

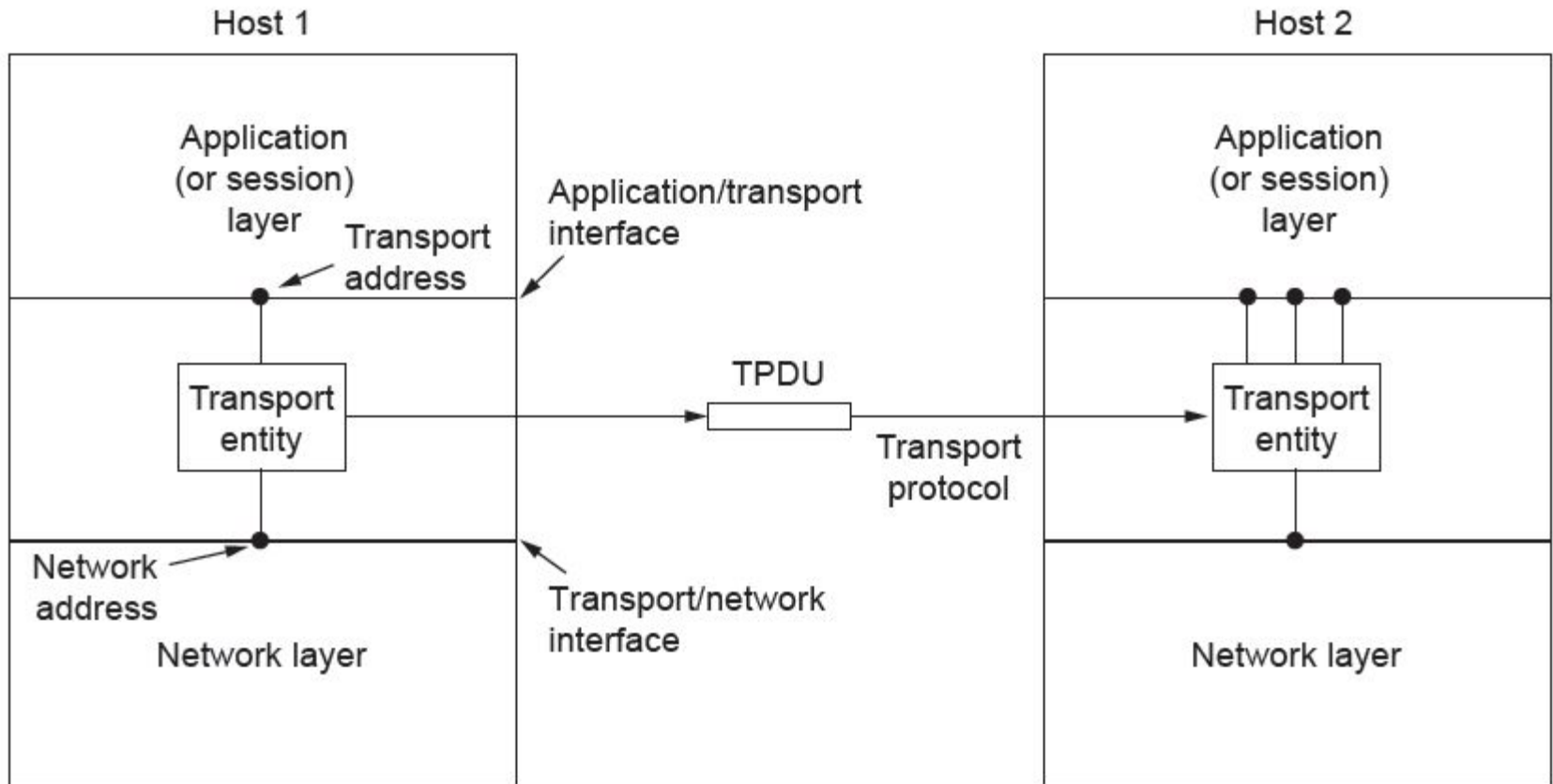
The Transport Layer

Chapter 6

Transport Service

- Upper Layer Services
 - Provides Services to applications
 - Takes services from Network Layer
 - Hides network imperfections
- Transport Service Primitives
 - Adds reliability to the data delivery process
- Berkeley Sockets
- Example of Socket Programming:
Internet File Server

Services Provided to the Upper Layers



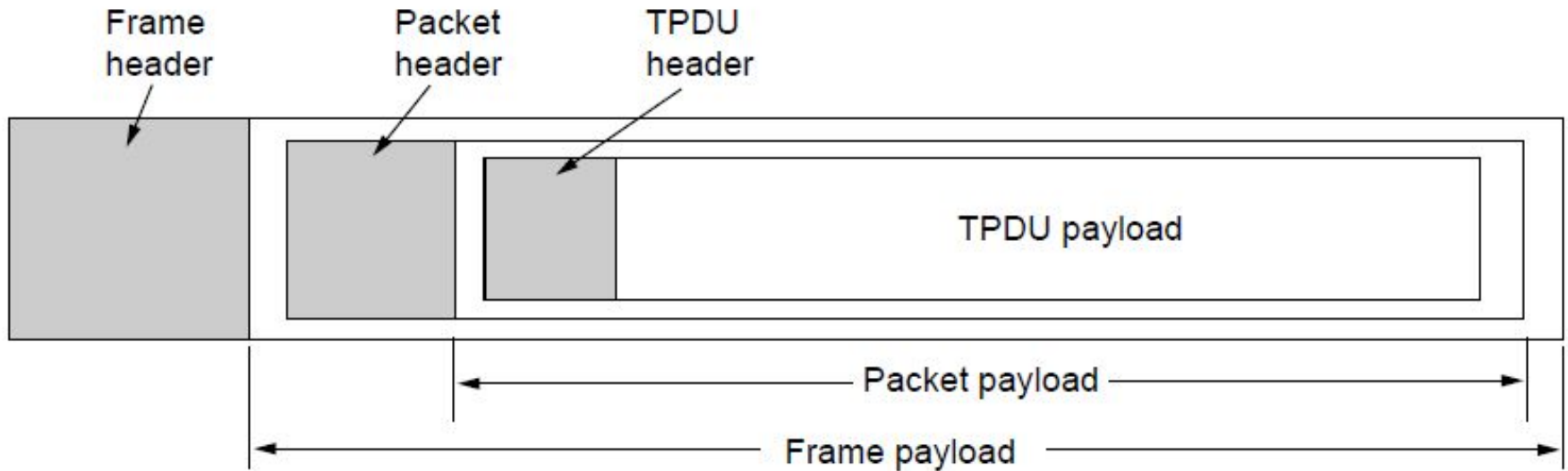
The network, transport, and application layers

Transport Service Primitives (1)

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

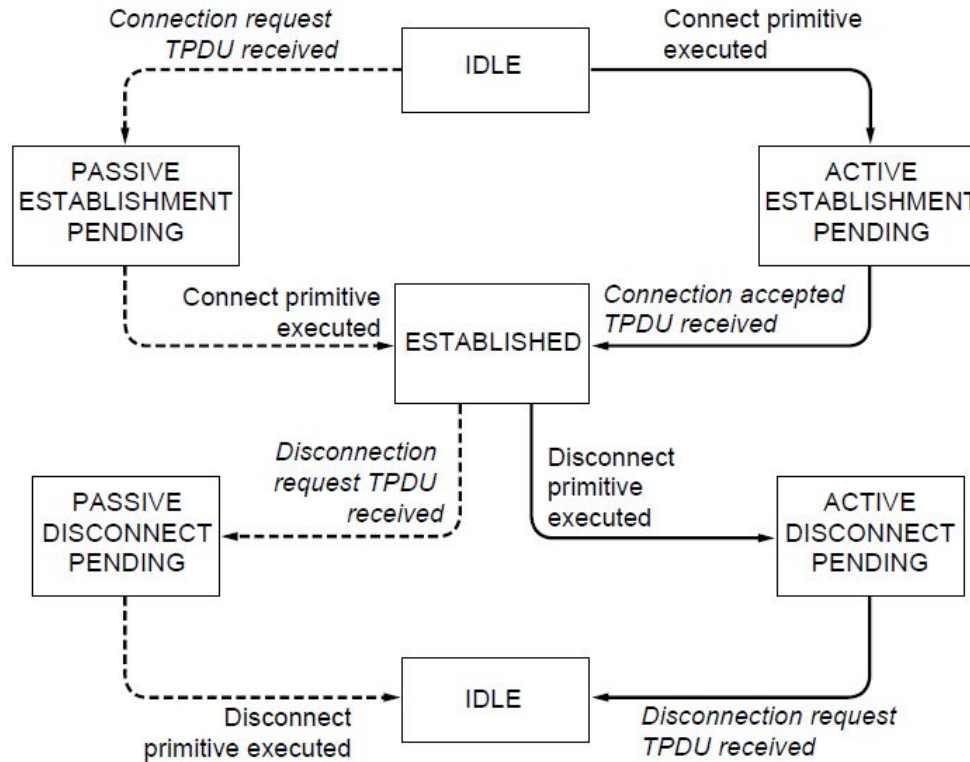
The primitives for a simple transport service

Transport Service Primitives (2)



Nesting of TPDUs, packets, and frames.

Berkeley Sockets (1)



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

Berkeley Sockets (2)

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

The socket primitives for TCP

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];
    struct hostent *h;
    struct sockaddr_in channel;

    /* buffer for incoming file */
    /* info about server */
    /* 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>
#include <sys/fcntl.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>

#define SERVER_PORT 12345
#define BUF_SIZE 4096
#define QUEUE_SIZE 10

int main(int argc, char *argv[])
{
    int s, b, l, fd, sa, bytes, on = 1;
    char buf[BUF_SIZE];
    struct sockaddr_in channel;

    . . .
```

/* This is the server code */

/* arbitrary, but client & server must agree */
/* block transfer size */

/* buffer for outgoing file */
/* 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 */
}
}
```

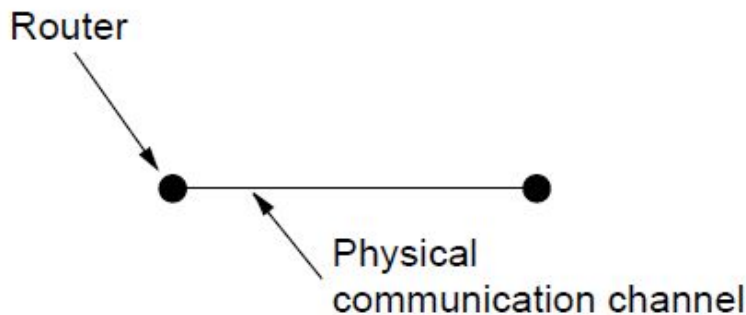
Server code

Elements of Transport Protocols (1)

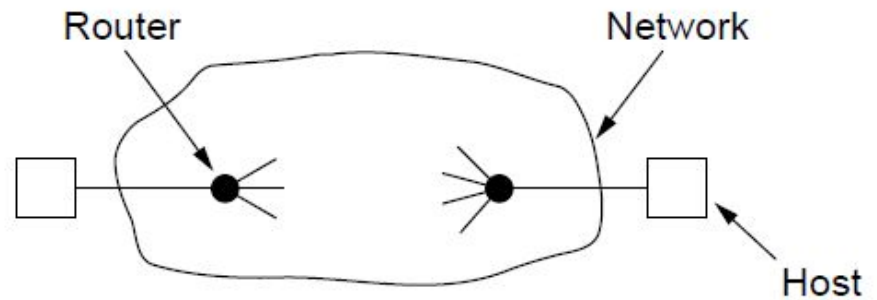
- Addressing
- Connection establishment
- Connection release
- Error control and flow control
- Multiplexing
- Crash recovery

Elements of Transport Protocols (2)

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



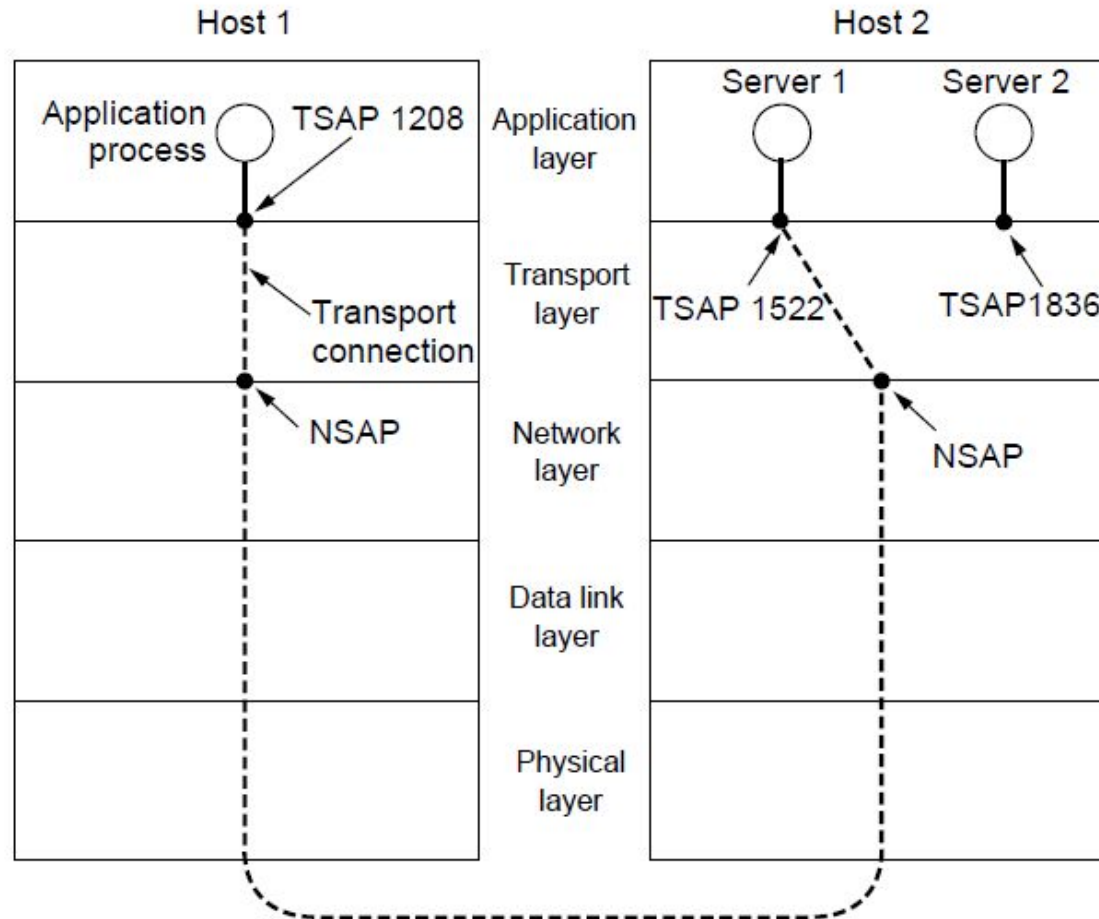
(a)



(b)

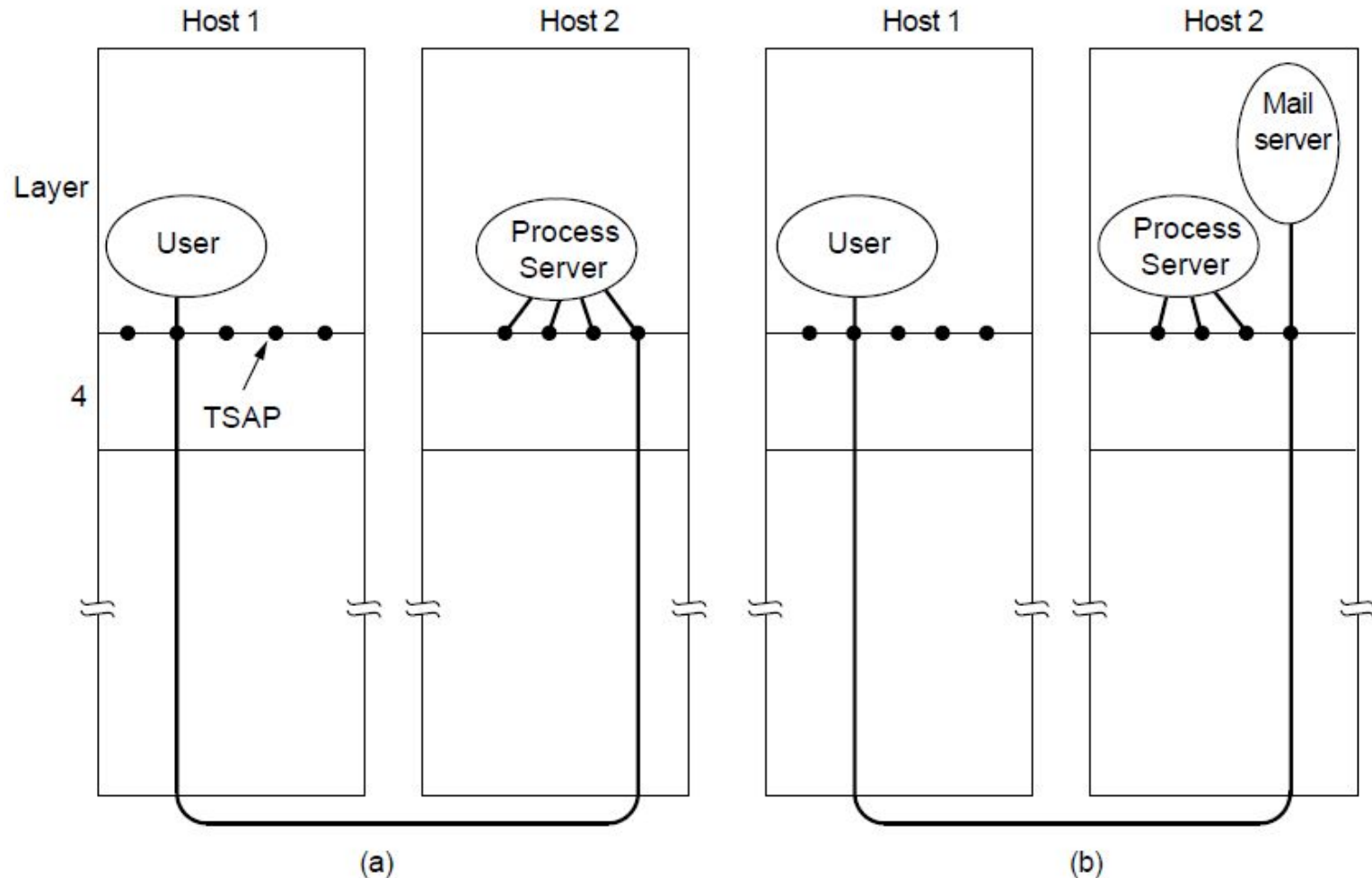
- Data Link Layer-router need not specify which router it want to talk to, Transport Layer- Required
- Connection Establishment- More Complex in Transport Layer
- Transport Layer needs to take care of storage capacity in the subnet
- Buffering and Flow Control – different for large and dynamically varying connections in the transport layer

Addressing (1)



TSAPs, NSAPs, and transport connections

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)

Difficulties : Network can

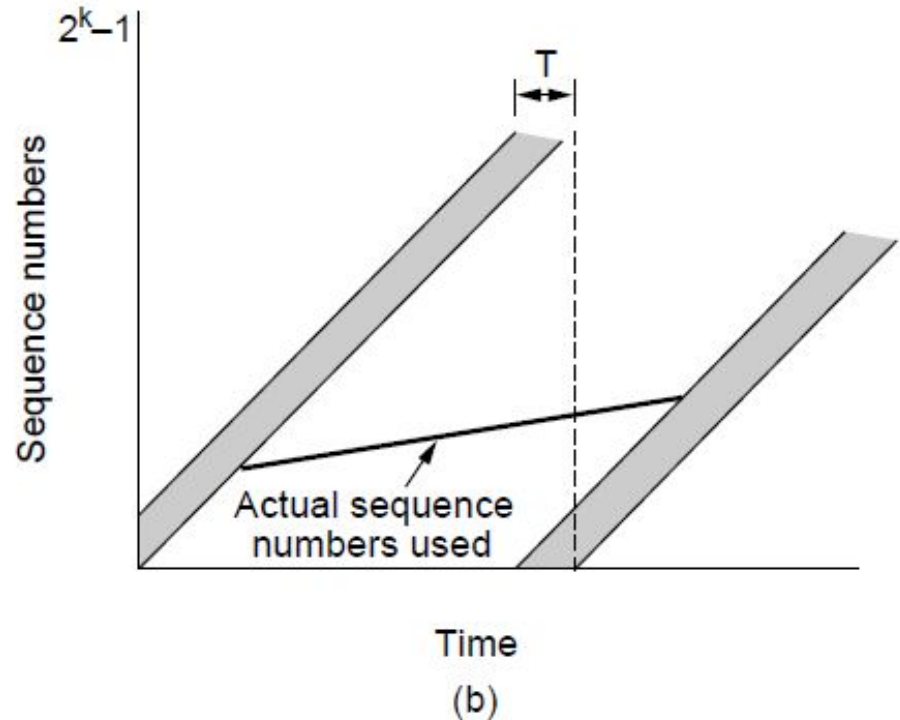
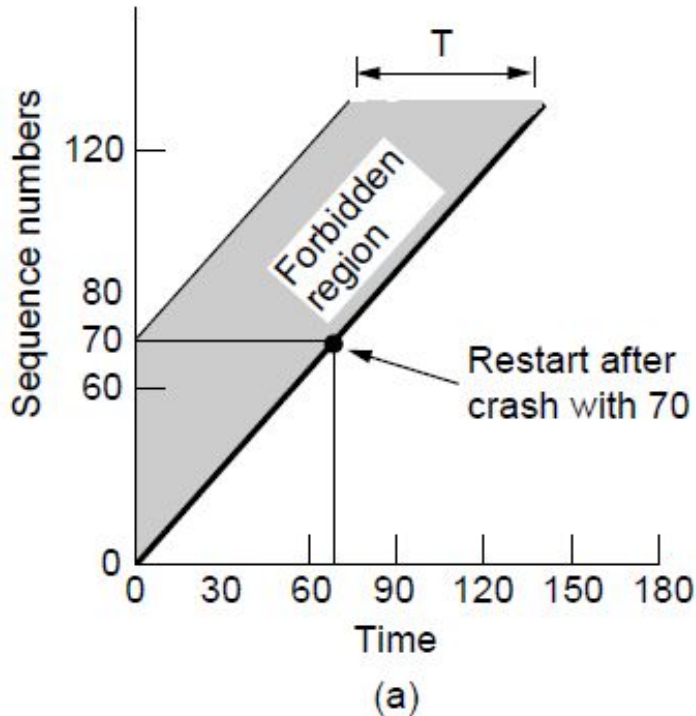
- Lose
- Store and
- Duplicate Packets including delayed duplicates

Possible Solutions

- Use throwaway Transport Addresses
- Using Connection Identifiers

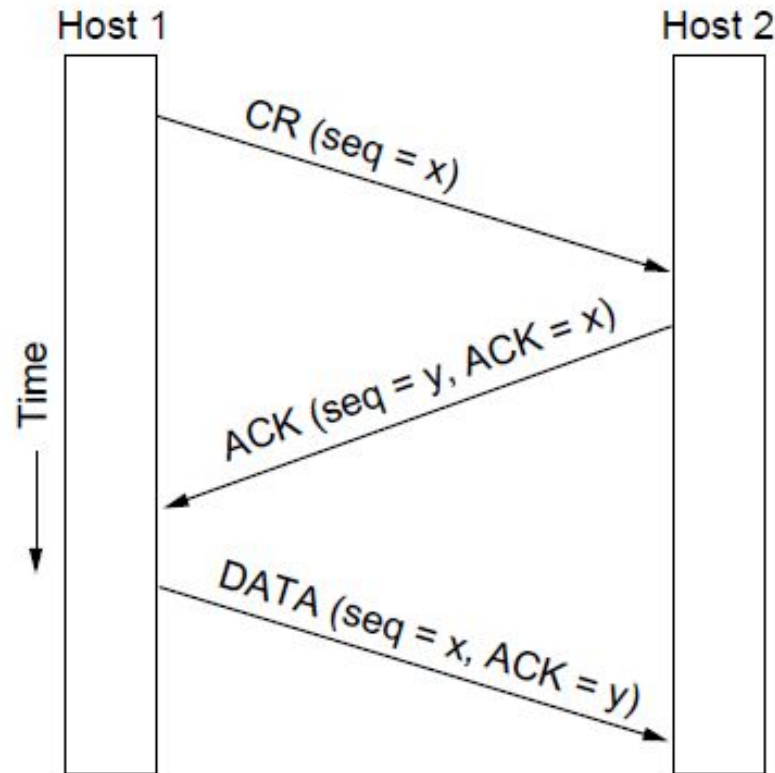
- Solution lies in not allowing packets to live forever within the subnet
- Techniques for restricting packet lifetime
- Restricted network design.
 - Prevent packets from looping combined with bounding congestion delay
- Putting a hop counter in each packet.
- Time-stamping each packet.

Connection Establishment (2)



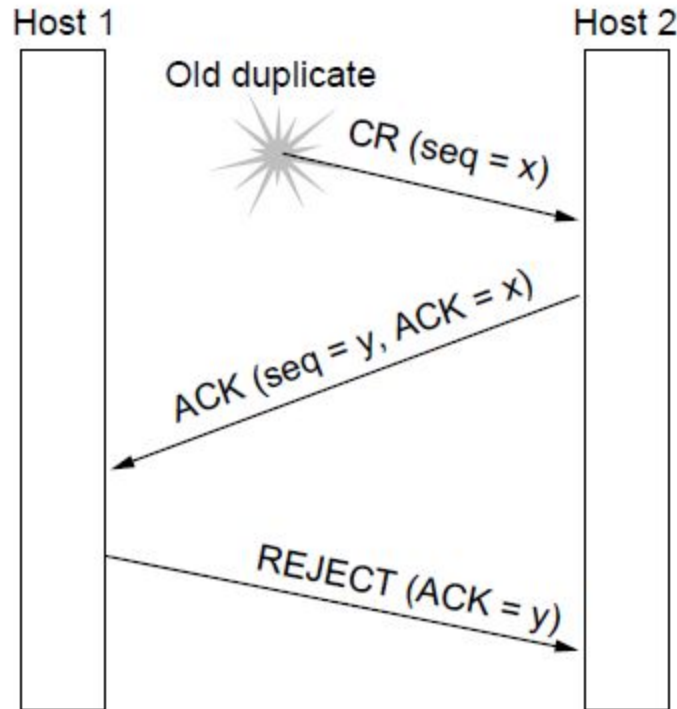
- (a) TPDUs may not enter the forbidden region.
- (b) The resynchronization problem.

Connection Establishment (3)



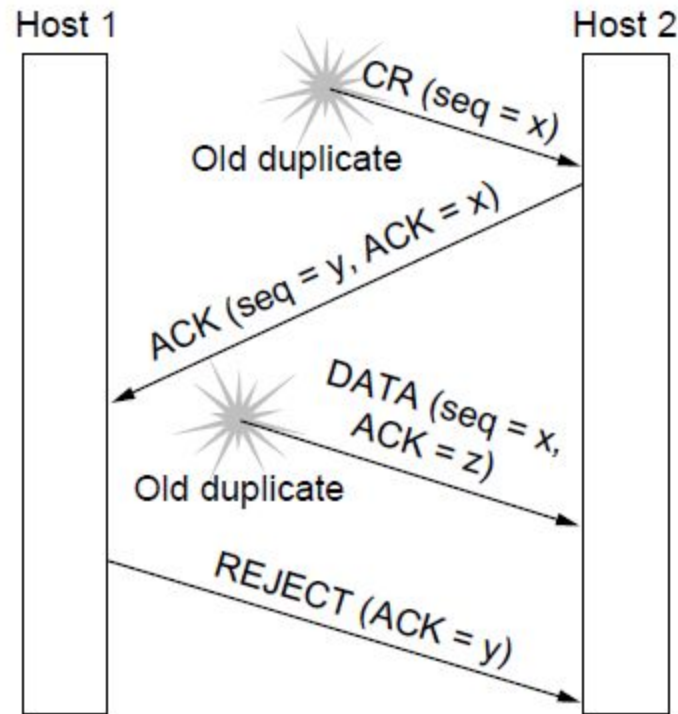
Three protocol scenarios for establishing a connection using a three-way handshake. CR denotes CONNECTION REQUEST. Normal operation.

Connection Establishment (4)



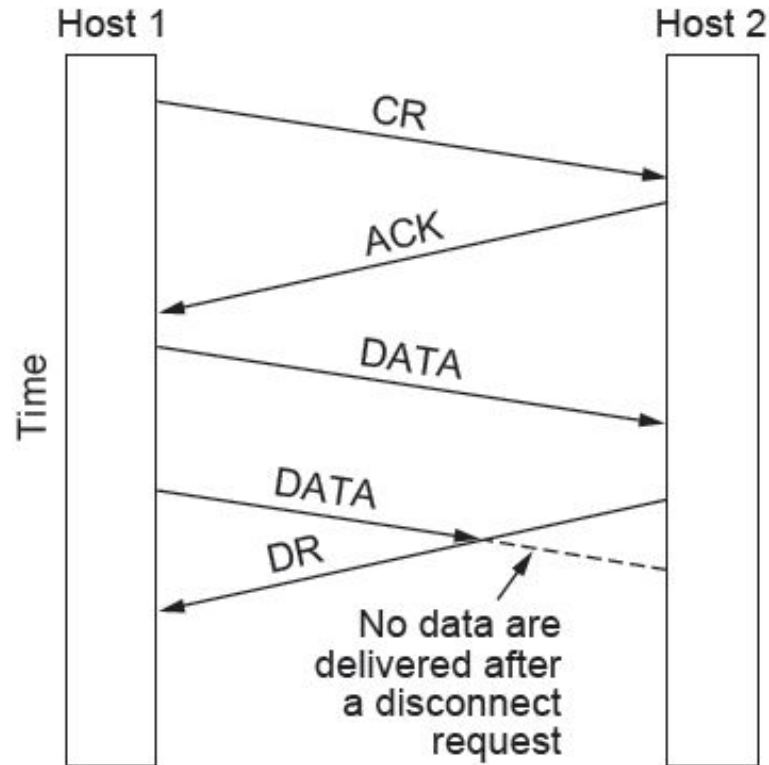
Three protocol scenarios for establishing a connection using a three-way handshake. CR denotes CONNECTION REQUEST. Old duplicate CONNECTION REQUEST appearing out of nowhere.

Connection Establishment (5)



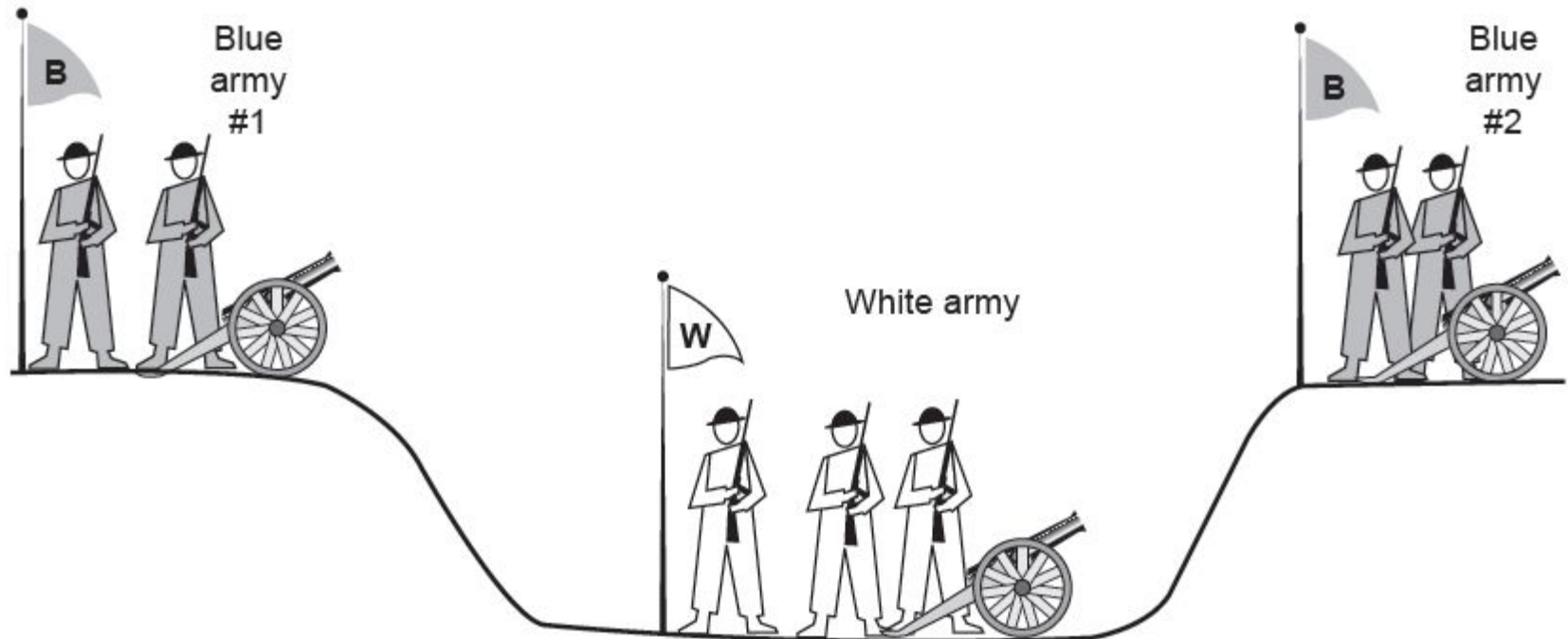
Three protocol scenarios for establishing a connection using a three-way handshake. CR denotes CONNECTION REQUEST. Duplicate CONNECTION REQUEST and duplicate ACK

Connection Release (1)



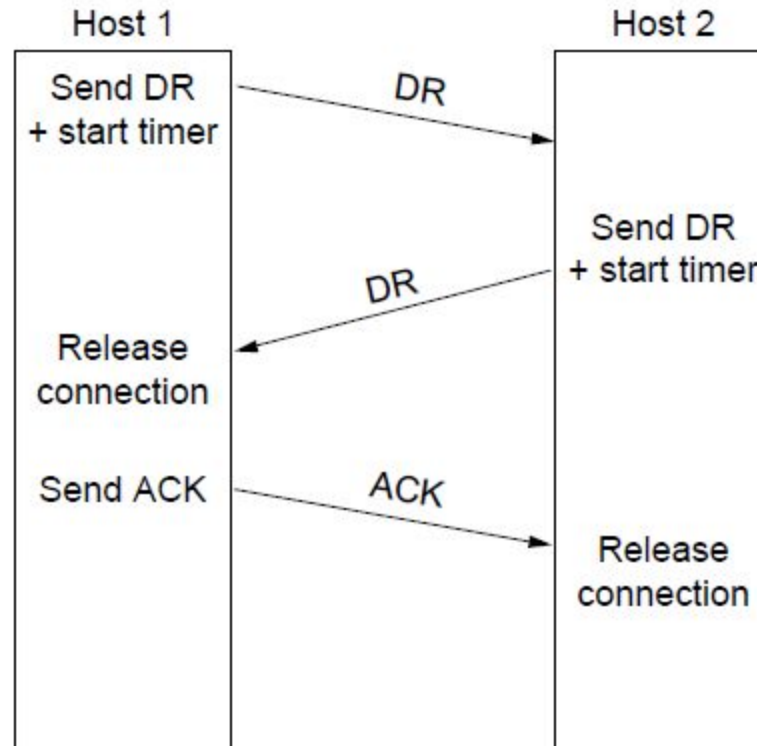
Abrupt disconnection with loss of data

Connection Release (2)



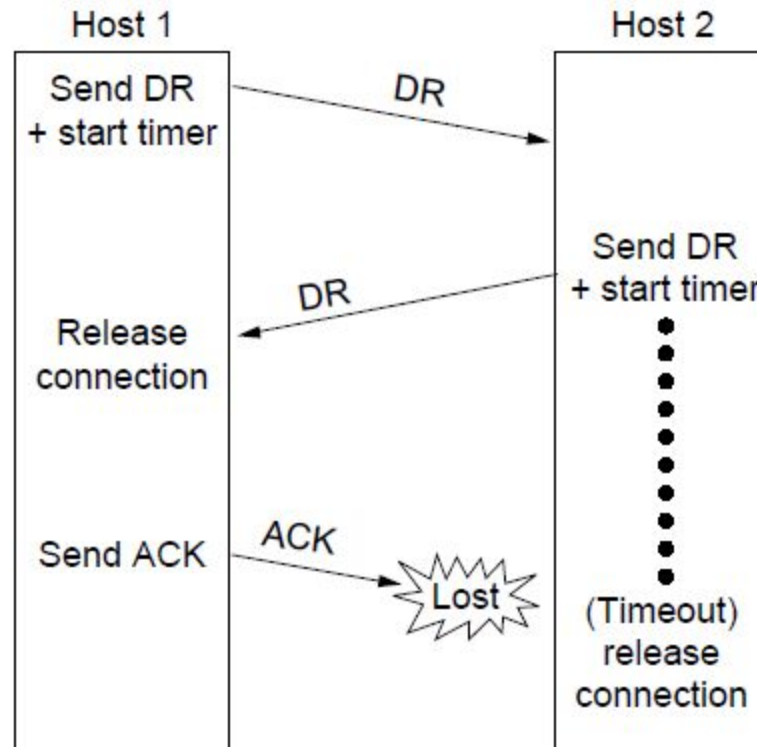
The two-army problem

Connection Release (3)



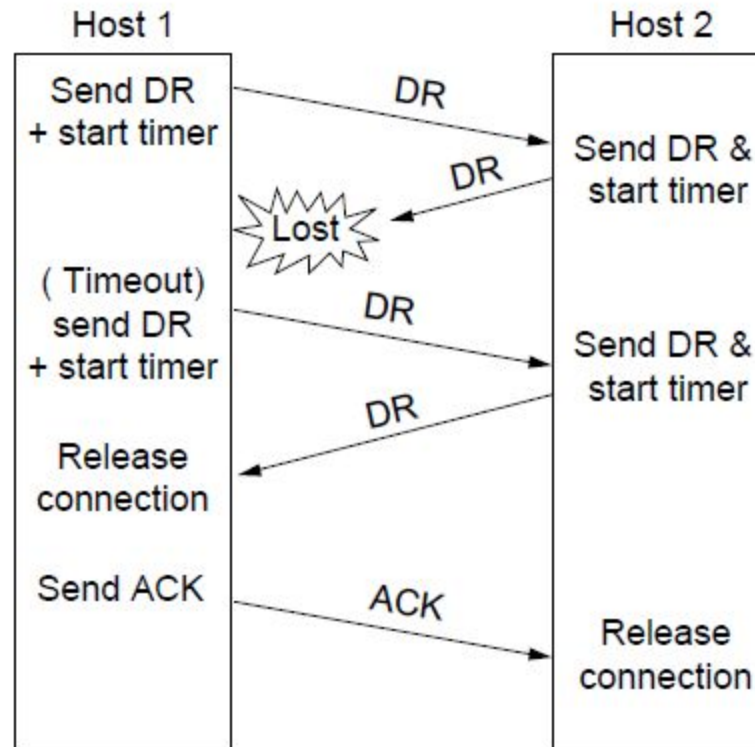
Four protocol scenarios for releasing a connection.
(a) Normal case of three-way handshake

Connection Release (4)



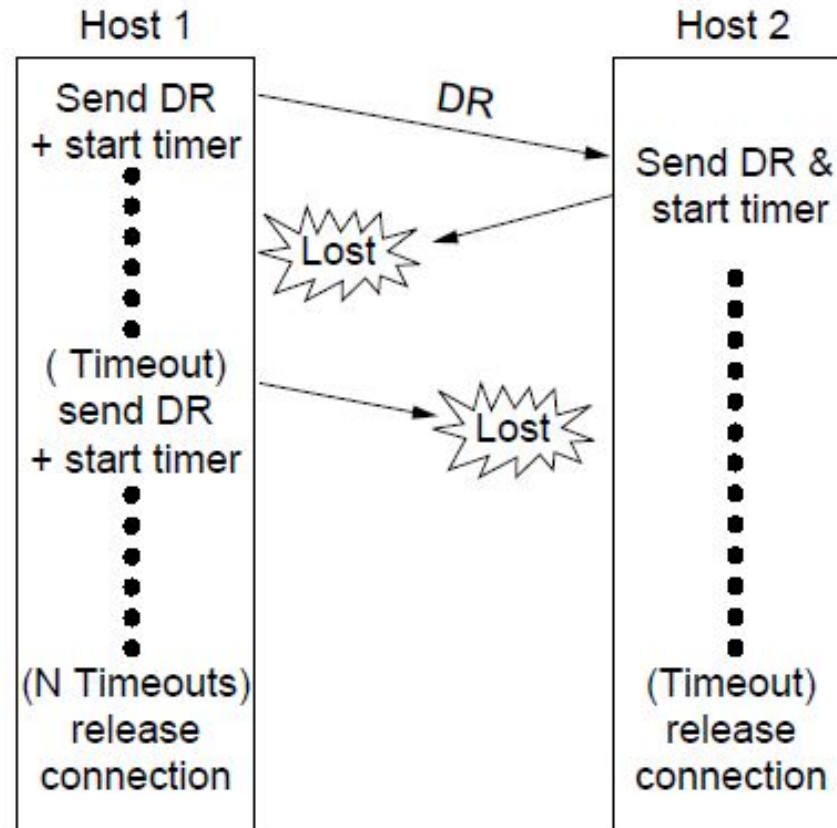
Four protocol scenarios for releasing a connection.
(b) Final ACK lost.

Connection Release (5)



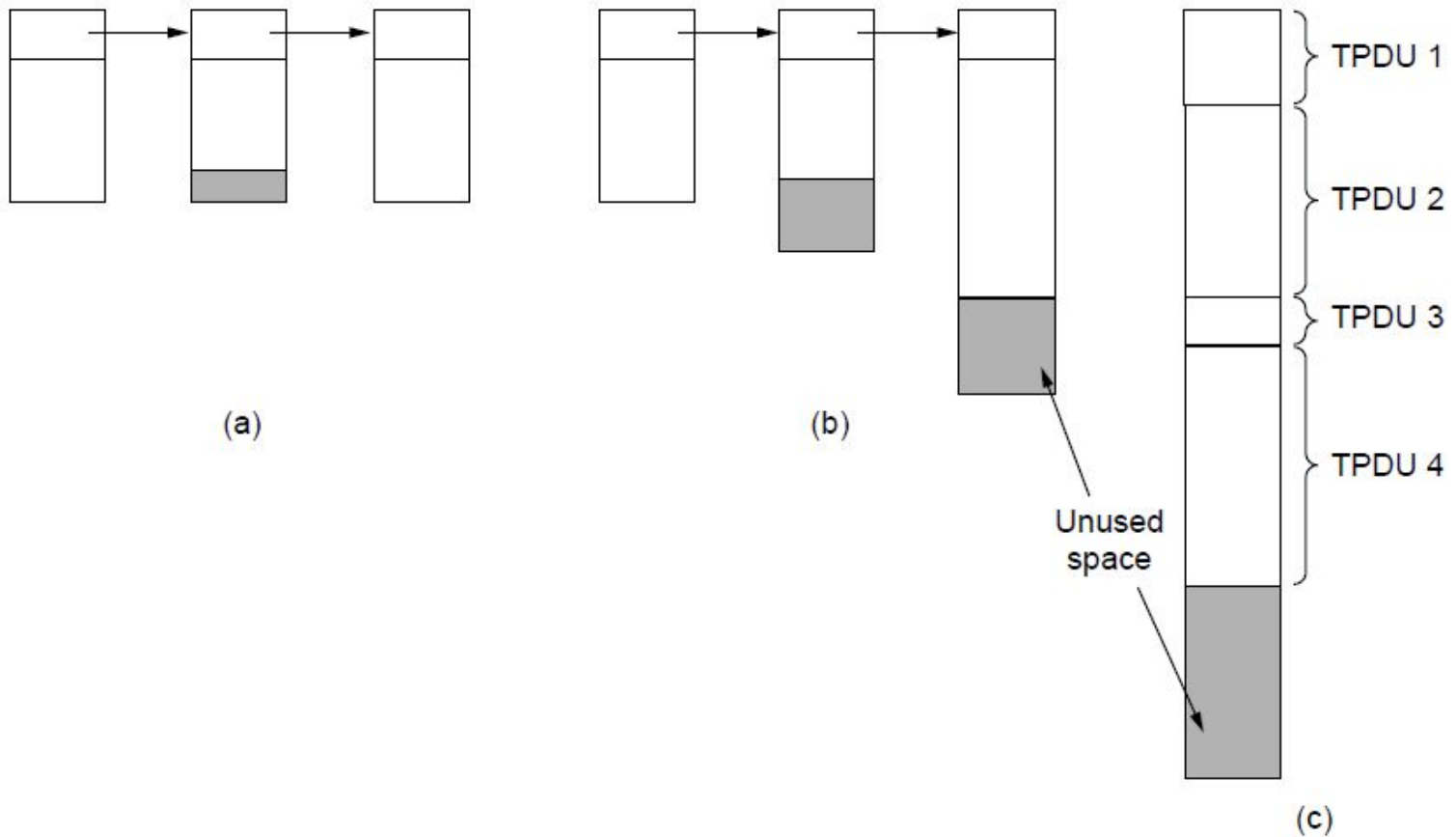
Four protocol scenarios for releasing a connection.
(c) Response lost

Connection Release (6)



Four protocol scenarios for releasing a connection.
(d) Response lost and subsequent DRs lost.

Error Control and Flow Control (1)



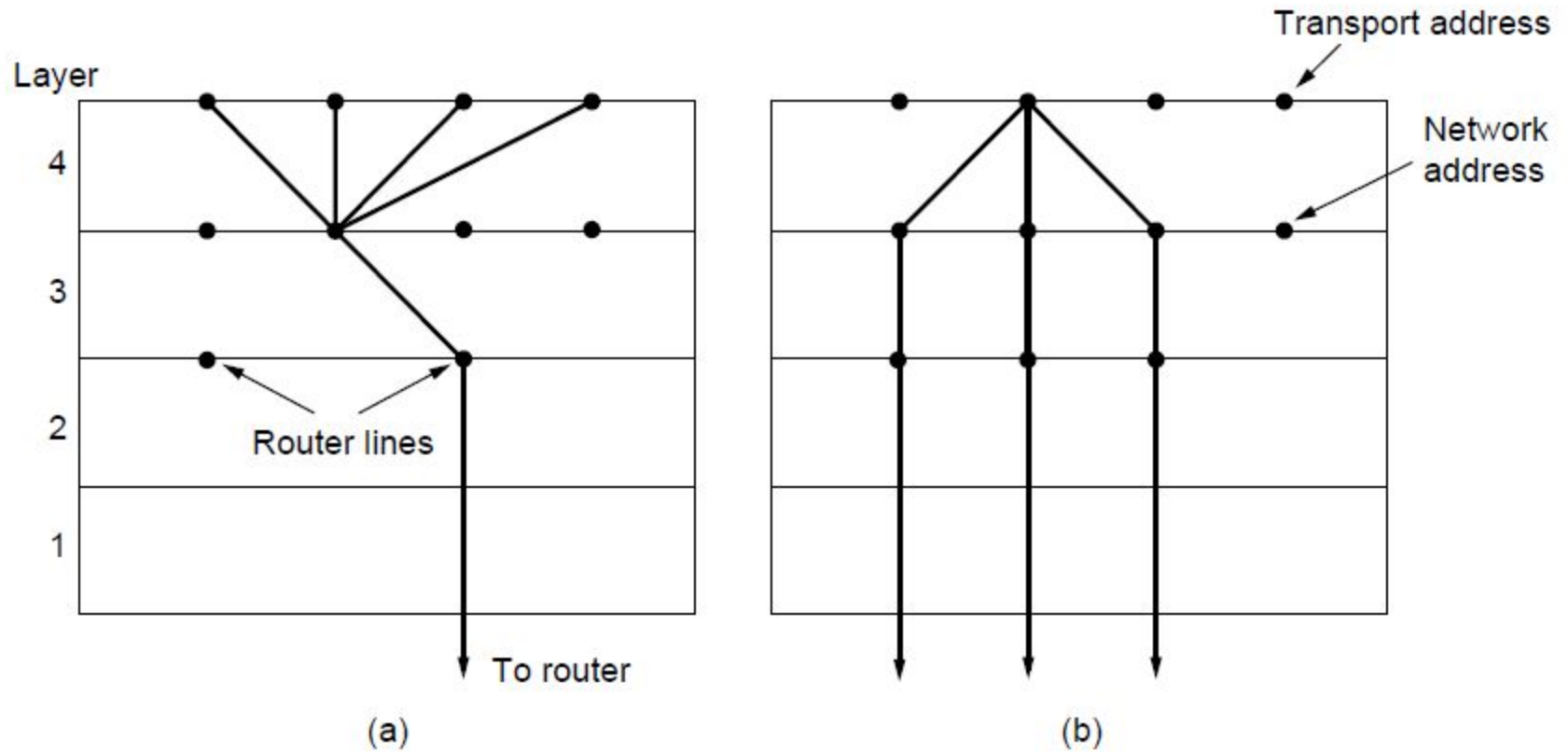
(a) Chained fixed-size buffers. **(b)** Chained variable-sized buffers. **(c)** One large circular buffer per connection.

Error Control and Flow Control (2)

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

Dynamic buffer allocation. The arrows show the direction of transmission. An ellipsis (...) indicates a lost TPDU

Multiplexing



(a) Multiplexing. (b) Inverse multiplexing.

Crash Recovery

Strategy used by sending host	Strategy used by receiving 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

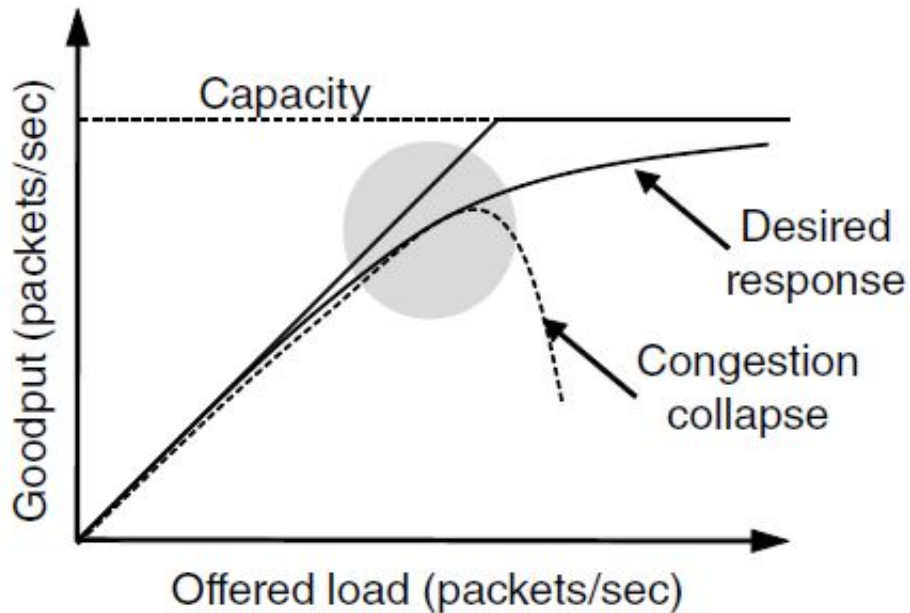
OK = Protocol functions correctly
 DUP = Protocol generates a duplicate message
 LOST = Protocol loses a message

Different combinations of client and server strategy

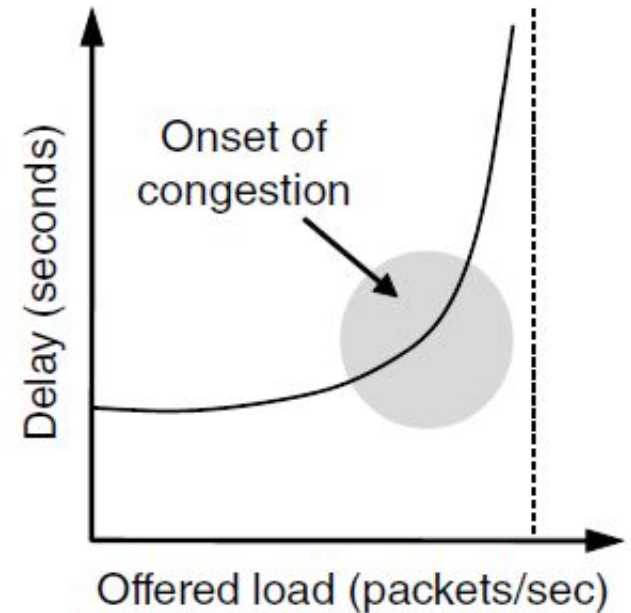
Congestion Control

- Desirable bandwidth allocation
- Regulating the sending rate

Desirable Bandwidth Allocation (1)



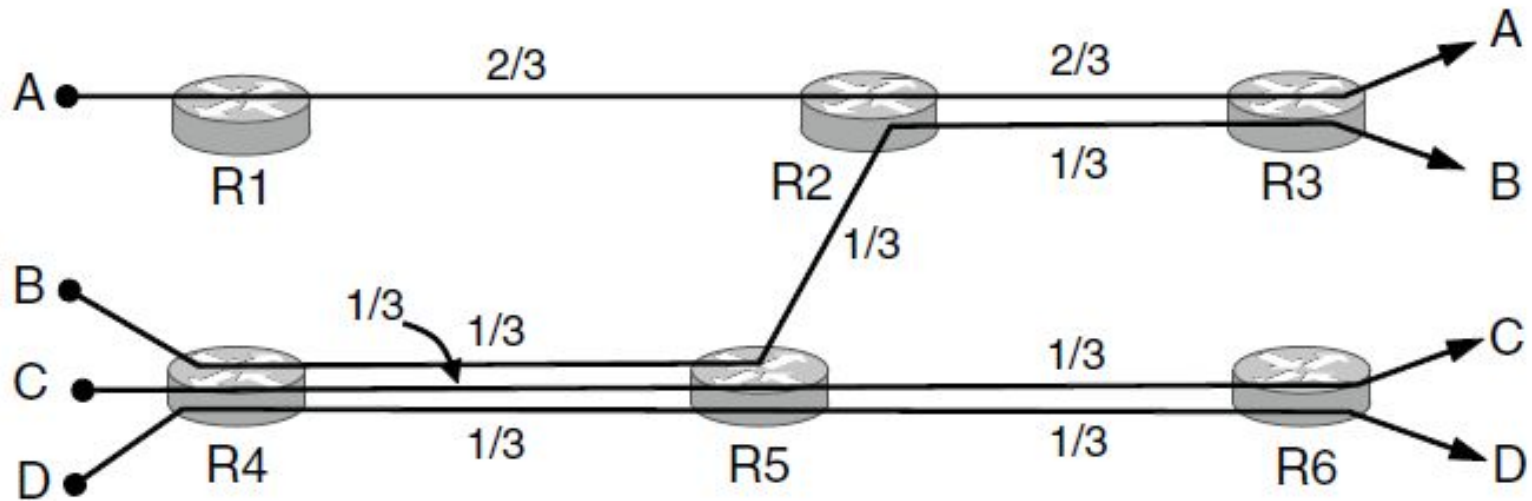
(a)



(b)

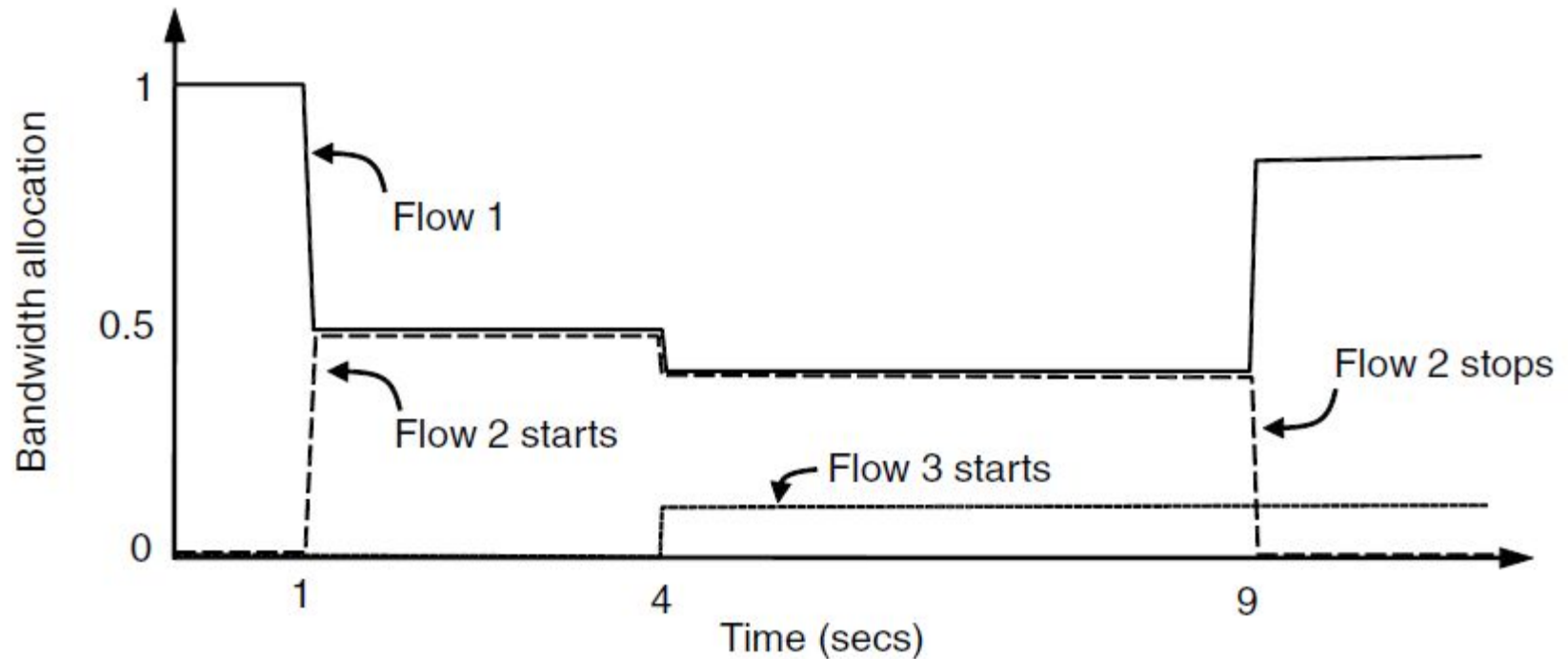
(a) Goodput and (b) delay as a function of offered load

Desirable Bandwidth Allocation (2)



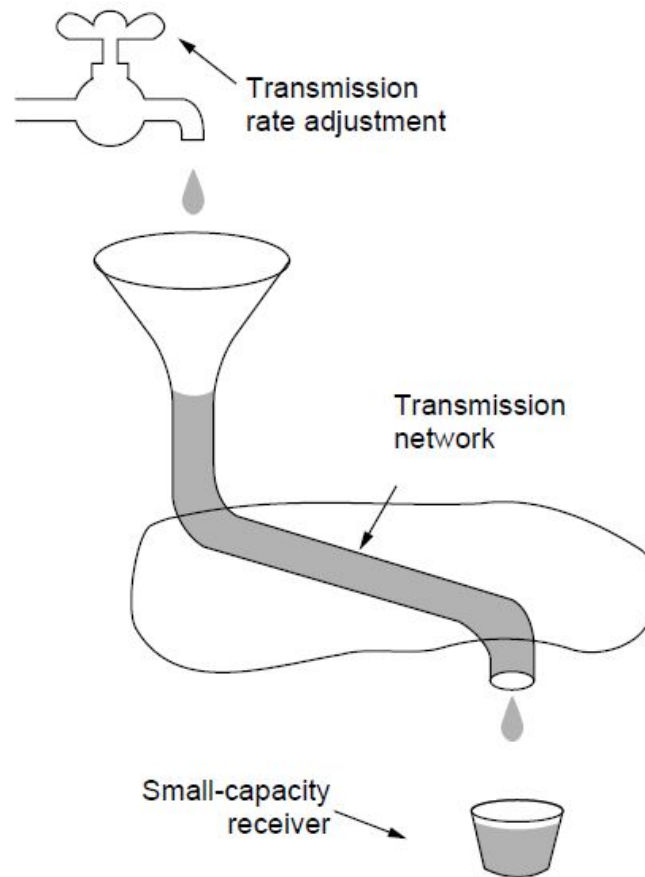
Max-min bandwidth allocation for four flows

Desirable Bandwidth Allocation (3)



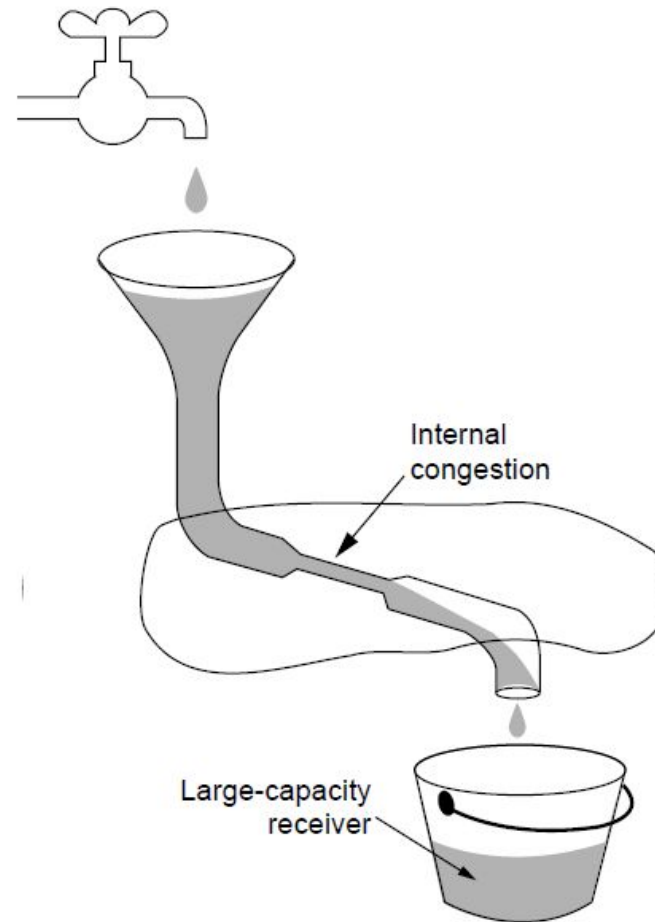
Changing bandwidth allocation over time

Regulating the Sending Rate (1)



A fast network feeding a low-capacity receiver

Regulating the Sending Rate (2)



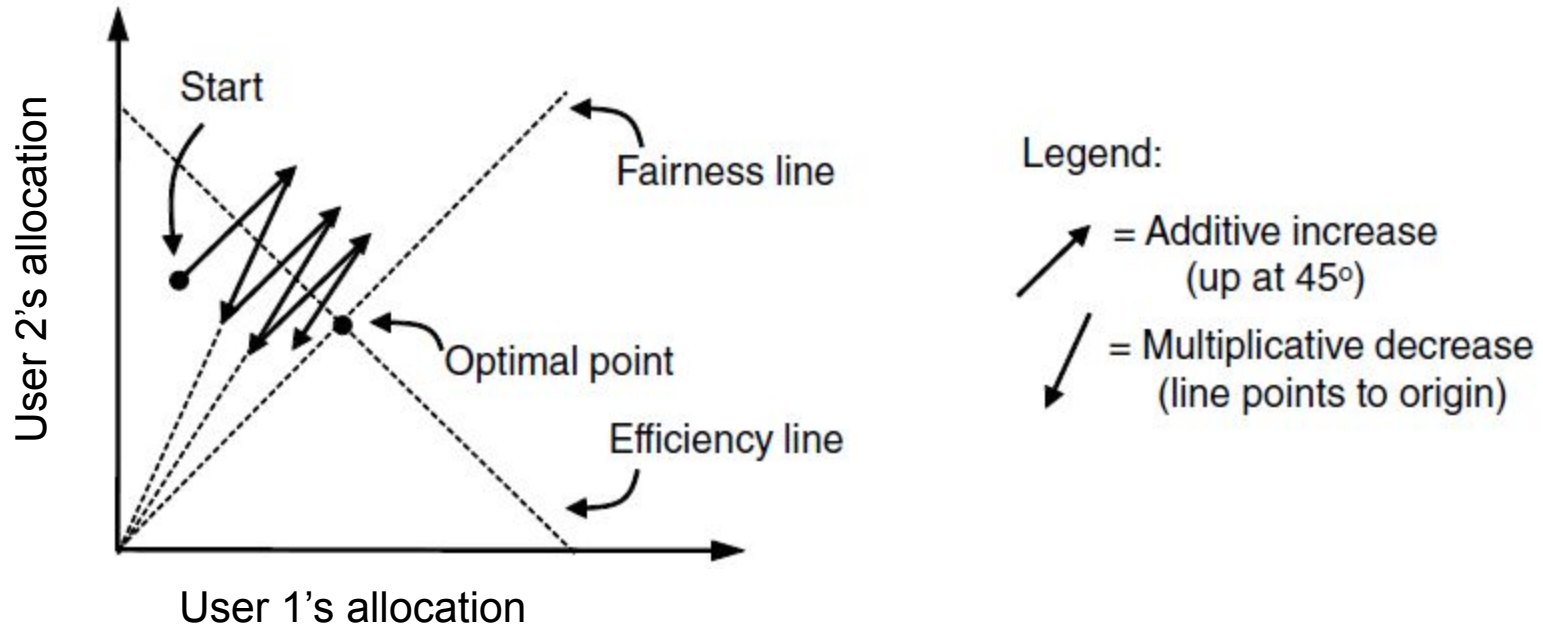
A slow network feeding a high-capacity receiver

Regulating the Sending Rate (3)

Protocol	Signal	Explicit?	Precise?
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

Some congestion control protocols

Regulating the Sending Rate (4)

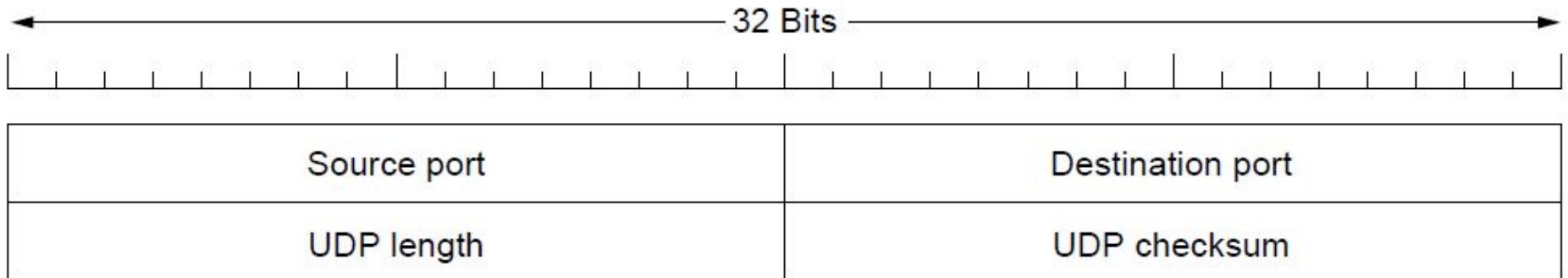


Additive Increase Multiplicative Decrease (AIMD) control law.

The Internet Transport Protocols: UDP

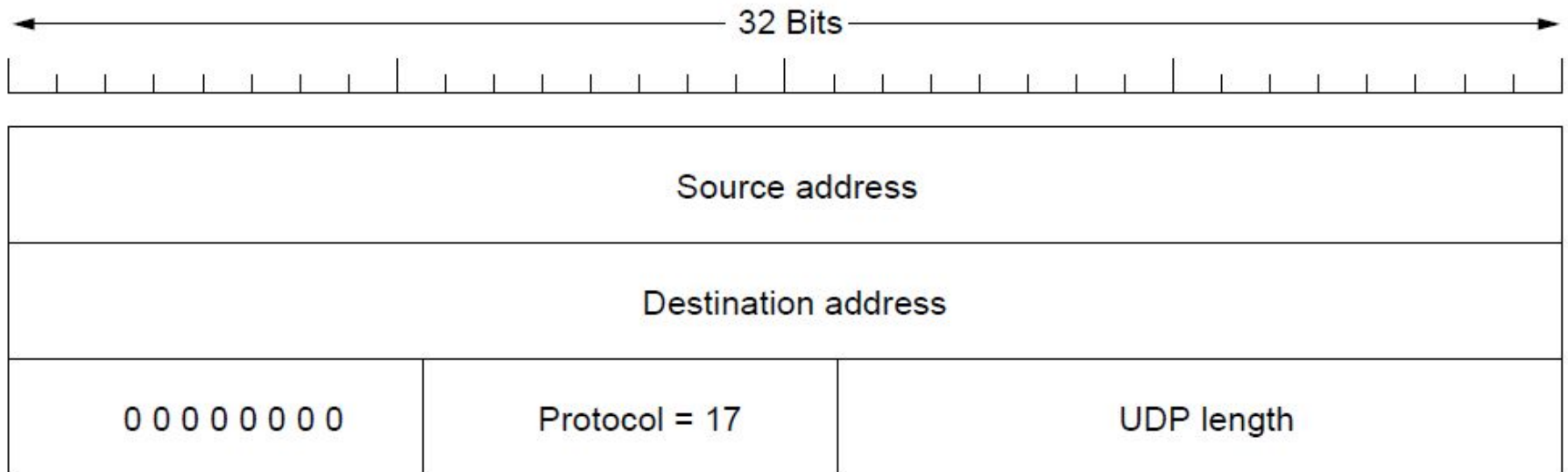
- Introduction to UDP
- Remote Procedure Call
- Real-Time Transport

Introduction to UDP (1)



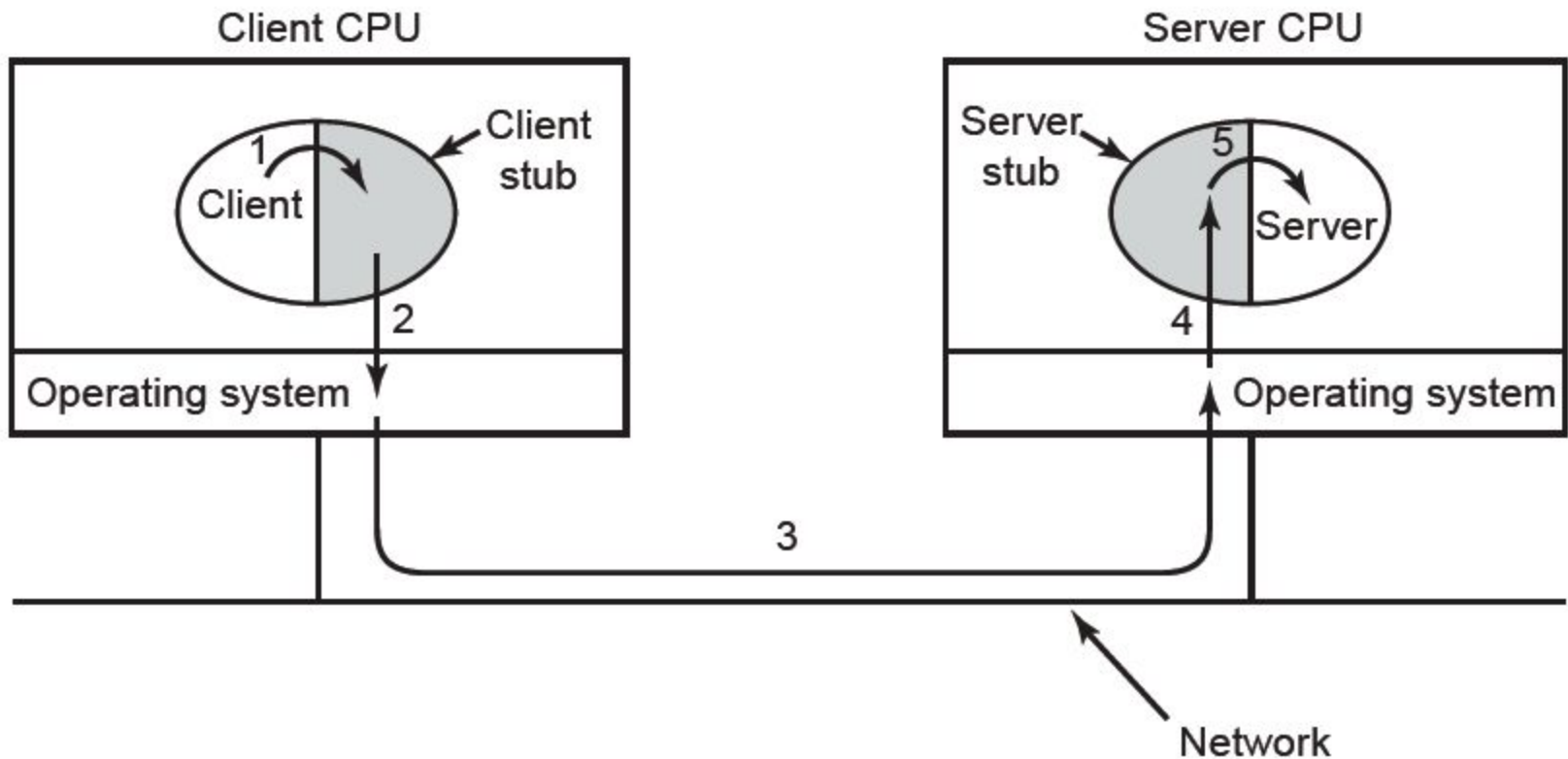
The UDP header.

Introduction to UDP (2)



The IPv4 pseudoheader included in the UDP checksum.

Remote Procedure Call

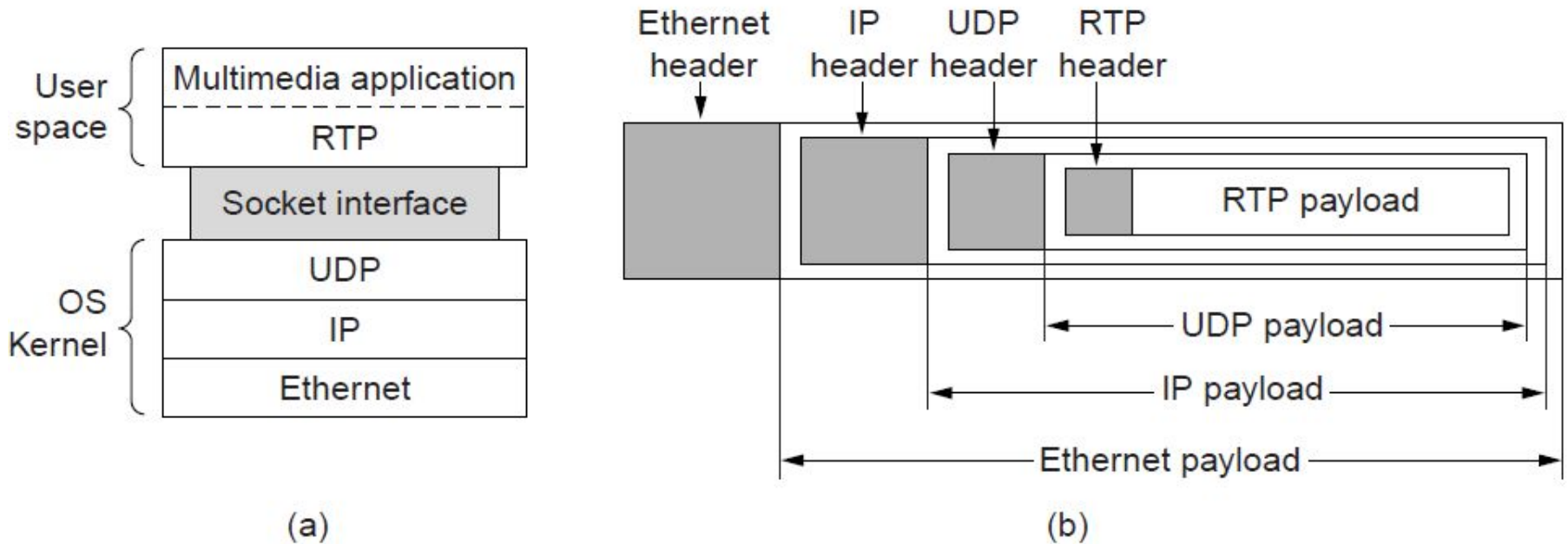


Steps in making a remote procedure call. The stubs are shaded.

Limitations of RPC

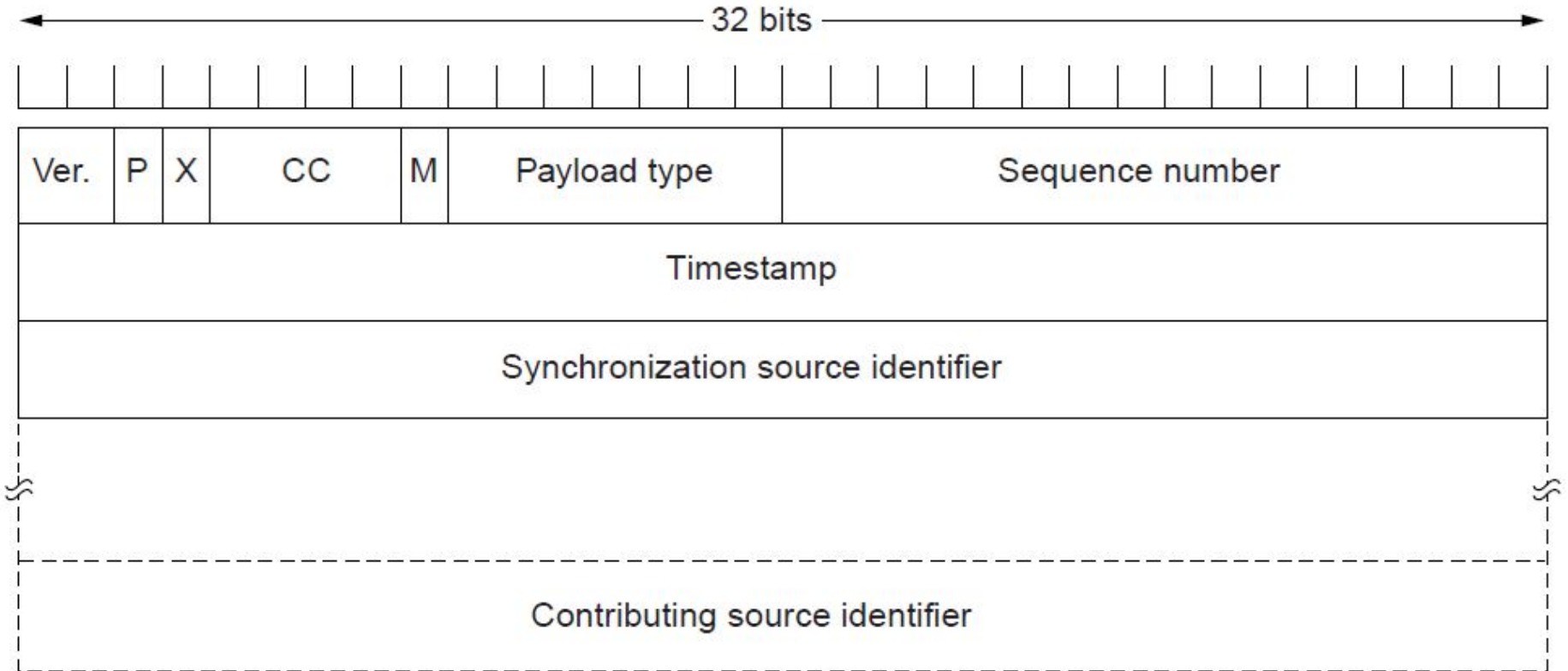
- a) Use of Pointers
- b) In languages like C, specifying size is not necessary e.g. you can multiply two arrays without specifying size
- c) Type of parameters may not be possible to deduce
 - *printf* in C
- d) Use of Global Variables- not shared

Real-Time Transport (1)



(a) The position of RTP in the protocol stack. (b) Packet nesting.

Real-Time Transport (2)



The RTP header

RTP Fields

The first twelve octets are present in every RTP packet, while the list of CSRC identifiers is present only when inserted by a mixer.

- **version (V) -: 2 bits**
 - This field identifies the version of RTP. The version defined by this specification is two (2).
- **padding (P) -: 1 bit**
 - If the padding bit is set, the packet contains one or more additional padding octets at the end which are not part of the payload.
- **extension (X)-: 1 bit**
 - If the extension bit is set, the fixed header is followed by exactly one header extension, with a format defined in Section 5.2.1.
- **CSRC count (CC) -: 4 bits**
 - The CSRC count contains the number of CSRC identifiers that follow the fixed header.

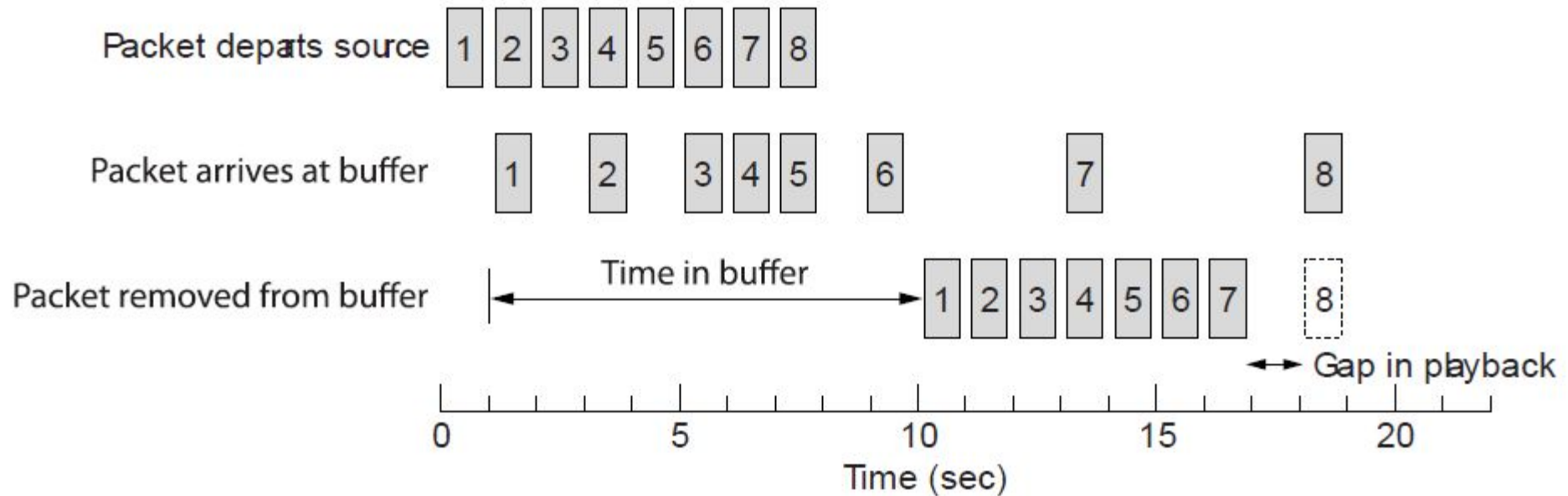
- **marker (M)-: 1 bit**
 - The interpretation of the marker is defined by a profile. It is intended to allow significant events such as frame boundaries to be marked in the packet stream.
- **payload type (PT)-: 7 bits**
 - This field identifies the format of the RTP payload and determines its interpretation by the application.
- **sequence number-: 16 bits**
 - The sequence number increments by one for each RTP data packet sent, and may be used by the receiver to detect packet loss and to restore packet sequence.
- **Timestamp-: 32 bits**
 - The timestamp reflects the sampling instant of the first octet in the RTP data packet. The sampling instant must be derived from a clock that increments monotonically and linearly in time to allow synchronization and jitter calculations

- **SSRC-:** 32 bits
 - The SSRC field identifies the synchronization source.
- **CSRC list-:** 0 to 15 items, 32 bits each
 - The CSRC list identifies the contributing sources for the payload contained in this packet. The number of identifiers is given by the CC field. If there are more than 15 contributing sources, only 15 may be identified. CSRC identifiers are inserted by mixers, using the SSRC identifiers of contributing sources.

a) RTCP

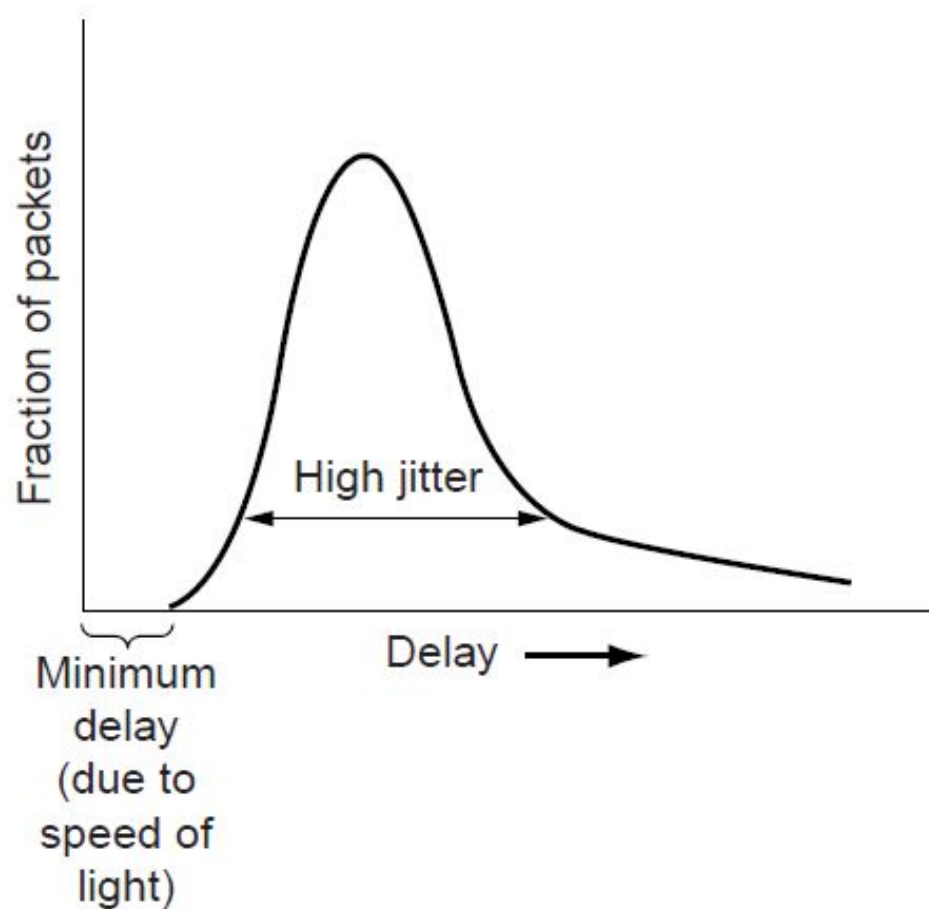
- A companion control protocol to RTP.
- RTCP is not a signaling protocol
- RTCP is used to collect end-to-end information about the quality of the session to each participant.
- RTCP packet types:
 - SR: sender report
 - RR: receiver report
 - SDES: Source DEscription
 - BYE: Hangs up from a session
 - APP: Application-Specific packet

Real-Time Transport (3)



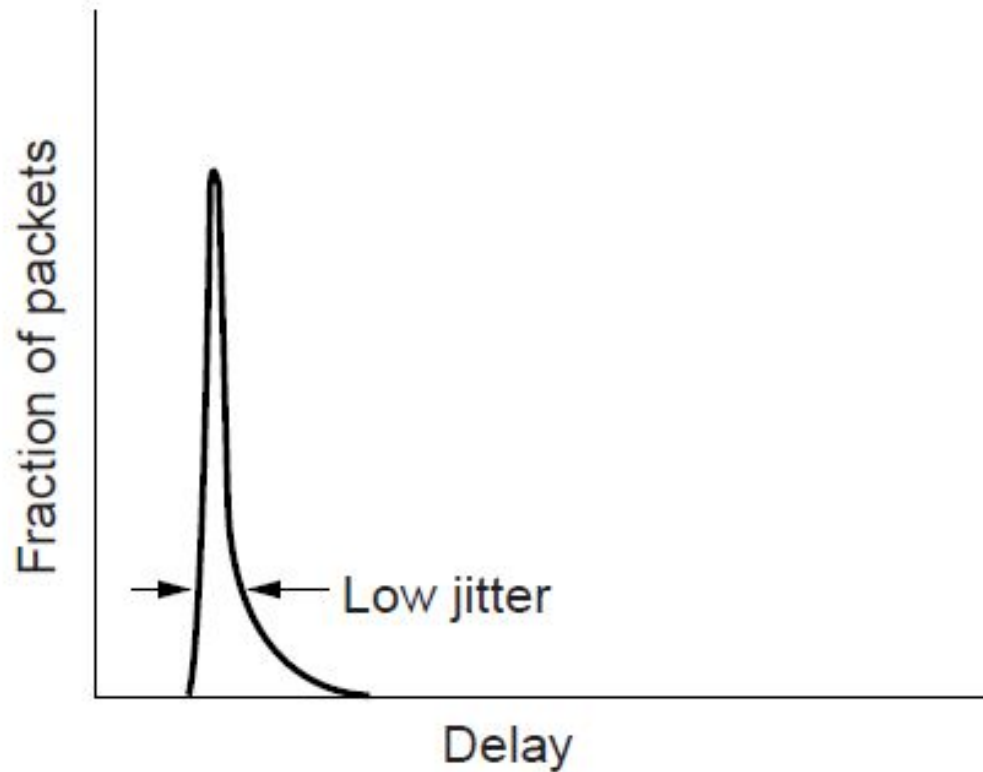
Smoothing the output stream by buffering packets

Real-Time Transport (3)



High jitter

Real-Time Transport (4)



Low jitter

The Internet Transport Protocols: TCP (1)

- Introduction to TCP
- The TCP service model
- The TCP protocol
- The TCP segment header
- TCP connection establishment
- TCP connection release

The Internet Transport Protocols: TCP (2)

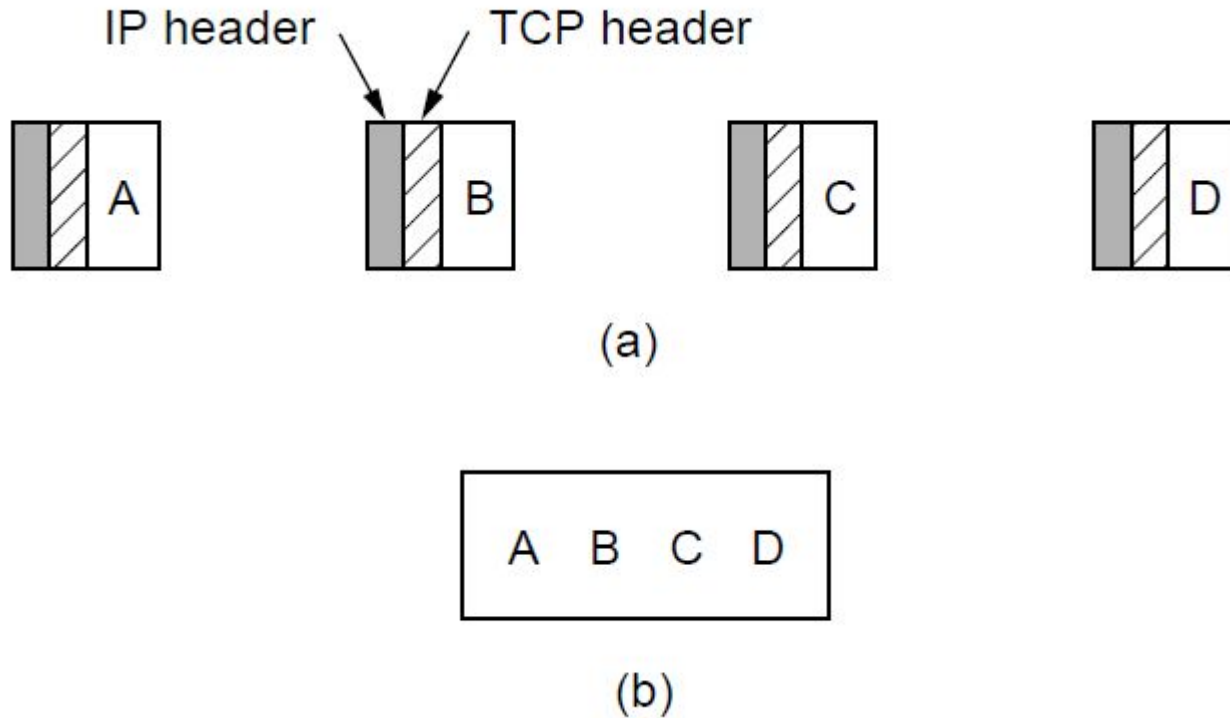
- TCP connection management modeling
- TCP sliding window
- TCP timer management
- TCP congestion control
- TCP futures

The TCP Service Model (1)

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

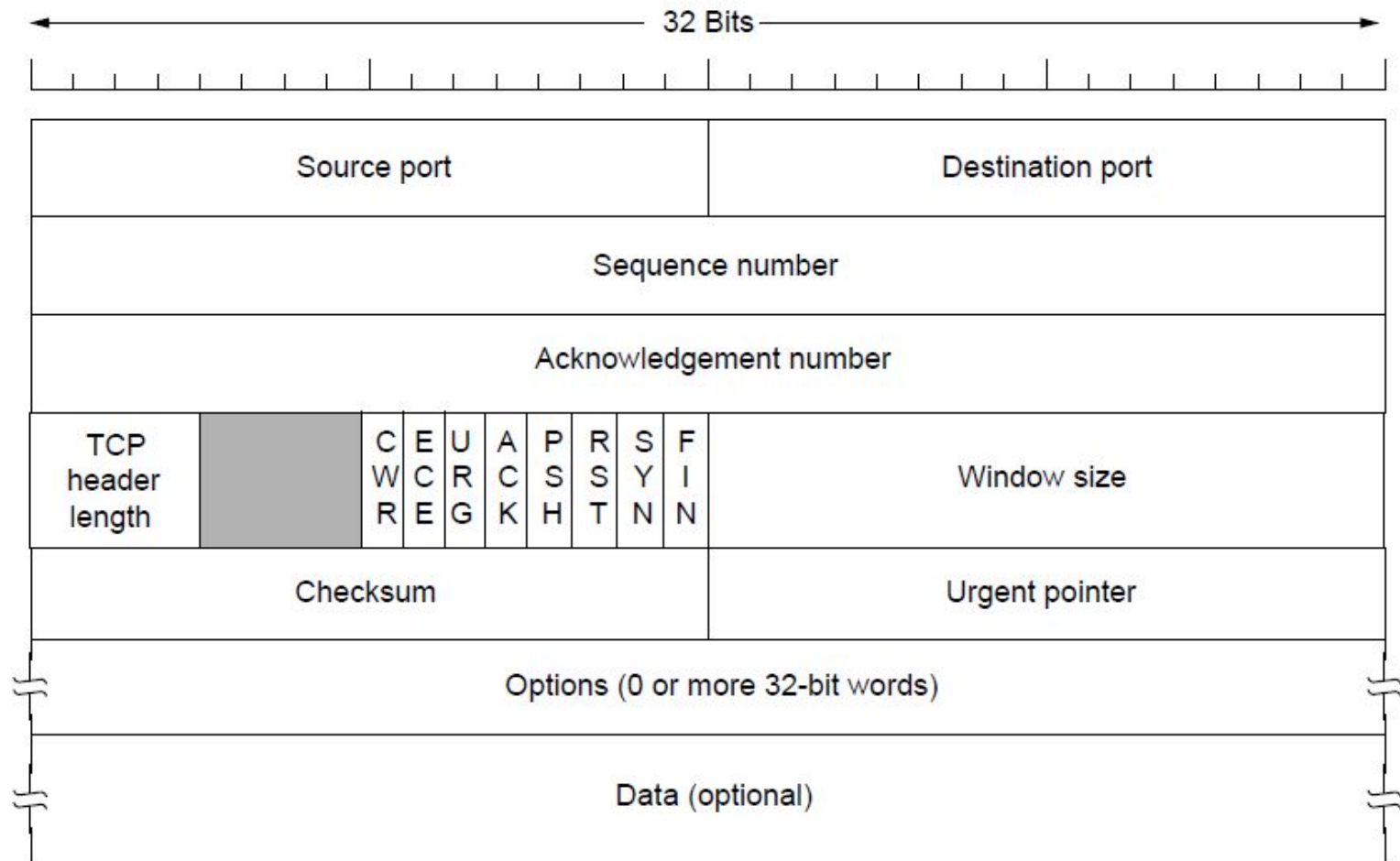
Some assigned ports

The TCP Service Model (2)



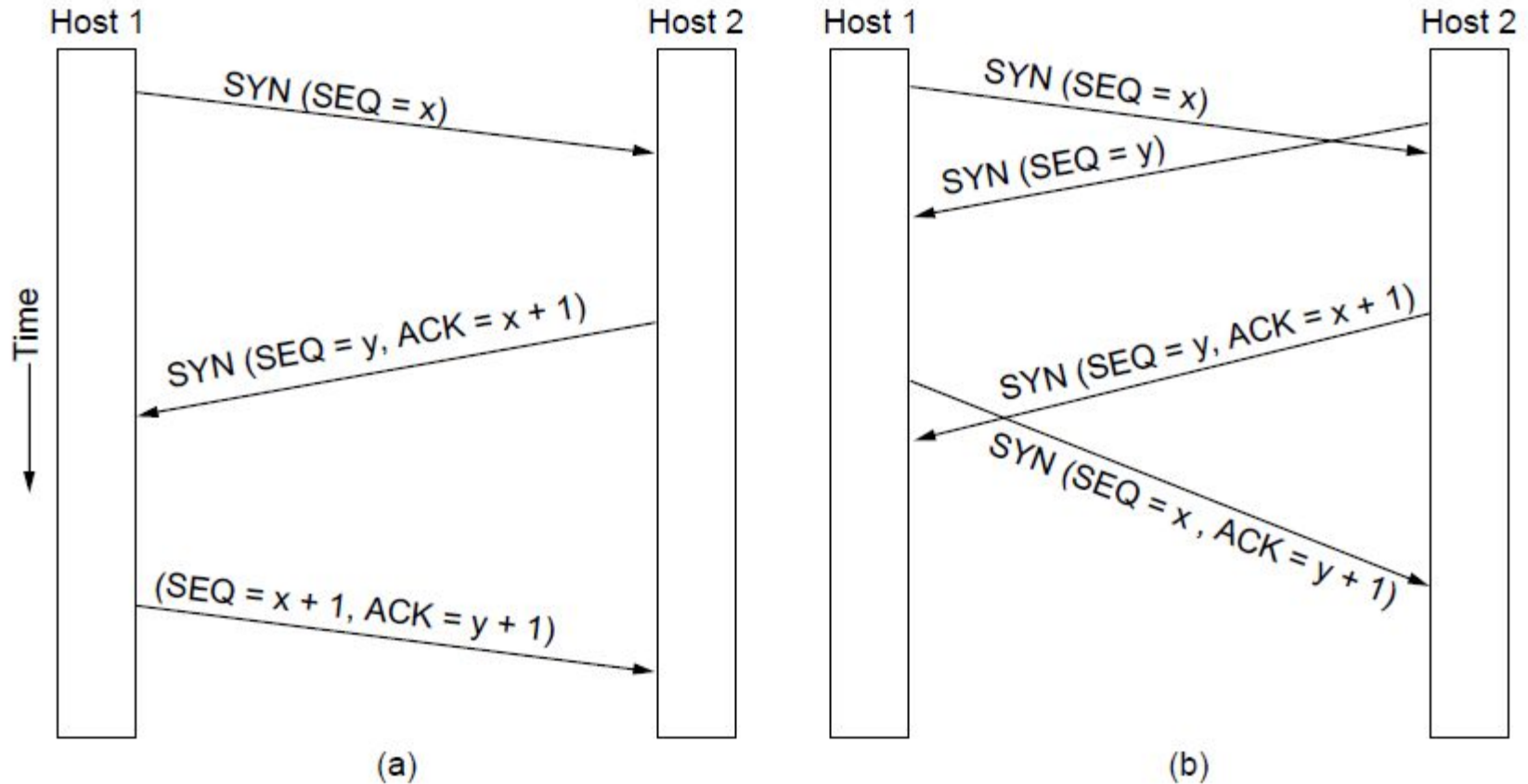
- (a) Four 512-byte segments sent as separate IP datagrams
- (b) The 2048 bytes of data delivered to the application in a single READ call

The TCP Segment Header



The TCP header.

TCP Connection Establishment



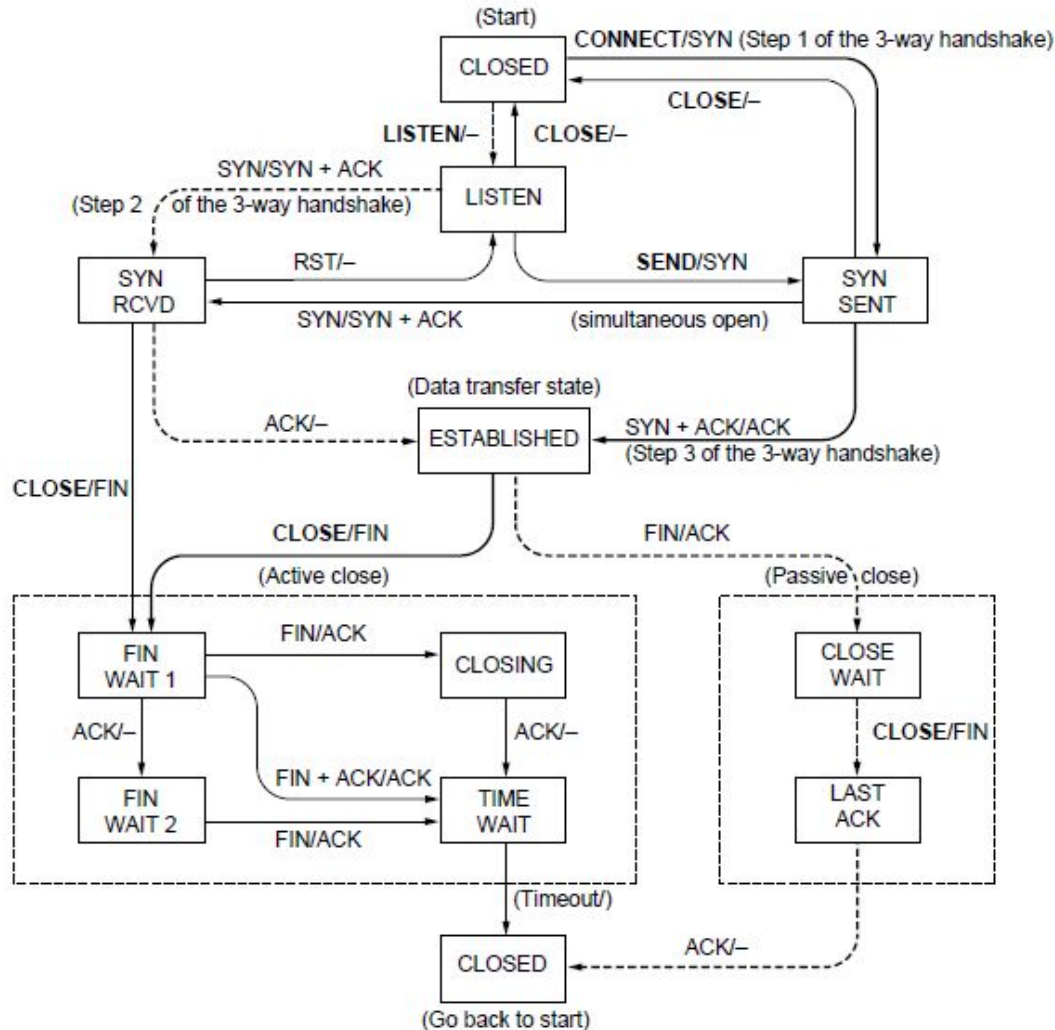
- (a) TCP connection establishment in the normal case.
- (b) Simultaneous connection establishment on both sides.

TCP Connection Management Modeling (1)

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

The states used in the TCP connection management finite state machine.

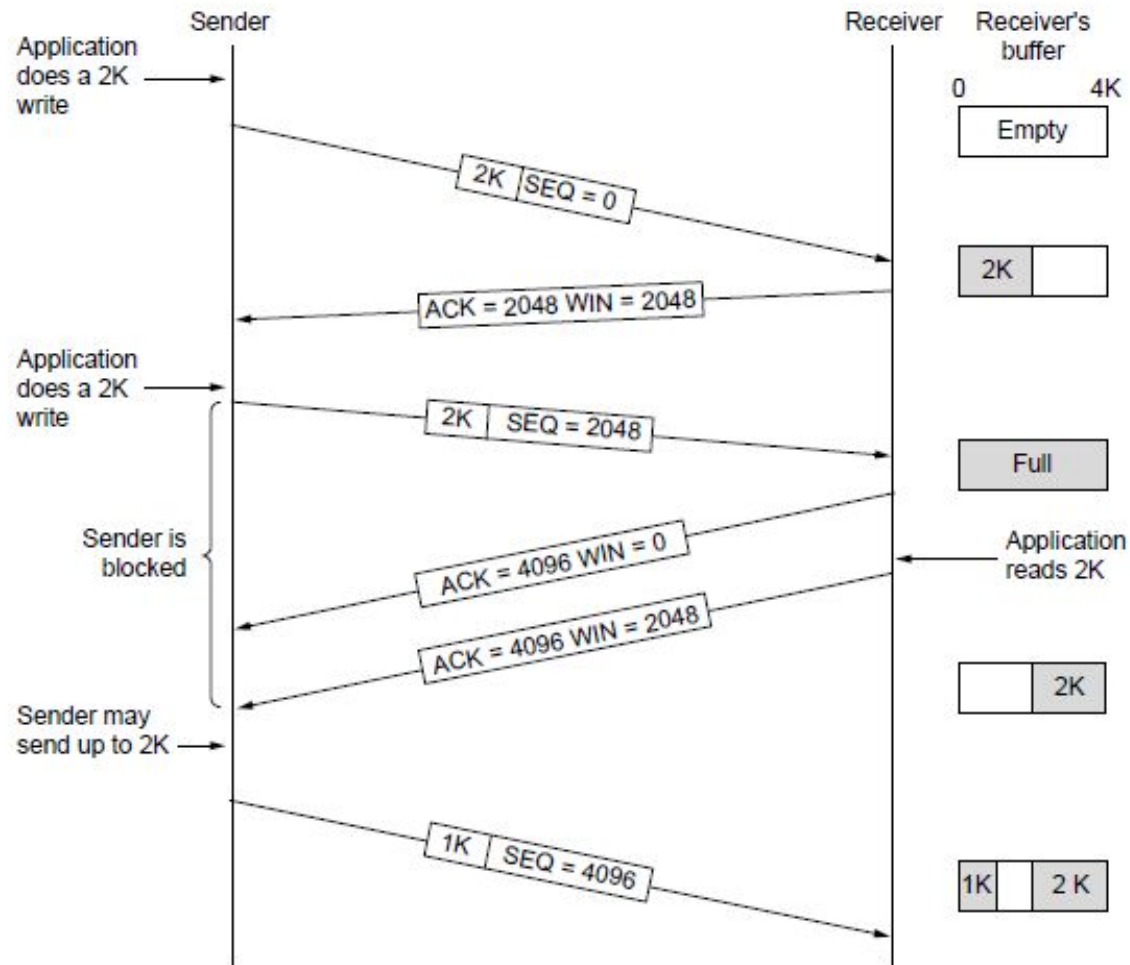
TCP Connection Management Modeling (2)



TCP connection management finite state machine.

The heavy solid line is the normal path for a client. The heavy dashed line is the normal path for a server. The light lines are unusual events. Each transition is labeled by the event causing it and the action resulting from it, separated by a slash.

TCP Sliding Window (1)



Window management in TCP

What if Window=0

No data can be sent by the sender except

1. Urgent data
2. 1-Byte segment to make the receiver re-announce the next byte expected and the window size-required to avoid deadlock in case window announcement is lost

Nagle's Algorithm

If application has one only byte to send continuously, 41 Byte is transmitted by the network layer (sender), receiver also sends 40 Byte acknowledgement- Poor Efficiency.

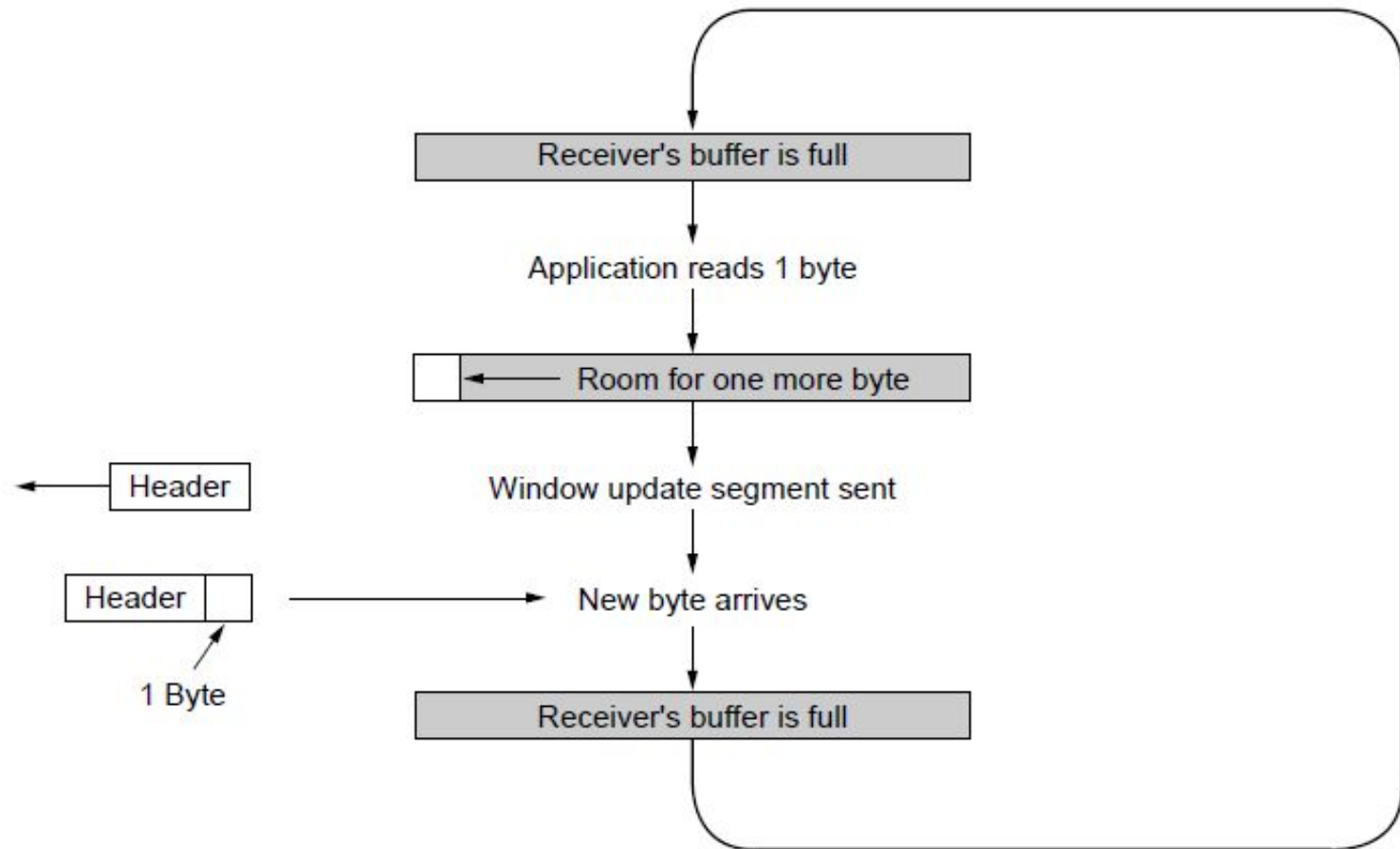
Nagle Suggested-

- Send first byte, buffer subsequent bytes till acknowledgement for the previous byte is received.
- Send buffered data in one segment, buffer till all the bytes are acknowledged.

Difficulties-

X-windows mouse clicks

Silly Window Syndrome



Silly window syndrome

Clark's Solution

Delay advertizing the window till

- It can handle the maximum segment size it advertized when the connection was established, or
- Its buffer is half empty

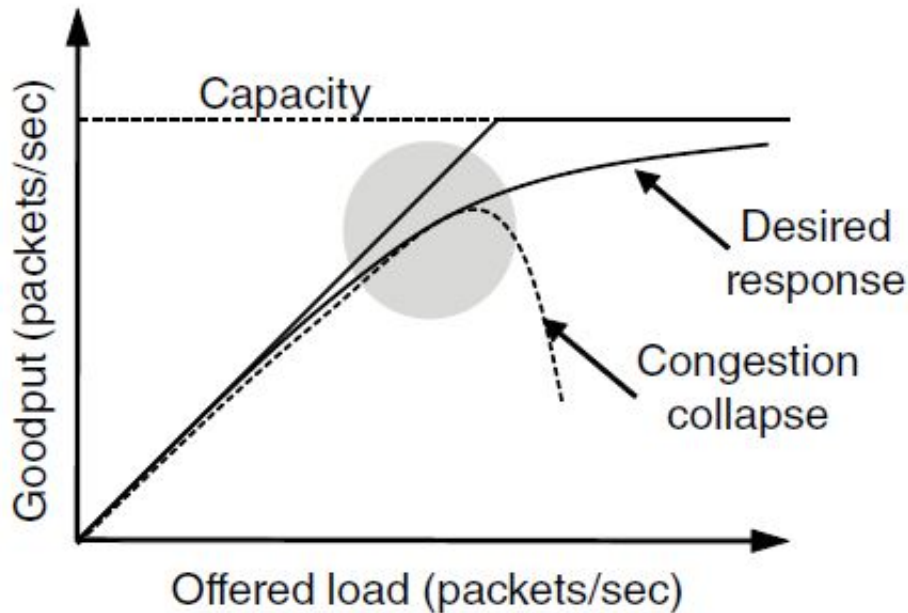
Whichever is smaller

Combine Nagle's Algo and Clark's Solution
improves efficiency and performance

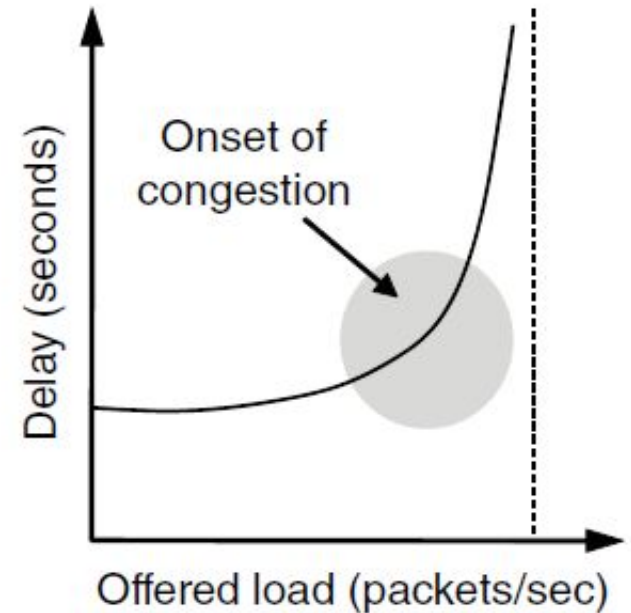
Congestion Control

- Desirable bandwidth allocation
- Regulating the sending rate

Desirable Bandwidth Allocation (1)



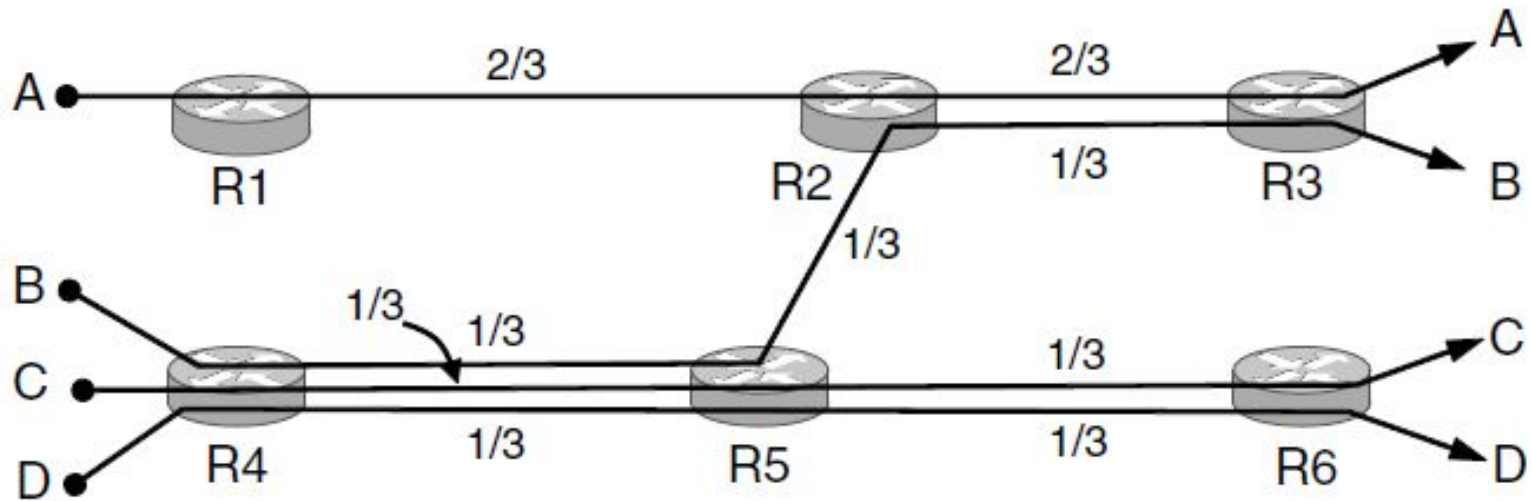
(a)



(b)

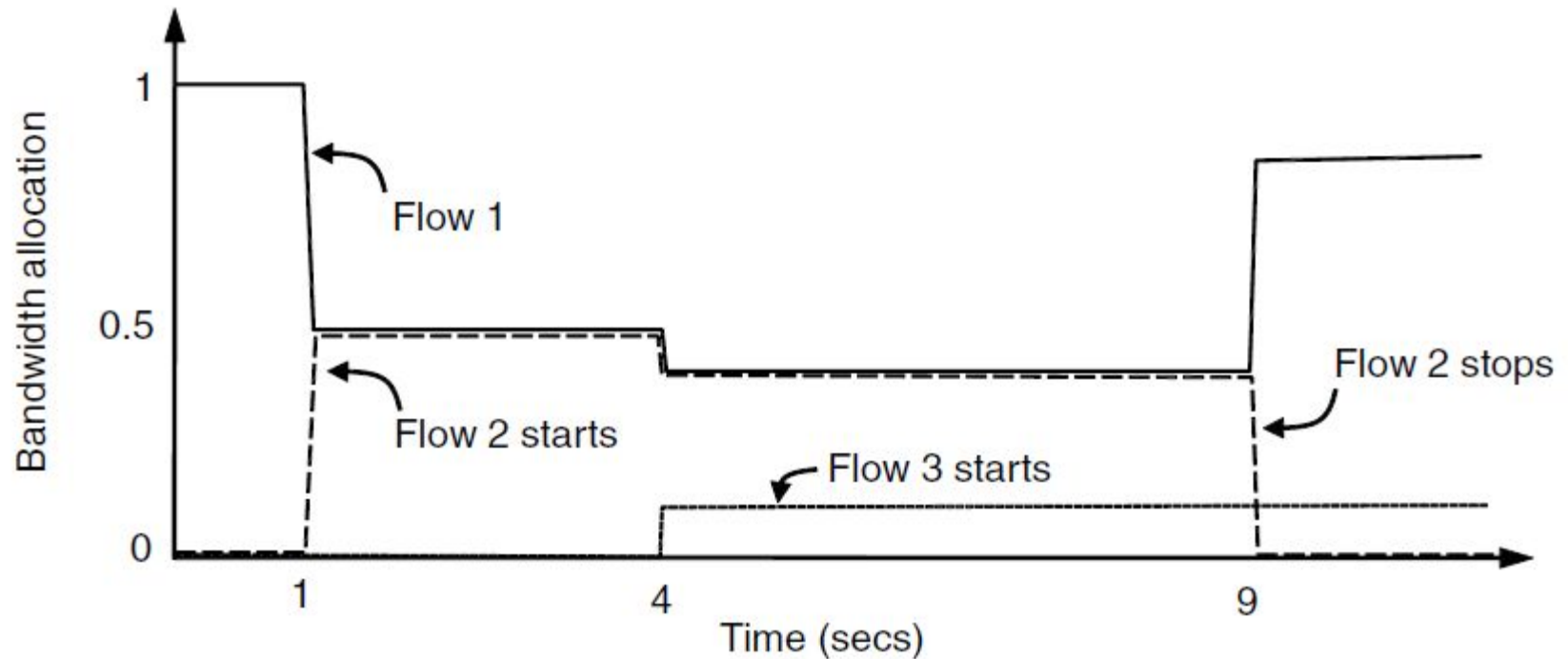
(a) Goodput and (b) delay as a function of offered load

Desirable Bandwidth Allocation (2)



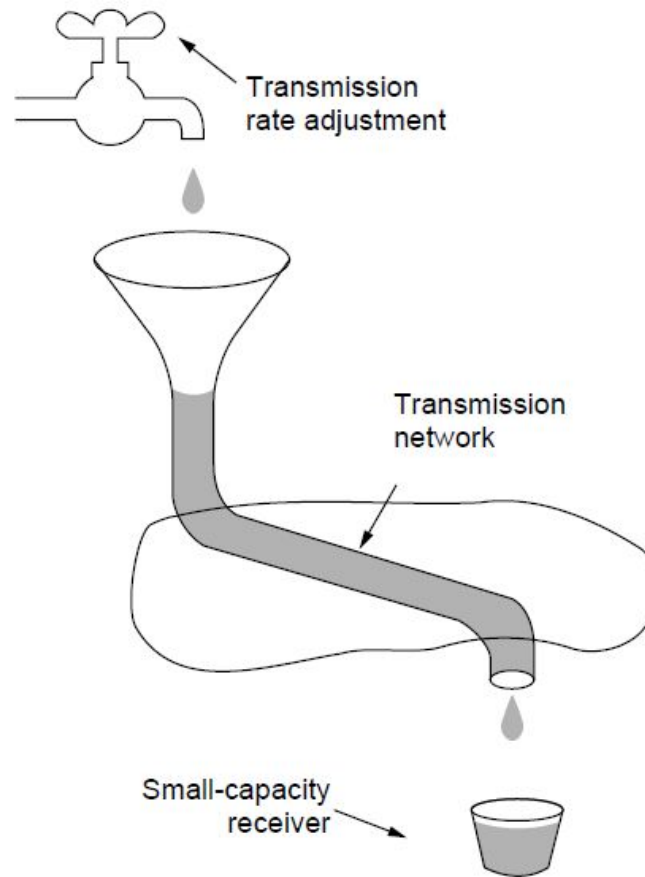
Max-min bandwidth allocation for four flows

Desirable Bandwidth Allocation (3)



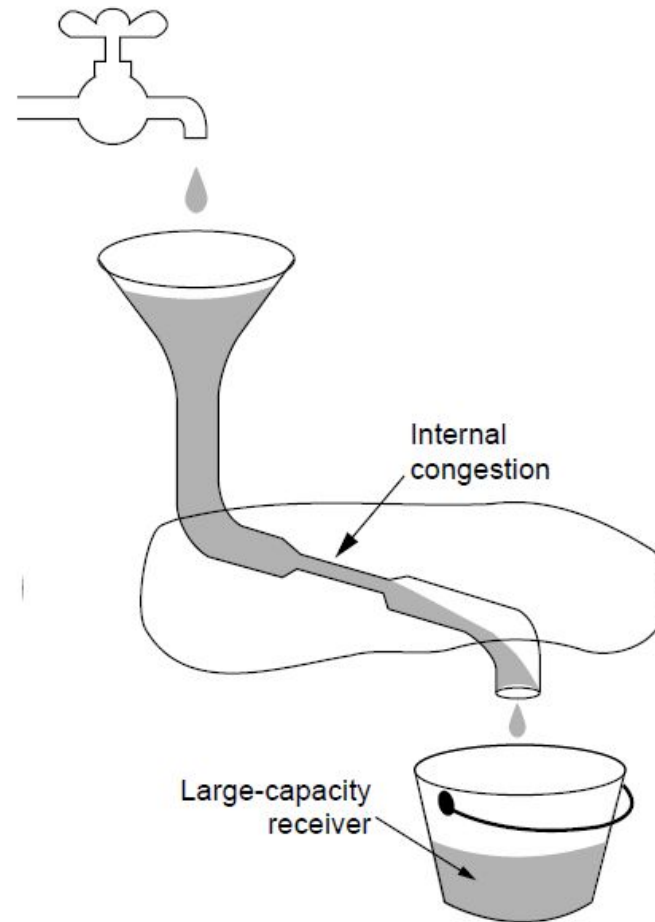
Changing bandwidth allocation over time

Regulating the Sending Rate (1)



A fast network feeding a low-capacity receiver

Regulating the Sending Rate (2)



A slow network feeding a high-capacity receiver

TCP Congestion Control

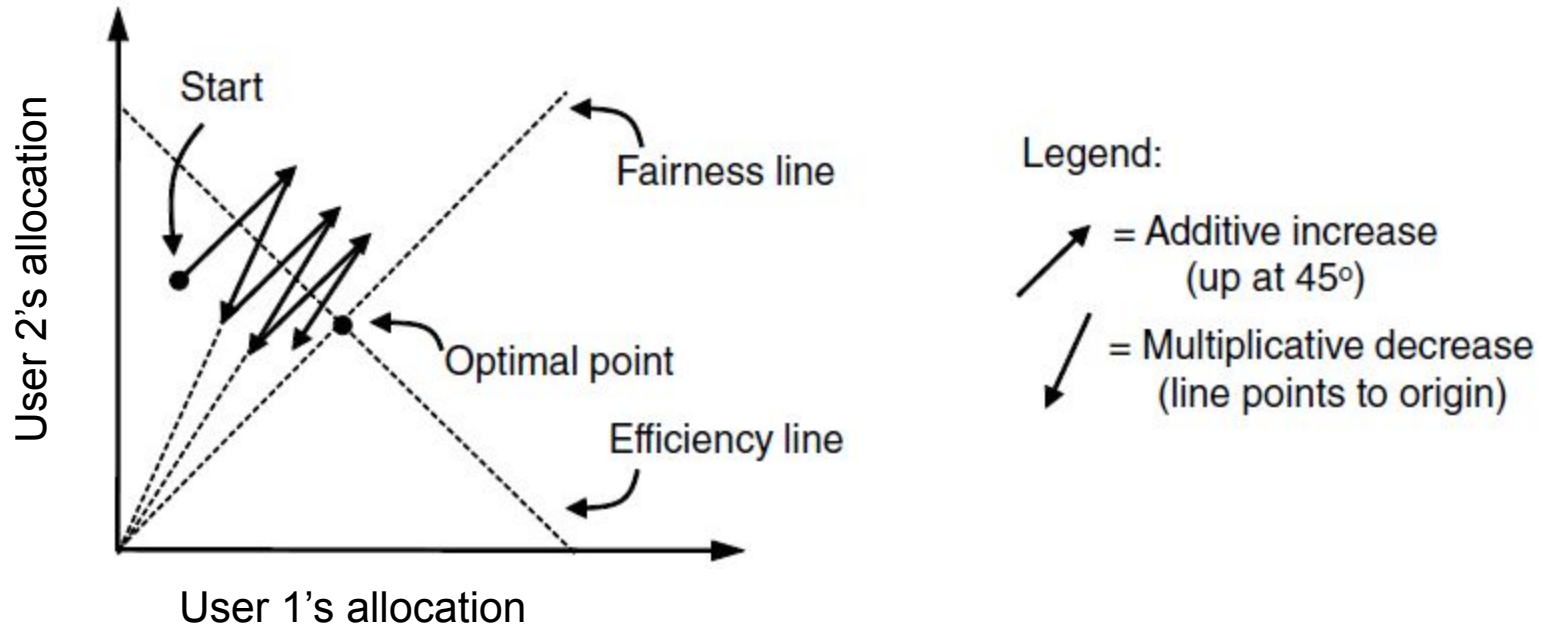
- a) Two windows
 - the window a receiver grants
 - Congestion window
- b) Only the smaller of the two is sent
- c) Initially, congestion window equals maximum segment size on the connection, if acknowledged, add the acknowledged window to congestion window
- d) Slow start – 1024 byte - move exponentially – SLOW START
- e) Set threshold – grow linearly after that

Regulating the Sending Rate (3)

Protocol	Signal	Explicit?	Precise?
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

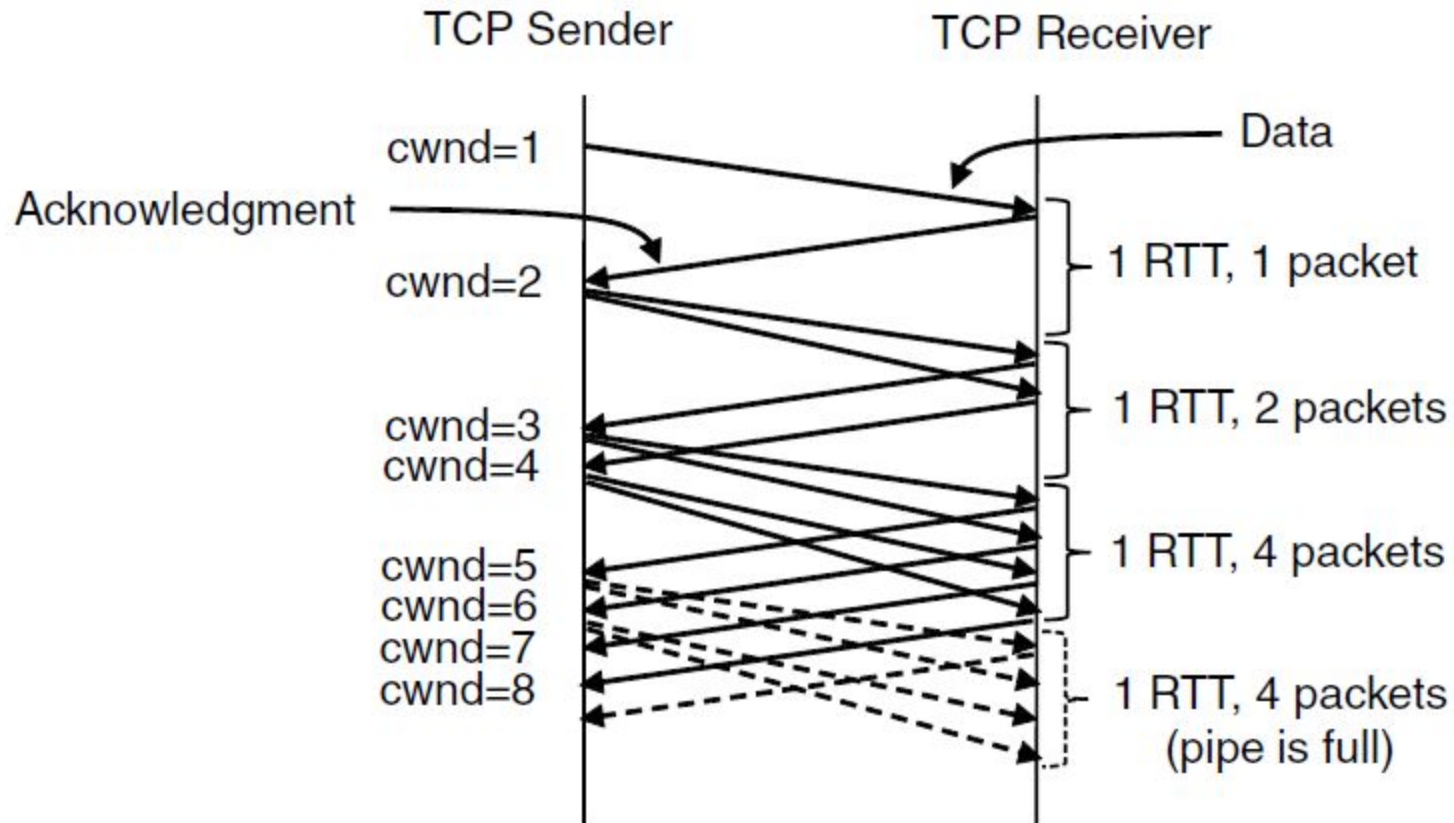
Some congestion control protocols

Regulating the Sending Rate (4)



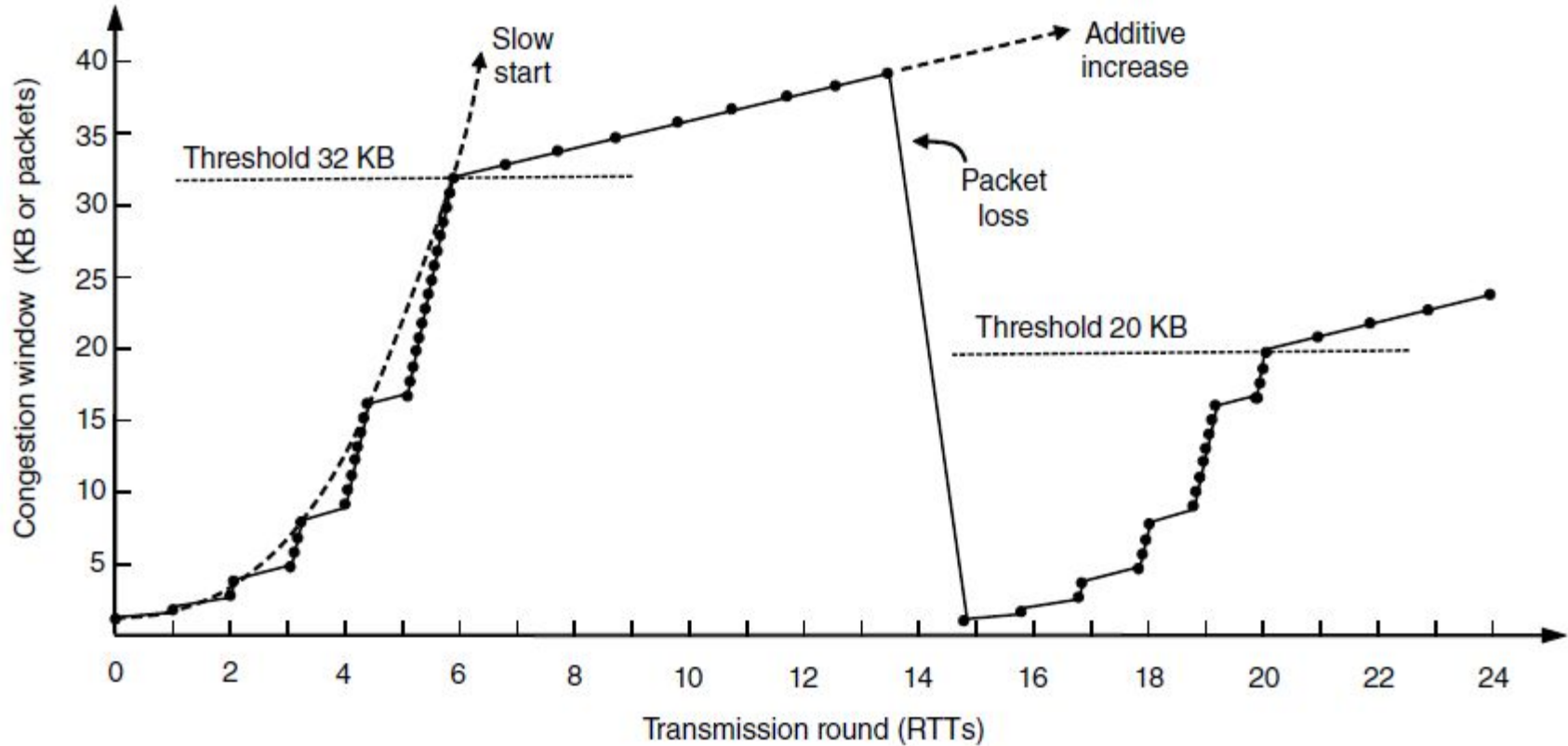
Additive Increase Multiplicative Decrease (AIMD) control law.

TCP Congestion Control (1)



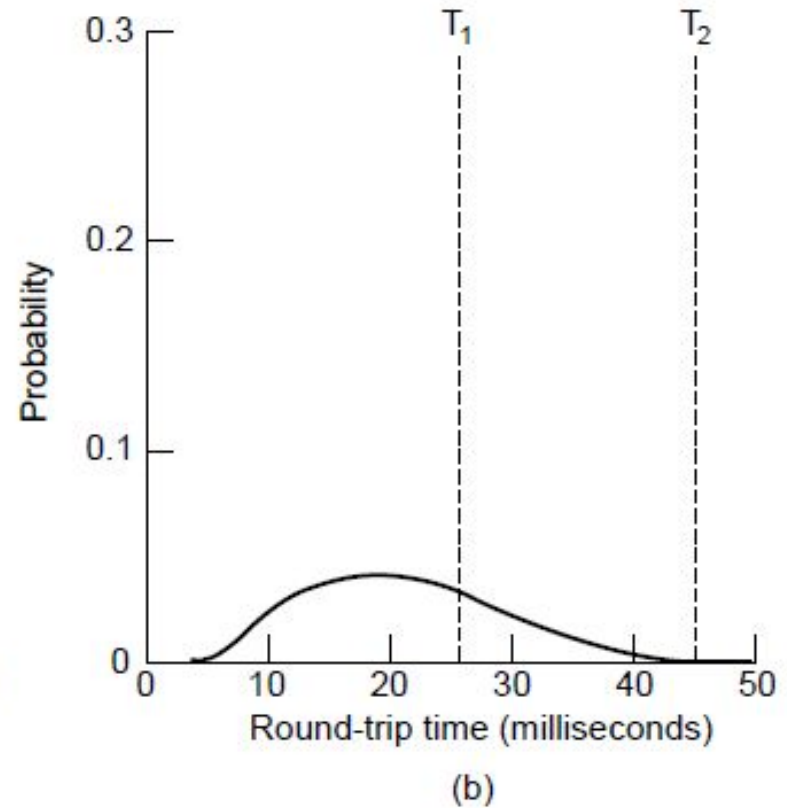
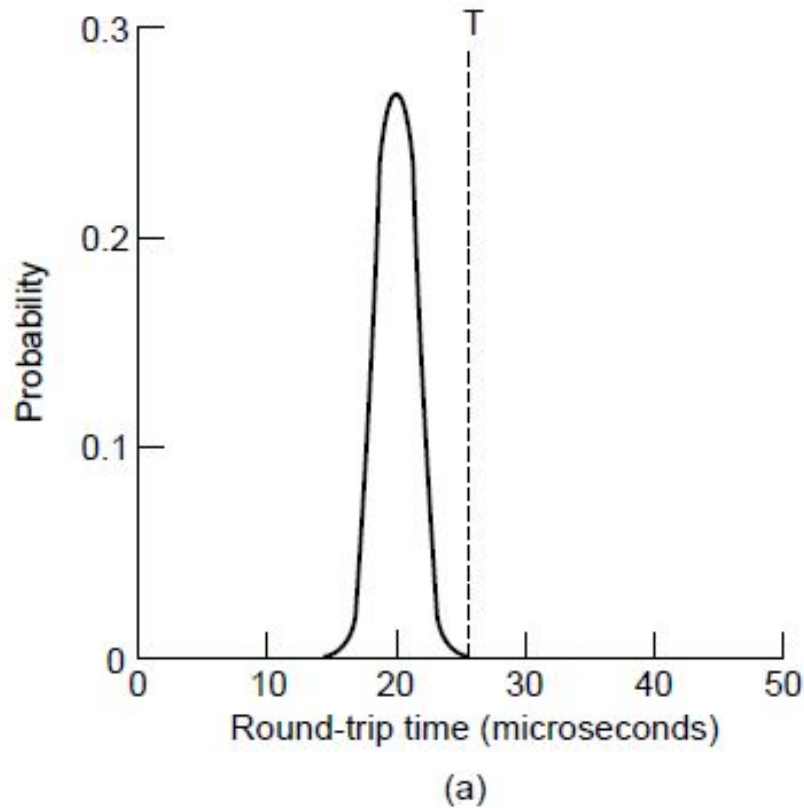
Slow start from an initial congestion window of 1 segment

TCP Congestion Control (3)



Slow start followed by additive increase in TCP Tahoe.

TCP Timer Management



a) Probability density of acknowledgment arrival times in data link layer. (b) ... for TCP

Retransmission Timer

- a) $RTT = \alpha RTT + (1 - \alpha) M$, $\alpha = 7/8$ (Smoothing Factor)
- b) $\text{Time-out} = \beta RTT$

Jacobson suggested making β roughly proportional to the SD of the acknowledgement arrival time probability density function. He suggested using mean deviation as an estimate to SD.

Modified, to use

- a) $D = \alpha D + (1 - \alpha) |RTT - M|$
- b) $\text{Timeout} = RTT + 4 \times D$

Retransmission Timer

- a) Karn's algorithm – In dynamic estimation of RTT when a segment times out and is resend.
- b) It is not clear whether the acknowledgement is from the original or resend.
- c) SO don't include resent packet's RTT into calculation.
- d) Each time failure double Time-out time

Persistence Timer

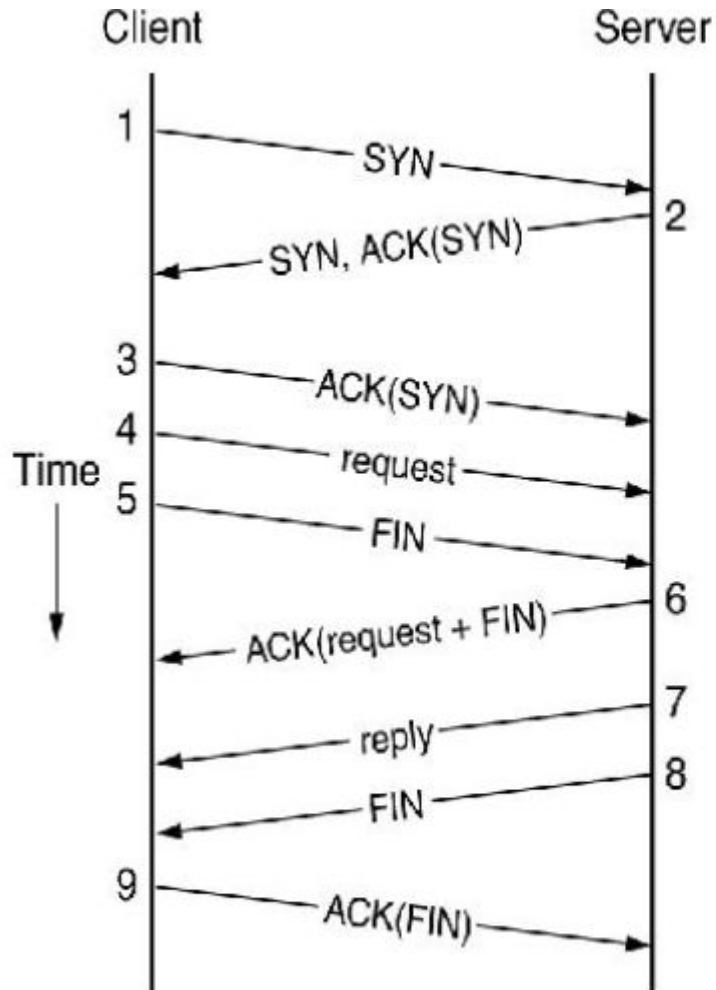
- a) When Persistence timer goes off the transmitter pings the receiver whether buffer space is available

Keepalive Timer

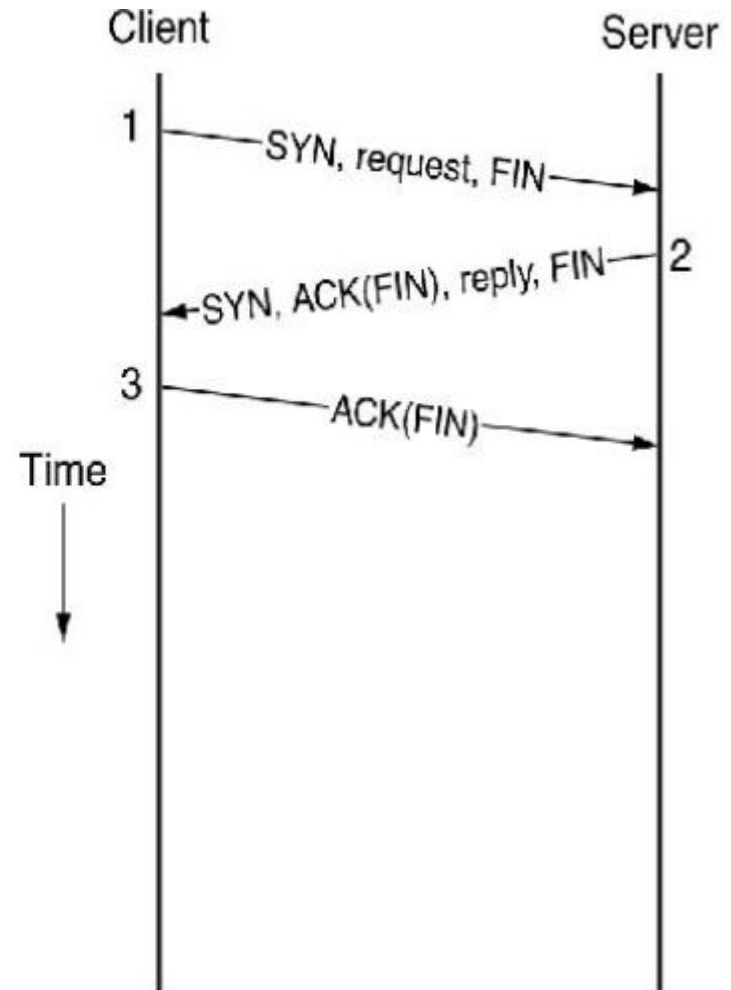
If idle □ checks whether the connection is active and then if not closes connection

CONTROVERSIAL – It may stop healthy connection due to transient network partitioning

RPC Over TCP-T/TCP



(a)

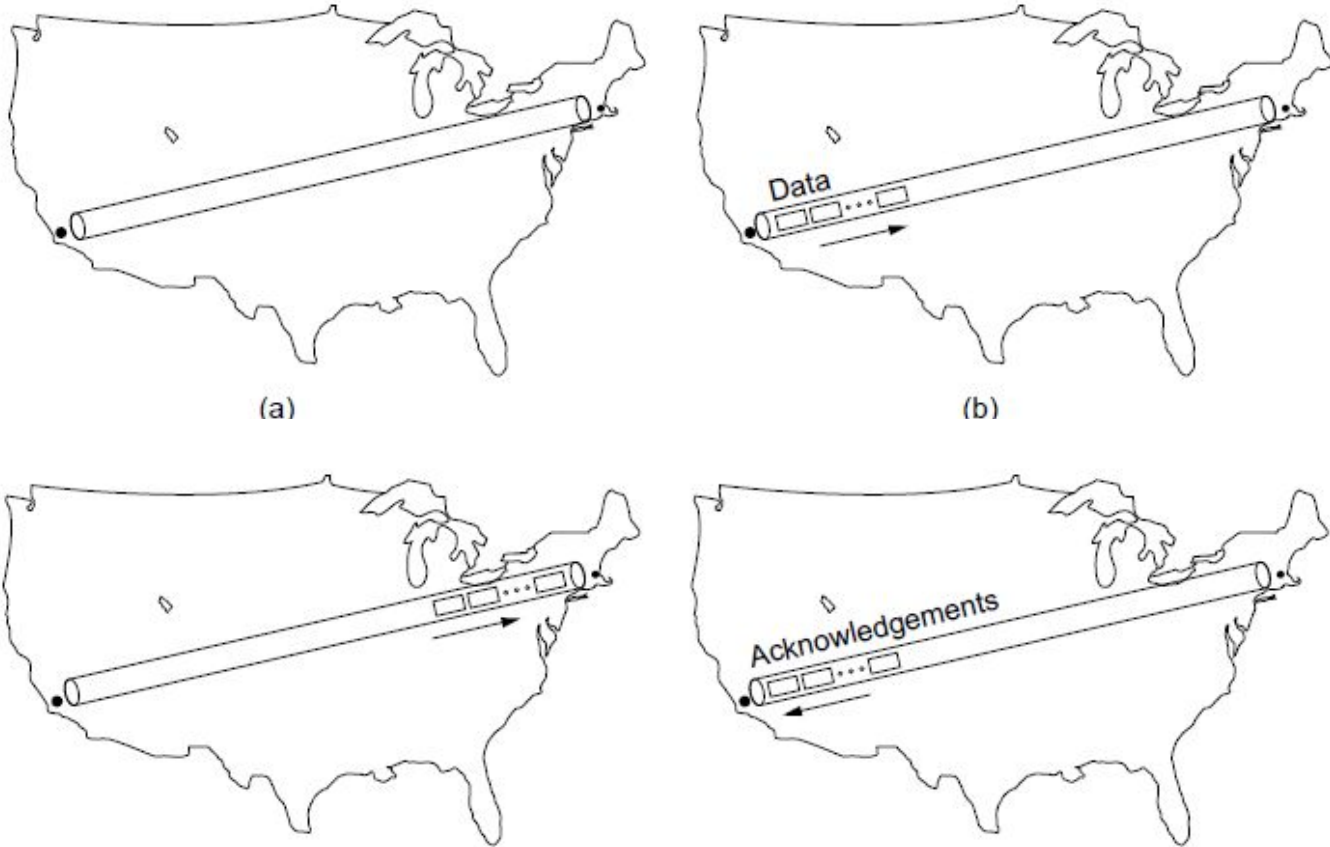


(b)

Performance Issues

- Performance problems in computer networks
- Network performance measurement
- System design for better performance
- Fast TPDU processing
- Protocols for high-speed networks

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 msec . (d) After 40 msec .

Network Performance Measurement (1)

Steps to performance improvement

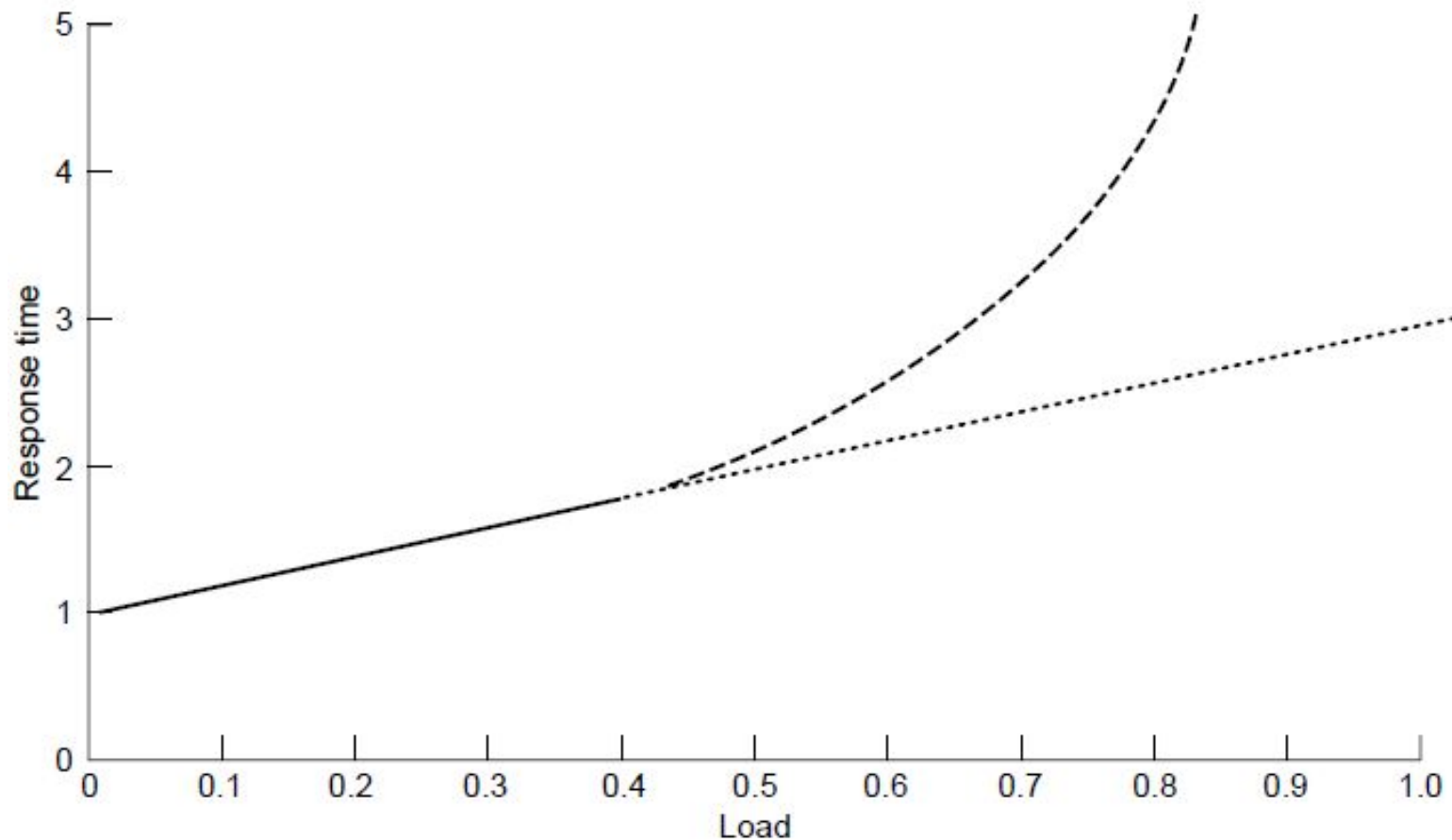
1. Measure relevant network parameters, performance.
2. Try to understand what is going on.
3. Change one parameter.

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

Network Performance Measurement (3)



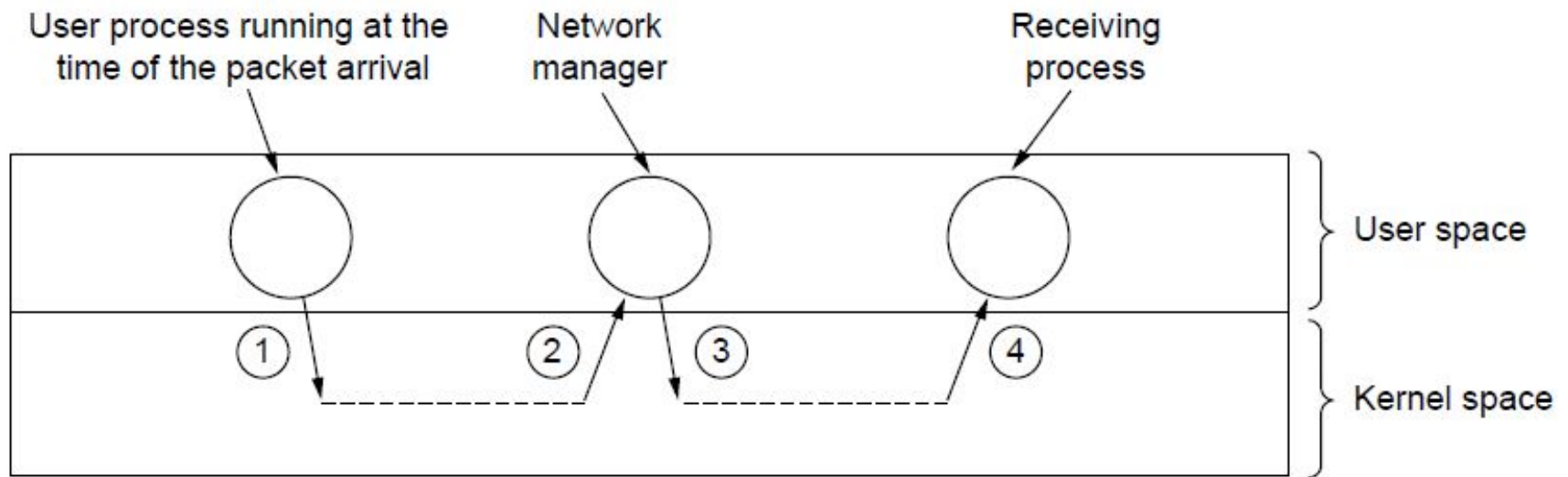
Response as a function of load.

System Design for Better Performance (1)

Rules of thumb

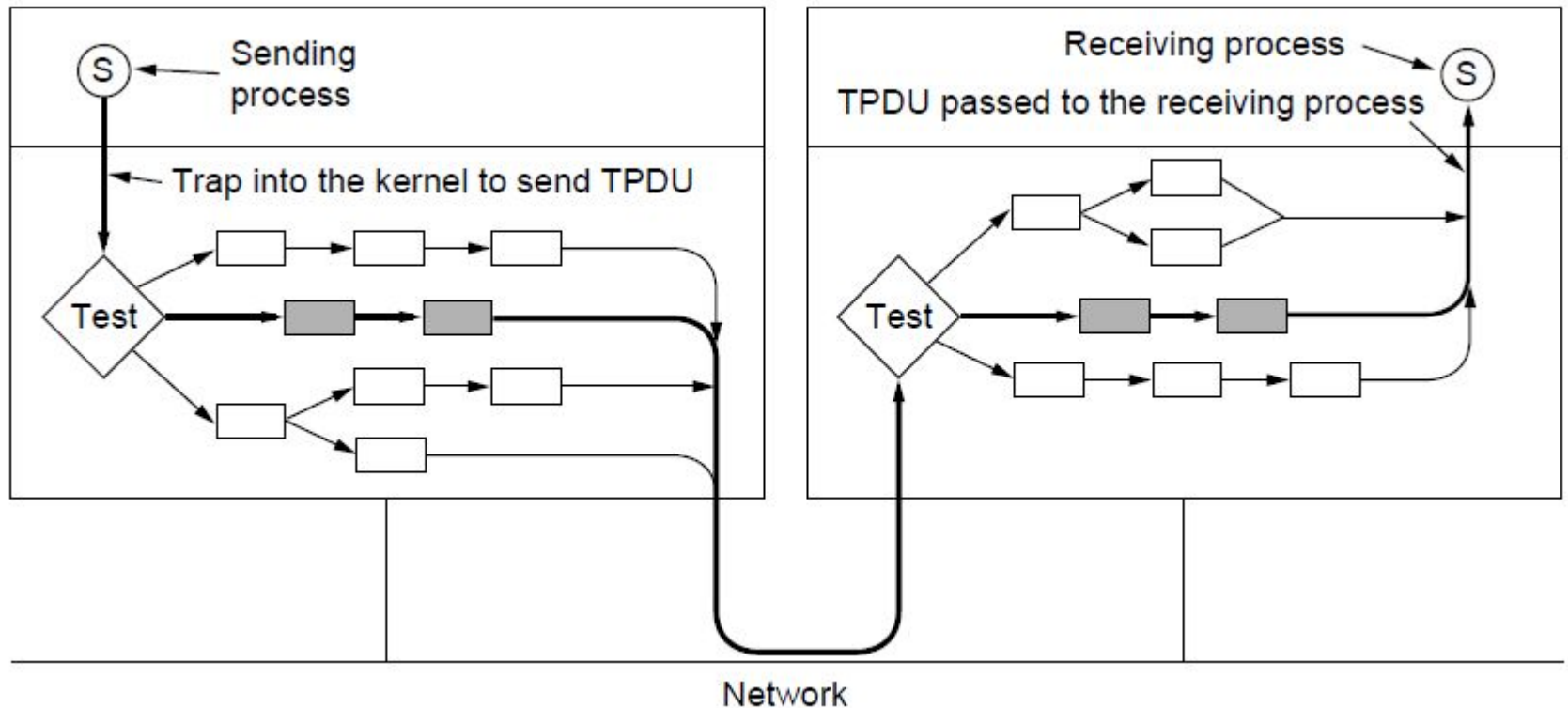
1. CPU speed more important than network speed
2. Reduce packet count to reduce software overhead
3. Minimize data touching
4. Minimize context switches
5. Minimize copying
6. You can buy more bandwidth but not lower delay
7. Avoiding congestion is better than recovering from it
8. Avoid timeouts

System Design for Better Performance (2)



Four context switches to handle one packet
with a user-space network manager.

Fast TPDU Processing (1)



The fast path from sender to receiver is shown with a heavy line. The processing steps on this path are shaded.

Fast TPDU Processing (2)

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

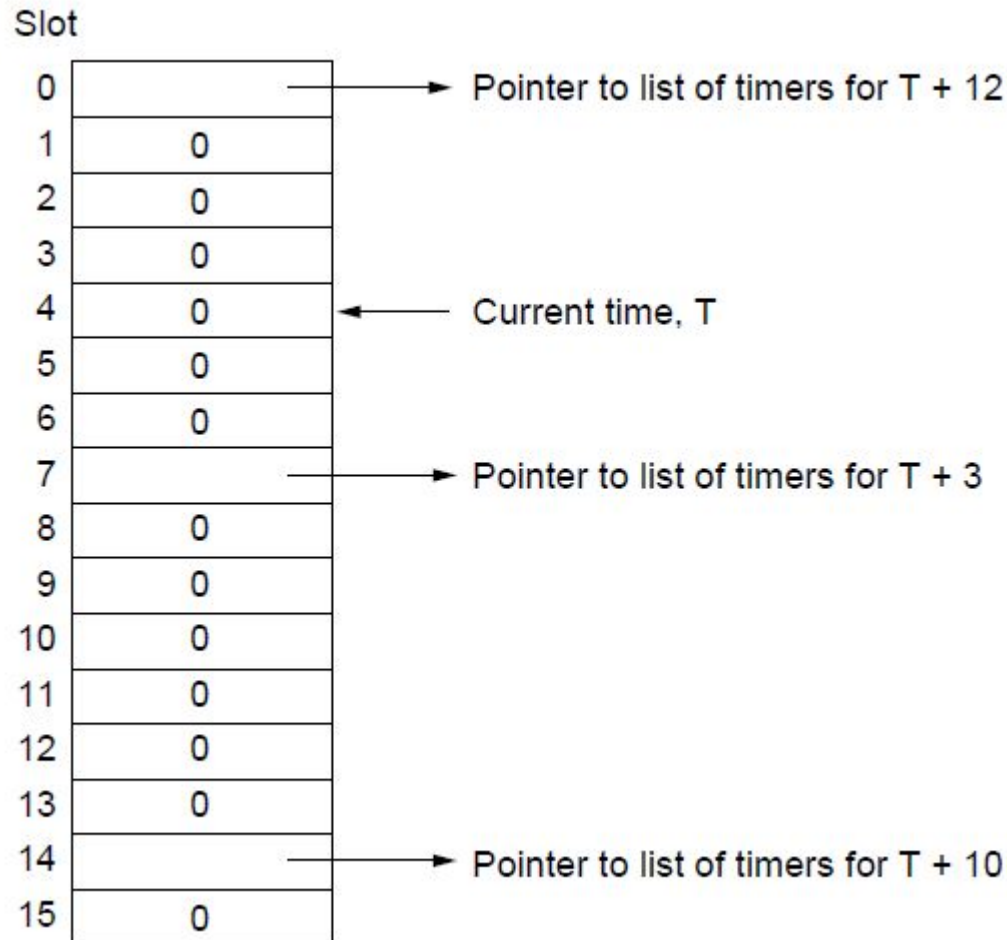
(a)

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

(b)

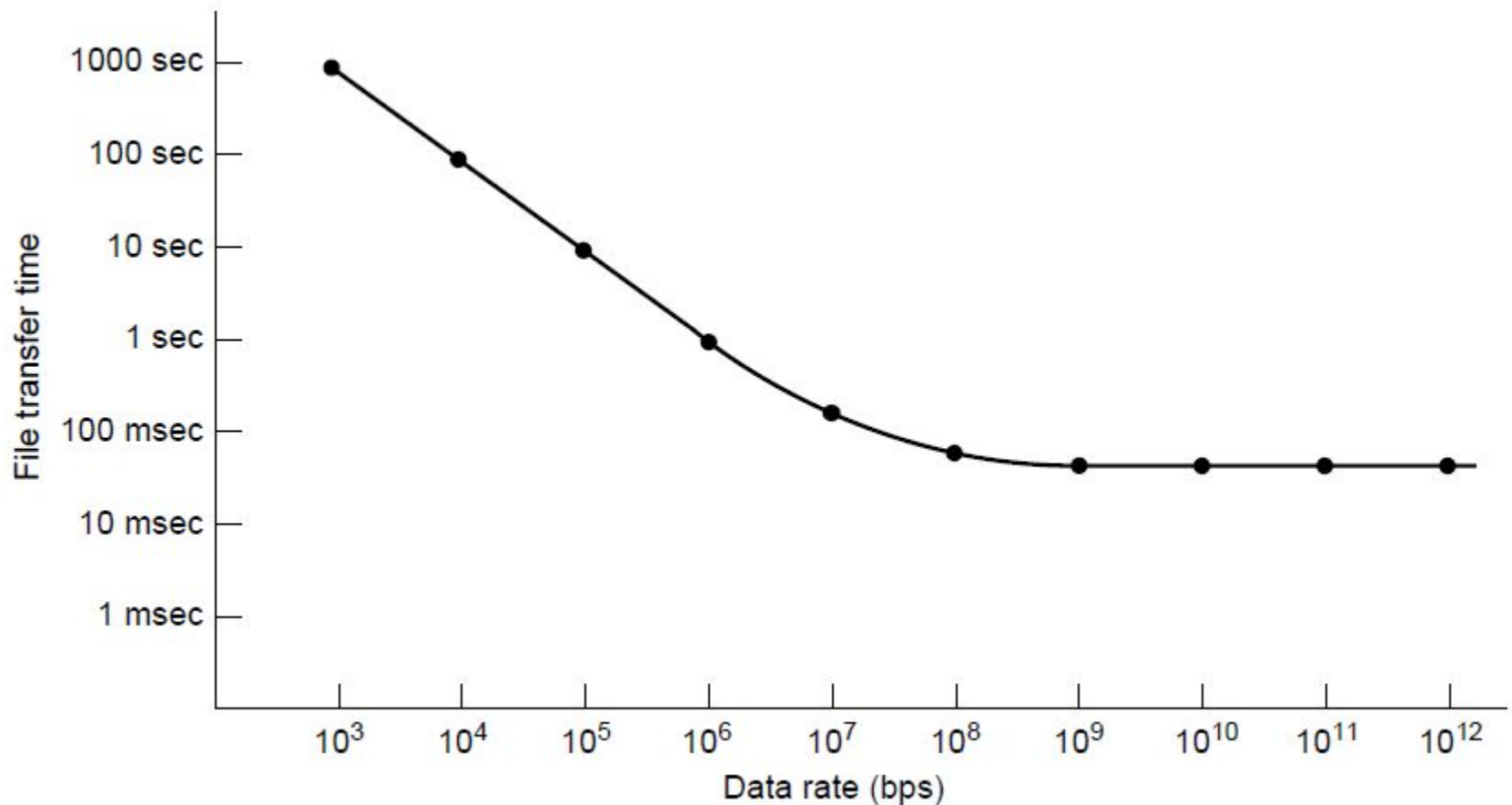
(a) TCP header. (b) IP header. In both cases, the shaded fields are taken from the prototype without change.

Protocols for High-Speed Networks (1)



A timing wheel

Protocols for High-Speed Networks (2)

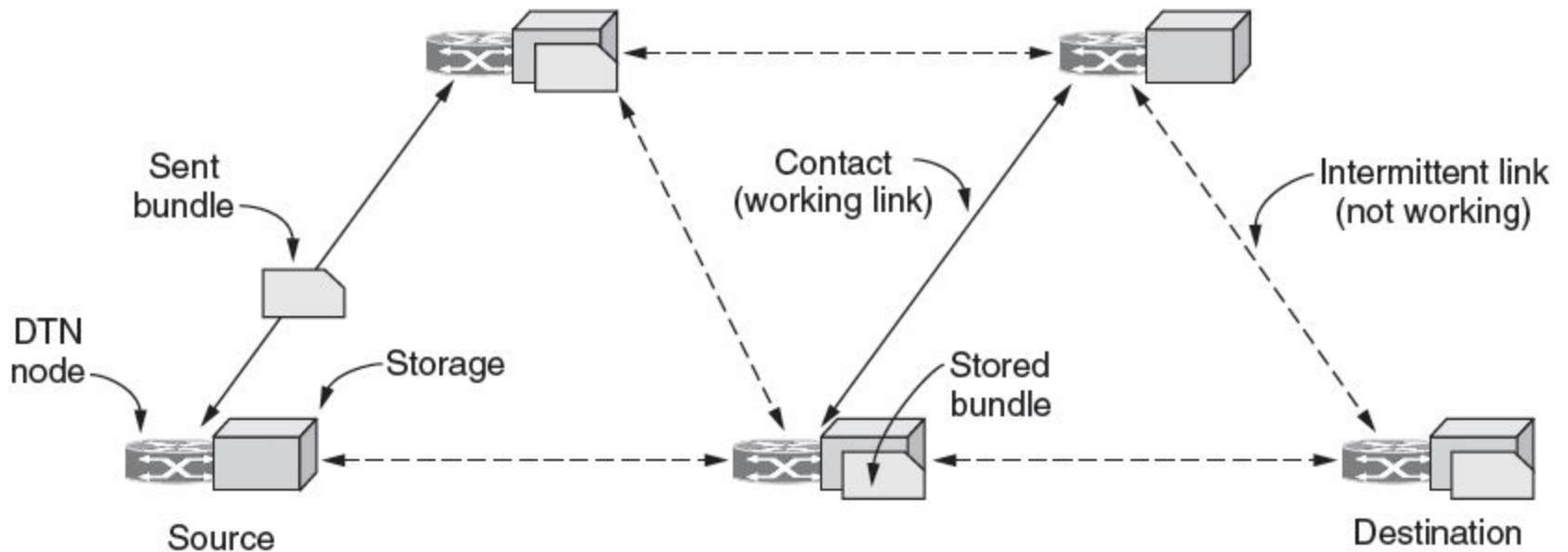


Time to transfer and acknowledge a
1-megabit file over a 4000-km line

Delay Tolerant Networking

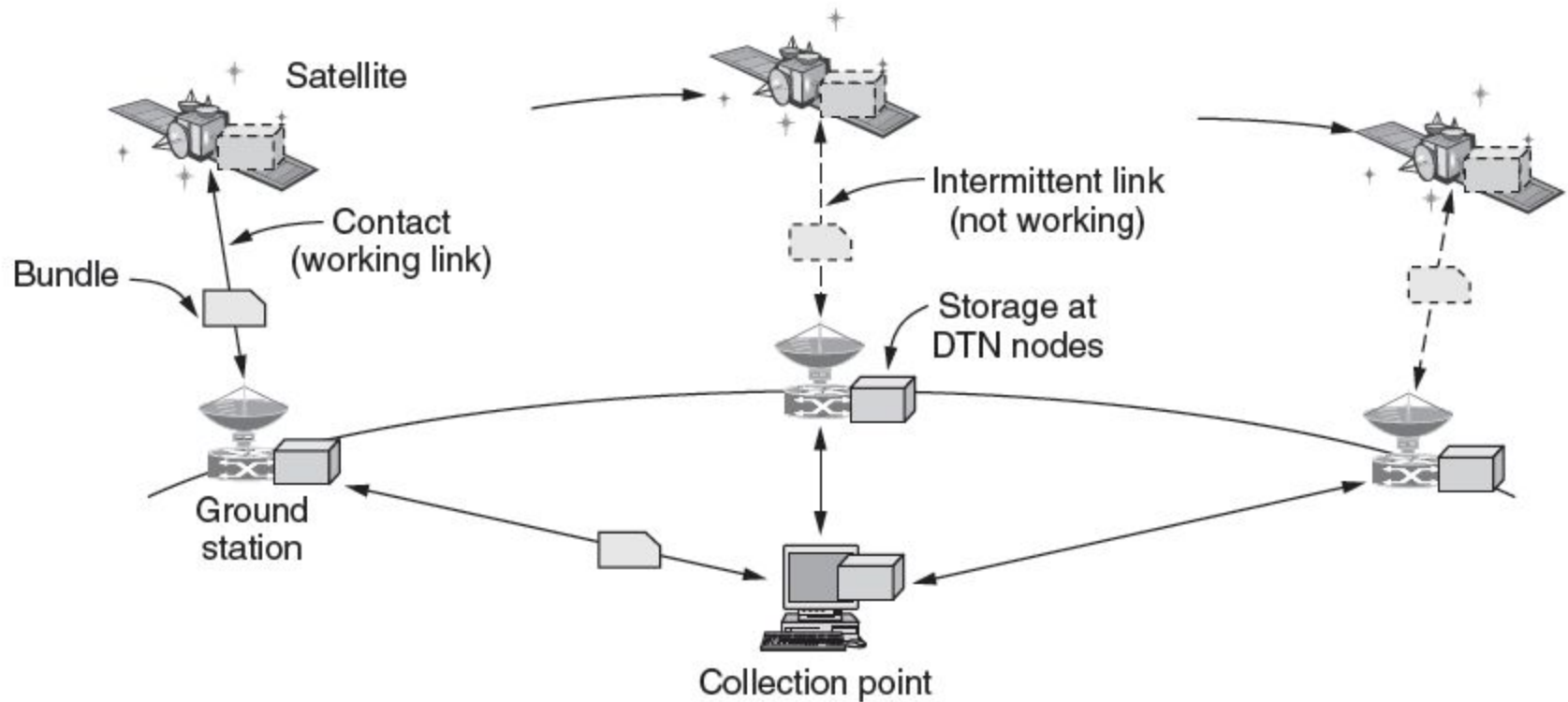
- DTN Architecture
- The Bundle Protocol

DTN Architecture (1)



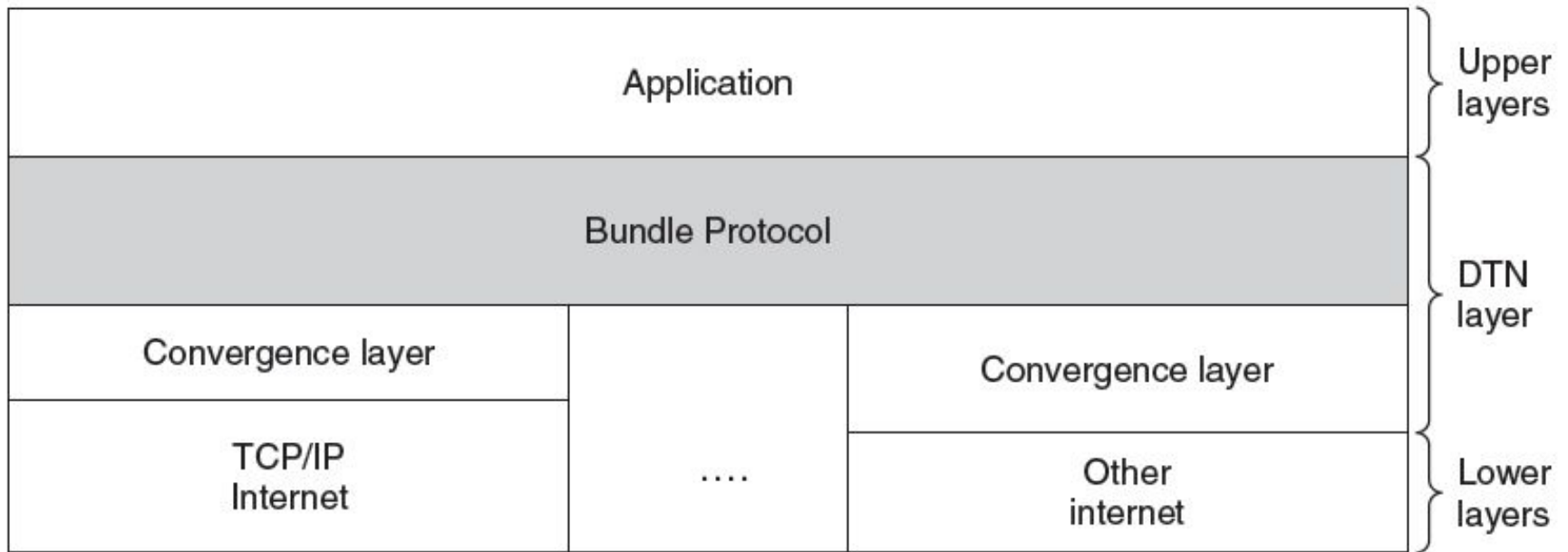
Delay-tolerant networking architecture

DTN Architecture (2)



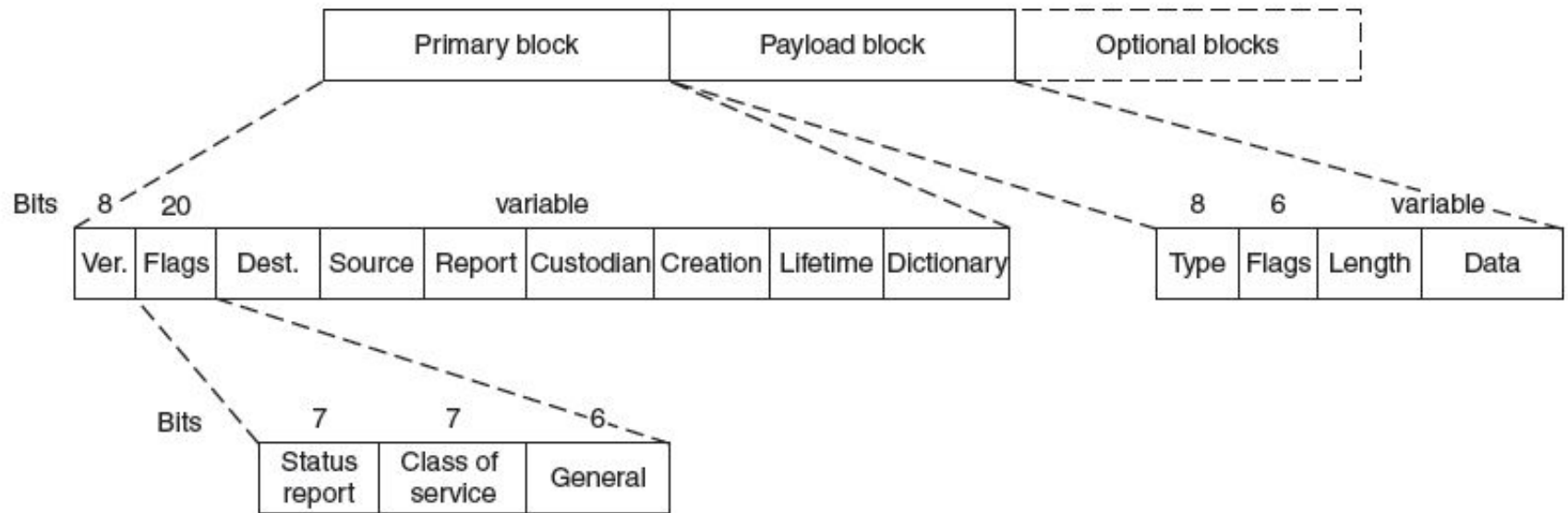
Use of a DTN in space.

The Bundle Protocol (1)



Delay-tolerant networking protocol stack.

The Bundle Protocol (2)



Bundle protocol message format.

End

Chapter 6