to send, this scheme works fine. When a connection is no longer needed, it must be released to free up table space within the two transport entities. Disconnection has two variants: <u>asymmetric and symmetric</u>.

In the <u>asymmetric</u> variant, either transport user can issue a DISCONNECT primitive, which results in a DISCONNECT segment being sent to the remote transport entity. Upon its arrival, the connection is released.

In the <u>symmetric</u> variant, each direction is closed separately,independently of the other one. When one side does a DISCONNECT, that means it has no more data to send but it is still willing to accept data from its partner. In this model, a connection is released when both sides have done a DISCONNECT



# **TransportServicePrimitives**



Astate diagram for a simple connection management scheme. Transitions labeled in italics are caused by packet arrivals.The solid lines show the client's state sequence. The dashed lines show the server's state sequence.

# ElementsofTransportProtoc ols

- Addressing
- ConnectionEstablishment
- ConnectionRelease
- FlowControlandBuffering
- Multiplexing
- CrashRecovery

# **TransportProtocol**



(a) Environmentofthedatalinklayer.(b) Environmentofthetransportlayer.

1 Over point-to-point links such as wires or optical fiber, it is usually not necessary for a router to specify which router it wants to talk to—each outgoing line leads directly to a particular router. In the transport layer, explicit addressing of destinations is required.

2 The process of establishing a connection over the wire of Fig(a) is simple: theotherendisalwaysthere(unlessithascrashed,inwhichcaseitisnot there).Eitherway,thereisnotmuchtodo.Evenonwirelesslinkstheprocess is not much different. Just sending a message is sufficient to have it reachallotherdestinations.Ifthemessageisnotacknowledgedduetoan error, it can be resent. In the transport layer, initial connection establishment is complicated, as we will see.

3 Another (exceedingly annoying) difference between the data link layer and thetransportlayeristhepotentialexistenceofstoragecapacityinthe network.Theconsequencesofthenetwork'sabilitytodelayandduplicate packetscansometimesbedisastrousandcanrequiretheuseofspecialprotocols to correctly transport information.

4.Buffering and flow control are needed in both layers, but the presence in the transport layer of a large and varying number of connections with bandwidth that fluctuates as the connections competewitheachothermayrequireadifferentapproachthanwe

# **Addressing**

When an application (e.g., a user) process wishes to set up a connection to a remote application process, it must specify which one to connect to. (Connectionless transport has the same problem: to whom should each message be sent?) The method normally used is to define transport addresses to which processes can listen for connection requests. In the Internet, these endpoints are called **ports.** We will use the generic term **TSAP(TransportServiceAccessPoint)** tomeanaspecificendpointinthe transport layer. The analogous endpoints in the network layer (i.e., network layer addresses) are not-surprisingly called **NSAPs** (**Network Service Access Points**). **IP addresses are examples of NSAPs**.



TSAPs,NSAPsandtransport connections.

Apossible scenario fora transportconnection is as follows:

1. Amail server process attaches itself to TSAP 1522 on host 2 to wait for an incoming call.Acall such as our LISTEN might be used, for example.

2. An application process on host 1 wants to send an email message, so it attaches itself to TSAP 1208 and issues a CONNECT request. Therequest specifies TSAP 1208 on host 1 as the source and TSAP 1522 on host 2 as the destination. This action ultimately results in a transport connection being established between the application process and the server.

- 3. Theapplicationprocesssends overthe mailmessage.
- 4. Themailserverrespondstosaythatitwilldeliver the message.
- 5. Thetransportconnectionisreleased.

# special process called a portmapper

### CONNECTIONESTABLISHMENT

Establishinga connection sounds easy, butitis actually surprisingly tricky. At first glance, it would seem sufficient for one transport entity to just send a CONNECTION REQUEST segment to the destination and wait for a CONNECTION ACCEPTED reply. The problem occurs when the network canlose, delay, corrupt, and duplicate packets. This behavior causes serious complications

Tosolvethisspecificproblem, (DELAYEDDUPLICATES) Tomlinson (1975) introduced the **three-way handshake.** This establishment protocol involves one peer checking with the other that the connection request is indeed current. The normal setup procedure when host 1 initiates is shown in Fig. (a). Host 1 chooses a sequence number, *x*, and sends a CONNECTION REQUEST segment containing it to host 2. Host 2 replies with an ACK segment acknowledging *x* and announcing its own initial sequence number, *y*. Finally, host 1 acknowledges host 2's choice of an initial sequence number in the first data segment that it sends.

# ConnectionEstablishment



Three protocol scenarios for establishing a connection using a three-way handshake. CRdenotesCONNECTIONREQUEST. (a) Normaloperation,

- (b) OldCONNECTIONREQUESTappearingoutofnowhere.
- (c) DuplicateCONNECTIONREQUESTandduplicateACK.

In Fig.(b), the first segment is a delayed duplicate CONNECTION REQUEST from an old connection. This segment arrives at host 2 without host 1's knowledge. Host 2 reacts to this segment by sending host 1 an ACK segment, in effect asking for verification that host 1 was indeed trying to set up a new connection. When host 1 rejects host 2's attempt to establish a connection, host 2 realizes that it was tricked by a delayed duplicate and abandons the connection. In this way, a delayed duplicate does no damage

The worst case is when both a delayed CONNECTION REQUEST and an ACK are floating around in the subnet. This case is shown in Fig. (c). As in the previous example, host 2 gets a delayed CONNECTION REQUEST and replies to it. Atthis point, it crucial to realize that host 2 has proposed using *y* as the initial sequence number for host 2 to host 1 traffic, knowing full well that no segments containing sequence number *y* or acknowledgements to *y* are still in existence. When the second delayed segment arrives at host 2, the fact that *z* has been acknowledged rather than *y* tells host 2 that this, too, is an old duplicate. The important thing to realize here is that there is no combination of old segments that can cause the protocol to fail and have a connection set up by accident when no one wants it.



Abruptdisconnectionwithlossofdata.

there are two styles of terminating a connection: asymmetric release and symmetric release Asymmetric release is the way the telephone system works: when one party hangs up, the connection is broken. Symmetric release treats the connection as two separate unidirectional connections and requires each one to be released separately

Asymmetric release is abrupt and may result in data loss. Consider thescenario of Fig.After the connection is established, host 1 sends a segment that arrives properly at host 2. Then host 1 sends another segment. Unfortunately, host 2 issues a DISCONNECT before the second segment arrives. The result is that the connection is released and data are lost.

Clearly, a more sophisticated release protocol is needed to avoid data loss. One way is to use symmetric release, in which each direction is released independently of the other one. Here, a host can continue to receive data even after it has sent a DISCONNECT segment.

Symmetric release does the job when each process has a fixed amount of data to send and clearly knows when it has sent it. One can envision a protocol in which host 1 says "I am done. Are you done too?" If host 2 responds: "I am done too. Goodbye, the connection can be safely released."

# Thetwo-armyproblem.





Four protocol scenarios for releasing a connection.(a) Normalcaseofathree-wayhandshake.(b)finalACKlost.



(c) Response lost. (d) Response lost and subsequent DRs lost. In Fig. (a), we see the normal case in which one of the users sends a DR (DISCONNECTION REQUEST) segment to initiate the connection release. When it arrives, the recipient sends back a DR segment and starts a timer, just in case its DR is lost. When this DR arrives, the original sender sends back an ACK segment and releases the connection. Finally, when the ACK segment arrives, thereceiveralsoreleasestheconnection.

If the final ACK segment is lost, as shown in Fig.(b), the situation is saved by the timer. When the timer expires, the connection is released anyway. Now consider the case of the second DR being lost. The user initiating the disconnection will not receive the expected response, will time out, and will start all over again.

In Fig.(c), we see how this works, assuming that the second time nosegmentsarelostandallsegmentsaredeliveredcorrectly and on time.

Last scenario, Fig.(d), is the same as Fig. (c) except that now we assume all the repeated attempts to retransmit the DR also fail duetolostsegments. After *Nretries*, these nderjust gives up and releases the connection. Meanwhile, there convertimes out and

# TCP

TCP is a connection oriented protocol; it creates a virtual connection between two TCPs to send data. In addition, TCP uses flow and error control mechanisms at the transport level. In brief, TCP is called a *connectionoriented, reliable transport protocol. It adds* connection-oriented and reliability features to the services of IP.

**Topicsdiscussedinthissection:** 

TCPServices TCPFeatures Segment ATCPConnection Flow Control Error Control

### **TCPServices**

#### <u>1 Process-to-ProcessCommunication</u>

TCP provides process-to-process communication using port numbers. Below Table lists some well-known port numbers used by TCP.

Port	Protocol	Description		
7	Echo	Echoes a received datagram back to the sender		
9	Discard	Discards any datagram that is received		
11	Users	Active users		
13	Daytime	Returns the date and the time		
17	Quote	Returns a quote of the day		
19	Chargen	Returns a string of characters		
20	FTP, Data	File Transfer Protocol (data connection)		
21	FTP, Control	File Transfer Protocol (control connection)		
23	TELNET	Terminal Network		
25	SMTP	Simple Mail Transfer Protocol		
53	DNS	Domain Name Server		
67	BOOTP	Bootstrap Protocol		
79	Finger	Finger		
80	HTTP	Hypertext Transfer Protocol		
111	RPC	Remote Procedure Call		

### <u>2 StreamDeliveryService</u>

TCP, on the other hand, allows the sending process to deliver data as a stream of bytes and allows the receiving process to obtain data as a stream of bytes. TCP creates an environment in which the two processes seem to be connected by an imaginary "tube" that carries their data across the Internet. This imaginary environment is showedin below Figure. The sendingprocess produces (writes to) the stream of bytes, and the receiving process consumes (reads from) them



<u>3</u> Sending and Receiving Buffers Because the sending and the receiving processes may not write or read data at the same speed, TCP needs buffers for storage. There are two buffers, the sending buffer and the receiving buffer, one for each direction. One way to implement a buffer is to use a circular array of I-byte locations as shown in Figure. For simplicity, we have shown two buffers of 20 bytes each. Normally the buffers are hundreds or thousands of bytes, depending on the implementation. Wealso show the buffers as the same size, which is not always the case.



Figureshows the movement of the data in one direction. At the <u>sendingsite</u>, the buffer has three types of chambers. The white section contains empty chambers that can be filled by the sending process (producer). The gray area holds bytes that have been sent but not yet acknowledged. TCP keeps these bytes in the buffer until it receives an acknowledgment. The colored area contains bytes to be sent by the sending TCP.

However, as we will see later in this chapter, TCP may be able to send only part of this colored section. This could be due to the slowness of the receiving process or perhaps to congestion in the network. Also note that after the bytes in the gray chambers are acknowledged, the chambers are recycled and available for use by the sending process.

Thisis whyweshowa circularbuffer.

The operation of the buffer at the <u>receiver site</u> is simpler. The circular buffer is divided into two areas (shown as white and colored). The white area contains empty chambers to be filled by bytes received from the network. The colored sections contain received bytes that can be read by the receiving process. When a byte is read by the receiving process, the chamber is recycled and added to the pool of empty chambers.

#### 4TCPsegments



At the transport layer, TCP groups a number of bytes together into a packet called as egment. TCP adds a header to each segment (for control purposes) and delivers the segment to the IP layer for transmission. These gments are encapsulated in IP datagrams and transmitted.

This entire operation is transparent to the receiving process. Later we will see that segments may be received out of order, lost, or corrupted and resent.All these are handled byTCPwith the receiving process unaware of any activities.Above fig shows howsegments are created from the bytes in the buffers

### 5Full-DuplexCommunication

TCP offers full-duplex service, in which data can flow in both directions at the same time. Each TCP then has a sending and receiving buffer, and segments move in both directions

### 6Connection-OrientedService

TCPisaconnection-orientedprotocol.WhenaprocessatsiteA wantsto sendandreceivedatafromanotherprocessatsiteB,thefollowingoccurs:

- 1. ThetwoTCPsestablishaconnectionbetween them.
- 2. Dataare exchanged in both directions.
- 3. The connection is terminated.

#### 7ReliableService

TCPis a reliable transport protocol. It uses an acknowledgment mechanism to check the safe and sound arrival of data. We will discuss this feature further in the section on error control.

## **TCPFeatures**

#### 1NumberingSystem

There are two fields called the sequence number and the acknowledgment number. These two fields refer to the byte number and not the segment number.

Byte Number The bytes of data being transferred in each connection are numbered by TCP.The numbering starts with a randomly generated number. For example, if the random number happens to be 1057 and the total data to be sent are 6000 bytes, the bytes are numbered from 1057 to 7056. We will see that byte numbering is used for flow and error control. Sequence Number After the bytes have been numbered, TCP assigns a sequence number to each segment that is being sent. The sequence number for each segment is the number of the first byte carried in that segment.

Acknowledgment Number The value of the acknowledgment field in a segment defines the number of the next byte a party expects to receive. The acknowledgment number is cumulative.

#### 2FlowControl

TCP, provides *flow control. The receiver of the data controls the amount of* data that are to be sent by the sender. This is done to prevent the receiver from being overwhelmed with data. The numbering system allows TCP to use a byte-oriented flow control.

#### <u>3ErrorControl</u>

Toprovide reliable service, TCPimplements an error control mechanism. Although error control considers a segment as the unit of data for error detection (loss or corrupted segments), error control is byte-oriented, as we will see later.

#### 4Congestion Control

TCP takes into account congestion in the network. The amount of data sent by a sender is not only controlled by the receiver (flow control), butis also determined by the level of congestion in the network

## **TCPsegment** format



Thesegmentconsistsofa20-to60-byteheader,.

Source port address. This is a 16-bit field that defines the port number of the application program in the host that is sending the segment.

Destination port address. This is a 16-bit field that defines the port number of the application program in the host that is receiving the segment.

Sequence number. This 32-bit field defines the number assigned to the first byte of data contained in this segment. As we said before, TCP is a stream transport protocol. To ensure connectivity, each byte to be transmitted is numbered. The sequencenumbertellsthedestinationwhichbyteinthissequencecomprises the first byte in the segment. During connection establishment, each party uses a random number generator to create an initial sequence number (ISN), which is usually different in each direction.

Acknowledgment number. This 32-bit field defines the byte number that the receiver of the segment is expecting to receive from the other party. If the receiver of the segment has successfully received byte number x from the other party, it defines x + I as the acknowledgment number. Acknowledgment and data can be piggybacked together.

Header length. This 4-bit field indicates the number of 4-byte words in the TCP header. The length of the header can be between 20 and 60 bytes. Therefore, the value of this field can be between 5 (5 x 4 = 20) and 15 (15 x 4 = 60).

Reserved. This is a 6-bit field reserved for future use.

Control.Thisfielddefines6differentcontrolbitsorflagsasshowninFigure.One or more of these bits can be set at a time.

URG: Urgent pointer is valid ACK: Acknowledgment is valid PSH: Request for push RST: Reset the connection SYN: Synchronize sequence numbers FIN: Terminate the connection

URG ACK	PSH	RST	SYN	FIN
---------	-----	-----	-----	-----

These bits enable flow control, connection establishment and termination, connection abortion, and the mode of data transfer in TCP.

Window size. This field defines the size of the window, in bytes, that the other party must maintain. Note that the length of this field is 16 bits, which means that the maximum size of the window is 65,535 bytes. This value is normally referred to as the receiving window (rwnd) and is determined by the receiver. The sender must obey the dictation of the receiver in this case. Checksum. This 16-bit field contains the checksum. The calculation of the checksumforTCPfollowsthesameprocedureastheonedescribedforUDP. However, the inclusion of the checksum in the UDP datagram is optional,

whereas the inclusion of the checksum for TCP is mandatory. The same pseudoheader, serving the same purpose, is added to the segment. For the TCP pseudoheader, the value for the protocol field is 6.
Urgent pointer. This I6-bit field, which is valid only if the urgent flag is set, is used when the segment contains urgent data. It defines the number that must be added to the sequence number to obtain the number of the last urgent byte in the data section of the segment. This will be discussed later in this chapter.

Options. There can be up to 40 bytes of optional information in the TCP header. We will not discuss these options here; please refer to the reference list for more information.

## **ATCPConnection**

TCP is connection-oriented. A connection-oriented transport protocol establishes a virtual path between the source and destination. All the segments belonging to a message are then sent over this virtual path. Using a single virtual pathway for the entire message facilitates the acknowledgment process as well as retransmission of damaged or lost frames.

InTCP, connection-oriented transmission requires three phases:

- 1. connectionestablishment,
- 2. data transfer,
- 3. connectiontermination.

### **TCPconnectionestablishment(3wayhandshaking)**

1 The client sends the first segment, a SYN segment, in which only the SYNflag is set. NOTE:A SYN segment cannot carry data, but it consumes one sequence number.

 The server sends the second segment, a SYN +ACK segment, with 2 flag bitsset:<u>SYNandACK</u>.Thissegmenthasadualpurpose.ItisaSYN segment for communication in the other direction and serves as the acknowledgment for the SYN segment. It consumes one sequence number. NOTE:ASYN+ACKsegmentcannot

carrydata, but does consume on esequence number

3. The client sends the third segment. This is just an ACK segment. It acknowledges the receipt of the second segment with the ACK flag and acknowledgment number field. Note that the sequence number in this segment is the same as the one in the SYN segment; the ACK segment does not consume any sequence numbers.

NOTE:AnACKsegment, if carrying nodata, consumes no sequence number



#### **SYNFloodingAttack**

This happens when a malicious attacker sends a large number of SYN segments to a server, pretending that each of them is corning from a different client by faking the source IP addresses in the datagram's.

The server, assuming that the clients are issuing an active open, allocates the necessary resources, such as creating communication tables and setting timers. The TCP server then sends the SYN +ACK segments to the fake clients, which are lost. During this time, however, a lot of resources are occupied without being used. If, during this short time, the number of SYN segments is large, the server eventually runs out of resources and may crash. This SYN flooding attack belongs to a type of security attackknown as a denial-of-service attack, in which an attacker monopolizes a system with so many service requests that the system collapses and denies service to every request.

SOLUTIONS:

1Somehaveimposedalimitonconnectionrequestsduringaspecified period of

time.

20thersfilteroutdatagramscomingfromunwantedsource addresses. 30nerecentstrategyistopostponeresourceallocationuntilthe

entireconnectionissetup usingwhatiscalledacookie

#### **DataTransfer**

After connection is established, bidirectional data transfer can take place. The client and server can both send data and acknowledgments. Data travelinginthesamedirectionasanacknowledgmentarecarriedon the same segment. The acknowledgment is piggybacked with the data

In this example, after connection is established (not shown in the figure), the client sends 2000 bytes of data in two segments. The server then sends 2000 bytes in one segment. The client sends one more segment. The first three segments carry both data and acknowledgment, but the last segment carries only an acknowledgment because there are no more data to be sent. Note the values of the sequence and acknowledgment numbers. The data segments sent by the client have the PSH (push) flag set so that theserver TCP knows to deliver data to the server process as soon as they are received.



**PUSHING DATA:** Delayed transmission and delayed delivery of data may not be acceptable by the application program.

TCP can handle such a situation. The application program at the sending site can request a *push operation. This means that the sending TCPmust not wait for the window* to be filled. It must create a segment and send it immediately. The sending TCP must also set the push bit (PSH) to let the receiving TCP know that the segment includes data that must bedelivered to the receiving application program as soon as possible andnot to wait for more data to come.

<u>Urgent Data :</u> TCP is a stream-oriented protocol. This means that the data are presented from the application program to TCP as a stream of bytes. Each byte of data has a position in the stream. However, sending application program wants a piece of data to be read out of order by the receiving application program.

**Connection Termination** (three-way handshaking and four-way handshaking with a half-close option.)

1. In a normal situation, the client TCP, after receiving a close command from the clientprocess, sends the firstsegment, <u>a FIN segmentinwhichtheFIN flag</u> is set.

Note that a FIN segment can include the last chunk of data sent by the client, or it can be just a control segment as shown in Figure. If it is only a control segment, it consumes only one sequence number.

NOTE: The FIN segment consumes one sequence number if it does not carry data.

2 The server TCP, after receiving the FIN segment, informs its process of the situation and sends the second segment, a <u>FIN +ACK segment</u>, to confirm the receipt of the FIN segment from the client and at the same time to announce the closing of the connection in the other direction. This segment can also contain the last chunk of data from the server. If it does not carry data, it consumes only one sequence number. NOTE:TheFIN+ACKsegmentconsumesonesequencenumberif

itdoesnotcarrydata.

3. The client TCP sends the last segment, an ACK segment, to confirm the receipt of the FIN segment from the TCPserver. This segment contains the acknowledgment number, which is 1 plus the sequence number received in the FIN segment from the server. This segment cannot carry data and consumes no sequence numbers.

<u>**Half-Close</u>**In TCP, one end can stop sending data while still receiving data. This is called a half-close.Although either end can issue a half-close, it is normally initiated by the client. It can occur when the server needs all the data before processing can begin.</u>

A good example is sorting. When the client sends data to the server to be sorted, the server needs to receive all the data before sorting can start. This means the client, after sending all the data, can close the connection in the outbound direction. However, the inbound direction must remain open to receive the sorted data. The server, after receiving the data, still needs time for sorting; its outbound direction must remain open

#### Connectionterminationusingthree-wayhandshaking



#### Figure23.21Half-close



#### FlowControlorTCPSliding Window

TCP uses a sliding window, to handle flow control. The sliding window protocolusedbyTCP,however,issomethingbetweenthe *Go-Back-N*and Selective Repeat sliding window.

The sliding window protocol in TCP looks like the Go-Back-N protocol because it does not use NAKs;

it looks like Selective Repeat because the receiver holds the out-of-order segments until the missing ones arrive.

Therearetwobigdifferencesbetweenthisslidingwindowandtheone we used at the data link layer.

1 the sliding window of TCP is byte-oriented; the one we discussed in the data link layer is frame-oriented.

2 theTCP'sslidingwindowisofvariablesize;theonewediscussedin the data link layer was of fixed size

Sliding window



The window is <u>opened, closed, or shrunk</u>. These three activities, as we will see, are in the control of the receiver (and depend on congestion in the network), not the sender.

The sender mustobey the commands of the receiver in this matter.

<u>Opening</u>awindowmeansmovingtherightwalltotheright.Thisallows more new bytes in the buffer that are eligible for sending.

<u>Closing</u> the window means moving the left wall to the right. This means that some bytes have been acknowledged and the sender need not worry about them anymore.

Shrinking thewindowmeansmovingtherightwalltotheleft.

The size of the window at one end is determined by the lesser of two values: <u>receiver window (rwnd)</u> or <u>congestion window (cwnd)</u>.

The <u>receiver window</u> is the value advertised by the opposite end in a segment containing acknowledgment. It is the number of bytes the other end can accept before its buffer overflows and data are discarded.

The <u>congestion window</u> is a value determined by the network to avoid congestion



WindowmanagementinTCP

When the window is 0, the sender may not normally send segments, with two exceptions.

- 1) urgentdata may be sent, for example, to allow the user to kill the process running on the remote machine.
- 2) the sender may send a 1-byte segment to force the receiver to reannounce the next byte expected and the window size. This packet is called a **window probe.**
- The TCP standard explicitly provides this option to prevent deadlock if a window update ever gets lost.

Senders are not required to transmit data as soon as they come in from the application. Neither are receivers required to send acknowledgements as soon as possible.

Forexample, in Fig. when the first 2KB of data came in, TCP, knowing that it had a 4-KB window, would have been completely correct in just buffering the data until another 2 KB came in, to be able to transmit a segment with a 4-KB payload. This freedom can be used to improve performance



Remoteterminalapplications(e.g.,Telnet)sendcharacterstoaserver.

Theserverinterpretsthecharacterandsendstheoutputattheserver to the client.

Foreach charactertyped,you see three packets: Client Server: Sendtypedcharacter Server Client: Echoof character (or user output) and acknowledgement for first packet Client Server: Acknowledgement for second packet

#### DelayedAcknowledgement

- TCPdelaystransmissionofACKsforupto 500ms
- AvoidtosendACKpacketsthatdonotcarrydata.
  - -The hope is that, within the delay, the receiver will have data ready to besenttothereceiver. Then, the ACK can be piggybacked with a data segment

#### **Exceptions:**

- ACKshould be sentfor everyfullsized segment
- DelayedACKisnotusedwhenpacketsarriveoutoforder

Although delayed acknowledgements reduce the load placed on the network by the receiver, a sender that sends multiple short packets (e.g., 41-byte packets containing 1 byte of data) is still operating inefficiently. A way to reduce this usage is known as **Nagle's algorithm (Nagle, 1984).** 

#### Nagel'sRule

Send one byte and buffer all subsequent bytes until acknowledgement is received. Then send all buffered bytes in a single TCP segment and start buffering again until the sent segment is acknowledged. Nagle's algorithm will put the many pieces in one segment, greatly <u>reducingthe</u> bandwidth used

Nagle's algorithm is widely used by TCP implementations, but there are times when it is better to disable it. In particular, in interactive games that are run over the Internet.

Amore subtle problem is that Nagle's algorithm can sometimes interact with delayed acknowledgements to cause a temporary deadlock: the receiver waits for data on which to piggyback an acknowledgement, and the sender waits on the acknowledgement to send more data.

Because of these problems, Nagle's algorithm can be disabled (which is called the *TCP NODELAY option*).

Anotherproblem that can degrade TCP performance is the **silly window syndrome** (Clark, 1982).



Clark's solution is to prevent the receiver from sending a window update for 1 byte. Instead, it is forced to wait until it has a decent amount of space available and advertise that instead. Specifically, the receiver should not send a window update until it can handle the maximum segment size it advertised when the connection was established or until its buffer is half empty, whichever is smaller.

Furthermore, the sender can also help by not sending tiny segments. Instead, it should wait until it can send a full segment, or at least one containing half of the receiver's buffer size.

The goal is for the sender not to send small segments and the receiver not to ask for them. (Nagel+ Clark). Both are used to improve TCP performance

The receiver will buffer the data until it can be passed up to the application in order (handling out of order segments)

#### Cumulativeacknowledgements

#### **ErrorControl**

TCP is a reliable transport layer protocol. This means that an application program that delivers a stream of data to TCP relies on TCP to deliver the entire stream to the application program on the other end in order, without error, and without any part lost or duplicated.

TCP provides reliability using error control. Error control includes mechanisms for detecting corrupted segments, lost segments, out-of-order segments, and duplicated segments. Error control also includes a mechanism for correcting errors after they are detected. Error detection and correction in TCP is achieved through the use of three simple tools: <u>checksum</u>, <u>acknowledgment</u>, and time-out.

#### <u>Checksum</u>

Each segment includes a checksum field which is used to check for a corrupted segment. If the segment is corrupted, it is discarded by the destination TCP and is considered as lost. TCP uses a 16-bit checksum that is mandatory in every segment

#### Figure 23.11 Checksum calculation of a simple UDP user data gram

153.18.8.105			
171.2.14.10			
All Os	17	15	
1087		13	
15		All Os	
Т	E	S	Т
I	N	G	All Os



#### <u>Acknowledgment</u>

TCP uses acknowledgments to confirm the receipt of data segments. Control segments that carry no data but consume a sequence number are also acknowledged. ACK segments are never acknowledged. ACK segments do not consume sequence numbers and are not

ACK segments do not consume sequence numbers and are not acknowledged.

#### <u>Retransmission</u>

The heart of the error control mechanism is the retransmission of segments. When a segment is corrupted, lost, or delayed, it is retransmitted. In modern implementations, a retransmission occurs if the retransmission timer expires or three duplicate ACK segments have arrived.

RetransmissionAfterRTO(retransmissiontimeout) RetransmissionAfterThreeDuplicateACKSegments(alsocalledfast retransmission)

#### Out-of-OrderSegments

Data may arrive out of order and be temporarily stored by the receiving TCP, but yet guarantees that no out-of-order segment is delivered to the process

### **TCPCongestionControl**

When the load offered to any network is more than it can handle, congestion builds up.

The network layer detects congestion when queues grow large at routers and tries to manage it, if only by dropping packets. It is up to the transportlayer to receive congestion feedback from the network layer and slow downthe rate of traffic that it is sending into the network.

For Congestion control, transport protocol uses an <u>AIMD</u> (Additive Increase Multiplicative Decrease) control law.

TCP congestion control is based on implementing this approach using a window called **congestion window.** TCP adjusts the size of the window according to the AIMD rule.

#### Thewindowsizeatthesenderissetasfollows:

## Send Window = MIN (flow control window, congestion window)

where

flowcontrolwindowisadvertisedbythereceiver(rwnd) congestionwindowisadiustedbasedonfeedbackfromthe Modern congestion control was added to TCP largely through the efforts of Van Jacobson (1988). It is a fascinating story.Starting in 1986, the growing popularity of the early Internet led to the first occurrence of what became known as a **congestion collapse**, a prolonged period during which goodput dropped suddenly (i.e., by more than a factor of 100) due to congestion in the <u>network</u>. Jacobson (and many others) set out to understand whatwas happening and remedy the situation.

To start, he observed that <u>packet loss is a suitable signal of congestion</u>. This signal comes a little late (as the network is already congested) but it is quite dependable

At the beginning how sender knows at what speed receiver can receive the packets?



The key observation is this: the acknowledgements return to the sender at about the rate that packets can be sent over the slowest link in the path. This is precisely the rate that the sender wants to use. If it injects newpackets into the network at this rate, they will be sent as fast as the slow link permits, but they will not queue up and congest any router along the path. This timing is known as an <u>ack clock</u>. It is an essential part ofTCP.By using an ack clock,TCPsmoothes outtraffic and avoids unnecessary queues at routers. This is first consideration A second consideration is that the AIMD rule will take a very long time to reach a good operating point on fast networks if the congestion window is started from a small size

Instead, the solution Jacobson chose to handle both of these considerations is a <u>mix of linear and multiplicative increase</u>.

#### **SLOW-START**



## **TCP Congestion Control**

### Slow Start

- Additive Increase / Multiplicative Decrease is only suitable for source, that is operating close to the available capacity of the network, but it takes too long to ramp up a connection when it is starting from scratch.
- slow start, that is used to increase the congestion window rapidly from a cold start.
- Slow start effectively increases the congestion window exponentially, rather than linearly.
  - the source starts out by setting CongestionWindow to one packet.
  - When the ACK for this packet arrives, TCP adds 1 to CongestionWindow and then sends two packets.
  - Upon receiving the corresponding two ACKs, TCP increments CongestionWindow by 2—one for each ACK—and next sends four packets.
  - The end result is that TCP effectively doubles the number of packets it has in transit every RTT.

Whenever a packet loss is detected, for example, by a timeout, the slow start threshold is set to be <u>half of the congestion window</u> and the entire process is restarted.

Congestion avoidance phase is started if cwnd has reached the slowstart threshold value

Whenever the slow start threshold is crossed, TCP switches from slow start to additive increase. In this mode, the congestion window is increased by one segment every round-trip time.



Figure 6-45. Additive increase from an initial congestion window of one segment.



Fig. 6-37. An example of the Internet congestion algorithm.

# ResponsestoCongestion

- So, TCP assumes there is congestion if it detects a packet loss
- ATCPsendercandetect lostpacketsvia:
  - Timeoutofaretransmissiontimer
  - ReceiptofaduplicateACK
- TCP interprets a Timeout as a binary congestion signal. When a timeout occurs, the sender performs:
  - cwndisresettoone:

cwnd = 1

- ssthreshissettohalfthecurrentsizeofthecongestionwindow:

ssthressh=cwnd/2

- andslow-startisentered

## FastRetransmit

- If three or more duplicate ACKs are received in a row, the TCP sender believes that a segment has been lost.
- Then TCP performs a retransmission of what seems to be the missing segment, without waiting for a timeout to happen.
- Enterslowstart: ssthresh=cwnd/2 cwnd = 1



# FlavorsofTCPCongestion Control

### • **TCPTahoe**(1988)

- SlowStart
- CongestionAvoidance
- FastRetransmit
- TCPReno(1990) (TCPTahoe+FR)
  - FastRecovery
- NewReno(1996)
- SACK (1996) (SACK (Selective ACKnowledgements))

RED(FloydandJacobson1993)



Figure 6-47. Fast recovery and the sawtooth pattern of TCP Reno.
The use of <u>ECN</u> (Explicit Congestion Notification) in addition to packet loss as a congestion signal. ECN is an IP layer mechanism to notify hosts of congestion.

The sender tells the receiver that it has heard the signal by using the <u>CWR</u> (Congestion Window Reduced) flag.

The User Datagram Protocol (UDP) is called a connectionless, unreliable transport protocol. It does not add anything to the services of IP except to provide process-to-process communication instead of host-to-host communication.

*Topicsdiscussedinthissection:* Well-KnownPortsforUDP User Datagram Checksum UDPOperation Use of UDP

## Table23.1Well-knownportsusedwithUDP

Port	Protocol	Description
7	Echo	Echoes a received datagram back to the sender
9	Discard	Discards any datagram that is received
11	Users	Active users
13	Daytime	Returns the date and the time
17	Quote	Returns a quote of the day
19	Chargen	Returns a string of characters
53	Nameserver	Domain Name Service
67	BOOTPs	Server port to download bootstrap information
68	BOOTPc	Client port to download bootstrap information
69	TFTP	Trivial File Transfer Protocol
111	RPC	Remote Procedure Call
123	NTP	Network Time Protocol
161	SNMP	Simple Network Management Protocol
162	SNMP	Simple Network Management Protocol (trap)

## Figure23.9Userdatagram format



## **Checksum**(OPTIONAL, IFNOTUSEDSETALL1'SDEFAULT)

The UDP checksum calculation is different from the one for IP and ICMP.Here the checksum includes three sections: <u>a pseudo header, the UDP header, and</u> <u>the data coming from the application layer.</u>

The pseudo header is the part of the header of the IP packet in which the user datagram is to be encapsulated with some fields filled with Os

If the checksum does not include the pseudo header, a user datagram may arrive safe and sound. However, if the IP header is corrupted, it may be delivered to the wrong host.

TheprotocolfieldisaddedtoensurethatthepacketbelongstoUDP, and not to other transport-layer protocols.

## Figure 23.10 Pseudoheader for checksum calculation



### **UDPOperation** *ConnectionlessServices*

UDP provides a connectionless service. This means that each user datagram sent by UDP is an independent datagram. There is no relationshipbetweenthedifferentuserdatagramseveniftheyarecoming fromthesamesourceprocessandgoingtothesamedestinationprogram. The user datagrams are not numbered. Also, there is no connection establishmentandnoconnectiontermination, asisthecaseforTCP. This meansthateachuserdatagramcantravelonadifferentpath.

#### <u>FlowandErrorControl</u>

UDP is a very simple, unreliable transport protocol. There is no flow control and hence no window mechanism. The receiver may overflow withincoming messages. There is no error control mechanism in UDP except for the checksum. This means that the sender does not know if a message has been lost or duplicated. When the receiver detects an error through the checksum, the user datagram is silently discarded. The lack of flow control and error control

### **EncapsulationandDecapsulation**

To send a message from one process to another, the UDP protocol encapsulates and decapsulates messages in an IP datagram.



## Figure 23.11 Checksum calculation of a simple UDP user data gram

153.18.8.105								
171.2.14.10								
All Os	17	15						
10	87	13						
1	5	All Os						
Т	E	S	Т					
I	N	G	All Os					



## Figure23.12Queues inUDP



## RemoteProcedureCall

The key work was done by Birrell and Nelson (1984). In a nutshell, what Birrell and Nelson suggested was allowing programs to call procedures located on remote hosts. When a process on machine 1 calls a procedure on machine 2, the calling process on 1 is suspended and execution of the called procedure takes place on 2. Information can be transported from the caller to the callee in the parameters and can come back in the procedure result. No message passing is visible to the application programmer. This technique is known as **RPC** (**Remote Procedure Call).** Traditionally, the calling procedure is known as the client and the called procedure is known as the server, and we will use those names here too.

to call a remote procedure, the client program must be bound with a small library procedure, called the **client stub**, that represents the server procedure in the client'saddress space. Similarly,the server is bound with a procedure called the **server stub**. These procedures hide the fact that the procedure call from the client to the server is not local Step 1 is the client calling the client stub. This call is a local procedure call, with the parameters pushed onto the stack in the normal way.

Step 2 is the client stub packing the parameters into a message and making a system call to send the message. Packing the parameters is called **marshaling.** 

Step 3 is the operating system sending the message from the client machine to the server machine.

Step 4 is the operating system passing the incoming packet to the server stub.

Finally, step 5 is the server stub calling the server procedure with the unmarshaled parameters.

Thereplytraces thesame path in the other direction.



Figure 6-29. Steps in making a remote procedure call. The stubs are shaded.

## ProblemswithRPC:

1 WithRPC, passing pointers is impossible because the client and server are in different address spaces.

- 2 Itisessentiallyimpossiblefortheclientstubtomarshaltheparameters: it has no way of determining how large they are.
- 3 A third problem is that it is not always possible to deduce the types of the parameters, not even from a formal specification or the code itself. (exa: printf)
- 4 A fourth problem relates to the use of global variables. Normally, the calling and called procedure can communicate by using global variables, in addition to communicating via parameters. But if the called procedure is moved to a remote machine, the code will fail because the global variables are no longer shared

## TCP

TCP is a connection oriented protocol; it creates a virtual connection between two TCPs to send data. In addition, TCP uses flow and error control mechanisms at the transport level. In brief, TCP is called a *connectionoriented, reliable transport protocol. It adds* connection-oriented and reliability features to the services of IP.

**Topicsdiscussedinthissection:** 

TCPServices TCPFeatures Segment ATCPConnection Flow Control Error Control

## **TCPServices**

### **1Process-to-ProcessCommunication**

TCP provides process-to-process communication using port numbers. Below Table lists some well-known port numbers used by TCP.

Port	Protocol	Description		
7	Echo	Echoes a received datagram back to the sender		
9	Discard	Discards any datagram that is received		
11	Users	Active users		
13	Daytime	Returns the date and the time		
17	Quote	Returns a quote of the day		
19	Chargen	Returns a string of characters		
20	FTP, Data	File Transfer Protocol (data connection)		
21	FTP, Control	File Transfer Protocol (control connection)		
23	TELNET	Terminal Network		
25	SMTP	Simple Mail Transfer Protocol		
53	DNS	Domain Name Server		
67	BOOTP	Bootstrap Protocol		
79	Finger	Finger		
80	HTTP	Hypertext Transfer Protocol		
111	RPC	Remote Procedure Call		

## <u>2StreamDeliveryService</u>

TCP, on the other hand, allows the sending process to deliver data as a stream of bytes and allows the receiving process to obtain data as a stream of bytes. TCP creates an environment in which the two processes seem to be connected by an imaginary "tube" that carries their data across the Internet. This imaginary environment is showedin below Figure. The sendingprocess produces (writes to) the stream of bytes, and the receiving process consumes (reads from) them



<u>3 Sending and Receiving Buffers</u> Because the sending and the receiving processes may not write or read data at the same speed, TCP needs buffers for storage. There are two buffers, the sending buffer and the receiving buffer, one for each direction. One way to implement a buffer is to use a circular array of I-byte locations as shown in Figure. For simplicity, we have shown two buffers of 20 bytes each. Normally the buffers are hundreds or thousands of bytes, depending on the implementation. Wealso show the buffers as the same size, which is not always the case.



Figureshows the movement of the data in one direction. At the <u>sendingsite</u>, the buffer has three types of chambers. The white section contains empty chambers that can be filled by the sending process (producer). The gray area holds bytes that have been sent but not yet acknowledged. TCP keeps these bytes in the buffer until it receives an acknowledgment. The colored area contains bytes to be sent by the sending TCP.

However, as we will see later in this chapter, TCP may be able to send only part of this colored section. This could be due to the slowness of the receiving process or perhaps to congestion in the network. Also note that after the bytes in the gray chambers are acknowledged, the chambers are recycled and available for use by the sending process.

Thisis whyweshowa circularbuffer.

The operation of the buffer at the <u>receiver site</u> is simpler. The circular buffer is divided into two areas (shown as white and colored). The white area contains empty chambers to be filled by bytes received from the network. The colored sections contain received bytes that can be read by the receiving process. When a byte is read by the receiving process, the chamber is recycled and added to the pool of empty chambers.

#### 4TCPsegments



At the transport layer, TCP groups a number of bytes together into a packet called as egment. TCP adds a header to each segment (for control purposes) and delivers the segment to the IP layer for transmission. These gments are encapsulated in IP datagrams and transmitted.

This entire operation is transparent to the receiving process. Later we will see that segments may be received out of order, lost, or corrupted and resent.All these are handled byTCPwith the receiving process unaware of any activities. Above fig shows howsegments are created from the bytes in the buffers

## 5 Full-Duplex Communication

TCP offers full-duplex service, in which data can flow in both directions at the same time. Each TCP then has a sending and receiving buffer, and segments move in both directions

## 6<u>Connection-OrientedService</u>

TCPisaconnection-orientedprotocol.WhenaprocessatsiteA wantsto sendandreceivedatafromanotherprocessatsiteB,thefollowingoccurs:

- 1. ThetwoTCPsestablishaconnectionbetween them.
- 2. Dataare exchanged in both directions.
- 3. The connection is terminated.

### 7<u>ReliableService</u>

TCPis a reliable transport protocol. It uses an acknowledgment mechanism to check the safe and sound arrival of data. We will discuss this feature further in the section on error control.

## **TCPFeatures**

### 1NumberingSystem

There are two fields called the sequence number and the acknowledgment number. These two fields refer to the byte number and not the segment number.

Byte Number The bytes of data being transferred in each connection are numbered by TCP.The numbering starts with a randomly generated number. For example, if the random number happens to be 1057 and the total data to be sent are 6000 bytes, the bytes are numbered from 1057 to 7056. We will see that byte numbering is used for flow and error control. Sequence Number After the bytes have been numbered, TCP assigns a sequence number to each segment that is being sent. The sequence number for each segment is the number of the first byte carried in that segment.

Acknowledgment Number The value of the acknowledgment field in a segment defines the number of the next byte a party expects to receive. The acknowledgment number is cumulative.

#### 2FlowControl

TCP, provides *flow control. The receiver of the data controls the amount of* data that are to be sent by the sender. This is done to prevent the receiver from being overwhelmed with data. The numbering system allows TCP to use a byte-oriented flow control.

### <u>3ErrorControl</u>

Toprovide reliable service, TCPimplements an error control mechanism. Although error control considers a segment as the unit of data for error detection (loss or corrupted segments), error control is byte-oriented, as we will see later.

### **4Congestion Control**

TCP takes into account congestion in the network. The amount of data sent by a sender is not only controlled by the receiver (flow control), butis also determined by the level of congestion in the network

## **TCPsegment** format



Thesegmentconsistsofa20-to60-byteheader,.

Source port address. This is a 16-bit field that defines the port number of the application program in the host that is sending the segment.

Destination port address. This is a 16-bit field that defines the port number of the application program in the host that is receiving the segment.

Sequence number. This 32-bit field defines the number assigned to the first byte of data contained in this segment. As we said before, TCP is a stream transport protocol. To ensure connectivity, each byte to be transmitted is numbered. The sequencenumbertellsthedestinationwhichbyteinthissequencecomprises the first byte in the segment. During connection establishment, each party uses a random number generator to create an initial sequence number (ISN), which is usually different in each direction.

Acknowledgment number. This 32-bit field defines the byte number that the receiver of the segment is expecting to receive from the other party. If the receiver of the segment has successfully received byte number x from the other party, it defines x + I as the acknowledgment number. Acknowledgment and data can be piggybacked together.

Header length. This 4-bit field indicates the number of 4-byte words in the TCP header. The length of the header can be between 20 and 60 bytes. Therefore, the value of this field can be between 5 (5 x 4 = 20) and 15 (15 x 4 = 60).

Reserved. This is a 6-bit field reserved for future use.

Control.Thisfielddefines6differentcontrolbitsorflagsasshowninFigure.One or more of these bits can be set at a time.

URG: Urgent pointer is valid ACK: Acknowledgment is valid PSH: Request for push RST: Reset the connection SYN: Synchronize sequence numbers FIN: Terminate the connection

URG ACK	PSH	RST	SYN	FIN
---------	-----	-----	-----	-----

These bits enable flow control, connection establishment and termination, connection abortion, and the mode of data transfer in TCP.

Window size. This field defines the size of the window, in bytes, that the other party must maintain. Note that the length of this field is 16 bits, which means that the maximum size of the window is 65,535 bytes. This value is normally referred to as the receiving window (rwnd) and is determined by the receiver. The sender must obey the dictation of the receiver in this case.

Checksum. This 16-bit field contains the checksum. The calculation of the checksumforTCPfollowsthesameprocedureastheonedescribedforUDP.

However, the inclusion of the checksum in the UDP datagram is optional, whereas the inclusion of the checksum for TCP is mandatory. The same pseudoheader, serving the same purpose, is added to the segment. For the TCP pseudoheader, the value for the protocol field is 6. Urgent pointer. This I6-bit field, which is valid only if the urgent flag is set, is used when the segment contains urgent data. It defines the number that must be added to the sequence number to obtain the number of the last urgent byte in the data section of the segment. This will be discussed later in this chapter.

Options. There can be up to 40 bytes of optional information in the TCP header. We will not discuss these options here; please refer to the reference list for more information.

# **ATCPConnection**

TCP is connection-oriented. A connection-oriented transport protocol establishes a virtual path between the source and destination. All the segments belonging to a message are then sent over this virtual path. Using a single virtual pathway for the entire message facilitates the acknowledgment process as well as retransmission of damaged or lost frames.

InTCP, connection-oriented transmission requires three phases:

- 1. connectionestablishment,
- 2. data transfer,
- 3. connectiontermination.

## **TCPconnectionestablishment(3wayhandshaking)**

1 The client sendsthe first segment, a SYN segment, in which only the SYNflag is set. NOTE:A SYN segment cannot carry data, but it consumes one

sequence number.

2. The server sends the second segment, a SYN +ACK segment, with 2 flag bitsset: <u>SYNandACK</u>. This segment has adual purpose. It is a SYN segment for communication in the other direction and serves as the acknowledgment for the SYN segment. It consumes one sequence number. NOTE: ASYN+ACK segment cannot

carrydata, but does consume on esequence number

3. The client sends the third segment. This is just an ACK segment. It acknowledges the receipt of the second segment with the ACK flag and acknowledgment number field. Note that the sequence number in this segment is the same as the one in the SYN segment; the ACK segment does not consume any sequence numbers.

NOTE:AnACKsegment, if carrying nodata, consumes no sequence number



## **SYNFloodingAttack**

This happens when a malicious attacker sends a large number of SYN segments to a server, pretending that each of them is corning from a different client by faking the source IP addresses in the datagram's.

The server, assuming that the clients are issuing an active open, allocates the necessary resources, such as creating communication tables and setting timers. The TCP server then sends the SYN +ACK segments to the fake clients, which are lost. During this time, however, a lot of resources are occupied without being used. If, during this short time, the number of SYN segments is large, the server eventually runs out of resources and may crash. This SYN flooding attack belongs to a type of security attackknown as a denial-of-service attack, in which an attacker monopolizes a system with so many service requests that the system collapses and denies service to every request.

SOLUTIONS:

1Somehaveimposedalimitonconnectionrequestsduringaspecified period of

time.

20thersfilteroutdatagramscomingfromunwantedsource addresses. 30nerecentstrategyistopostponeresourceallocationuntilthe

entireconnectionissetup usingwhatiscalledacookie

### **DataTransfer**

After connection is established, bidirectional data transfer can take place. The client and server can both send data and acknowledgments. Data travelinginthesamedirectionasanacknowledgmentarecarriedon the same segment. The acknowledgment is piggybacked with the data

In this example, after connection is established (not shown in the figure), the client sends 2000 bytes of data in two segments. The server then sends 2000 bytes in one segment. The client sends one more segment. The first three segments carry both data and acknowledgment, but the last segment carries only an acknowledgment because there are no more data to be sent. Note the values of the sequence and acknowledgment numbers. The data segments sent by the client have the PSH (push) flag set so that theserver TCP knows to deliver data to the server process as soon as they are received.



**PUSHING DATA:** Delayed transmission and delayed delivery of data may not be acceptable by the application program.

TCP can handle such a situation. The application program at the sending site can request a *push operation. This means that the sending TCPmust not wait for the window* to be filled. It must create a segment and send it immediately. The sending TCP must also set the push bit (PSH) to let the receiving TCP know that the segment includes data that must bedelivered to the receiving application program as soon as possible andnot to wait for more data to come.

<u>Urgent Data :</u> TCP is a stream-oriented protocol. This means that the data are presented from the application program to TCP as a stream of bytes. Each byte of data has a position in the stream. However, sending application program wants a piece of data to be read out of order by the receiving application program.

**Connection Termination** (three-way handshaking and four-way handshaking with a half-close option.)

1. In a normal situation, the client TCP, after receiving a close command from the clientprocess, sends the firstsegment, <u>a FIN segmentinwhichtheFIN flag</u> is set.

Note that a FIN segment can include the last chunk of data sent by the client, or it can be just a control segment as shown in Figure. If it is only a control segment, it consumes only one sequence number.

NOTE: The FIN segment consumes one sequence number if it does not carry data.

2 The server TCP, after receiving the FIN segment, informs its process of the situation and sends the second segment, a <u>FIN +ACK segment</u>, to confirm the receipt of the FIN segment from the client and at the same time to announce the closing of the connection in the other direction. This segment can also contain the last chunk of data from the server. If it does not carry data, it consumes only one sequence number.

NOTE: The FIN+ACK segment consumes one sequence number if it does not carry data.

3. The client TCP sends the last segment, an ACK segment, to confirm the receipt of the FIN segment from the TCPserver. This segment contains the acknowledgment number, which is 1 plus the sequence number received in the FIN segment from the server. This segment cannot carry data and consumes no sequence numbers.

<u>**Half-Close</u>**In TCP, one end can stop sending data while still receiving data. This is called a half-close.Although either end can issue a half-close, it is normally initiated by the client. It can occur when the server needs all the data before processing can begin.</u>

A good example is sorting. When the client sends data to the server to be sorted, the server needs to receive all the data before sorting can start. This means the client, after sending all the data, can close the connection in the outbound direction. However, the inbound direction must remain open to receive the sorted data. The server, after receiving the data, still needs time for sorting; its outbound direction must remain open
### Connectionterminationusingthree-wayhandshaking



## Figure23.21Half-close



### FlowControlorTCPSliding Window

TCP uses a sliding window, to handle flow control. The sliding window protocolusedbyTCP, however, is something between the *Go-Back-N* and Selective Repeat sliding window.

The sliding window protocol in TCP looks like the Go-Back-N protocol because it does not use NAKs;

it looks like Selective Repeat because the receiver holds the out-of-order segments until the missing ones arrive.

Therearetwobigdifferencesbetweenthisslidingwindowandtheone we used at the data link layer.

1 the sliding window of TCP is byte-oriented; the one we discussed in the data link layer is frame-oriented.

2 theTCP'sslidingwindowisofvariablesize;theonewediscussedin the data link layer was of fixed size

Sliding window



The window is <u>opened, closed, or shrunk</u>. These three activities, as we will see, are in the control of the receiver (and depend on congestion in the network), not the sender.

The sender mustobey the commands of the receiver in this matter.

<u>Opening</u>awindowmeansmovingtherightwalltotheright. Thisallows more new bytes in the buffer that are eligible for sending.

<u>Closing</u> the window means moving the left wall to the right. This means that some bytes have been acknowledged and the sender need not worry about them anymore.

Shrinking thewindowmeansmovingtherightwalltotheleft.

The size of the window at one end is determined by the lesser of two values: <u>receiver window (rwnd)</u> or <u>congestion window (cwnd)</u>.

The <u>receiver window</u> is the value advertised by the opposite end in a segment containing acknowledgment. It is the number of bytes the other end can accept before its buffer overflows and data are discarded.

The <u>congestion window</u> is a value determined by the network to avoid congestion



WindowmanagementinTCP

When the window is 0, the sender may not normally send segments, with two exceptions.

- 1) urgentdata may be sent, for example, to allow the user to kill the process running on the remote machine.
- 2) the sender may send a 1-byte segment to force the receiver to reannounce the next byte expected and the window size. This packet is called a **window probe.**
- The TCP standard explicitly provides this option to prevent deadlock if a window update ever gets lost.

Senders are not required to transmit data as soon as they come in from the application. Neither are receivers required to send acknowledgements as soon as possible.

Forexample, in Fig. when the first 2KB of data came in, TCP, knowing that it had a 4-KB window, would have been completely correct in just buffering the data until another 2 KB came in, to be able to transmit a segment with a 4-KB payload. This freedom can be used to improve performance



Remoteterminalapplications(e.g.,Telnet)sendcharacterstoaserver.

Theserverinterpretsthecharacterandsendstheoutputattheserver to the client.

Foreach charactertyped,you see three packets: Client Server: Sendtypedcharacter Server Client: Echoof character (or user output) and acknowledgement for first packet Client Server: Acknowledgement for second packet

## DelayedAcknowledgement

- TCPdelaystransmissionofACKsforupto 500ms
- AvoidtosendACKpacketsthatdonotcarrydata.
  - -The hope is that, within the delay, the receiver will have data ready to besenttothereceiver. Then, the ACK can be piggybacked with a data segment

### **Exceptions:**

- ACKshould be sentfor everyfullsized segment
- DelayedACKisnotusedwhenpacketsarriveoutoforder

Although delayed acknowledgements reduce the load placed on the network by the receiver, a sender that sends multiple short packets (e.g., 41-byte packets containing 1 byte of data) is still operating inefficiently. A way to reduce this usage is known as **Nagle's algorithm (Nagle, 1984).** 

### Nagel'sRule

Send one byte and buffer all subsequent bytes until acknowledgement is received. Then send all buffered bytes in a single TCP segment and start buffering again until the sent segment is acknowledged. Nagle's algorithm will put the many pieces in one segment, greatly <u>reducingthe</u> bandwidth used

Nagle's algorithm is widely used by TCP implementations, but there are times when it is better to disable it. In particular, in interactive games that are run over the Internet.

Amore subtle problem is that Nagle's algorithm can sometimes interact with delayed acknowledgements to cause a temporary deadlock: the receiver waits for data on which to piggyback an acknowledgement, and the sender waits on the acknowledgement to send more data.

Because of these problems, Nagle's algorithm can be disabled (which is called the *TCP NODELAY option*).

Anotherproblem that can degrade TCP performance is the **silly window syndrome** (Clark, 1982).



Clark's solution is to prevent the receiver from sending a window update for 1 byte. Instead, it is forced to wait until it has a decent amount of space available and advertise that instead. Specifically, the receiver should not send a window update until it can handle the maximum segment size it advertised when the connection was established or until its buffer is half empty, whichever is smaller.

Furthermore, the sender can also help by not sending tiny segments. Instead, it should wait until it can send a full segment, or at least one containing half of the receiver's buffer size.

The goal is for the sender not to send small segments and the receiver not to ask for them. (Nagel+ Clark). Both are used to improve TCP performance

The receiver will buffer the data until it can be passed up to the application in order (handling out of order segments)

### Cumulativeacknowledgements

### **ErrorControl**

TCP is a reliable transport layer protocol. This means that an application program that delivers a stream of data to TCP relies on TCP to deliver the entire stream to the application program on the other end <u>in order</u>, <u>withouterror</u>, and without any part lost or duplicated.

TCP provides reliability using error control. Error control includes mechanisms for detecting corrupted segments, lost segments, out-of-order segments, and duplicated segments. Error control also includes a mechanism for correcting errors after they are detected. Error detection and correction in TCP is achieved through the use of three simple tools: <u>checksum</u>, <u>acknowledgment</u>, and time-out.

#### <u>Checksum</u>

Each segment includes a checksum field which is used to check for a corrupted segment. If the segment is corrupted, it is discarded by the destination TCP and is considered as lost. TCP uses a 16-bit checksum that is mandatory in every segment

### Figure 23.11 Checksum calculation of a simple UDP user data gram

153.18.8.105						
171.2.14.10						
All Os	17	15				
1087		13				
15		All Os				
Т	E	S	Т			
I	N	G	All Os			



### <u>Acknowledgment</u>

TCP uses acknowledgments to confirm the receipt of data segments. Control segments that carry no data but consume a sequence number are also acknowledged. ACK segments are never acknowledged. ACK segments do not consume sequence numbers and are not

ACK segments do not consume sequence numbers and are not acknowledged.

### <u>Retransmission</u>

The heart of the error control mechanism is the retransmission of segments. When a segment is corrupted, lost, or delayed, it is retransmitted. In modern implementations, a retransmission occurs if the retransmission timer expires or three duplicate ACK segments have arrived.

RetransmissionAfterRTO(retransmissiontimeout) RetransmissionAfterThreeDuplicateACKSegments(alsocalledfast retransmission)

### Out-of-OrderSegments

Data may arrive out of order and be temporarily stored by the receiving TCP, but yet guarantees that no out-of-order segment is delivered to the process

# **TCPCongestionControl**

When the load offered to any network is more than it can handle, congestion builds up.

The network layer detects congestion when queues grow large at routers and tries to manage it, if only by dropping packets. It is up to the transportlayer to receive congestion feedback from the network layer and slow downthe rate of traffic that it is sending into the network.

For Congestion control, transport protocol uses an <u>AIMD</u> (Additive Increase Multiplicative Decrease) control law.

TCP congestion control is based on implementing this approach using a window called **congestion window.** TCP adjusts the size of the window according to the AIMD rule.

### Thewindowsizeatthesenderissetasfollows:

# Send Window = MIN (flow control window, congestion window)

where

flowcontrolwindowisadvertisedbythereceiver(rwnd) congestionwindowisadiustedbasedonfeedbackfromthe Modern congestion control was added to TCP largely through the efforts of Van Jacobson (1988). It is a fascinating story.Starting in 1986, the growing popularity of the early Internet led to the first occurrence of what became known as a **congestion collapse**, a prolonged period during which goodput dropped suddenly (i.e., by more than a factor of 100) due to congestion in the <u>network</u>. Jacobson (and many others) set out to understand whatwas happening and remedy the situation.

To start, he observed that <u>packet loss is a suitable signal of congestion</u>. This signal comes a little late (as the network is already congested) but it is quite dependable

At the beginning how sender knows at what speed receiver can receive the packets?



The key observation is this: the acknowledgements return to the sender at about the rate that packets can be sent over the slowest link in the path. This is precisely the rate that the sender wants to use. If it injects newpackets into the network at this rate, they will be sent as fast as the slow link permits, but they will not queue up and congest any router along the path. This timing is known as an <u>ack clock</u>. It is an essential part ofTCP.By using an ack clock,TCPsmoothes outtraffic and avoids unnecessary queues at routers. This is first consideration A second consideration is that the AIMD rule will take a very long time to reach a good operating point on fast networks if the congestion window is started from a small size

Instead, the solution Jacobson chose to handle both of these considerations is a <u>mix of linear and multiplicative increase</u>.

### **SLOW-START**



# **TCP Congestion Control**

# Slow Start

- Additive Increase / Multiplicative Decrease is only suitable for source, that is operating close to the available capacity of the network, but it takes too long to ramp up a connection when it is starting from scratch.
- slow start, that is used to increase the congestion window rapidly from a cold start.
- Slow start effectively increases the congestion window exponentially, rather than linearly.
  - the source starts out by setting CongestionWindow to one packet.
  - When the ACK for this packet arrives, TCP adds 1 to CongestionWindow and then sends two packets.
  - Upon receiving the corresponding two ACKs, TCP increments CongestionWindow by 2—one for each ACK—and next sends four packets.
  - The end result is that TCP effectively doubles the number of packets it has in transit every RTT.

Whenever a packet loss is detected, for example, by a timeout, the slow start threshold is set to be <u>half of the congestion window</u> and the entire process is restarted.

Congestion avoidance phase is started if cwnd has reached the slowstart threshold value

Whenever the slow start threshold is crossed, TCP switches from slow start to additive increase. In this mode, the congestion window is increased by one segment every round-trip time.



Figure 6-45. Additive increase from an initial congestion window of one segment.



Fig. 6-37. An example of the Internet congestion algorithm.

# ResponsestoCongestion

- So, TCP assumes there is congestion if it detects a packet loss
- ATCPsendercandetect lostpacketsvia:
  - Timeoutofaretransmissiontimer
  - ReceiptofaduplicateACK
- TCP interprets a Timeout as a binary congestion signal. When a timeout occurs, the sender performs:
  - cwndisresettoone:

cwnd = 1

- ssthreshissettohalfthecurrentsizeofthecongestionwindow:

ssthressh=cwnd/2

- andslow-startisentered

# FastRetransmit

- If three or more duplicate ACKs are received in a row, the TCP sender believes that a segment has been lost.
- Then TCP performs a retransmission of what seems to be the missing segment, without waiting for a timeout to happen.
- Enterslowstart: ssthresh=cwnd/2 cwnd = 1



# FlavorsofTCPCongestion Control

# • **TCPTahoe**(1988)

- SlowStart
- CongestionAvoidance
- FastRetransmit
- TCPReno(1990) (TCPTahoe+FR)
  - FastRecovery
- NewReno(1996)
- SACK (1996) (SACK (Selective ACKnowledgements))

RED(FloydandJacobson1993)



Figure 6-47. Fast recovery and the sawtooth pattern of TCP Reno.

The use of <u>ECN</u> (Explicit Congestion Notification) in addition to packet loss as a congestion signal. ECN is an IP layer mechanism to notify hosts of congestion.

The sender tells the receiver that it has heard the signal by using the <u>CWR</u> (Congestion Window Reduced) flag.

The User Datagram Protocol (UDP) is called a connectionless, unreliable transport protocol. It does not add anything to the services of IP except to provide process-to-process communication instead of host-to-host communication.

*Topicsdiscussedinthissection:* Well-KnownPortsforUDP User Datagram Checksum UDPOperation Use of UDP

# Table23.1Well-knownportsusedwithUDP

Port	Protocol	Description
7	Echo	Echoes a received datagram back to the sender
9	Discard	Discards any datagram that is received
11	Users	Active users
13	Daytime	Returns the date and the time
17	Quote	Returns a quote of the day
19	Chargen	Returns a string of characters
53	Nameserver	Domain Name Service
67	BOOTPs	Server port to download bootstrap information
68	BOOTPc	Client port to download bootstrap information
69	TFTP	Trivial File Transfer Protocol
111	RPC	Remote Procedure Call
123	NTP	Network Time Protocol
161	SNMP	Simple Network Management Protocol
162	SNMP	Simple Network Management Protocol (trap)

### Figure23.9Userdatagram format



# **Checksum**(OPTIONAL, IFNOTUSEDSETALL1'SDEFAULT)

The UDP checksum calculation is different from the one for IP and ICMP.Here the checksum includes three sections: <u>a pseudo header, the UDP header, and</u> <u>the data coming from the application layer.</u>

The pseudo header is the part of the header of the IP packet in which the user datagram is to be encapsulated with some fields filled with Os

If the checksum does not include the pseudo header, a user datagram may arrive safe and sound. However, if the IP header is corrupted, it may be delivered to the wrong host.

TheprotocolfieldisaddedtoensurethatthepacketbelongstoUDP, and not to other transport-layer protocols.

### Figure 23.10 Pseudoheader for checksum calculation



### **UDPOperation** *ConnectionlessServices*

UDP provides a connectionless service. This means that each user datagram sent by UDP is an independent datagram. There is no relationshipbetweenthedifferentuserdatagramseveniftheyarecoming fromthesamesourceprocessandgoingtothesamedestinationprogram. The user datagrams are not numbered. Also, there is no connection establishmentandnoconnectiontermination, asisthecaseforTCP. This meansthateachuserdatagramcantravelonadifferentpath.

### <u>FlowandErrorControl</u>

UDP is a very simple, unreliable transport protocol. There is no flow control and hence no window mechanism. The receiver may overflow withincoming messages. There is no error control mechanism in UDP except for the checksum. This means that the sender does not know if a message has been lost or duplicated. When the receiver detects an error through the checksum, the user datagram is silently discarded. The lack of flow control and error control

### **EncapsulationandDecapsulation**

To send a message from one process to another, the UDP protocol encapsulates and decapsulates messages in an IP datagram.



### Figure 23.11 Checksum calculation of a simple UDP user data gram

153.18.8.105					
171.2.14.10					
All Os	17	15			
1087		13			
15		All Os			
Т	E	S	Т		
I	Ν	G	All Os		



## Figure23.12Queues inUDP


# RemoteProcedureCall

The key work was done by Birrell and Nelson (1984). In a nutshell, what Birrell and Nelson suggested was allowing programs to call procedures located on remote hosts. When a process on machine 1 calls a procedure on machine 2, the calling process on 1 is suspended and execution of the called procedure takes place on 2. Information can be transported from the caller to the callee in the parameters and can come back in the procedure result. No message passing is visible to the application programmer. This technique is known as **RPC** (**Remote Procedure Call)**. Traditionally, the calling procedure is known as the client and the called procedure is known as the server, and we will use those names here too.

to call a remote procedure, the client program must be bound with a small library procedure, called the **client stub**, that represents the server procedure in the client'saddress space. Similarly,the server is bound with a procedure called the **server stub**. These procedures hide the fact that the procedure call from the client to the server is not local Step 1 is the client calling the client stub. This call is a local procedure call, with the parameters pushed onto the stack in the normal way.

Step 2 is the client stub packing the parameters into a message and making a system call to send the message. Packing the parameters is called **marshaling.** 

Step 3 is the operating system sending the message from the client machine to the server machine.

Step 4 is the operating system passing the incoming packet to the server stub.

Finally, step 5 is the server stub calling the server procedure with the unmarshaled parameters.

Thereplytraces thesame path in the other direction.



Figure 6-29. Steps in making a remote procedure call. The stubs are shaded.

## ProblemswithRPC:

1 WithRPC, passing pointers is impossible because the client and server are in different address spaces.

- 2 Itisessentiallyimpossiblefortheclientstubtomarshaltheparameters: it has no way of determining how large they are.
- 3 A third problem is that it is not always possible to deduce the types of the parameters, not even from a formal specification or the code itself. (exa: printf)
- 4 A fourth problem relates to the use of global variables. Normally, the calling and called procedure can communicate by using global variables, in addition to communicating via parameters. But if the called procedure is moved to a remote machine, the code will fail because the global variables are no longer shared

## Real-TimeTransportProtocols

- Client-serverRPCisoneareainwhichUDPiswidelyused. Another one is for real-time multimedia applications.
  - Internet radio,
  - Internet telephony,
  - music-on-demand,
  - videoconferencing,
  - video-on-demand,

and other multimedia applications became more commonplace, people have discovered that each application was reinventing more or less the same real-time transport protocol.

- Itgraduallybecameclearthathavingagenericreal-timetransportprotocol for multiple applications would be a good idea.
- Thus was **RTP (Real-time Transport Protocol)** born. It is described in RFC 3550 and is now in widespread use for multimedia applications. We will describe two aspects of real-time transport.
- The first is the RTP protocol for transporting audio and video data in packets. The second is the processing that takes place, mostly at the receiver, to play out the audio and video at the right time..



Figure 6-30. (a) The position of RTP in the protocol stack. (b) Packet nesting.

RTPnormallyrunsin userspaceoverUDP(inthe operating system).

It operates as follows. The multimedia application consists of multiple audio, video, text, and possibly other streams. These are fed into the RTP library, which is in user space along with the application. This library multiplexes the streams and encodes them in RTP packets, which it stuffs into a socket. On the operating system side of the socket, UDP packets are generated to wrap the RTP packets and handed to IP for transmission over a link such as Ethernet.

The reverse process happens at the receiver. The multimedia application eventually receives multimedia data from the RTP library. It is responsible for playing out the media. The protocol stack for this situation is shown in Fig. 6-30(a). The packet nesting is shown in Fig. 6-30(b).

### **RTP**—TheReal-timeTransportProtocol

The basic function of RTP is to multiplex several real-time data streamsonto a single stream of UDP packets. The UDP stream can be sent to a single destination (unicasting) or to multiple destinations (multicasting).

Because RTPjust uses normal UDP, its packets are not treated specially by the routers unless some normal IP quality-of-service features are enabled. In particular, there are no special guarantees about delivery, and packets may be lost, delayed, corrupted, etc.

TheRTPformatcontainsseveralfeatures.

Each packet sent in an RTP stream is given a number one higher than its predecessor. This <u>numbering</u> allows the destination to determine if any packets are missing.

RTP has no acknowledgements, and no mechanism to request retransmissions.

Each RTP payload may contain multiple samples, and they may be coded any way that the application wants. To allow for interworking, RTP defines several profiles (e.g., a single audio stream), and for each profile, <u>multipleencoding formats may be allowed</u>

Another facility many real-time applications need is <u>time stamping</u>. Not only doestimestamping<u>reducetheeffectsofvariationinnetworkdelay</u>, but it

-										
Ver.	P	x	сс	М	Payload type	Sequence number				
	Timestamp									
Synchronization source identifier										
						\$				
Contributing source identifier										

Figure 6-31. The RTP header.

It consists of three 32-bit words and potentially some extensions.

The first word contains the *Version field*, *which* is already at 2.

The *P* bitindicatesthatthe packethasbeenpaddedtoamultiple of 4 bytes. The last padding byte tells how many bytes were added.

The Xbit indicates that an extension header is present.

The CC field tells how many contributing sources are present, from 0 to 15

The *M* bit is an application-specific marker bit. It can be used to mark the start of a video frame, the start of a word in an audio channel, or something else that the application understands.

The *Payload type field tells which encoding* algorithm has been used (e.g., uncompressed 8-bit audio, MP3, etc.). Since every packetcarries this field, theencodingcanchangeduring transmission.

The Sequence number is just a counter that is incremented on each RTP packet sent. It is used to detect lost packets.

The *Timestamp*, *this* value can help reduce timing variability called jitter at the receiver by decoupling the playback from the packet arrival time. The *Synchronization source identifier tells which stream the packet belongs to*. It is the method used to multiplex and demultiplex multiple data streams onto a single stream of UDP packets.

Finally, the Contributing source identifiers, if any, are used when

### RTCP—TheReal-timeTransportControlProtocol

RTPhas a little sister protocol (little sibling protocol?) called RTCP(Real time Transport Control Protocol). It is defined along with RTP in RFC 3550 and handles <u>feedback</u>, synchronization, and the user interface. It does not transport any media samples.

The first function can be used to provide <u>feedback on delay,variation in delayor</u> <u>jitter, bandwidth, congestion, and other network properties</u> to the sources. This information can be used by the encoding process to increase the data rate (and give better quality) when the network is functioning well and to cut back the data rate when there is trouble in the network.By providing continuous feedback, It provides the best quality

The *Payload type field is used to tell the destination* what encoding algorithm is used for the current packet, making it possible to vary it on demand.

RTCP also handles <u>inter stream synchronization</u>. The problem is that different streams may use different clocks, with different granularities and different drift rates. RTCP can be used to keep them in sync.

Finally, RTCP provides a way for <u>naming the various sources</u> (e.g., in ASCII text). This information can be displayed on the receiver's screen to indicate who is talking at the moment.

### **Playoutwith Bufferingand Jitter Control**

Once the media information reaches the receiver, it must be played out at the right time. Even if the packets are injected with exactly the right intervals between them at the sender, they will reach the receiver with different relative times. This variation in delay is called **jitter**. Even a small amount of packet jitter can cause distracting media artifacts, such as jerky video frames and unintelligible audio, if the media is simply played out as itarrives.

Thesolutiontothisproblemisto **buffer**packetsatthereceiverbeforethey are played out to reduce the jitter.



Figure 6-32. Smoothing the output stream by buffering packets.

Akeyconsiderationforsmoothplayoutisthe**playbackpoint**,or how long to wait at the receiver for media before playing it out. Deciding how long to wait depends on the jitter. The difference between a low-jitter and high-jitter connection is shown in Fig. The average delay may not differ greatly between the two, but if there is high jitter the playback point may need to be much further out to capture 99% of the packets than if there is low jitter



Figure 6-33. (a) High jitter. (b) Low jitter.

One way to avoid this problem for audio is to adapt the playback pointbetween**talkspurts**, inthe gaps in a conversation. Noone will notice the difference between a short and slightly longer silence

## <u>TELNET</u>

It is client/server application program. TELNET is an abbreviation for *TErminaLNETwork*. TELNETenables theestablishmentofa connection to a remote system in such a way that the local terminal appears to be a terminal at the remote system.

#### **TimesharingEnvironment**

A large computer supports multiple users. The interaction between auser and the computer occurs through a terminal, which is usually a combination of keyboard, monitor, and mouse.

### Logging

Toaccessthesystem, the user logs into the system with a user idor log-in name. The system also includes password checking to prevent an unauthorized user from accessing the resources.

> Local login Remotelogin



a. Local log-in



b. Remote log-in

When a userlogs into a localtimesharing system, it is called <u>locallog-in</u>. As a user types at a terminal or at a workstation running a terminal emulator, the keystrokes are accepted by the terminal driver. The terminal driverpassesthecharacterstotheoperatingsystem. The operating system, in turn, interprets the combination of characters and invokes the desired application program or utility.

When a user wants to access an application program or utility located on a remote Machine, it is called **remote log-in**. Here the TELNET client and server programs come into use. The user sends the keystrokes to the terminal driver, where the local operating system accepts the characters but does not interpret them. The characters are sent to the TELNET client, whichtransformsthecharacterstoauniversalcharactersetcalled <u>network virtual terminal(NVT)</u> characters and delivers them to the local TCP/IP protocol stack.

The commands or text, in NVT form, travel through the Internet and arrive at the TCP/IP stack at the remote machine. Here the characters are delivered to the operating system and passed to the TELNET server, which changes the characters to the corresponding characters understandable by the remote computer. However, the characters cannot be passed directlyto the operating system because the remote operating system is not designedtoreceivecharactersfromaTELNETserver:Itisdesignedto receivecharactersfromaterminaldriver.Thesolutionistoaddapieceof

### ConceptofNVT(networkvirtual terminal)



26.75



Data Communications and Networking Fourth Edition



# WWWandHTTP

27.76

Copyright @The McGraw-Hill Companies, Inc. Permission required for reproduction or display.

### ARCHITECTURE

The WWW todayisadistributed client/serverservice, in which a client using a browser can access a service using a server. However, the service provided is distributed over many locations called sites as shown in fig.

**Topicsdiscussedinthissection:** 

Client(Browser)

Server

UniformResourceLocator

Cookies

## Figure27.1ArchitectureofWWW



27.79

Client(Browser)

A variety of vendors offer commercial browsers that interpret and displayaWebdocument, and all use nearly the same architecture.

Each browser usually consists of three parts: a controller, client protocol, and interpreters.

The controller receives input from the keyboard or the mouse and uses theclient programs to access thedocument.

After the document has been accessed, the controller uses one of the interpreters to display the document on the screen. The interpreter can be HTML, Java, or JavaScript, depending on the type of document

The client protocol can be one of the protocols described previously such as FTP or HTTP.

Server

The Web page is stored at the server. Each time a client request arrives, the corresponding document is sent to the client. To improve efficiency, servers normally store requested files in a cache in memory; memory is faster to access than disk. A server can also become more efficient through multithreading or multiprocessing.Inthiscase,aservercananswermorethanone

## Figure27.2Browser



27.80

UniformResourceLocator

A client that wants to access a Web page needs the address. To facilitate the access of documents distributed throughout the world, HTTP uses locators. The uniform resource locator (URL) is a standard for specifying any kind of information on the Internet. The URL defines four things: protocol, host computer, port, and path.

The protocol is the client/server program used to retrieve the document. Many different protocols can retrieve a document; among them are FTP or HTTP. The most common today is HTTP.

The host is the computer on which the information is located, although the name of the computer can be an alias. Web pages are usually stored in computers, and computers are given alias names that usually begin with the characters "www".

The URL can optionally contain the port number of the server. If the *port is* included, it is inserted between the host and the path, and it is separated from the host by a colon.

Path is the pathname of the file where the information is located. Note that the path can itself contain slashes that, in the UNIX operating system, separate the directories from the subdirectories and files.



AnHTTPcookie(alsocalledwebcookie,Internetcookie, browsercookieor simplycookie, the latter which is not to be confused with the literal definition), is a small piece of data sent from a website and stored in a user's web browser while the user is browsing that website The documents in the WWW can be grouped into three broad categories: static, dynamic, and active. The category is based on the time at which the contents of the document are determined.

<u>Topicsdiscussedinthissection:</u> Static Documents DynamicDocuments Active Documents

### StaticDocuments

Static documents are fixed-content documents that are created and stored in a server. The client can get only a copy of the document. When a client accesses the document, a copy of the document is sent. The user can then use a browsing program to display the document



Figure27.5Boldface tags

### HTML

HypertextMarkupLanguage(HTML)isalanguageforcreating Web pages.



# Figure27.7Beginningandendingtags

< TagName	Attribute = Value	Attribute = Value	••• >
a. Beginning tag			

< /TagName >	

b. Ending tag

DynamicDocuments

A dynamic document is created by a Web server whenever a browserrequeststhedocument.Whenarequestarrives,theWeb server runs an application program or a script that creates the dynamicdocument.Theserverreturnstheoutputoftheprogram or script as a response to the browser that requested the document. A very simple example of a dynamic document is the retrieval of the time and date from a server. Time and date are kinds of information that are dynamic in that they change from moment to moment. The client can ask the server to run a program such as the*dateprograminUNIXandsendtheresultoftheprogramto* theclient.

CommonGatewayInterface(CGI)

TheCommonGatewayInterface(CGI)isatechnologythat

createsandhandlesdynamicdocuments.

Hypertext Preprocessor (pHP), which uses the Perl language; Java Server Pages (JSP), which uses the Java language for scripting; Active Server Pages (ASP), a Microsoft product which uses Visual Basic language for scripting; and ColdFusion, which embeds SQL database queries in the HTML document

## Figure 27.8Dynamic documentusingCGI



## Figure27.9Dynamicdocumentusingserver-site script





Dynamic documents are sometimes referred to as server-site dynamic documents.

## Figure 27.10 Active documentusing Javaapplet

ActiveDocuments Formanyapplications, we need a program or a script to be run at the client site. These are called active documents


# Figure27.11Activedocumentusingclient-sitescript





Active documents are sometimes referred to as client-site dynamic documents.



HTTPversion1.1specifiesapersistentconnectionbydefault.

### DNS(DomainNameSystem)

To identify an entity, TCP/IP protocols use the IP address, which uniquely identifies the connection of a host to the Internet. However, people prefer to use names instead of numeric addresses. Therefore, we need a system that can <u>map a name to an address or an address to a name</u>.



# **NAMESPACE**

Aname space that maps each address to a unique name can be organized in two ways: fiat or hierarchical.

### FlatNameSpace

Inaflatnamespace, an ameisassigned to an address. An ameinthis space is a sequence of characters without structure.

#### Hierarchical NameSpace

In a hierarchical name space, each name is made of several parts. The first part can define the nature of the organization, the second part can define the name of an organization, the third part can define departments in the organization, and so on.

Exa:challenger.jhda.edu, challenger.smart.com

challenger.berkeley.edu, and

# DOMAINNAMESPACE

Tohaveahierarchicalnamespace, adomainnamespacewas designed. In this design the names are defined in an inverted-tree structure with the root at the top. The tree can have only 128 levels: level 0 (root) to level 127.



#### Label

Each node in the tree has a label, which is a string with a maximum of 63 characters. The root label is a null string (empty string). DNS requires that children of a node (nodes that branch from the same node) have different labels, which guarantees the uniqueness of the domain names. <u>DomainName</u>

Eachnodeinthetreehasadomainname.Afulldomainnameisa sequenceoflabelsseparatedbydots(.).Thedomainnamesarealways readfromthenodeuptotheroot.Thelastlabelisthelabeloftheroot(null). Thismeansthatafulldomainnamealwaysendsinanulllabel,which meansthelastcharacterisadotbecausethenullstringisnothing.Below Figure shows some domain names



challenger.atc.fhda.edu. cs.hmme.com. www.funny.int.



#### Domainnames andlabels



### Domain

A domain is a subtree of the domain name space. The name of the domain is the domain name of the node at the top of the subtree.



### DISTRIBUTIONOFNAME SPACE:

The information contained in the domain name space must be stored. However, it is very inefficient and also unreliable to have just one computer store such a huge amount of information. In this section, we discuss the distribution of the domain name space

#### <u>1HierarchyofName Servers</u>

distribute the information among many computers called DNS servers. we let the root stand alone and create as many domains (subtrees) as there are first-level nodes



#### <u>2 Zone</u>

Since the complete domain name hierarchy cannot be stored on a single server, it is divided among manyservers. What a server is responsible foror has authority over is called a zone. We can define a zone as a contiguous part of the entire tree



### <u>3Root Server</u>

A root server is a server whose zone consists of the whole tree. A root server usually does not store any information about domains but delegates its authorityto otherservers,keepingreferencestothoseservers.Thereare several root servers, each covering the whole domain name space. The servers are distributed all around the world.

#### **4Primaryand Secondary Servers**

Aprimary server is a server that stores a file about the zone for which it is an authority. It is responsible for creating, maintaining, and updating the zone file. It stores the zone file on a local disk

Asecondary server is a server that transfers the complete information about a zone from another server (primary or secondary) and stores the file on its local disk. The secondary server neither creates nor updates the zone files

# **DNSINTHEINTERNET**

DNS is a protocol that can be used in different platforms. In the Internet, the domain name space (tree) is divided into three different sections: genericdomains, country domains, and the inverse domain



# 1 GenericDomains

The generic domains define registered hosts according to their generic behavior. Each node in the tree defines a domain, which is an index to the domain name space database



Label	Description
aero	Airlines and aerospace companies
biz	Businesses or firms (similar to "com")
com	Commercial organizations
соор	Cooperative business organizations
edu	Educational institutions
gov	Government institutions
info	Information service providers
int	International organizations
mil	Military groups
museum	Museums and other nonprofit organizations
name	Personal names (individuals)
net	Network support centers
org	Nonprofit organizations
pro	Professional individual organizations

### 2 CountryDomains

The country domains section uses two-character country abbreviations(e.g., us for United States). Second labels can be organizational, or theycanbemorespecific, national designations. The United States, for example, uses state abbreviations as a subdivision of us (e.g., ca.us.).



#### 3InverseDomain

Theinverse domain isused tomapan address toaname.



# **RESOLUTION**

Mappinganametoanaddressoranaddresstoanameiscalledname-address resolution

# <u>1Resolver</u>

DNS is designed as a client/server application. A host that needs to map an address to a name or a name to an address calls a DNSclient called a resolver. The resolver accesses the closest DNSserver with a mapping request. If the server has the information, it satisfies the resolver; otherwise, it either refers the resolver to other servers or asks other servers to provide the information.

### 2MappingNamestoAddresses

In this case, the server checks the generic domains or the country domains to find the mapping.

### <u>3MappingAddressestoNames</u>

To answer queries of this kind, DNS uses the inverse domain <u>4</u> <u>Recursive Resolution</u>

The client (resolver) can ask for a recursive answer from a name server. This means that the resolver expects the server to supply the final answer.. When the query is finally resolved, the response travels back until it finally reaches the requesting client. This is called recursive resolution and is shown in FIG



### 5IterativeResolution

If the client does not ask for a recursive answer, the mapping can be done iteratively. If the server is an authority for the name, it sends the answer. If it is not, it returns (to the client) the IP address of the server that it thinks can resolve the query



### 6Caching

Each time a server receives a query for a name that is not in its domain, it needs to search its database for a server IP address. Reduction of this search time would increase efficiency. DNS handles this with a mechanismcalled caching

### DNS MESSAGES

DNS has two types of messages: query and response. Both types have the same format.

Thequerymessageconsistsofaheader,

andquestionrecords; the

response message consists of a header,

question records,

answer records,

authoritative records,

andadditionalrecords



#### Header

Both query and response messages have the same header format with some fields set to zero for the query messages. The header is 12 bytes,

Identification	Flags
Number of question records	Number of answer records (all 0s in query message)
Number of authoritative records (all 0s in query message)	Number of additional records (all 0s in query message)

### **TYPESOF RECORDS**

The question records are used in the question section of the query and response messages. The resource records are used in the answer, authoritative, and additional information sections of the response message.

#### QuestionRecord

Aquestionrecordisusedbytheclienttogetinformationfromaserver...

#### Resource Record

Each domain name (each node on the tree) is associated with a record called the resource record. The server database consists of resource records. Resource records are also what is returned by the server to the client.

# **REGISTRARS**

How are new domains added to DNS? This is done through a registrar, a commercial entity accredited by ICANN. A registrar first verifies that the requested domain name is unique and then enters it into the DNS database. Afeeischarged. Today, there are many registrars; their names and addresses can be found at http://www.intenic.net

# DYNAMICDOMAINNAMESYSTEM(DDNS)

In DNS, when there is a change, such as adding a new host, removing a host, or changing an IP address, the change must be made to the DNS master file. The size of today's Internet does not allow for this kind of manual operation.

The DNS master file must be updated dynamically. The Dynamic Domain Name System (DDNS) therefore was devised to respond to this need.

# **ENCAPSULATION**

DNS can use either UDPor TCP.In both cases the well-known port used by the server is port 53.