# CSCI 460 Networks and Communications

## **Network Layer**

#### Humayun Kabir

Professor, CS, Vancouver Island University, BC, Canada

## Outline

- Store and Forward Packet
   (Datagram) Switching
- Routers
- Routing Algorithms
  - Shortest Path Routing
  - Distance Vector Routing
  - Link State Routing
- Internet Protocol (IP)
  - IP Packet
  - IP Address, Subnet, and CIDR
  - Network Address Translation (NAT)
  - Internet Control Message Protocol (ICMP)

- Address Resolution Protocol (ARP)
- Dynamic Host Configuration Protocol (DHCP)

#### **The Network Layer**

Responsible for routing datagrams (packets) from source to destination network (eventually source to destination nodes) over multiple hops (networks).

Application
Transport
Network
Link
Physical

#### Store-and-Forward Packet Switching

# <u>Hosts</u> send <u>packets</u> into the network; packets are <u>forwarded</u> by <u>routers</u>



## **Connectionless Service – Datagrams**

Packet is forwarded using destination address inside it

Different packets may take different paths



## **Routing and Forwarding**

Routing is the process of discovering network paths

- Model the network as a graph of nodes and links
- Decide what to optimize (e.g., fairness vs efficiency)
- Build routing tables in each router
- Update routes for changes in topology (e.g., failures)

A's tab	ole (i	initia	lly)	A's ta	able (	(later)	C's	Tabl	е	E's1	Table	<b>;</b>
[	Α			Α		1	Α	Α		Α	С	
	В	В		В	В	1	В	Α		В	D	
[	С	С		С	С	]	С			С	С	
[	D	В		D	В		D	Е		D	D	
	Е	С		E	D		E	Е		Е		
	F	С		F	D	]	F	E		F	F	

## **Routing and Forwarding**

<u>Forwarding</u> is the sending of packets along a path using the routing table



### Flooding

A simple method to send a packet to all network nodes

Each node floods a new packet received on an incoming link by sending it out all of the other links

Nodes need to keep track of flooded packets to stop the explosion; even using a hop limit can blow up exponentially

## **Distance Vector Routing**

Distance vector is a distributed routing algorithm

- Shortest path computation is split across nodes
- Often used in Internet (RIP)

Algorithm:

- Each node knows distance of links to its neighbors
- Each node advertises vector of lowest known distances to all neighbors
- Each node uses received vectors to update its own
- Repeat periodically

#### **Distance Vector Routing**



Vectors received at J from Neighbors A, I, H and K

#### The Count-to-Infinity Problem

Failures can cause DV to "count to infinity" while seeking a path to an unreachable node



Good news of a path to *A* spreads quickly



Bad news of no path to A is learned slowly

<u>Dijkstra</u>'s algorithm computes a sink tree on the graph:

- Each link is assigned a non-negative weight/distance
- Shortest path is the one with lowest total weight
- Using weights of 1 gives paths with fewest hops

Algorithm:

- Start with sink, set distance at other nodes to infinity
- Relax distance to other nodes
- Pick the lowest distance node, add it to sink tree
- Repeat until all nodes are in the sink tree





C (9, B)

H̃ (∞, -)

)⊃D(∞,1)





A network and first five steps in computing the shortest paths from A to D. Pink arrows show the sink tree so far.



From A	Α	В	С	D	Е	F	G	Н
	0, -	2,A	∞, -	∞, -	∞, -	∞, -	6, A	∞, -
	•	D	C	D	Б	Б	C	II
From B	A	D		D	Ľ	Г	G	п
	0, -	2,A	9, B	∞, -	4, B	∞, -	6, A	∞, -
From F	Δ	R	C	D	F	F	G	н
FIOILIE	1	D	C	D		T	U	
	0, -	2,A	9, B	∞, -	4, B	6, E	5, E	∞, -
From G	Α	В	С	D	Ε	F	G	Н
	0, -	2,A	9, B	∞, -	4, B	6, E	5, E	9, G



From F	Α	В	С	D	E	F	G	Н
	0, -	2,A	9, B	∞, -	4, B	6, E	5, E	8, F
From H	Α	В	С	D	E	F	G	Н
	0, -	2,A	9, B	10,H	4, B	6, E	5, E	8, F
From C	Α	В	С	D	Е	F	G	Н
	0, -	2,A	9, B	10,H	4, B	6, E	5, E	8, F
From D	Α	В	С	D	E	F	G	Н
	0, -	2,A	9, B	10,H	4, B	6, E	5, E	8, F

#### Routing Table from Shortest Path Algorithm



A's Routing Table

Destination	Cost	Next Hop
А	0	-
В	2	В
С	9	В
D	10	В
E	4	В
F	6	В
G	5	В
Н	8	В

## Link State Routing

Link state is an alternative to distance vector

- More computation but simpler dynamics
- Widely used in the Internet (OSPF, ISIS)

Algorithm:

- Each node floods information about its neighbors in LSPs (Link State Packets); all nodes learn the full network graph
- Each node runs Dijkstra's algorithm to compute the path to take for each destination

#### Link State Routing – LSPs

LSP (Link State Packet) for a node lists neighbors and weights of links to reach them



ŀ	4	E	3	0	0		)	E		F	
Se	q.	Se	eq.	Se	eq.	Se	eq.	Se	eq.	Se	q.
A	ge	Ag	ge	Ag	ge	Aç	ge	Ag	je	Ag	je
В	4	Α	4	В	2	С	3	Α	5	В	6
Е	5	С	2	D	3	F	7	С	1	D	7
		F	6	F	1			F	8	F	8

Network

LSP for each node

## Link State Routing – Reliable Flooding

Seq. number and age are used for reliable flooding

- New LSPs are acknowledged on the lines they are received and sent on all other lines
- Example shows the LSP database at router B

			Ser	nd fla	ags ACK flags			gs	
Source	Seq.	Age	Á	С	F	Á	С	F	Data
А	21	60	0	1	1	1	0	0	
F	21	60	1	1	0	0	0	1	
E	21	59	0	1	0	1	0	1	
С	20	60	1	0	1	0	1	0	
D	21	59	1	0	0	0	1	1	

#### Internet Protocol (IP)

# Internet is an interconnected collection of many networks that is held together by the IP protocol



#### **IP Version 4 Protocol Header**

IPv4 (Internet Protocol) header is carried on all packets and has fields for the key parts of the protocol:



#### **IP** Addresses – Prefixes

Addresses are allocated in blocks called prefixes

- Prefix is determined by the network portion
- Has 2<sup>L</sup> addresses aligned on 2<sup>L</sup> boundary
- Written address/length, e.g., 18.0.31.0/24



 Prefix length is sometime represented by subnet mask, e.g., 18.0.31.0, 255.255.255.0

#### IP Addresses – Classful Addressing

Old addresses came in blocks of fixed size (A, B, C)

- Carries size as part of address, but lacks flexibility
- Called class full (vs. classless) addressing



#### IP Addresses – CIDR

- In Classless Inter-Domain Routing (CIDR), IP address prefixes are not of fixed sizes but varying sizes.
- Routers have the corresponding prefix information of an IP address to make routing decision, e.g., 128.120.0.0/16 or 128.120.0.0, 255.255.0.0

#### CIDR IP Addresses – Subnets

Subnetting splits up IP prefix to help with management

• Looks like a single prefix outside the network



Network divides it into subnets internally

#### IP Addresses – CIDR

- 128.208.0.0/16
- 16 bit network prefix 128.208
- 16 bits left for host numbering, 2<sup>16</sup> hosts if placed in a single network.
- 1 or more left most host bits are used to divide the network into subnets.
- Using **1 left most host** bit **two subnets** are created as follows and one subnet is **assigned** to **CS** department.

128.208. <b>0</b> 0000000.00000000	128.208. <b>0</b> .0/17	Other
128.208.10000000.00000000	128.208. <b>1</b> 28.0/17	CS

#### IP Addresses – CIDR

 Using 1 left most host bit of the Other subnet two more subnets are created as follows and one subnet is assigned to EE department.

<b>128.208.00</b> 000000.00000000	128.208. <b>0</b> .0/18	EE
128.208. <b>01</b> 000000.000000000	128.208. <b>64</b> .0/18	<b>New Other</b>

 Using 1 left most host bit of the New Other subnet two more subnets are created as follows and one subnet is assigned to EE department.

128.208. <b>010</b> 00000.000000000	128.208. <b>64</b> .0/19	Available
128.208. <b>011</b> 00000.000000000	128.208. <b>96</b> .0/19	Arts

#### **CIDR IP Address: ISP Assignments**



#### **CIDR IP Address: ISP Assignments**

- 192.24.0.0/16
- 16 bit network prefix 192.24
- 16 bits left for host numbering, 2<sup>16</sup> hosts if placed in a single network.
- Using 4 left most host bit 2<sup>4</sup> or 16 subnets are created as follows and the second subnet is assigned to Oxford.

192.24. <b>0000</b> 0000.00000000	192.24. <b>0</b> .0/20	First subnet
192.24. <b>0001</b> 0000.00000000	192.24. <b>1</b> 6.0/20	Oxford
192.24. <b>0010</b> 0000.00000000	192.24. <b>32.0/20</b>	Third subnet

192.24.**1111**0000.00000000 192.24.**240**.0/20

Sixteenth subnet

#### **CIDR IP Address: ISP Assignments**

 Using 1 left most host bit of the first subnet two more subnets are created as follows and the first subnet is assigned to Cambridge.

192.24. <b>00000</b> 000.00000000	192.24. <b>0</b> .0/21	Cambridge
192.24. <b>00001</b> 000.00000000	192.24. <b>8.0/21</b>	Second sub-subnet

 Using 1 left most host bit of the second sub-subnet two more subnets are created as follows and one sub-sub-subnet is assigned to Edinburgh and the other sub-sub-subnet is left available.

192.24. <b>000010</b> 00.00000000	192.24. <b>8.0/22</b>	Edinburgh
192.24. <b>000011</b> 00.00000000	192.24. <b>12</b> .0/22	Available

## **CIDR IP Addresses – Aggregation**

CIDR enables route aggregation to reduce the number of routing table entries. Aggregation joins multiple IP prefixes into a single larger prefix to reduce routing table size.



ISP customers have different prefixes

#### **CIDR IP Address: Aggregation**

192.24. <b>00000</b> 000.000000000	192.24. <b>0</b> .0/21	Cambridge
192.24. <b>000</b> 10000.000000000	192.24. <b>1</b> 6.0/20	Oxford
192.24. <b>000010</b> 00.000000000	192.24. <b>8.0/22</b>	Edinburgh

Aggregated Prefix by London router

A router may have multiple entries with common prefix but with different prefix lengths.

Packets are forwarded to the entry with the longest matching prefix

Complicates forwarding but adds flexibility



- Oxford, Cambridge, and Edinburgh subnets are connected by London router.
- London router aggregates 3 subnets into 192.24.0.0/19 aggregated prefix and advertises it to New York router.
- One of the subnet 192.24.12.0/22 available before is assigned to San Francisco. San Francisco router advertises it to New York router.
- **New York** router has two entries having with common prefix but varying prefix lengths (19 and 22).

192.24. <b>000</b> 00000.00000000	192.24.0.0/19	London
192.24. <b>000011</b> 00.00000000	192.24. <b>12.0/22</b>	San Francisco

192.24.00000000000000192.24.0.0/19London192.24.00001100.00000000192.24.12.0/22San Francisco

- When a packet comes to New York router its destination address is compared with the longest prefix entry 192.24.12.0/22 first and forward the packet to San Francisco if matches.
- New York router compares packet destination address with the next longest prefix entry 192.24.0.0/19 next and forward the packet to London if matches.

- If a packet comes with the **destination address** 192.24.15.252
- It is compared with the **longest prefix entry** 192.24.12.0/22 first and matches.

 192.24.000011
 00.00000000
 192.24.12.0/22
 San Francisco

 192.24.000011
 11.1111100
 192.24.15.252

• The packet is forwarded to **San Francisco**.

- If another packet comes with the **destination address** 192.24.20.248
- It is compared with the **longest prefix entry** 192.24.12.0/22 first and **does not match**.

192.24.00001100.00000000192.24.12.0/22San Francisco192.24.00010100.11111000192.24.20.248

 It is compared with the next longest prefix entry 192.24.0.0/19 next and matches.

• The packet is forwarded to **London**.

#### IP Addresses – NAT

NAT (Network Address Translation) box maps one external IP address to many internal IP addresses

- Uses TCP/UDP port to tell connections apart
- Violates layering; very common in homes, etc.



#### IP Addresses – NAT

Three range of IP addresses are declared private.

10.0.0.0	- 10.255.255.255/8	(16,777,216 hosts)
172.16.0.0	- 172.31.255.255/12	(1,048,576 hosts)
192.168.0.0	- 192.168.255.255/16	(65,536 hosts)

Internet routers do not forward any IP packet with these private address as the destination.

#### IP Addresses – NAT

NAT table is used for translating private-to-public and public-to-private IP addresses.

Index	Source Port	Source IP
0		
1		
2		
3344	5544	10.0.0.1

#### **Internet Control Protocols**

IP works with the help of several control protocols:

- <u>ICMP</u> is a companion to IP that returns error info
  - Required, and used in many ways, e.g., for traceroute
- <u>ARP</u> finds Ethernet address of a local IP address
  - Glue that is needed to send any IP packets
  - Host queries an address and the owner replies
- <u>DHCP</u> assigns a local IP address to a host
  - Gets host started by automatically configuring it
  - Host sends request to server, which grants a lease

#### **ICMP**

The Internet Control Message Protocol (ICMP) is a helper protocol that supports IP with facility for

- Error reporting
- Simple queries

#### ICMP messages are encapsulated as IP datagrams:



## **ICMP Error Reporting**

If, in the destination host, the IP module cannot deliver the datagram because the indicated protocol module or process port is not active, the destination host may send a destination unreachable message to the source host.



#### **ICMP** Request and Reply

Ping's are handled directly by the kernel

Each Ping is translated into an Echo Request

The Ping'ed host responds with an Echo Reply



#### **ICMP** Time Exceeded

If IP packet's TTL reaches to zero



ARP (Address Resolution Protocol) operates below the network layer as a part of the interface between the network (IPv4) and the data link layer (Ethernet).

ARP lets nodes find target Ethernet addresses from their IP addresses.



- ARP **Request** and **Reply** messages (broadcast).
- ARP Cache
  - Reduces ARP broadcast.
- Gratitude ARP
  - ARP request against self IP address
- **Proxy** ARP
  - Gateway (router) machine replies against out of network IP address.



Destination IP	<b>Destination Ethernet</b>
192.32.65.5	E2
192.32.65.1	E3 (router/gateway)
192.32.63.3	E3 (proxy)
192.32.63.8	E3 (proxy)

ARP Cache at Host 1 in CS Network



Frame	Source IP	Source Eth.	Destination IP	Destination Eth.
Host 1 to 2, on CS net	IP1	E1	IP2	E2
Host 1 to 4, on CS net	IP1	E1	IP4	E3
Host 1 to 4, on EE net	IP1	E4	IP4	E6

#### DHCP





## Summary

- Store and Forward Packet
   Switching
- Datagrams
- Routers
- Routing Algorithms
  - Shortest Path Routing
  - Distance Vector Routing
  - Link State Routing
- Internet Protocol (IP)
  - IP Packet
  - IP Address, Subnet, and CIDR
  - Network Address Translation (NAT)

- Internet Control Message Protocol (ICMP)
- Address Resolution Protocol (ARP)
- Dynamic Host Configuration Protocol (DHCP)

#### Next

#### Transport Layer

- User Datagram Protocol (UDP)
- Transport Control Protocol (TCP)
  - TCP Segment Header
  - TCP Connection
  - TCP Flow Control
  - TCP Congestion Control
- TCP Retransmission Timer