

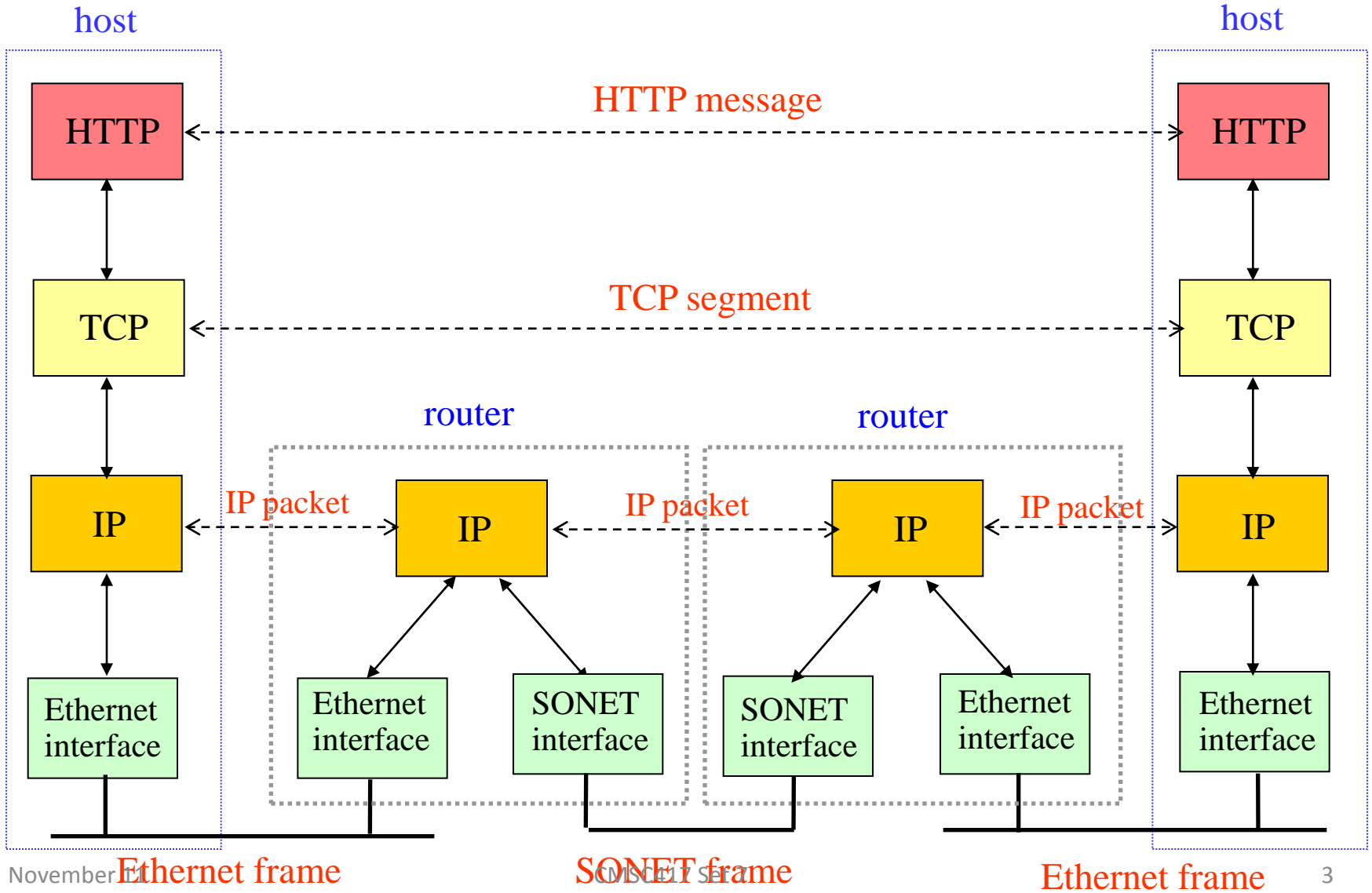
# CMSC 417

## Computer Networks Prof. Ashok K Agrawala

© 2011 Ashok Agrawala  
Set 7

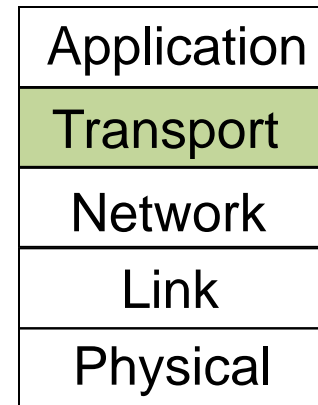
# The Transport Layer

# Message, Segment, Packet, and Frame



# The Transport Layer

Responsible for delivering data across networks with the desired reliability or quality



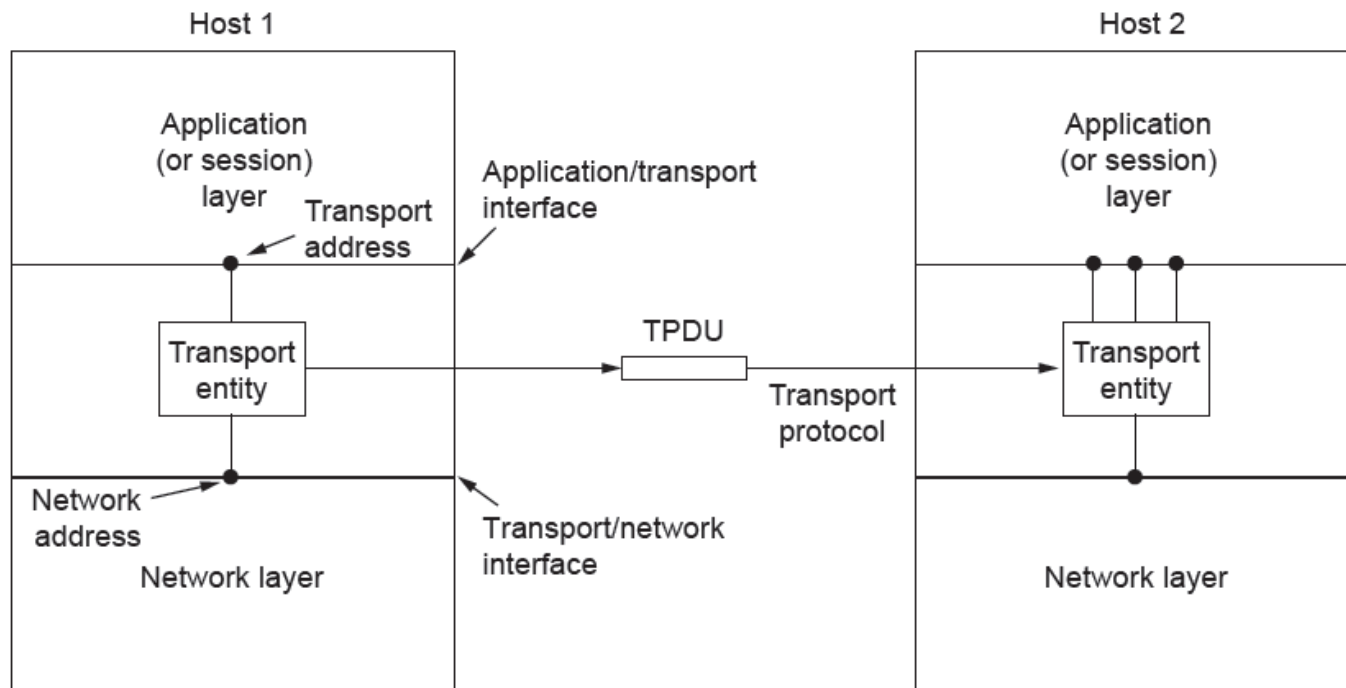
# The Transport Service

- Services Provided to the Upper Layers
- Transport Service Primitives
- Berkeley Sockets
- An Example of Socket Programming:
  - An Internet File Server

# Services Provided to the Upper Layers (1)

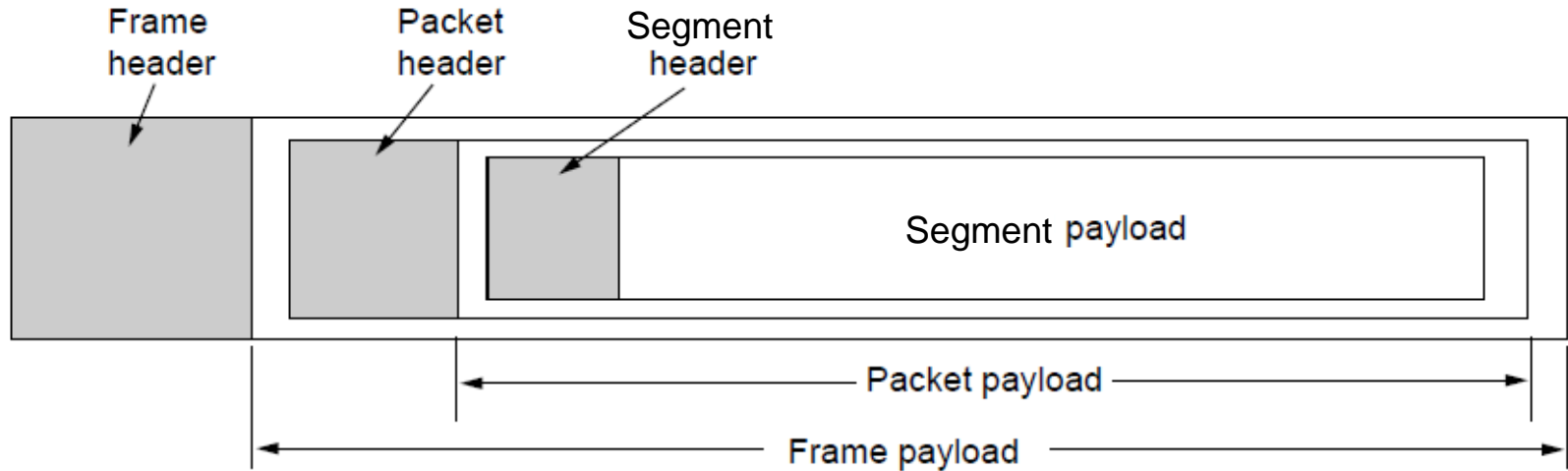
Transport layer adds reliability to the network layer

- Offers connectionless (e.g., UDP) and connection-oriented (e.g, TCP) service to applications



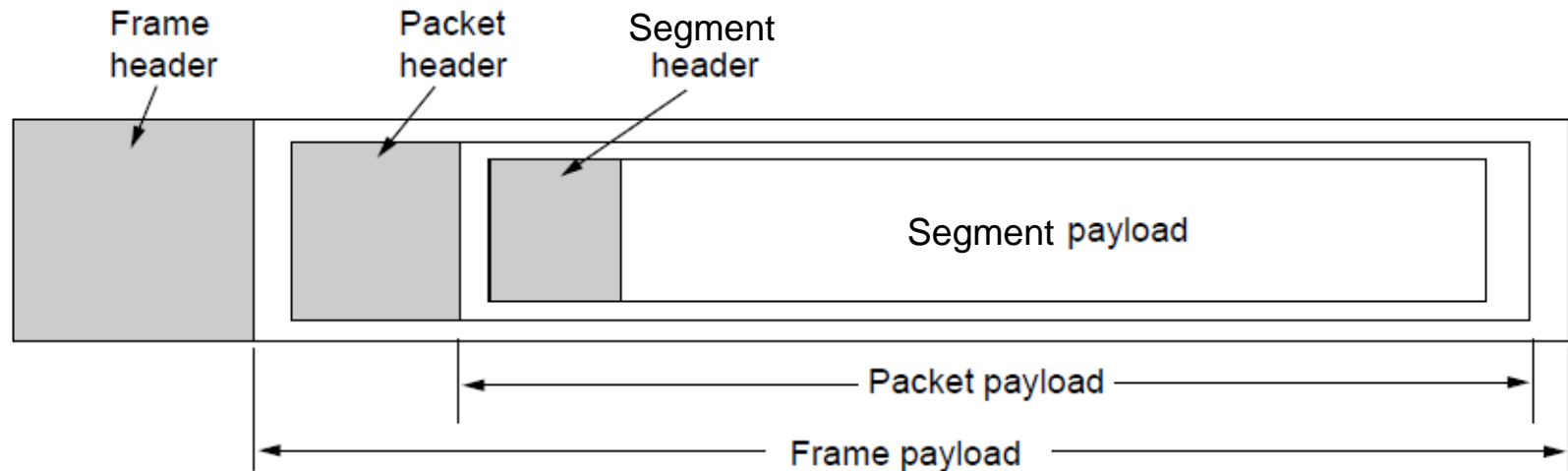
# Services Provided to the Upper Layers (2)

Transport layer sends segments in packets (in frames)



# Services Provided to the Upper Layers (2)

Transport layer sends segments in packets (in frames)



# Berkeley Sockets

Very widely used primitives started with TCP on UNIX

- Notion of “sockets” as transport endpoints
- Like simple set plus SOCKET, BIND, and ACCEPT

Primitive	Meaning
SOCKET	Create a new communication end point
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 connection
SEND	Send some data over the connection
RECEIVE	Receive some data from the connection
CLOSE	Release the connection

# Example of Socket Programming: An Internet File Server (1)

```
/* This page contains a client program that can request a file from the server program
 * on the next page. The server responds by sending the whole file.
 */

#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>

#define SERVER_PORT 12345          /* arbitrary, but client & server must agree */
#define BUF_SIZE 4096            /* block transfer size */

int main(int argc, char **argv)
{
    int c, s, bytes;
    char buf[BUF_SIZE];           /* buffer for incoming file */
    struct hostent *h;           /* info about server */
    struct sockaddr_in channel;  /* holds IP address */

    . . .
```

Client code using sockets

# Example of Socket Programming: An Internet File Server (2)

. . .

```
if (argc != 3) fatal("Usage: client server-name file-name");
h = gethostbyname(argv[1]);           /* look up host's IP address */
if (!h) fatal("gethostbyname failed");

s = socket(PF_INET, SOCK_STREAM, IPPROTO_TCP);
if (s < 0) fatal("socket");
memset(&channel, 0, sizeof(channel));
channel.sin_family= AF_INET;
memcpy(&channel.sin_addr.s_addr, h->h_addr, h->h_length);
channel.sin_port= htons(SERVER_PORT);

c = connect(s, (struct sockaddr *) &channel, sizeof(channel));
if (c < 0) fatal("connect failed");
```

. . .

Client code using sockets

# Example of Socket Programming: An Internet File Server (3)

...

```
c = connect(s, (struct sockaddr *) &channel, sizeof(channel));
if (c < 0) fatal("connect failed");

/* Connection is now established. Send file name including 0 byte at end. */
write(s, argv[2], strlen(argv[2])+1);

/* Go get the file and write it to standard output. */
while (1) {
    bytes = read(s, buf, BUF_SIZE);          /* read from socket */
    if (bytes <= 0) exit(0);                 /* check for end of file */
    write(1, buf, bytes);                    /* write to standard output */
}
}

fatal(char *string)
{
    printf("%s\n", string);
    exit(1);
}
```

Client code using sockets

# Example of Socket Programming: An Internet File Server (4)

```
#include <sys/types.h>                /* This is the server code */
#include <sys/fcntl.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>

#define SERVER_PORT 12345              /* arbitrary, but client & server must agree */
#define BUF_SIZE 4096                 /* block transfer size */
#define QUEUE_SIZE 10

int main(int argc, char *argv[])
{
    int s, b, l, fd, sa, bytes, on = 1;
    char buf[BUF_SIZE];                /* buffer for outgoing file */
    struct sockaddr_in channel;        /* holds IP address */
```

. . .

Server code

# Example of Socket Programming: An Internet File Server (5)

...

```
/* Build address structure to bind to socket. */
memset(&channel, 0, sizeof(channel));    /* zero channel */
channel.sin_family = AF_INET;
channel.sin_addr.s_addr = htonl(INADDR_ANY);
channel.sin_port = htons(SERVER_PORT);

/* Passive open. Wait for connection. */
s = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP); /* create socket */
if (s < 0) fatal("socket failed");
setsockopt(s, SOL_SOCKET, SO_REUSEADDR, (char *) &on, sizeof(on));

b = bind(s, (struct sockaddr *) &channel, sizeof(channel));
if (b < 0) fatal("bind failed");

l = listen(s, QUEUE_SIZE);                /* specify queue size */
if (l < 0) fatal("listen failed");
```

...

Server code

# Example of Socket Programming: An Internet File Server (6)

...

```
/* Socket is now set up and bound. Wait for connection and process it. */
while (1) {
    sa = accept(s, 0, 0);                /* block for connection request */
    if (sa < 0) fatal("accept failed");

    read(sa, buf, BUF_SIZE);           /* read file name from socket */

    /* Get and return the file. */
    fd = open(buf, O_RDONLY);           /* open the file to be sent back */
    if (fd < 0) fatal("open failed");

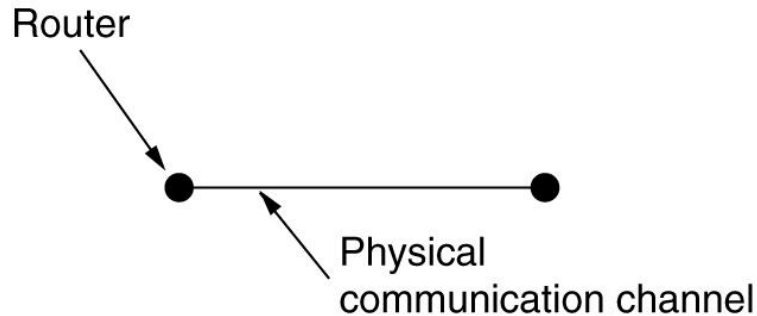
    while (1) {
        bytes = read(fd, buf, BUF_SIZE); /* read from file */
        if (bytes <= 0) break;           /* check for end of file */
        write(sa, buf, bytes);           /* write bytes to socket */
    }
    close(fd);                           /* close file */
    close(sa);                            /* close connection */
}
}
```

Server code

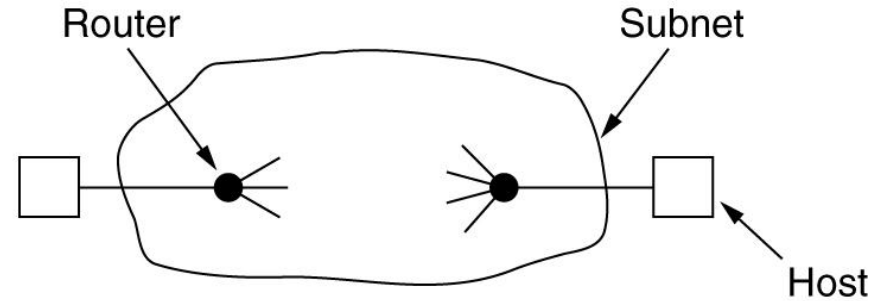
# Elements of Transport Protocols

- Addressing »
- Connection establishment »
- Connection release »
- Error control and flow control »
- Multiplexing »
- Crash recovery »

# Transport Protocol



(a)

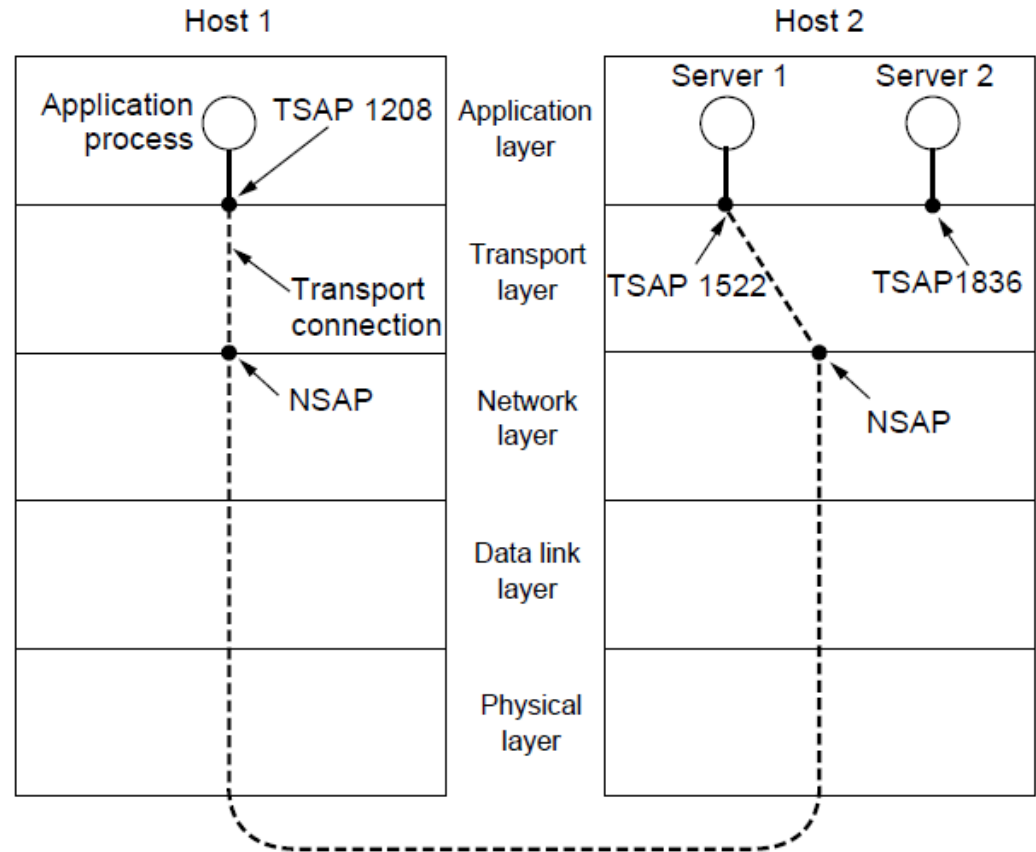


(b)

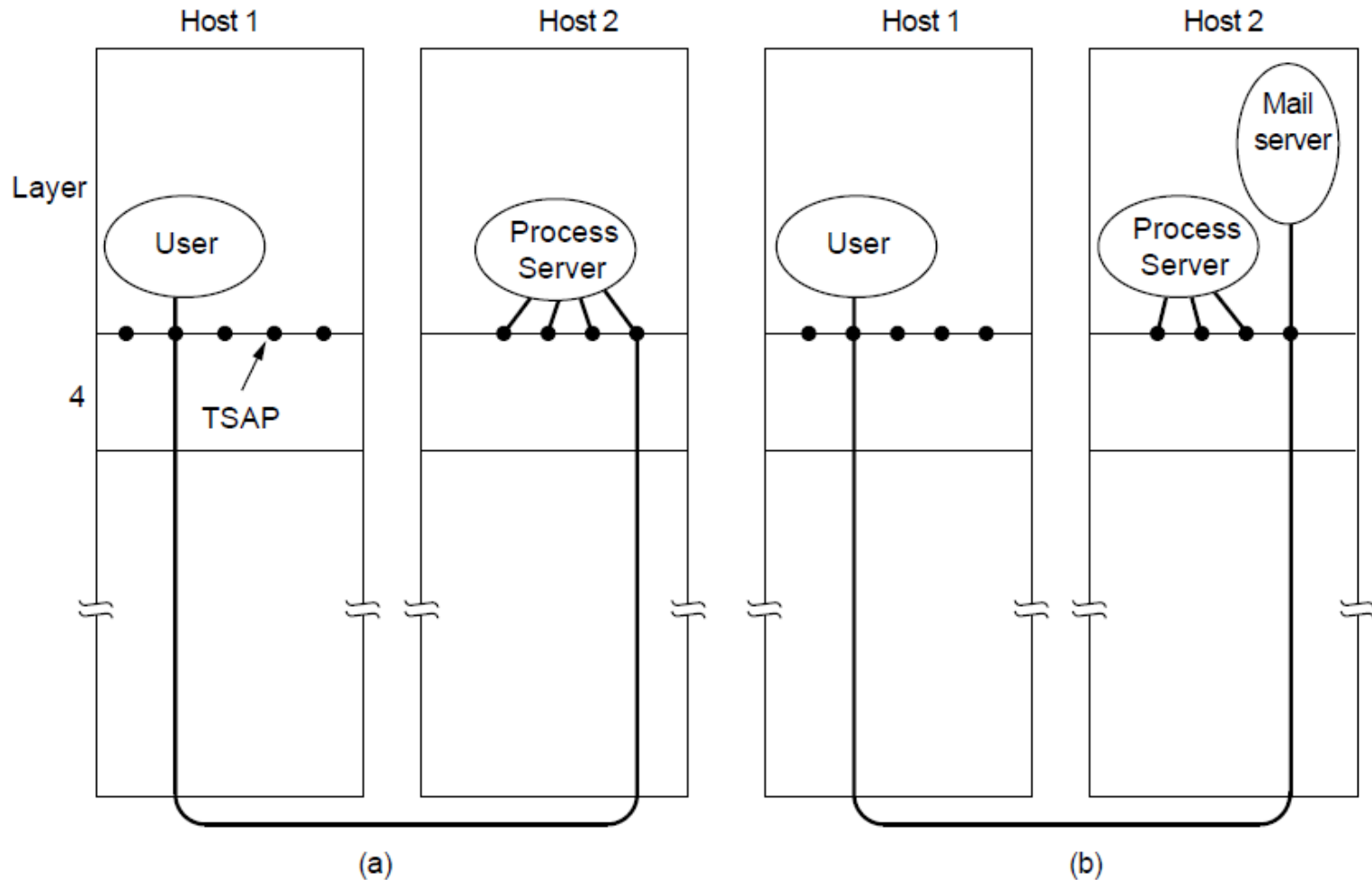
- (a) Environment of the data link layer.
- (b) Environment of the transport layer.

# Addressing

- Transport layer adds TSAPs
- Multiple clients and servers can run on a host with a single network (IP) address
- TSAPs are ports for TCP/UDP



# Addressing (2)



How a user process in host 1 establishes a connection with a mail server in host 2 via a process server.

# Connection Establishment (1)

Key problem is to ensure reliability even though packets may be lost, corrupted, delayed, and duplicated

- Don't treat an old or duplicate packet as new
- (Use ARQ and checksums for loss/corruption)

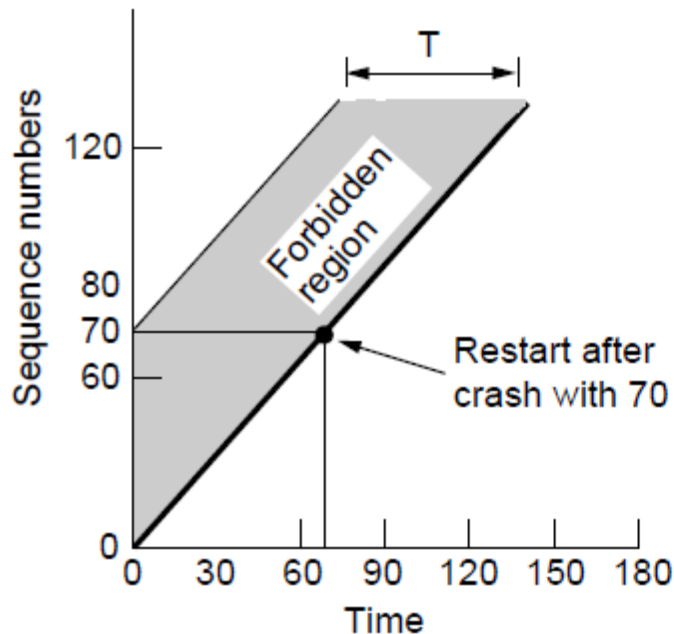
Approach:

- Don't reuse sequence numbers within twice the MSL (Maximum Segment Lifetime) of  $2T=240$  secs
- Three-way handshake for establishing connection

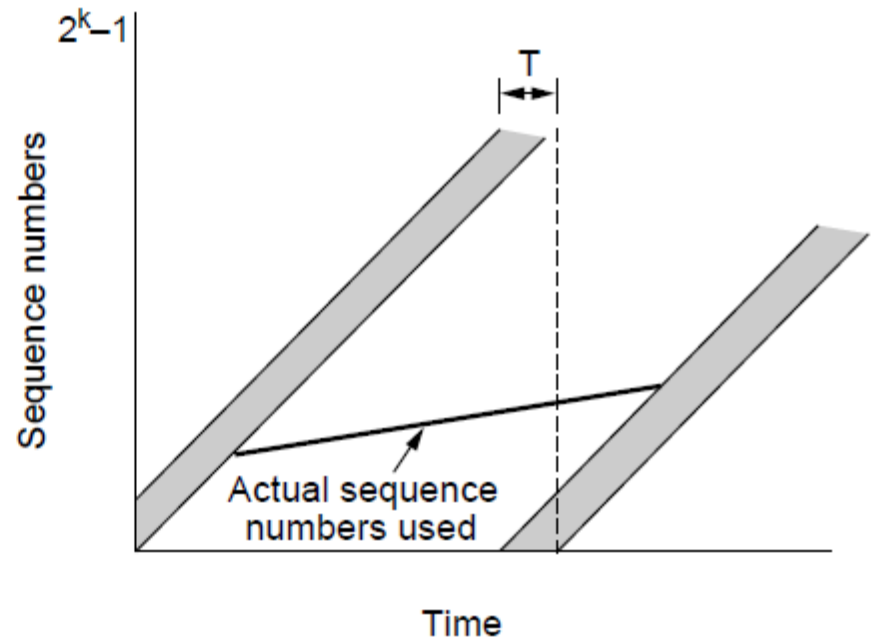
# Connection Establishment (2)

Use a sequence number space large enough that it will not wrap, even when sending at full rate

- Clock (high bits) advances & keeps state over crash



Need seq. number not to wrap within T seconds

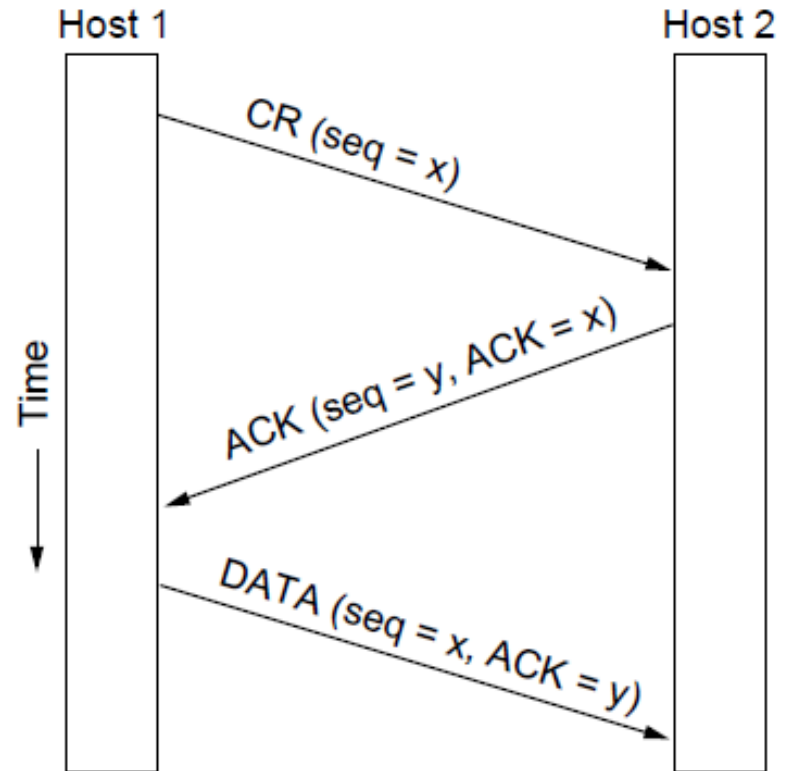


Need seq. number not to climb too slowly for too long

# Connection Establishment (3)

Three-way handshake used for initial packet

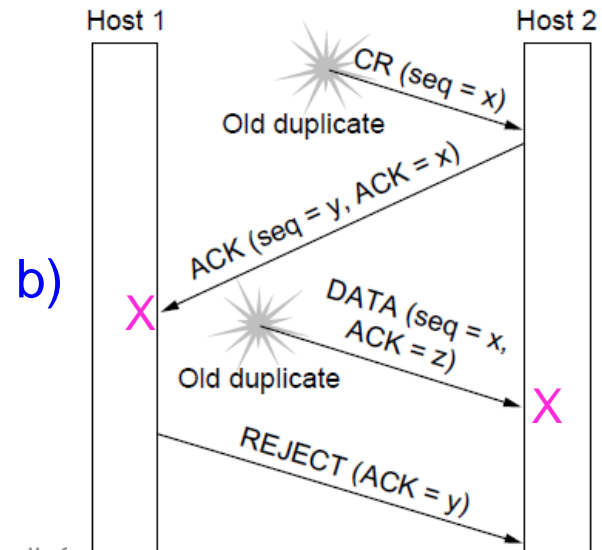
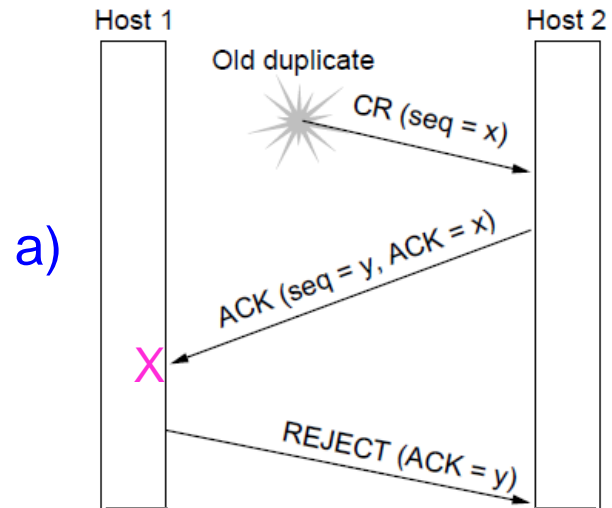
- Since no state from previous connection
- Both hosts contribute fresh seq. numbers
- CR = Connect Request



# Connection Establishment (4)

Three-way handshake  
protects against odd  
cases:

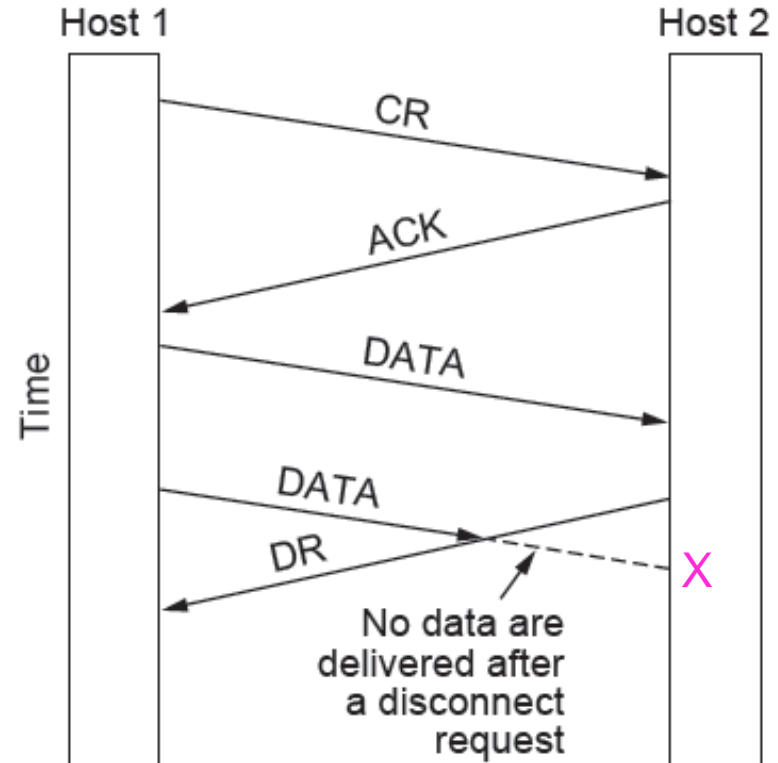
- a) Duplicate CR. Spurious ACK does not connect
- b) Duplicate CR and DATA. Same plus DATA will be rejected (wrong ACK).



# Connection Release (1)

Key problem is to ensure reliability while releasing connection

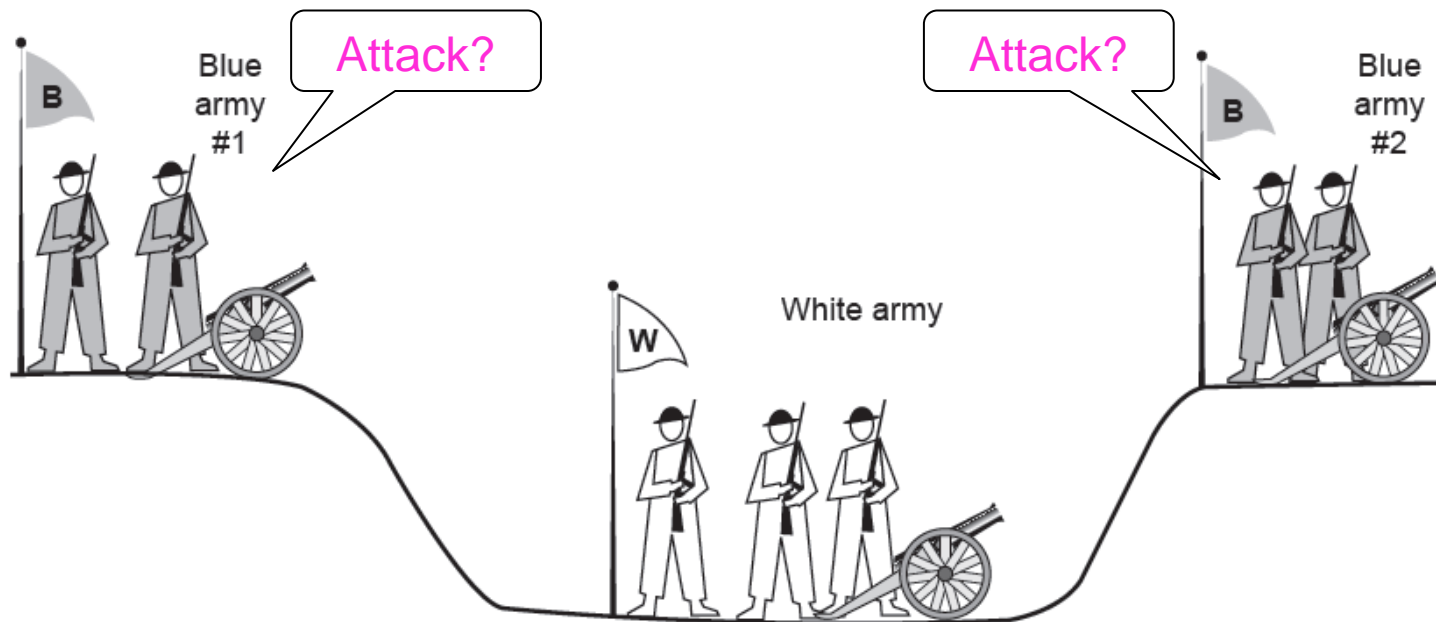
Asymmetric release (when one side breaks connection) is abrupt and may lose data



# Connection Release (2)

Symmetric release (both sides agree to release)  
can't be handled solely by the transport layer

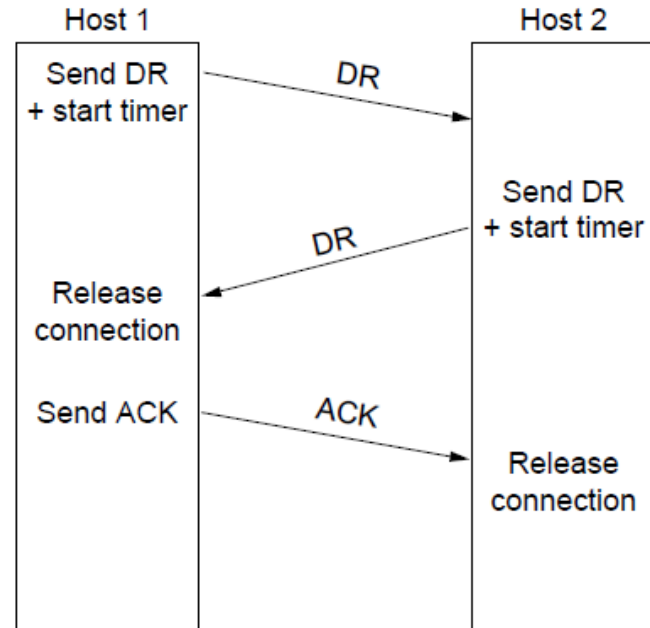
- Two-army problem shows pitfall of agreement



# Connection Release (3)

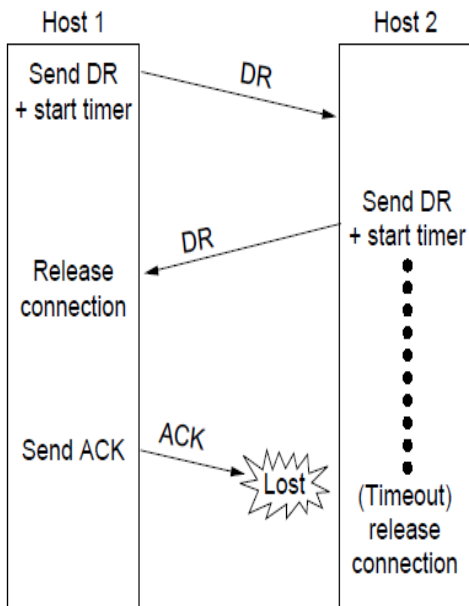
Normal release sequence,  
initiated by transport  
user on Host 1

- DR=Disconnect Request
- Both DRs are ACKed by the other side

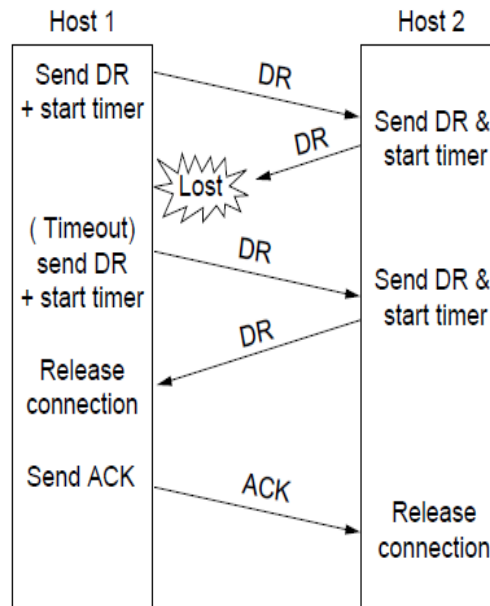


# Connection Release (4)

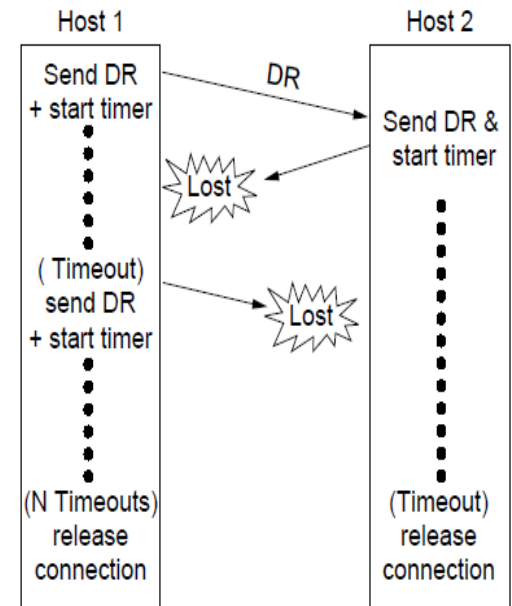
Error cases are handled with timer and retransmission



Final ACK lost,  
Host 2 times out



Lost DR causes  
retransmissions



Extreme: Many lost  
DRs cause both  
hosts to timeout

# Flow Control

- Use Sliding Window
- Buffering
  - Sender buffers all TPDU's until acknowledged
  - TPDU lost by the network
    - Unreliable service
    - Receiver not having buffer
- How should buffers be managed
  - Dedicate
  - Acquire when needed
- Traffic
  - Low bandwidth, Bursty – buffer at sender – acquire at receiver
  - High bandwidth, smooth – Buffer at both ends
- Exchange Buffer information

# Error Control and Flow Control (1)

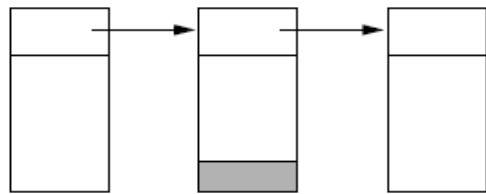
Foundation for error control is a sliding window (from Link layer) with checksums and retransmissions

Flow control manages buffering at sender/receiver

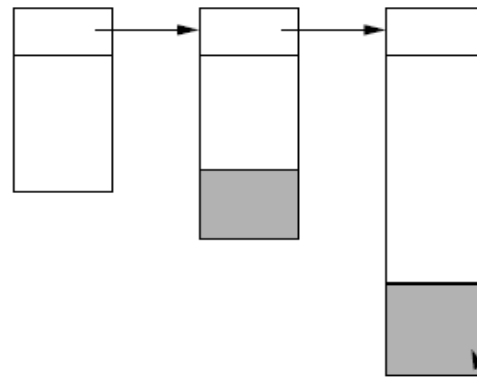
- Issue is that data goes to/from the network and applications at different times
- Window tells sender available buffering at receiver
- Makes a variable-size sliding window

# Error Control and Flow Control (2)

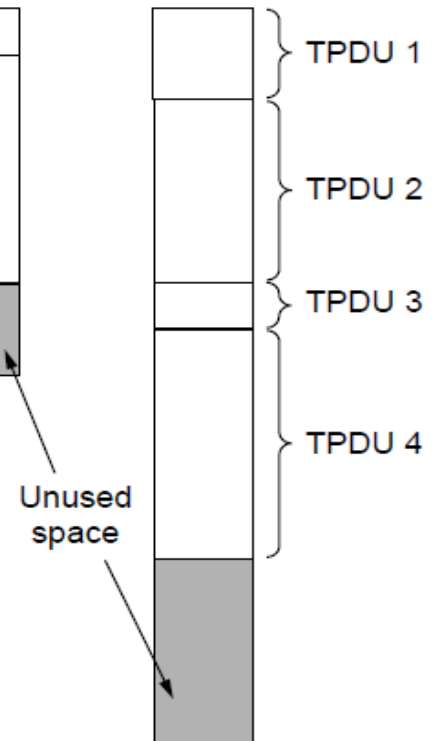
Different buffer strategies trade efficiency / complexity



a) Chained fixed-size buffers



b) Chained variable-size buffers



c) One large circular buffer

# Error Control and Flow Control (3)

Flow control example: A's data is limited by B's

A	Message	B	B's Buffer	Comments
1	→ < request 8 buffers >	→		A wants 8 buffers
2	← <ack = 15, buf = 4 >	←	0 1 2 3	B grants messages 0-3 only
3	→ <seq = 0, data = m0 >	→	0 1 2 3	A has 3 buffers left now
4	→ <seq = 1, data = m1 >	→	0 1 2 3	A has 2 buffers left now
5	→ <seq = 2, data = m2 >	...	0 1 2 3	Message lost but A thinks it has 1 left
6	← <ack = 1, buf = 3 >	←	1 2 3 4	B acknowledges 0 and 1, permits 2-4
7	→ <seq = 3, data = m3 >	→	1 2 3 4	A has 1 buffer left
8	→ <seq = 4, data = m4 >	→	1 2 3 4	A has 0 buffers left, and must stop
9	→ <seq = 2, data = m2 >	→	1 2 3 4	A times out and retransmits
10	← <ack = 4, buf = 0 >	←	1 2 3 4	Everything acknowledged, but A still blocked
11	← <ack = 4, buf = 1 >	←	2 3 4 5	A may now send 5
12	← <ack = 4, buf = 2 >	←	3 4 5 6	B found a new buffer somewhere
13	→ <seq = 5, data = m5 >	→	3 4 5 6	A has 1 buffer left
14	→ <seq = 6, data = m6 >	→	3 4 5 6	A is now blocked again
15	← <ack = 6, buf = 0 >	←	3 4 5 6	A is still blocked
16	... <ack = 6, buf = 4 >	←	7 8 9 10	Potential deadlock

# Sliding Window Protocol

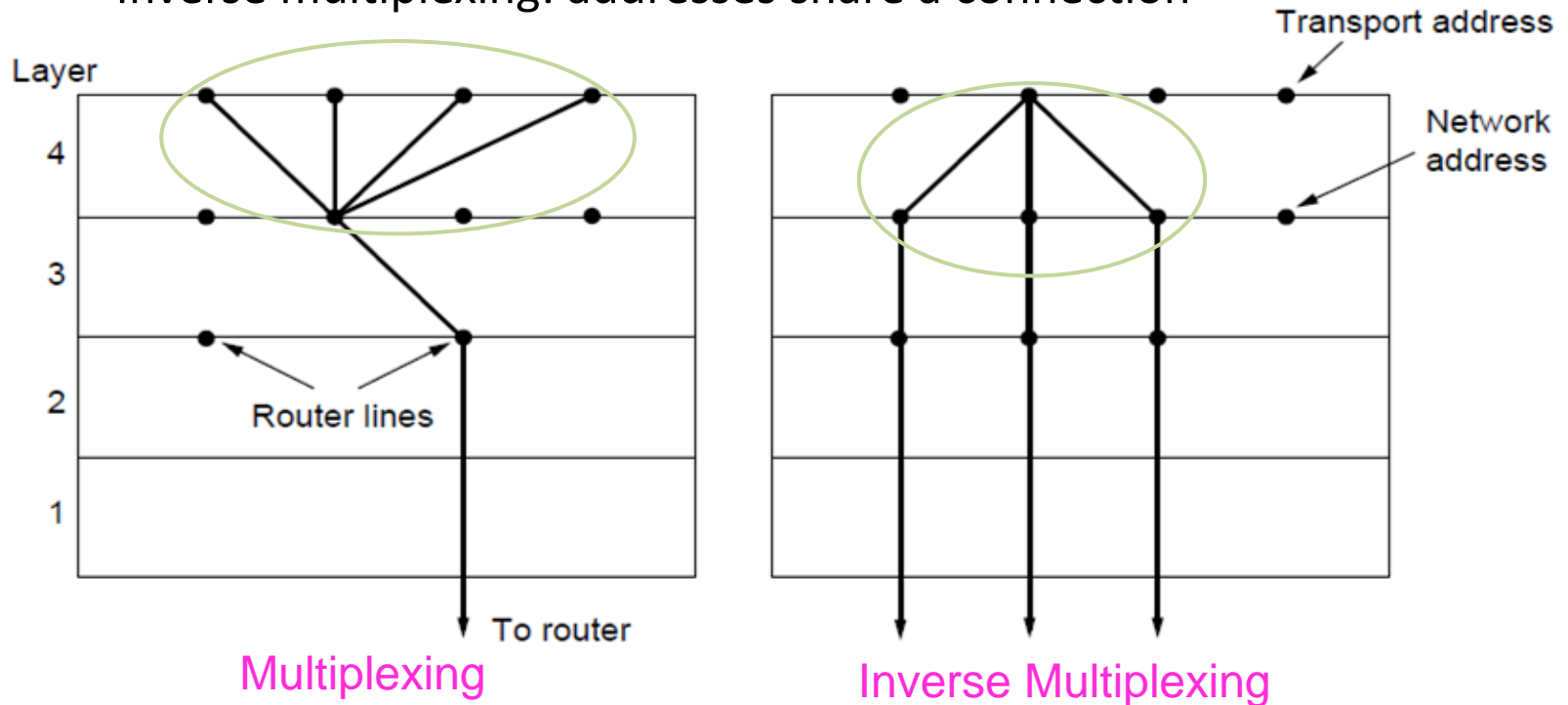
- Animations
  - [http://www.site.uottawa.ca/~elsaddik/abedweb/applets/Applets/Sliding\\_Window/sliding\\_window.html](http://www.site.uottawa.ca/~elsaddik/abedweb/applets/Applets/Sliding_Window/sliding_window.html)
  - <http://www.osischool.com/protocol/Tcp/slidingWindow/index.php>

# Throughput limits

- Buffers
- Bandwidth – subnet's carrying capacity
  - $K$  TPDU's per second
  - $X$  paths then total of  $XK$
- Flow control to manage
  - Manage window size
    - If network can handle  $c$  TPDU's/sec and Cycle time is  $r$  then the window size should be  $cr$

# Multiplexing

- Kinds of transport / network sharing that can occur:
  - Multiplexing: connections share a network address
  - Inverse multiplexing: addresses share a connection



# Crash Recovery

- Network Failures
  - Transport layer handles
    - Connectionless
    - Connection oriented
- Host Crashes
  - Server crash and may reboot
    - Send broadcast asking clients to inform of prior connections ( stop and wait protocol)
      - Client – one TPDU outstanding or none outstanding

# Crash Recovery

Application needs to help recovering from a crash

- Transport can fail since A(ck) / W(rite) not atomic

Strategy used by receiving host

← First ACK, then write →

← First write, then ACK →

Strategy used by sending host	First ACK, then write			First write, then ACK		
	AC(W)	AWC	C(AW)	C(WA)	W AC	WC(A)
Always retransmit	OK	DUP	OK	OK	DUP	DUP
Never retransmit	LOST	OK	LOST	LOST	OK	OK
Retransmit in S0	OK	DUP	LOST	LOST	DUP	OK
Retransmit in S1	LOST	OK	OK	OK	OK	DUP

OK = Protocol functions correctly  
 DUP = Protocol generates a duplicate message  
 LOST = Protocol loses a message

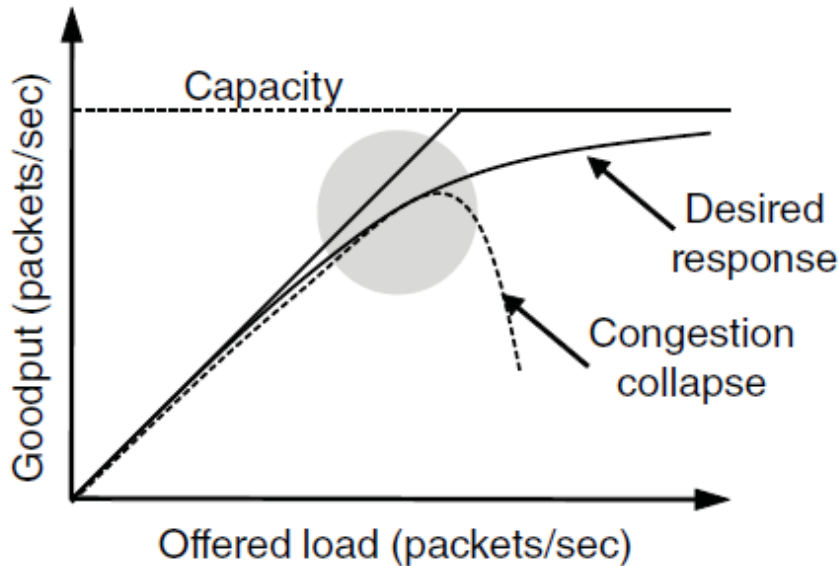
# Congestion Control

Two layers are responsible for congestion control:

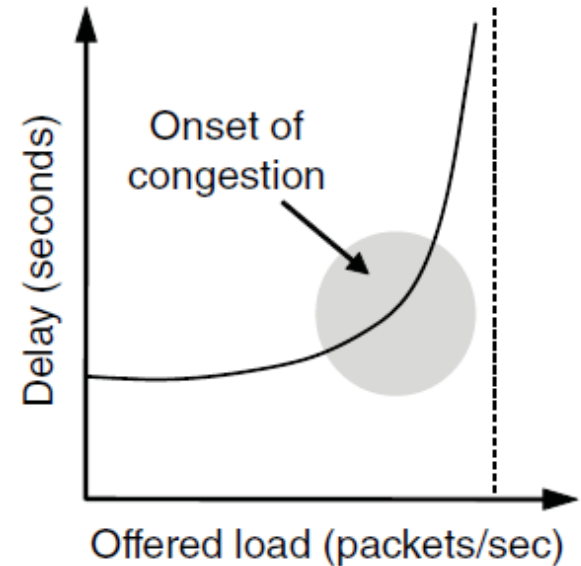
- Transport layer, controls the offered load [here]
  - Network layer, experiences congestion [previous]
- 
- Desirable bandwidth allocation »
  - Regulating the sending rate »
  - Wireless issues »

# Desirable Bandwidth Allocation (1)

Efficient use of bandwidth gives high goodput, low delay



Goodput rises more slowly than load when congestion sets in

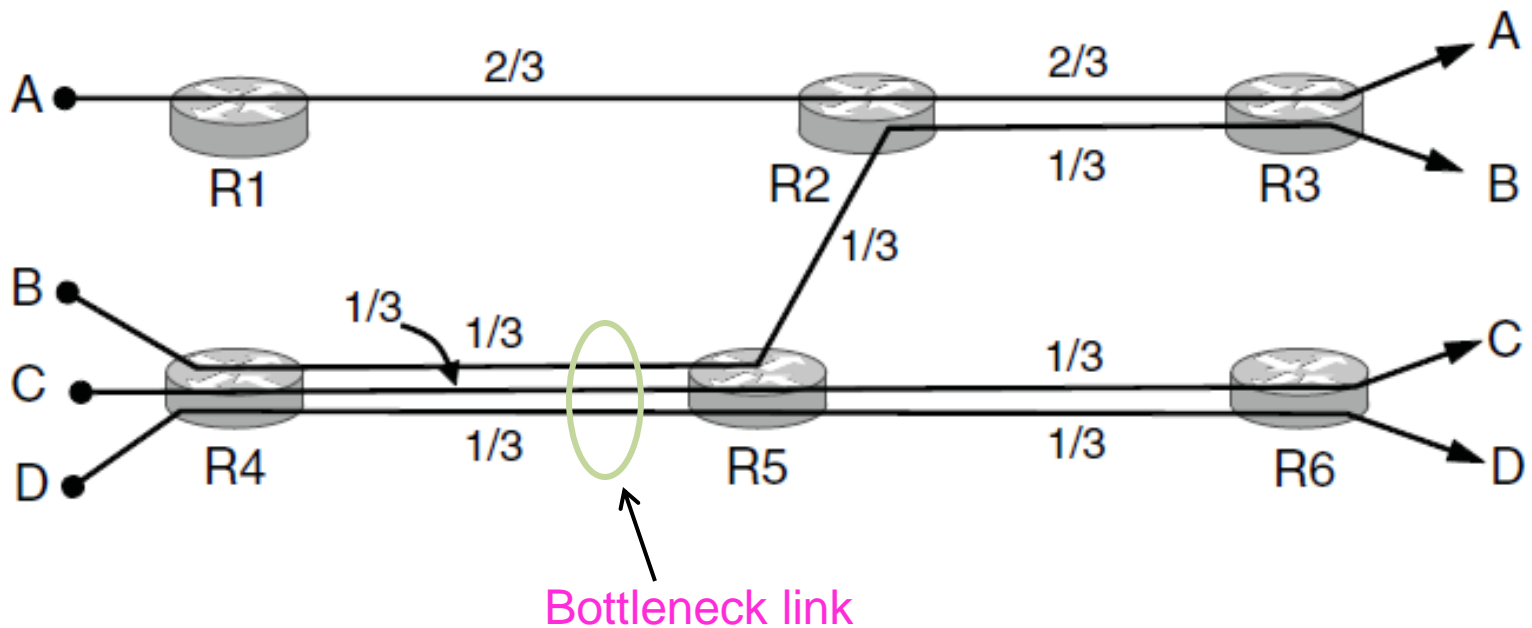


Delay begins to rise sharply when congestion sets in

# Desirable Bandwidth Allocation (2)

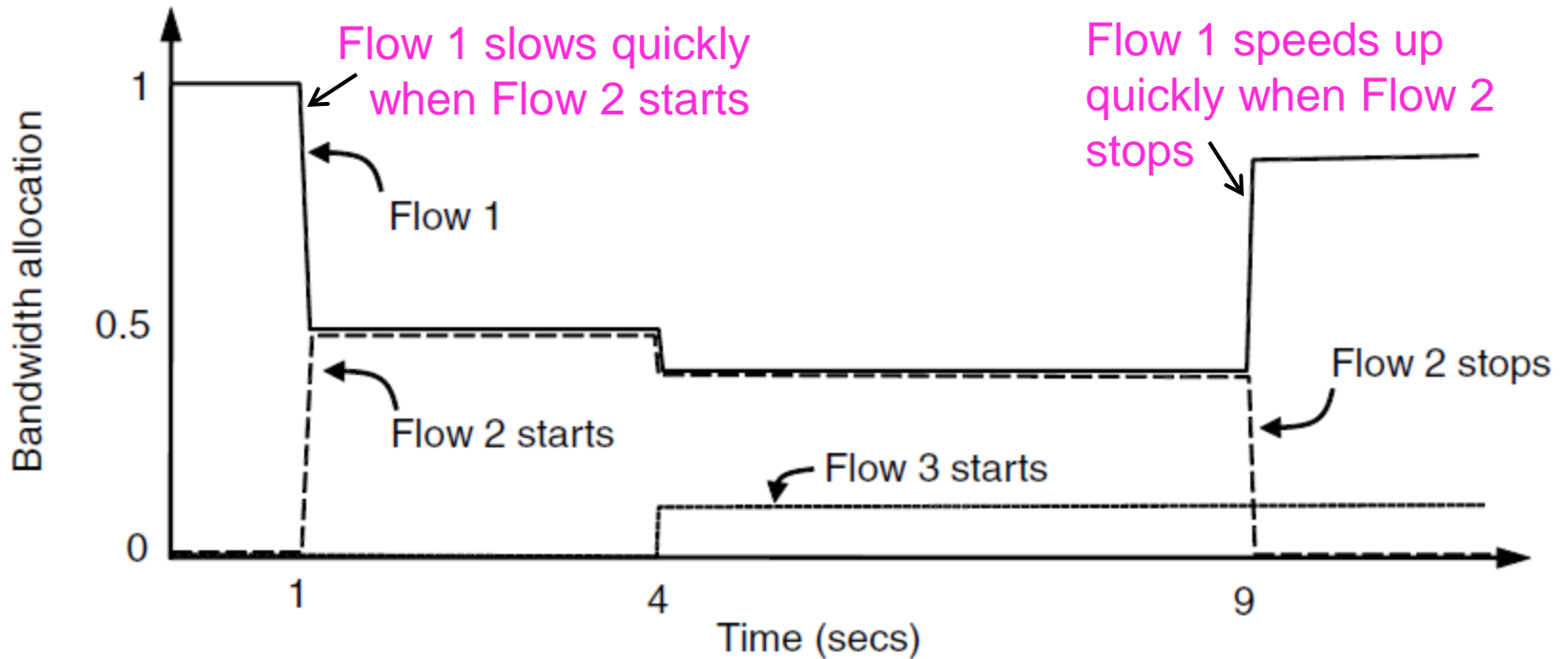
Fair use gives bandwidth to all flows (no starvation)

- Max-min fairness gives equal shares of bottleneck



# Desirable Bandwidth Allocation (3)

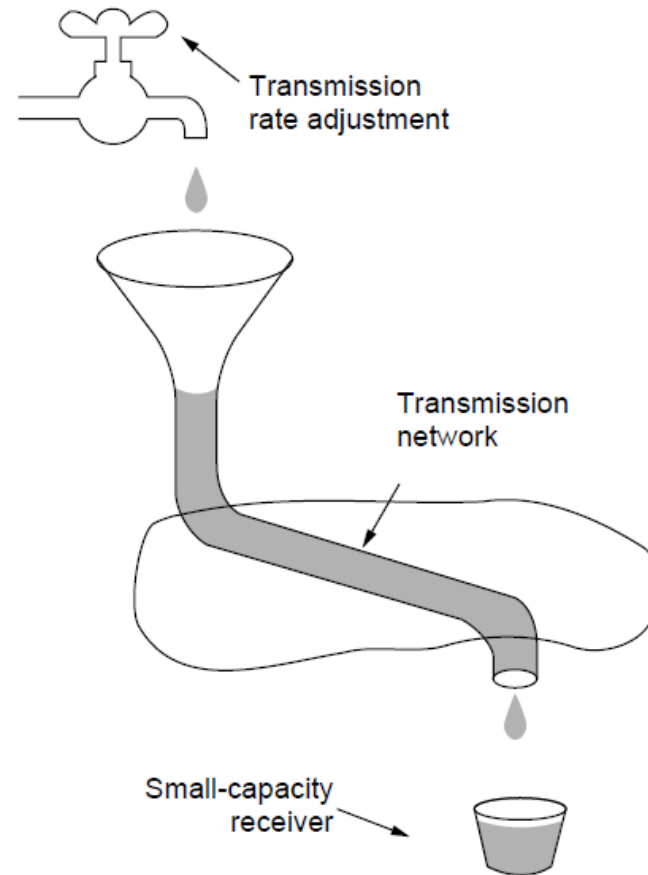
We want bandwidth levels to converge quickly when traffic patterns change



# Regulating the Sending Rate (1)

Sender may need to slow down for different reasons:

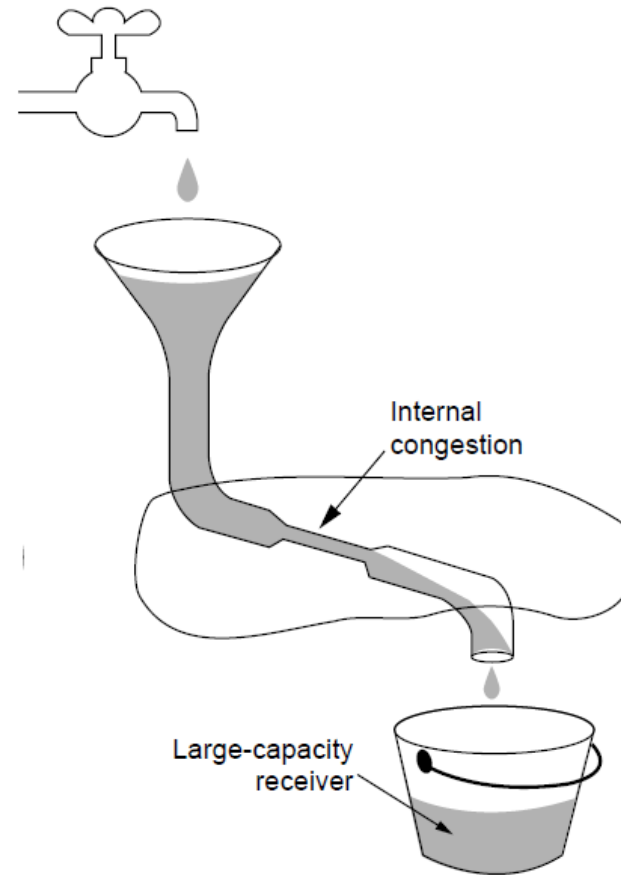
- Flow control, when the receiver is not fast enough [right]
- Congestion, when the network is not fast enough [over]



A fast network feeding a low-capacity receiver  
→ flow control is needed

# Regulating the Sending Rate (2)

Our focus is dealing with this problem – congestion



A slow network feeding a high-capacity receiver  
→ congestion control is needed

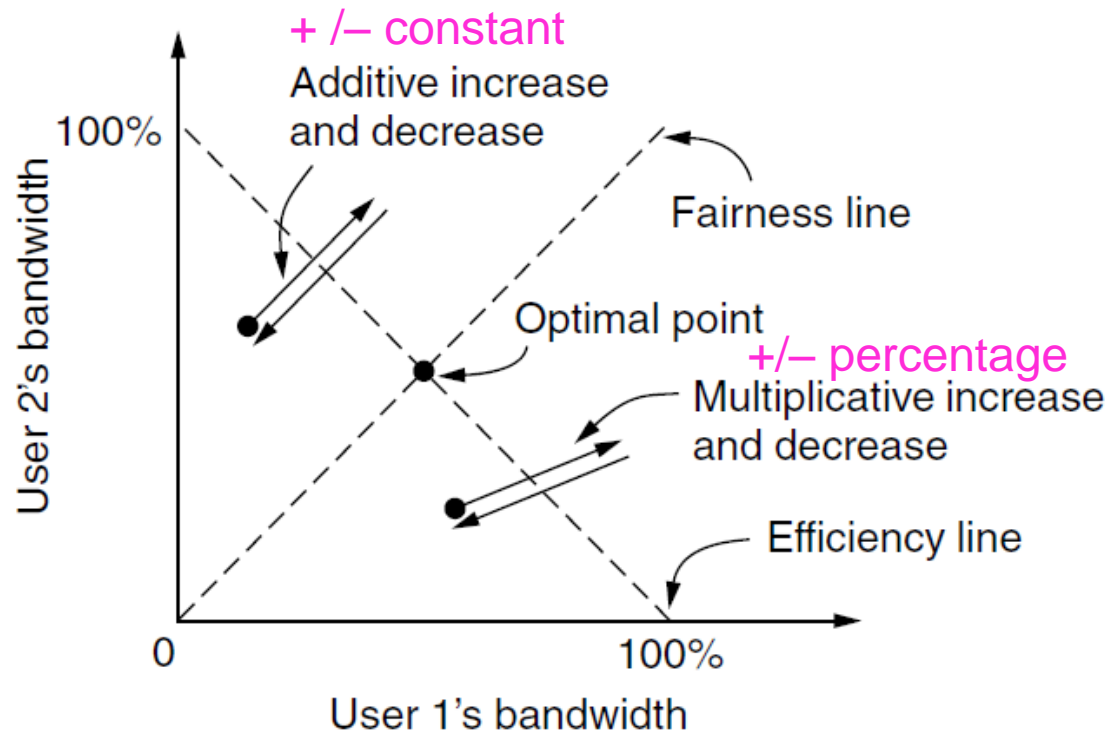
# Regulating the Sending Rate (3)

Different congestion signals the network may use to tell the transport endpoint to slow down (or speed up)

<b>Protocol</b>	<b>Signal</b>	<b>Explicit?</b>	<b>Precise?</b>
XCP	Rate to use	Yes	Yes
TCP with ECN	Congestion warning	Yes	No
FAST TCP	End-to-end delay	No	Yes
CUBIC TCP	Packet loss	No	No
TCP	Packet loss	No	No

# Regulating the Sending Rate (3)

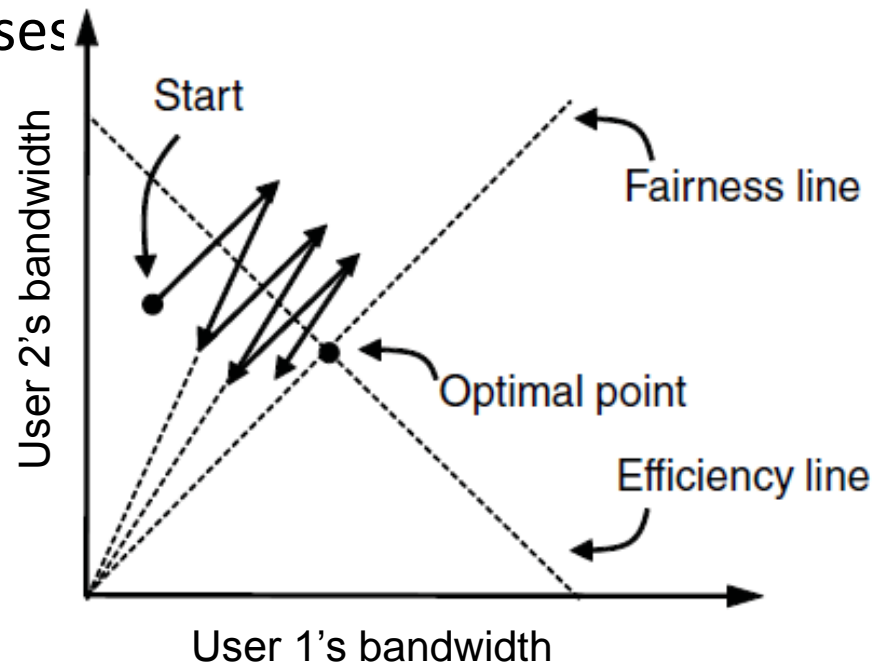
If two flows increase/decrease their bandwidth in the same way when the network signals free/busy they will not converge to a fair allocation



# Regulating the Sending Rate (4)

The AIMD (Additive Increase Multiplicative Decrease) control law does converge to a fair and efficient point!

– TCP uses



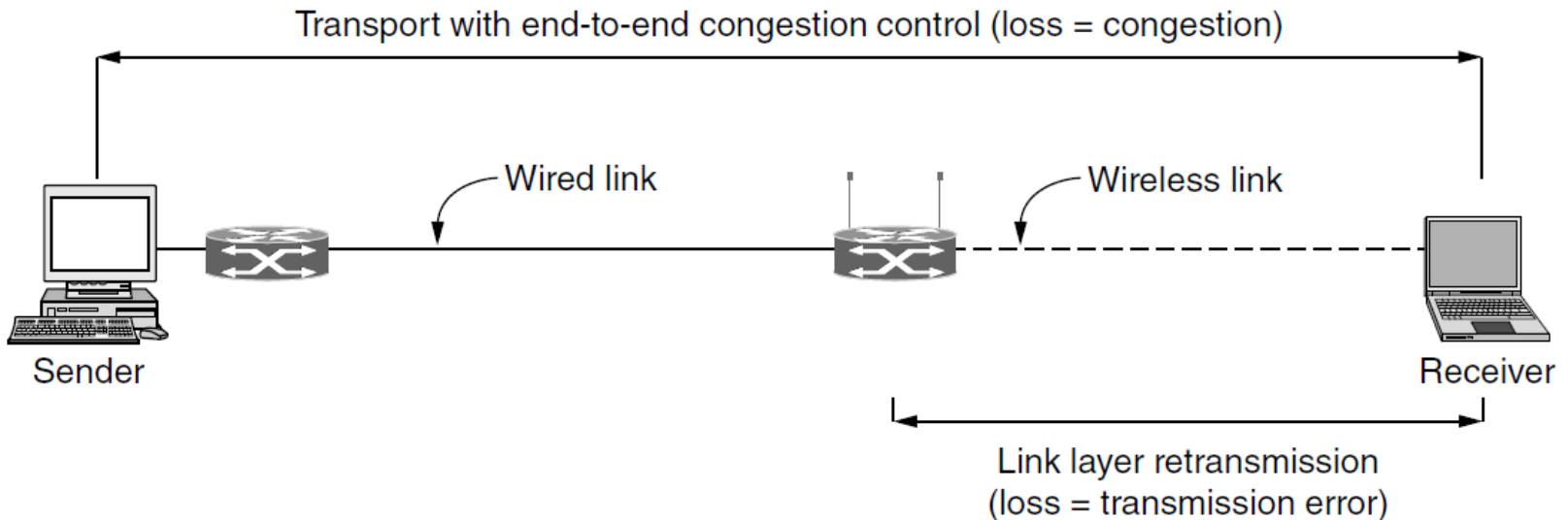
# Wireless Issues

Wireless links lose packets due to transmission errors

- Do not want to confuse this loss with congestion
- Or connection will run slowly over wireless links!

Strategy:

- Wireless links use ARQ, which masks errors



# Internet Protocols – UDP

- Introduction to UDP »
- Remote Procedure Call »
- Real-Time Transport »

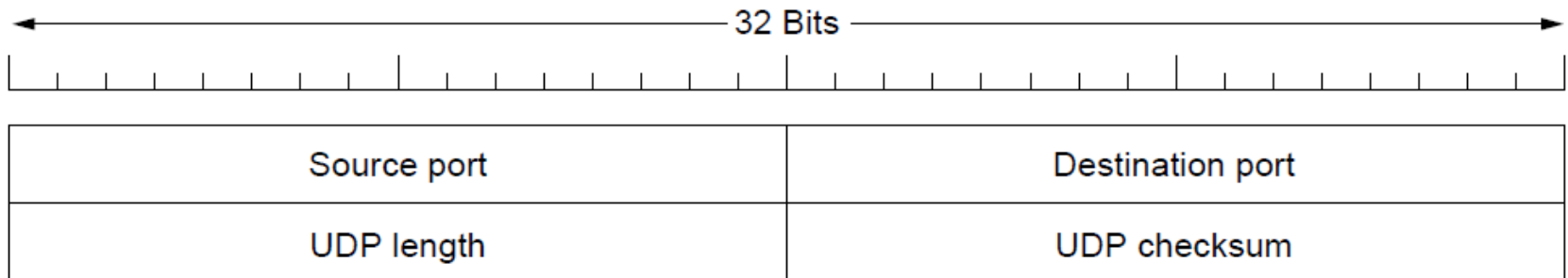
# User Datagram Protocol

- Connectionless
- Does not do
  - Flow control
  - Error control
  - Retransmissions
- Useful in client-server situations
- Sends segments consisting of an 8-byte header followed by the payload

# Introduction to UDP (1)

UDP (User Datagram Protocol) is a shim over IP

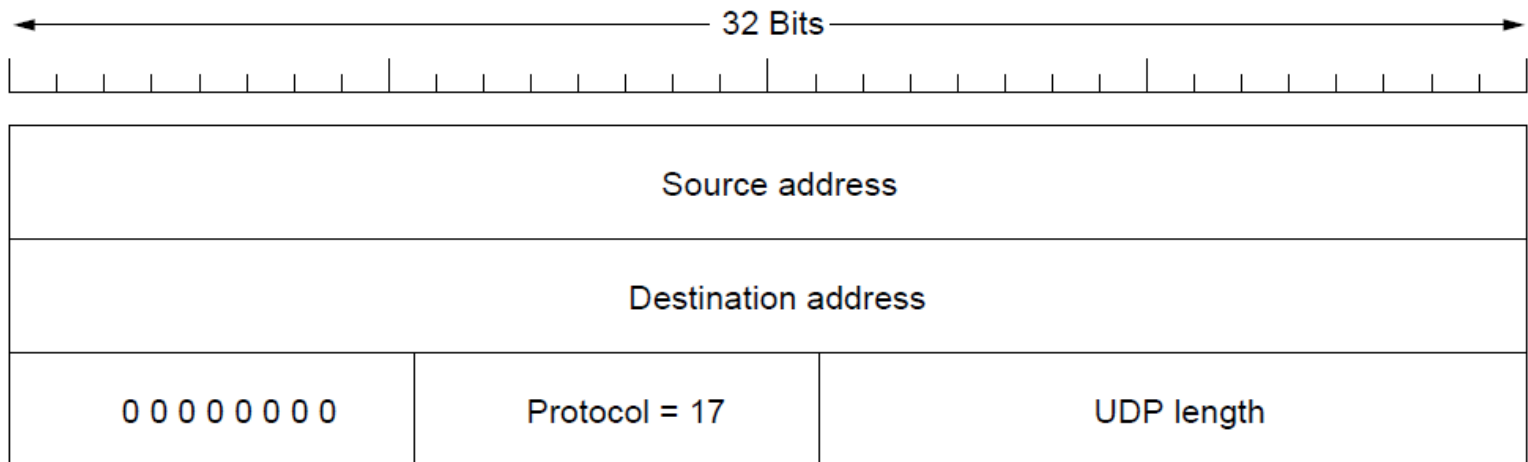
- Header has ports (TSAPs), length and checksum.



# Introduction to UDP (2)

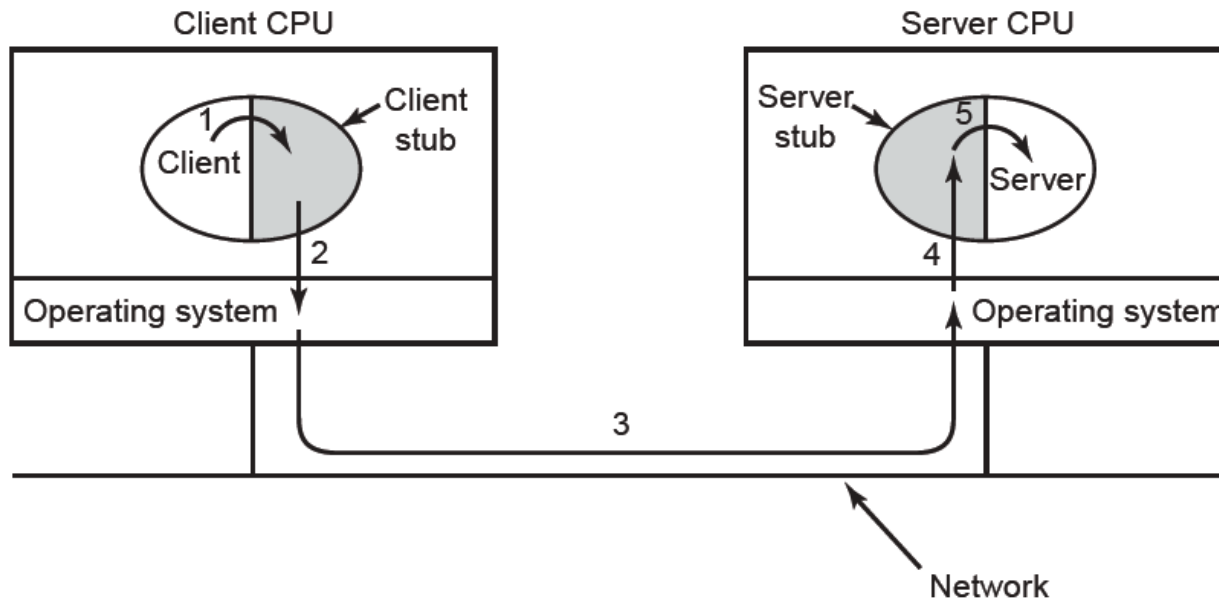
Checksum covers UDP segment and IP pseudoheader

- Fields that change in the network are zeroed out
- Provides an end-to-end delivery check



# RPC (Remote Procedure Call)

- RPC connects applications over the network with the familiar abstraction of procedure calls
  - Stubs package parameters/results into a message
  - UDP with retransmissions is a low-latency transport



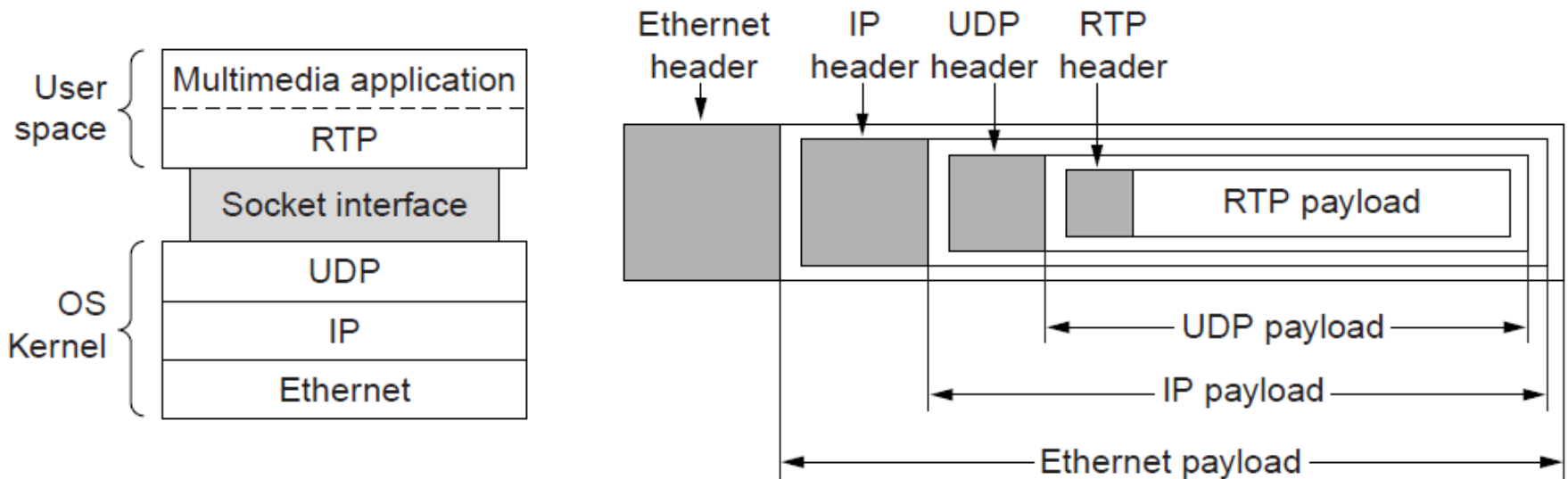
# Limitations of RPC

- Pointers
- Weakly Typed languages – variable length arrays
- Not possible always to deduce parameter types
- Global variables

# Real-Time Transport (1)

RTP (Real-time Transport Protocol) provides support for sending real-time media over UDP

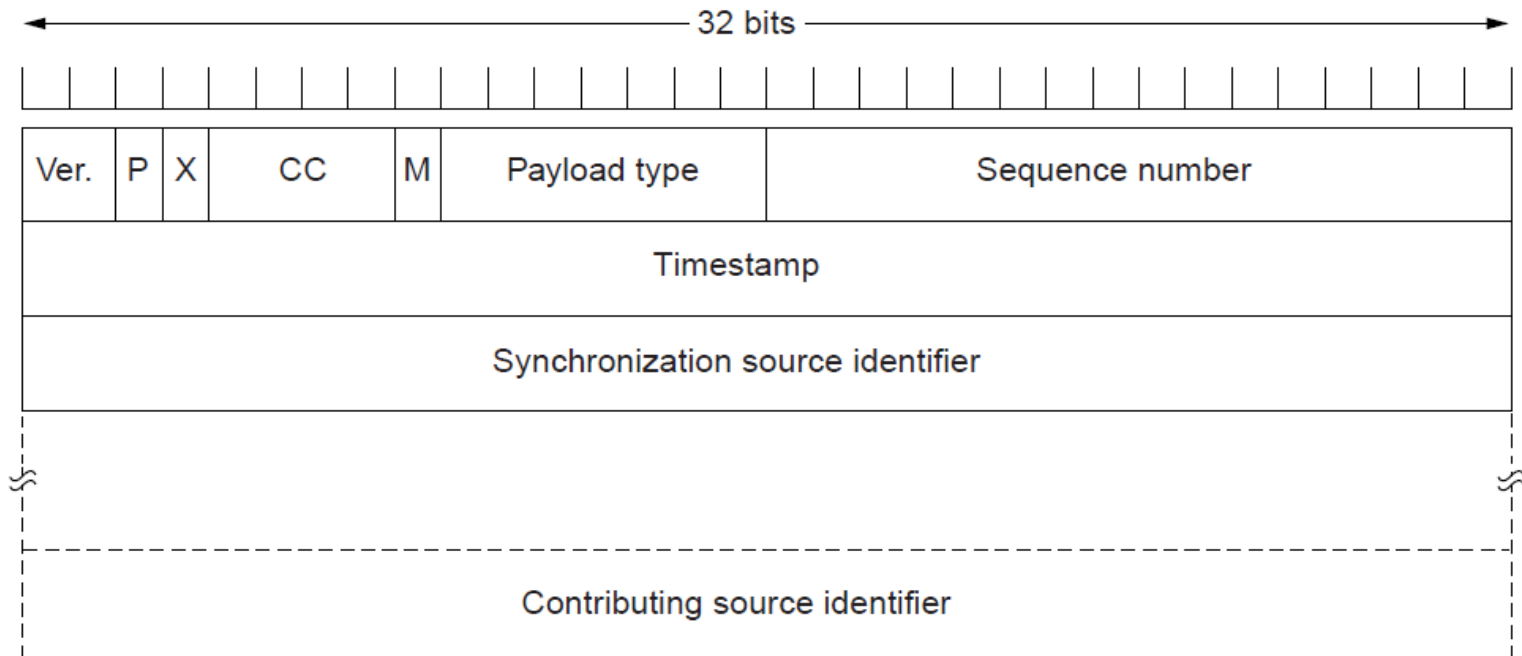
- Often implemented as part of the application



# Real-Time Transport (2)

RTP header contains fields to describe the type of media and synchronize it across multiple streams

- RTCP sister protocol helps with management tasks



# RTP Header Fields

- Ver – 2
- P – Packet padded to multiple of 4 bytes
- X – extension header present
- CC – number of contributing sources
- M bit – Application specific marker
- Payload Type – encoding used
- Sequence Number
- Time stamp – produced by the source
- Synchronizations Source Identifier – which stream the packet belongs to

# RTP Profiles

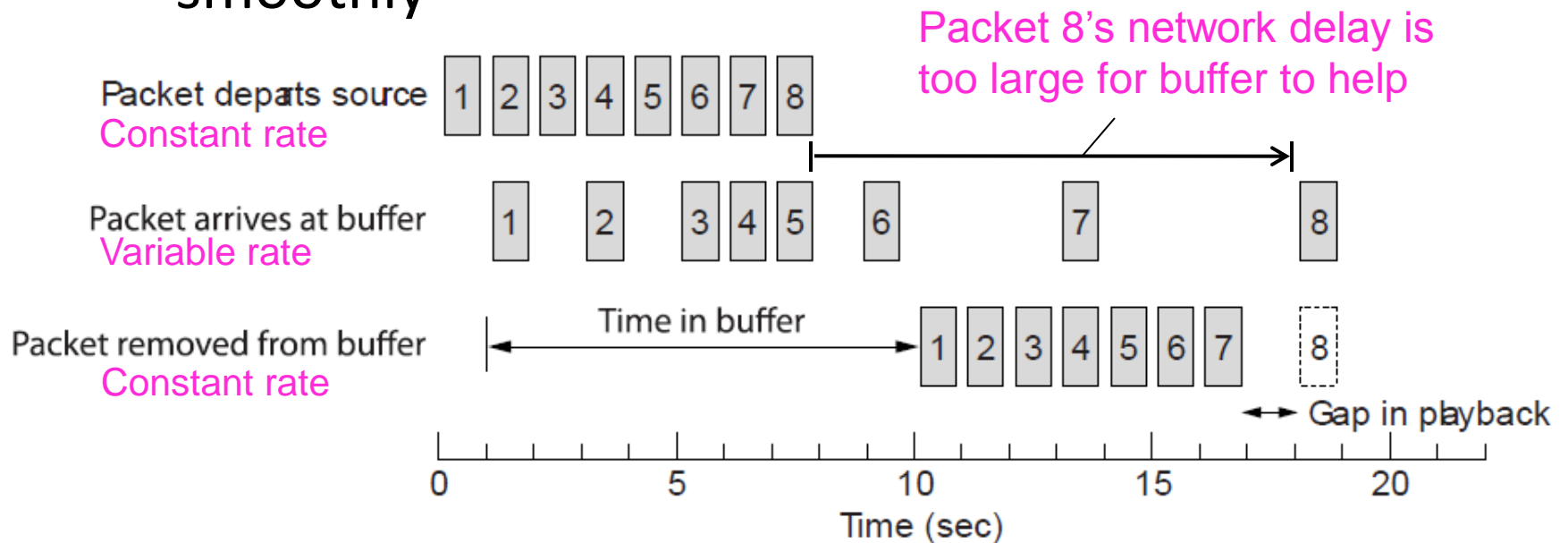
- RTP payloads may contain multiple samples coded in any way the application wants
- Profiles – to support interworking
  - Single Audio Stream
  - Multiple encoding formats may be supported
    - 8-bit pcm samples at 8KHz
    - Delta encoding
    - Predictive encoding
    - MP3
    - ...

# RTCP – Real-time Transport Control Protocol

- Control Protocol for RTP
- Does not transport any data
- Handles:
  - Feedback
    - Delay
    - Jitter
    - Bandwidth
    - Congestion, etc.
  - Synchronization
    - Interstream Synchronization – Different clocks, drifts, etc.
  - User Interface

# Real-Time Transport (3)

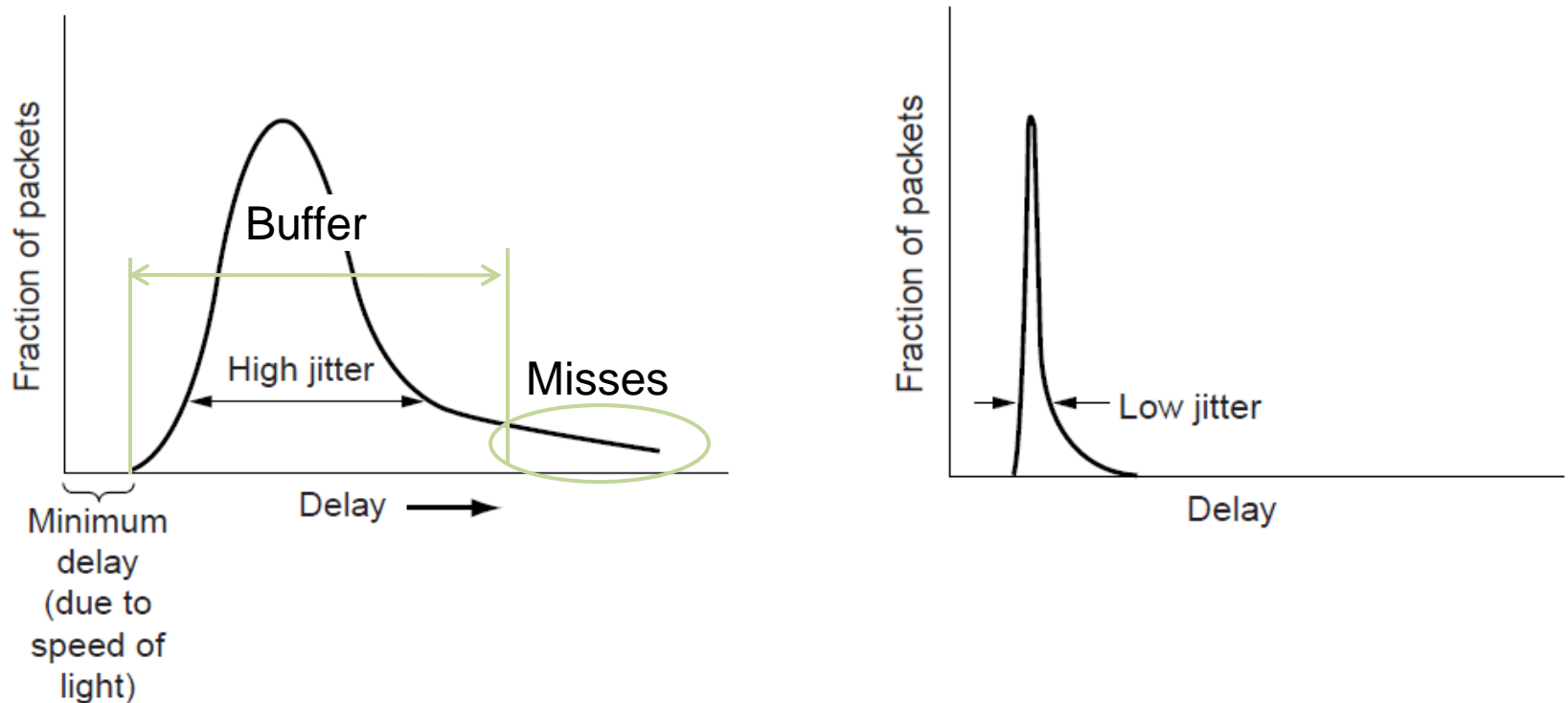
Buffer at receiver is used to delay packets and absorb jitter so that streaming media is played out smoothly



# Real-Time Transport (3)

High jitter, or more variation in delay, requires a larger playout buffer to avoid playout misses

- Propagation delay does not affect buffer size



# Internet Protocols – TCP

- The TCP service model »
- The TCP segment header »
- TCP connection establishment »
- TCP connection state modeling »
- TCP sliding window »
- TCP timer management »
- TCP congestion control »

# TCP

- Reliable end-to-end byte stream over an unreliable internetwork
- Dynamically adapt to the properties of the network
- Robust
  
- RFC 793 – 1122 and 1323

# TCP Transport Entity

- Implemented as
  - Library procedures
  - User process, or
  - Part of the kernel
- Accepts data streams from local processes
- Breaks them into segments of >64KB (usually 1460 bytes)
- Sends each piece in a separate IP packet
- On receive side – reconstruct the original byte stream and give to a process.
- Must recover from errors – time outs, retransmissions, etc.

# TCP Service Model

- End points called sockets
  - Socket number
    - IP address of the host
    - 16-bit number (called the *port*)
- Connection Oriented – Full Duplex, Point-to-Point
  - Establish a connection between sockets
  - An socket may be used for multiple connections at the same time
  - Connection (*socket1, socket2*)

# Berkeley Sockets

The socket primitives for TCP.

<b>Primitive</b>	<b>Meaning</b>
SOCKET	Create a new communication end point
BIND	Attach a local address to a socket
LISTEN	Announce willingness to accept connections; give queue size
ACCEPT	Block the caller until a connection attempt arrives
CONNECT	Actively attempt to establish a connection
SEND	Send some data over the connection
RECEIVE	Receive some data from the connection
CLOSE	Release the connection

# The TCP Service Model (1)

TCP provides applications with a reliable byte stream between processes; it is the workhorse of the Internet

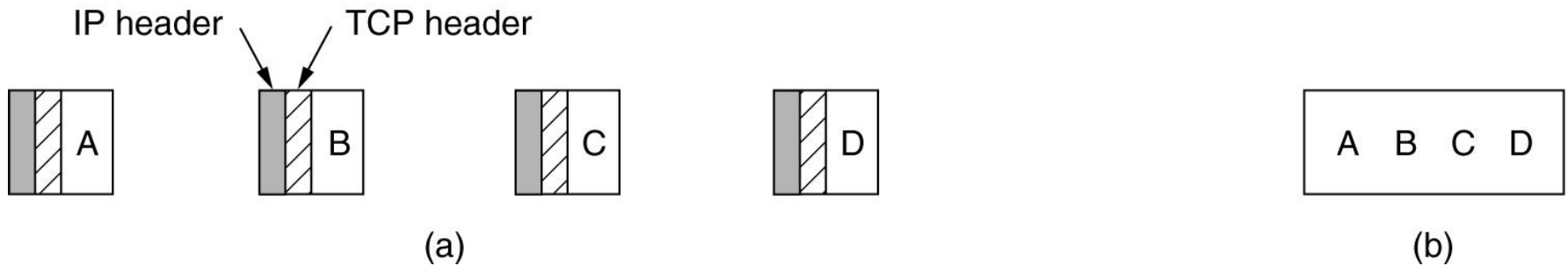
- Popular servers run on well-known ports

Port	Protocol	Use
20, 21	FTP	File transfer
22	SSH	Remote login, replacement for Telnet
25	SMTP	Email
80	HTTP	World Wide Web
110	POP-3	Remote email access
143	IMAP	Remote email access
443	HTTPS	Secure Web (HTTP over SSL/TLS)
543	RTSP	Media player control
631	IPP	Printer sharing

# Internet Daemon

- Attaches to multiple well-known ports and waits
- When a connection comes in it forks off a new process and executes the appropriate daemon
- That daemon handles the request

# The TCP Service Model



- (a) Four 512-byte segments sent as separate IP datagrams.
- (b) The 2048 bytes of data delivered to the application in a single READ CALL.

# TCP Service

- No message boundaries are preserved
- Send data as received or buffer it
- PUSH Flag
  - Send data now
  - Useful in sending command from terminal
- Urgent Data
  - DEL or CTRL-C to break off a remote computation
  - Use URGENT flag – Transmit everything right now
    - Receiving application is interrupted

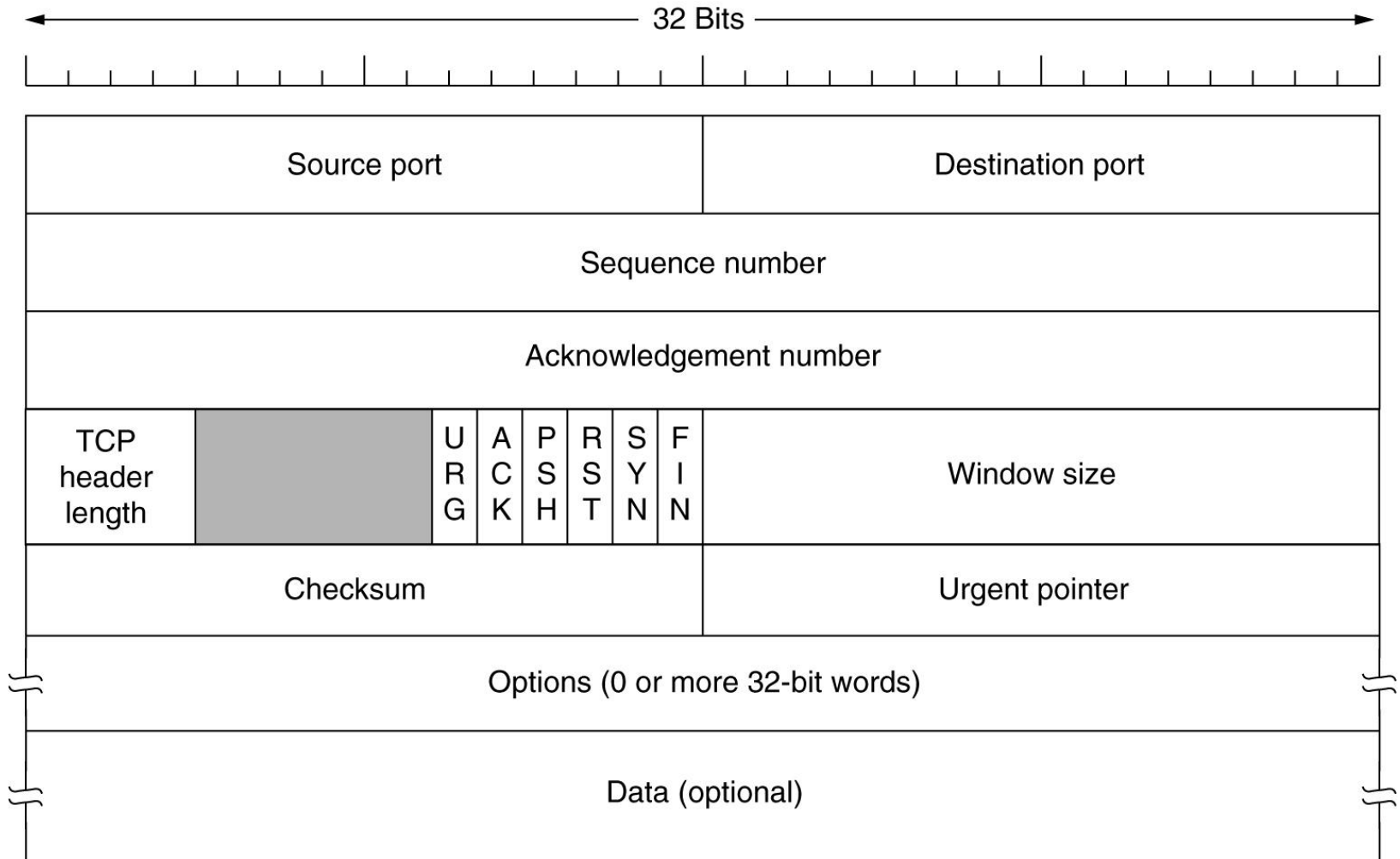
# TCP Protocol

- Exchange segments
  - 20 byte header (plus optional parts)
  - 0 or more data bytes
- Accumulate data from several writes into one segment
- May split data from one write into multiple segments
- Each segment, including the header, must be <65515 byte IP payload
- Each network has MTU- Maximum Transfer Unit
  - Each segment must be less than or equal to MTU

# TCP Protocol

- TCP uses sliding window protocol
- Sequence numbers are for bytes not segments
- Sending – start a timer
- Receiving – send an ack with sequence number = next sequence number expected

# The TCP Segment Header

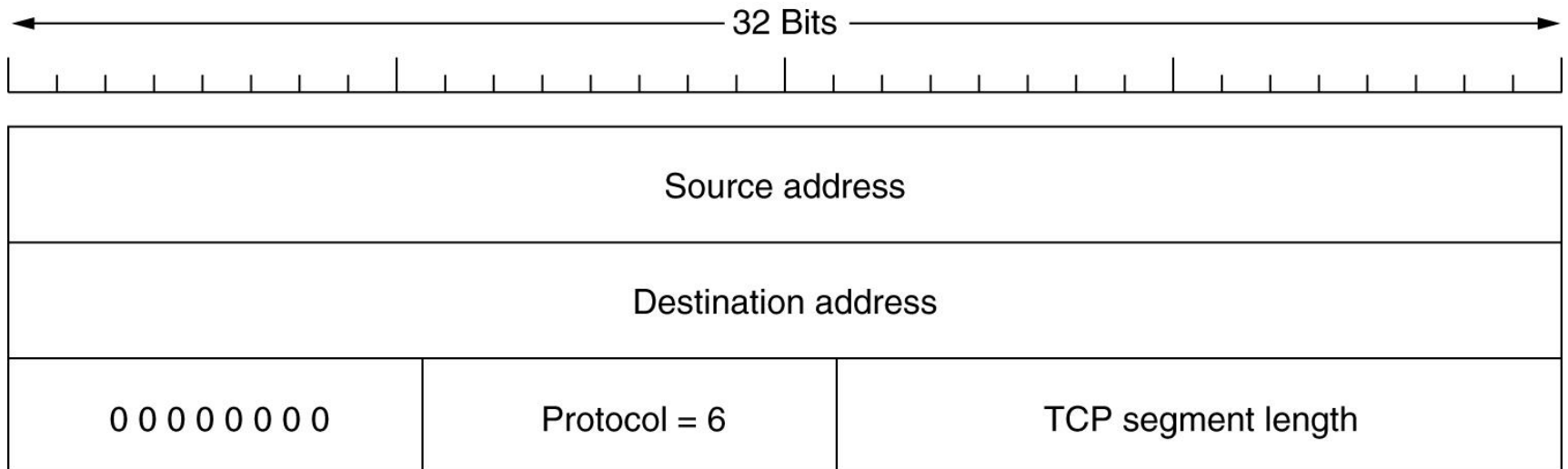


# TCP Segment Header

- Source and Destination Port
- Sequence number
- Acknowledgement Number
  - Next byte expected
- TCP Header Length – Number of 32 bit words in TCP header
- Flags
  - URG – set to 1 if Urgent Pointer is in use
    - Used to indicate a byte offset from the current seq no at which urgent data care there
  - ACK – set to 1 when ack no is valid
  - PSH – Push bit
  - RST – Reset
  - SYN – Used for connection establishment
  - FIN – Used to close a connection

# The TCP Segment Header (2)

- The pseudoheader included in the TCP checksum
- Checksum the header, the data and the pseudoheader
- Add all 16-bit words in 1's complement and then take 1's complement of the sum
- To check calculate on the entire segment and result should be 0.



# TCP Window

- Window size tells how many bytes may be sent starting at the byte acknowledged
  - If 0 means do not send now.
  - May send a segment with same ack no and non-zero window size.

# Maximum segment Size

- All hosts are required to accept TCP segments of  $536+20 = 556$  bytes
- May negotiate max segment size using options.
- Another negotiable parameter – Window Scale
  - May shift to the left by up to 14 bits
  - Giving a max window size of  $2^{30}$  bytes

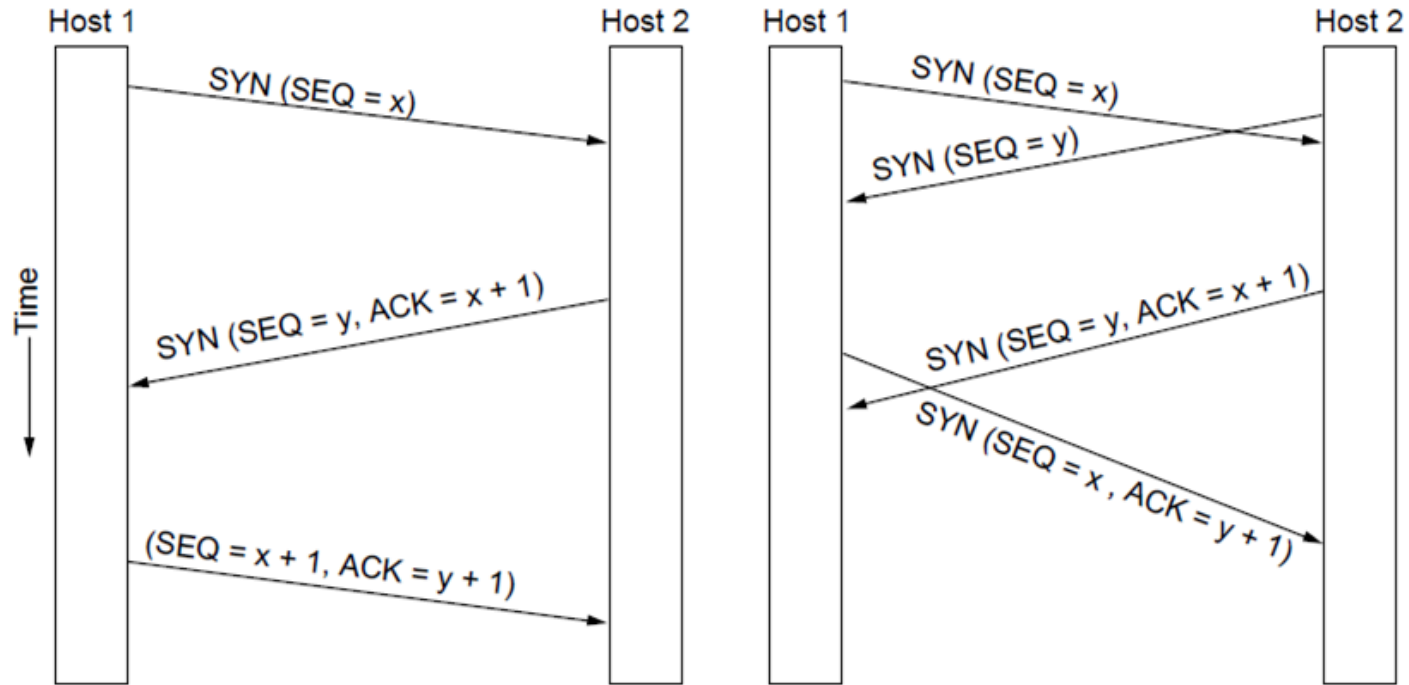
# Connection Establishment

- Uses three-way handshake
- Server passively listens
  - LISTEN and ACCEPT primitives
- Client executes CONNECT
  - Send a TCP Segment with SYN bit on and ACK bit off.
- Check to see if there is a process listening
  - If not send a RST
  - If yes, then give the segment to the process
  - If accepted – send an ACK message

# TCP Connection Establishment

TCP sets up connections with the three-way handshake

- Release is symmetric, also as described before



Normal case

Simultaneous connect

# TCP Connection State Modeling (1)

The TCP connection finite state machine has more states than our simple example from earlier.

State	Description
CLOSED	No connection is active or pending
LISTEN	The server is waiting for an incoming call
SYN RCVD	A connection request has arrived; wait for ACK
SYN SENT	The application has started to open a connection
ESTABLISHED	The normal data transfer state
FIN WAIT 1	The application has said it is finished
FIN WAIT 2	The other side has agreed to release
TIME WAIT	Wait for all packets to die off
CLOSING	Both sides have tried to close simultaneously
CLOSE WAIT	The other side has initiated a release
LAST ACK	Wait for all packets to die off

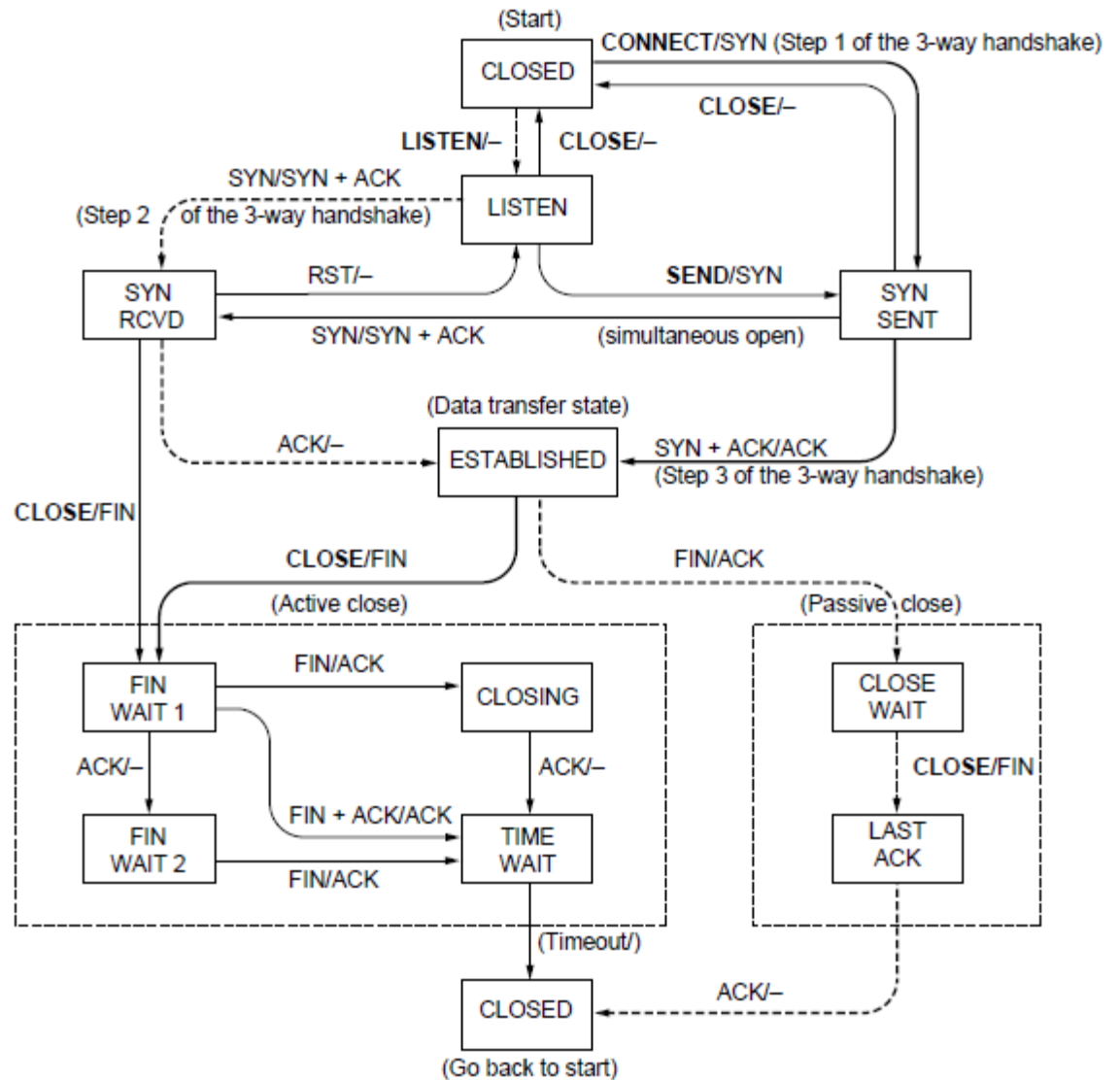
# TCP Connection State Modeling (2)

Solid line is the normal path for a client.

Dashed line is the normal path for a server.

Light lines are unusual events.

Transitions are labeled by the cause and action, separated by a slash.



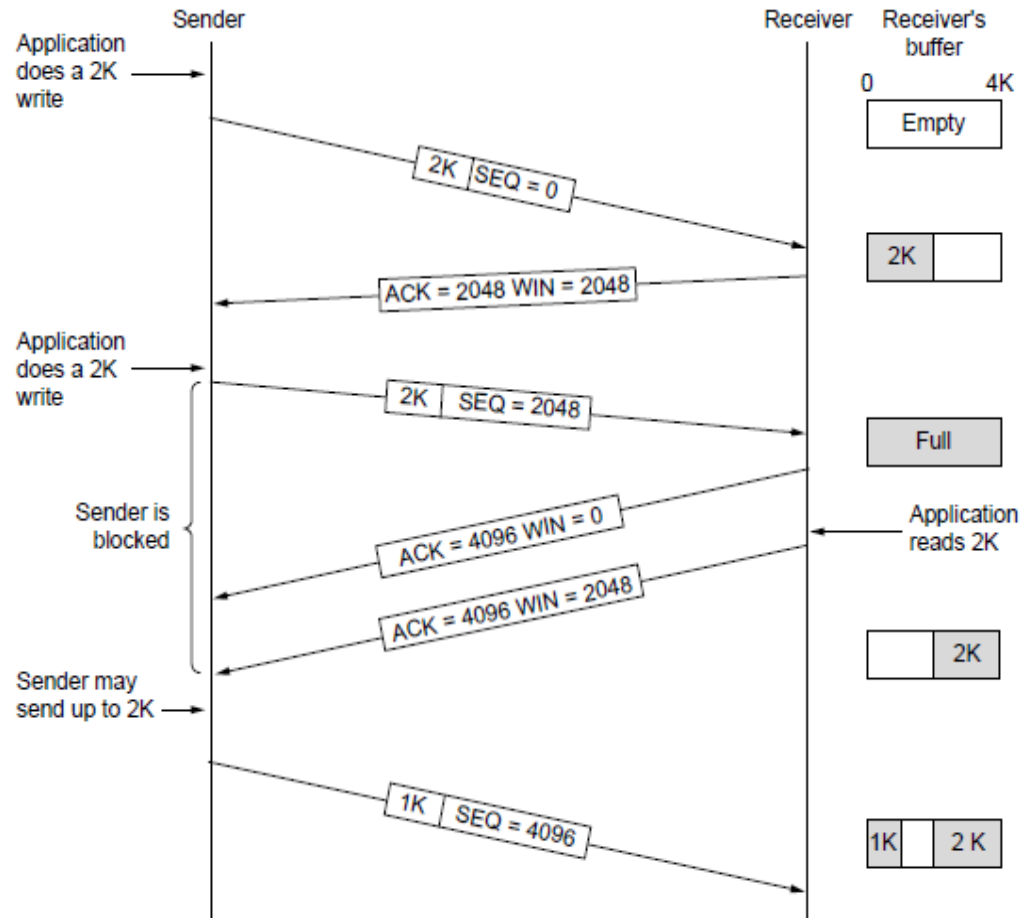
# Connection Release

- Think of the connection as a pair of simplex connections
  - Each is released independently
- Either party sends a segment with FIN bit set
- When FIN acked that direction is shut down

# TCP Sliding Window (1)

TCP adds flow control to the sliding window as before

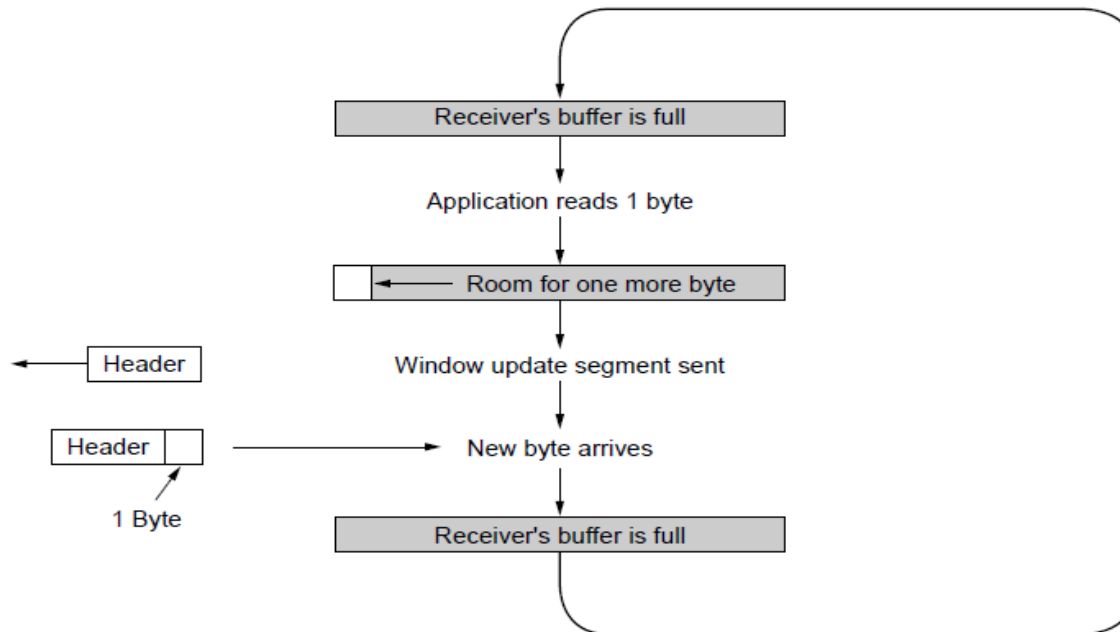
- ACK + WIN is the sender's limit



# TCP Sliding Window (2)

Need to add special cases to avoid unwanted behavior

- E.g., silly window syndrome [below]



Receiver application reads single bytes, so sender always sends one byte segments

# Timer Management

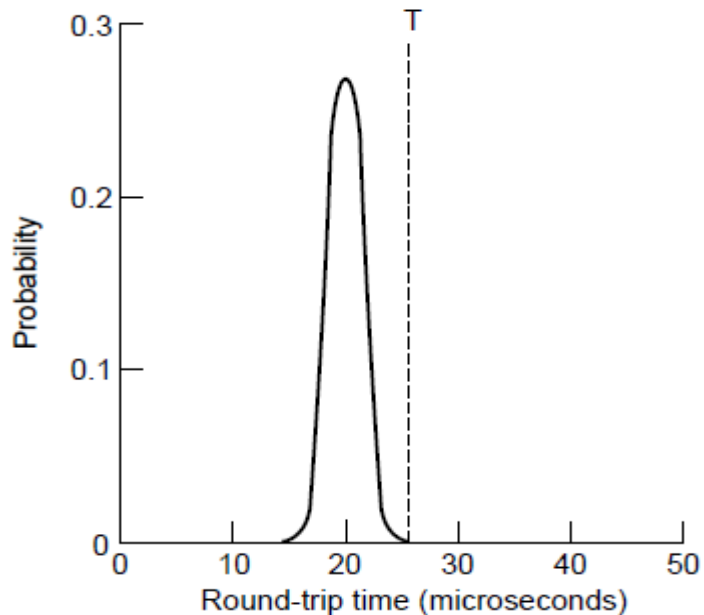
- TCP uses multiple timers
  - Most important is the *retransmission timer*
  - What value to set it at??
- Round Trip time
  - Highly variable
  - Varies over time
  - Have to track it
  - Estimate it
  - M is a new measurement
  - $\alpha = 7/8$
  - Use  $\beta RTT$  for retransmission timer
  - Initial values of  $\beta$  were 2 – make is proportional to standard deviation of M

$$RTT = \alpha RTT + (1 - \alpha)M$$

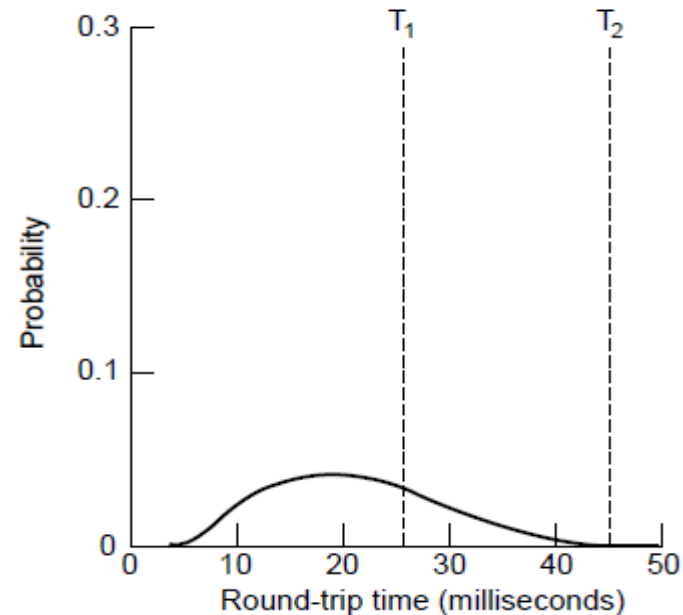
# TCP Timer Management

TCP estimates retransmit timer from segment RTTs

- Tracks both average and variance (for Internet case)
- Timeout is set to average plus 4 x variance



LAN case – small,  
regular RTT



Internet case –  
large, varied RTT

# Timer Management

- Jacobson Approach
- Mean deviation estimate
- $D = \gamma D + (1-\gamma) |RTT - M|$
- Timeout = RTT + 4 D
  
- What to do on retransmissions
  - Do not know if the ack is for the first or second
- Karn Algorithm
  - Do not update RTT on any segments that have been retransmitted

# Persistence Timer

- Receiver sends a window of 0
- Later sends a window size but that packet is lost
  - Both wait
- Persistence Timer
  - When it goes off – sender sends a probe request to receiver to get a window size
  - If still zero – continue to wait and reset persistence timer
- Keepalive Timer
  - When a connection is idle for a long time – check if the other side is still there

# TCP Congestion Control

- Congestion – a function of total number of packets in the network, and where they are
- First step – detection
  - Is packet loss an indication of congestion??
  - All TCP algorithms assume timeouts are caused by congestion
- Initial steps
  - When connection is established – use suitable window size
    - Loss will not occur due to buffers at receiver
- Two issues
  - Network Capacity
  - Receiver Capacity

# TCP Congestion Control

- Network Capacity and Receiver Capacity
- Maintain two windows
  - Receiver window
  - Congestion window
  - Use the min ( Receiver window and Congestion window)
- Initially
  - Sender sets congestion window to MSS (Max Seg Size)
  - If acked add one more MSS – 2 now
  - Repeat for each acked MSS
  - Congestion window grows exponentially
  - If timeout – go back to previous window size
  - SLOW START

# Internet Congestion Control

- Use a Threshold – initially 64 KB
- When a timeout occurs set threshold to half the current congestion window and reset congestion window to 1 MSS
- Use slow start till the threshold is reached
- Then successful transmissions grow congestion window linearly

# TCP Congestion Control (1)

TCP uses AIMD with loss signal to control congestion

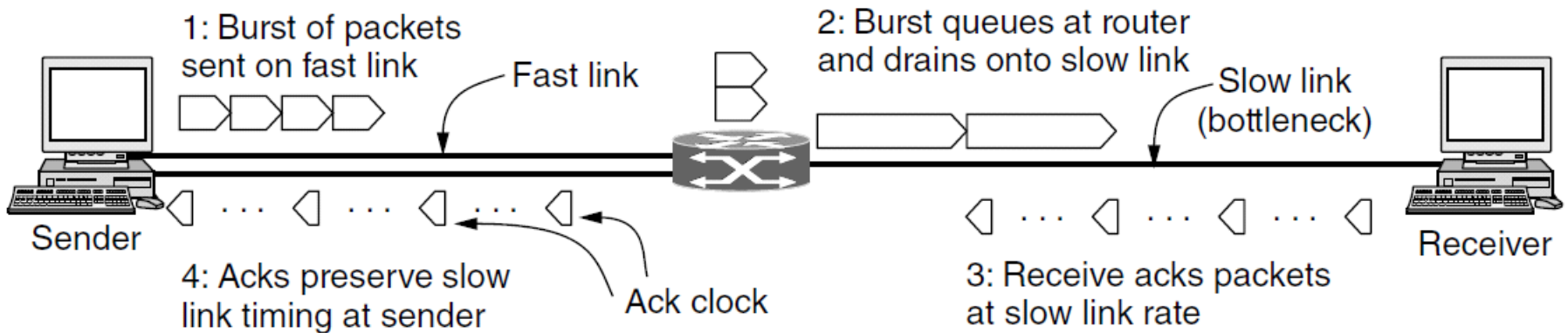
- Implemented as a congestion window (cwnd) for the number of segments that may be in the network
- Uses several mechanisms that work together

Name	Mechanism	Purpose
ACK clock	Congestion window (cwnd)	Smooth out packet bursts
Slow-start	Double cwnd each RTT	Rapidly increase send rate to reach roughly the right level
Additive Increase	Increase cwnd by 1 packet each RTT	Slowly increase send rate to probe at about the right level
Fast retransmit / recovery	Resend lost packet after 3 duplicate ACKs; send new packet for each new ACK	Recover from a lost packet without stopping ACK clock

# TCP Congestion Control (2)

Congestion window controls the sending rate

- Rate is  $cwnd / RTT$ ; window can stop sender quickly
- ACK clock (regular receipt of ACKs) paces traffic and smoothes out sender bursts

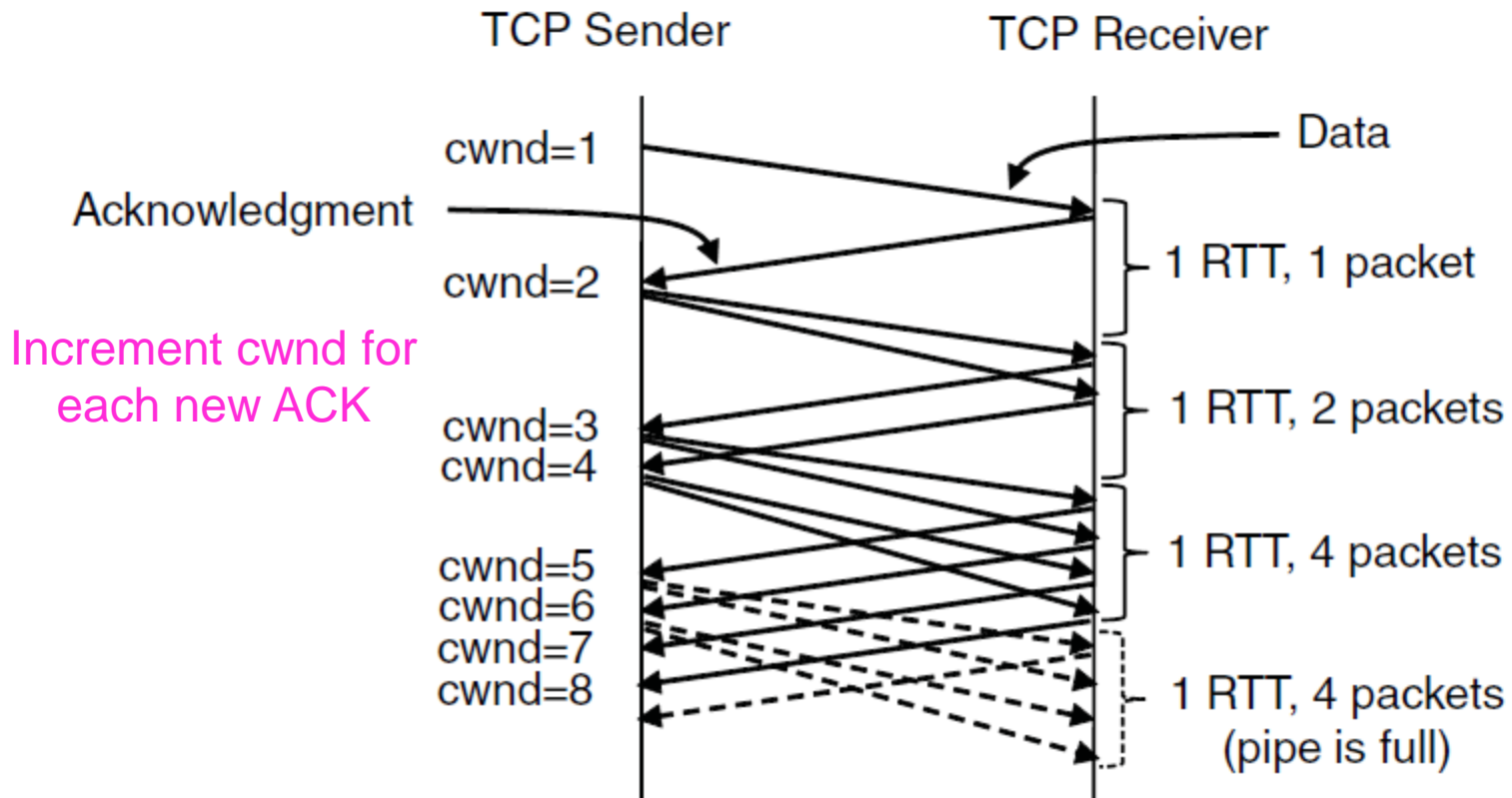


ACKs pace new segments into the network and smooth bursts

# TCP Congestion Control (3)

Slow start grows congestion window exponentially

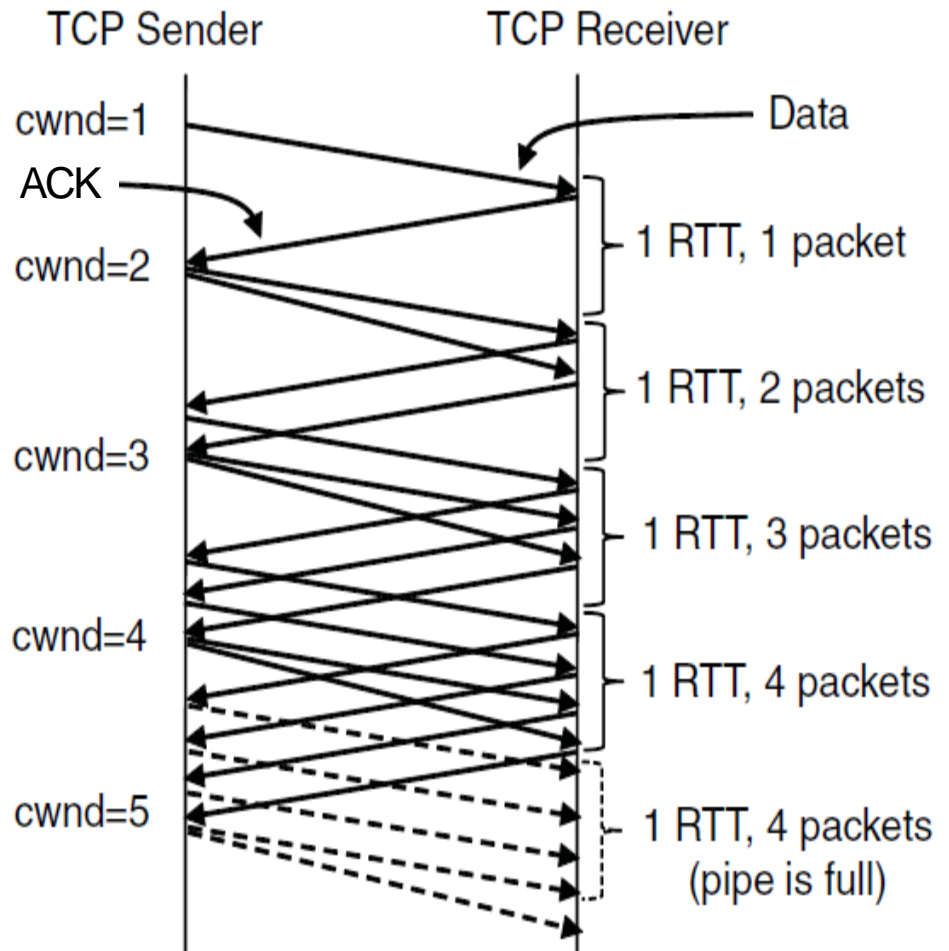
- Doubles every RTT while keeping ACK clock going



# TCP Congestion Control (4)

Additive increase grows  
cwnd slowly

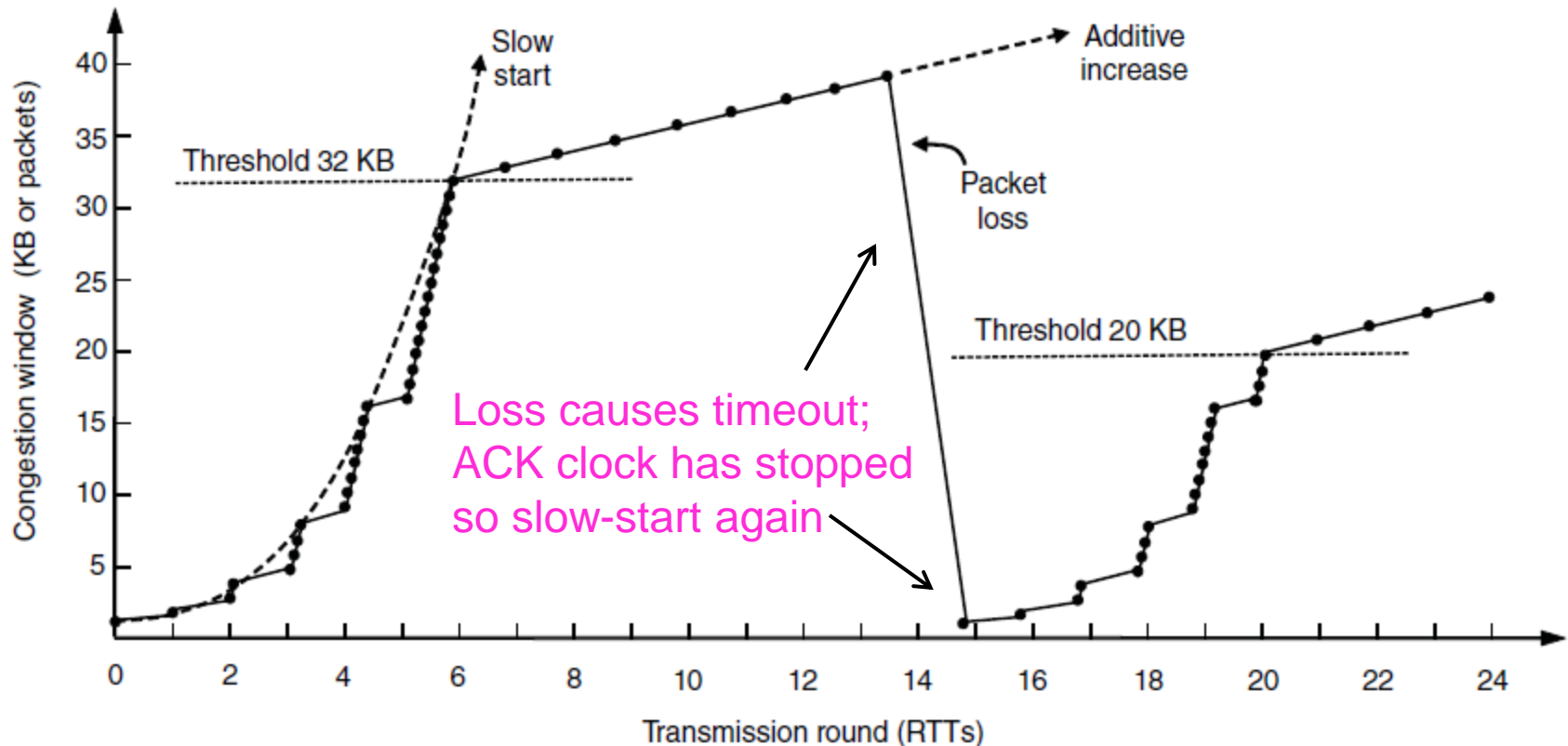
- Adds 1 every RTT
- Keeps ACK clock



# TCP Congestion Control (5)

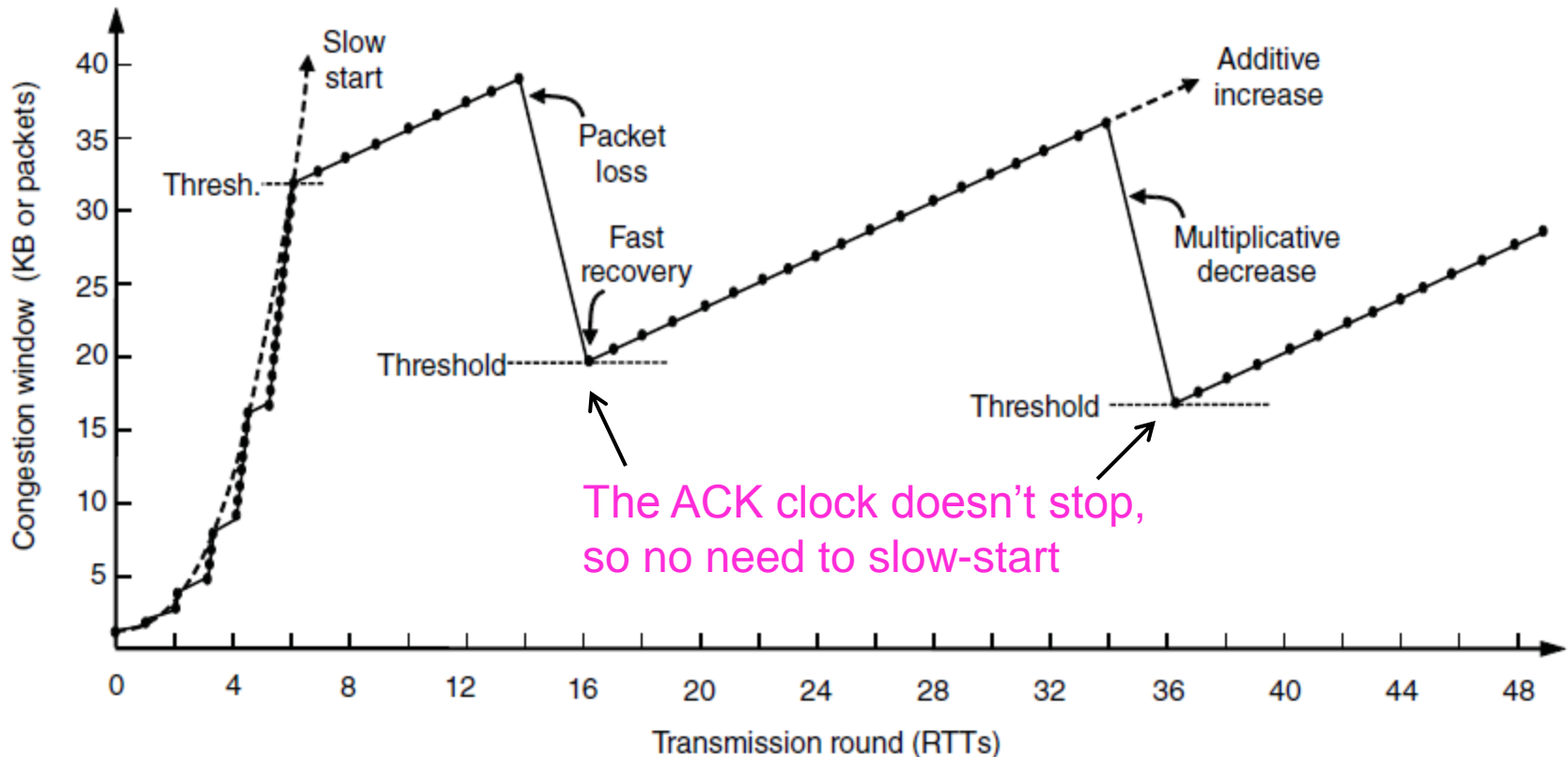
Slow start followed by additive increase (TCP Tahoe)

- Threshold is half of previous loss cwnd



# TCP Congestion Control (6)

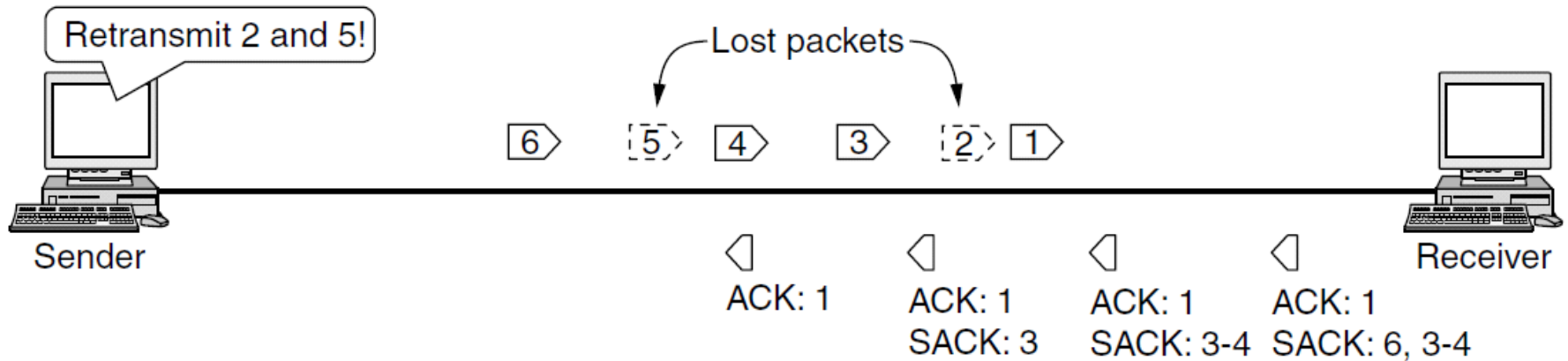
- With fast recovery, we get the classic sawtooth (TCP Reno)
  - Retransmit lost packet after 3 duplicate ACKs



# TCP Congestion Control (7)

SACK (Selective ACKs) extend ACKs with a vector to describe received segments and hence losses

- Allows for more accurate retransmissions / recovery



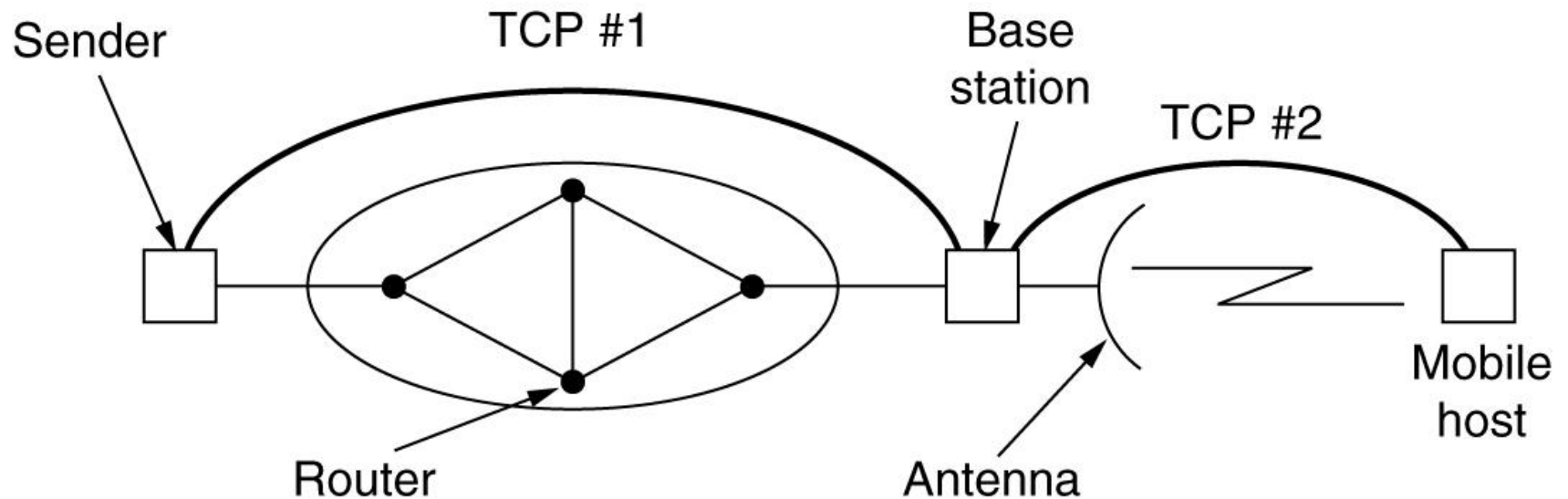
No way for us to know that 2 and 5 were lost with only ACKs

# Handling Silly Window Problem

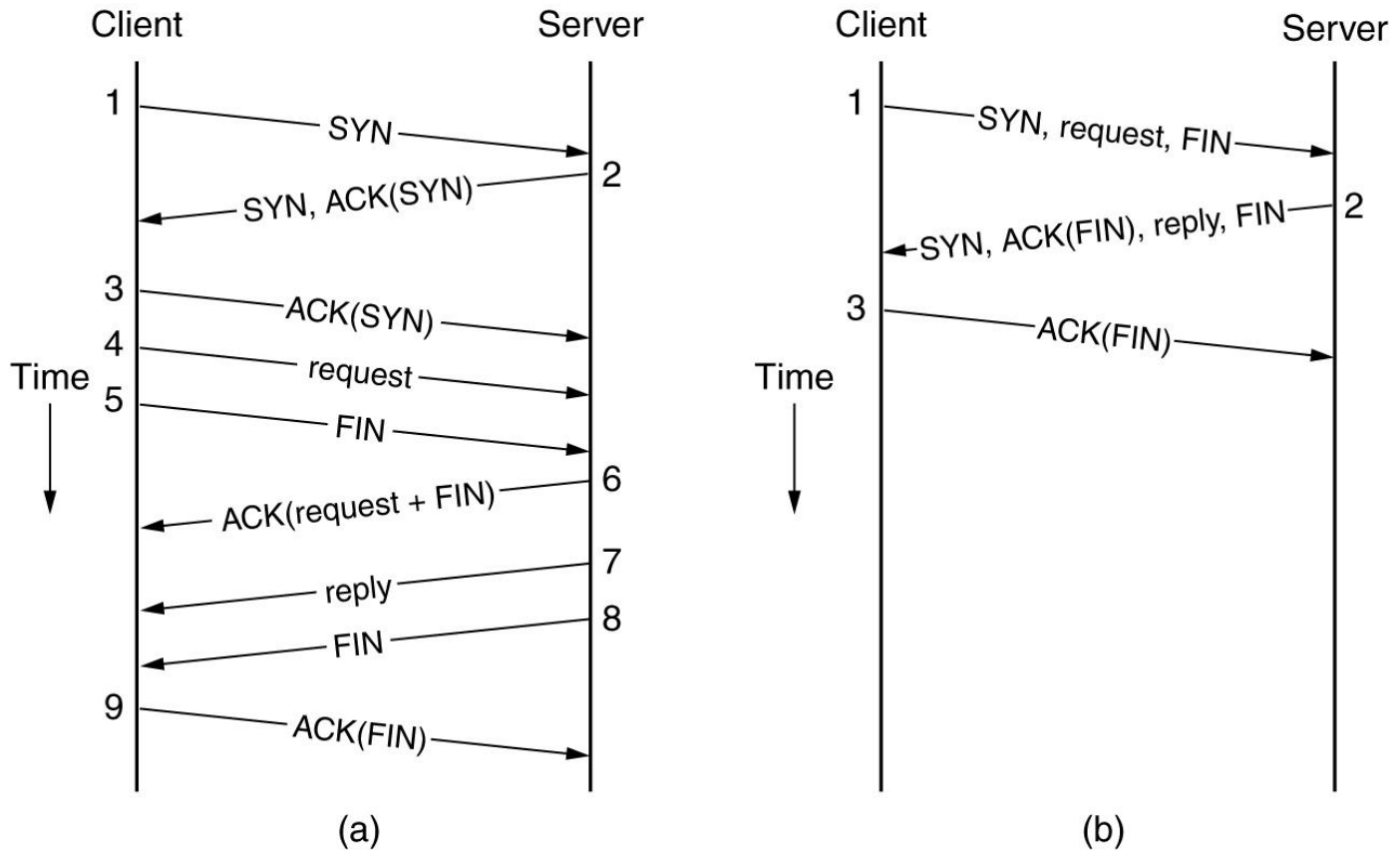
- Delay ack and window updates for 500 ms.
- Nagle's Algorithm
  - When data comes in one byte at a time
    - Send the first byte and buffer the rest till the outstanding byte is acknowledged
    - Then send all the buffered characters in one TCP segment
  - Mouse movements have to be sent – Burst does not work well.
- Clark's Solution
  - Wait until decent amount of space available then advertise
    - Max segment size or half buffer
  - Sender not send tiny segments

# Wireless TCP and UDP

Splitting a TCP connection into two connections.



# Transactional TCP



**(a)** RPC using normal TCP.

**(b)** RPC using T/TCP.

# Performance Issues

Many strategies for getting good performance have been learned over time

- Performance problems »
- Measuring network performance »
- Host design for fast networks »
- Fast segment processing »
- Header compression »
- Protocols for “long fat” networks »

# Performance Problems

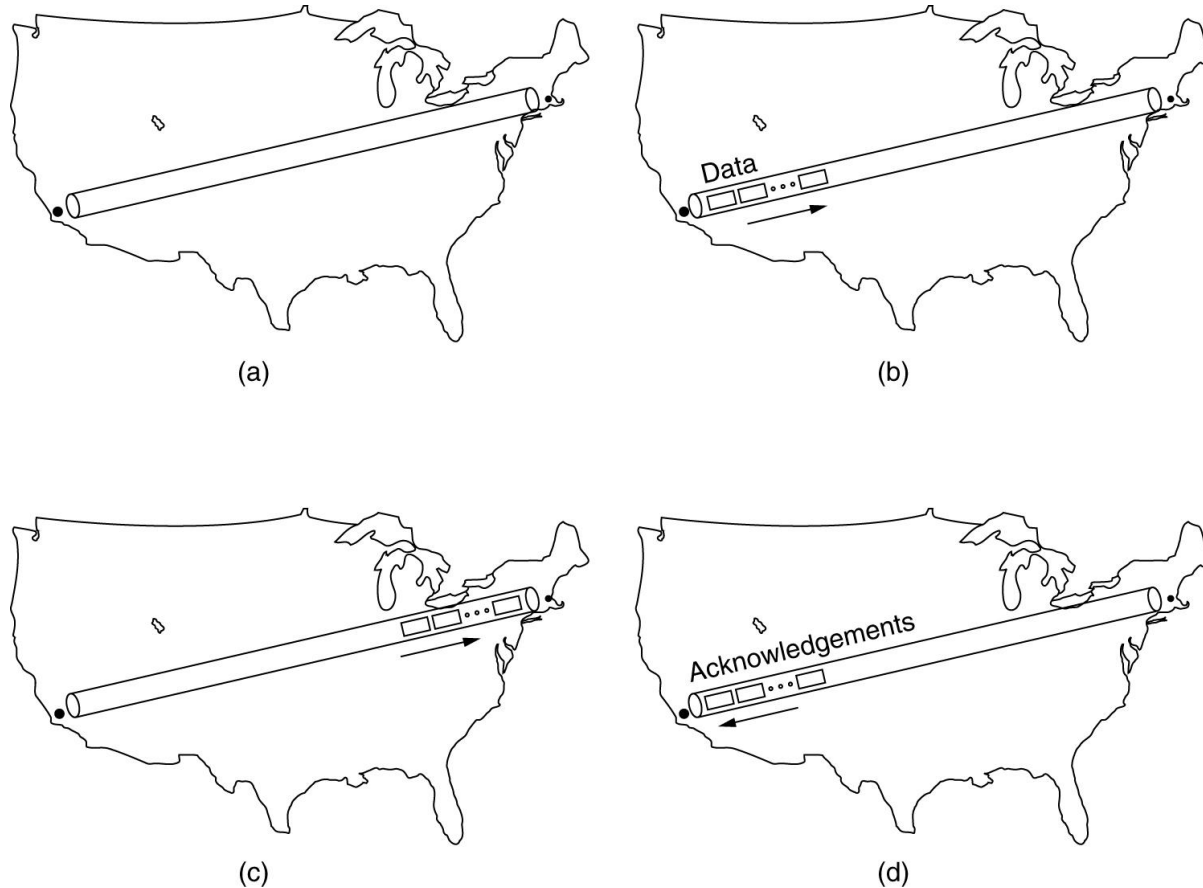
Unexpected loads often interact with protocols to cause performance problems

- Need to find the situations and improve the protocols

Examples:

- Broadcast storm: one broadcast triggers another
- Synchronization: a building of computers all contact the DHCP server together after a power failure
- Tiny packets: some situations can cause TCP to send many small packets instead of few large ones

# Performance Problems in Computer Networks



The state of transmitting one megabit from San Diego to Boston  
(a) At  $t = 0$ , (b) After  $500 \mu\text{sec}$ , (c) After  $20 \text{ msec}$ , (d) after  $40 \text{ msec}$ .

# Measuring Network Performance

Measurement is the key to understanding performance – but has its own pitfalls.

Example pitfalls:

- Caching: fetching Web pages will give surprisingly fast results if they are unexpectedly cached
- Timing: clocks may over/underestimate fast events
- Interference: there may be competing workloads

# Network Performance Measurement

The basic loop for improving network performance.

1. Measure relevant network parameters, performance.
2. Try to understand what is going on.
3. Change one parameter.

# Network Performance Measurement (2)

## Issues in measuring performance

- Sufficient sample size
- Representative samples
- Clock accuracy
- Measuring typical representative load
- Beware of caching
- Understand what you are measuring
- Extrapolate with care

# Host Design for Fast Networks

Poor host software can greatly slow down networks.

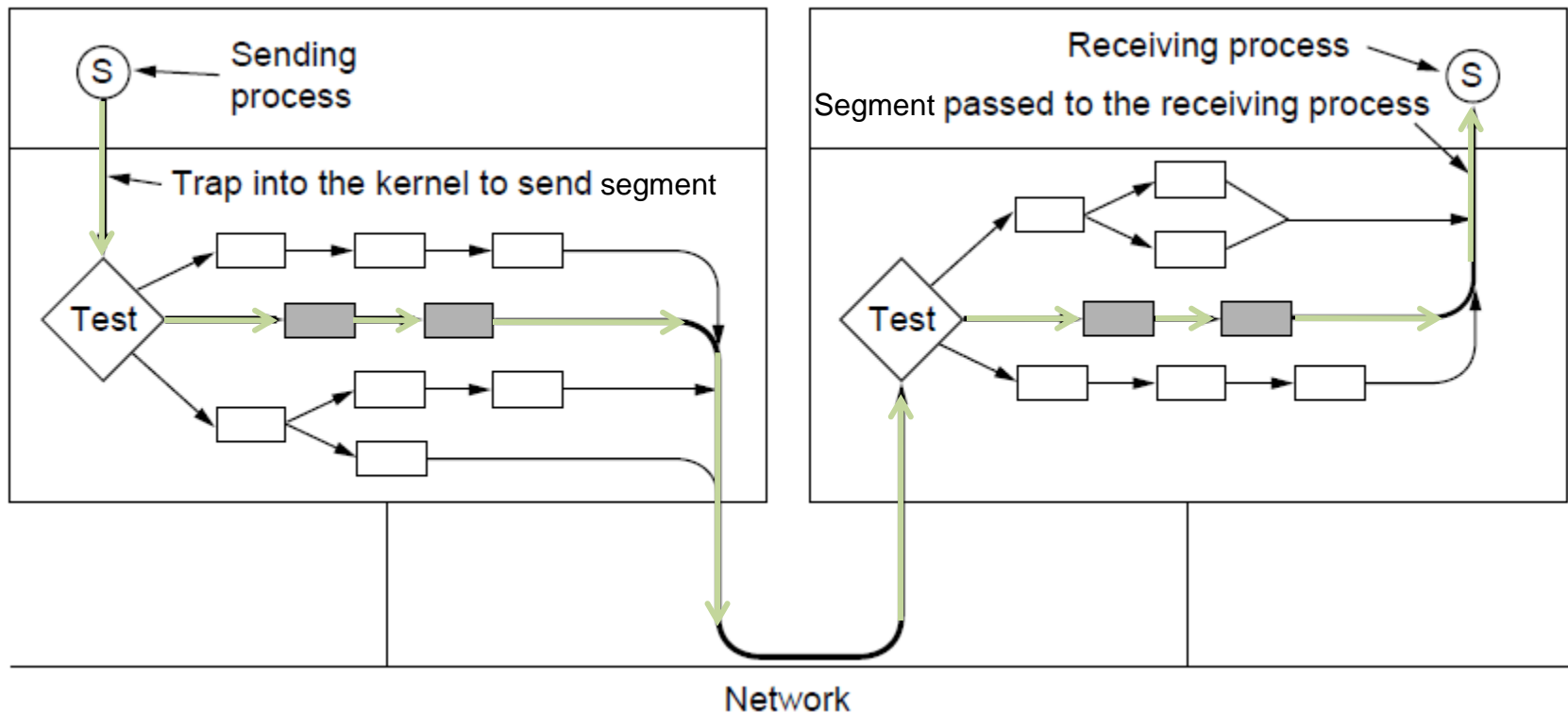
Rules of thumb for fast host software:

- Host speed more important than network speed
- Reduce packet count to reduce overhead
- Minimize data touching
- Minimize context switches
- Avoiding congestion is better than recovering from it
- Avoid timeouts

# Fast Segment Processing (1)

Speed up the common case with a fast path [pink]

- Handles packets with expected header; OK for others to run slowly



# Fast Segment Processing (2)

Header fields are often the same from one packet to the next for a flow; copy/check them to speed up processing

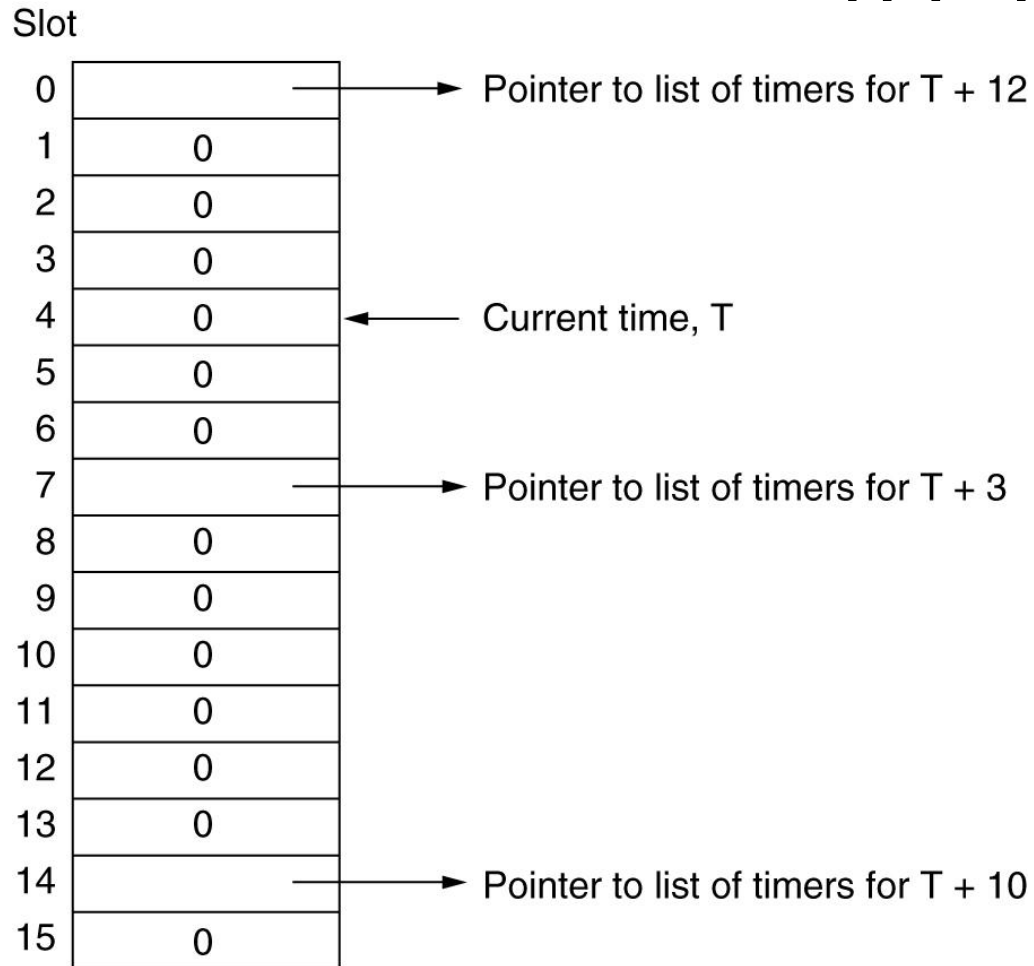
Source port		Destination port	
Sequence number			
Acknowledgement number			
Len	Unused	Window size	
Checksum		Urgent pointer	

TCP header fields that stay the same for a one-way flow (shaded)

VER.	IHL	TOS	Total length	
Identification			Fragment offset	
TTL	Protocol		Header checksum	
Source address				
Destination address				

IP header fields that are often the same for a one-way flow (shaded)

# Fast TPDU Processing (3)



A timing wheel.

# Header Compression

Overhead can be very large for small packets

- 40 bytes of header for RTP/UDP/IP VoIP packet
- Problematic for slow links, especially wireless

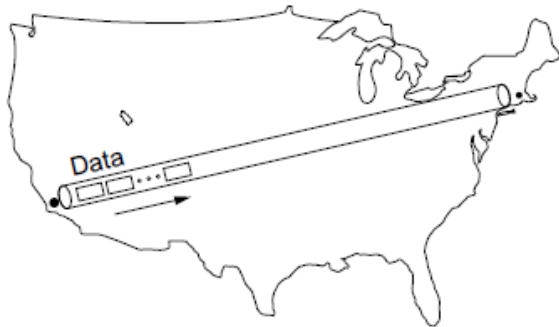
Header compression mitigates this problem

- Runs between Link and Network layer
- Omits fields that don't change or change predictably
  - 40 byte TCP/IP header → 3 bytes of information
- Gives simple high-layer headers and efficient links

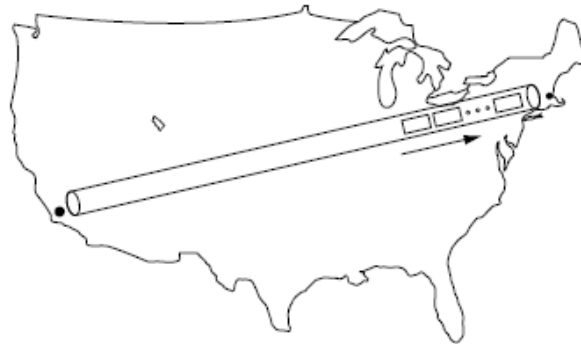
# Protocols for “Long Fat” Networks (1)

Networks with high bandwidth (“Fat”) and high delay (“Long”) can store much information inside the network

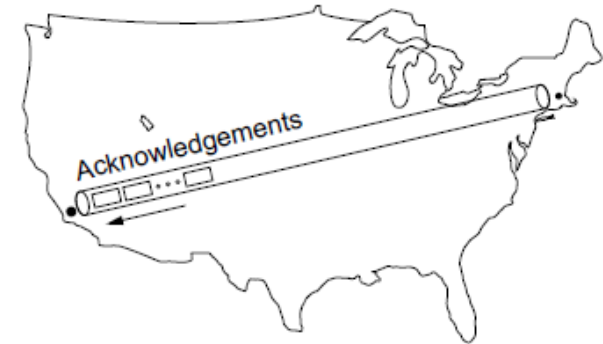
- Requires protocols with ample buffering and few RTTs.



Starting to send 1 Mbit  
San Diego → Boston



20ms after start

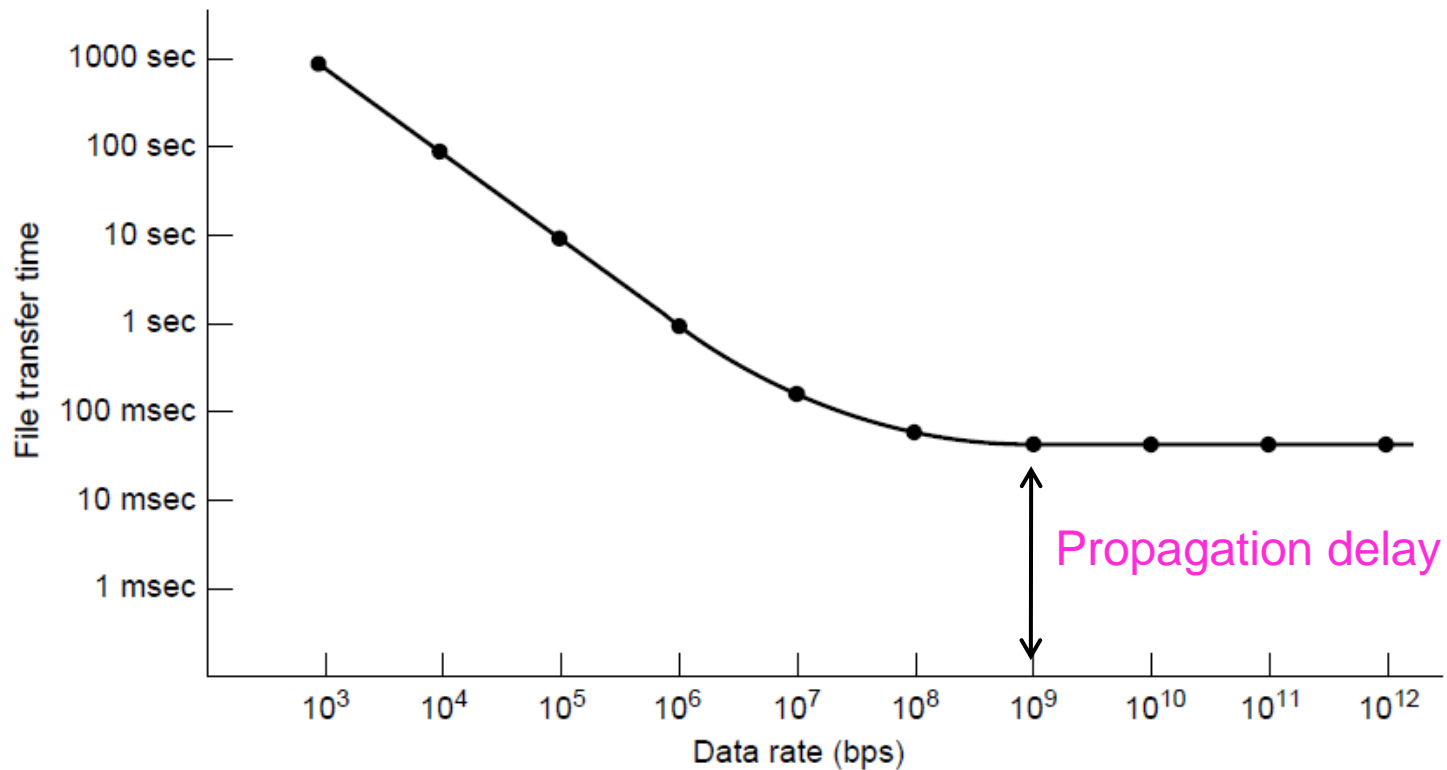


40ms after start

# Protocols for “Long Fat” Networks (2)

You can buy more bandwidth but not lower delay

- Need to shift ends (e.g., into cloud) to lower further



Minimum time to send and ACK a 1-Mbit file over a 4000-km line

# Delay Tolerant Networking

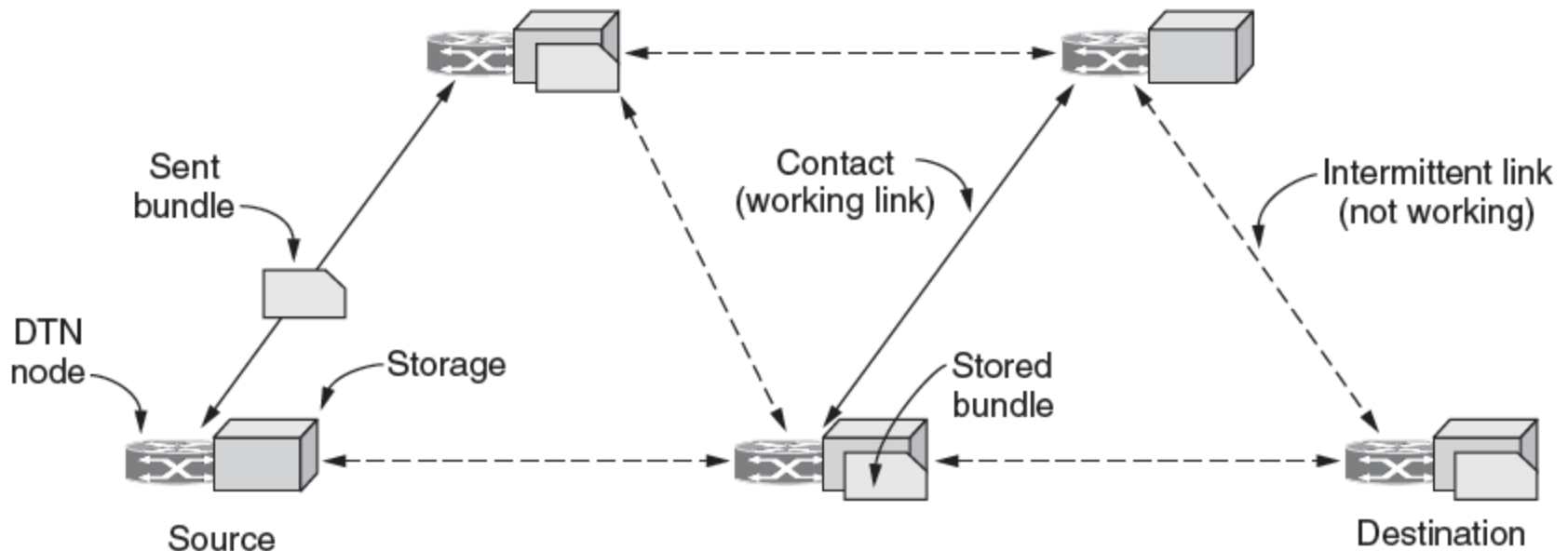
DTNs (Delay Tolerant Networks) store messages inside the network until they can be delivered

- DTN Architecture »
- Bundle Protocol »

# DTN Architecture (1)

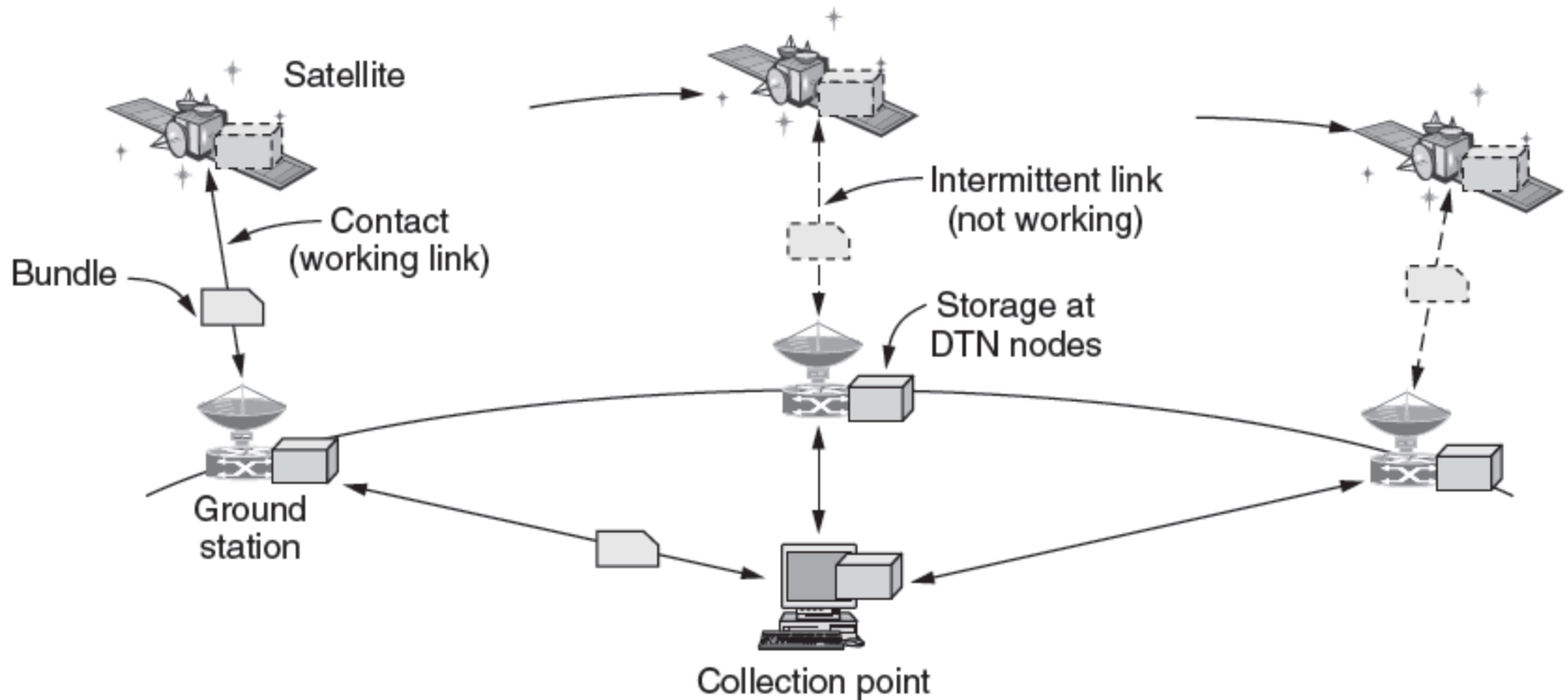
Messages called bundles are stored at DTN nodes while waiting for an intermittent link to become a contact

- Bundles might wait hours, not milliseconds in routers



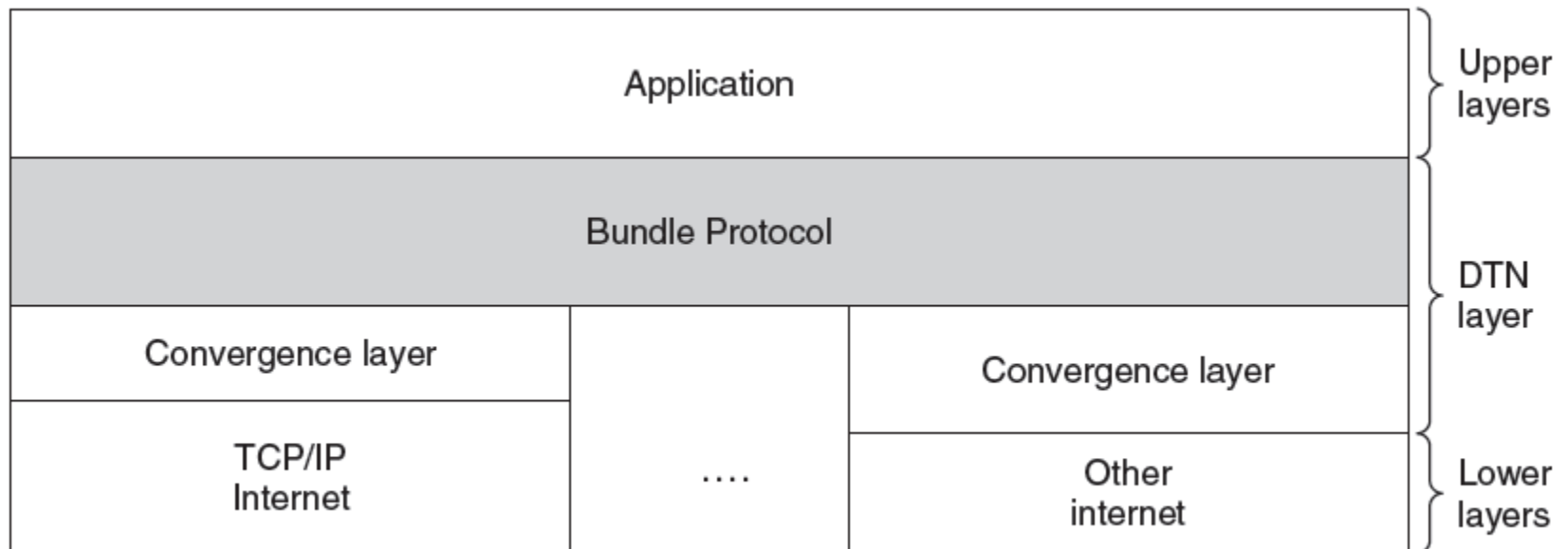
# DTN Architecture (2)

Example DTN connecting a satellite to a collection point



# Bundle Protocol (1)

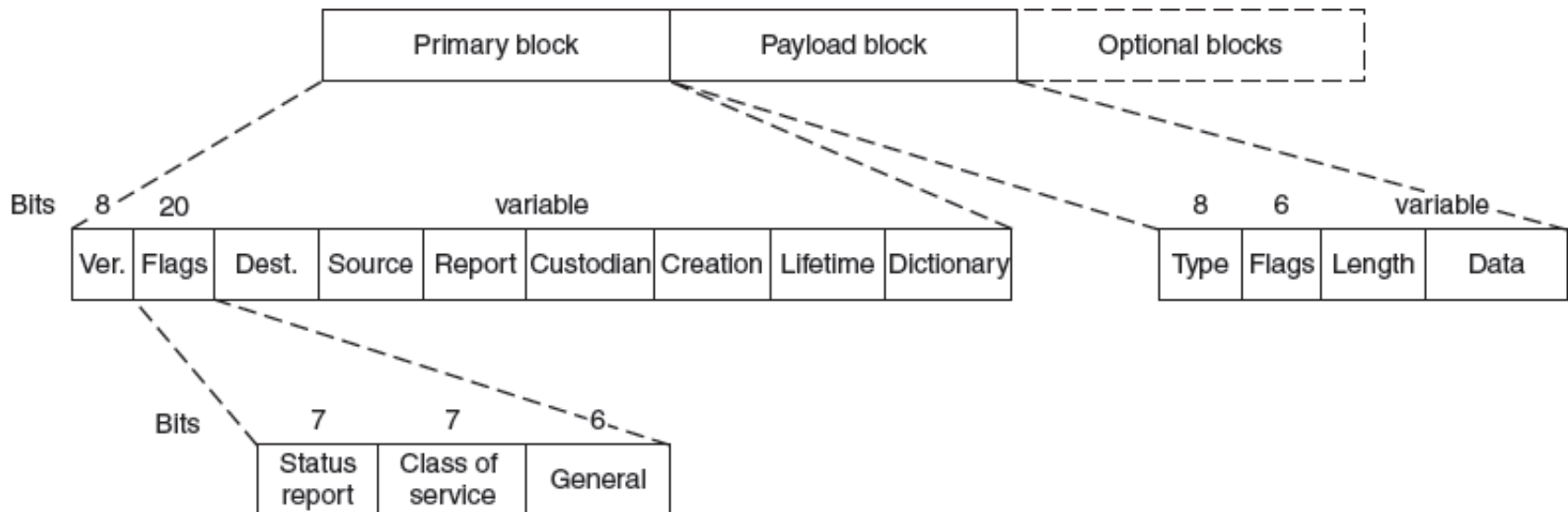
The Bundle protocol uses TCP or other transports and provides a DTN service to applications



# Bundle Protocol (2)

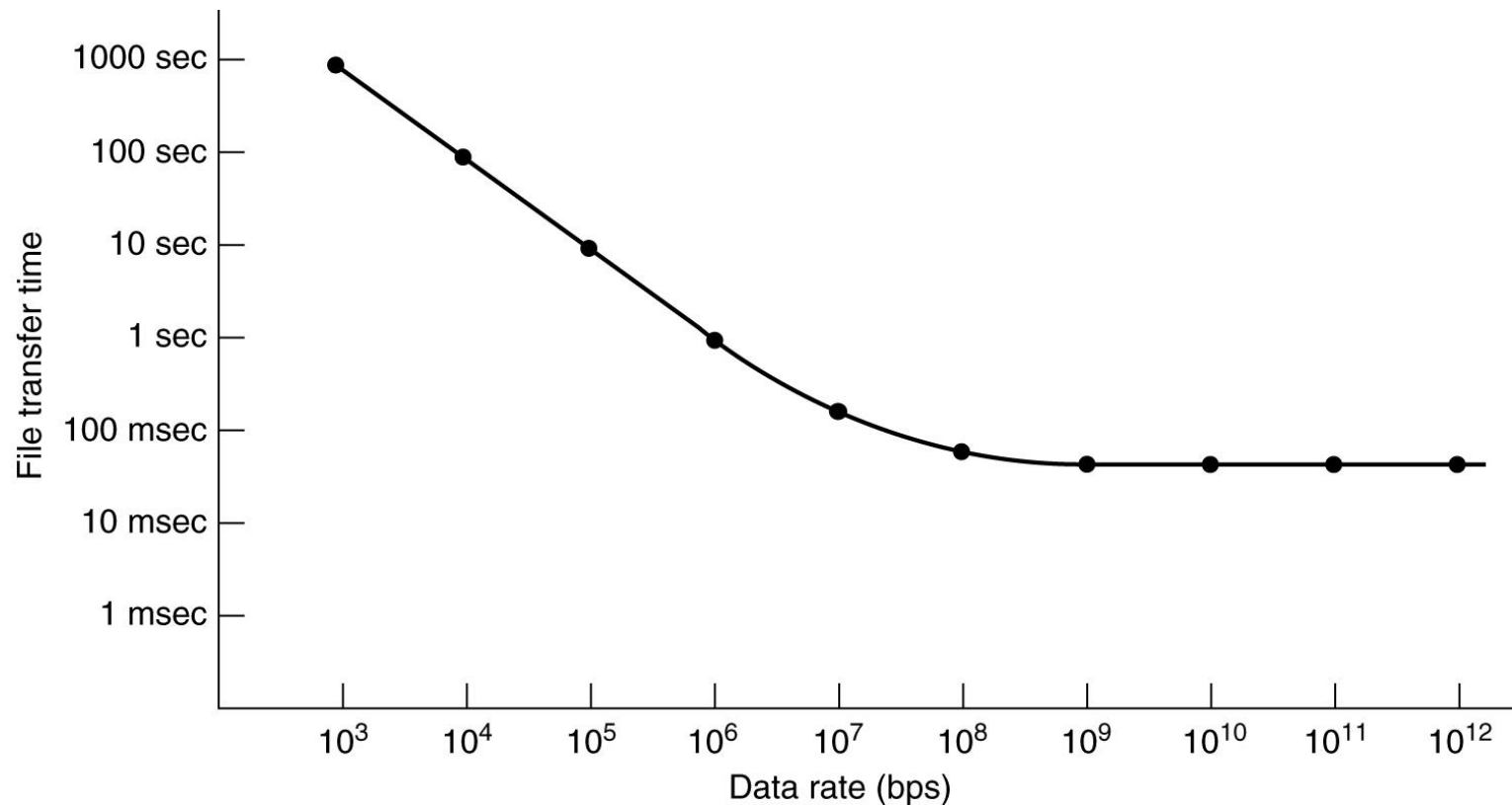
## Features of the bundle message format:

- Dest./source add high-level addresses (not port/IP)
- Custody transfer shifts delivery responsibility
- Dictionary provides compression for efficiency



# Protocols for Gigabit Networks

Time to transfer and acknowledge a 1-megabit file over a 4000-km line.



# A Simple Transport Protocol

- The Example Service Primitives
- The Example Transport Entity
- The Example as a Finite State Machine

Similar to TCP but simpler

# Service Primitives

- Connect
  - Parameters – local and remote TSAPs
  - Caller is blocked
  - If connection succeeds the caller is unblocked and transmission starts
- Listen – specifies a TSAP to listen to
- Disconnect
- Send
- Receive
- \*\* Library procedures

# Service Primitives

- Connnum=LISTEN(local)
- Connnum=Connect(local,remote)
- Status = Send(Connnum, buffer,bytes)
  - No Connection, illegal buffer address, negative count
- Status = Receive(Connnum, buffer, bytes)
- Status = Disconnect(Connnum)

# The Transport Entity

- Use connection oriented, reliable network service
- Transport Entity is part of the user process
- Network Layer interface
  - To\_net and from\_net
  - Parameters –
    - Connection Identifier
    - Q bit – control message
    - M bit – more data from this message to follow
    - Packet Type
    - Pointer to data

# The Example Transport Entity

The network layer packets used in our example.

<b>Network packet</b>	<b>Meaning</b>
CALL REQUEST	Sent to establish a connection
CALL ACCEPTED	Response to CALL REQUEST
CLEAR REQUEST	Sent to release a connection
CLEAR CONFIRMATION	Response to CLEAR REQUEST
DATA	Used to transport data
CREDIT	Control packet for managing the window

# The Example Transport Entity (2)

Each connection is in one of seven states:

1. Idle – Connection not established yet.
2. Waiting – CONNECT has been executed, CALL REQUEST sent.
3. Queued – A CALL REQUEST has arrived; no LISTEN yet.
4. Established – The connection has been established.
5. Sending – The user is waiting for permission to send a packet.
6. Receiving – A RECEIVE has been done.
7. DISCONNECTING – a DISCONNECT has been done locally.

# State Transitions

- A primitive is executed
- A packet arrives
- A timer expires

# The Example Transport Entity (3)

```
#define MAX_CONN 32                /* max number of simultaneous connections */
#define MAX_MSG_SIZE 8192          /* largest message in bytes */
#define MAX_PKT_SIZE 512          /* largest packet in bytes */
#define TIMEOUT 20
#define CRED 1
#define OK 0

#define ERR_FULL -1
#define ERR_REJECT -2
#define ERR_CLOSED -3
#define LOW_ERR -3

typedef int transport_address;
typedef enum {CALL_REQ,CALL_ACC,CLEAR_REQ,CLEAR_CONF,DATA_PKT,CREDIT} pkt_type;
typedef enum {IDLE,WAITING,QUEUED,ESTABLISHED,SENDING,RECEIVING,DISCONN} cstate;

/* Global variables. */
transport_address listen_address; /* local address being listened to */
int listen_conn;                 /* connection identifier for listen */
unsigned char data[MAX_PKT_SIZE]; /* scratch area for packet data */

struct conn {
    transport_address local_address, remote_address;
    cstate state; /* state of this connection */
    unsigned char *user_buf_addr; /* pointer to receive buffer */
    int byte_count; /* send/receive count */
    int clr_req_received; /* set when CLEAR_REQ packet received */
    int timer; /* used to time out CALL_REQ packets */
    int credits; /* number of messages that may be sent */
} conn[MAX_CONN + 1]; /* slot 0 is not used */
```

# The Example Transport Entity (4)

```
void sleep(void);                /* prototypes */
void wakeup(void);
void to_net(int cid, int q, int m, pkt_type pt, unsigned char *p, int bytes);
void from_net(int *cid, int *q, int *m, pkt_type *pt, unsigned char *p, int *bytes);

int listen(transport_address t)
{ /* User wants to listen for a connection. See if CALL_REQ has already arrived. */
  int i, found = 0;

  for (i = 1; i <= MAX_CONN; i++)          /* search the table for CALL_REQ */
    if (conn[i].state == QUEUED && conn[i].local_address == t) {
      found = i;
      break;
    }

  if (found == 0) {
    /* No CALL_REQ is waiting. Go to sleep until arrival or timeout. */
    listen_address = t; sleep(); i = listen_conn ;
  }
  conn[i].state = ESTABLISHED;             /* connection is ESTABLISHED */
  conn[i].timer = 0;                       /* timer is not used */
}
```

# The Example Transport Entity (5)

```
listen_conn = 0; /* 0 is assumed to be an invalid address */
to_net(i, 0, 0, CALL_ACC, data, 0); /* tell net to accept connection */
return(i); /* return connection identifier */
}
```

```
int connect(transport_address l, transport_address r)
{ /* User wants to connect to a remote process; send CALL_REQ packet. */
  int i;
  struct conn *cptr;

  data[0] = r; data[1] = l; /* CALL_REQ packet needs these */
  i = MAX_CONN; /* search table backward */
  while (conn[i].state != IDLE && i > 1) i = i - 1;
  if (conn[i].state == IDLE) {
    /* Make a table entry that CALL_REQ has been sent. */
    cptr = &conn[i];
    cptr->local_address = l; cptr->remote_address = r;
    cptr->state = WAITING; cptr->clr_req_received = 0;
    cptr->credits = 0; cptr->timer = 0;
    to_net(i, 0, 0, CALL_REQ, data, 2);
    sleep(); /* wait for CALL_ACC or CLEAR_REQ */
    if (cptr->state == ESTABLISHED) return(i);
    if (cptr->clr_req_received) {
      /* Other side refused call. */
      cptr->state = IDLE; /* back to IDLE state */
      to_net(i, 0, 0, CLEAR_CONF, data, 0);
      return(ERR_REJECT);
    }
  }
  else return(ERR_FULL); /* reject CONNECT: no table space */
}
```

# The Example Transport Entity (6)

```
int send(int cid, unsigned char bufptr[], int bytes)
{ /* User wants to send a message. */
  int i, count, m;
  struct conn *cptr = &conn[cid];

  /* Enter SENDING state. */
  cptr->state = SENDING;
  cptr->byte_count = 0; /* # bytes sent so far this message */
  if (cptr->clr_req_received == 0 && cptr->credits == 0) sleep();
  if (cptr->clr_req_received == 0) {
    /* Credit available; split message into packets if need be. */
    do {
      if (bytes - cptr->byte_count > MAX_PKT_SIZE) { /* multipacket message */
        count = MAX_PKT_SIZE; m = 1; /* more packets later */
      } else { /* single packet message */
        count = bytes - cptr->byte_count; m = 0; /* last pkt of this message */
      }
      for (i = 0; i < count; i++) data[i] = bufptr[cptr->byte_count + i];
      to_net(cid, 0, m, DATA_PKT, data, count); /* send 1 packet */
      cptr->byte_count = cptr->byte_count + count; /* increment bytes sent so far */
    } while (cptr->byte_count < bytes); /* loop until whole message sent */
  }
}
```

# The Example Transport Entity (7)

```
    cptr->credits -- ;                /* each message uses up one credit */
    cptr->state = ESTABLISHED;
    return(OK);
} else {
    cptr->state = ESTABLISHED;
    return(ERR_CLOSED);              /* send failed: peer wants to disconnect */
}
}
```

```
int receive(int cid, unsigned char bufptr[], int *bytes)
{ /* User is prepared to receive a message. */
    struct conn *cptr = &conn[cid];

    if (cptr->clr_req_received == 0) {
        /* Connection still established; try to receive. */
        cptr->state = RECEIVING;
        cptr->user_buf_addr = bufptr;
        cptr->byte_count = 0;
        data[0] = CRED;
        data[1] = 1;
        to_net(cid, 1, 0, CREDIT, data, 2);    /* send credit */
        sleep();                               /* block awaiting data */
        *bytes = cptr->byte_count;
    }
    cptr->state = ESTABLISHED;
    return(cptr->clr_req_received ? ERR_CLOSED : OK);
}
```

# The Example Transport Entity (8)

```
int disconnect(int cid)
{ /* User wants to release a connection. */
  struct conn *cptr = &conn[cid];

  if (cptr->clr_req_received) { /* other side initiated termination */
    cptr->state = IDLE; /* connection is now released */
    to_net(cid, 0, 0, CLEAR_CONF, data, 0);
  } else { /* we initiated termination */
    cptr->state = DISCONN; /* not released until other side agrees */
    to_net(cid, 0, 0, CLEAR_REQ, data, 0);
  }
  return(OK);
}

void packet_arrival(void)
{ /* A packet has arrived, get and process it. */
  int cid; /* connection on which packet arrived */
  int count, i, q, m;
  pkt_type ptype; /* CALL_REQ, CALL_ACC, CLEAR_REQ, CLEAR_CONF, DATA_PKT, CREDIT */
  unsigned char data[MAX_PKT_SIZE]; /* data portion of the incoming packet */
  struct conn *cptr;

  from_net(&cid, &q, &m, &ptype, data, &count); /* go get it */
  cptr = &conn[cid];
```

# The Example Transport Entity (9)

```
switch (ptype) {
  case CALL_REQ: /* remote user wants to establish connection */
    cptr->local_address = data[0]; cptr->remote_address = data[1];
    if (cptr->local_address == listen_address) {
      listen_conn = cid; cptr->state = ESTABLISHED; wakeup();
    } else {
      cptr->state = QUEUED; cptr->timer = TIMEOUT;
    }
    cptr->clr_req_received = 0; cptr->credits = 0;
    break;

  case CALL_ACC: /* remote user has accepted our CALL_REQ */
    cptr->state = ESTABLISHED;
    wakeup();
    break;

  case CLEAR_REQ: /* remote user wants to disconnect or reject call */
    cptr->clr_req_received = 1;
    if (cptr->state == DISCONN) cptr->state = IDLE; /* clear collision */
    if (cptr->state == WAITING || cptr->state == RECEIVING || cptr->state == SENDING) wakeup();
    break;

  case CLEAR_CONF: /* remote user agrees to disconnect */
    cptr->state = IDLE;
    break;

  case CREDIT: /* remote user is waiting for data */
    cptr->credits += data[1];
    if (cptr->state == SENDING) wakeup();
    break;

  case DATA_PKT: /* remote user has sent data */
    for (i = 0; i < count; i++) cptr->user_buf_addr[cptr->byte_count + i] = data[i];
    cptr->byte_count += count;
    if (m == 0) wakeup();
}
}
```

# The Example Transport Entity (10)

```
}  
void clock(void)  
{ /* The clock has ticked, check for timeouts of queued connect requests. */  
  int i;  
  struct conn *cptr;  
  for (i = 1; i <= MAX_CONN; i++) {  
    cptr = &conn[i];  
    if (cptr->timer > 0) { /* timer was running */  
      cptr->timer--;  
      if (cptr->timer == 0) { /* timer has now expired */  
        cptr->state = IDLE;  
        to_net(i, 0, 0, CLEAR_REQ, data, 0);  
      }  
    }  
  }  
}
```

# The Example as a Finite State Machine

The example protocol as a finite state machine. Each entry has an optional predicate, an optional action, and the new state. The tilde indicates that no major action is taken. An overbar above a predicate indicate the negation of the predicate. Blank entries correspond to impossible or invalid events.

		State						
		Idle	Waiting	Queued	Established	Sending	Receiving	Dis- connecting
Primitives	LISTEN	P1: ~/Idle P2: A1/Estab P2: A2/Idle		~/Estab				
	CONNECT	P1: ~/Idle P1: A3/Wait						
	DISCONNECT				P4: A5/Idle P4: A6/Disc			
	SEND				P5: A7/Estab P5: A8/Send			
	RECEIVE				A9/Receiving			
Incoming packets	Call_req	P3: A1/Estab P3: A4/Queue'd						
	Call_acc		~/Estab					
	Clear_req		~/Idle		A10/Estab	A10/Estab	A10/Estab	~/Idle
	Clear_conf							~/Idle
	DataPkt						A12/Estab	
Clock	Credit				A11/Estab	A7/Estab		
	Timeout			~/Idle				

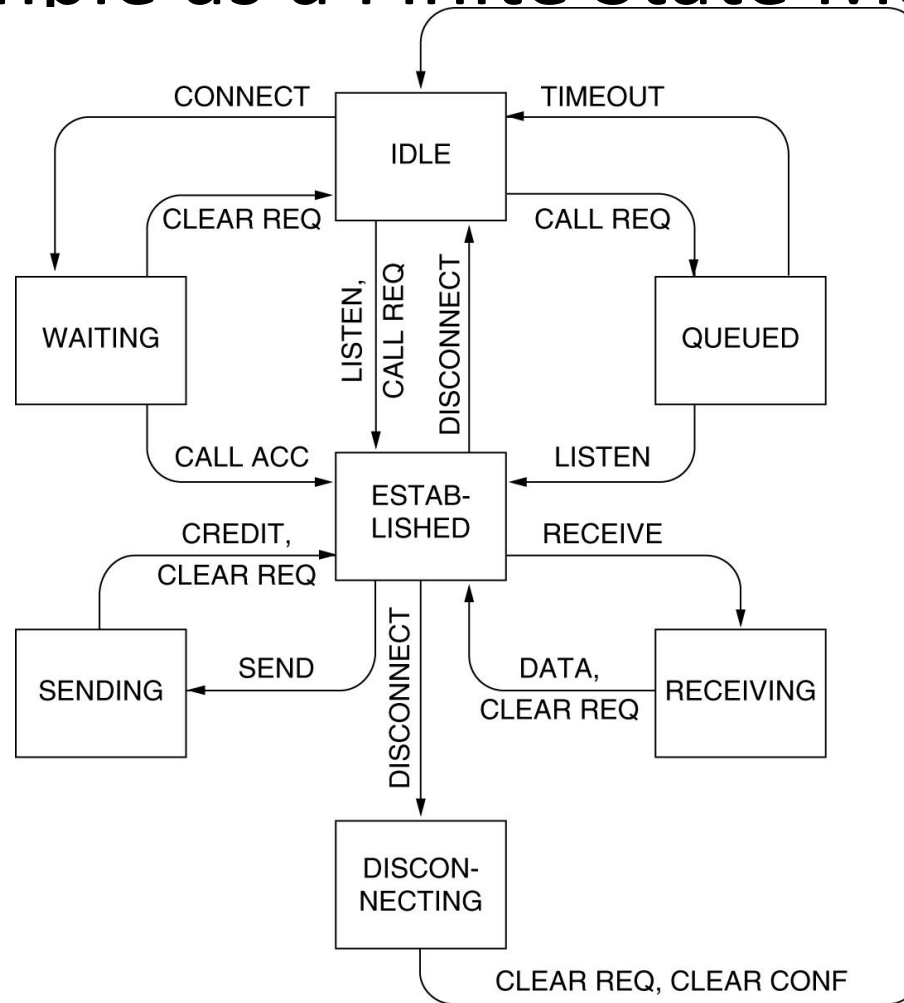
## Predicates

P1: Connection table full  
 P2: Call\_req pending  
 P3: LISTEN pending  
 P4: Clear\_req pending  
 P5: Credit available

## Actions

A1: Send Call\_acc      A7: Send message  
 A2: Wait for Call\_req    A8: Wait for credit  
 A3: Send Call\_req      A9: Send credit  
 A4: Start timer        A10: Set Clr\_req\_received flag  
 A5: Send Clear\_conf    A11: Record credit  
 A6: Send Clear\_req     A12: Accept message

# The Example as a Finite State Machine (2)



The example protocol in graphical form. Transitions that leave the connection state unchanged have been omitted for simplicity.