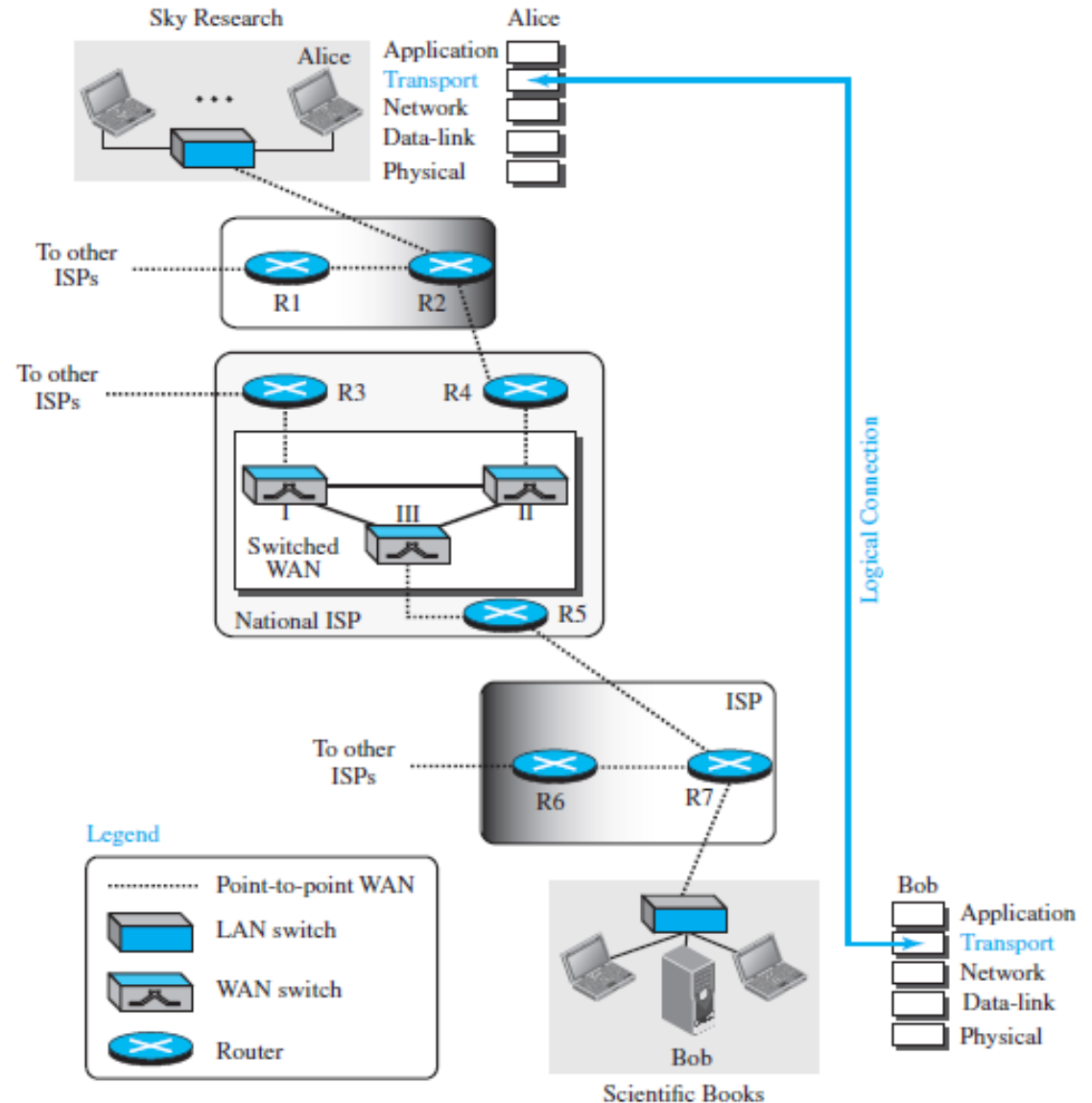# Transport Layer

Dr. Subrat K. Dash

# Transport Layer – The Link

- Application layer
  – Communication for specific applications
  – E.g., HyperText Transfer Protocol (HTTP), File Transfer Protocol (FTP), Network News Transfer Protocol (NNTP)
- Transport layer
  – Communication between processes (e.g., socket)
  – Relies on network layer and serves the application layer
  – E.g., TCP and UDP
- Network layer
  – Logical communication between nodes
  – Hides details of the link technology
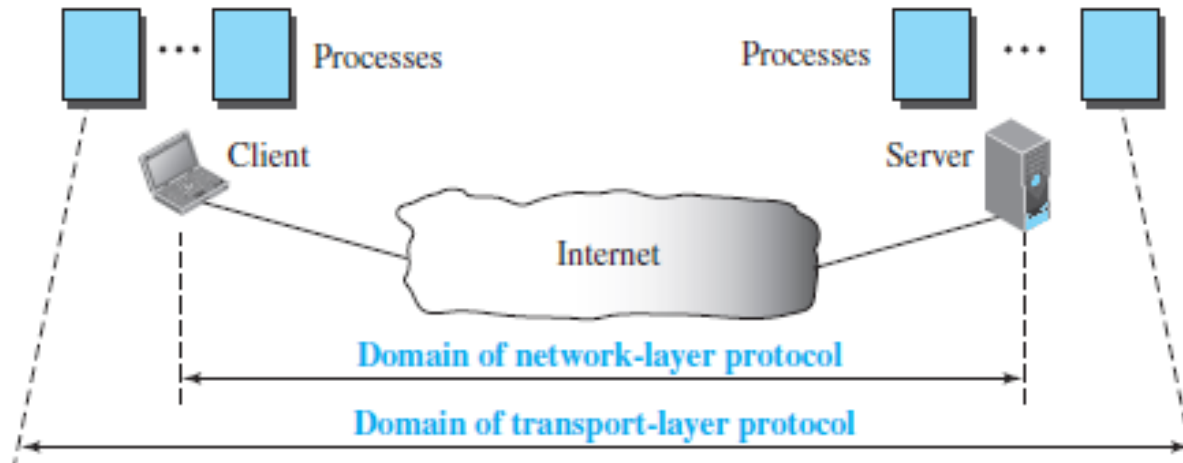  – E.g., IP

# Transport Layer

# Transport Protocols

- Run on end hosts
  - Sender: breaks application messages into segments,
    and passes to network layer
  - Receiver: reassembles segments into messages, passes to application layer
- Multiple transport protocol available to applications
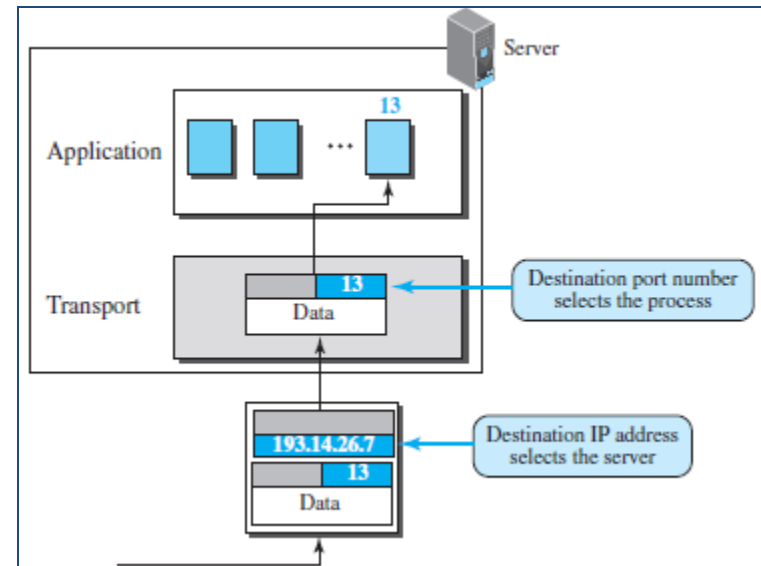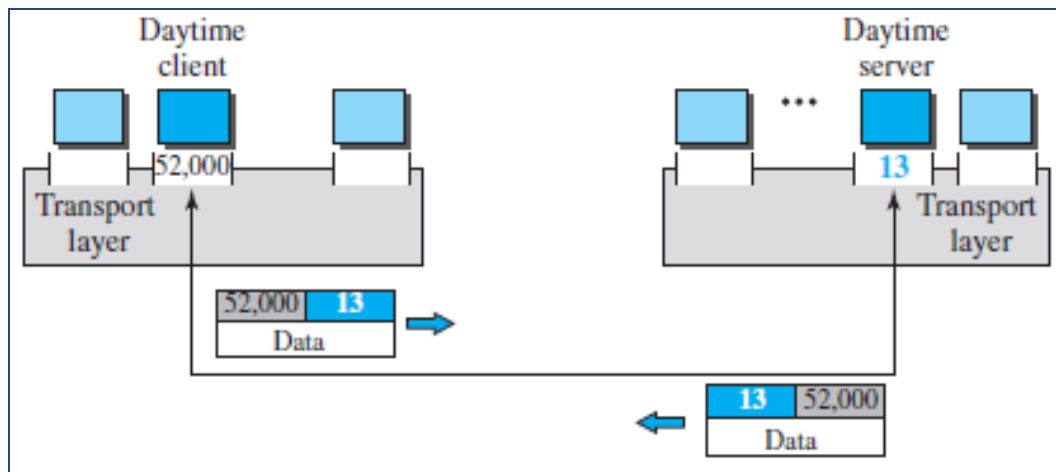  - Internet: TCP and UDP

# Transport Layer Services
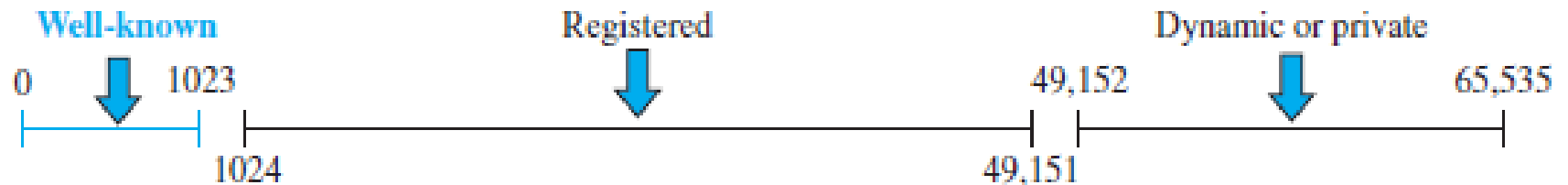
- Process-to-Process Communication

# Transport Layer Services

- Addressing: Port Numbers (16 bits)
  - Client-server paradigm
  - Port numbers
    - Client: Ephemeral port (>1023)
    - Server: Well-known port

# ICANN Ranges



| Well-known | Registered | Dynamic or private |
|---|---|---|
| 0    1023 | 49,152 | 65,535 |
| 1024 | 49,151 | |

```
$grep    tftp/etc/services
tftp 69/tcp
tftp 69/udp
```
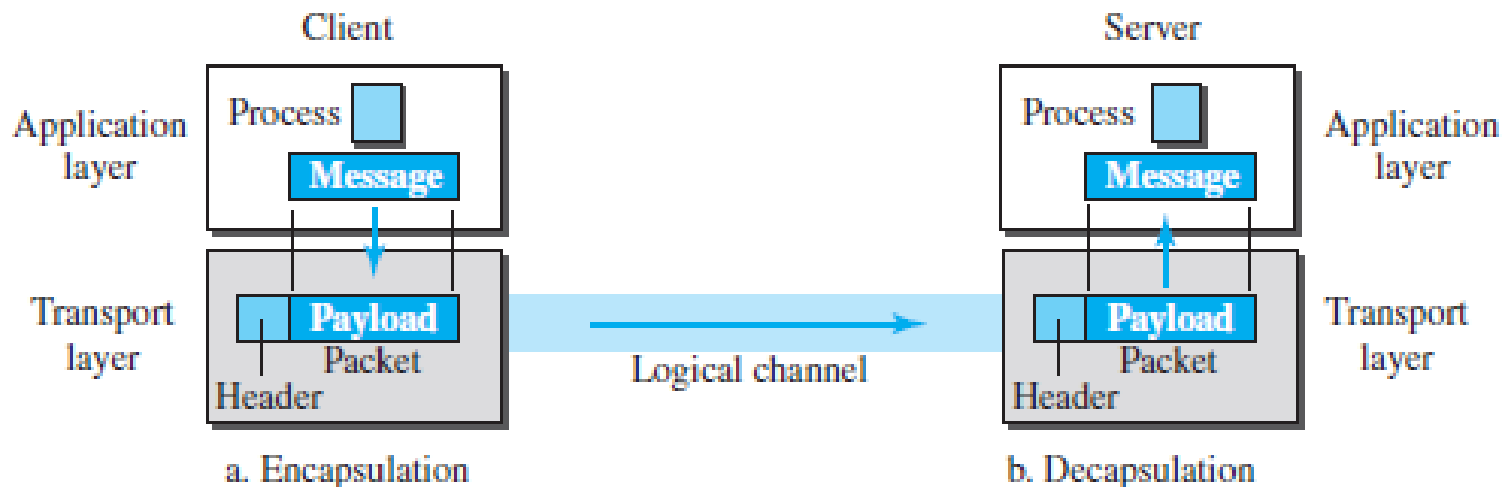
```
$grep    snmp/etc/services
snmp161/tcp#Simple Net Mgmt Proto
snmp161/udp#Simple Net Mgmt Proto
snmptrap162/udp#Traps for SNMP
```
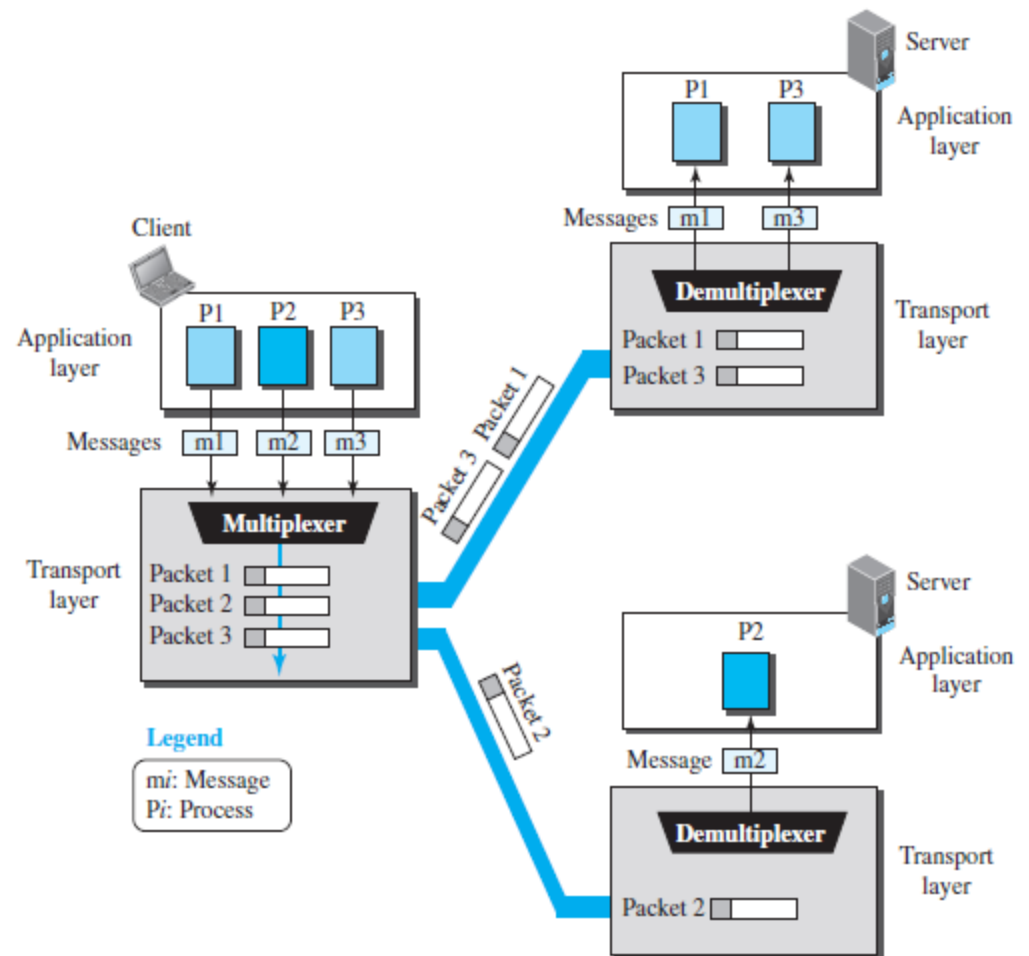
# Transport Layer Services

- Encapsulation and Decapsulation
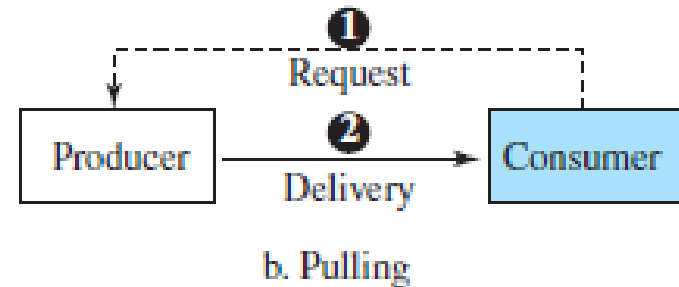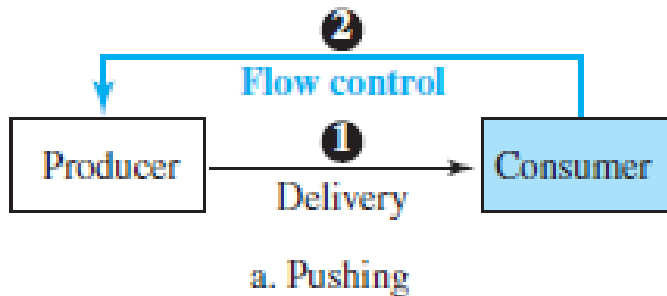  - User datagrams, Segments, Packets

# Transport Layer Services

- Multiplexing and Demultiplexing

# Transport Layer Services

- Flow Control
  - Items are produced faster than they can be consumed.



a. Pushing

b. Pulling

# Transport Layer Services

- Flow Control
  - Items are produced faster than they can be consumed.



a. Pushing

b. Pulling



Buffers

Buffers

# Transport Layer Services

- Error Control
  - Detecting and discarding corrupted packets.
  - Keeping track of lost and discarded packets and resending them.
  - Recognizing duplicate packets and discarding them.
  - Buffering out-of-order packets until the missing packets arrive.

# Transport Layer Services

- Error Control
  - Sequence Number
  - Acknowledgement



a. Four packets have been sent.

b. Five packets have been sent.

c. Seven packets have been sent; window is full.

d. Packet 0 has been acknowledged; window slides.

# Transport Layer Services

- Congestion Control
  - Congestion
    - Results due to imbalance between load of network and capacity of network
    - Waiting
  - TCP

# Protocol Types

- Connectionless
  - Independency between packets (TL)
  - Different paths for different datagrams belonging to the same message (NL)
- Connection-oriented
  - Dependency between packets (TL)
  - Same path for different datagrams belonging to the same message (NL)

# Connectionless vs. Connection-oriented

# Transport-Layer Protocols in TCP/IP Protocol Suite

# Sample list of well-known ports

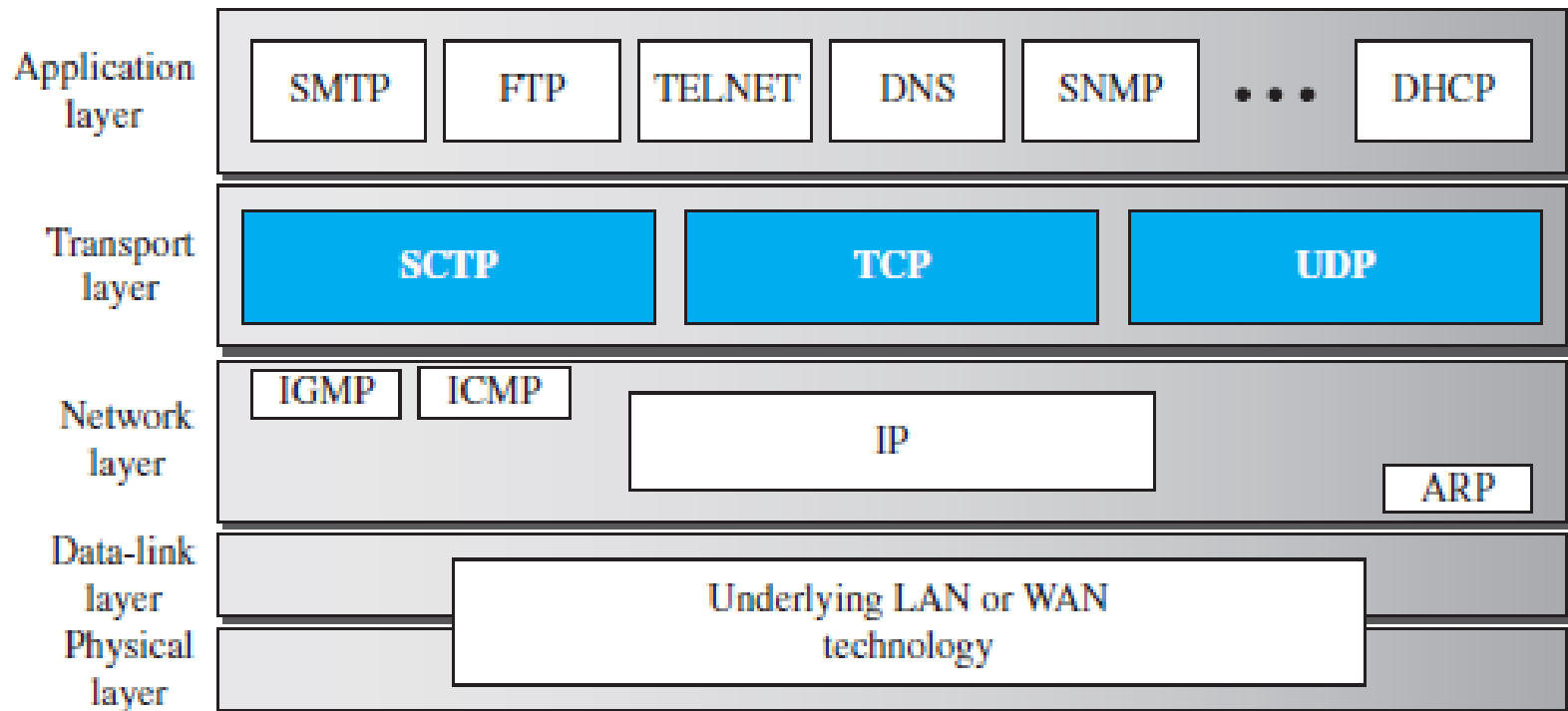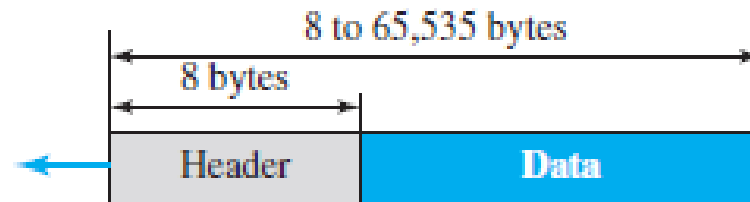| Port | Protocol | UDP | TCP | SCTP | Description |
|------|----------|-----|-----|------|-------------|
| 7 | Echo | √ | √ | √ | Echoes back a received datagram |
| 9 | Discard | √ | √ | √ | Discards any datagram that is received |
| 11 | Users | √ | √ | √ | Active users |
| 13 | Daytime | √ | √ | √ | Returns the date and the time |
| 17 | Quote | √ | √ | √ | Returns a quote of the day |
| 19 | Chargen | √ | √ | √ | Returns a string of characters |
| 20 | FTP-data | | √ | √ | File Transfer Protocol |
| 21 | FTP-21 | | √ | √ | File Transfer Protocol |
| 23 | TELNET | | √ | √ | Terminal Network |
| 25 | SMTP | | √ | √ | Simple Mail Transfer Protocol |
| 53 | DNS | √ | √ | √ | Domain Name Service |
| 67 | DHCP | √ | √ | √ | Dynamic Host Configuration Protocol |
| 69 | TFTP | √ | √ | √ | Trivial File Transfer Protocol |
| 80 | HTTP | | √ | √ | HyperText Transfer Protocol |
| 111 | RPC | √ | √ | √ | Remote Procedure Call |
| 123 | NTP | √ | √ | √ | Network Time Protocol |
| 161 | SNMP-server | √ | | | Simple Network Management Protocol |
| 162 | SNMP-client | √ | | | Simple Network Management Protocol |

# User Datagram Protocol (UDP)

- Connectionless, unreliable TL Protocol

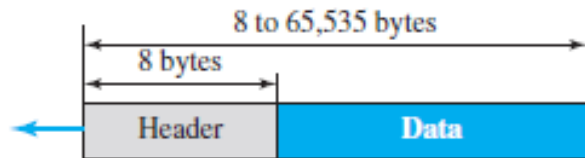- User datagram packet format



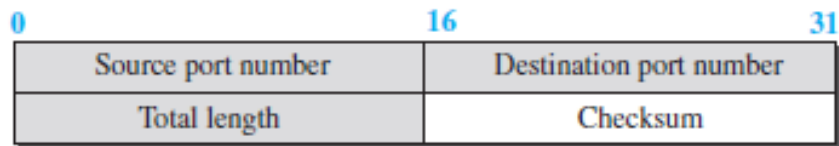a. UDP user datagram

b. Header format

# UDP Services

- Process-to-process communication
- Encapsulation and decapsulation
- Multiplexing and demultiplexing
- Flow control
- Error control
- Congestion control
- Connectionless services
- Checksum

# UDP: Checksum



8 to 65,535 bytes

8 bytes

| Header | Data |
|--------|------|

a. UDP user datagram

| 0 | 16 | 31 |
|---|----|----|

| Source port number | Destination port number |
|--------------------|-------------------------|
| Total length | Checksum |

b. Header format

**Pseudoheader**

| 32-bit source IP address | | |
|---|---|---|
| 32-bit destination IP address | | |
| All 0s | 8-bit protocol | 16-bit UDP total length |

**Header**

| Source port address 16 bits | Destination port address 16 bits |
|------------------------------|----------------------------------|
| UDP total length 16 bits | Checksum 16 bits |

**Data**
(Padding must be added to make the data a multiple of 16 bits)

# TCP

- Connection oriented
  - Explicit set-up and tear-down of TCP session
- Stream-of-bytes service
  - Sends and receives a stream of bytes, not messages
- Reliable, in-order delivery
  - Checksums to detect corrupted data
  - Acknowledgments & retransmissions for reliable delivery
  - Sequence numbers to detect losses and reorder data
- Flow control
  - Prevent overflow of the receiver's buffer space
- Congestion control
  - Adapt to network congestion for the greater good

# An Analogy: Talking on a Cell Phone

- Alice and Bob on their cell phones
  - Both Alice and Bob are talking
- What if Alice couldn't understand Bob?
  - Bob asks Alice to repeat what she said
- What if Bob hasn't heard Alice for a while?
  - Is Alice just being quiet?
  - Or, have Bob and Alice lost reception?
  - How long should Bob just keep on talking?
  - Maybe Alice should periodically say "uh huh"
  - … or Bob should ask "Can you hear me now?"

# Some Take-Aways from the Example

- Acknowledgments from receiver
  - Positive: "okay" or "ACK"
  - Negative: "please repeat that" or "NACK"
- Timeout by the sender ("stop and wait")
  - Don't wait indefinitely without receiving some response
  - … whether a positive or a negative acknowledgment
- Retransmission by the sender
  - After receiving a "NACK" from the receiver
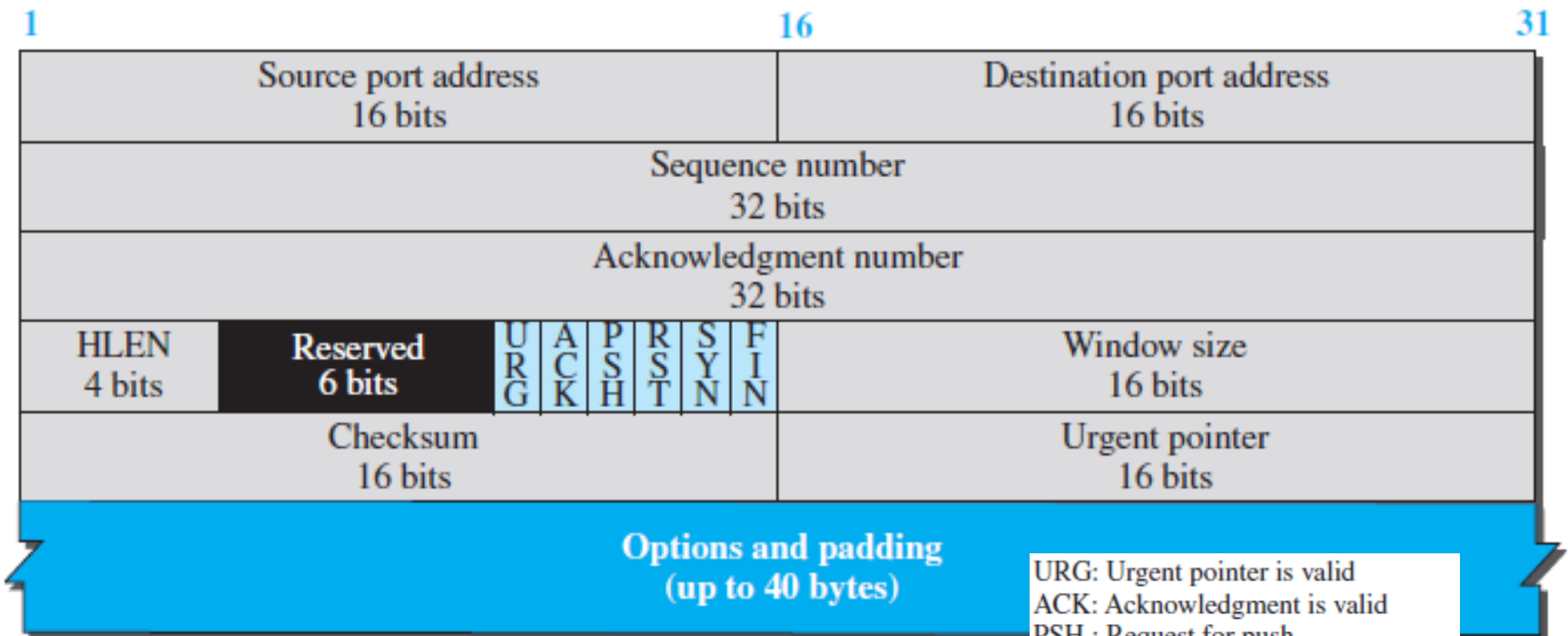  - After receiving no feedback from the receiver

# TCP Services

- Process-to-process communication

- Stream delivery service

- Full-Duplex communication

- Multiplexing and demultiplexing

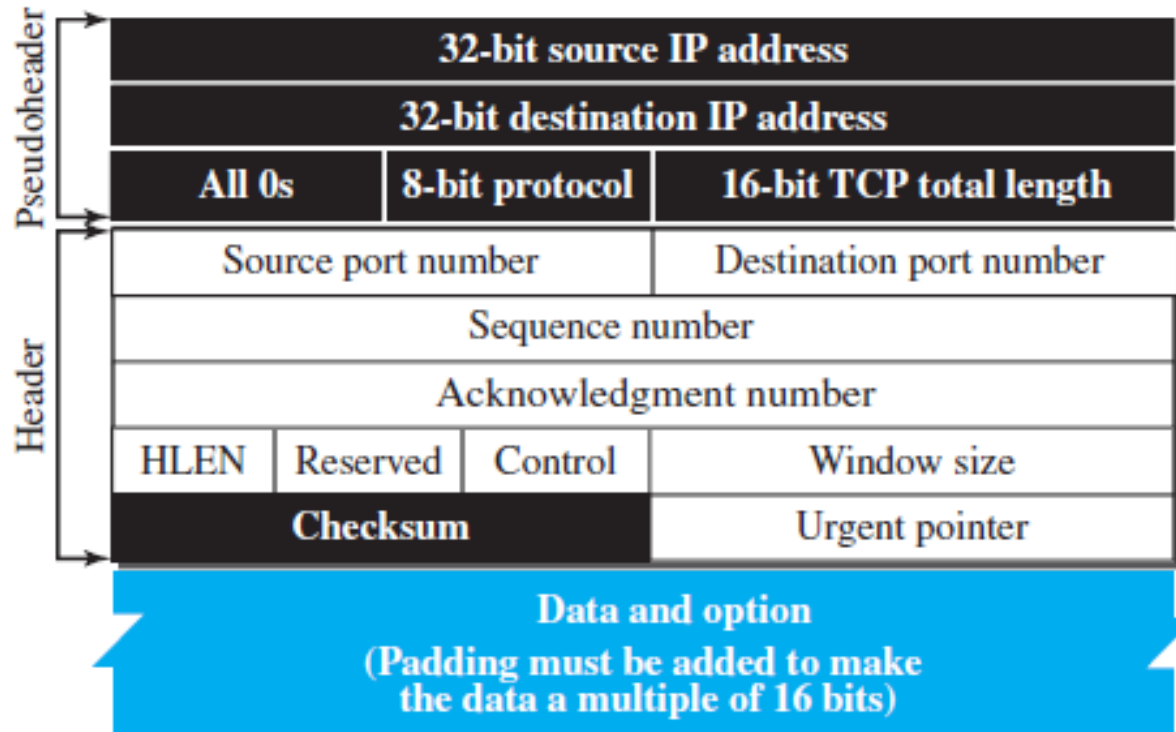- Connection-oriented service

- Reliable service

# TCP Segment



20 to 60 bytes

| Header | Data |

a. Segment

| 1 | 16 | 31 |
|---|---|---|
| Source port address<br>16 bits | | Destination port address<br>16 bits |
| Sequence number<br>32 bits | | |
| Acknowledgment number<br>32 bits | | |
| HLEN<br>4 bits | Reserved<br>6 bits | U R G / A C K / P S H / R S T / S Y N / F I N | Window size<br>16 bits |
| Checksum<br>16 bits | | Urgent pointer<br>16 bits |
| Options and padding<br>(up to 40 bytes) | | |

b. Header

URG: Urgent pointer is valid
ACK: Acknowledgment is valid
PSH : Request for push
RST : Reset the connection
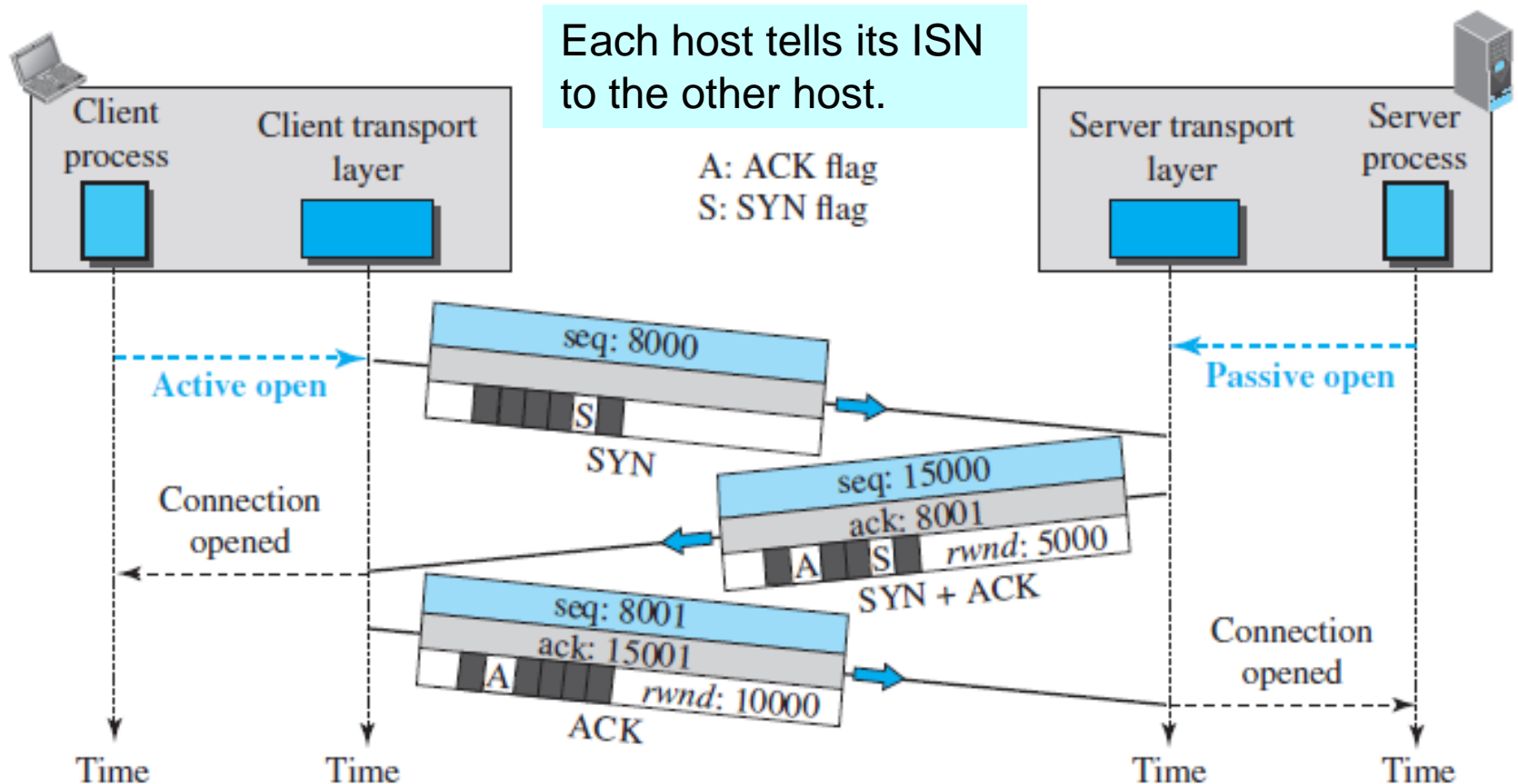SYN: Synchronize sequence numbers
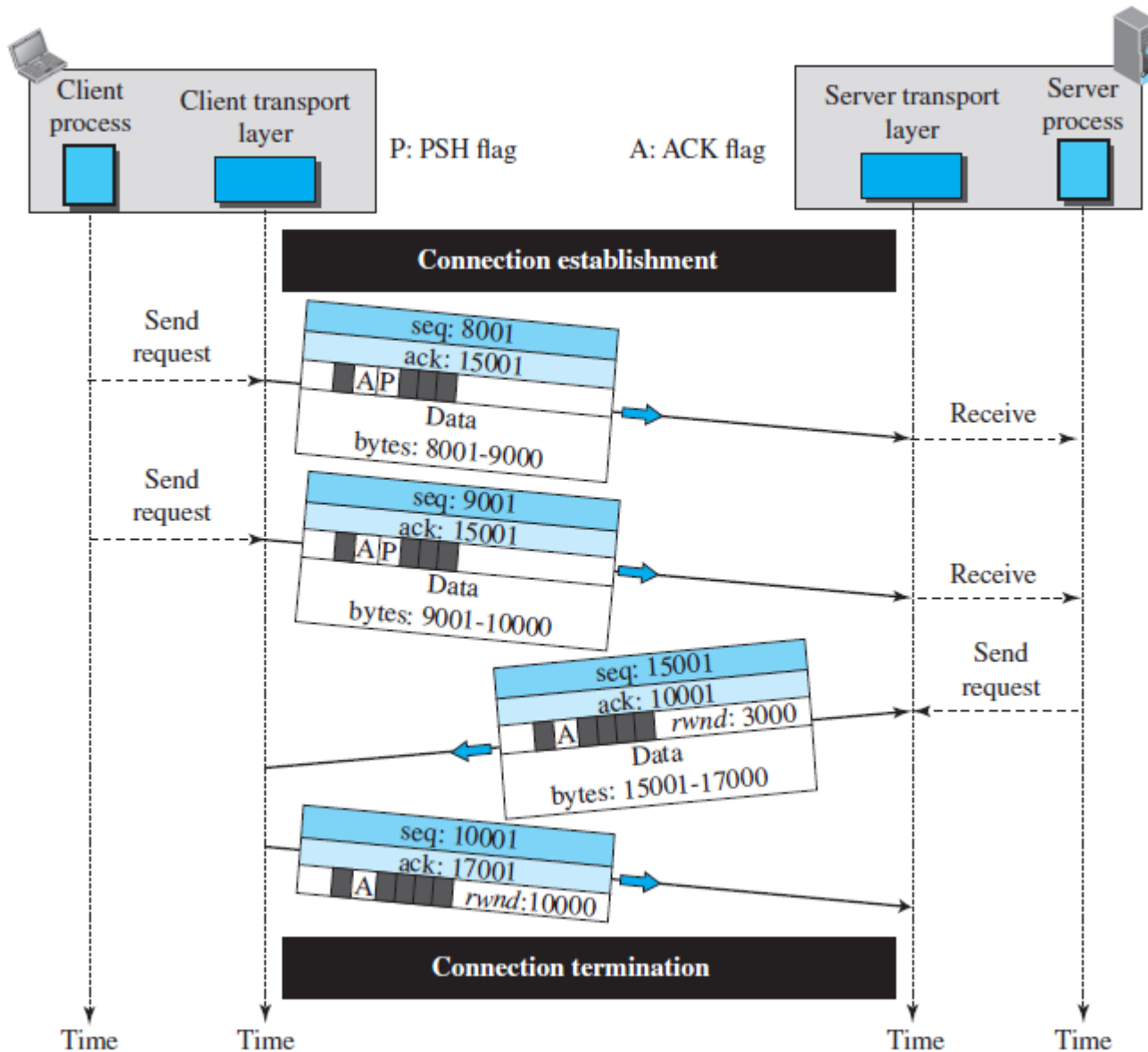FIN : Terminate the connection

# TCP – Checksum



The use of the checksum in TCP is mandatory.

# TCP Connection

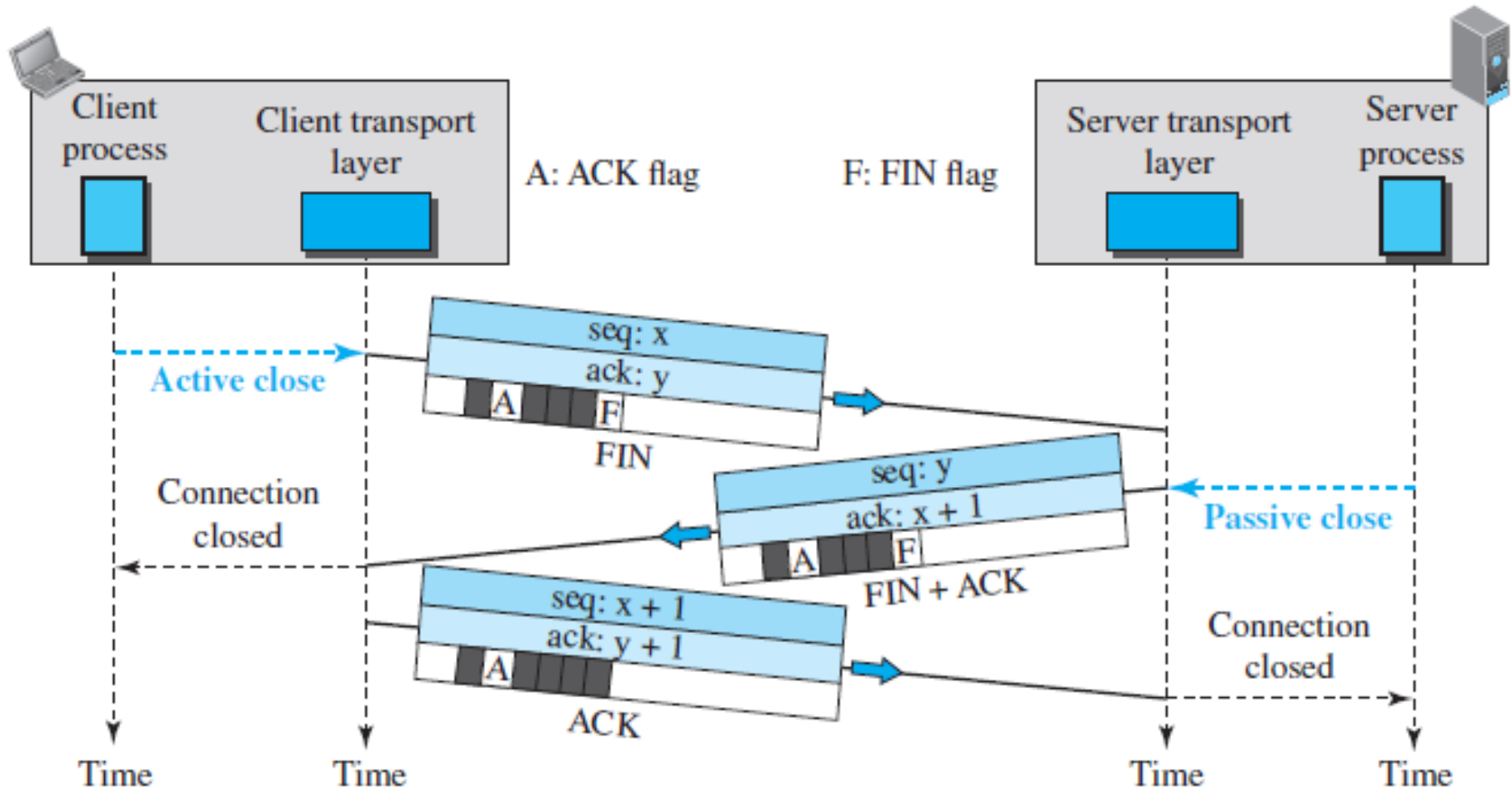- Connection establishment
  - Three-way handshaking

SYN Flooding Attack
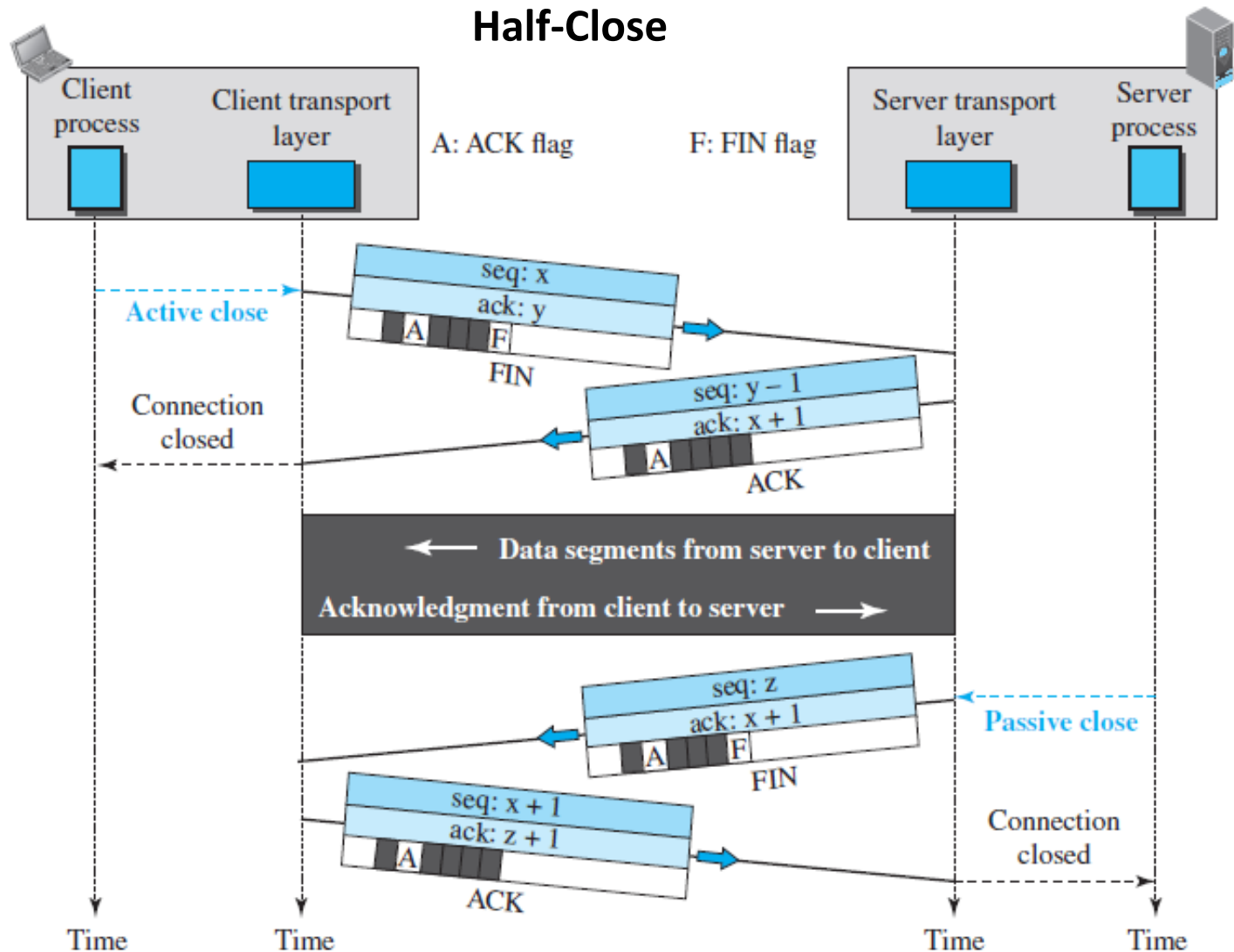


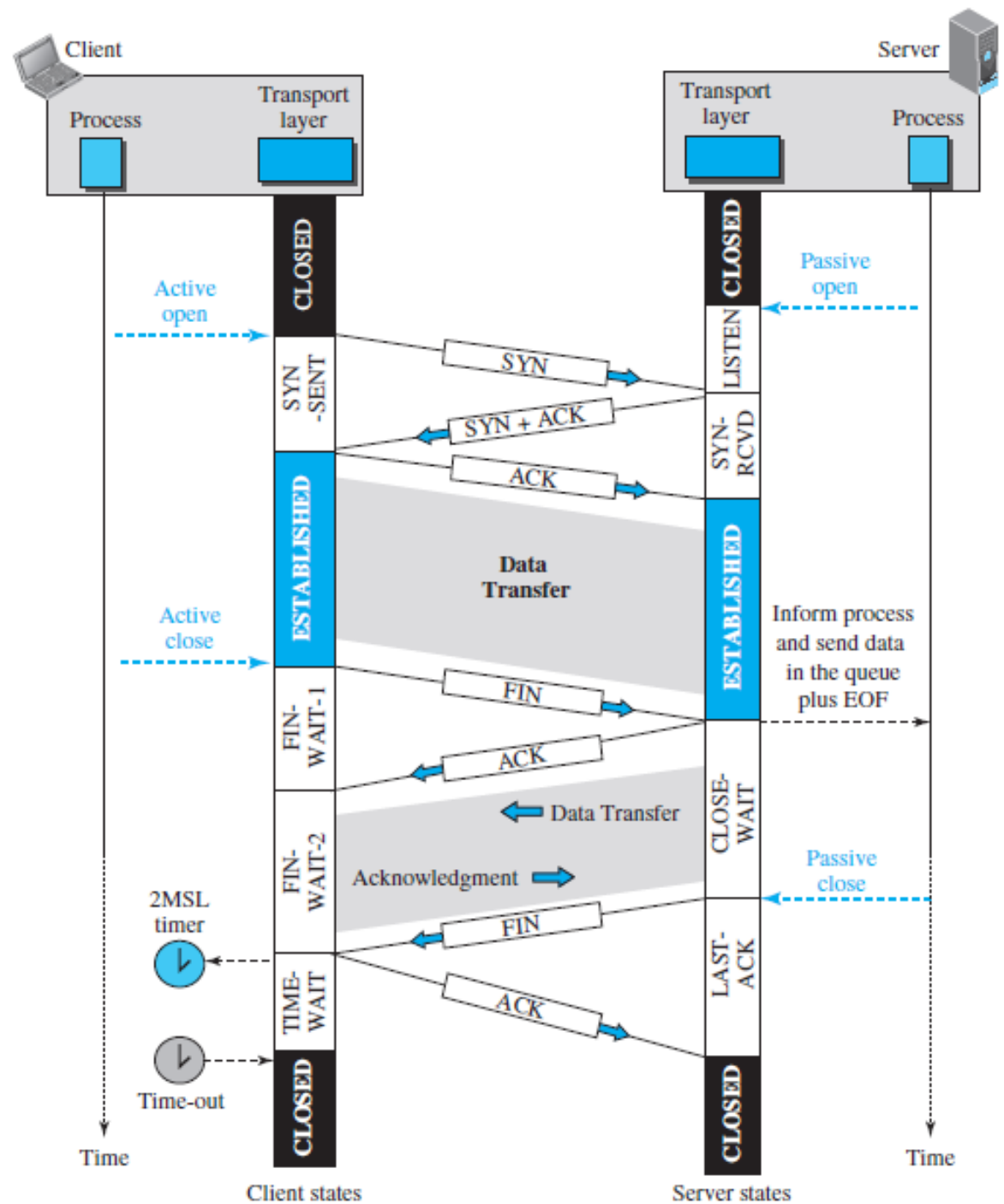Each host tells its ISN to the other host.

A: ACK flag
S: SYN flag

# Data Transfer

# Connection Termination

**Three-Way Handshaking**
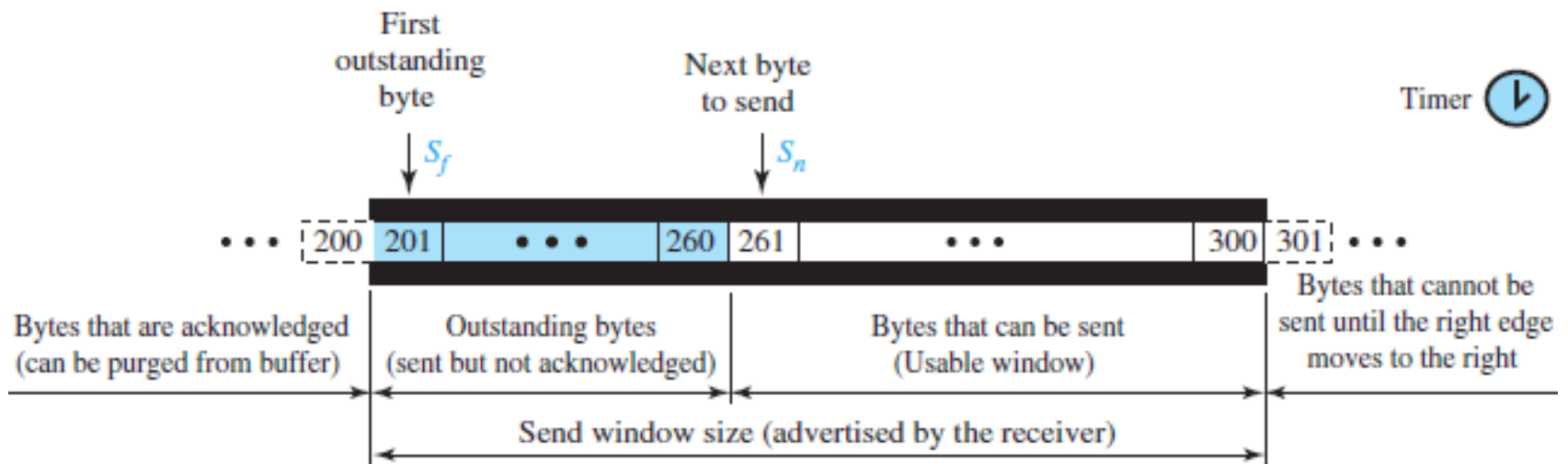
# Connection Termination



**Half-Close**

Time-line diagram for a common scenario

# Windows in TCP

**Send Window**



a. Send window

b. Opening, closing, and shrinking send window
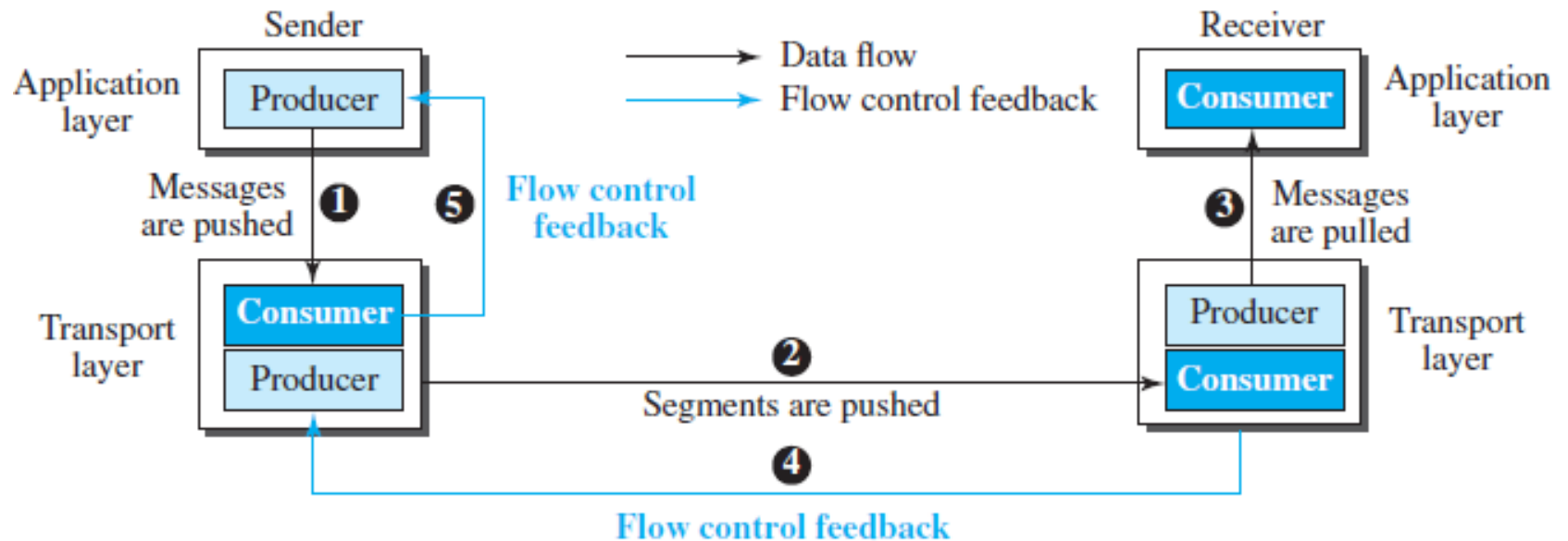
# Windows in TCP

**Receive Window**



Next byte to be pulled by the process

Next byte expected to be received

$R_n$

... 200 | 201 | ... | 260 | 261 | ... | 300 | 301 ...

Bytes that have already been pulled by the process

Bytes received and acknowledged, waiting to be consumed by process

Bytes that can be received from sender
Receive window size (*rwnd*)

Bytes that cannot be received from sender

Allocated buffer

a. Receive window and allocated buffer

Left wall

Right wall

Closes

Opens

... 200 | 201 | ... | 260 | 261 | ... | 300 | 301 ...

b. Opening and closing of receive window

# TCP – Flow Control



Data flow and flow control feedbacks in TCP

**Client**

SYN
① seqNo: 100

SYN + ACK
② seqNo: 1000
ackNo: 101
rwnd: 800

Size = 800
101                     901
Send window is set.

ACK
③ ackNo: 1001
rwnd: 2000

Size = 800
101  301                901
Sender sends 200 bytes.

Data
④ seqNo: 101
Data: 200 bytes

Size = 600
301                     901
Bytes acknowledged, window closes.

ACK
⑤ ackNo: 301
rwnd: 600

Size = 600
301  601                901
Sender sends 300 bytes.

Data
⑥ seqNo: 301
Data: 300 bytes

Size = 400
601        1001
Window closes and opens.

ACK
⑦ ackNo: 601
rwnd: 400

Size = 600
601                     1201
Window opens.

ACK
⑧ ackNo: 601
rwnd: 600

**Server**

rwnd = 800
101                     901
Receive window is set.

rwnd = 600
101  301                901
200 bytes received, window closes.

rwnd = 400
201  601        1001
300 bytes received, 100 bytes consumed.

rwnd = 600
401  601                1201
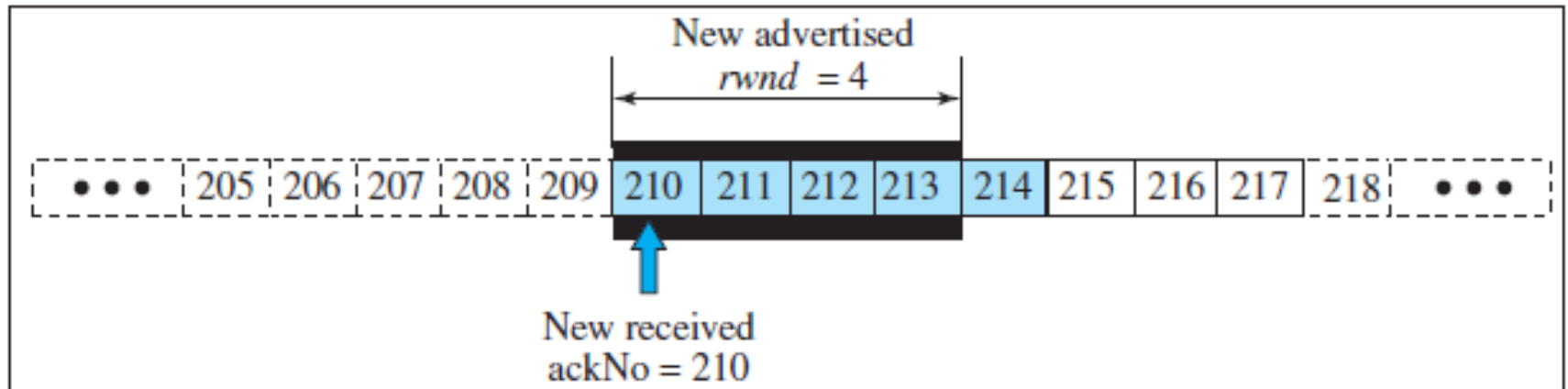200 bytes consumed, window opens.

# Shrinking of Window



a. The window after the last advertisement

b. The window after the new advertisement; window has shrunk

**Window Shutdown**

# Silly Window Syndrome

- The sending application program creates data slowly or the receiving application program consumes data slowly, or both.

  - Creates problem in sliding window operation and leads to overhead.

- Syndrome created by

  - Sender

  - Receiver

# Syndrome created by the Sender

- The sending TCP may create a silly window syndrome if it is serving an application program that creates data slowly
- Solution: prevent the sending TCP from sending the data in small quantity
  - Nagle's Algorithm
    - Send first piece of data received from sending app.
    - Accumulate data and wait until an ACK is received or until max data is accumulated

# Syndrome created by the Receiver

- The receiving TCP may create a silly window syndrome if it is serving an application program that consumes data slowly

- Two Solutions:
  - send an ACK as soon as the data arrive, but to announce a window size of zero until either there is enough space to accommodate a segment of maximum size or until at least half of the receive buffer is empty (Clark's Solution)
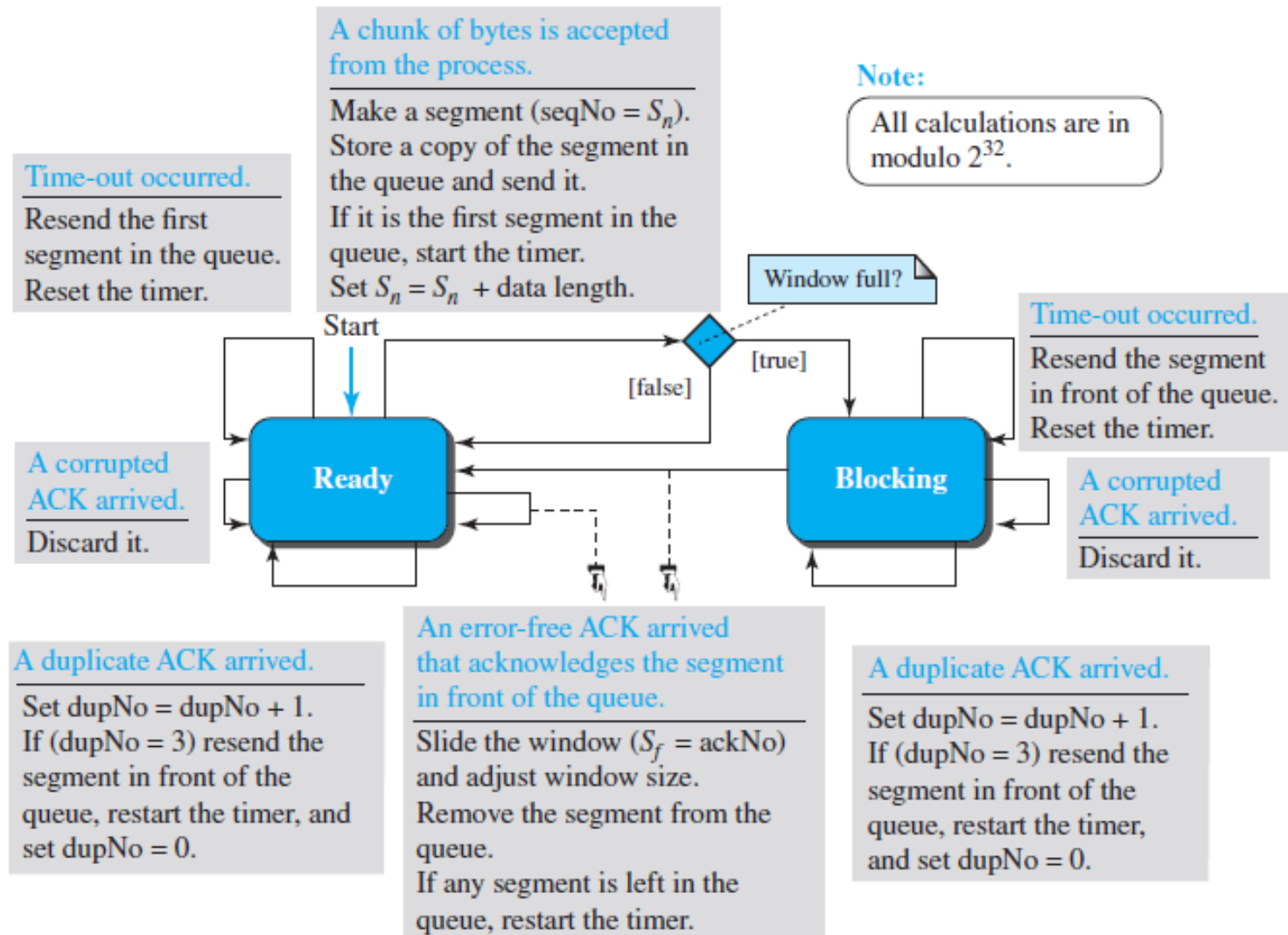  - delay sending the acknowledgment

# TCP – Error Control

- Detecting and resending corrupted segments
- Resending lost segments
- Storing out-of order segments until missing segments arrive
- Detecting and discarding duplicated segments

- Solution
  - checksum, acknowledgment (Cumulative, Selective), and time-out

# TCP Acknowledgement

- Acknowledgement
  - Cumulative (ACK)
  - Selective (SACK)

- Retransmission
  - After RTO (Retransmission time-out)
  - After three duplicate ACK segments
    - Fast retransmission

# FSM for the TCP sender side

# FSM for the TCP receiver side

**An expected error-free segment arrived.**

Buffer the message.
$R_n = R_n$ + data length.
If the ACK-delaying timer is running,
stop the timer and send a cumulative ACK.
Otherwise, start the ACK-delaying timer.

**Note:**

All calculations are in
modulo $2^{32}$.

**A request for delivery of**
**$k$ bytes of data from**
**process came.**

Deliver the data.
Slide the window and
adjust window size.

Start

**Ready**

**ACK-delaying timer expired.**
Send the delayed ACK.

**An error-free, but out-of-**
**order segment arrived.**

Store the segment if not duplicate.
Send an ACK with ackNo equal
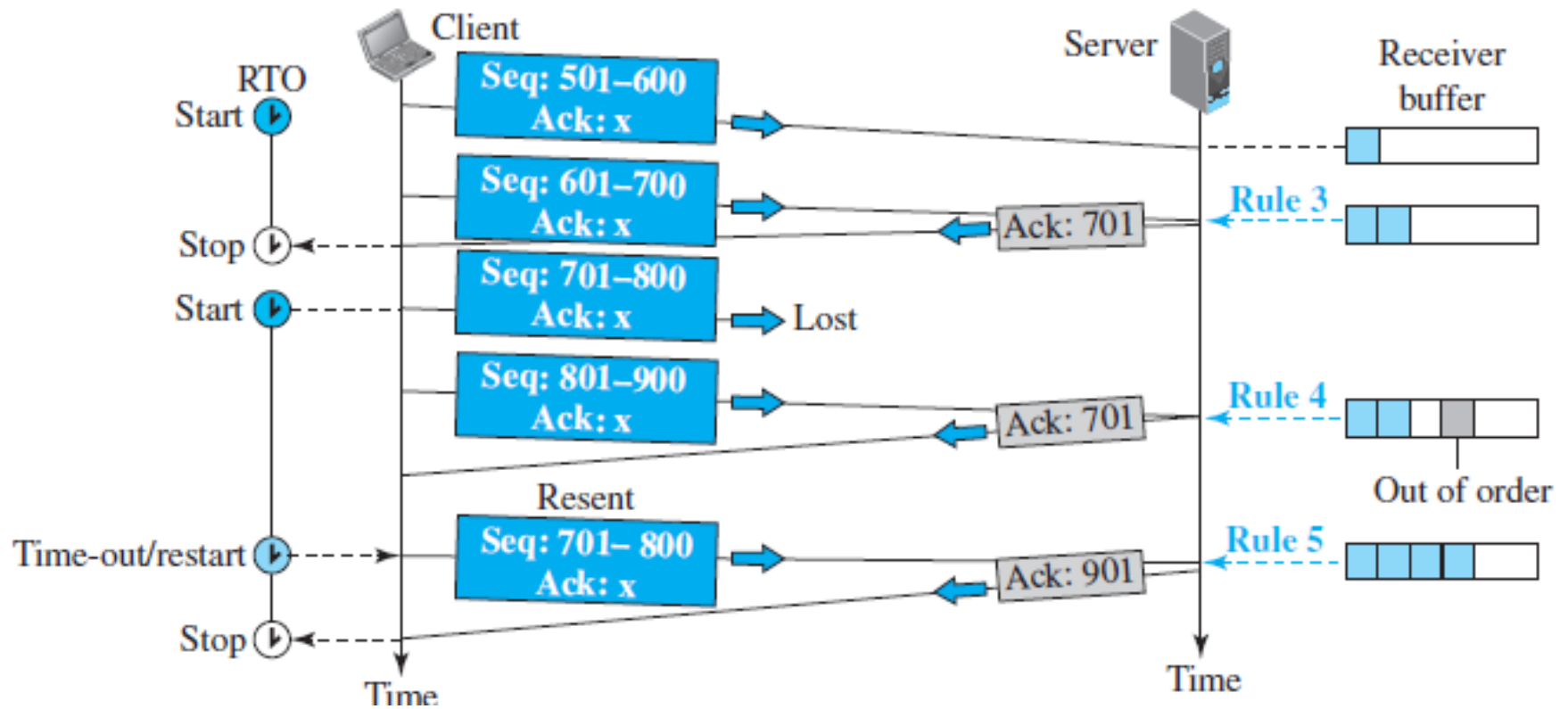to the sequence number of expected
segment (duplicate ACK).

**An error-free duplicate segment**
**or an error-free segment with**
**sequence number outside**
**window arrived.**

Discard the segment.
Send an ACK with ackNo equal
to the sequence number of expected
segment (duplicate ACK).
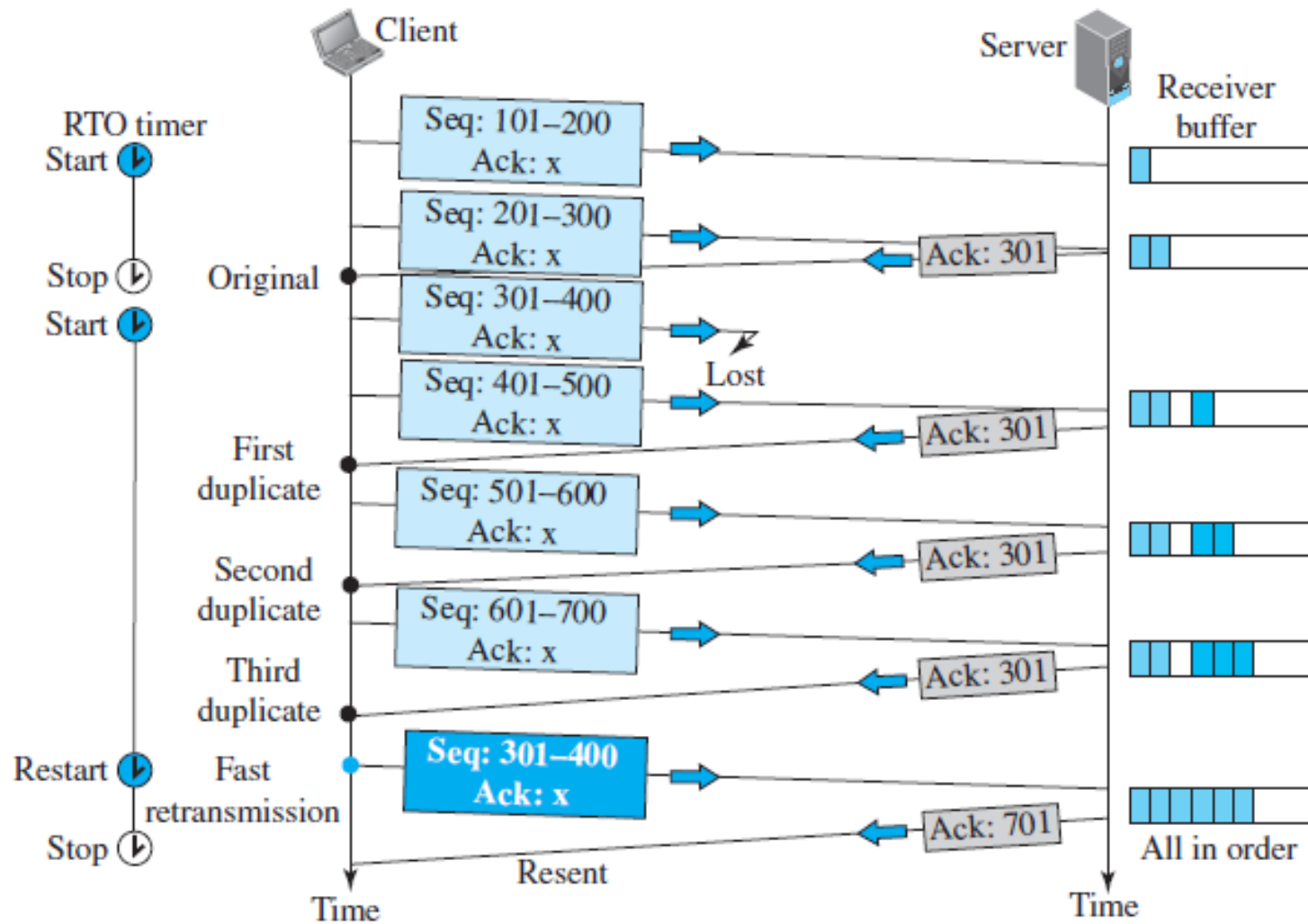
**A corrupted segment arrived.**
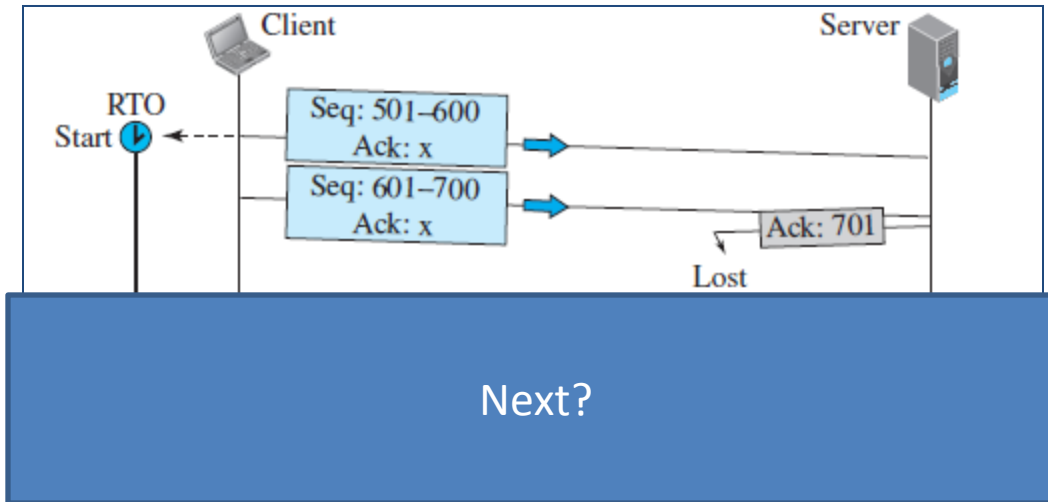
Discard the segment.

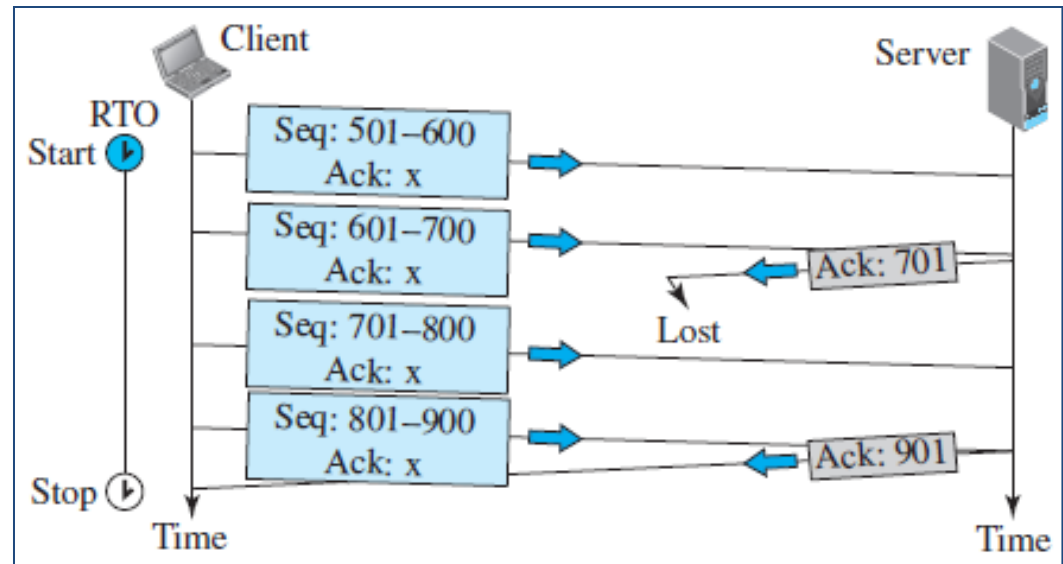# Normal Operations

# Lost Segment

# Fast retransmission

# Lost acknowledgement



Next?

Deadlock due to Lost acknowledgement

# TCP – Congestion Control

- Implemented at the Sender side
- Send window = MIN(*rwnd, cwnd*)
  - *rwnd*: advertised by the receiver
  - *cwnd*: adjusted based on feedback from the network

- Congestion Detection
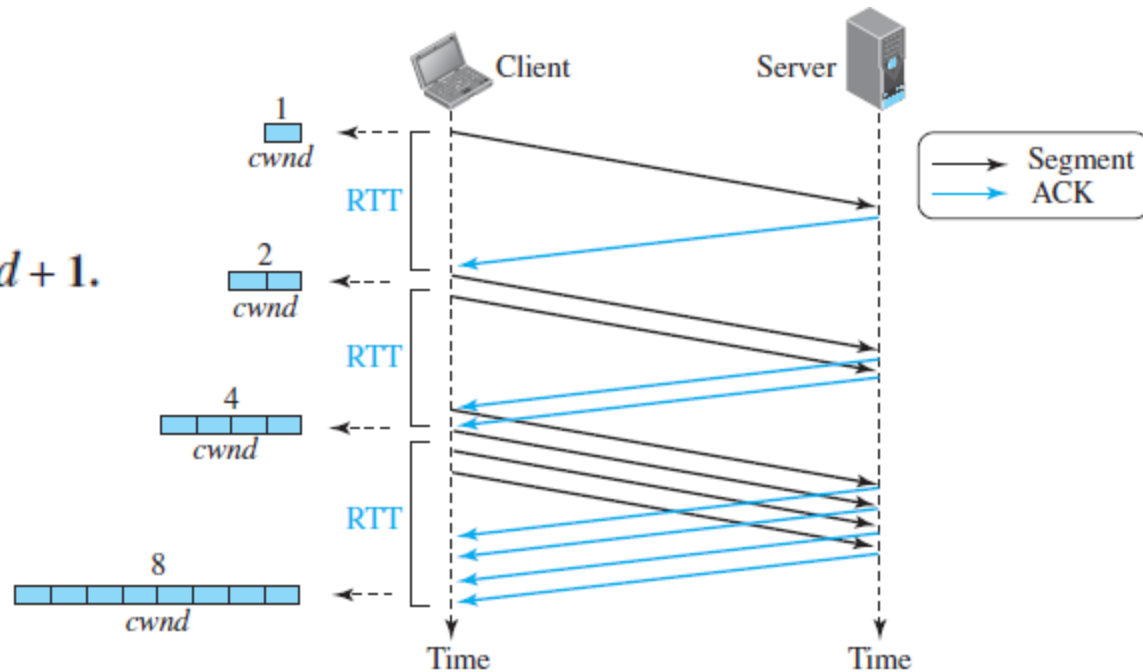  - Time-out – Strong congestion
  - Three duplicate ACKs – Weak congestion

# Congestion Policies

- Slow start: Exponential increase
- Congestion avoidance
- Fast recovery

# Slow Start

**Initially, *cwnd* = 1 MSS**

**If an ACK arrives, $cwnd = cwnd + 1$.**



| Start | $\rightarrow$ | $cwnd = 1 \rightarrow 2^0$ |
|---|---|---|
| **After 1 RTT** | $\rightarrow$ | $cwnd = cwnd + 1 = 1 + 1 = 2 \rightarrow 2^1$ |
| **After 2 RTT** | $\rightarrow$ | $cwnd = cwnd + 2 = 2 + 2 = 4 \rightarrow 2^2$ |
| **After 3 RTT** | $\rightarrow$ | $cwnd = cwnd + 4 = 4 + 4 = 8 \rightarrow 2^3$ |

# Congestion Avoidance
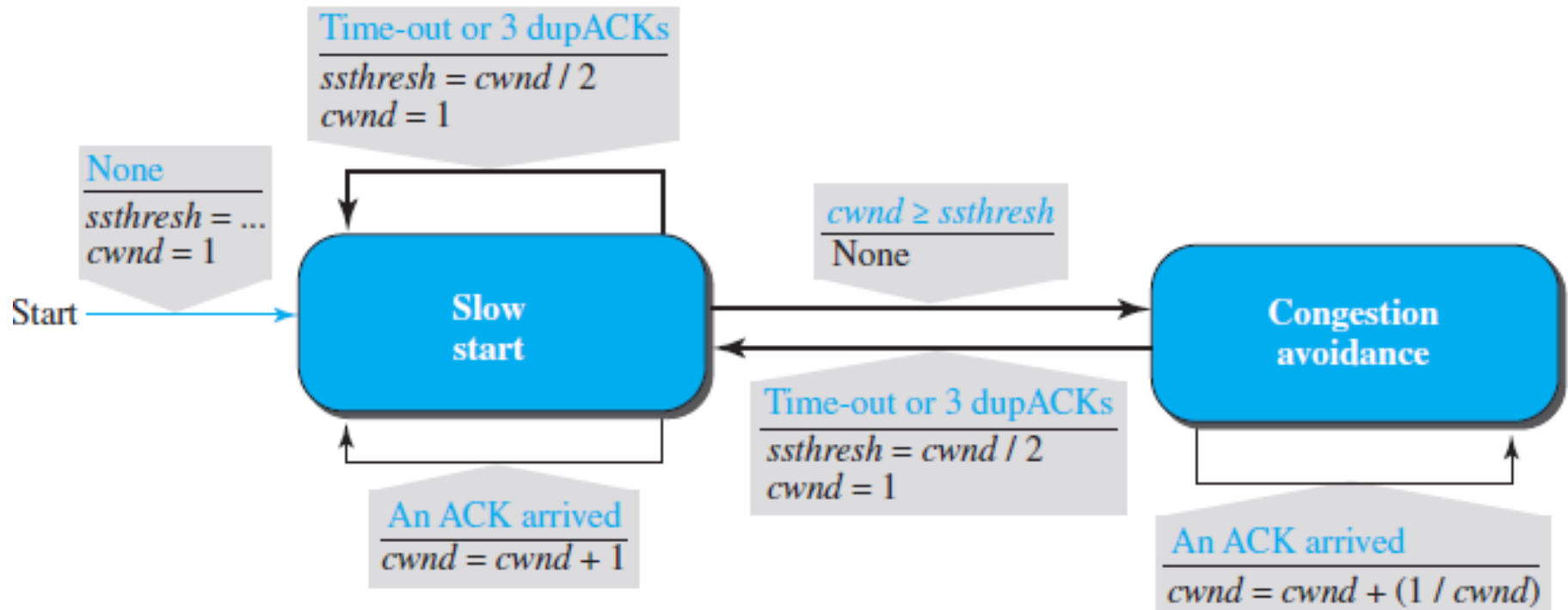
- Congestion avoidance phase is started if *cwnd* has reached the *slow-start threshold* value

- If <span style="color:red">cwnd >= ssthresh</span> then each time an ACK is received, increment *cwnd* as follows:
  - cwnd = cwnd + 1/cwnd

- So *cwnd* is increased by one segment (=MSS bytes) only if all segments have been acknowledged.
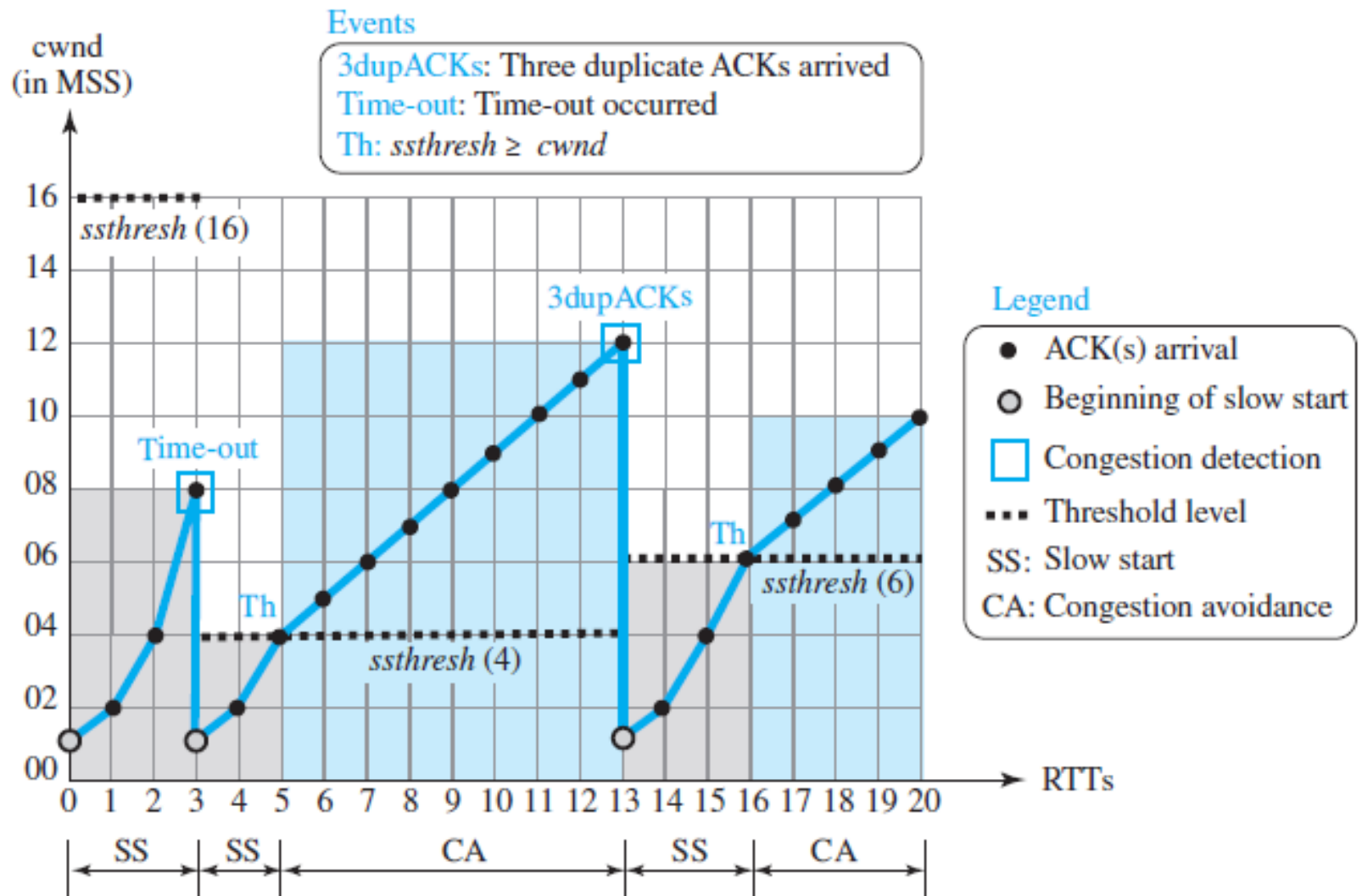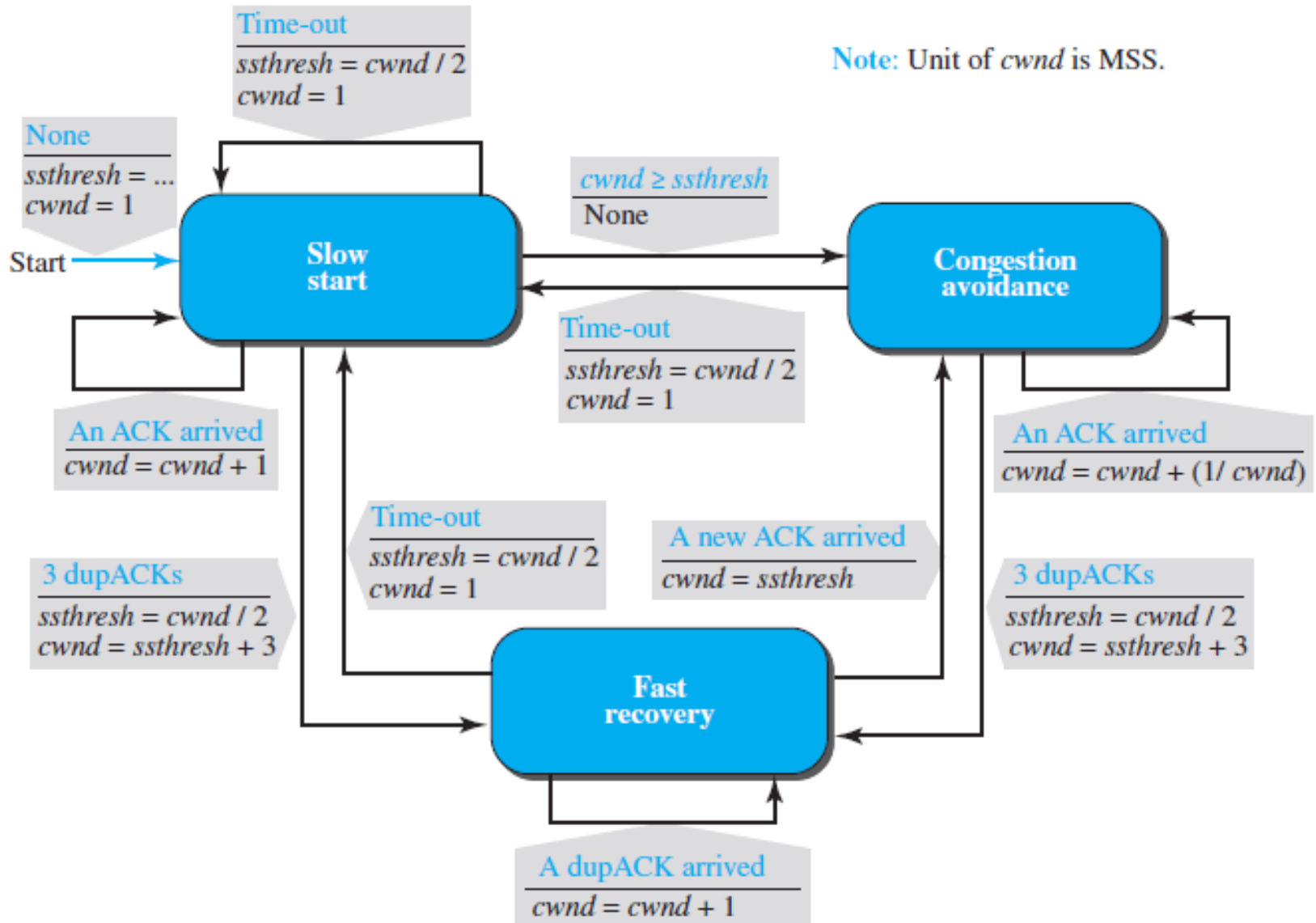
# Congestion Avoidance

# TCP – Congestion Control

- Taho TCP
  - Slow start and Congestion Avoidance

# Taho TCP – Example

# Reno TCP

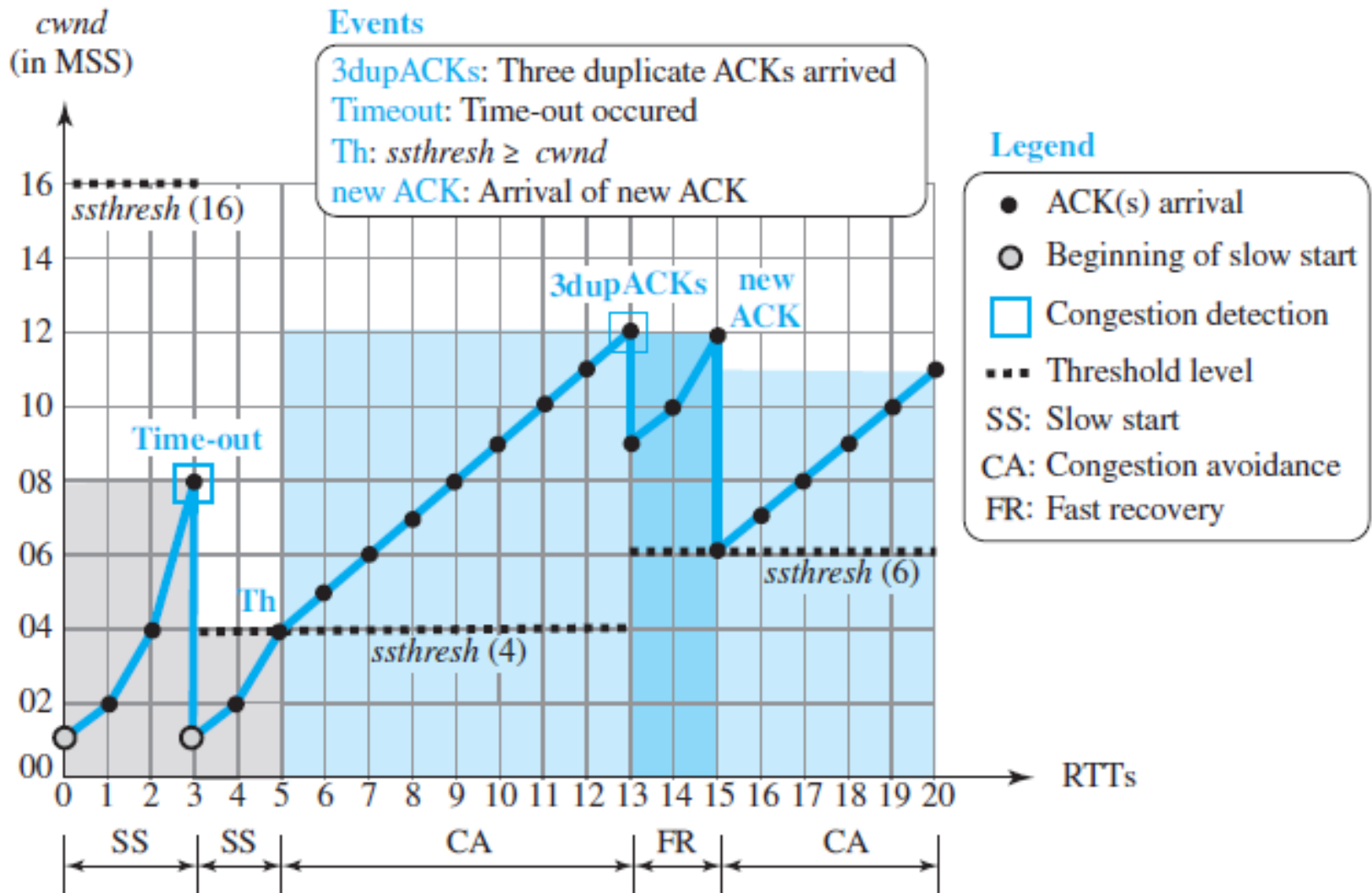# Reno TCP – Example

# AIMD

- Additive increase, multiplicative decrease (AIMD)