

Network Layer

November 5, 2019

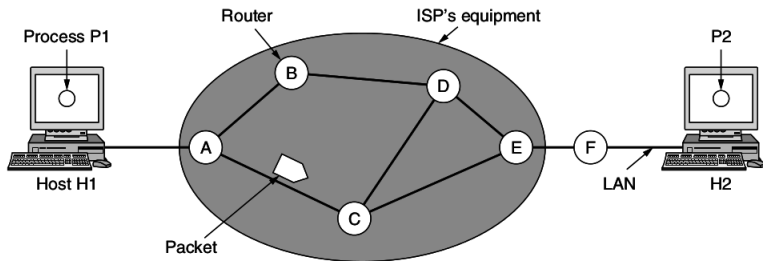
- Getting packets from the source all the way to the destination.
- Getting to the destination may require making many hops at intermediate routers along the way.

How to achieve this?

- To achieve its goals, N/W layer must know about the topology of the network (i.e., the set of all routers and links) and choose appropriate paths through it.
- Routes chosen should avoid overloading some of the communication lines and routers while leaving others idle.
- Finally, when the source and destination are in different networks, new problems occur. It is up to the network layer to deal with them.

- Basically we are concerned with design of network layer and how/what services are provided to transport layer.

Store and Forward Packet Switching



- A host with a packet to send transmits it to the nearest router, either on its own LAN or over a point-to-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.

Services provided to Transport Layer

- The network layer provides services to the transport layer at the network layer/transport layer interface.
- Question/Dilemma: Should N/W layer provide connection oriented service or connectionless service.

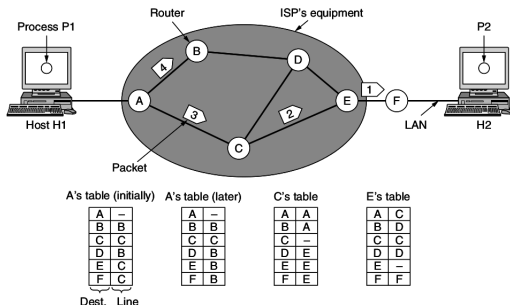
- Argument: Routers' job is moving packets around and nothing else.
- The network is inherently unreliable, no matter how it is designed. Therefore, the hosts should accept this fact and do error control (i.e., error detection and correction) and flow control themselves.
- Should comprise SEND PACKET and RECEIVE PACKET primitives, and each packet must carry the full destination address, because each packet sent is carried independently of its predecessors, if any.

Connection oriented service

- Argument: Success of telephone network. Reliable and provides QoS.
- X.25, Frame Relay (CO): ARPANET, Internet (CL).

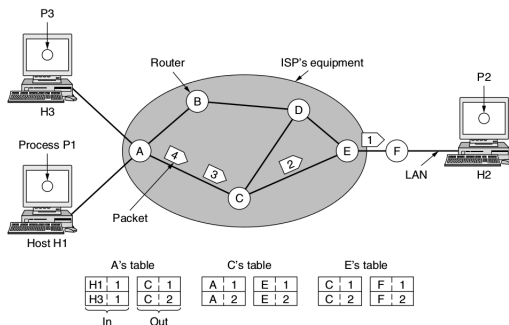
Implementation of Connection Less Service

- Packets are injected into the network individually and routed independently of each other. No advance setup is needed.
- Packets are frequently called datagrams (in analogy with telegrams) and the network is called a datagram network.
- The algorithm that manages the tables and makes the routing decisions is called the routing algorithm.



Implementation of Connection Oriented Service

- Path (VC) from the source router to destination router must be established before any data packets can be sent. Each packet carries an identifier telling which virtual circuit it belongs to.
- 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.



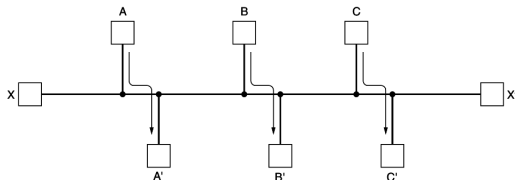
VC & Datagram Networks

- Resources (e.g., buffers, band-width, and CPU cycles) can be reserved in advance, when the connection is established.

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

Routing Algorithms

- Routing algorithm is that part of the network layer software responsible for deciding which output line an incoming packet should be transmitted on.
- One can think of a router as having two processes inside it. One of them handles each packet as it arrives, looking up the outgoing line to use for it in the routing tables. This process is forwarding.
- The other process is responsible for filling in and updating the routing tables.
- Desirable properties of a routing algorithm: correctness, simplicity, robustness, stability, fairness, and efficiency.



Classes of Routing Algorithms

- Non-adaptive: Routes are computed in advance, offline, and downloaded to the routers when the network is booted(static routing).
- Adaptive algorithms: Change their routing decisions to reflect changes in the topology, and sometimes changes in the traffic as well.
- Dynamic routing algorithms differ in where they get their information, when they change the routes, and what metric is used for optimization

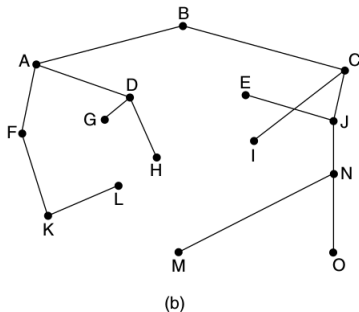
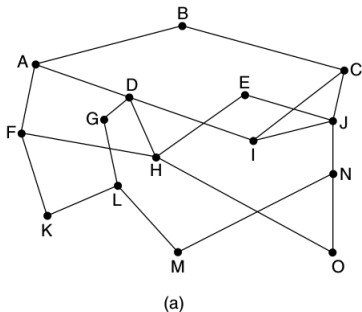
Optimality Principle

- 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 route.

Optimality Principle

- 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 route.
- To see this, call the part of the route from I to J r_1 and the rest of the route r_2 .
- If a route better than r_2 existed from J to K, it could be concatenated with r_1 to improve the route from I to K, contradicting our statement that $r_1 r_2$ is optimal.

Optimality Principle



- Based on optimality principle, the set of optimal routes from all sources to a given destination form a tree rooted at the destination (**sink tree**) where the distance metric is the number of hops.
- The goal of all routing algorithms is to discover and use the sink trees for all routers.

Optimality Principle

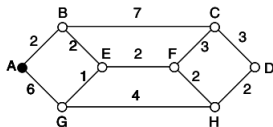
- A sink tree is not necessarily unique; other trees with the same path lengths may exist.
- Since a sink tree is indeed a tree, it does not contain any loops, so each packet will be delivered within a finite and bounded number of hops.

Shortest Path Algorithm

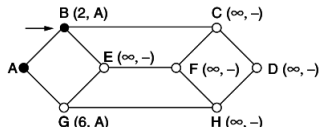
Let us compute optimal paths given a complete picture of the network.

- Dijkstra's algorithm (1959) finds the shortest paths between a source and all destinations in the network. Each node is labelled with its distance from the source node along the best known path.
- Initially, no paths are known, so all nodes are labelled with infinity.
- As the algorithm proceeds and paths are found, the labels may change, reflecting better paths.
- A label may be either tentative or permanent. Initially, all labels are tentative.
- When it is discovered that a label represents the shortest possible path from the source to that node, it is made permanent and never changed thereafter.

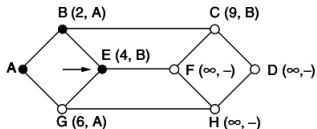
Shortest Path Algorithm



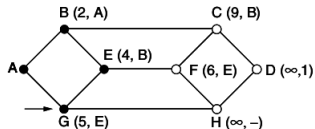
(a)



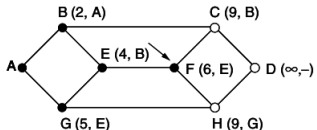
(b)



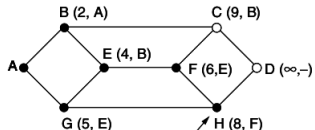
(c)



(d)



(e)



(f)

When a routing algorithm is implemented, each router must make decisions based on local knowledge, not the complete picture of the network.

- Flooding is the technique in which every incoming packet is sent out on every outgoing line except the one it arrived on.
- Flooding generates vast numbers of duplicate packets, in fact, an infinite number that can be limited by using hop-counter.
- A better technique for damming the flood is to have routers keep track of which packets have been flooded, to avoid sending them out a second time.

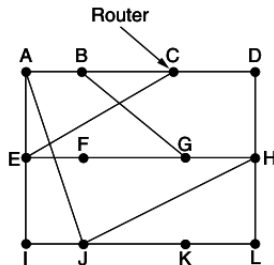
Distance Vector Routing

A distance vector routing algorithm operates by having each router maintain a table (i.e., a vector) giving the best known distance to each destination and which link to use to get there.

- These tables are updated by exchanging information with the neighbors. Eventually, every router knows the best link to reach each destination.
- It was the original ARPANET routing algorithm and was also used in the Internet under the name RIP.

DV Routing

Suppose that J has measured or estimated its delay to its neighbors, A, I, H, and K, as 8, 10, 12, and 6 msec, respectively.



					New estimated delay from J	
					↓ Line	
To	A	I	H	K		
A	0	24	20	21	8	A
B	12	36	31	28	20	A
C	25	18	19	36	28	I
D	40	27	8	24	20	H
E	14	7	30	22	17	I
F	23	20	19	40	30	I
G	18	31	6	31	18	H
H	17	20	0	19	12	H
I	21	0	14	22	10	I
J	9	11	7	10	0	–
K	24	22	22	0	6	K
L	29	33	9	9	15	K
					New routing table for J	

JA delay is 8
JI delay is 10
JH delay is 12
JK delay is 6

Vectors received from
its four neighbors

Count to Infinity Problem

- The settling of routes to best paths across the network is called convergence.
- DV reacts rapidly to good news, but leisurely to bad news.

A	B	C	D	E	
•	•	•	•	•	Initially
	1	•	•	•	After 1 exchange
	1	2	•	•	After 2 exchanges
	1	2	3	•	After 3 exchanges
	1	2	3	4	After 4 exchanges

A	B	C	D	E	
•	1	2	3	4	Initially
	3	2	3	4	After 1 exchange
	3	4	3	4	After 2 exchanges
	5	4	5	4	After 3 exchanges
	5	6	5	6	After 4 exchanges
	7	6	7	6	After 5 exchanges
	7	8	7	8	After 6 exchanges
	•	•	•	•	

Link State Routing

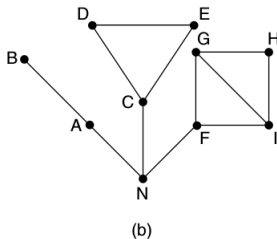
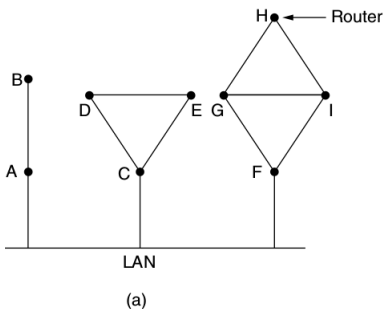
The idea behind link state routing is fairly simple and can be stated as five parts. Each router must do the following things to make it work:

- Discover its neighbors and learn their network addresses.
- Set the distance or cost metric to each of its neighbors.
- Construct a packet telling all it has just learned.
- Send this packet to and receive packets from all other routers.
- Compute the shortest path to every other router.

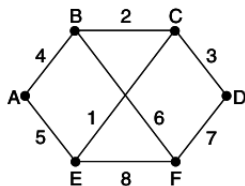
In effect, the complete topology is distributed to every router. Then Dijkstra's algorithm can be run at each router to find the shortest path to every other router.

Link State Routing: Discovering Neighbors

- It accomplishes this goal by sending a special HELLO packet on each point-to-point line. The router on the other end is expected to send back a reply giving its name.
- The router names must be globally unique.



Link State Routing: Building Link State Packets



(a)

		Link		State		Packets	
A		B		C		D	
Seq.		Seq.		Seq.		Seq.	
Age		Age		Age		Age	
B	4	A	4	B	2	C	3
E	5	C	2	D	3	F	7
		F	6	E	1		
						F	8
						E 8	

(b)

Link State Routing: Distributing Link State Packets

The fundamental idea is to use flooding to distribute the link state packets to all routers. To keep the flood in check, each packet contains a sequence number that is incremented for each new packet sent.

The Age field is decremented by each router during the initial flooding process, to make sure no packet can get lost and live for an indefinite period of time.

Source	Seq.	Age	Send flags			ACK flags			Data
			A	C	F	A	C	F	
A	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	
C	20	60	1	0	1	0	1	0	
D	21	59	1	0	0	0	1	1	

Figure: The packet buffer for router B

Link State Routing: Computing Routes

Once a router has accumulated a full set of link state packets, it can construct the entire network graph because every link is represented.

Dijkstra's algorithm can be run locally to construct the shortest paths to all possible destinations.

OSPF and IS-IS are some examples.

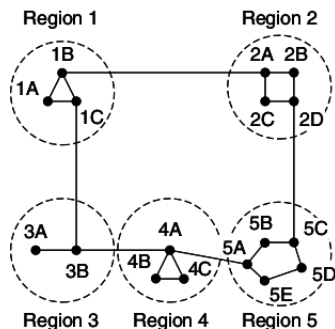
Hierarchical Routing

Can you guess why hierarchical routing?

Hierarchical Routing

- As networks grow in size, the router routing tables grow proportionally.
- More router memory and more CPU time 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.

Hierarchical Routing



(a)

Full table for 1A

Dest.	Line	Hops
1A	—	—
1B	1B	1
1C	1C	1
2A	1B	2
2B	1B	3
2C	1B	3
2D	1B	4
3A	1C	3
3B	1C	2
4A	1C	3
4B	1C	4
4C	1C	4
5A	1C	4
5B	1C	5
5C	1B	5
5D	1C	6
5E	1C	5

(b)

Hierarchical table for 1A

Dest.	Line	Hops
1A	—	—
1B	1B	1
1C	1C	1
2	1B	2
3	1C	2
4	1C	3
5	1C	4

(c)

Broadcast Routing

- Sending a packet to all destinations simultaneously is called broadcasting.
- Simple **broadcast** and **multi-destination** broadcasting are naive techniques.
- Flooding is one of the better broadcasting techniques. When implemented with a sequence number per source, flooding uses links efficiently with a decision rule at routers that is relatively simple.
- However, we can do better still once the shortest path routes for regular packets have been computed (**Reverse path forwarding**:)

Reverse path forwarding:

- When a broadcast packet arrives at a router, the router checks to see if the packet arrived on the link that is **normally used for sending packets toward the source of the broadcast**.
- If so, there is an excellent chance that the **broadcast packet itself followed the best route from the router** and is therefore the first copy to arrive at the router.
- This being the case, the router forwards copies of it onto all links except the one it arrived on.
- If, however, the broadcast packet arrived on a link other than the preferred one for reaching the source, the packet is discarded as a likely duplicate.

Reverse Path Forwarding

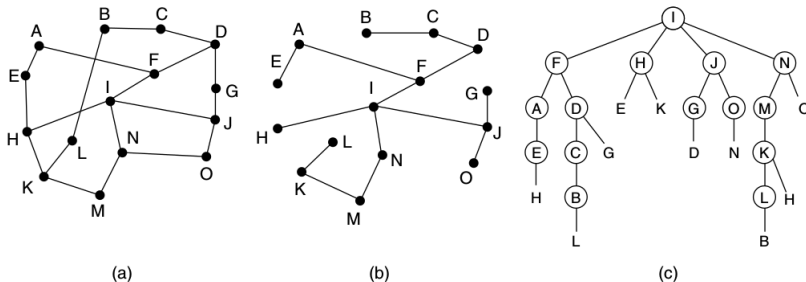


Figure 5-15. Reverse path forwarding. (a) A network. (b) A sink tree. (c) The tree built by reverse path forwarding.

Reverse path forwarding:

- The principal advantage of reverse path forwarding is that it is efficient while **being easy to implement**.
- It sends the broadcast packet over each link only once in each direction, **just as in flooding**, yet it requires only that routers know how to reach all destinations, without needing to remember sequence numbers (or use other mechanisms to stop the flood) or list all destinations in the packet.

Spanning Tree:

- A spanning tree is a subset of the network that includes all the routers but contains no loops. Sink trees are spanning trees.
- If each router knows which of **its lines belong to the spanning tree**, it can copy an incoming broadcast packet onto all the spanning tree lines except the one it arrived on.
- This method makes excellent use of bandwidth, generating the absolute minimum number of packets necessary to do the job.

Multicast Routing

- Some applications, such as a multiplayer game or live video of a sports event streamed to many viewing locations, send packets to multiple receivers.
- Unless the group is very small, **sending a distinct packet** to each receiver is **expensive**.
- On the other hand, broadcasting a packet is wasteful if the group consists of, say, **1000 machines on a million-node network**, so that most receivers are not interested in the message (or they are not supposed to see it).

Multicast Routing

- Sending a message to a group is called multicasting, and the routing algorithm used is called multicast routing.
- All multicasting schemes require some way to **create and destroy groups** and to **identify which routers are members** of a group.
- For now, we will assume that **each group is identified by a multicast address** and that **routers know the groups to which they belong**.

Multicast Routing

- Multicast routing schemes build on the broadcast routing schemes, sending packets along spanning trees to deliver the packets to the members of the group while making efficient use of bandwidth.
- However, the best spanning tree to use depends on whether the group is dense, with receivers scattered over most of the network, or sparse, with much of the network not belonging to the group.

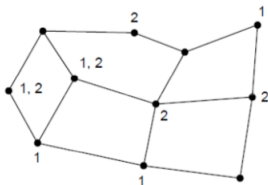
Multicast Routing

- If the group is dense, broadcast is a good start because it efficiently gets the packet to all parts of the network.
- But broadcast will reach some routers that are not members of the group, which is wasteful.
- The solution explored by **Deering and Cheriton (1990)** is to **prune the broadcast spanning tree** by removing links that do not lead to members.
- The result is an efficient multicast spanning tree.

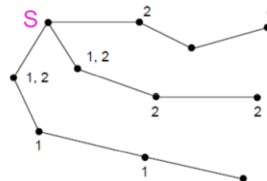
Multicast Routing

Multicast sends to a subset of the nodes called a group

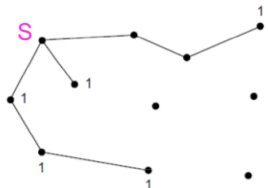
- Uses a different tree for each group and source



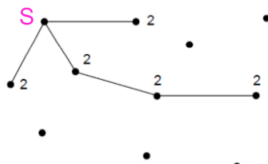
Network with groups 1 & 2



Spanning tree from source S



Multicast tree from S to group 1

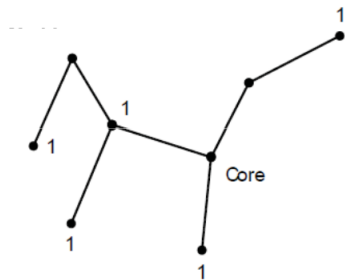


Multicast tree from S to group 2

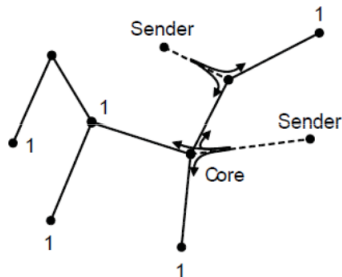
Multicast Routing

CBT (Core-Based Tree) uses a single tree to multicast

- Tree is the sink tree from core node to group members
- Multicast heads to the core until it reaches the CBT



Sink tree from core to group 1



Multicast is send to the core then down when it reaches the sink tree

Core based trees

- All of the routers agree on a root (**called the core or rendezvous point**) and build the tree by **sending a packet from each member to the root**.
- The tree is the union of the paths traced by these packets.
- To send to this group, a sender sends a packet to the core. When the packet reaches the core, it is forwarded down the tree.
- Packets destined for the group do not need to reach the core before they are multicast. As soon as a packet reaches the tree, it can be forwarded up toward the root, as well as down all the other branches.

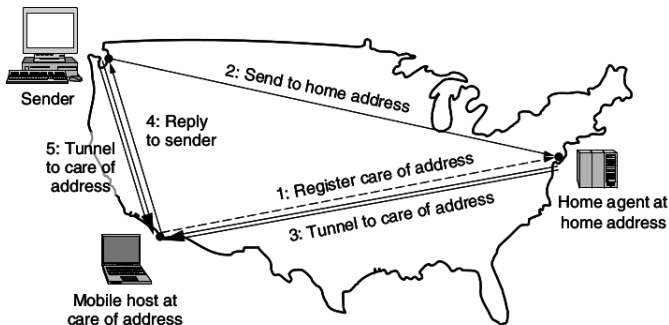
Core based trees

- Shared trees can be a major savings in storage costs, messages sent, and computation.
- Each router has to keep only one tree per group, instead of m trees.
- Shared tree approaches like core-based trees are used for multicasting to sparse groups in the Internet as part of popular protocols such as PIM (Protocol Independent Multicast).

Routing for Mobile Hosts

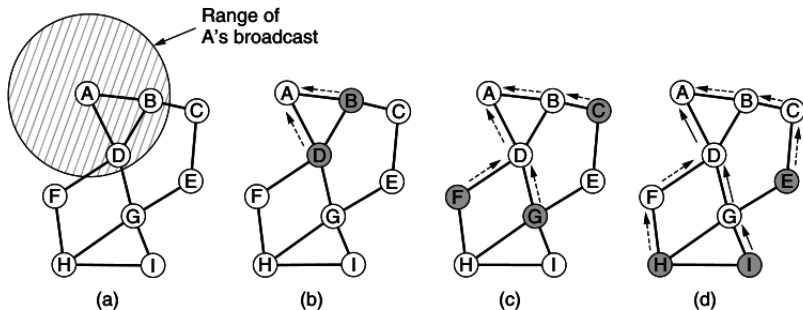
Mobile hosts introduce a new complication: to route a packet to a mobile host, the network first has to find it.

Assumption: all hosts are assumed to have a permanent home location that never changes.



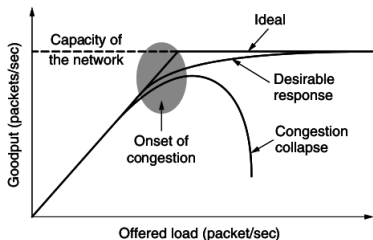
Routing in Ad-hoc Networks

Routers themselves are mobile



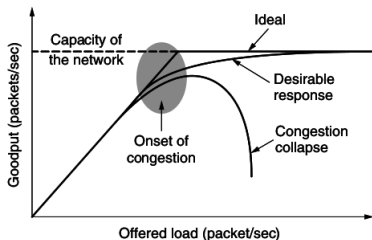
Congestion Control Algorithms

- 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.
- The most effective way to control congestion is to reduce the load that the transport layer is placing on the network.



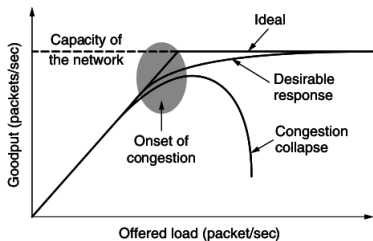
Congestion Control Algorithms

- 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.
- The most effective way to control congestion is to reduce the load that the transport layer is placing on the network.



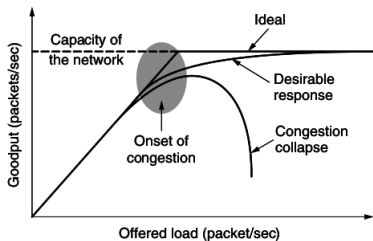
Congestion Control Algorithms

- 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.
- 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 number of delivered packets falls below the ideal curve. The network is now congested.



Congestion Control Algorithms

- **Congestion collapse**, in which performance plummets as the offered load increases beyond the capacity.
- This can happen because packets can be **sufficiently delayed inside the network** that they are no longer useful when they leave the network.
- A different failure mode occurs when **senders retransmit packets** that are greatly delayed, thinking that they have been lost. In this case, **copies of the same packet will be delivered by the network**, again wasting its capacity.



Congestion Control Algorithms

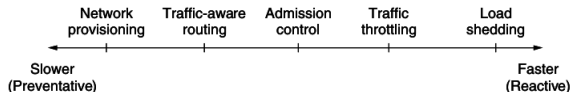
- Designing networks that avoid congestion where possible and do not suffer from congestion collapse if they do become congested.
- Adding more memory may help up to a point, if routers have an infinite amount of memory, congestion gets worse, not better.
- Low-bandwidth links or routers that process packets more slowly than the line rate can also become congested.

Congestion Control Algorithms

- Congestion control VS Flow Control:

Approaches to control congestion

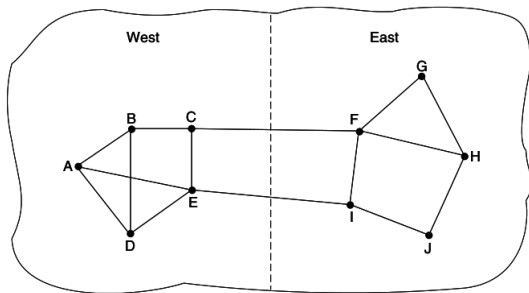
Two solutions come to mind: increase the resources or decrease the load.



Traffic Aware Routing

Routing based on load

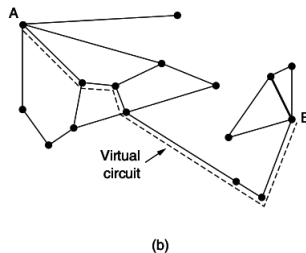
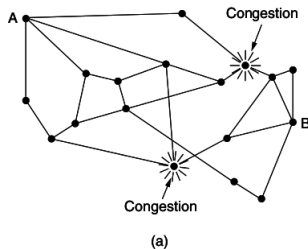
- But frequent oscillation of routes
- Multipath routes
- The routing scheme to shift traffic across routes slowly enough that it is able to converge.



Admission Control

Do not set up a new virtual circuit unless the network can carry the added traffic without becoming congested.

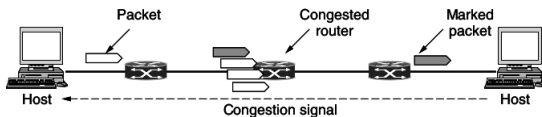
- Admission control can also be combined with traffic-aware routing by considering routes around traffic hotspots as part of the setup procedure.



Traffic Throttling

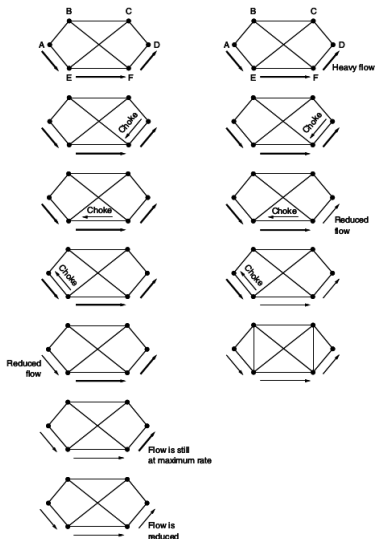
When congestion is imminent, the senders throttle back their transmissions and slow down. How to know?

- Each router can continuously monitor the resources it is using.
- The buffering of queued packets inside the router, and the number of packets that are lost due to insufficient buffering.
- Choke packets, explicit congestion notification



Traffic Throttling

Hop-by-hop back-pressure



Load Shedding

When routers can't take any more, they just throw packets away.

- What packets to drop?

Random Early Detection

The only reliable indication of congestion that hosts get from the network is packet loss.

- To determine when to start discarding, routers maintain a running average of their queue lengths. When the average queue length on some link exceeds a threshold, the link is said to be congested and a small fraction of the packets are dropped at random.
- How many many packets should be dropped?

Applications demand stronger performance guarantees from the network than "the best that could be done under the circumstances."

- Multimedia applications in particular, often need a minimum throughput and maximum latency to work.
- One way is over provisioning.

Four issues must be addressed to ensure quality of service:

- What applications need from the network.
- How to regulate the traffic that enters the network.
- How to reserve resources at routers to guarantee performance.
- Whether the network can safely accept more traffic.

QoS: Application Requirements

A stream of packets from a source to a destination is called a flow.

- The needs of each flow can be characterized by four primary parameters: bandwidth, delay, jitter, and loss.

QoS: Application Requirements

Different applications care about different properties

- We want all applications to get what they need

Application	Bandwidth	Delay	Jitter	Loss
Email	Low	Low	Low	Medium
File sharing	High	Low	Low	Medium
Web access	Medium	Medium	Low	Medium
Remote login	Low	Medium	Medium	Medium
Audio on demand	Low	Low	High	Low
Video on demand	High	Low	High	Low
Telephony	Low	High	High	Low
Videoconferencing	High	High	High	Low

“High” means a demanding requirement, e.g., low delay

Traffic Shaping

- Traffic shaping is a technique for regulating the average rate and burstiness of a flow of data that enters the network.
- In the telephone network, this characterization is simple. For example, a voice call needs 64 kbps and consists of one 8-bit sample every 125 μ sec.

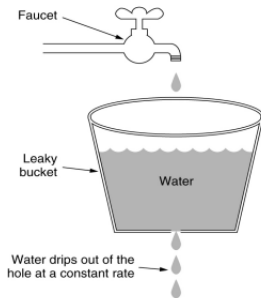
Traffic Shaping

- It is about regulating average rate of data flow.
- It is a method of congestion control by providing shape to data flow before entering the packet into the network.
- At connection set-up time, the sender and carrier negotiate a traffic pattern (shape)
- There are two types of Traffic shaping algorithm :-
 - 1. **Leaky Bucket Algorithm.**
 - 2. **Token Bucket Algorithm**

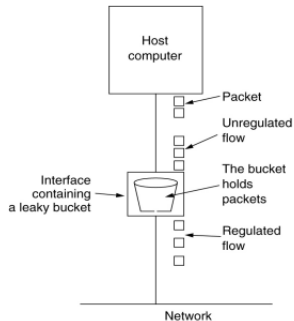
Leaky Bucket Algorithm

- The **Leaky Bucket Algorithm** used to control rate in a network.
- It is implemented as a single-server queue with constant service time.
- If the bucket (buffer) overflows then packets are discarded.
- In this algorithm the input rate can vary but the output rate remains constant.
- This algorithm saves bursty traffic into fixed rate traffic by averaging the data rate.

Leaky Bucket Algorithm



(a) A leaky bucket with water.

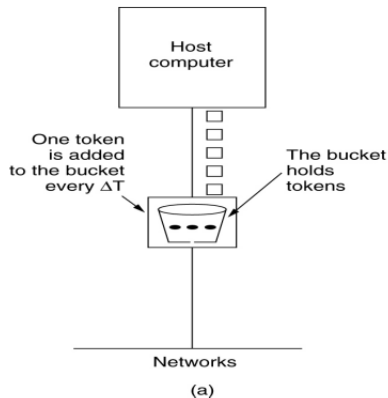


(b) a leaky bucket with packets.

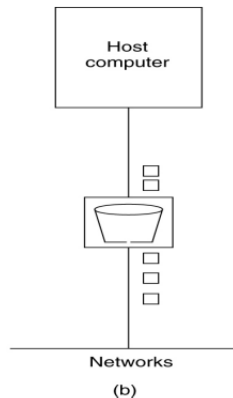
Token Bucket Algorithm

- The **Token Bucket Algorithm** compare to Leaky Bucket Algorithm allow the output rate vary depending on the size of burst.
- In this algorithm the buckets holds token to transmit a packet, the host must capture and destroy one token.
- Tokens are generated by a clock at the rate of one token every Δt sec.
- Idle hosts can capture and save up tokens (up to the max. size of the bucket) in order to send larger bursts later.

Token Bucket Algorithm



(a) Before



(b) After

Leaky/Token Buckets

- Leaky and token buckets limit the long-term rate of a flow but allow short-term bursts up to a maximum regulated length to pass through unaltered and without suffering any artificial delays.
- Large bursts will be smoothed by a leaky bucket traffic shaper to reduce congestion in the network.

Token Bucket Algorithm

A computer on a 6-Mbps network is regulated by a token bucket. The token bucket is filled at a rate of 1 Mbps. It is initially filled to capacity with 8 megabits. How long can the computer transmit at the full 6 Mbps?

Leaky/Token Buckets

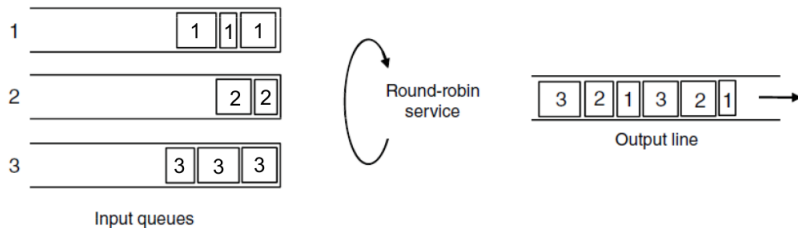
If we call the burst length **S** sec., the maximum output rate **M** bytes/sec, the token bucket capacity **B** bytes, and the token arrival rate **R** bytes/sec, we can see that an output burst contains a maximum of **B + RS** bytes. We also know that the number of bytes in a maximum speed burst of length *S* seconds is *MS*. Hence, we have **B + RS = MS**.

Packet Scheduling

- Using traffic shaping, we are able to regulate the shape of the offered traffic.
- However, to provide a performance guarantee, we must reserve sufficient resources along the route that the packets take through the network.
- Algorithms that allocate router resources among the packets of a flow and between competing flows are called packet scheduling algorithms.
- 3 different kinds of resources can potentially be reserved for different flows: 1. Bandwidth. 2. Buffer space. 3. CPU cycles.

Packet Scheduling

Packet scheduling divides router/link resources among traffic flows with alternatives to FIFO (First In First Out)

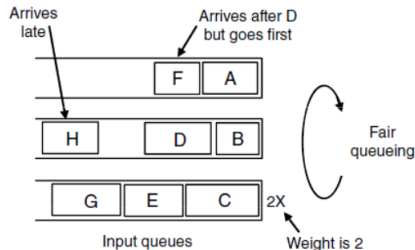


Example of round-robin queuing

Packet Scheduling

Fair Queueing approximates bit-level fairness with different packet sizes; weights change target levels

- Result is WFQ (Weighted Fair Queueing)



Packets may be sent
out of arrival order

Packet	Arrival time	Length	Finish time	Output order
A	0	8	8	1
B	5	6	11	3
C	5	10	10	2
D	8	9	20	7
E	8	8	14	4
F	10	6	16	5
G	11	10	19	6
H	20	8	28	8

$$F_i = \max(A_i, F_{i-1}) + L_i/W$$

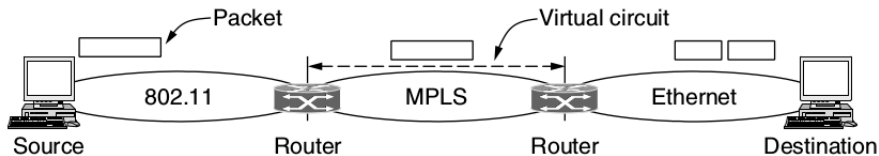
Finish virtual times determine
transmission order

- What issues arise when two or more networks are connected to form an internetwork, or more simply an internet?

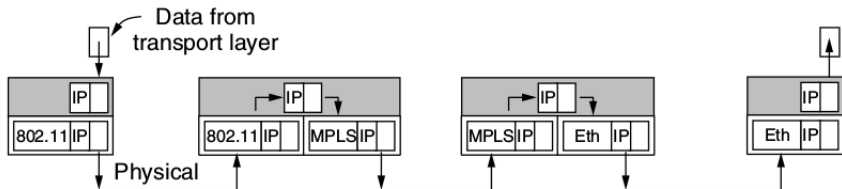
How networks differ?

Item	Some Possibilities
Service offered	Connectionless versus connection oriented
Addressing	Different sizes, flat or hierarchical
Broadcasting	Present or absent (also multicast)
Packet size	Every network has its own maximum
Ordering	Ordered and unordered delivery
Quality of service	Present or absent; many different kinds
Reliability	Different levels of loss
Security	Privacy rules, encryption, etc.
Parameters	Different timeouts, flow specifications, etc.
Accounting	By connect time, packet, byte, or not at all

How networks can be connected?



(a)

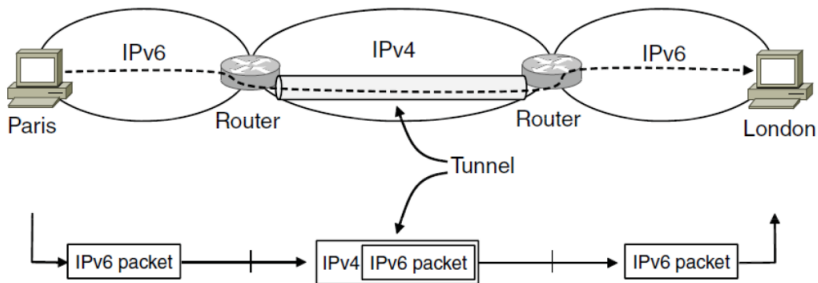


(b)

Tunneling

Connects two networks through a middle one

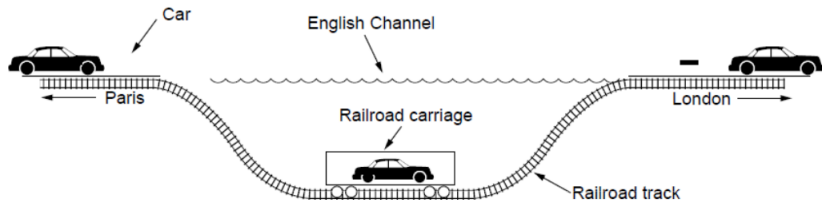
- Packets are encapsulated over the middle



Tunneling

Tunneling analogy:

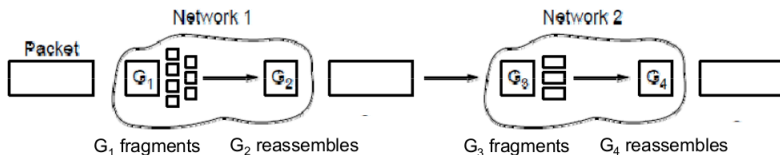
- tunnel is a link; packet can only enter/exit at ends



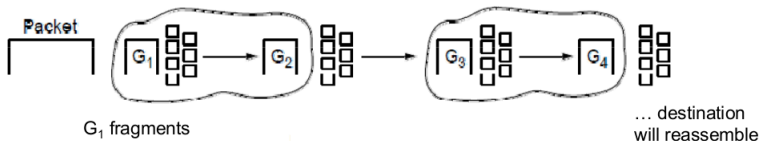
Packet Fragmentation

Networks have different packet size limits for many reasons

- Large packets sent with fragmentation & reassembly



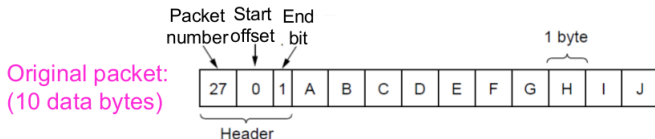
Transparent – packets fragmented / reassembled in each network



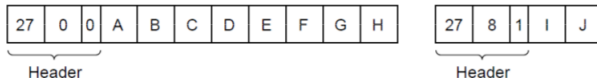
Non-transparent – fragments are reassembled at destination

Packet Fragmentation

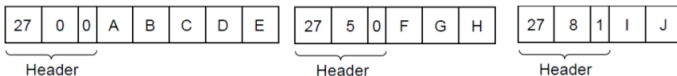
Example of IP-style fragmentation:



Fragmented:
(to 8 data bytes)



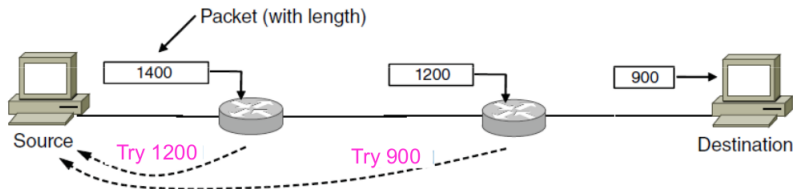
Re-fragmented:
(to 5 bytes)



Path MTU Discovery

Path MTU Discovery avoids network fragmentation

- Routers return MTU (Max. Transmission Unit) to source and discard large packets



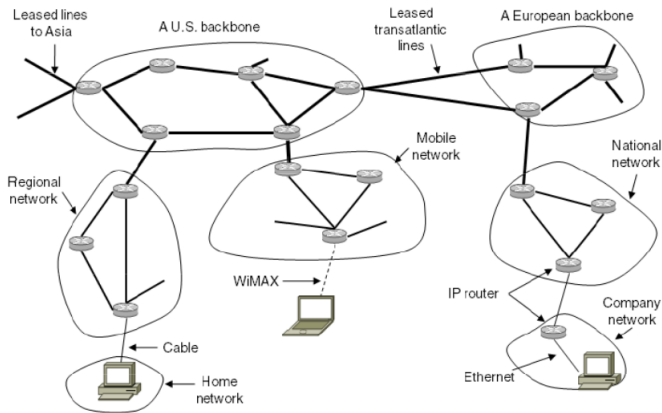
Network Layer in the Internet

IP has been shaped by guiding principles:

- Make sure it works
- Keep it simple
- Make clear choices
- Exploit modularity
- Expect heterogeneity
- Avoid static options and parameters
- Look for good design (not perfect)
- Strict sending, tolerant receiving
- Think about scalability
- Consider performance and cost

Network Layer in the Internet

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



IPv4 Protocol

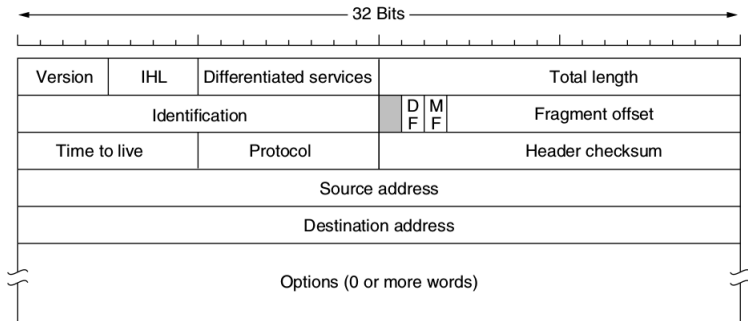


Figure 5-46. The IPv4 (Internet Protocol) header.

IPv4 Header Options

Option	Description
Security	Specifies how secret the datagram is
Strict source routing	Gives the complete path to be followed
Loose source routing	Gives a list of routers not to be missed
Record route	Makes each router append its IP address
Timestamp	Makes each router append its address and timestamp

Figure 5-47. Some of the IP options.

- A defining feature of IPv4 is its 32-bit addresses.
- Every host and router on the Internet has an IP address that can be used in the Source address and Destination address fields of IP packets.
- It is important to note that an IP address does not actually refer to a host. It really refers to a network interface, so if a host is on two networks, it must have two IP addresses.
- However, in practice, most hosts are on one network and thus have one IP address. In contrast, routers have multiple interfaces and thus multiple IP addresses.

- IP addresses are hierarchical, unlike Ethernet addresses.
- Each 32-bit address is comprised of a **variable-length network portion in the top bits** and a **host portion in the bottom bits**.
- The network portion has the same value for all hosts on a single network, such as an Ethernet LAN.
- This means that a **network** corresponds to a **contiguous block of IP address space**. This block is called a **prefix**.

- IP addresses are hierarchical, unlike Ethernet addresses.
- Each 32-bit address is comprised of a **variable-length network portion in the top bits** and a **host portion in the bottom bits**.
- The network portion has the **same value for all hosts on a single network**, such as an Ethernet LAN.
- This means that a **network** corresponds to a **contiguous block of IP address space**. This block is called a **prefix**.

- IP addresses are written in **dotted decimal notation**. In this format, each of the 4 bytes is written in decimal, from 0 to 255 (128.208.2.151).
- Prefixes are written by giving the **lowest IP address in the block** and the **size of the block**. The size is determined by the number of bits in the network portion; the remaining bits in the host portion can vary.
- By convention, it is written after the prefix IP address as a **slash(/)** followed by the length in bits of the network portion.
- If the prefix contains 2^8 addresses, it leaves 24 bits for the network portion, it is written as 128.208.2.0/24.

- Since the prefix length cannot be inferred from the IP address alone, routing protocols **must carry the prefixes** to routers.
- The length of the prefix corresponds to a binary mask of **1s** in the network portion.
- When written out this way, it is called a **subnet mask**.
- It can be ANDed with the IP address to extract only the network portion. For our example, the subnet mask is 255.255.255.0.

- **Address** - The unique number ID assigned to one host or interface in a network.
- A **network** corresponds to a **contiguous block of IP address space**. This block is called a **prefix**.
- **Subnet** - A portion of a network that shares a particular subnet address.
- **Subnet mask** - A 32-bit combination used to describe which portion of an address refers to the subnet and which part refers to the host.
- **Interface** - A network connection.

Hierarchical IP Addresses

- The key advantage of **prefixes** is that routers can forward packets based on only the **network portion** of the address, as long as each of the networks has a unique address block.
- The host portion does not matter to the routers because all hosts on the same network will be sent in the same direction.
- It is only when the packets reach the network for which they are destined that they are forwarded to the correct host.
- If the number of hosts on the Internet one billion, then by using a hierarchy, routers need to keep routes for only around 300,000 prefixes.

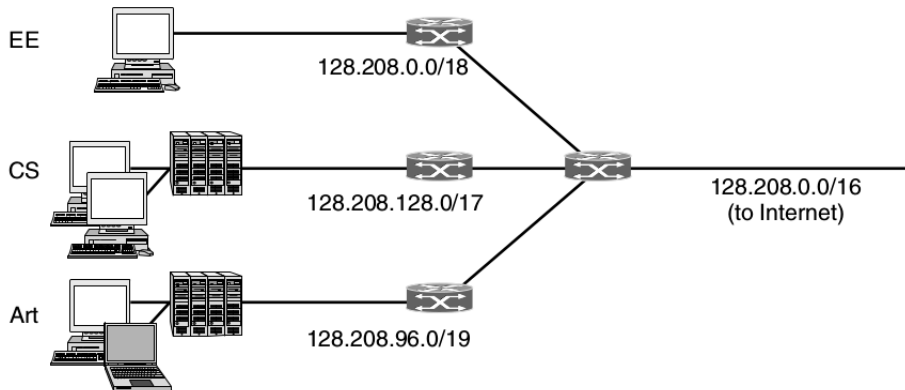
- Network numbers are managed by a non-profit corporation called **ICANN (Internet Corporation for Assigned Names and Numbers)**, to avoid conflicts.
- In turn, **ICANN** has delegated parts of the address space to various **regional authorities**, which dole out IP addresses to ISPs and other companies.

- Routing by prefix requires all the hosts in a network to have the same network number. This property can cause problems as networks grow.
- For example, consider a university that started out with **/16** prefix for use by the Computer Science Dept. for the computers on its Ethernet.
- A year later, the Electrical Engineering Dept. wants to get on the Internet. The Art Dept. soon follows suit. **What IP addresses should these departments use?**

- Getting further blocks requires going outside the university and may be expensive or inconvenient.
- Moreover, the /16 already allocated has enough addresses for over **60,000** hosts.
- It might be intended to allow for significant growth, but until that happens, it is wasteful to allocate further blocks of IP addresses to the same university. **A different organization is required.**

- The solution is to allow the block of addresses to be split into several parts for internal use as multiple networks, while still acting like a single network to the outside world.
- This is called **subnetting** and the networks (such as Ethernet LANs) that result from dividing up a larger network are called subnets.

subnets



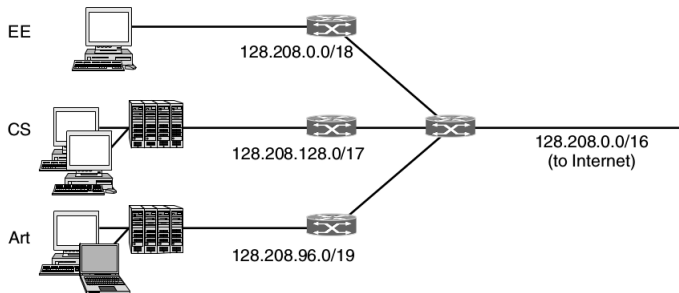
- The single **/16** has been split into pieces. This split does not need to be even, but each piece must be aligned so that any bits can be used in the lower host portion.

- A different way to see how the block was divided is to look at the resulting prefixes when written in binary notation:
- Here, the vertical bar (**—**) shows the boundary between the subnet number and the host portion.

Computer Science:	10000000	11010000	1 xxxxxxx	xxxxxxx
Electrical Eng.:	10000000	11010000	00 xxxxxx	xxxxxxx
Art:	10000000	11010000	011 xxxxx	xxxxxxx

subnets

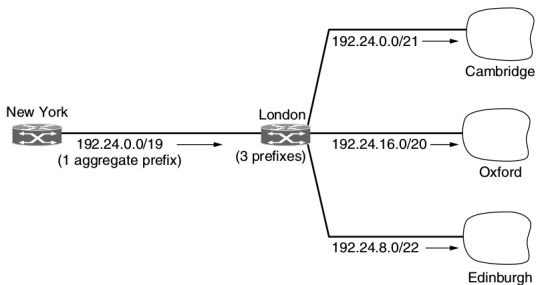
- When a packet comes into the main router, how does the router know which subnet to give it to?
- One way would be for each router to have a table with 65,536 entries telling it which outgoing line to use for each host on campus.
- But this would undermine the main scaling benefit we get from using a hierarchy. Instead, the routers simply need to know the subnet masks for the networks on campus.



- When a packet arrives, the router looks at the destination address of the packet and checks which subnet it belongs to.
- The router can do this by **ANDing** the **destination address** with the **mask for each subnet** and checking to see if the result is the corresponding prefix.
- For example, consider a packet destined for IP address 128.208.2.151. To see if it is for the Computer Science Dept., we AND with 255.255.128.0 to take the first 17 bits (which is 128.208.0.0) and see if they match the prefix address (which is 128.208.128.0). They do not match.

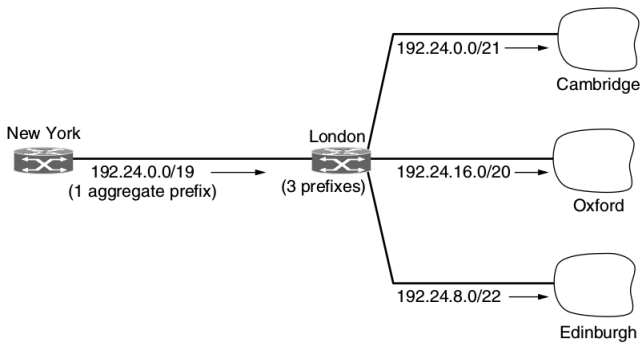
- Checking the first 18 bits for the Electrical Engineering Dept., we get 128.208.0.0 when ANDing with the subnet mask. This does match the prefix address, so the packet is forwarded onto the interface which leads to the Electrical Engineering network.
- The subnet divisions can be changed later if necessary, by updating all subnet masks at routers inside the university.

CIDR: Classless InterDomain Routing



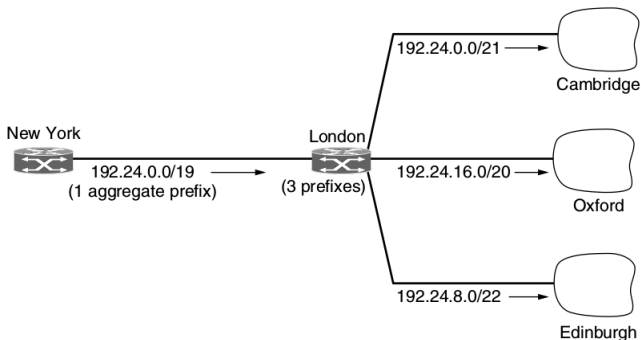
- Routers in organizations at the edge of a network, such as a university, need to have an entry for each of their subnets, telling the router which line to use to get to that network.
- For routes to destinations outside of the organization, they can use the simple default rule of sending the packets on the line toward the ISP that connects the organization to the rest of the Internet.

CIDR: Classless InterDomain Routing



- Routers in ISPs and backbones in the middle of the Internet must know which way to go to get to every network and no simple default will work.
- These core routers are said to be in the **default-free zone** of the Internet.

CIDR: Classless InterDomain Routing



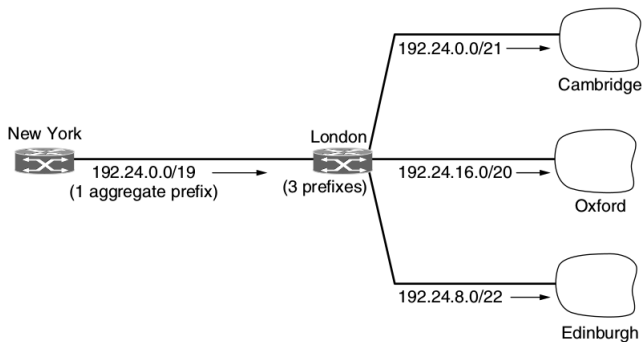
- **Solution:** Combine multiple small prefixes into a single larger prefix. This process is called route aggregation.
- This design works with subnetting and is called CIDR.

CIDR: Example

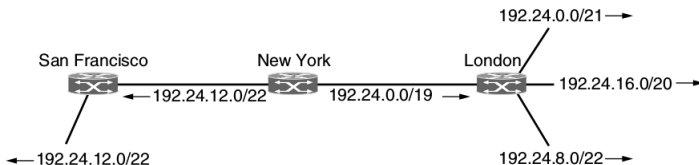
- Let us consider an example in which a block of 8192 IP addresses is available starting at 194.24.0.0.
- Suppose, Cambridge University needs 2048 addresses. Next, Oxford University asks for 4096 addresses. Finally, the University of Edinburgh asks for 1024 addresses.

University	First address	Last address	How many	Prefix
Cambridge	194.24.0.0	194.24.7.255	2048	194.24.0.0/21
Edinburgh	194.24.8.0	194.24.11.255	1024	194.24.8.0/22
(Available)	194.24.12.0	194.24.15.255	1024	194.24.12.0/22
Oxford	194.24.16.0	194.24.31.255	4096	194.24.16.0/20

- All of the routers in the default-free zone are now told about the IP addresses in the three networks.
- All of the IP addresses in the three prefixes should be sent from New York (or the U.S. in general) to London.



- Prefixes are allowed to overlap. The rule is that packets are sent in the direction of the most specific route, or the longest matching prefix.



Understanding IP Addresses

Before 1993, IP addresses were divided into the five categories. These are five different classes of networks, A to E.

Classful IP Addresses

Given an IP address, its class can be determined from the three high-order bits (the three left-most bits in the first octet).

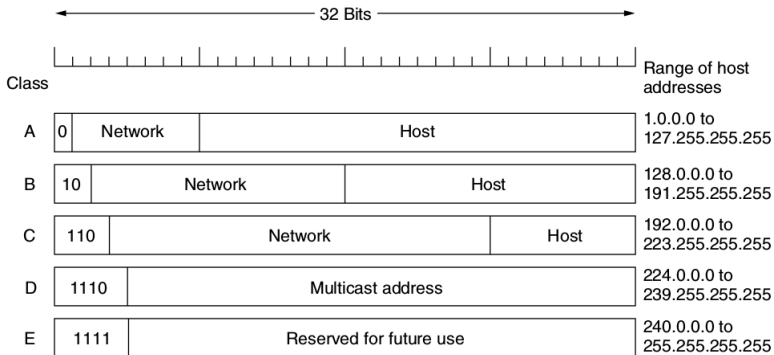


Figure 5-53. IP address formats.

Classful IP Addresses

- In a Class A address, the first octet is the network portion, so the Class A has a major network address of 1.0.0.0 - 127.255.255.255. Class A addresses are used for networks that have more than 65,536 hosts (actually, up to 16777214 hosts!).
- In a Class B address, the first two octets are the network portion, so the Class B has a major network address of 128.0.0.0 - 191.255.255.255. Class B addresses are used for networks that have between 256 and 65534 hosts.
- The Class C address has a major network address of 192.0.0.0 - 223.255.255.255. Octet 4 (8 bits) is for local subnets and hosts - perfect for networks with less than 254 hosts.

Understanding Subnetting

If you have network 172.16.0.0, then you know that its natural mask is 255.255.0.0 or 172.16.0.0/16. Extending the mask to anything beyond 255.255.0.0 means you are subnetting.

If you use a mask of 255.255.248.0 (/21), how many subnets and hosts per subnet does this allow for?