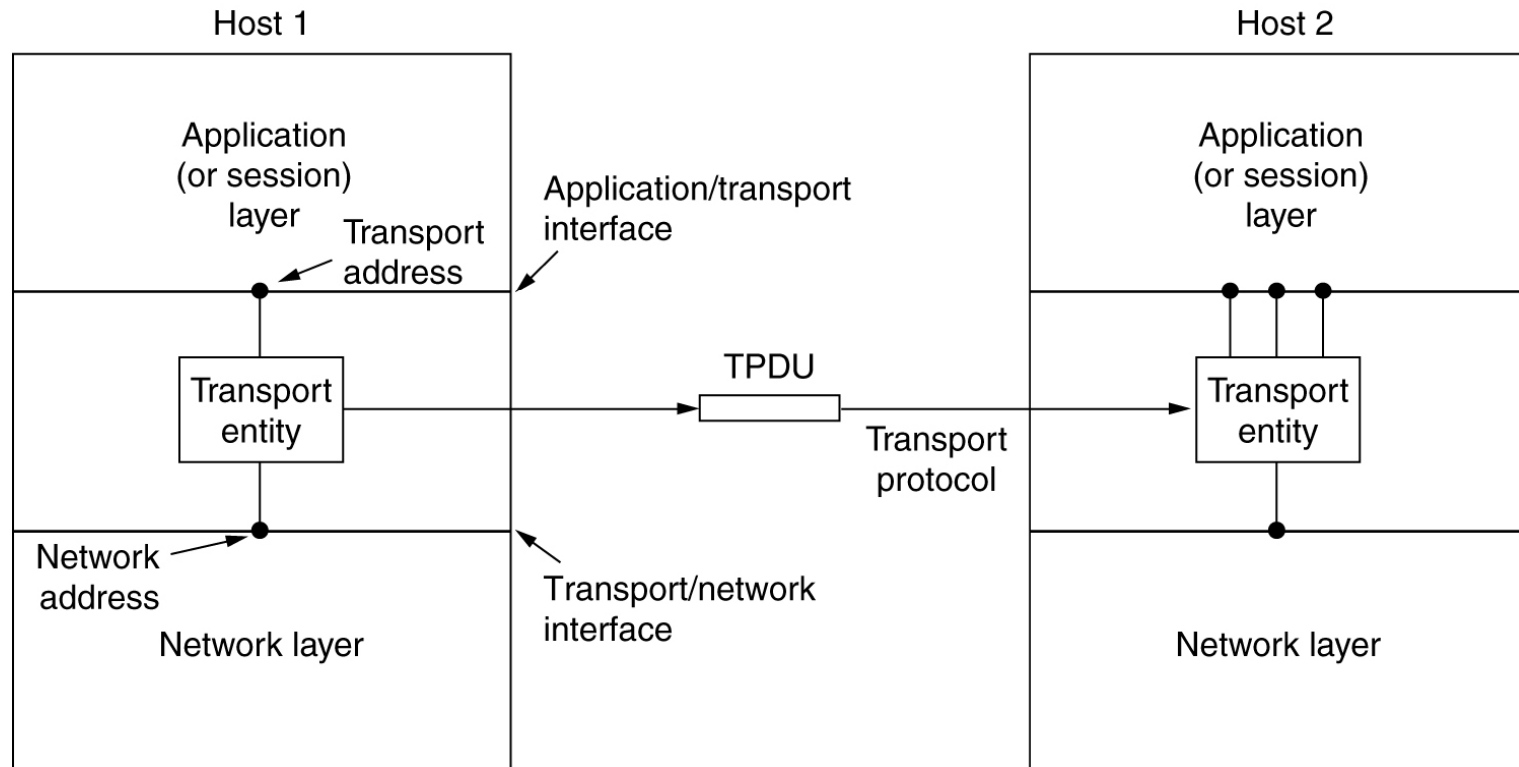# Chapter 6

# The Transport Layer

# The Transport Service

- Services Provided to the Upper Layers

- Transport Service Primitives

- Berkeley Sockets

- An Example of Socket Programming:
  - An Internet File Server
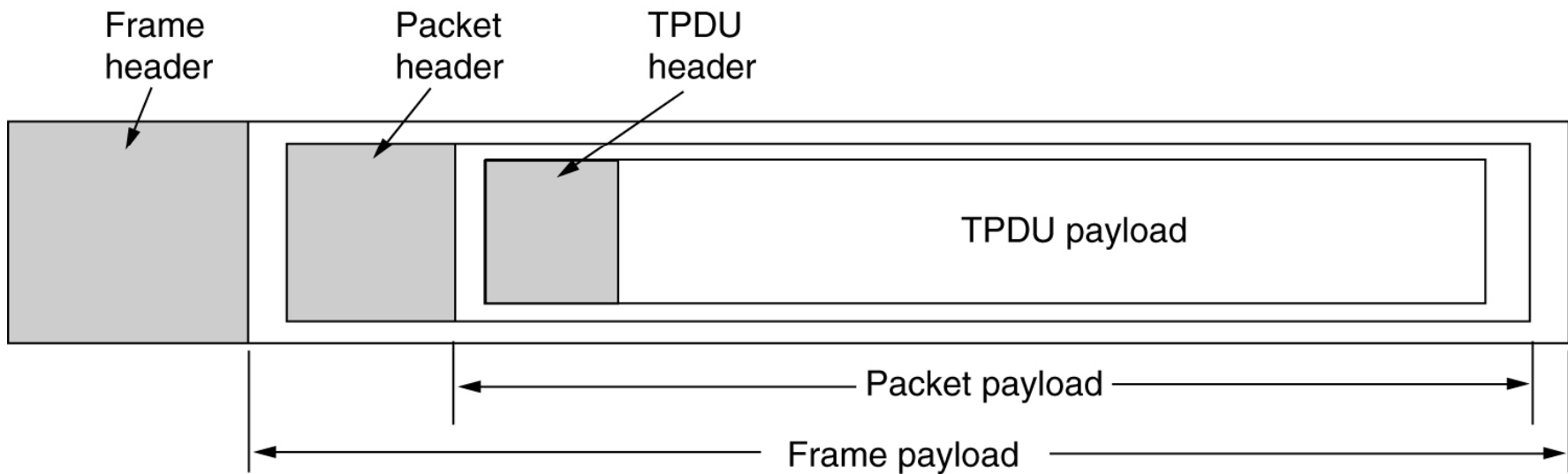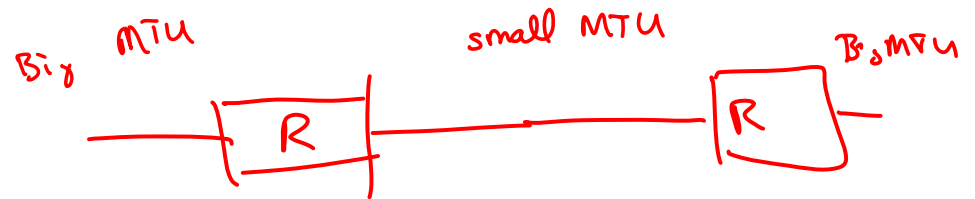
# Services Provided to the Upper Layers



The network, transport, and application layers.

# Transport Service Primitives

| 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.

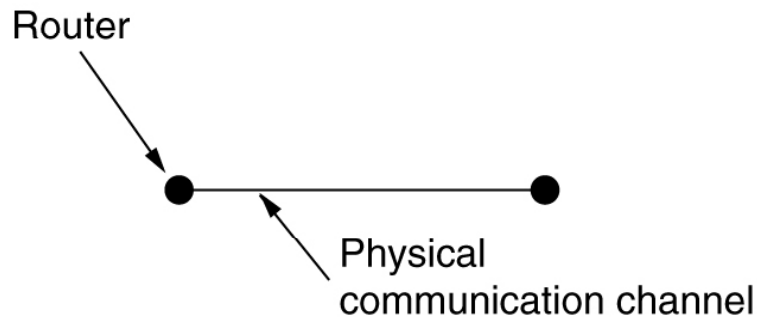# Transport Service Primitives (2)

Big MTU    small MTU    BoMTU



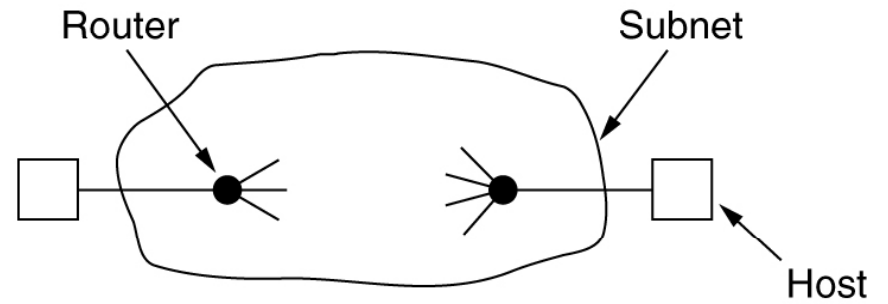The nesting of TPDUs, packets, and frames.

# Elements of Transport Protocols

- The environment in which Transport Operates

- Addressing
- Connection Establishment
- Connection Release
- Flow Control and Buffering
- Multiplexing
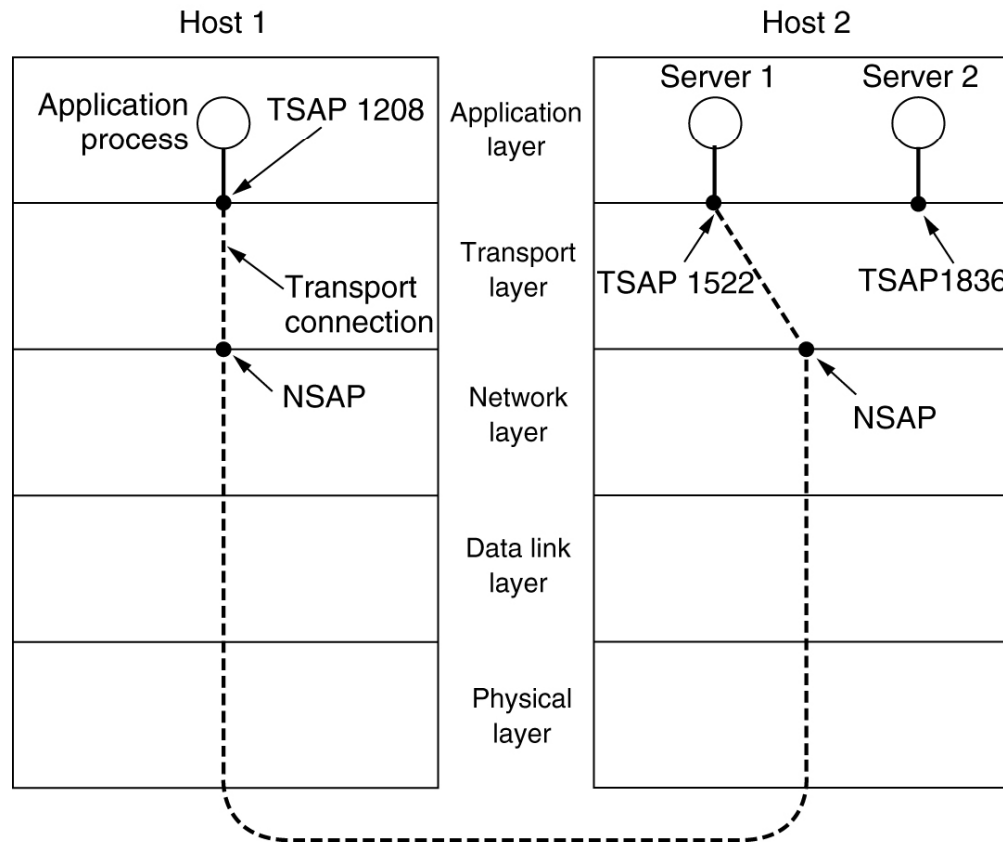- Crash Recovery

# Transport Protocol



(a)

(b)

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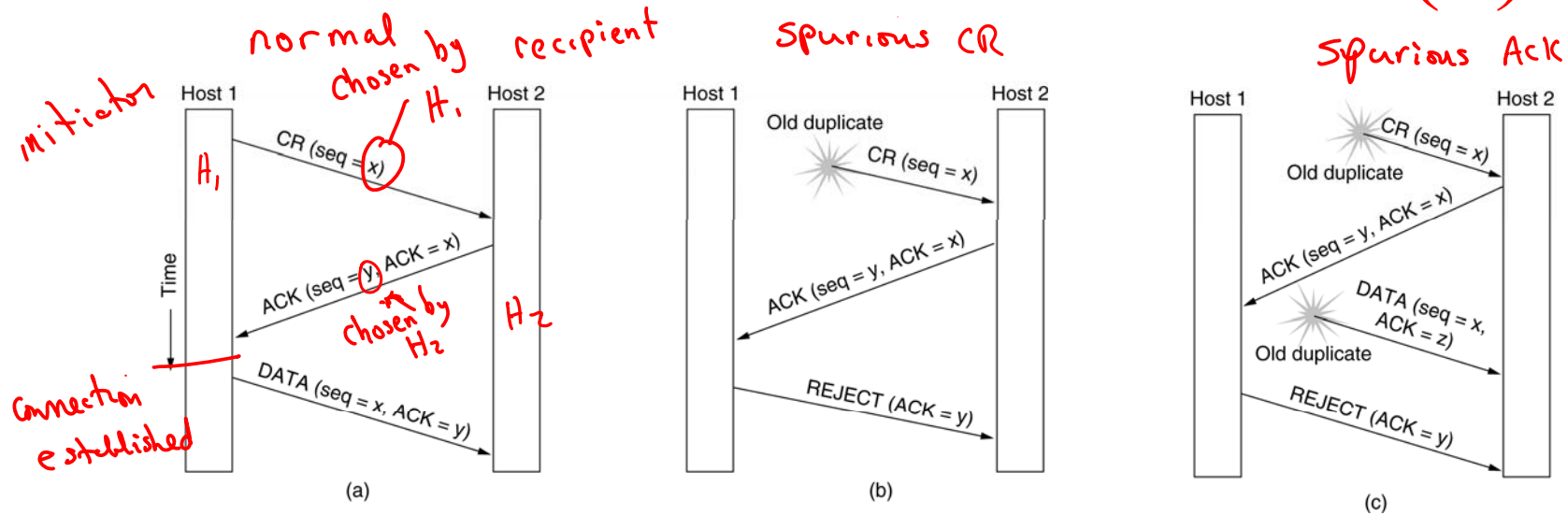(b) Environment of the transport layer.

# Addressing



TSAPs, NSAPs and transport connections.

# Connection Establishment (3)

Three-way handshake



(a) Normal — chosen by recipient

Host 1 — Initiator ($H_1$)

CR (seq = x) — chosen by $H_1$

ACK (seq = y, ACK = x) — chosen by $H_2$

Connection established

DATA (seq = x, ACK = y)

(b) Spurious CR

Host 1 — Host 2

Old duplicate — CR (seq = x)

ACK (seq = y, ACK = x)

REJECT (ACK = y)

(c) Spurious Ack

Host 1 — Host 2

CR (seq = x)

Old duplicate

ACK (seq = y, ACK = x)
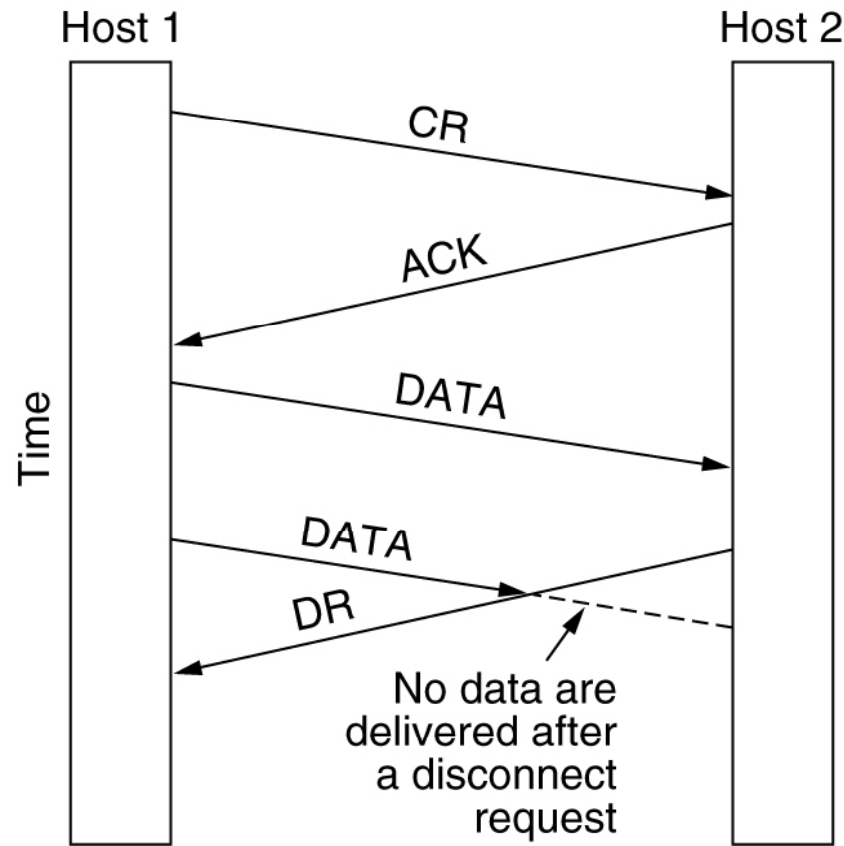
DATA (seq = x, ACK = z)

Old duplicate

REJECT (ACK = y)

Three protocol scenarios for establishing a connection using a three-way handshake. CR denotes CONNECTION REQUEST.
(a) Normal operation,
(b) Old CONNECTION REQUEST appearing out of nowhere.
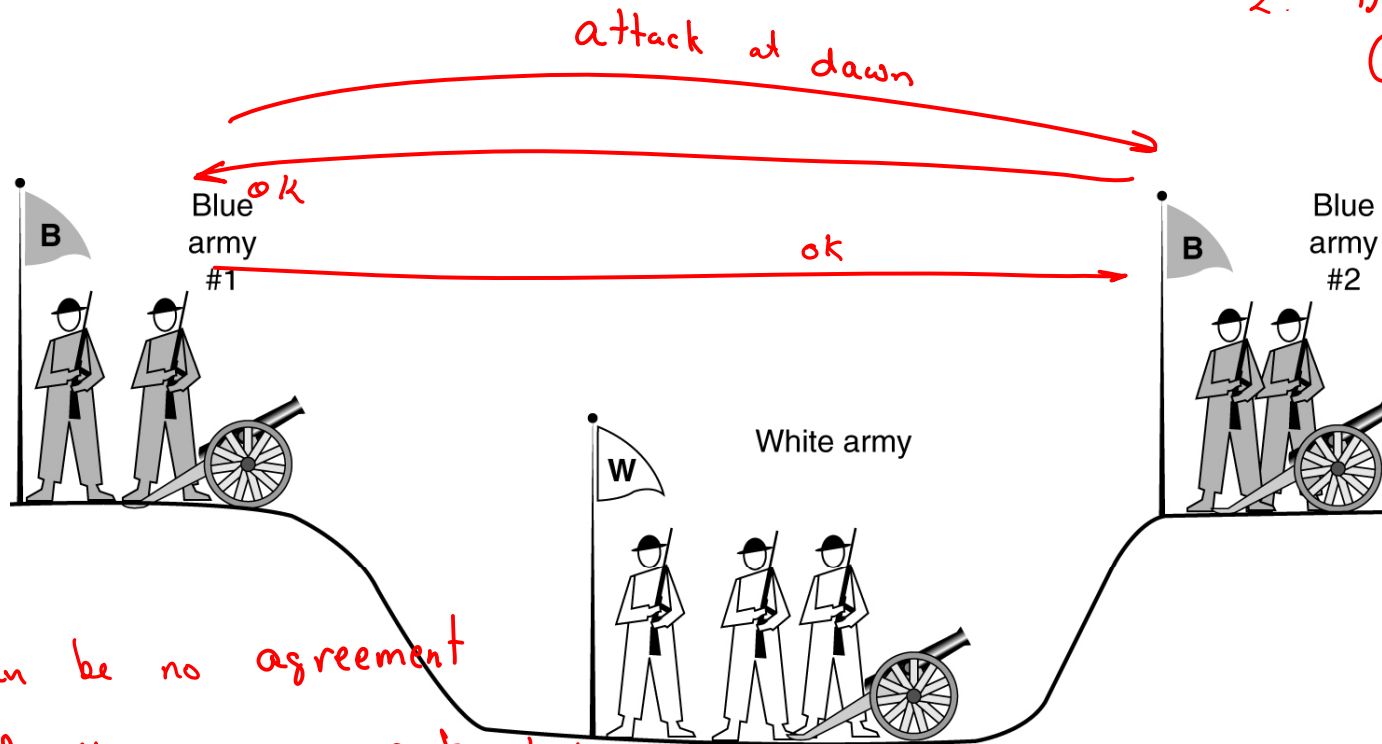(c) Duplicate CONNECTION REQUEST and duplicate ACK.

# Connection Release

Abrupt disconnection with loss of data.

1. Connection Release.

2. Distributed Commit of database transactions.

attack at dawn



Blue army #1

ok

ok

ok

Blue army #2

White army

There can be no agreement protocol if messages may be lost.
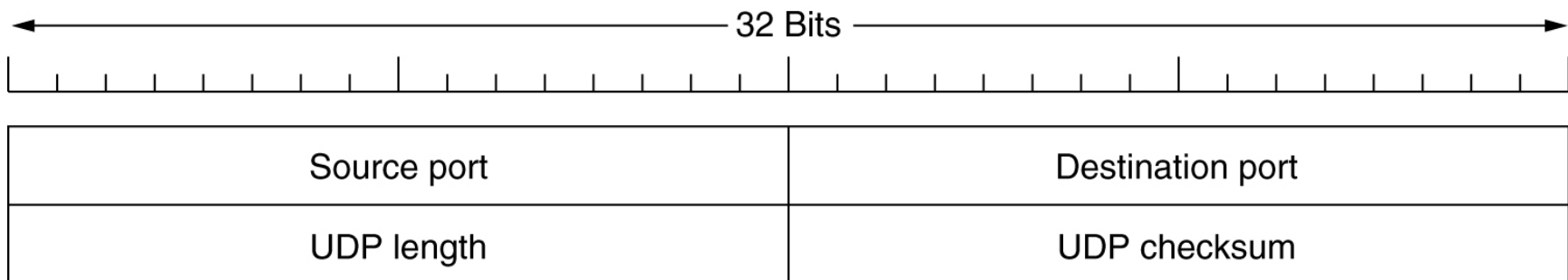
Proof idea:
message is essential if its receipt affects a party's behavior.

The two-army problem.

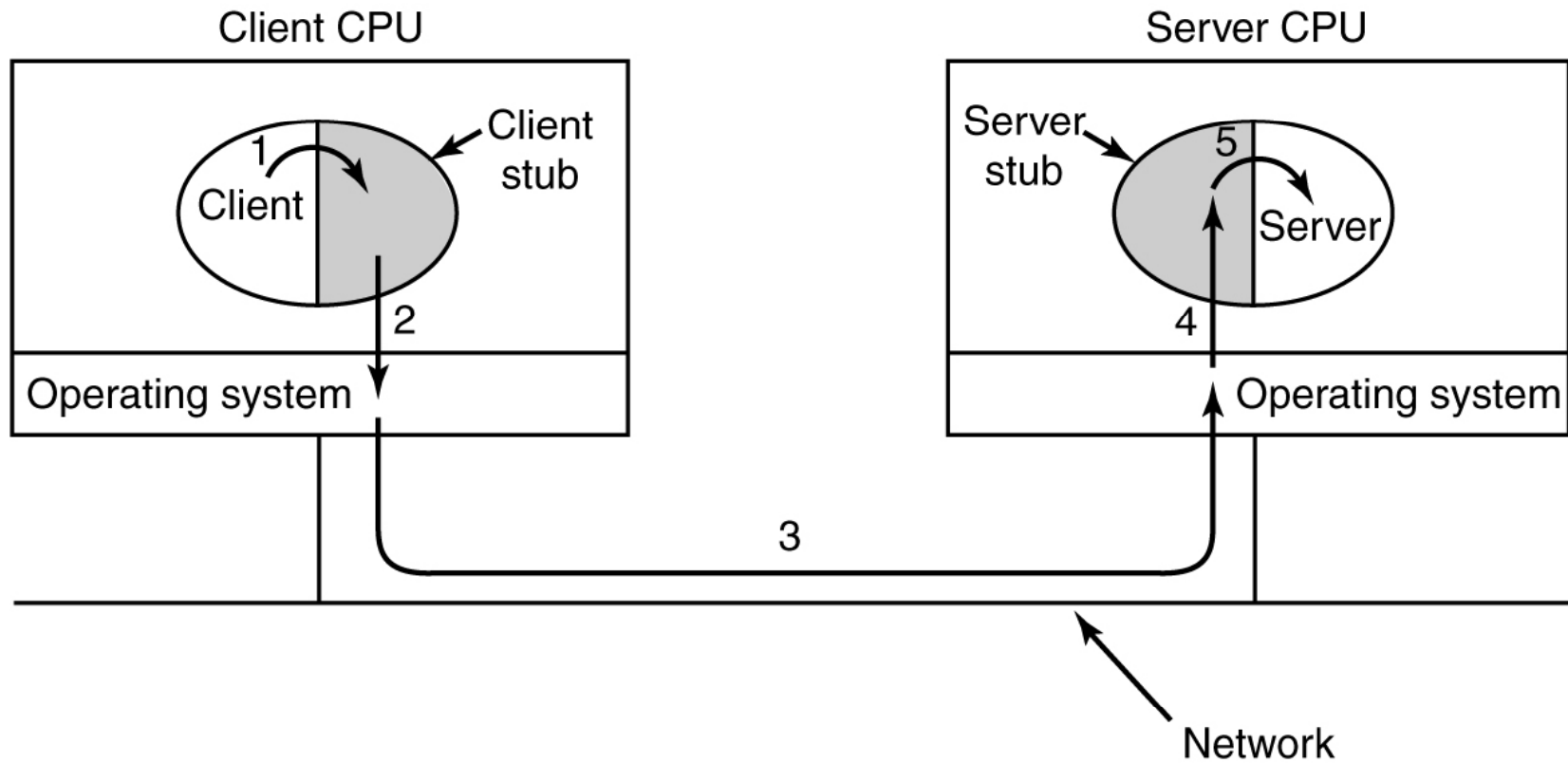# The Internet Transport Protocols: UDP

- Introduction to UDP

- Remote Procedure Call

- The Real-Time Transport Protocol

# Introduction to UDP



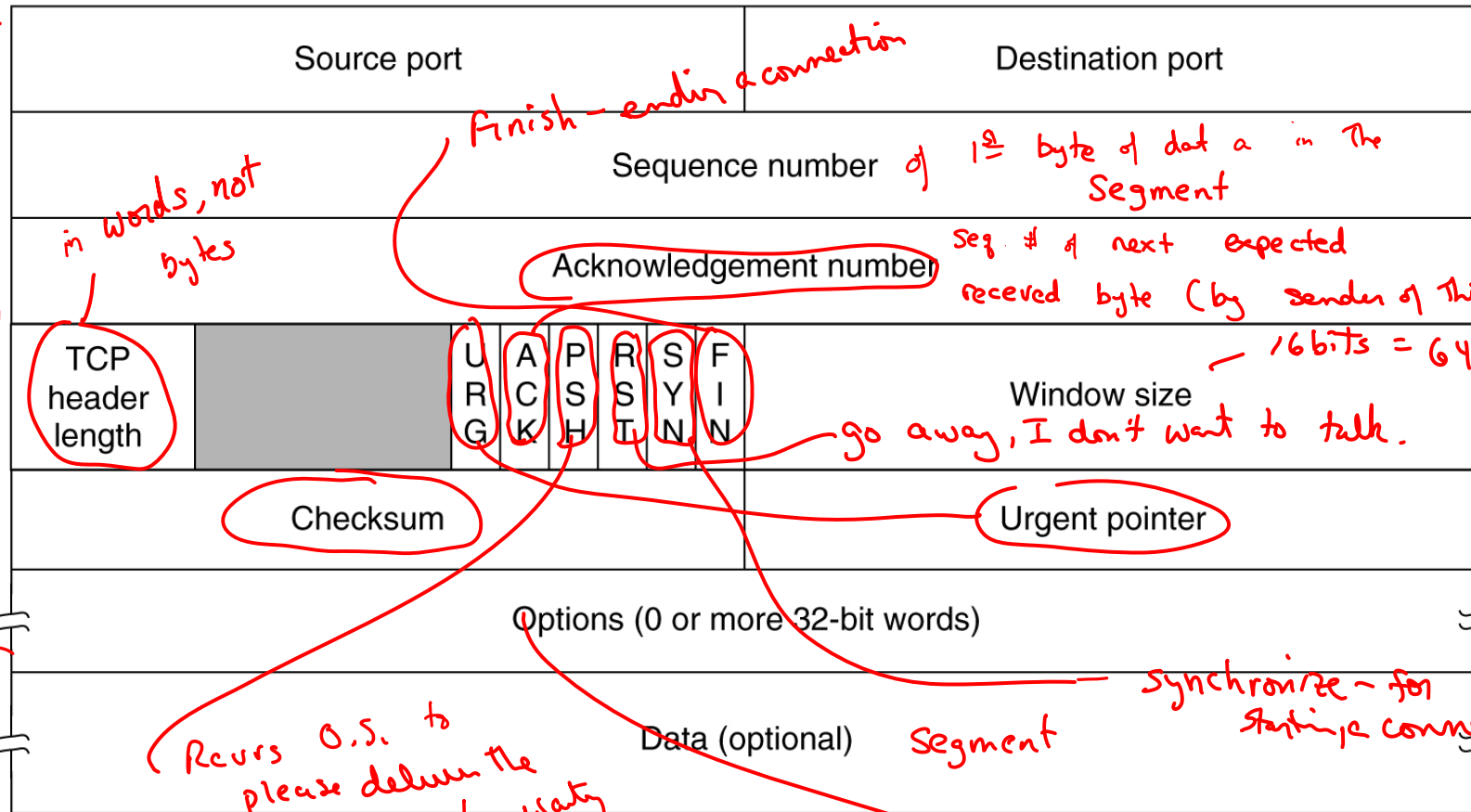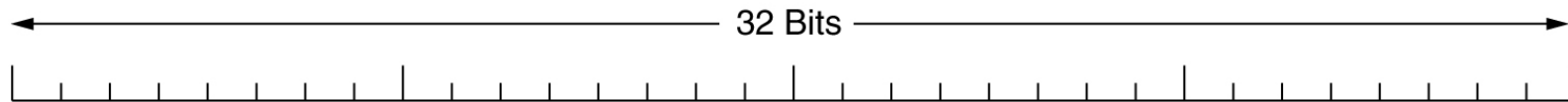| 32 Bits | |
|---|---|
| Source port | Destination port |
| UDP length | UDP checksum |

The UDP header.

# Remote Procedure Call



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

# The TCP Segment Header

TCP connection identified by (< src ip, src port>, <dest ip, dest port>)

←——————————— 32 Bits ———————————→

| Source port | Destination port |
|---|---|

Sequence number — of 1st byte of data in the Segment

Acknowledgement number — Seq # of next expected received byte (by sender of this seg).

TCP header length — in words, not bytes

| U R G | A C K | P S H | R S T | S Y N | F I N | Window size — 16 bits = 64 Kbyte |

Finish — ending a connection

go away, I don't want to talk.

| Checksum | Urgent pointer |

Options (0 or more 32-bit words)
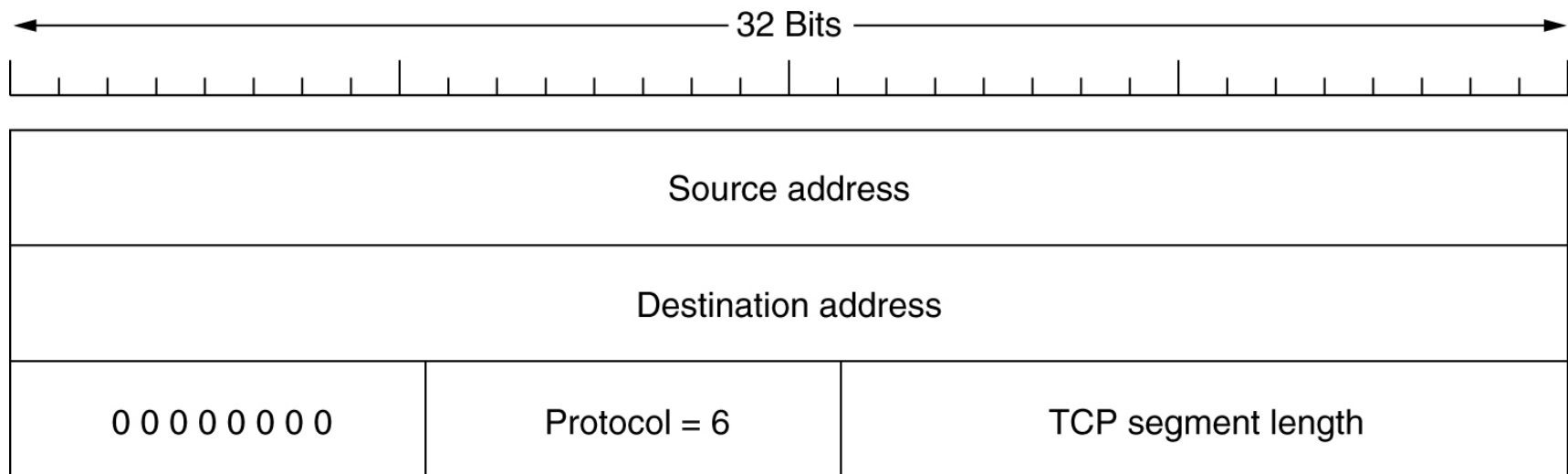
Data (optional)

Rcvrs O.S. to please deliver the data w/o waiting for more.

Segment

Synchronize — for Starting connection

Selective acks
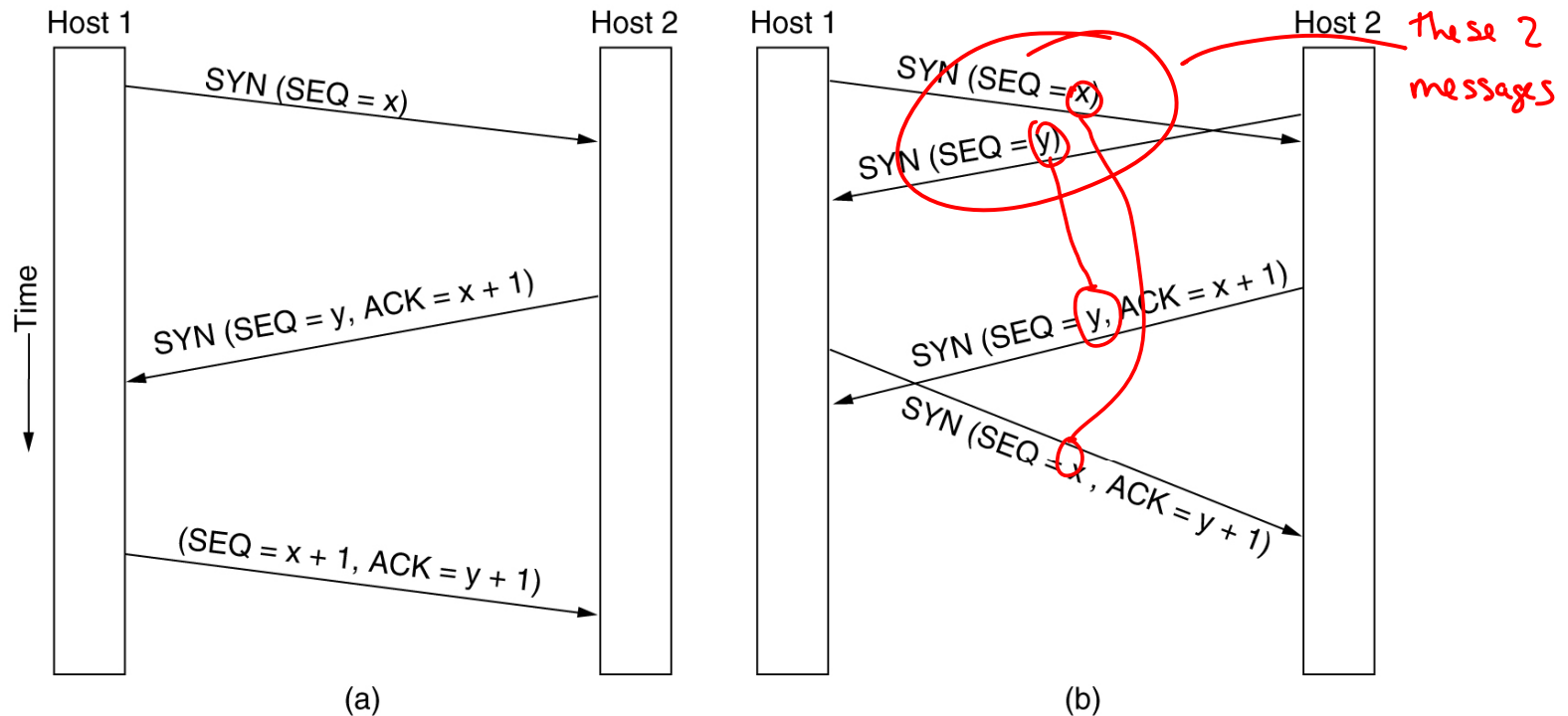
Window Scaling.

## TCP Header.

# The TCP Segment Header (2)

← ————————————— 32 Bits ————————————— →

| Source address |
|:---:|
| Destination address |

| 0 0 0 0 0 0 0 0 | Protocol = 6 | TCP segment length |
|:---:|:---:|:---:|

*violates layering*

The ~~IP~~ pseudoheader included in the TCP checksum.

# TCP Connection Establishment

only an issue of $H_1$ port and $H_2$ port are identical in these 2 messages

3-way handshake.

End up w/ 1 connection betw $H_1$ port and $H_2$ port.



(a) TCP connection establishment in the normal case.
(b) Call collision.

# TCP Connection Management Modeling

| 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 |
| TIMED 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 |

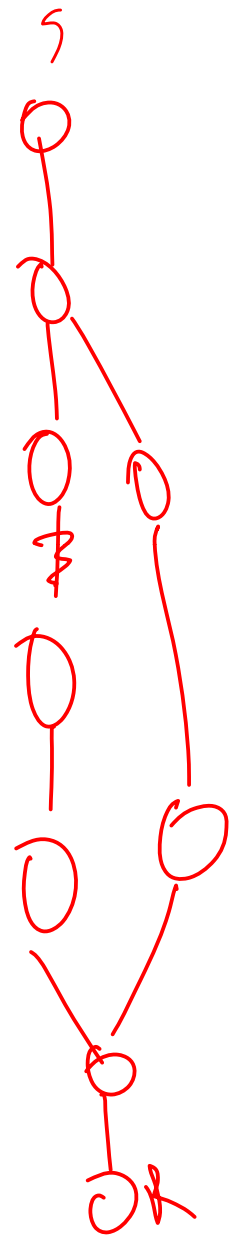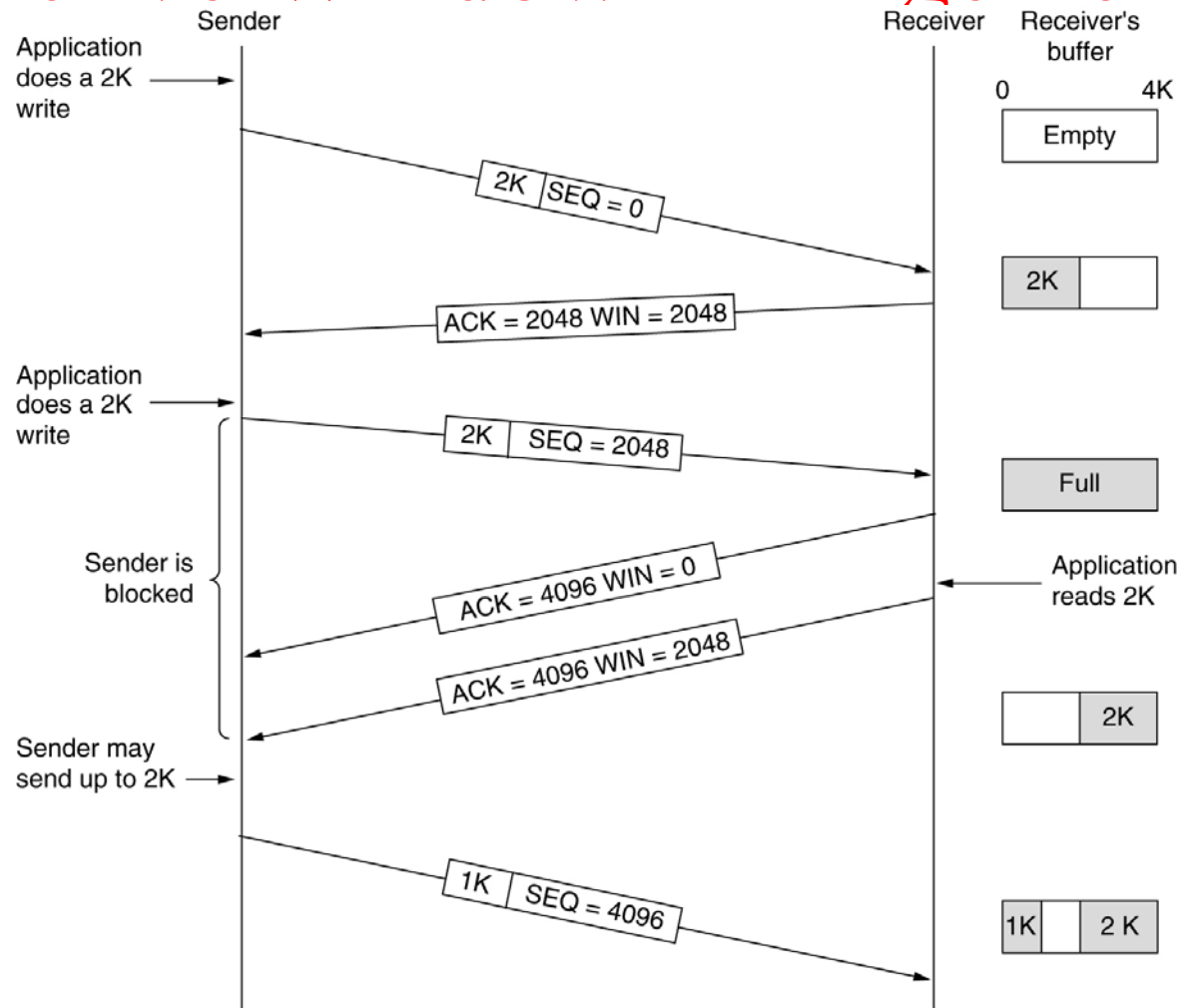The states used in the TCP connection management finite state machine.

# 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.
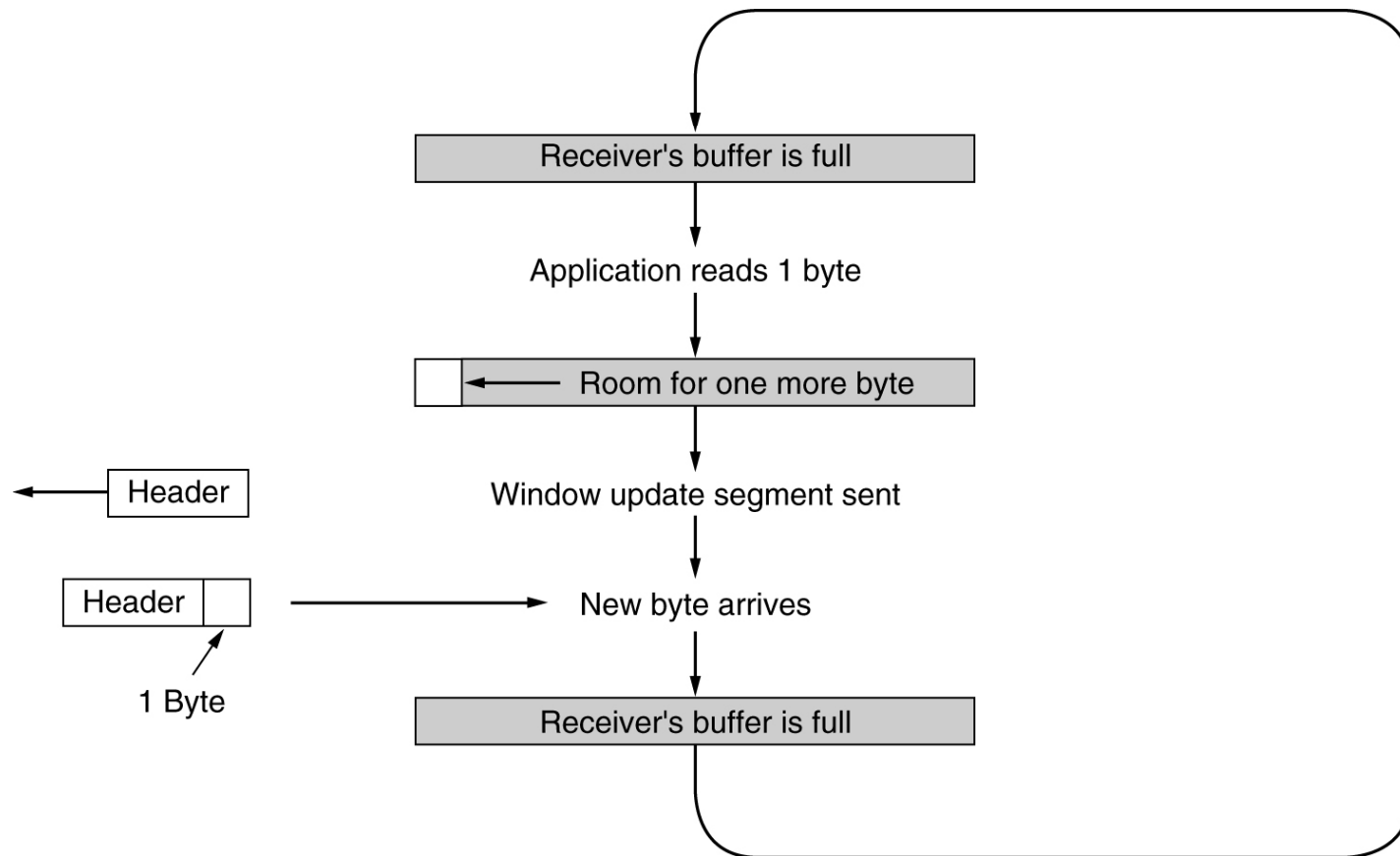
(Start)

CONNECT/SYN (Step 1 of the 3-way handshake)

CLOSED

CLOSE/–

LISTEN/–    CLOSE/–

SYN/SYN + ACK

(Step 2   of the 3-way handshake)    LISTEN

SYN RCVD    RST/–    SEND/SYN    SYN SENT

SYN/SYN + ACK    (simultaneous open)

(Data transfer state)

ACK/–    ESTABLISHED    SYN + ACK/ACK
(Step 3 of the 3-way handshake)

CLOSE/FIN

CLOSE/FIN    FIN/ACK

(Active close)    (Passive close)

FIN/ACK

FIN WAIT 1    CLOSING    CLOSE WAIT

ACK/–    ACK/–    CLOSE/FIN

FIN WAIT 2    FIN + ACK/ACK    TIME WAIT    LAST ACK

FIN/ACK

(Timeout/)

CLOSED    ACK/–

(Go back to start)

# TCP Transmission Policy
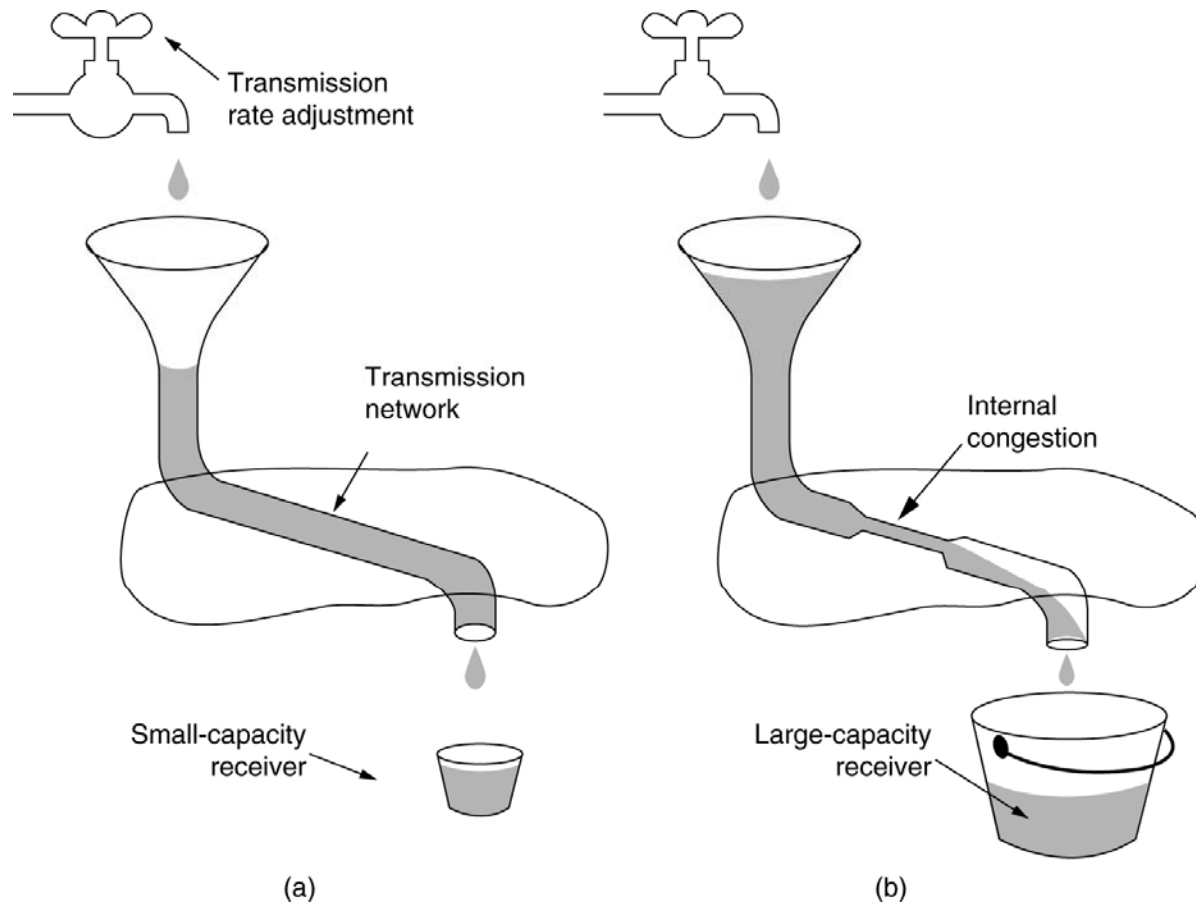# Active Window Management



Window management in TCP.

# TCP Transmission Policy (2)



Solving the silly window syndrome
Nagle's algorithm for transmission
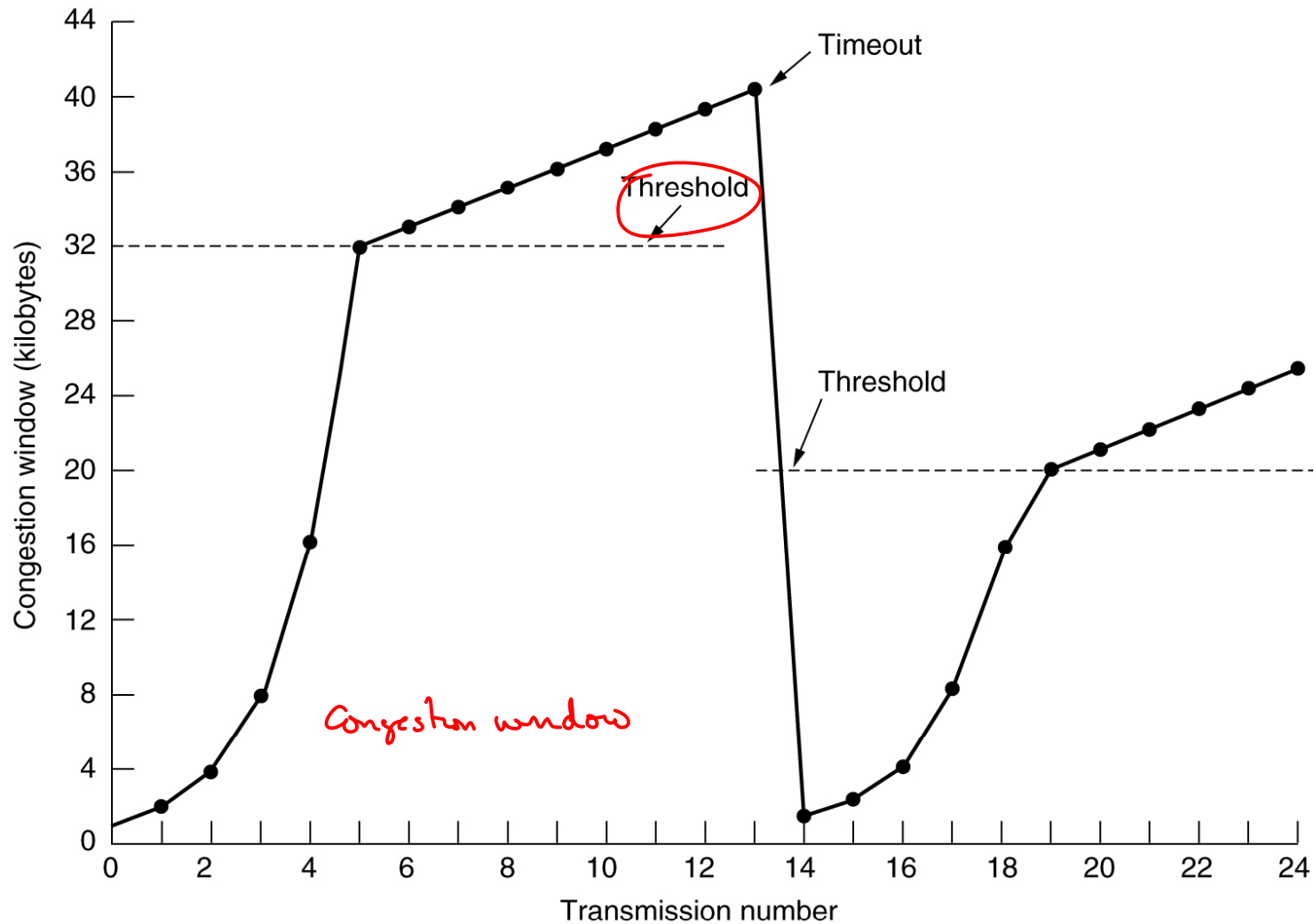
# TCP Congestion Control



(a) A fast network feeding a low capacity receiver.
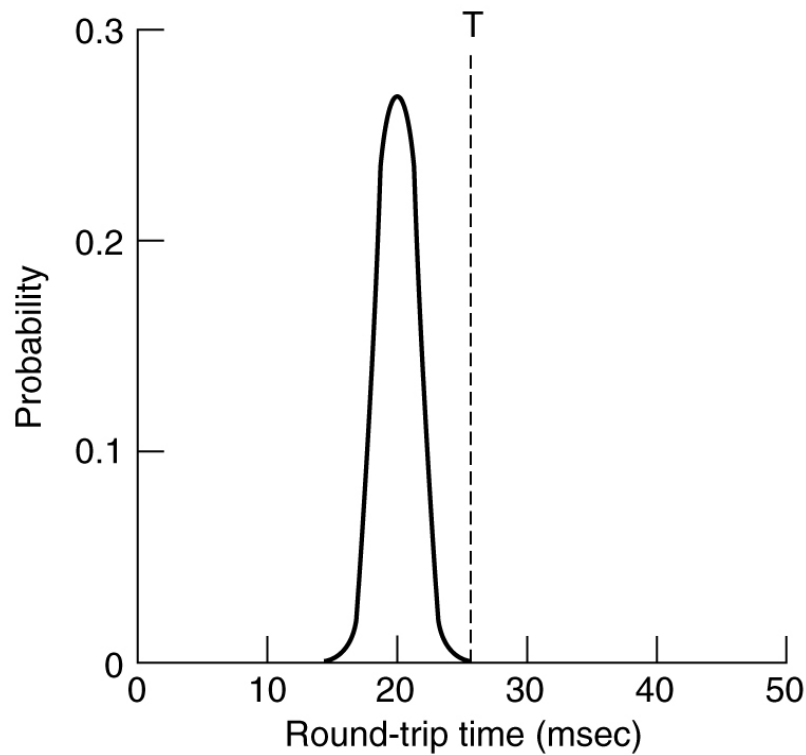(b) A slow network feeding a high-capacity receiver.
Outstanding data must be limited by both the network AND the receiver
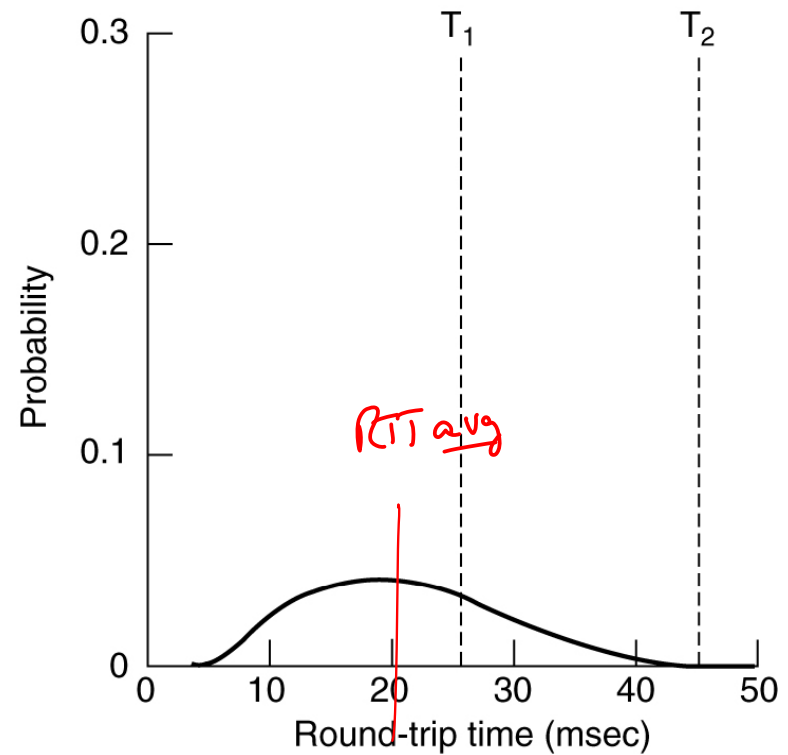
# TCP Congestion Control (2)



An example of the Internet congestion algorithm
TCP Slow-start (Slow wrt to what since it is exponential!?)

# TCP Timer Management



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

# Round-trip and Variance Estimation

$$0 \leq a \leq 1.0 \qquad a_{typical} = 7/8$$

(a)   RTT = a(RTT) + (1-a)M  (exponential smoothing)

(b)   Dev = a(Dev) + (1-a) |RTT-M|

(c)   Retransmission Timeout = RTT + 4*Dev

    (a)   Used to use RTO = 2*RTT

(d)   What about retransmitted segments?

    (a)   Ignore them in the RTT calculation
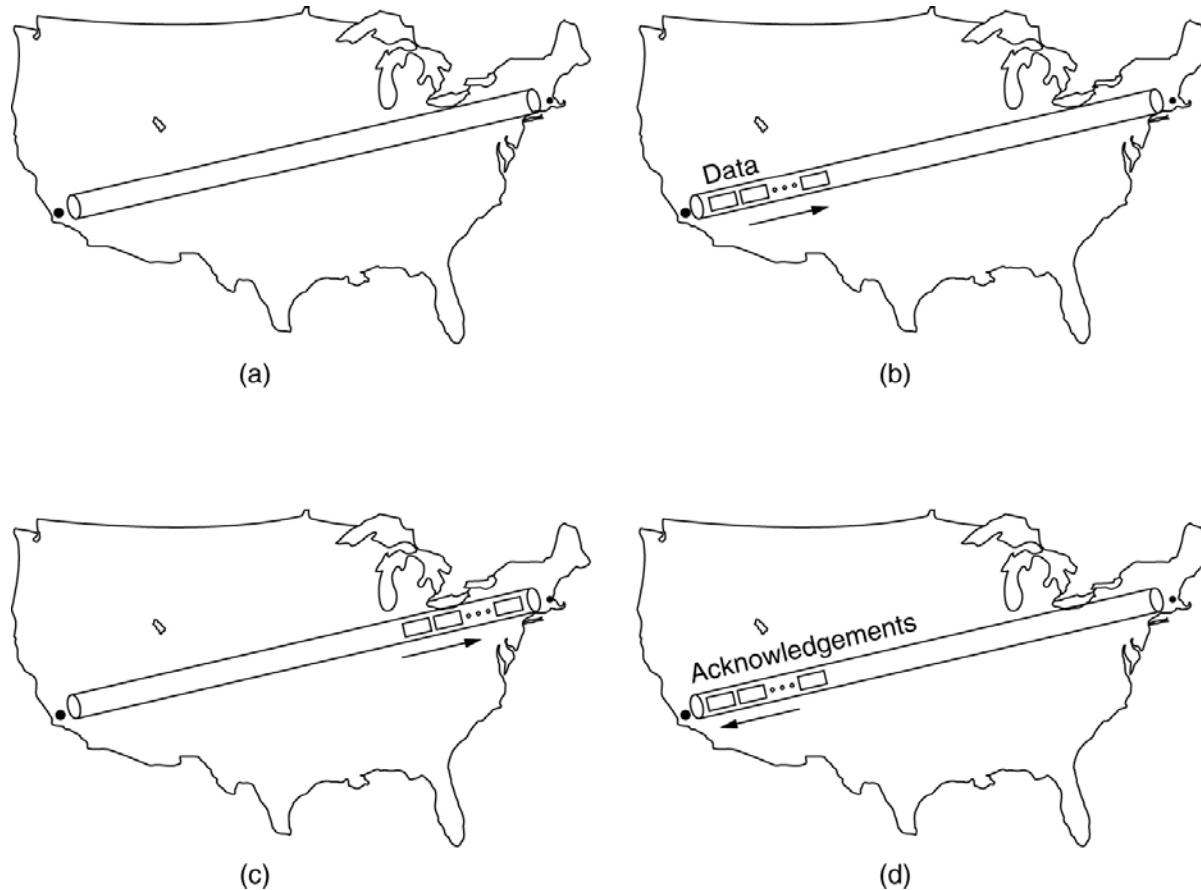
# Wireless TCP

- Missing acks can mean two things with opposite required responses
    - In wired networks: congestion; response: slow down!
    - In wireless networks: dropped packets due to noise; try again, real soon!
- Approaches for wireless
    - Split TCP connection in two (not common)
    - Acks and retransmission at the link layer

# Performance Issues

- Performance Problems in Computer Networks

- Network Performance Measurement

- System Design for Better Performance

- Fast TPDU Processing

- Protocols for Gigabit Networks

# Performance Problems in Computer Networks



The state of transmitting one megabit from San Diego to Boston
(a) At t = 0,   (b) After 500 µsec, (c) After 20 msec,  (d) after 40 msec.

# Network Performance Measurement

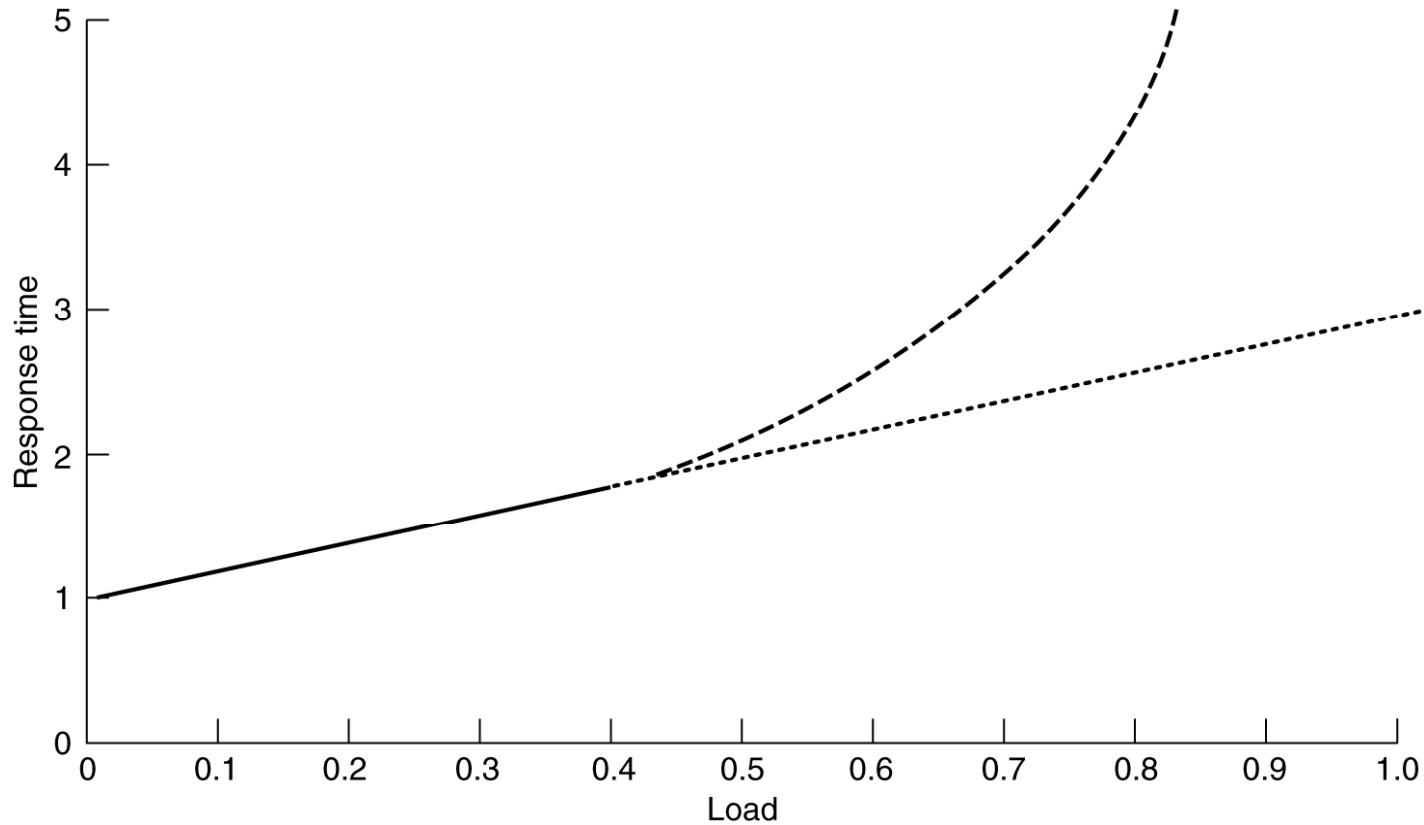The basic loop for improving network performance.

1. Measure relevant network parameters, performance.
2. Try to understand what is going on.
3. Change one parameter.
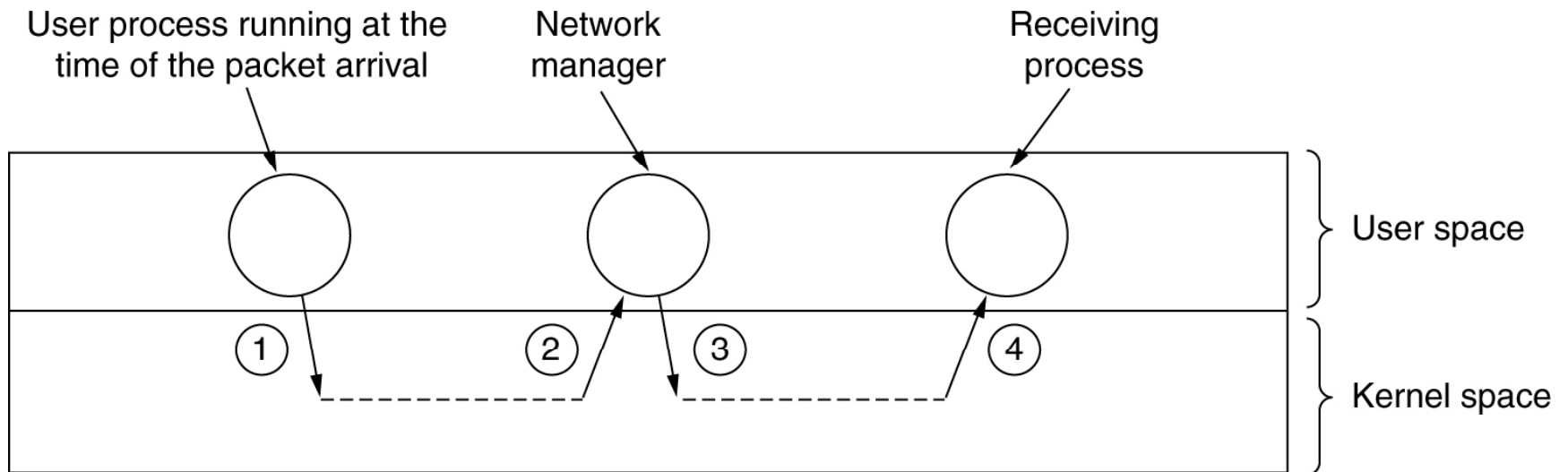
# System Design for Better Performance

Rules:

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

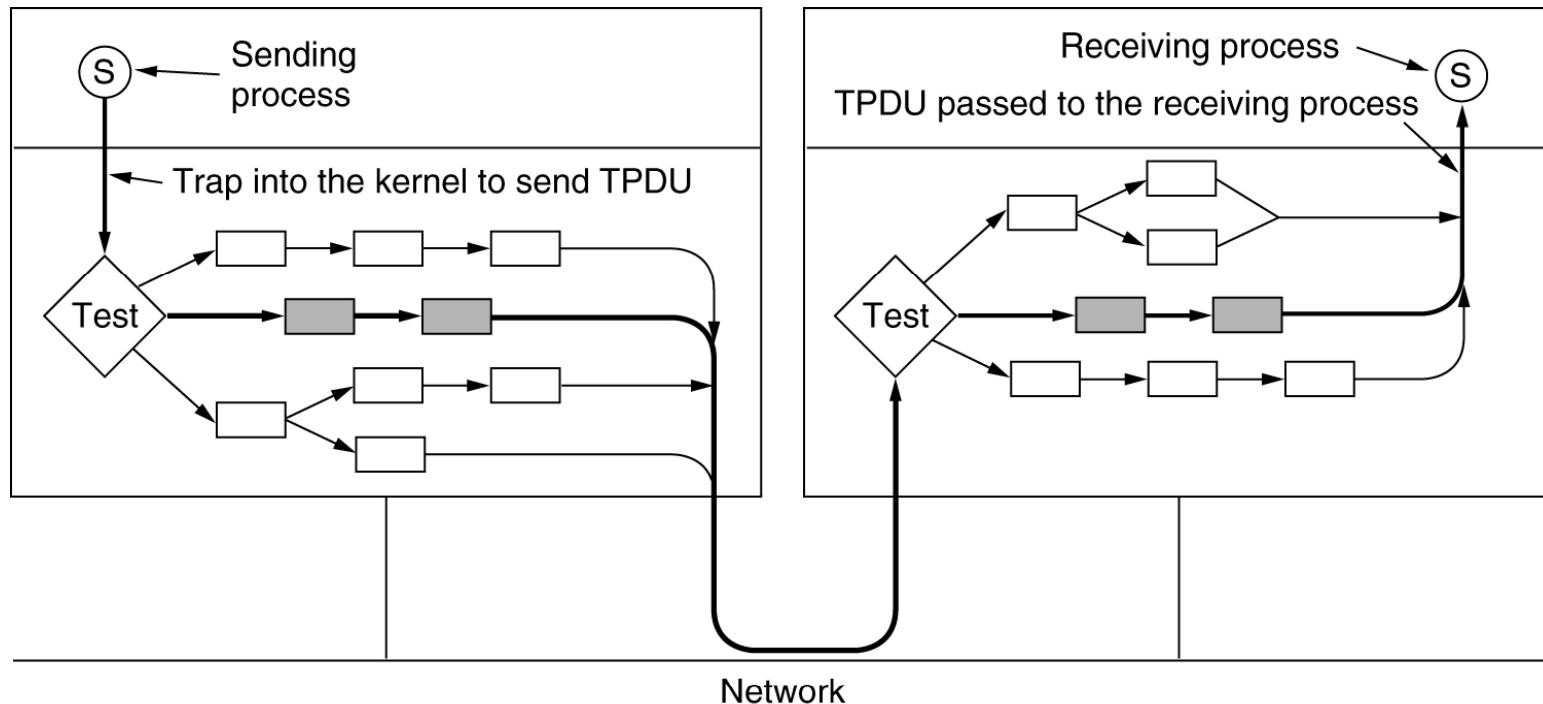# System Design for Better Performance (2)



Response as a function of load.

# System Design for Better Performance (3)



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

# Fast TPDU Processing



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)



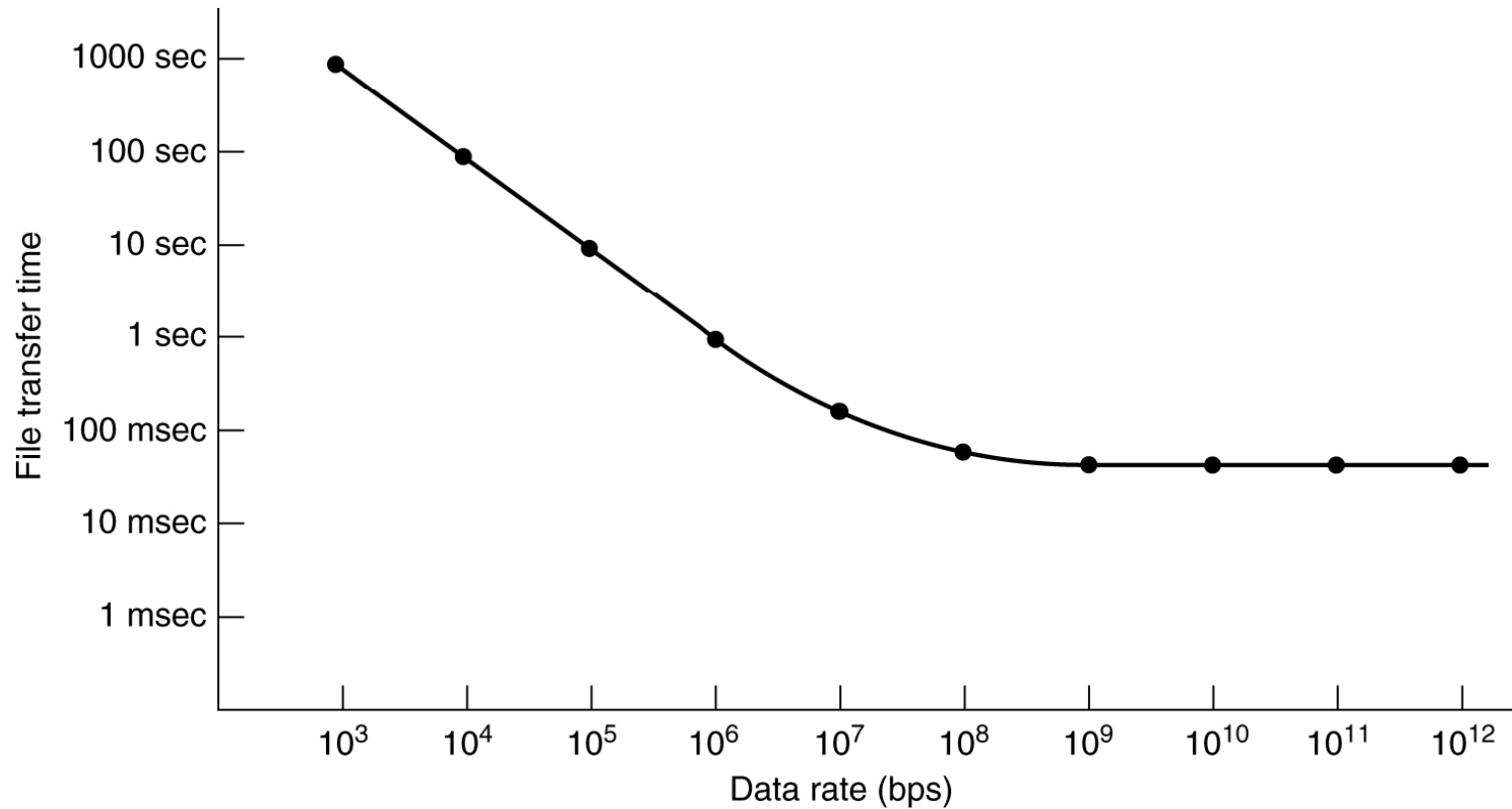(a) TCP header. (b) IP header. In both cases, the shaded fields are taken from the prototype without change.

# Fast TPDU Processing (3)

Slot

| | |
|---|---|
| 0 | → Pointer to list of timers for T + 12 |
| 1 | 0 |
| 2 | 0 |
| 3 | 0 |
| 4 | 0 ← Current time, T |
| 5 | 0 |
| 6 | 0 |
| 7 | → Pointer to list of timers for T + 3 |
| 8 | 0 |
| 9 | 0 |
| 10 | 0 |
| 11 | 0 |
| 12 | 0 |
| 13 | 0 |
| 14 | → Pointer to list of timers for T + 10 |
| 15 | 0 |

A timing wheel.

# Protocols for Gigabit Networks



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