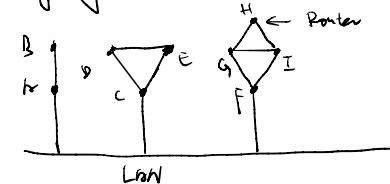


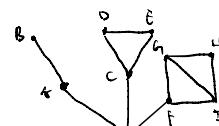
5.2.5.1 Learning about the Neighbors

• When a router is booted, its first task is to learn who its neighbors are. 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. These names must be globally unique.

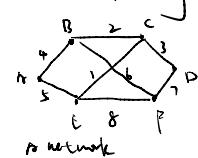


(a) Nine routers and a broadcast LAN



(b) A graph model of (a). Here we have introduced a new, artificial node N, to which A, C, and F are.

5.2.5.2 Building Link State Packets



Link	
A	B
B	C
C	D
D	E
E	F
F	G
G	H
H	I
I	J
J	K
K	L
L	M
M	N
N	P

to network

Link	
A	B
B	C
C	D
D	E
E	F
F	G
G	H
H	I
I	J
J	K
K	L
L	M
M	N
N	P

(b)

The link state packets for this network.

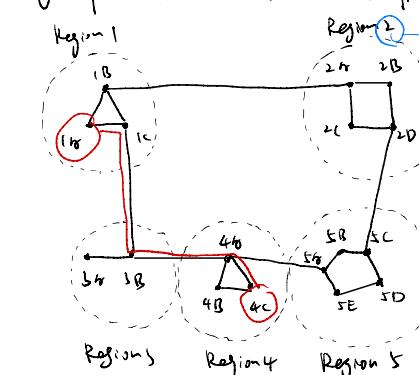
- The packet starts with the identity of the sender, followed by a sequence number and age and a list of neighbors. The cost to each neighbor is also given.
- The hard part is determining when to build packets. One possibility is to build them periodically at regular intervals. Another possibility when some specific event occurs, such as a line or neighbor going down or coming back up again or changing its properties.
- Drawback: 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.

5.2.6 Hierarchical Routing within a Network

• When hierarchical routing is used, the routers are divided into what we will call regions or areas.

• Drawback: There is a penalty to be paid: increased path length. For example, the best route from 1A to 5C is via region 2, but with hierarchical routing all traffic to region 5 goes via region 3, because that is better for most destinations in region 5.

e.g. If 1A \rightarrow 5C then it can read from hierarchical table 1A \rightarrow 1C $\rightarrow \dots \rightarrow$ 5C



1A \rightarrow 4C

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

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

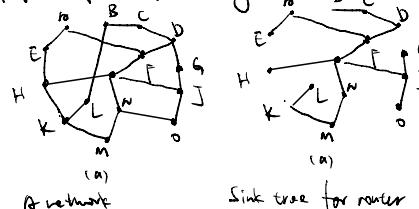
Bandwidth
(packets/sec)

5.2.7 Broadcast Routing

• sending a packet to all destinations simultaneously is called broadcasting.

- The principal advantage of reverse path forwarding is that it is efficient while 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 list all destinations.
- Drawback: Each router must have knowledge of some spanning tree.

reverse path forwarding:

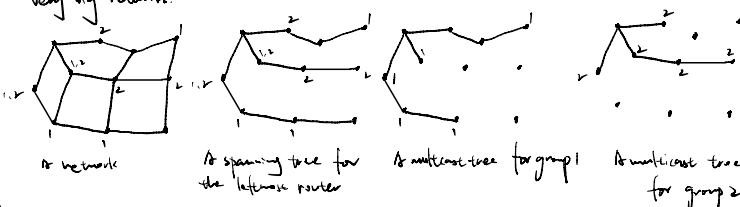


Sink tree for router

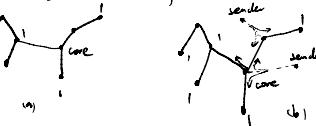
The tree built by reverse path forwarding from I.

5.2.8 Multicast Routing

- Sending a message to such a group is called multicasting, and the routing algorithm used is called multicast routing.
- The solution is to prune the broadcast spanning tree by removing links that do not lead to members. The result is an efficient multicast spanning tree.
- Drawback: One potential disadvantage is that it is lots of work for routers, especially for very big networks.

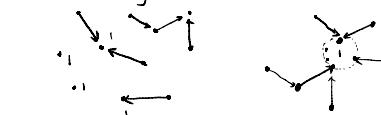


- An alternative design uses core-based trees to compute a single spanning tree for the group. All of the routers agree on a root (called the core) and build the tree by sending a packet from each member to the root. The tree is the union of the paths traced by these packets. Hence, to send to the group, a sender sends a packet to the core. When the packet reaches the core, it is forwarded down the tree.



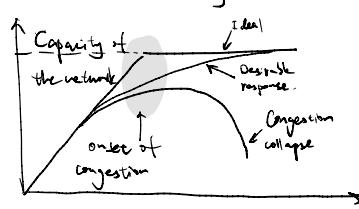
5.2.9 Anycast Routing

- In anycast, a packet is delivered to the nearest member of a group. Schemes that find these paths are called anycast routing.
- Suppose we want to anycast to the members of group 1. They will all be given the address "1", instead of different addresses. Distance Vector routing will distribute vectors as usual, and nodes will choose the shortest path to destination 1. This will result in nodes sending to the nearest instance of destination 1.



5.3.1 Congestion

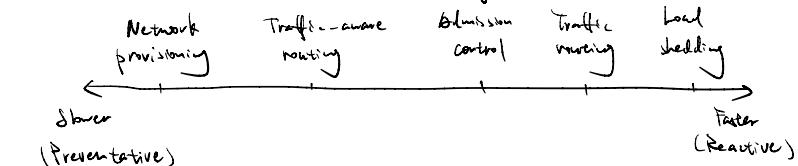
- Traffic management or engineering has to do with making sure the network is able to carry the offered traffic. It can be performed by devices in the network, or by the senders of traffic (often through mechanisms in the transport protocol), which are often referred to as congestion control.
- Congestion management and control concern the behavior of all the hosts and routers.
- Flow control relates to the traffic between a particular sender and a particular receiver and is generally concerned with making sure that the sender is not transmitting data faster than the receiver can process it.



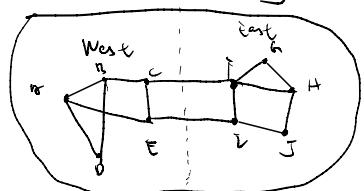
offered load (packets/sec)
Performance drops significantly in the presence of congestion:
packet loss rates increase,
and latency also increases
as router queues fill with packets

5.3.2 Approaches to Traffic management.

Timescales of approaches to traffic and congestion management



5.3.2.1 Traffic-Aware Routing

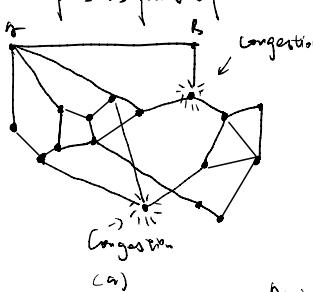


A network in which the East and West parts are connected by two links.

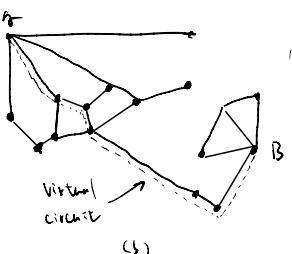
- The goal in taking load into account when computing routes is to shift traffic away from hotspots that will be the first places in the network to experience congestion.
- Set the link weight to be a function of the (fixed) link bandwidth and propagation delay plus the (variable) measured load or average queuing delay. least-weight paths will then favor paths that are more lightly loaded, all else being equal.
- Disadvantage: Routing tables may oscillate wildly, leading to erratic routing and many potential problems.

5.3.2.2 Admission Control

- The technique that is widely used in virtual-circuit networks to keep congestion at bay is admission control. The idea is simple: do not set up a new virtual circuit unless the network can carry the added traffic without becoming congested.
- Admission control can be combined with traffic-aware routing by considering routes around traffic hotspots as part of the setup procedure.

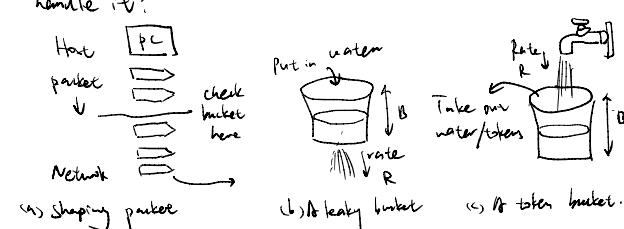


(a) A congested network (b) The portion of the network that is not congested.



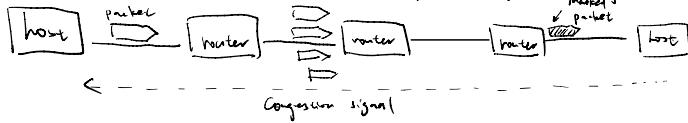
5.3.2.4 Traffic shaping

Traffic shaping is a technique for regulating the average rate and burstiness of a flow of data that enters the networks. The goal is to allow applications to transmit traffic that suits their needs. When a flow is set up, the user and the network agree on a certain traffic pattern for that flow. In effect, the user says to the provider "My transmission pattern will look like this: can you handle it?"



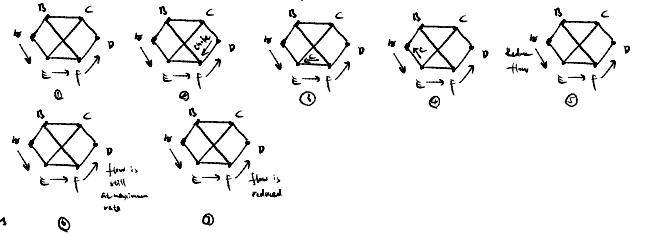
5.3.2.5 Explicit Congestion Notification

Instead of generating packets to warn of congestion, a sender can tag any packet it forwards (by setting a bit in the packet's header) to signal that it is experiencing congestion. When the network delivers the packet, the destination can note that there is congestion and inform the sender when it sends a reply packet. The design is called ECN (Explicit Congestion Notification).

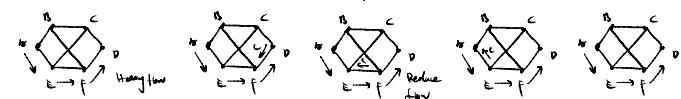


5.3.2.6 Hop-by-Hop Bandwidth

A choke packet that affects only the source.



A choke packet that affects each hop it passes through.



An alternative approach is to have the choke packet take effect at every hop it passes through, as shown in the sequence of (b). Here, as soon as the choke packet reaches F, F is required to reduce the flow to 0. Doing so will require F to devote more buffers to the connection, since the source is still sending away at full blast, but it gives immediate relief.

The next effect of this hop-by-hop scheme is to provide quick relief at the point of congestion, at the price of using up more buffers upstream.

5.4.1 Application QoS 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.
- Together, these four determine the QoS (Quality of Service) the flow requires.

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

Stringency of applications' quality-of-service requirements.

5.4.2 Overprovisioning

Helps to provide good quality of service

Ensures the network has the capacity for all traffic - expensive solution.

- Network uses to address for quality of service
- Addressing applications needs
- Regulating traffic entering the network
- Reserving resources at routers to guarantee performance
- Safely accepting more traffic

5.4.3 Packet Scheduling

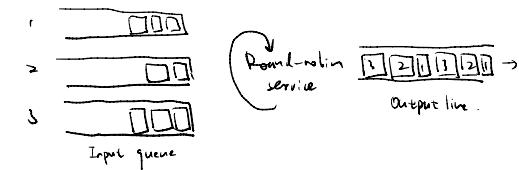
- Algorithms that allocate network resources among the packets of a flow and between competing flows are called packet scheduling algorithms.
- Routers reserve resources for different flows: Bandwidth, Buffer spaces, CPU cycles.
- Algorithms: FIFO, Fair queuing, weighted fair queuing.

5.4.3.1 FIFO

Each router buffers packets in a queue for each output line until they can be sent; they are sent in the same order that they arrived; This algorithm is known a FIFO.

Disadvantages: Processing packets in the order of their arrival means that the aggressive sender can hog most of the capacity of the routers its packets traverse, starving the other flows and reducing their quality of service.

5.4.3.2 Fair Queuing

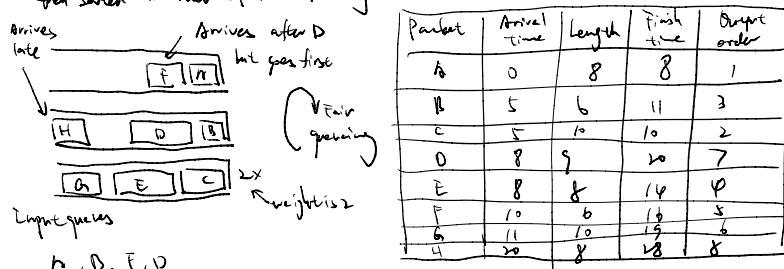


The essence of this algorithm is that routers have separate queues, one for each flow for a given output line; When the line becomes idle, the router scans the queues round robin. In this way, with n hosts competing for the output line, each host gets to send one out of every n packets; It's fair in the sense that all flows get to send packets at the same rate. Sending more packets will not improve this rate.

Dis: It gives more bandwidth to hosts that use large packets than to hosts that use small packets.

5.4.3.3 Improved Fair Queuing

- Fair robin is done in such a way as to simulate a byte-by-byte round robin, instead of a packet-by-packet round robin.
- The trick is to compute a virtual time that is the number of the round of which each packet would finish being sent. Each round derives a byte from all of the queues that have data to send. The packets are then sorted in order of their finishing times and sent in that order.



Dis: It gives all hosts the same priority.

5.4.3.4 Weighted Fair Queuing

In many situations, it is desirable to give, for example, video servers more bandwidth than, say, file servers.

This is easily by giving the video server two or more bytes per round. This modified algorithm is called WFA (Weighted Fair Queuing).

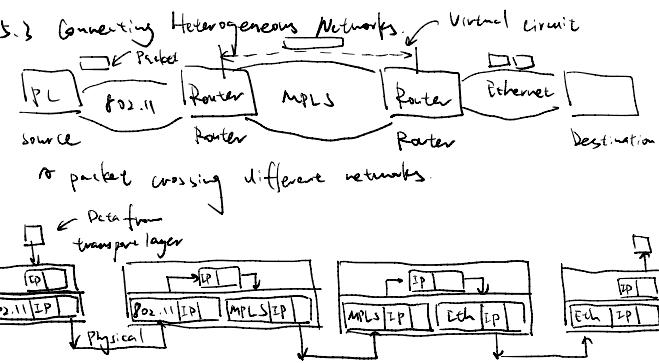
5.5.1 Internetworks:

- When two or more networks are connected, they form an internetwork, or more simply an Internet.
- How networks differ; Connecting heterogeneous networks;
- Connecting endpoints across heterogeneous networks; Supporting different packet size: packet fragmentation.

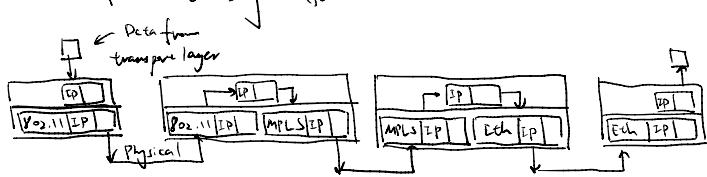
5.5.2 How Networks Differ

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.

5.5.3 Connecting Heterogeneous Networks



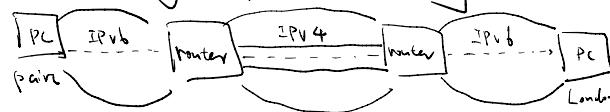
→ packet crossing different networks.



Network and link layer protocol processing.

- Because different networks may, in general, have different forms of addressing the packet carries a network layer address that can identify any host across the three networks;
- Once the packet has traveled along the virtual circuit, it will reach the Ethernet network;
- The packet may be too large to be carried, since 802.11 can work with larger frames than Ethernet;
- To handle this problem, the packet is divided into fragments, and each fragment is sent separately. When the fragments reach the destination, they are reassembled.

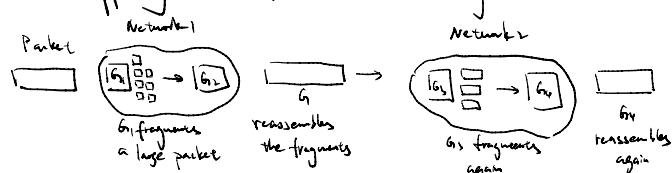
5.5.4 Connecting Endpoints Across Heterogeneous Networks



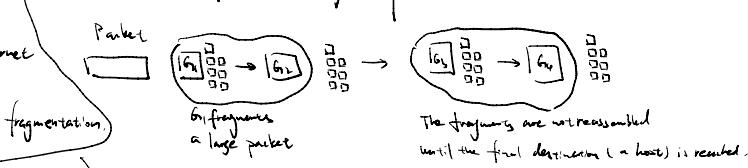
→ IPv6 packet → IPv4 [IPv6 packet] → IPv6 packet

- A common special case is where the source and destination hosts are on the same type of network, but there is a different network in between.
- The solution to this problem is a technique called tunneling.
- To send an IP packet to a host in the London office, a host in the Paris constructs the packet containing an IPv6 address in London and sends it to the multi-protocol router that connects the Paris IPv6 network to the IPv4 Internet.
- When this is between routers gets the IPv6 packet, the router puts the IPv6 packet inside a IPv4 packet. When this wrapped packet arrives the London router removes the original IPv6 packet and sends it onward to the destination host.

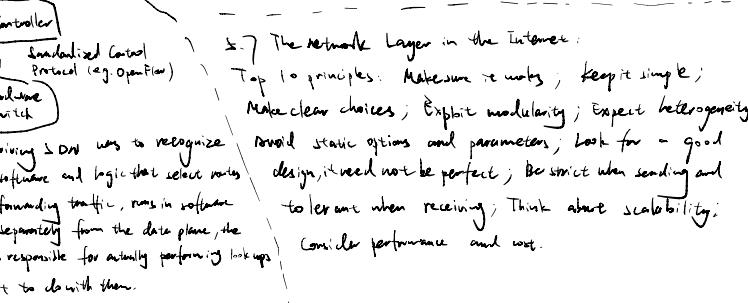
5.5.5 Supporting Different Packet Size: Packet Fragmentation



(a) Transparent fragmentation: is straightforward but has some problems



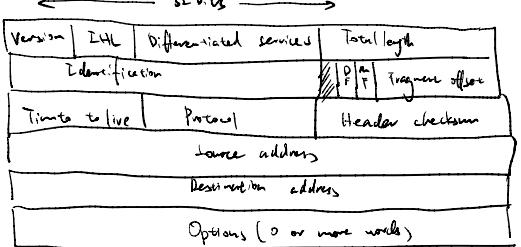
(b) Non-transparent fragmentation: The advantage of this is it requires routers to do less work.



5.5.6 The Network Layer in the Internet

- Top 10 principles: Make sure it works; keep it simple; Make clear choices; Exploit modularity; Expect heterogeneity; Avoid static options and parameters; Look for a good design, need not be perfect; Be strict when sending and tolerant when receiving; Think about scalability; Consider performance and cost.

I.7.1 The IP Version 4 Protocol



I.8 Policy at the Network Layer

- Peering disputes
 - A breakdown in negotiations over paying for transit.
- Traffic prioritization
 - Generally agreed upon bright-line rules
 - No blocking; - No throttling; - No paid prioritization;
 - Disclosure of any prioritization practices.

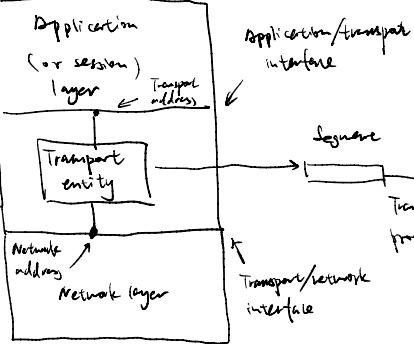
b.1 The Transport Service

- Services provided to the upper layers
- Transport service primitives
- Berkeley sockets
- Example of socket programming: An internet file server.

b.1.1 Services provided to the Upper layers

- The goal of the transport layer is to provide efficient, reliable, and cost-effective data transmission service to its users, normally processes in the application layer.
- The transport layer makes use of the services provided by the network layer. The software and/or hardware within the transport layer that does work is called the transport entity.
- The transport entity can be located in the operating system kernel, in a library package bound into network applications, in a separate user process, or even on the network interface card.

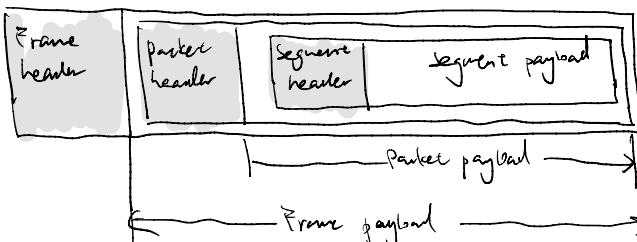
Host 1



The network, transport, and application layers

b.1.2 Transport Service Primitives

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

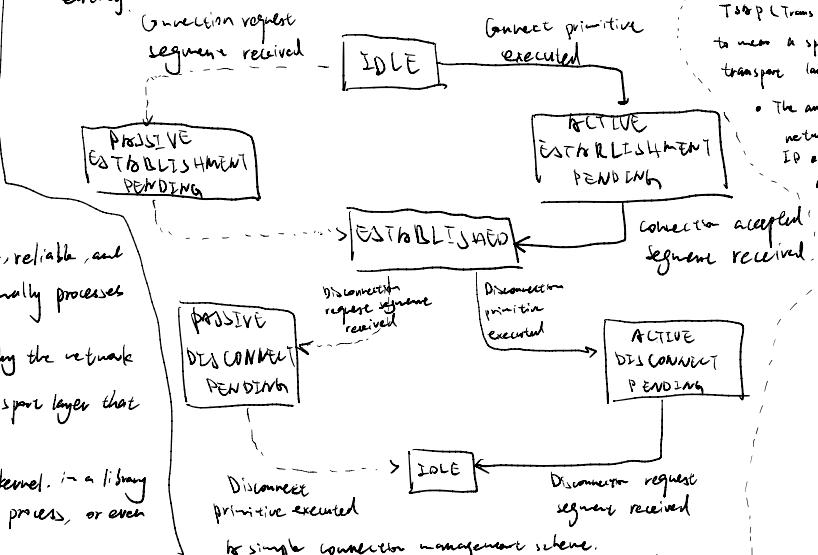


Nesting of segments, packets and frames.

- We will use the term segment for messages sent from transport entity to transport entity.

- Segments (exchanged by the transport layer) are contained in packets (which are exchanged by the network layer). In turn, these packets are contained in frames (exchanged by the data link layer).

- When a frame arrives, the data link layer processes the frame header and, if the destination address matches for local delivery, passes the contents of the frame payload field up to the network entity. The network entity similarly processes the packet header and then passes the contents of the packet payload up to the transport entity.



To simplify connection management scheme.
The solid lines show the client's state sequence. The dashed lines show the server's state sequence.

b.1.3 Berkeley sockets

Socket primitives were first released as part of the Berkeley UNIX 4.2BSD software distribution in 1983. They quickly became popular. The primitives are now widely used for Internet programming on many operating systems.

The socket primitives for TCP

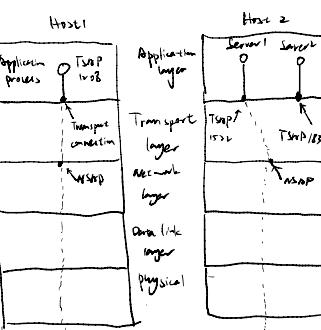
Primitive	Meaning
SOCKET	Create a new communication endpoint
BIND	Associate a local address with a socket
LISTEN	Announce willingness to accept connections; give queue size
ACCEPT	Passively establish an incoming connection
CONNECT	Actively attempt to establish a connecting
SEND	Send some data over the connection
RECEIVE	Receive some data from the connection
CLOSE	Release the connection

b.1.4 An Example of Socket Programming: An Internet File Server

b.2 Elements of Transport Protocols

- Addressing
- Error control and flow control
- Multiplexing
- Crash recovery.

b.2.1 Addressing

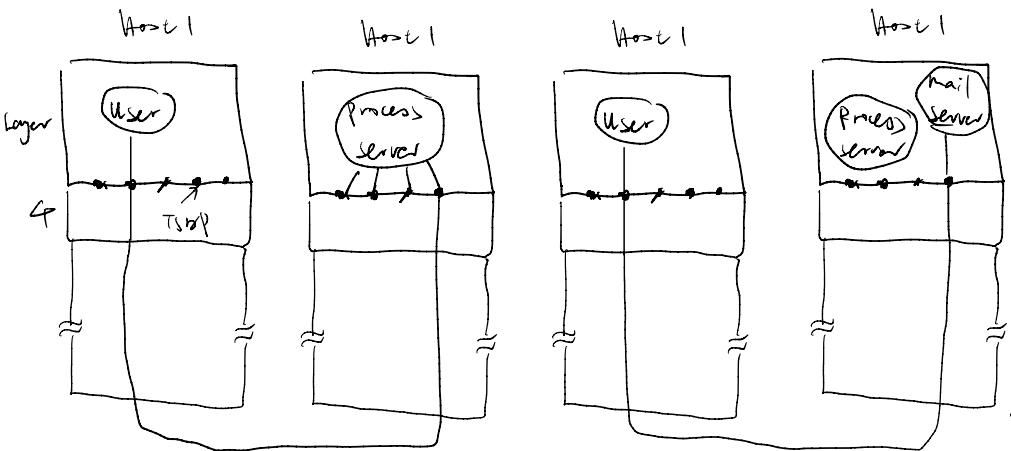


TCP/IPv4, TCP/IPv6, and transport connection.

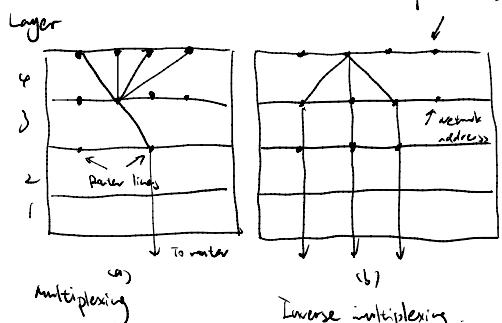
- The transport addresses to which processes can listen for connection requests are called ports.

- We will use the generic term TSP (Transport Service Access Point) to mean a specific endpoint in the transport layer.

- The analogous endpoints in the network layer are called ports. IP addresses are example of NISPs.



b.2.5 Multiplexing



(a) If only one network address is available on a host, all transport connections on that machine have to use it. When a segment comes in, some way is needed to tell which process to give it to. This situation is called multiplexing.

(b) A host has multiple network ports that it can use. If a user needs more bandwidth or more reliability than one of the network ports can provide, a way out is to have a connection that distributes the traffic among multiple network paths on a round-robin basis. This is called inverse multiplexing.

b.2.6 Crash Recovery

Strategy used by receiving host

Strategy used by sending host	First tick, then write			First write, then tick		
	OK	DUP	OK	OK	DUP	DUP
Always retransmit	OK	LOST	OK	OK	OK	OK
Never retransmit	LOST	OK	LOST	OK	OK	OK
Retransmit in SO	OK	DUP	LOST	OK	OK	OK
Retransmit in SJ	LOST	OK	OK	OK	OK	DUP

OK = protocol finishes correctly
DUP = protocol generates a duplicate message
LOST = protocol loses a message.

Each client can be one of two types: one segment outstanding SJ, or no segments outstanding SO.

b.3 Congestion Control

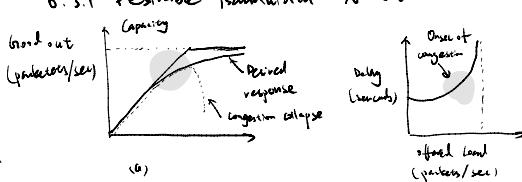
If the transport entities on many machines send too many packets into the network too quickly, the network will become congested, with performance degraded as packets are delayed and lost.

Controlling congestion to avoid this problem is the combined responsibility of the network and transport layers.

Congestion occurs at routers, so it is detected at the network layer. However, congestion is ultimately caused by traffic sent into the network by the transport layer.

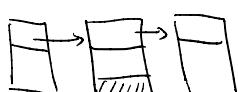
The only effective way to control congestion is for the transport protocols to send packets into the network more slowly.

b.3.1 Desirable Bandwidth Allocation

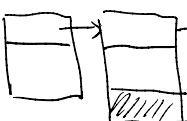


b.2.4 Error Control and Flow Control

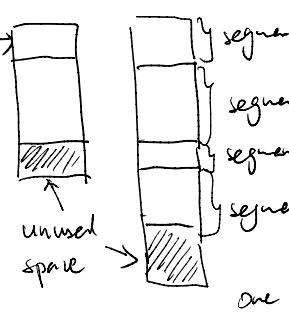
- Error Control is ensuring that the data is delivered with the desired level of reliability, usually that all of the data is delivered without any errors.
- Flow Control is keeping a fast transmitter from overrunning a slow receiver.
- A frame carries an error-detecting code (e.g. a CRC or checksum) that is used to check if the information was correctly received.
- A frame carries a sequence number to identify itself and is retransmitted by the sender until it receives an acknowledgement of successful receipt from the receiver. This is called ARQ (Automatic Repeat reQuest).
- There is a maximum number of frames that the sender will allow to be outstanding at any time, pausing if the receiver is not acknowledging frames quickly enough. If this maximum is one packet the protocol is called stop-and-wait.
- The sliding window protocol combines these features and is also used to support bidirectional data transfer.



(a)
Chained fixed-size buffers

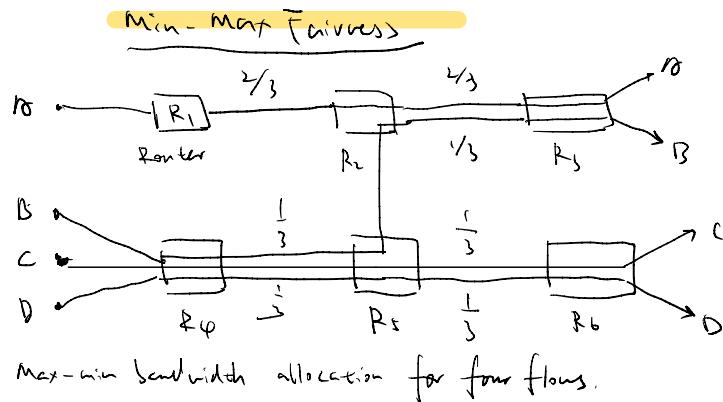


(b)
Chained variable-sized buffers



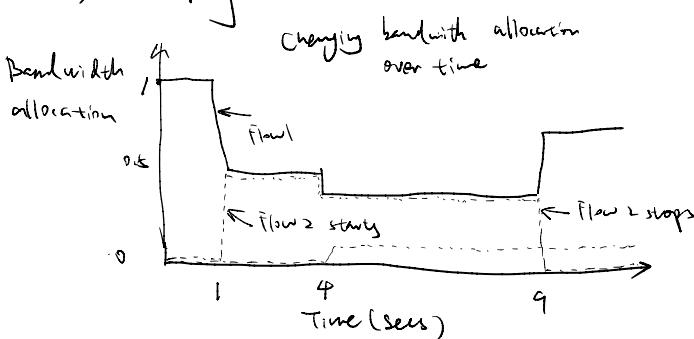
(c)
One large circular buffer per connection

6.3.1.2 Desirable Bandwidth Allocation:



- The form of fairness that is often desired for network usage is max-min fairness. An allocation is max-min fair if the bandwidth given to one flow cannot be increased without decreasing the bandwidth given to another flow with an allocation that is no larger.
 - Each of the links between routers has the same capacity, taken to be 1 unit, though in the general case the links will have different capacities.
 - If more of the bandwidth on the link between R2 and R3 is given to flow B, there will be less for flow A. This is reasonable as flow A already has more bandwidth. However, the capacity of flow C or D (or both) must be decreased to give more bandwidth to B, and these flows will have less bandwidth than B.
- Thus, the allocation is max-min fair.

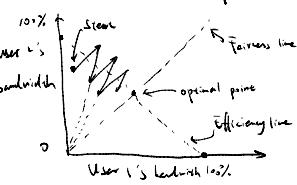
6.3.1.3 Convergence



Because of the variation in demand, the ideal operating point for the network varies over time. A good congestion control algorithm should rapidly converge to the ideal operating point, and it should track this point as it changes over time. Initially, flow 1 has all of the bandwidth. One second later, flow 2 starts. It needs bandwidth. At 4 seconds, a third flow joins. However, this flow uses only 50% of the bandwidth, which is less than its fair share (which is a third). Flow 1 can quickly adjust, dividing the available bandwidth to each flow 66.6% of the bandwidth. At 9 seconds, the second flow leaves, and the third flow remains unchanged. The first flow quickly captures 83.3% of the bandwidth.

6.3.4 Regulating the sending rate (AIMD)

Additive increase multiplicative decrease (AIMD) control law.



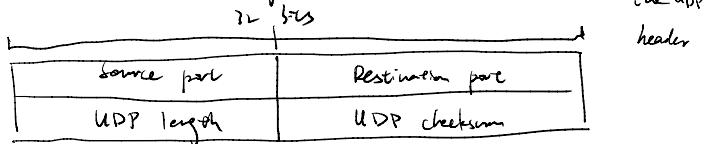
Legend: \nearrow = additive increase ($\approx 45^\circ$)

\searrow = multiplicative decrease
(use points to origin)

The users additively increase their bandwidth allocations and then multiplicatively decrease them when congestion is sensed. This behavior is the AIMD control law. It can be seen that the path traced by this behavior does converge to the optimal point that is both fair and efficient.

6.4.1 Introduction to UDP

- The connectionless protocol is UDP (User Datagram Protocol)
- It does almost nothing beyond sending packets between applications, letting applications build their own protocols on top as needed.
- UDP provides a way for applications to send encapsulated IP datagrams without having to establish a connection.



- The source port is primarily needed when a reply must be sent back to the source.
- By copying the source port field from the incoming segment into the Destination port field of the outgoing segment, the process sending the reply can specify which process on the sending machine is to get it.
- The UDP length field includes the 8-byte header and the data.
- An optional checksum is also provided for extra reliability.