

UNIT - 4

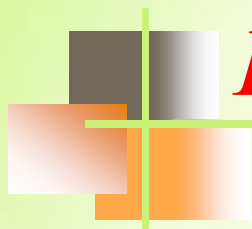
Transport Layer

3.1 INTRODUCTION

3.2 TRANSPORT-LAYER PROTOCOLS

3.3 USER DATAGRAM PROTOCOL

3.4 TRANSMISSION CONTROL PROTOCOL



INTRODUCTION

- ❑ Process-to-Process Communication

- ❑ Addressing: Port Numbers

- ❑ ICANN Ranges

 - ❖ Well-known ports

 - ❖ Registered ports

 - ❖ Dynamic ports

- ❑ Encapsulation and Decapsulation

- ❑ Multiplexing and Demultiplexing



INTRODUCTION

- Flow Control

- ❖ Pushing or Pulling
- ❖ Flow Control at Transport Layer
- ❖ Buffers

- Error Control

- ❖ Sequence Numbers
- ❖ Acknowledgment

- Combination of Flow and Error Control

- ❖ Sliding Window



INTRODUCTION

- ❑ Congestion Control

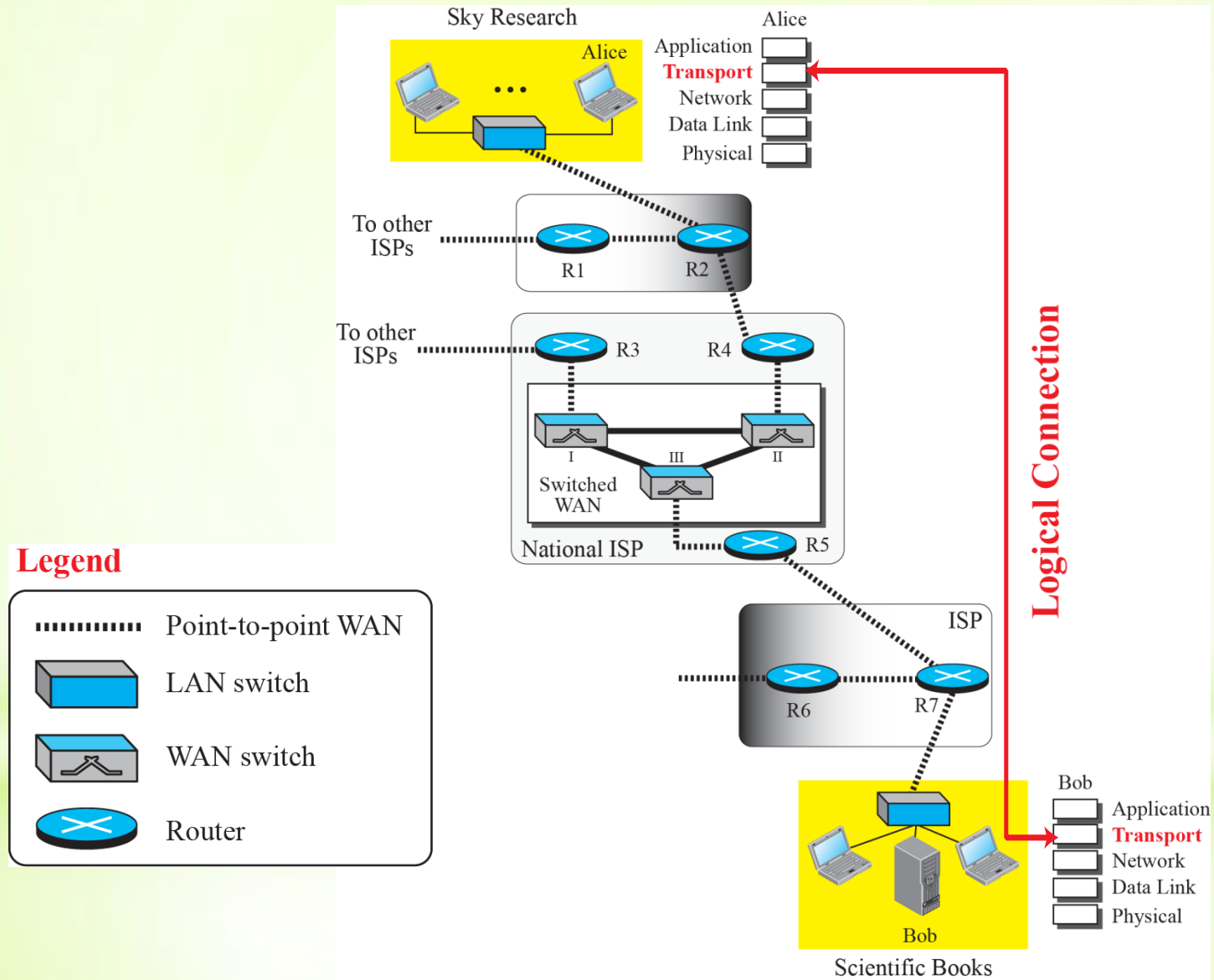
- ❑ Connectionless and Connection-Oriented

- ❖ Connectionless Service
- ❖ Connection-Oriented Service
- ❖ Finite State Machine

INTRODUCTION

- The transport layer provides a process-to-process communication between two application layers.
- Communication is provided using a logical connection, which means that the two Transport layers assume that there is an imaginary direct connection through which they can send and receive messages.
- The transport layer is located between the network layer and the application layer. The transport layer is responsible for providing services to the application layer; it receives services from the network layer.

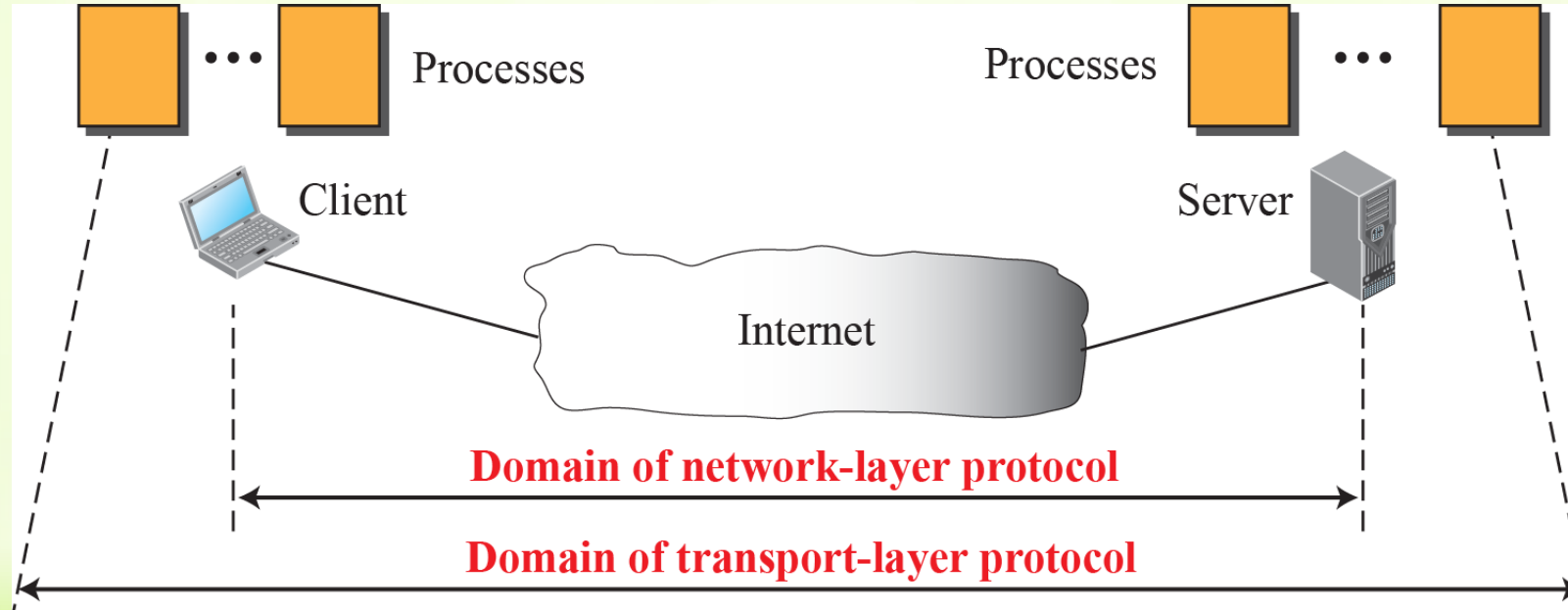
Logical connection at the transport layer



Process-to-Process Communication

- The first duty of a transport-layer protocol is to provide process-to-process communication. A process is an application-layer entity (running program) that uses the services of the transport layer.
- A network-layer protocol can deliver the message only to the destination computer. However, this is an incomplete delivery. The message still needs to be handed to the correct process. This is where a transport-layer protocol takes over. A transport-layer protocol is responsible for delivery of the message to the appropriate process.

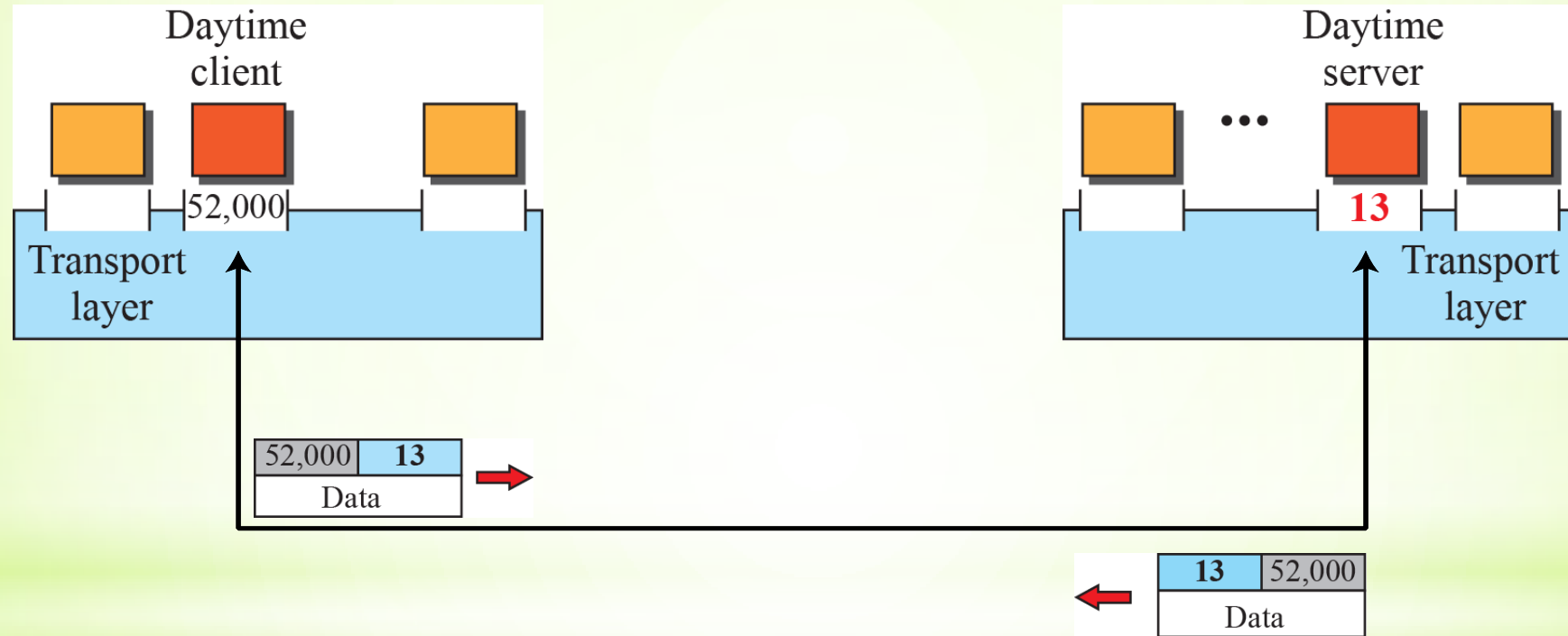
Network layer versus transport layer



Addressing: Port Numbers

- A process on the local host, called a client, needs services from a process usually on the remote host, called a server.
- The local host and the remote host are defined using IP addresses
- To define the processes, we need second identifiers, called port numbers.
- In the TCP/IP protocol suite, the port numbers are integers between 0 and 65,535 (16 bits).
- ICANN has divided the port numbers into three ranges:
 - **Well-known ports.** The ports ranging from 0 to 1,023 are assigned and controlled by ICANN. These are the well-known ports.
 - **Registered ports.** The ports ranging from 1,024 to 49,151 are not assigned or controlled by ICANN. They can only be registered with ICANN to prevent duplication.
 - **Dynamic ports.** The ports ranging from 49,152 to 65,535 are neither controlled nor registered. They can be used as temporary or private port numbers

Port numbers



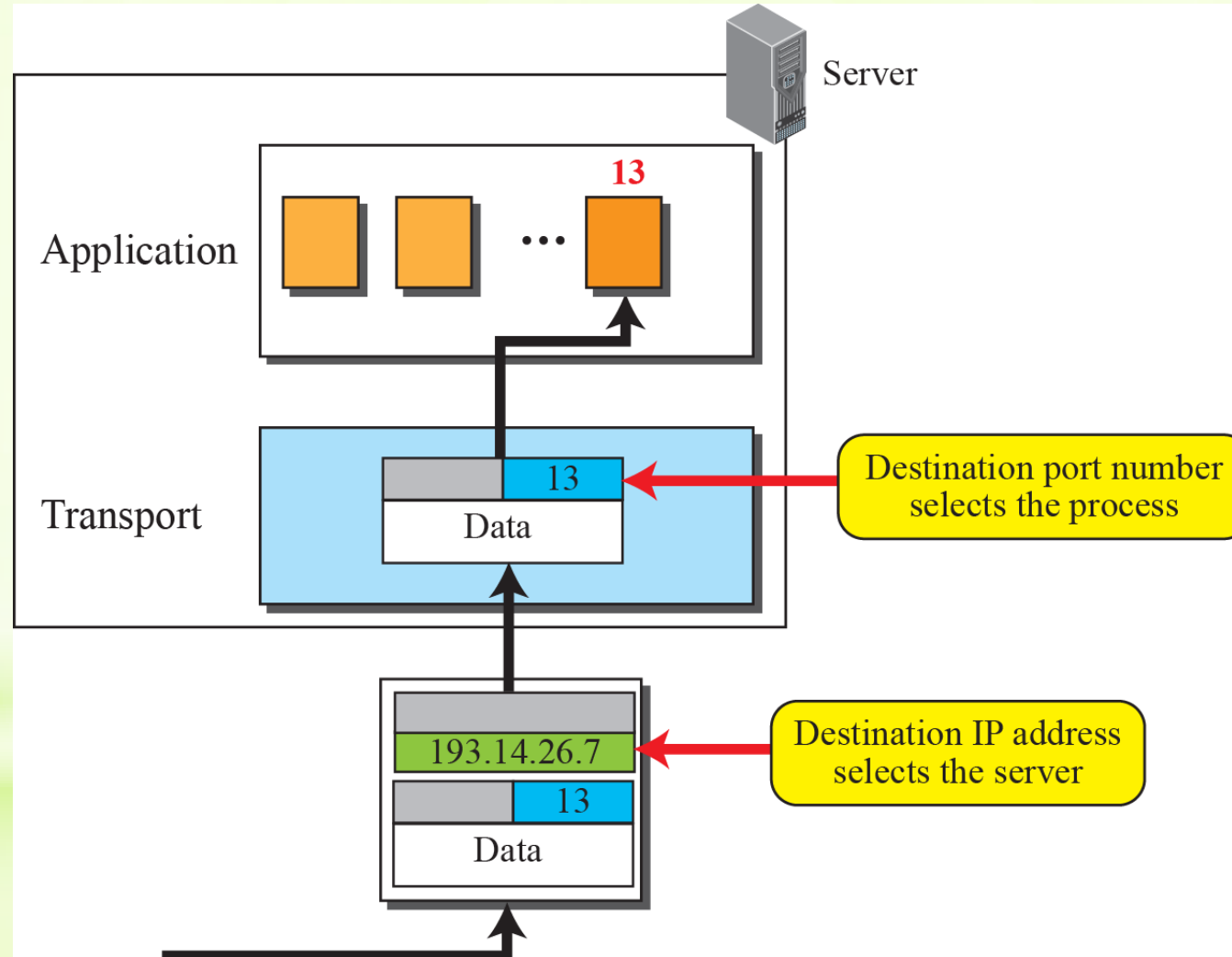
ICANN ranges



The client program defines itself with a port number, called the ephemeral port number.

The server process must also define itself with a port number. This port number, however, cannot be chosen randomly. TCP/IP has decided to use universal port numbers for servers; these are called well-known port numbers (for standardized Client-server applications)

IP addresses versus port numbers



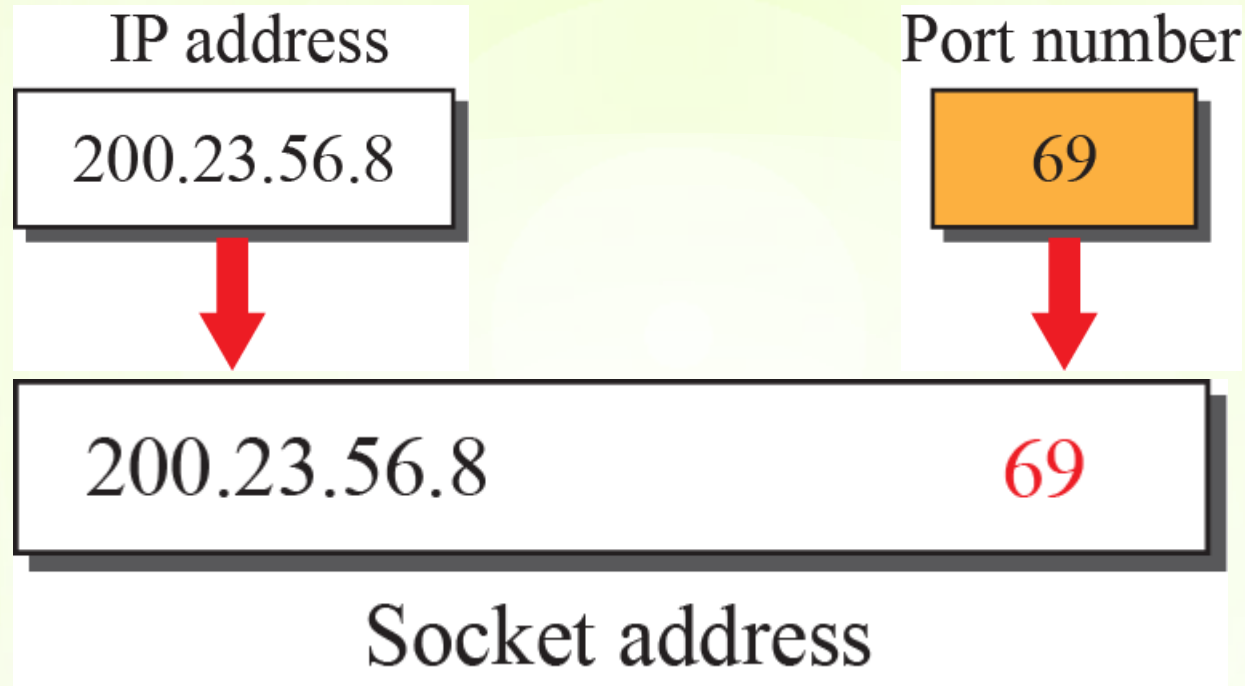
In UNIX, the well-known ports are stored in a file called `/etc/services`. We can use the *grep* utility to extract the line corresponding to the desired application.

```
$grep tftp/etc/services  
tftp 69/tcp  
tftp 69/udp
```

SNMP (see Chapter 9) uses two port numbers (161 and 162), each for a different purpose.

```
$grep snmp/etc/services  
snmp161/tcp#Simple Net Mgmt Proto  
snmp161/udp#Simple Net Mgmt Proto  
snmptrap162/udp#Traps for SNMP
```

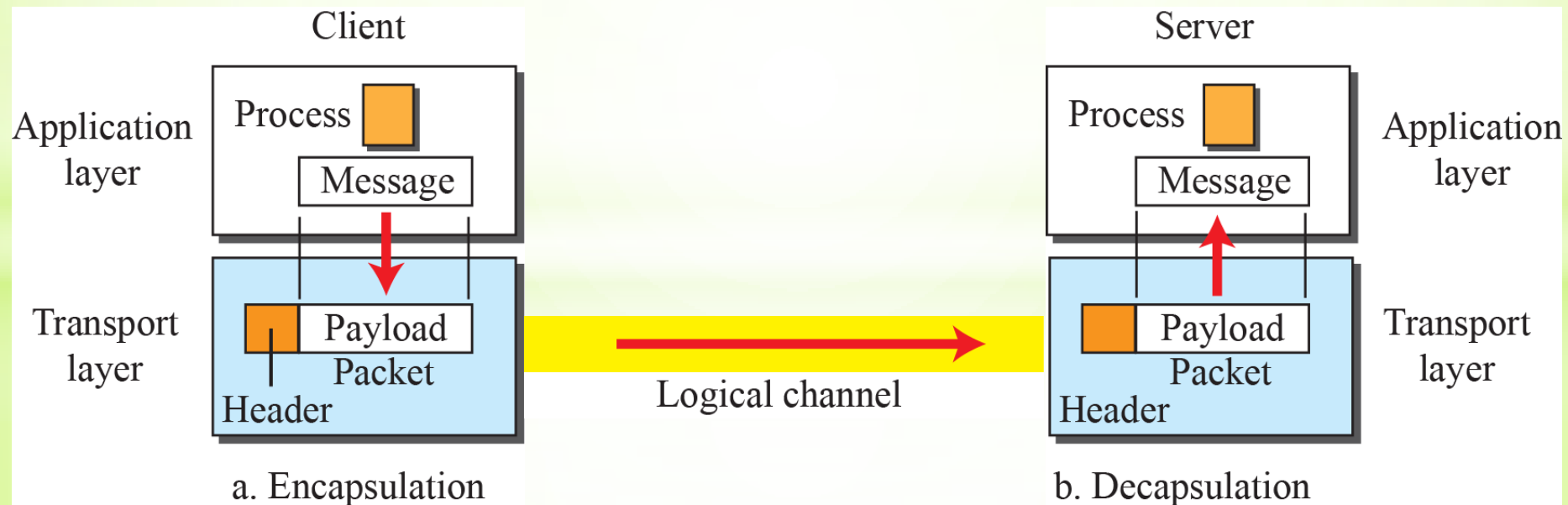
Socket address



- To use services in the Internet, we need a pair of socket addresses: the client socket address and the server socket address at each end, to make a connection.
- The combination of an IP address and a port number is called a socket address

Encapsulation and decapsulation

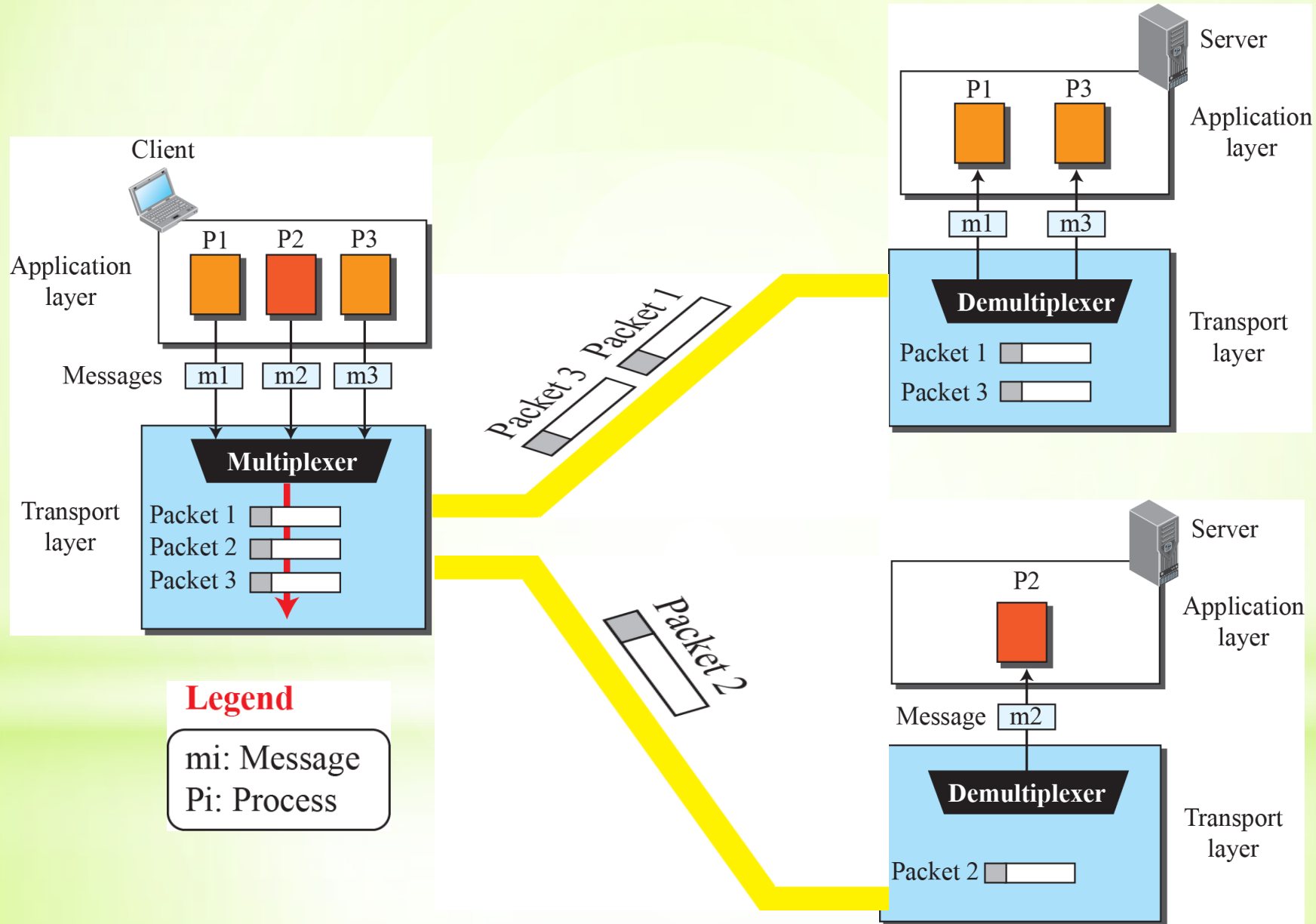
- Encapsulation happens at the sender site. The transport layer receives the data and adds the transport-layer header(i.e. port address, and other information). The packets at the transport layers in the Internet are called **user datagrams, segments, or packets**, depending on what transport-layer protocol we use.
- Decapsulation happens at the receiver site. When the message arrives at the destination transport layer, the header is dropped and the transport layer delivers the message to the process running at the application layer. The sender socket address is passed to the process in case it needs to respond to the message received.



Multiplexing and demultiplexing

- Whenever an entity accepts items from more than one source, this is referred to as multiplexing (many to one); whenever an entity delivers items to more than one source, this is referred to as demultiplexing (one to many).
- The transport layer at the source performs multiplexing; the transport layer at the destination performs demultiplexing.
- The transport layer at the client site accepts three messages from the three processes and creates three packets. It acts as a multiplexer.
- When they arrive at the server, the transport layer does the job of a demultiplexer and distributes the messages to two different processes.

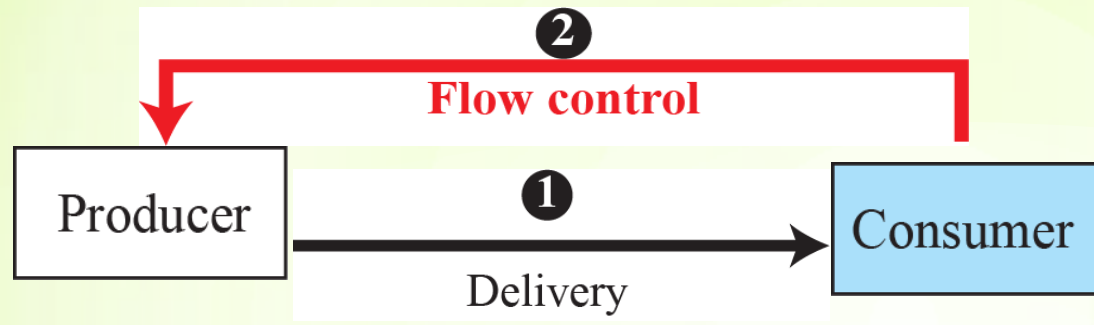
Multiplexing and demultiplexing



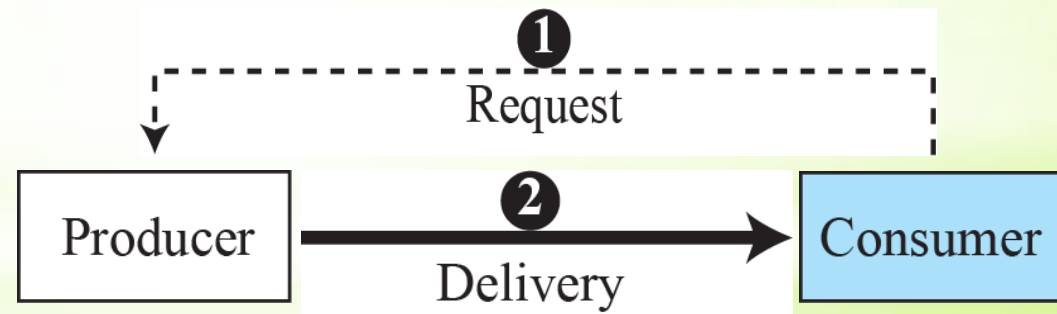
Flow Control

- In communication at the transport layer, we are dealing with four entities: sender process, sender transport layer, receiver transport layer, and receiver process.
- The sending process at the application layer is only a producer. It produces message chunks and pushes them to the transport layer. The sending transport layer has a double role: it is both a consumer and a producer. It consumes the messages pushed by the producer. It encapsulates the messages in packets and pushes them to the receiving transport layer.
- The receiving transport layer also has a double role, it is the consumer for the packets received from the sender and the producer that decapsulates the messages and delivers them to the application layer. The last delivery, however, is normally a pulling delivery; the transport layer waits until the application-layer process asks for messages.

Pushing or pulling

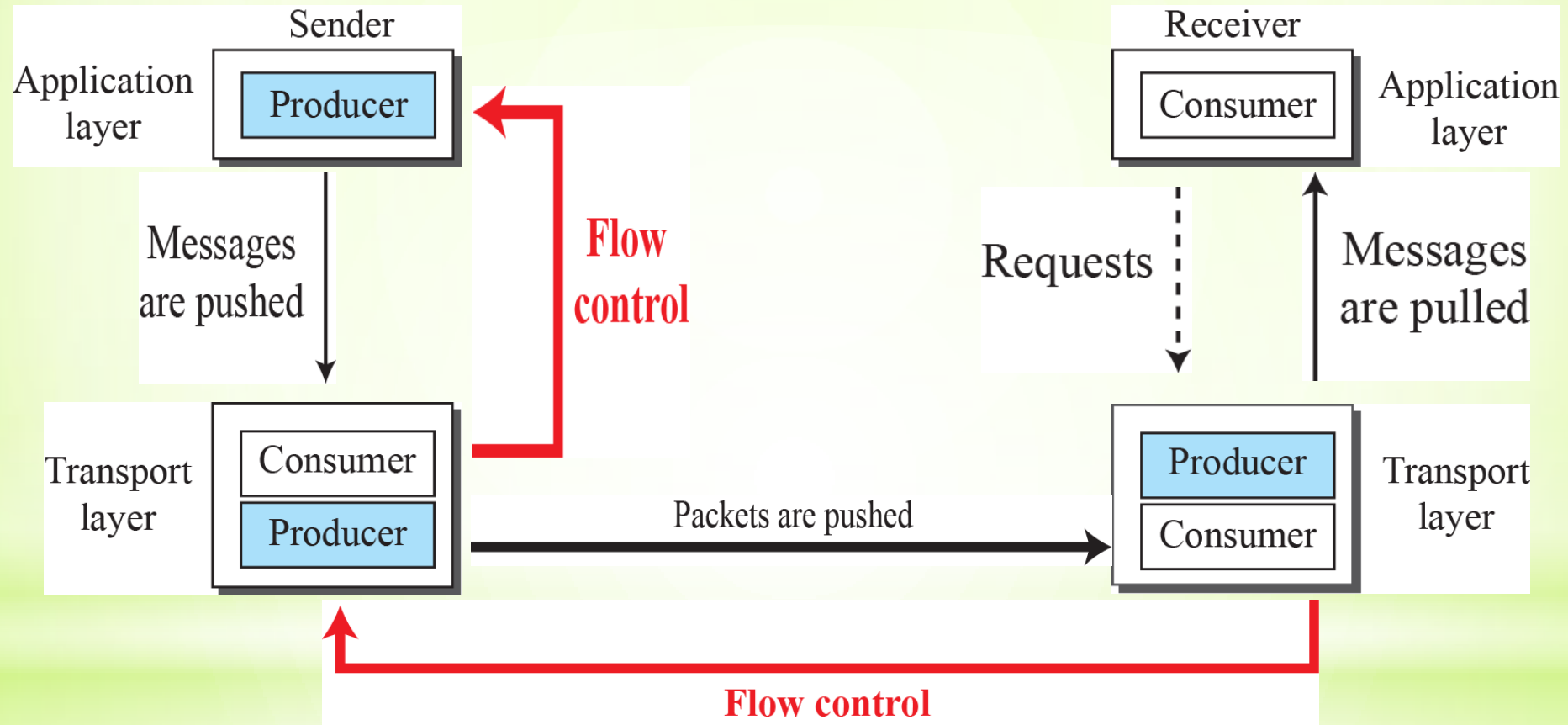


a. Pushing



b. Pulling

Flow control at the transport layer



Buffers

- Although flow control can be implemented in several ways, one of the solutions is normally to use two buffers: one at the sending transport layer and the other at the receiving transport layer.
- A buffer is a set of memory locations that can hold packets at the sender and receiver.
- When the buffer of the sending transport layer is full, it informs the application layer to stop passing chunks of messages; when there are some vacancies, it informs the application layer that it can pass message chunks again.
- When the buffer of the receiving transport layer is full, it informs the sending transport layer to stop sending packets. When there are some vacancies, it informs the sending transport layer that it can send packets again.

Error Control

- In the Internet, since the underlying network layer (IP) is unreliable, we need to make the transport layer reliable if the application requires reliability.
- Reliability can be achieved to add error control services to the transport layer. Error control at the transport layer is responsible for:
 1. Detecting and discarding corrupted packets.
 2. Keeping track of lost and discarded packets and resending them.
 3. Recognizing duplicate packets and discarding them.
 4. Buffering out-of-order packets until the missing packets arrive.

SEQUENCE NUMBER

- To perform error control, the packets are numbered. We can add a field to the transport-layer packet to hold the sequence number of the packet. Packets are numbered sequentially. If the header of the packet allows m bits for the sequence number, the sequence numbers range from 0 to $2^m - 1$.

ACKNOWLEDGMENT NUMBER

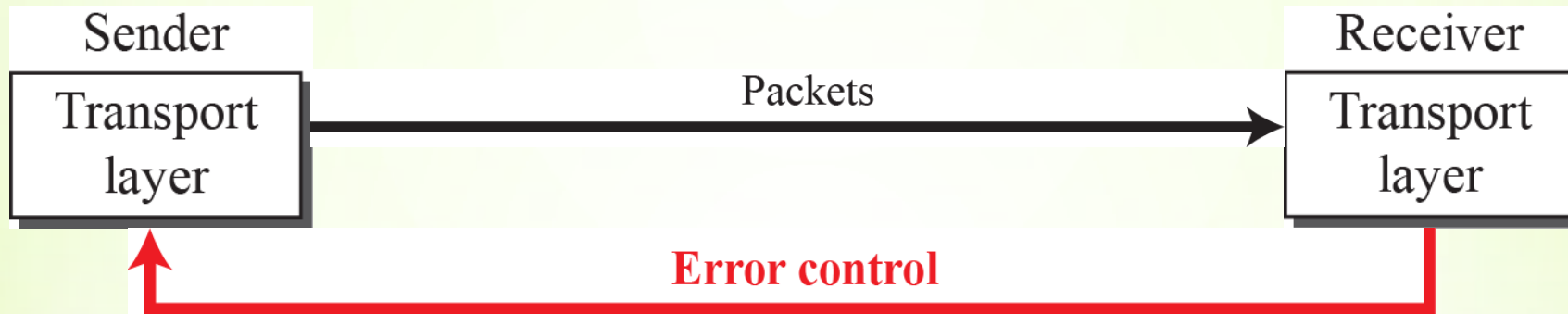
- We can use both positive and negative signals as error control
- The receiver side can send an acknowledgment (ACK) for each of a collection of packets that have arrived safe and sound. The receiver can simply discard the corrupted packets.

TIMERS

- The sender can detect lost packets if it uses a timer.
- When a packet is sent, the sender starts a timer. If an ACK does not arrive before the timer expires, the sender resends the packet.

Duplicate packets can be silently discarded by the receiver. Out-of-order packets can be either discarded (to be treated as lost packets by the sender), or stored until the missing ones arrives.

Error control at the transport layer



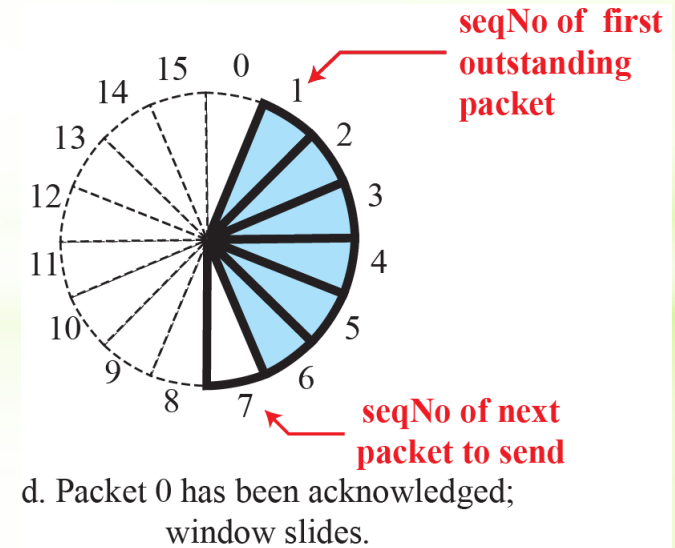
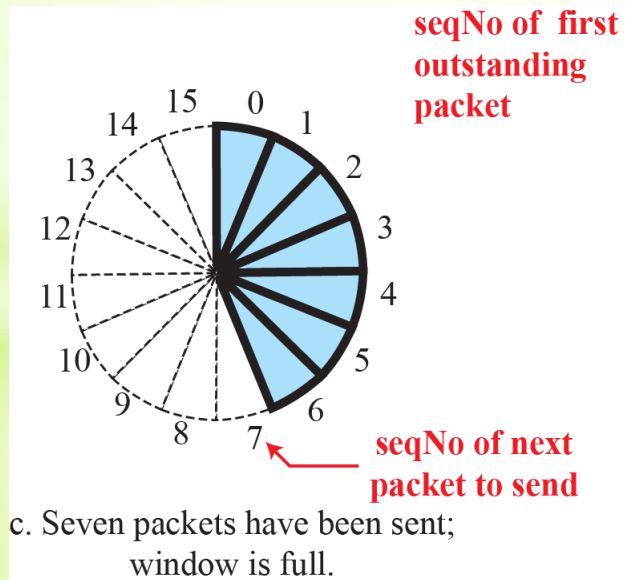
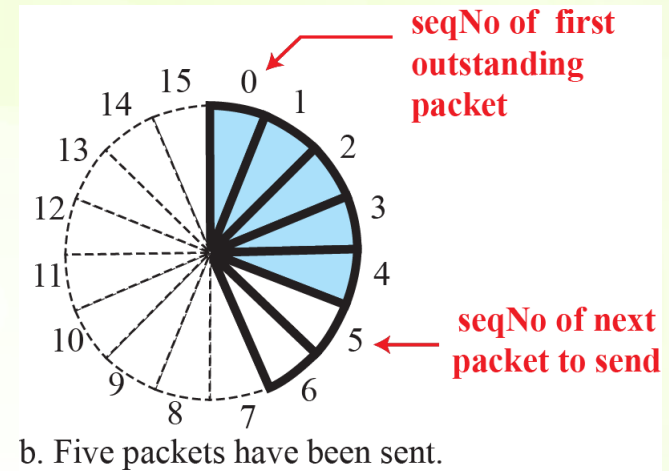
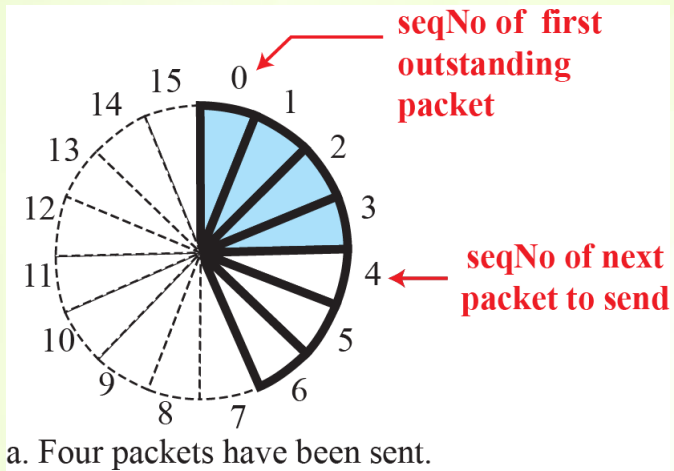
COMBINATION OF FLOW AND ERROR CONTROL

- Flow control requires the use of two buffers, one at the sender site and the other at the receiver site.
- Error control requires the use of sequence and acknowledgment numbers by both sides.
- These two requirements can be combined if we use two numbered buffers, one at the sender, one at the receiver.
- At the sender, when a packet is prepared to be sent, we use the number of the next free location, x , in the buffer as the sequence number of the packet. When the packet is sent, a copy is stored at memory location x , awaiting the acknowledgment from the other end. When an acknowledgment related to a sent packet arrives, the packet is purged and the memory location becomes free.
- At the receiver, when a packet with sequence number y arrives, it is stored at the memory location y until the application layer is ready to receive it. An acknowledgment can be sent to announce the arrival of packet y .

SLIDING WINDOW

- Since the sequence numbers used modulo 2^m , a circle can represent the sequence numbers from 0 to $2^m - 1$
- The buffer is represented as a set of slices, called the sliding window, that occupies part of the circle at any time.
- At the sender site, when a packet is sent, the corresponding slice is marked. When all the slices are marked, it means that the buffer is full and no further messages can be accepted from the application layer.
- When an acknowledgment arrives, the corresponding slice is unmarked.
- If some consecutive slices from the beginning of the window are unmarked, the window slides over the range of the corresponding sequence numbers to allow more free slices at the end of the window.
- Most protocols show the sliding window using linear representation.

Sliding window in circular format



Sliding window in linear format



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.

Congestion Control

- An important issue in a packet-switched network, such as the Internet, is congestion.
- Congestion in a network may occur if the load on the network—the number of packets sent to the network—is greater than the capacity of the network—the number of packets a network can handle.
- Congestion control refers to the mechanisms and techniques that control the congestion and keep the load below the capacity.
- Congestion at the transport layer is actually the result of congestion at the network layer, which manifests itself at the transport layer
- TCP, assuming that there is no congestion control at the network layer, implements its own congestion control mechanism

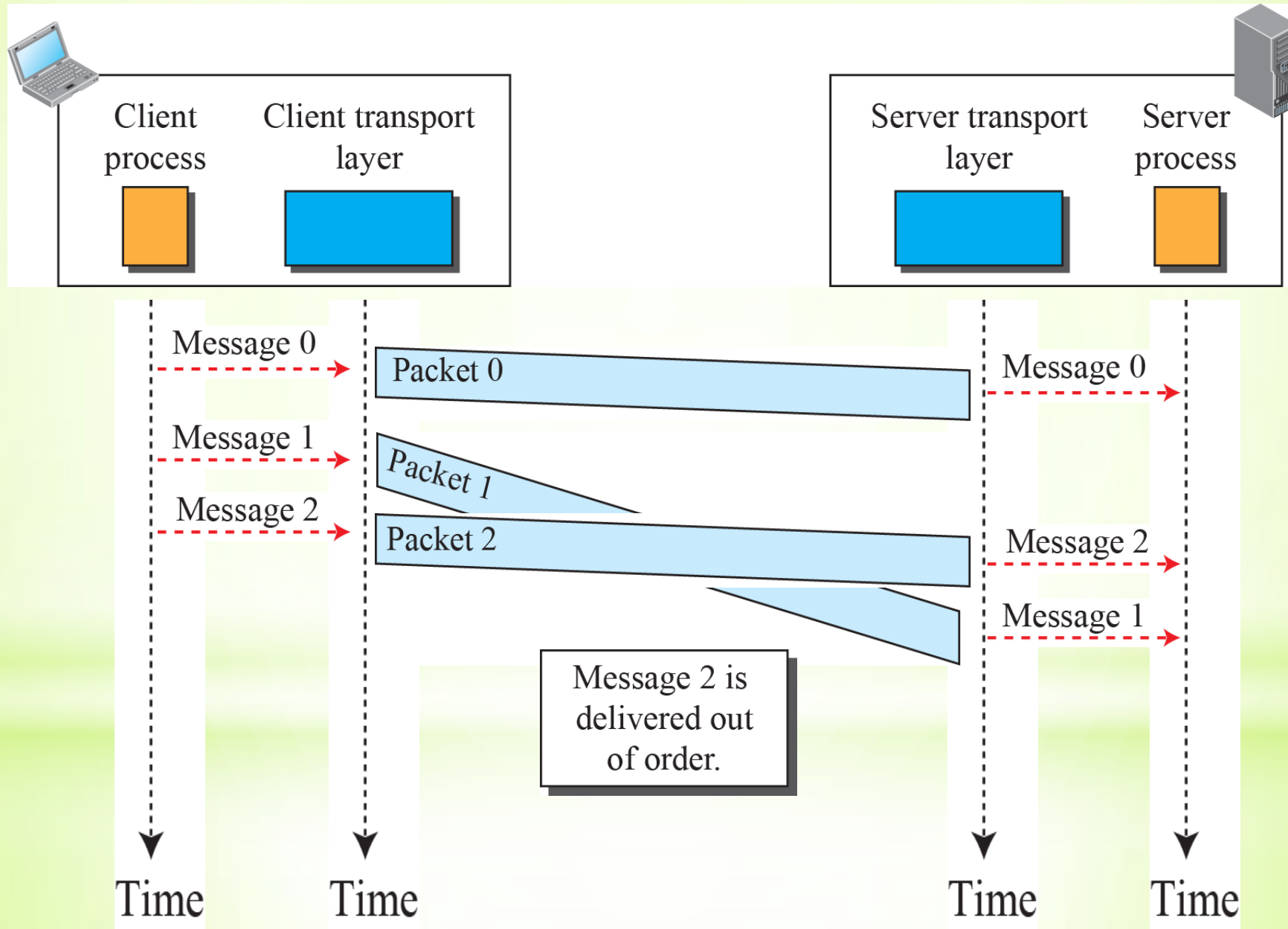
Connectionless and Connection-Oriented Services

- A transport-layer protocol, like a network-layer protocol, can provide two types of services: connectionless and connection-oriented.
- At the transport layer, we are not concerned about the physical paths of packets (we assume a logical connection between two transport layers).
- Connectionless service at the transport layer means independency between packets; connection-oriented means dependency.

CONNECTIONLESS SERVICE

- In a connectionless service, the source process (application program) needs to divide its message into chunks of data of the size acceptable by the transport layer and deliver them to the transport layer one by one.
- The transport layer treats each chunk as a single unit without any relation between the chunks. When a chunk arrives from the application layer, the transport layer encapsulates it in a packet and sends it.
- Problems such as lost packets, out of order delivery exists in such mechanism.
- We can say that no flow control, error control, or congestion control can be effectively implemented in a connectionless service.
- A well known real time protocol used for Connectionless service is UDP.

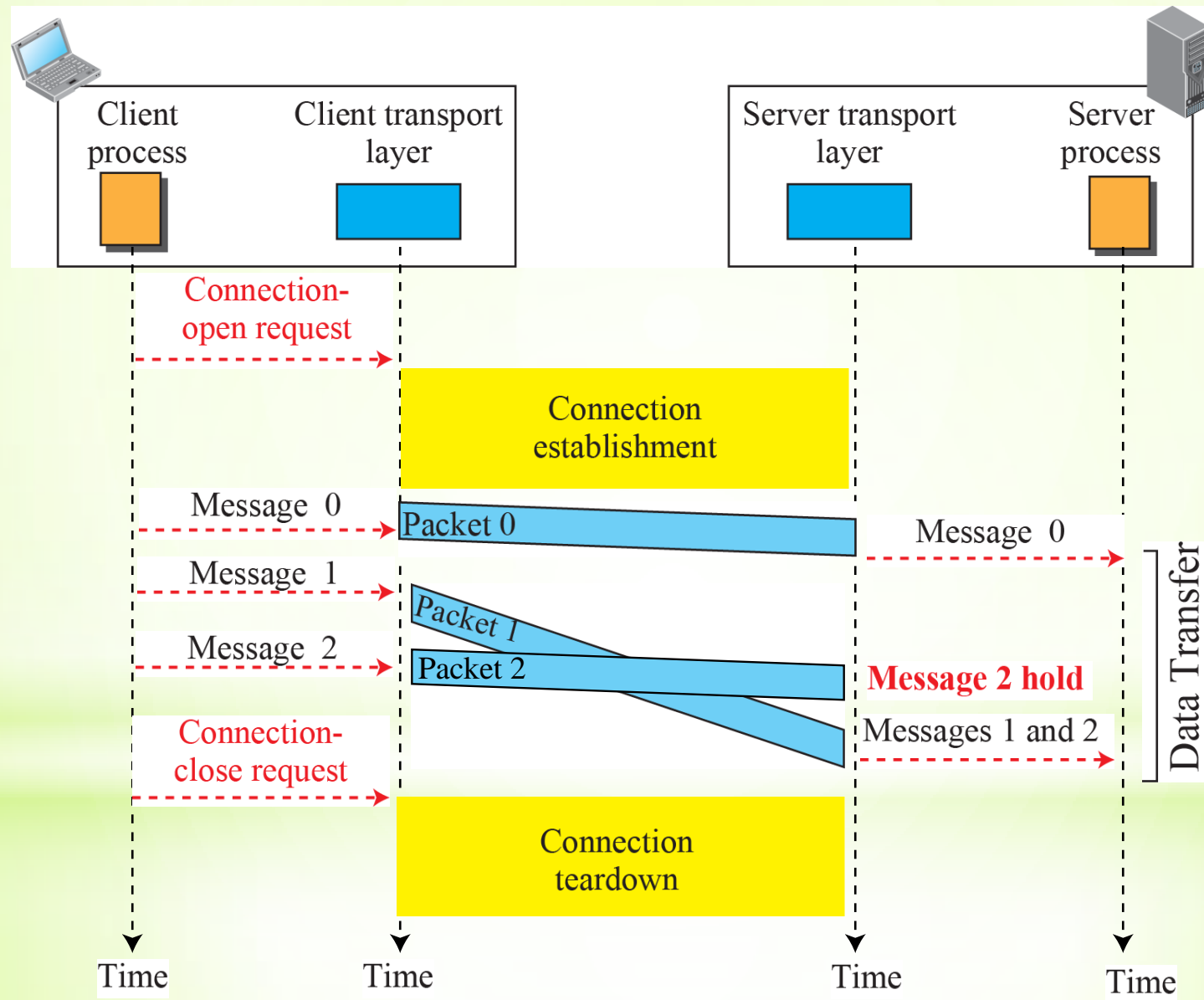
Connectionless service



CONNECTION-ORIENTED SERVICE

- In a connection-oriented service, the client and the server first need to establish a logical connection between themselves.
- The data exchange can only happen after the connection establishment.
- Once completed the data exchange the connection is closed
- We can implement flow control, error control, and congestion control in a connection oriented protocol.
- A well known real time protocol used for Connection oriented service is TCP.

Connection-oriented service

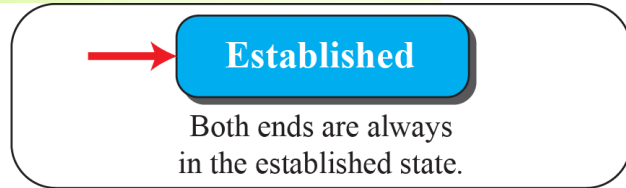


FINITE STATE MACHINE

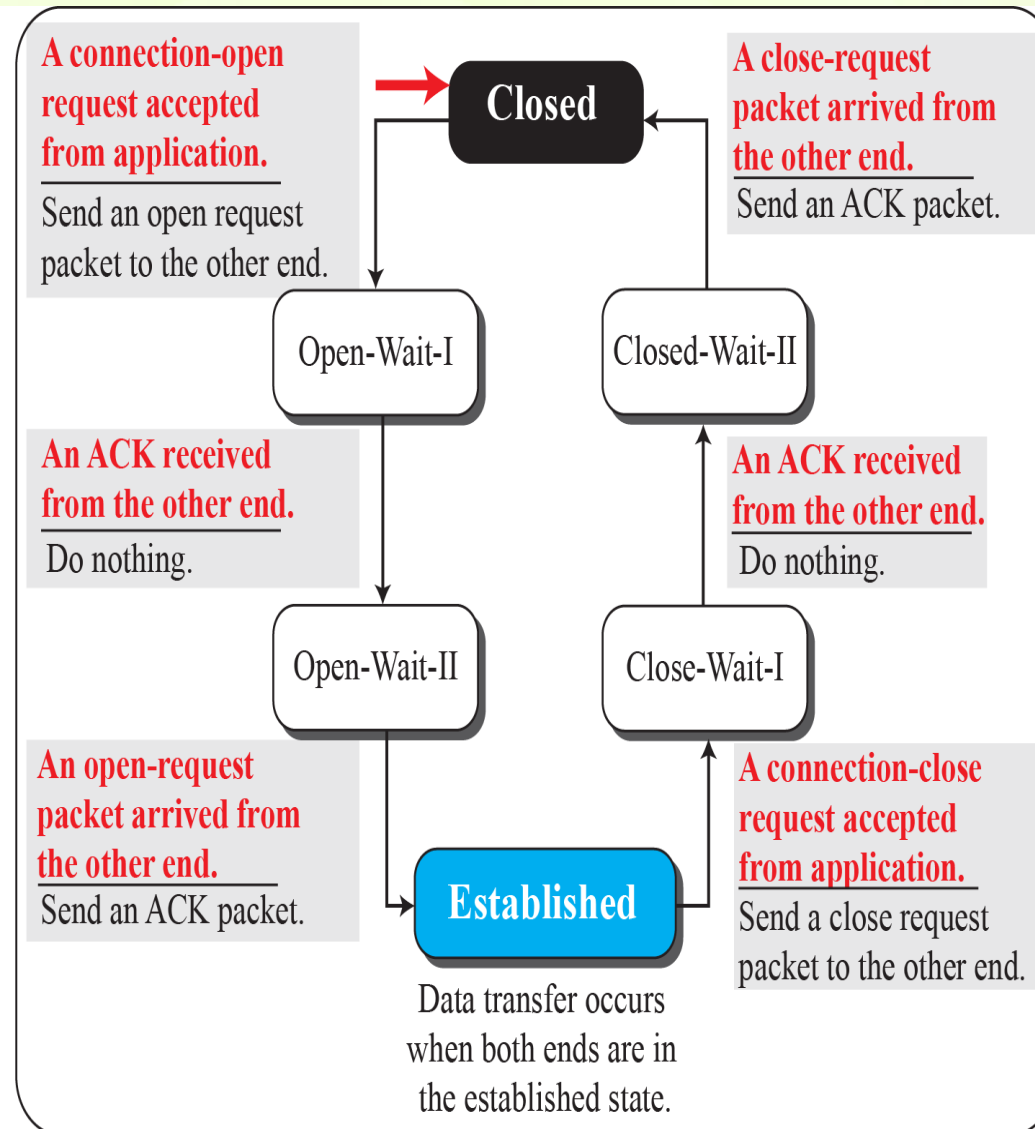
- The behavior of a transport-layer protocol, both when it provides a connectionless and when it provides a connection-oriented protocol, can be better shown as a **finite state machine (FSM)**.
- In a connectionless service mechanism, there is only one state i.e. both ends are always in the established state.
- In a connection oriented mechanism, there are totally six states:
 - Open-wait-I
 - Open-wait-II
 - Established
 - Close-wait-I
 - Close-wait-II
 - Closed

Connectionless and connection-oriented service represented as FSMs

FSM for
connectionless
transport layer



FSM for
connection-oriented
transport layer



Note:

The colored arrow shows the starting state.

UNIT - 4

Transport Layer

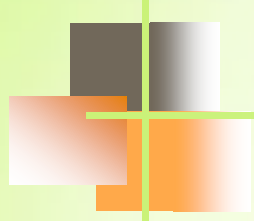
3.1 INTRODUCTION

3.2 TRANSPORT-LAYER PROTOCOLS

- Simplex Protocol
- Stop-and-Wait Protocol
- Go-Back-N Protocol (GBN)
- Selective-Repeat Protocol
- Bidirectional Protocols: Piggybacking

3.3 INTERNET TRANSPORT-LAYER PROTOCOLS

- User Datagram Protocol (UDP)
- Transmission Control Protocol (TCP)

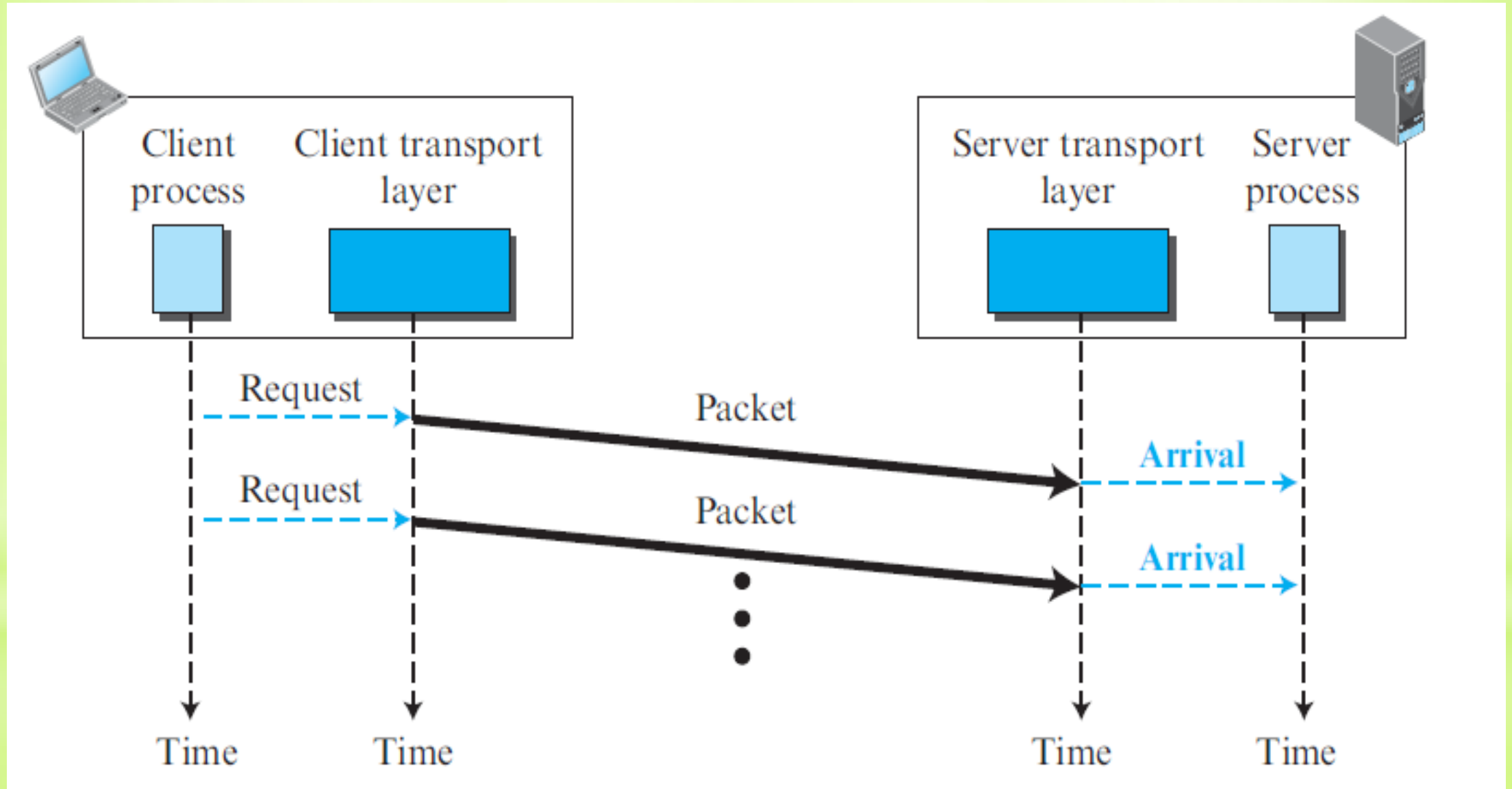


- ☐ Sequence Numbers
- ☐ Acknowledgment Numbers
- ☐ Send Window
- ☐ Receive Window
- ☐ Timers
- ☐ Resending packets
- ☐ FSMs (Sender & Receiver)

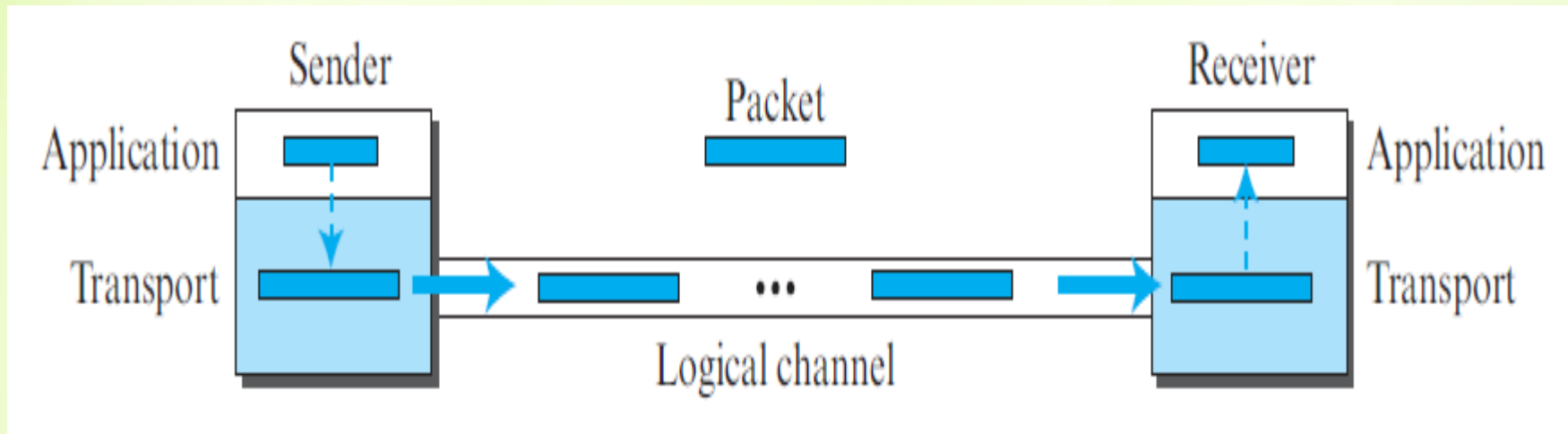
Simplex Protocol

- Used in Noiseless channels
- Simple connectionless protocol with neither flow nor error control
- The transport layer at the sender gets a message from its application layer, makes a packet out of it, and sends the packet.
- The transport layer at the receiver receives a packet from its network layer, extracts the message from the packet, and delivers the message to its application layer.
- No acknowledgement and no sequence number.
- Each FSM has only one state, the ready state.

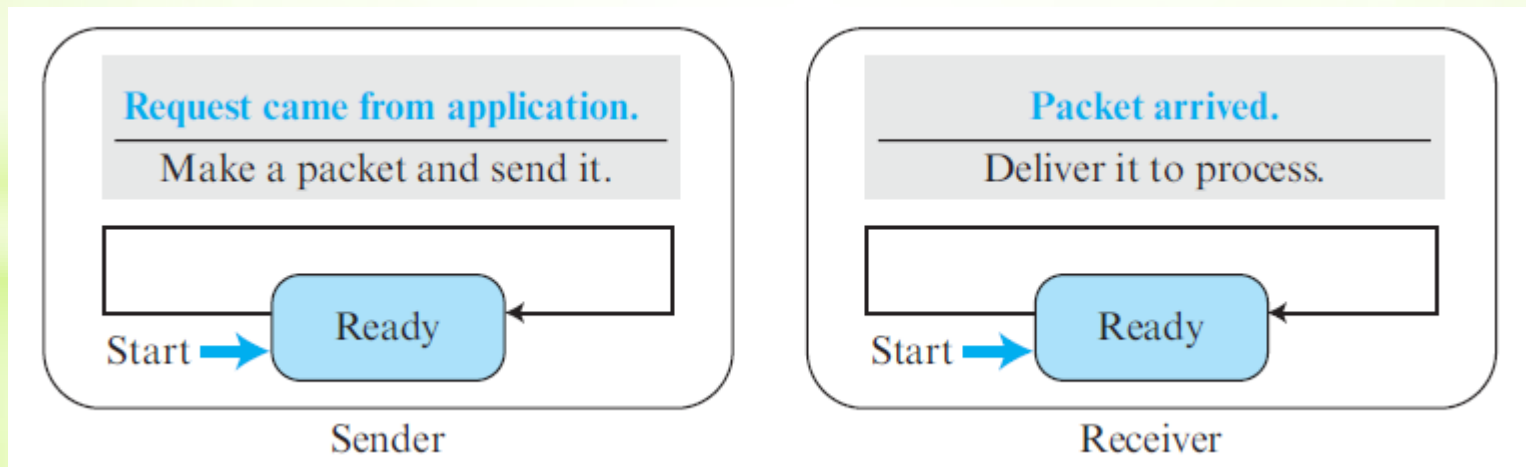
Flow diagram



Simplex protocol



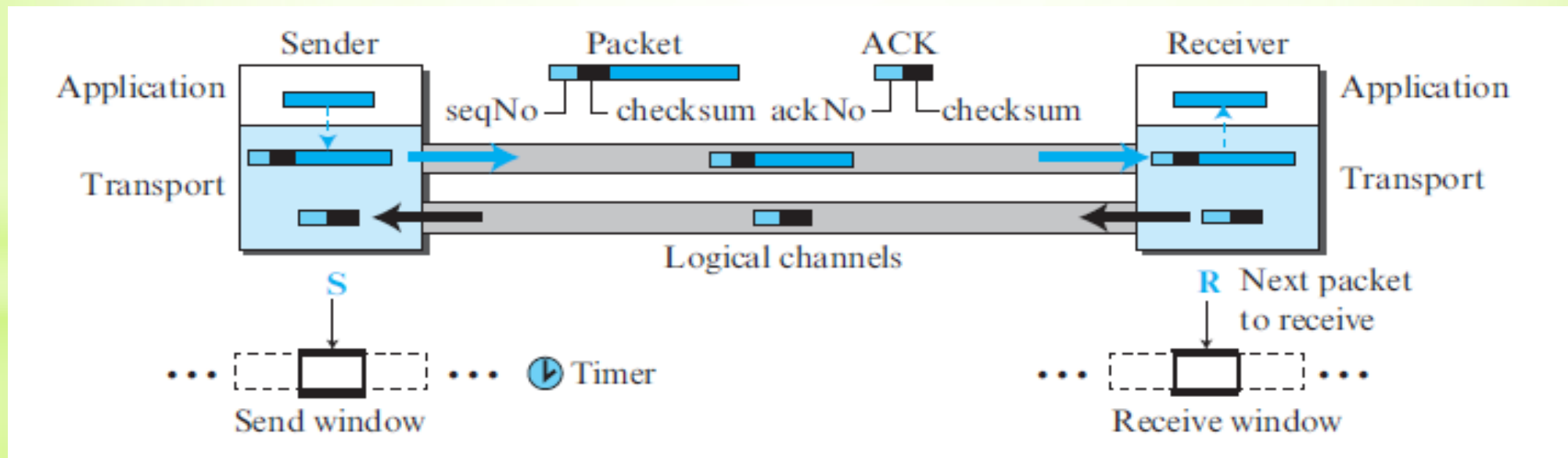
FSMs for the simplex protocol



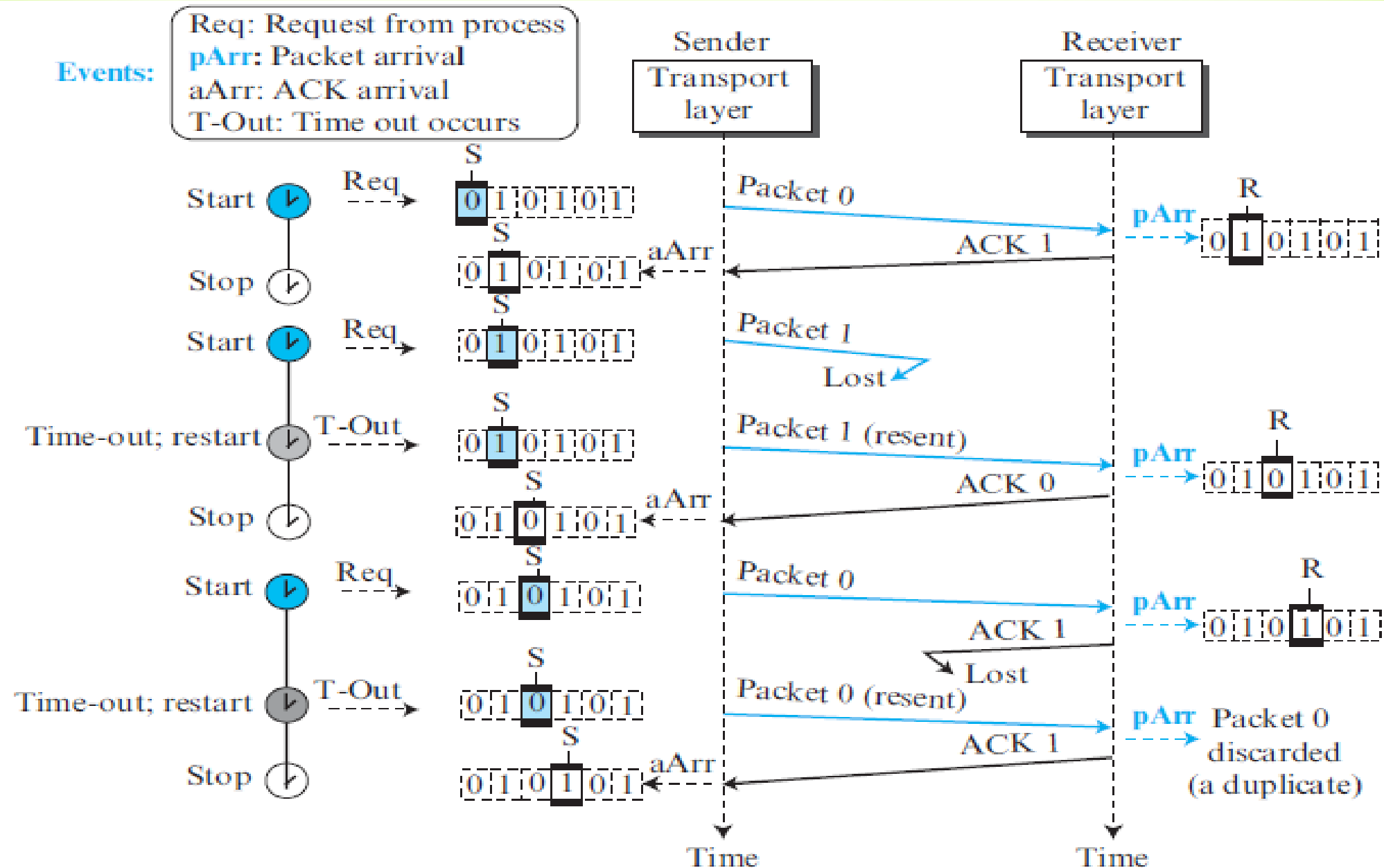
Stop and Wait protocol

- It is a connection-oriented protocol
- Uses both flow and error control
- Both the sender and the receiver use a sliding window of size 1.
- The sender sends one packet at a time and waits for an acknowledgment before sending the next one.
- To detect corrupted packets, we need to add a checksum to each data packet, if checksum not correct the packet is discarded.
- Every time the sender sends a packet, it starts a timer. If an acknowledgment arrives before the timer expires, the timer is stopped and the sender sends the next packet (if it has one to send).
- If the timer expires, the sender resends the previous packet, assuming that the packet was either lost or corrupted. This means that the sender needs to keep a copy of the packet until its acknowledgment arrives.

- To prevent duplicate packets, the protocol uses sequence numbers and acknowledgment Numbers
- The sender is initially in the **ready state**, but it can move between the ready and **blocking state**.
- The receiver is always in the **ready state**

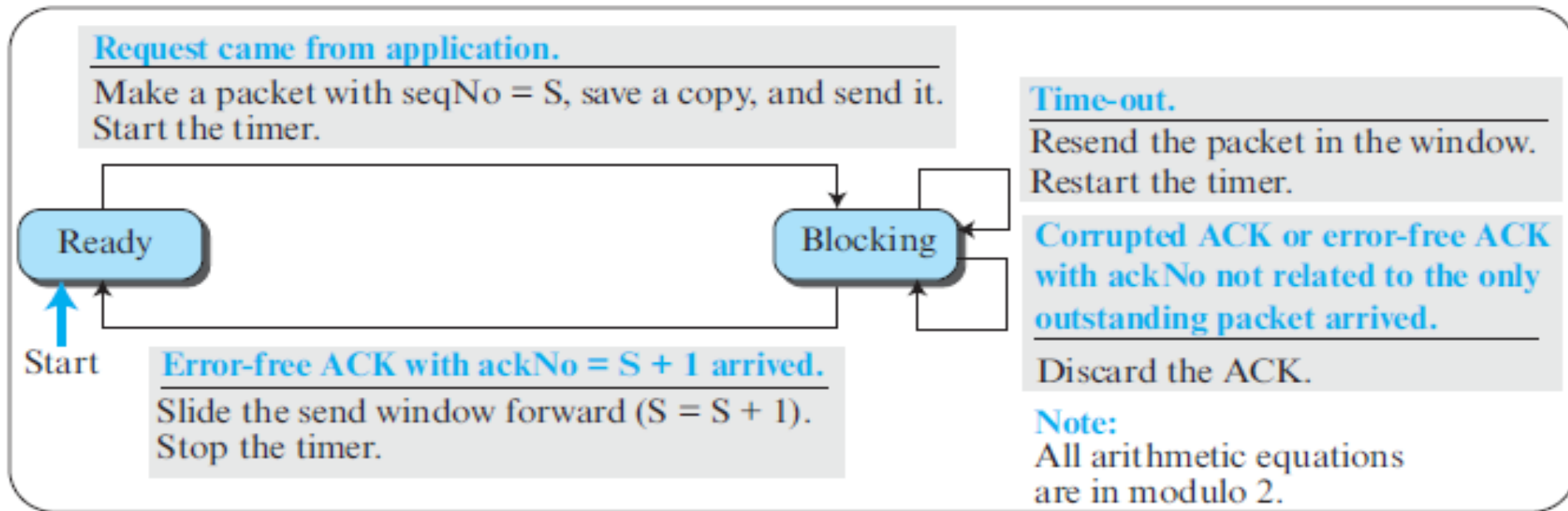


Flow diagram of stop and wait

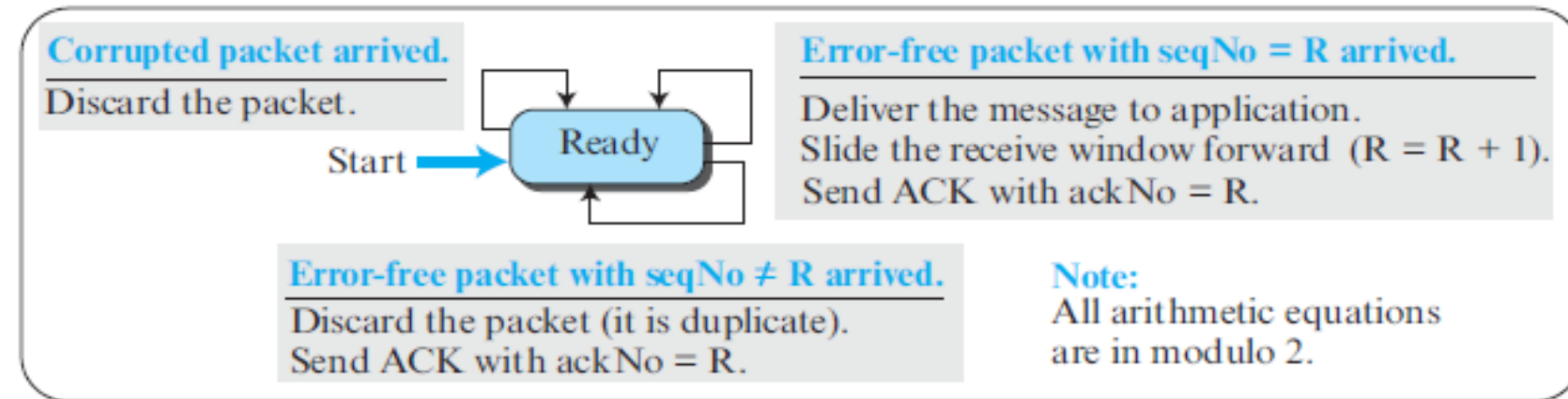


FSM for the Stop-and-Wait protocol

Sender



Receiver



Shortfalls of Stop and Wait Protocol:

- The Stop-and-Wait protocol is very inefficient if our channel is thick and long.
- By **thick**, we mean that our channel has a **large bandwidth** (high data rate); by **long**, we mean the **round-trip delay** is long.
- The product of these two is called the **bandwidth-delay product**.

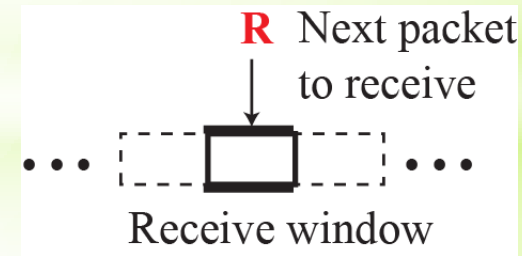
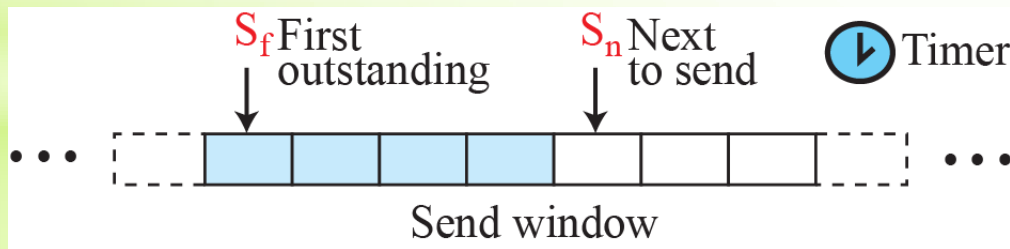
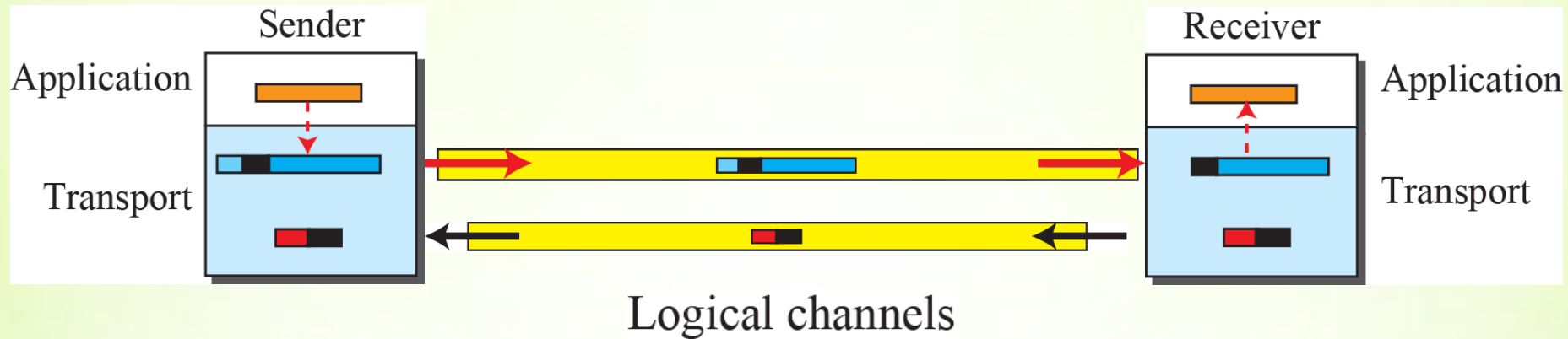
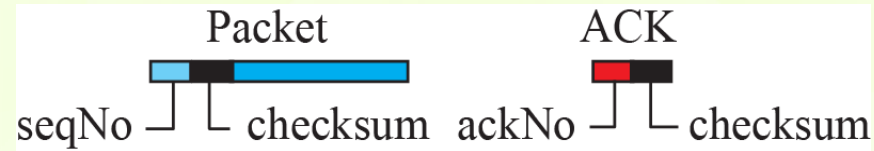
Pipelining:

- In networking and in other areas, a task is often begun before the previous task has ended. This is known as pipelining.
- There is no pipelining in the Stop-and-Wait protocol because a sender must wait for a packet to reach the destination and be acknowledged before the next packet can be sent.

Go-Back-N protocol

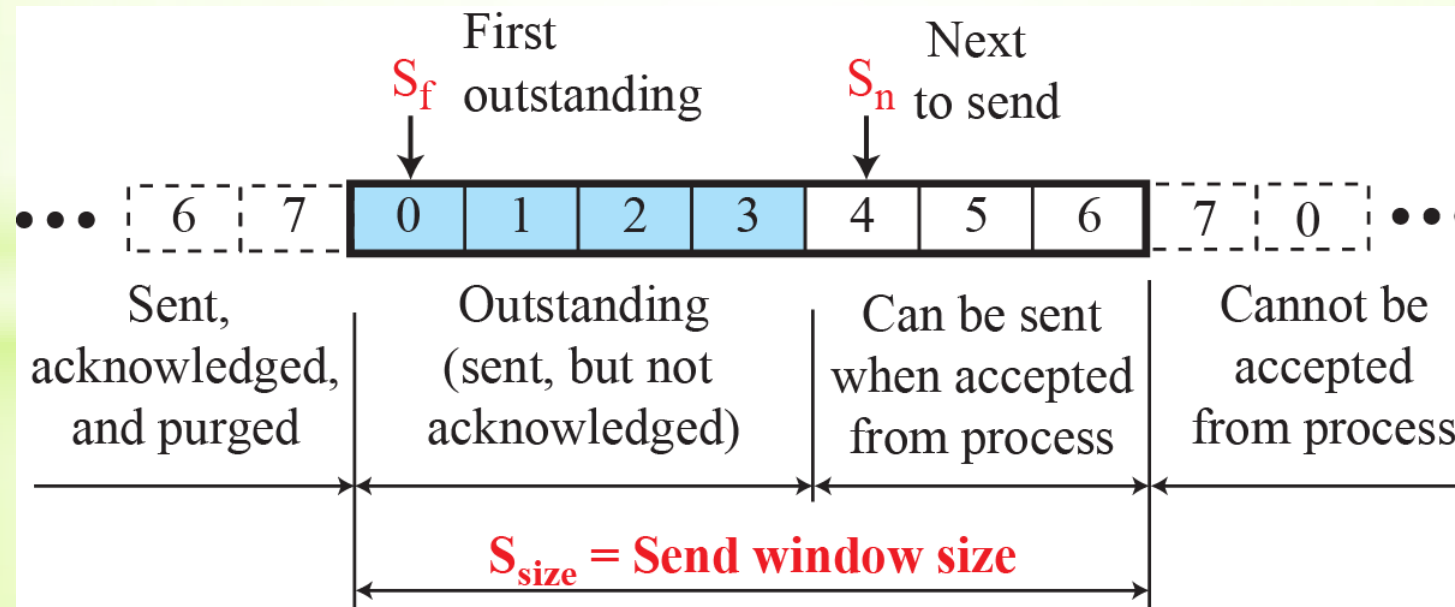
- To improve the efficiency of transmission, multiple packets must be in transition while the sender is waiting for acknowledgment. i.e. using the concept of pipelining.
- A simple protocol that can achieve this goal is called Go-Back-N (GBN)
- The key to Go-back-N is that we can send several packets before receiving acknowledgments, but the receiver can only buffer one packet.
- We keep a copy of the sent packets until the acknowledgments arrive.
- The sequence numbers are modulo 2^m , where 'm' is the size of the sequence number field in bits.
- An acknowledgment number in this protocol is cumulative and defines the sequence number of the next packet expected

Go-Back-N protocol

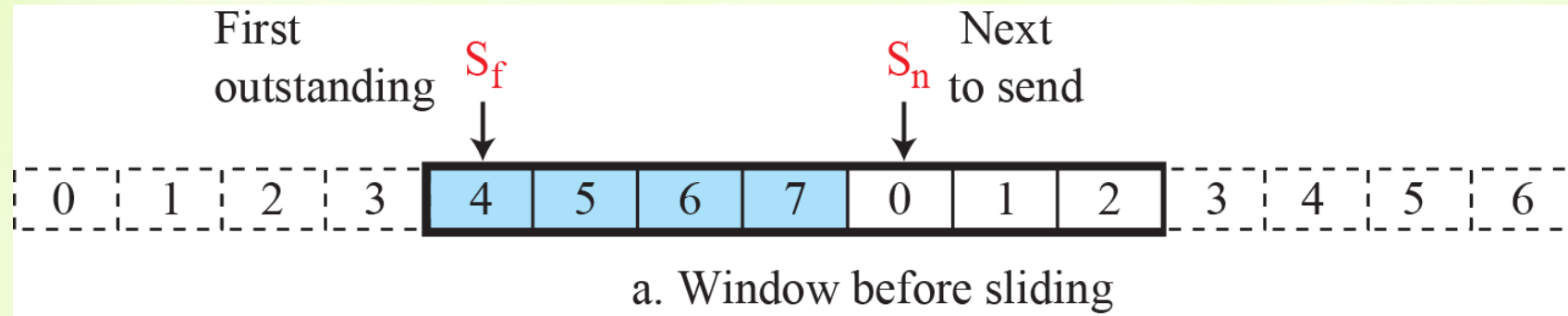


Send window for Go-Back-N

- The maximum **size of the senders window is $2^m - 1$** .
- The sender needs to wait to find out if the packets sent have been received or were lost. We call these **outstanding packets**.
- Three variables are defined:
 - S_f (send window, the first outstanding packet),
 - S_n (send window, the next packet to be sent), and
 - S_{size} (send window, size) = $2^m - 1$.

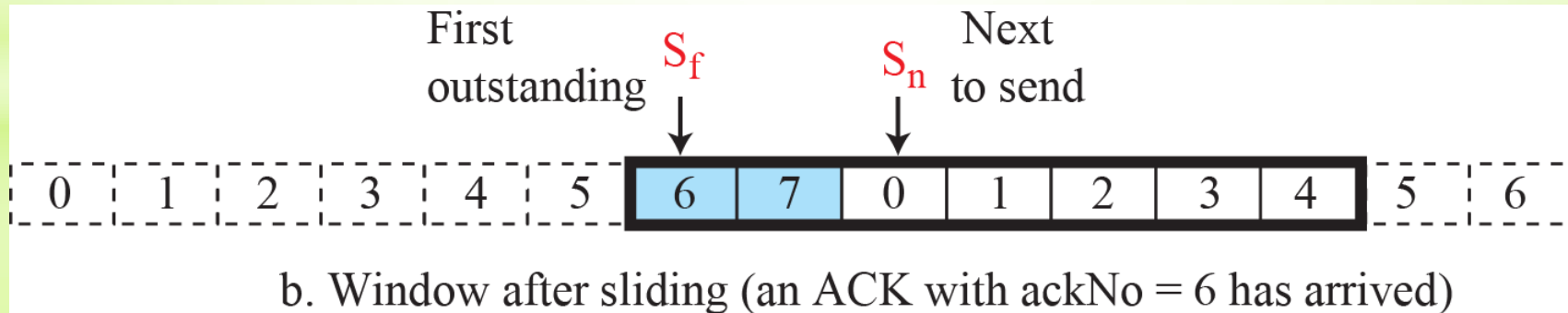


Sliding the send window



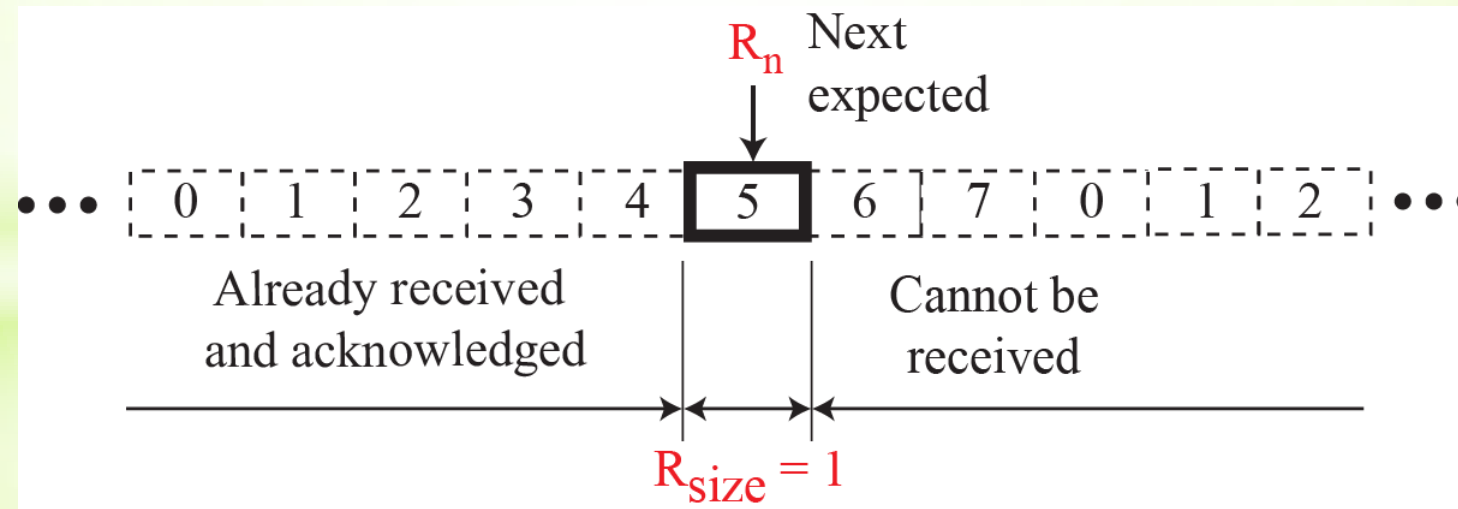
The send window can slide one or more slots when an error-free ACK with ackNo greater than or equal S_f and less than S_n (in modular arithmetic) arrives.

→ *Sliding direction*

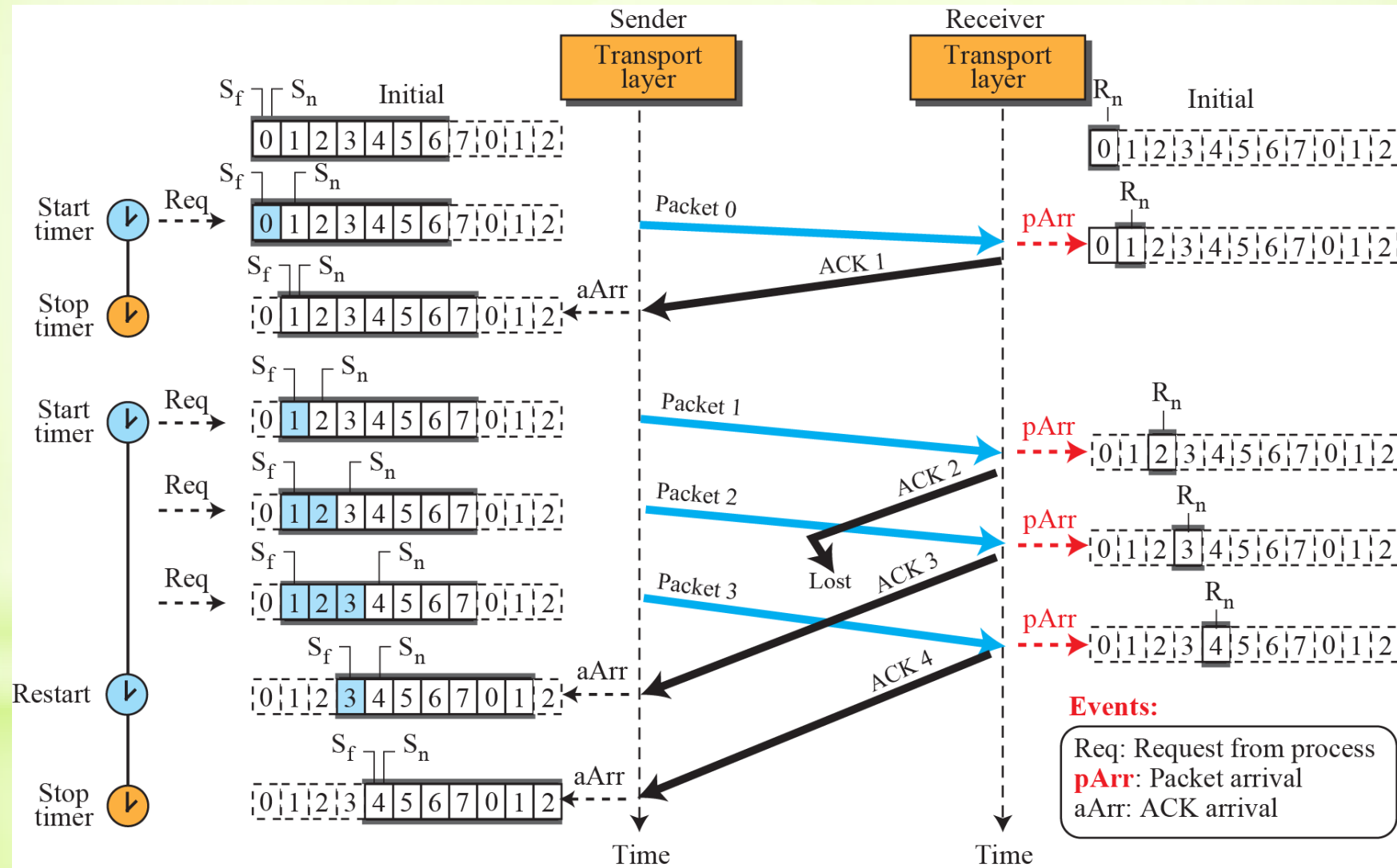


Receive window for Go-Back-N

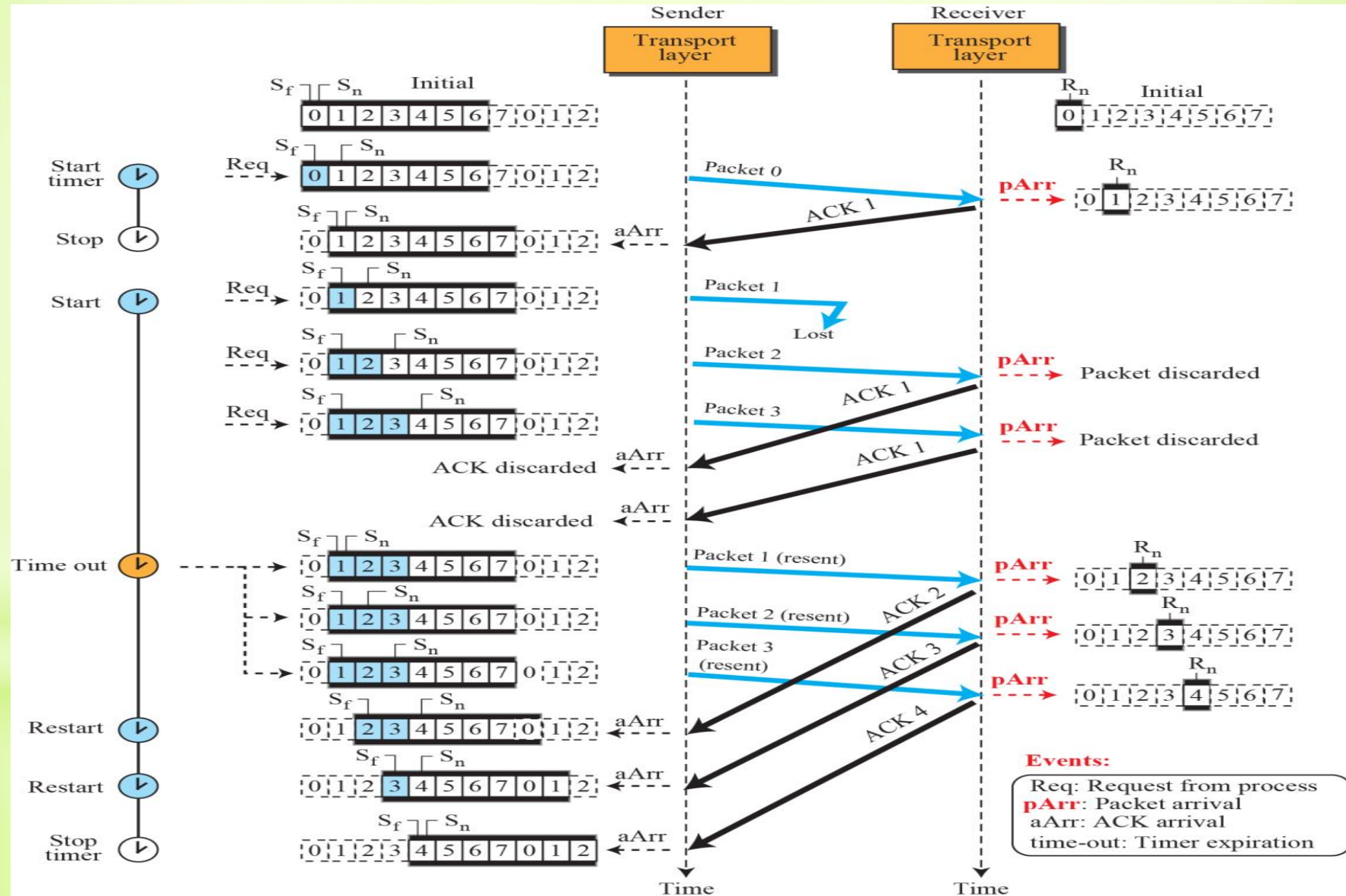
- The receive window makes sure that the correct data packets are received and that the correct acknowledgments are sent. In Go-back-N, the size of the **receive window is always 1**.
- Only one variable: R_n (receive window, next packet expected).
- Only a packet with a sequence number matching the value of R_n is accepted and acknowledged
- When a correct packet is received, the window slides, $R_n = (R_n + 1) \text{ modulo } 2^m$



Flow diagram for *reliability* on senders side



Flow diagram for *unreliability* in senders site



FSMs for the Go-Back-N protocol

Sender

Note:

All arithmetic equations are in modulo 2^m .

Time-out.

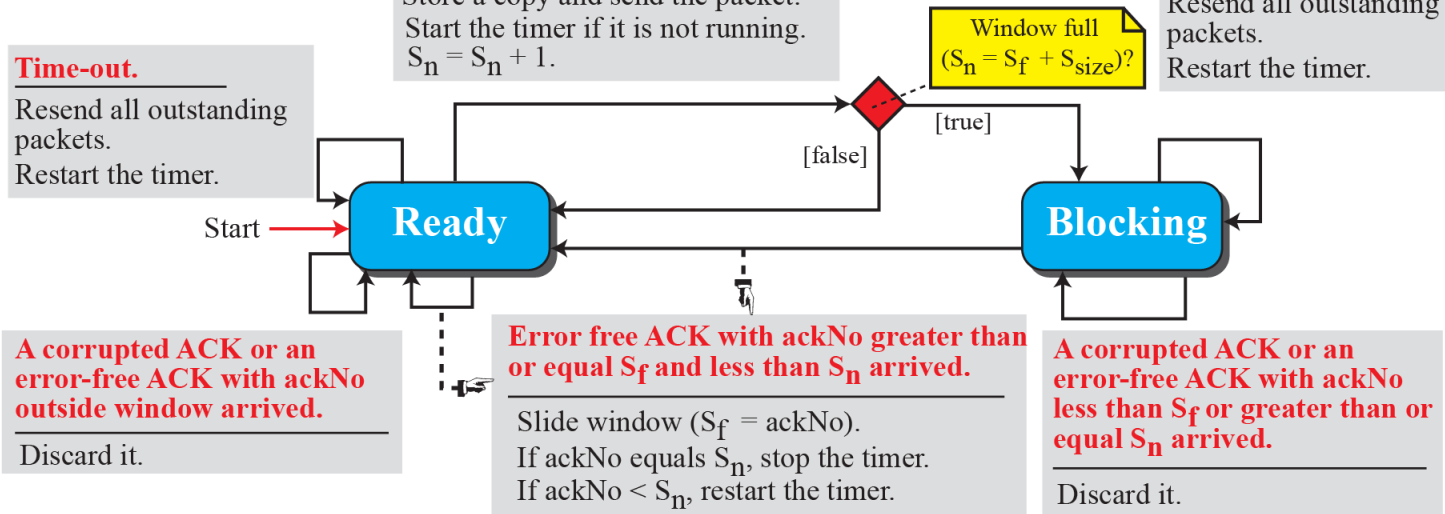
Resend all outstanding packets.
Restart the timer.

Request from process came.

Make a packet ($\text{seqNo} = S_n$).
Store a copy and send the packet.
Start the timer if it is not running.
 $S_n = S_n + 1$.

Time-out.

Resend all outstanding packets.
Restart the timer.



Receiver

Note:

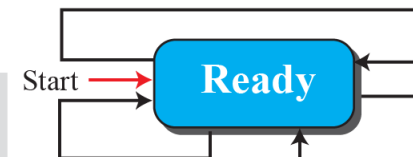
All arithmetic equations are in modulo 2^m .

Error-free packet with $\text{seqNo} = R_n$ arrived.

Deliver message.
Slide window ($R_n = R_n + 1$).
Send ACK ($\text{ackNo} = R_n$).

Corrupted packet arrived.

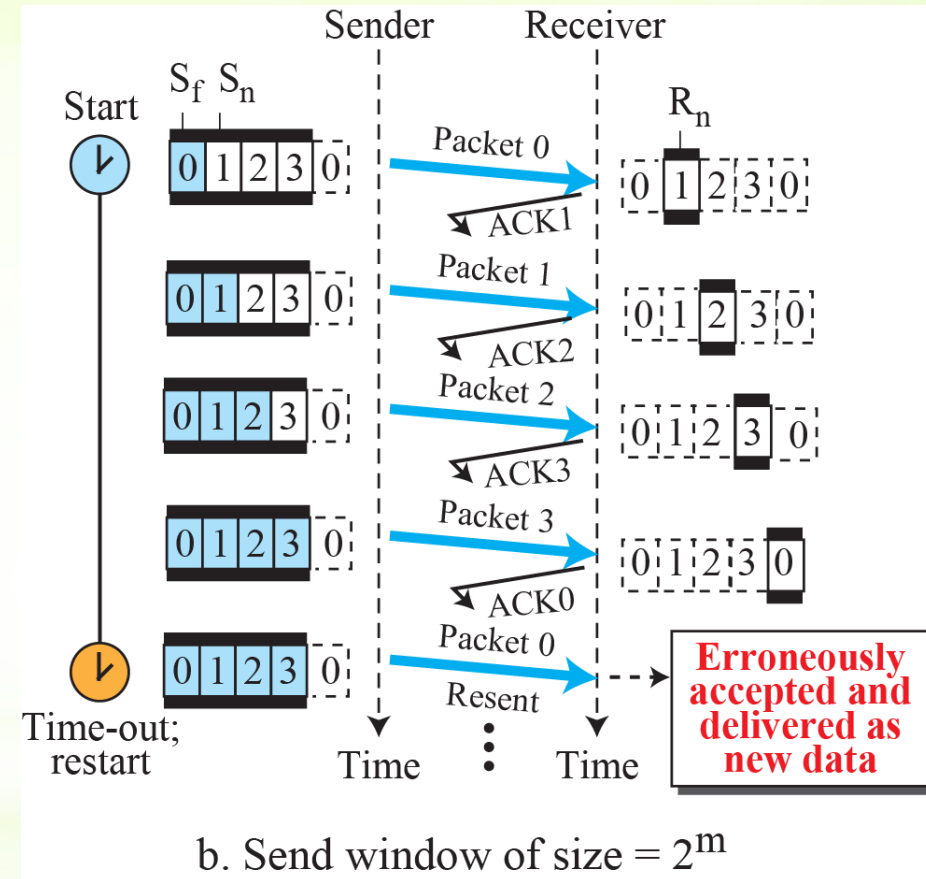
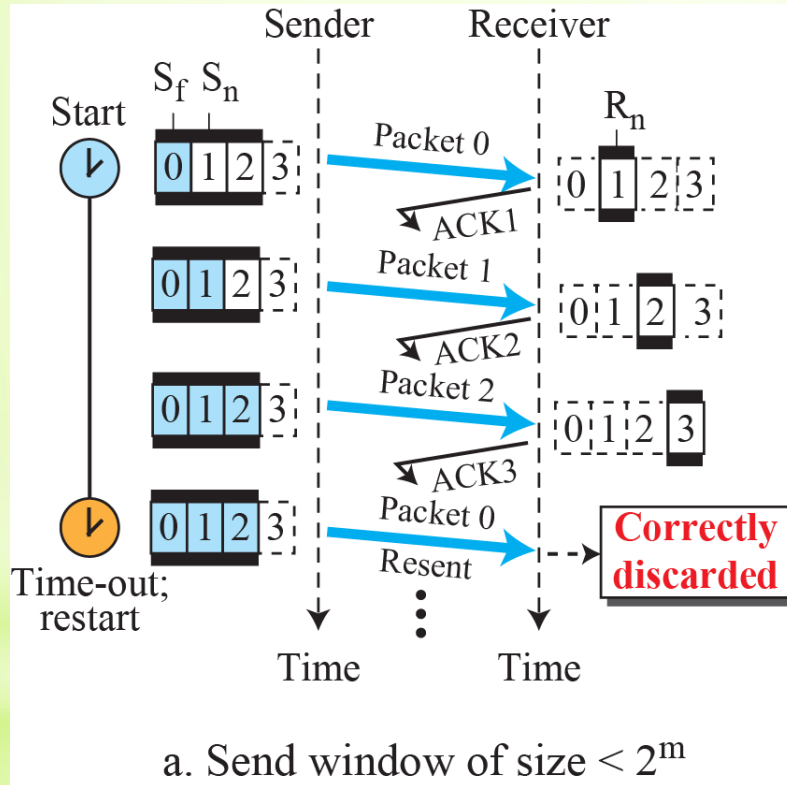
Discard packet.



Error-free packet with $\text{seqNo} \neq R_n$ arrived.

Discard packet.
Send an ACK ($\text{ackNo} = R_n$).

Send window size for Go-Back-N



- Cumulative acknowledgments can help if acknowledgments are delayed or lost

Shortfalls of Go-Back-N Protocol:

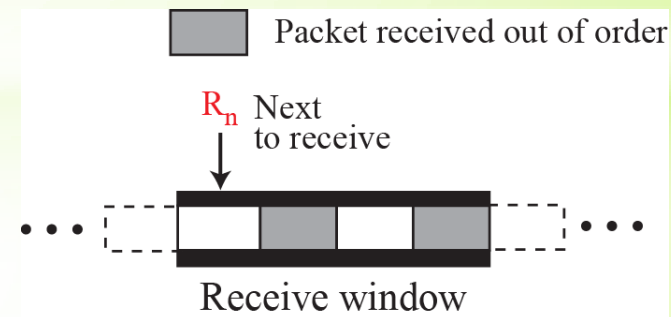
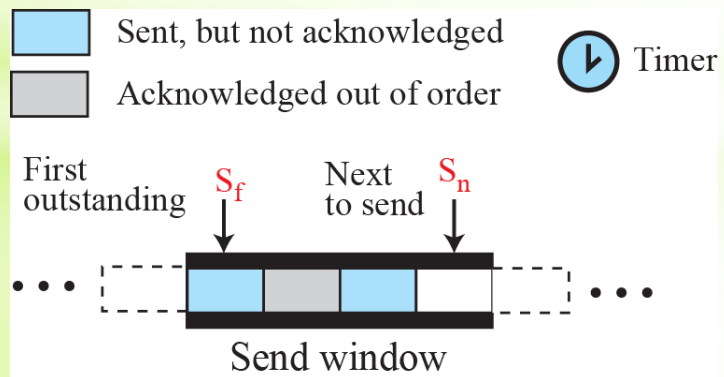
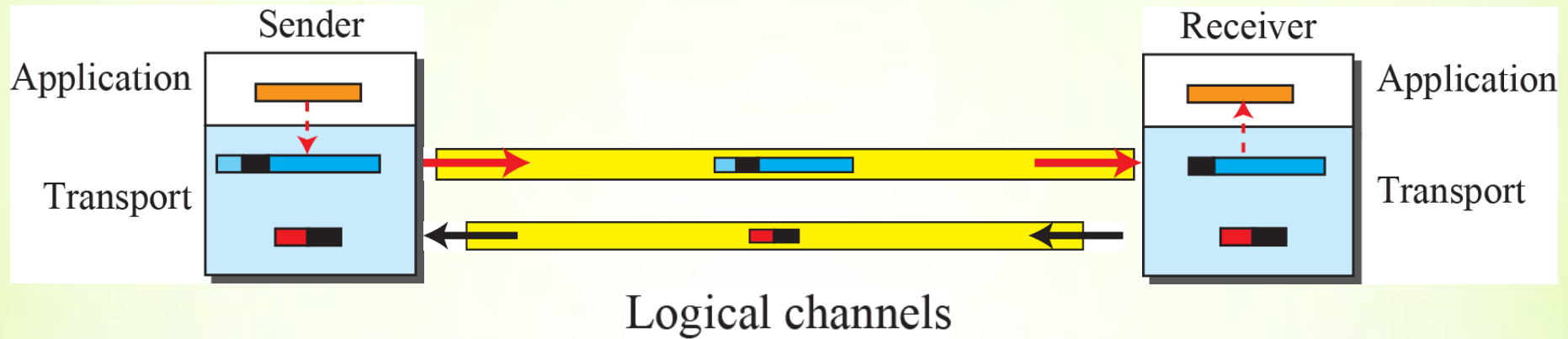
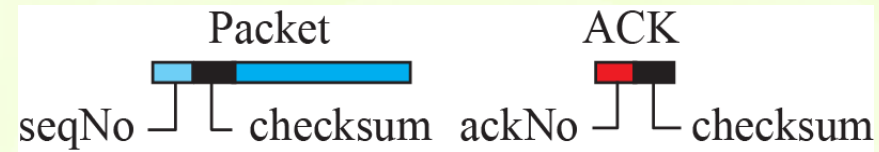
- There is no need to buffer out-of-order packets; they are simply discarded. However, this protocol is inefficient if the underlying network protocol loses a lot of packets.
- Each time a single packet is lost or corrupted, the sender resends all outstanding packets, even though some of these packets may have been received safe and sound but out of order.
- If the network layer is losing many packets because of congestion in the network, the resending of all of these outstanding packets makes the congestion worse, and eventually more packets are lost. This has an avalanche effect that may result in the total collapse of the network.

Selective Repeat protocol

The Go-Back-N protocol simplifies the process at the receiver. The receiver keeps track of only one variable, and there is no need to buffer out-of-order packets; they are simply discarded.

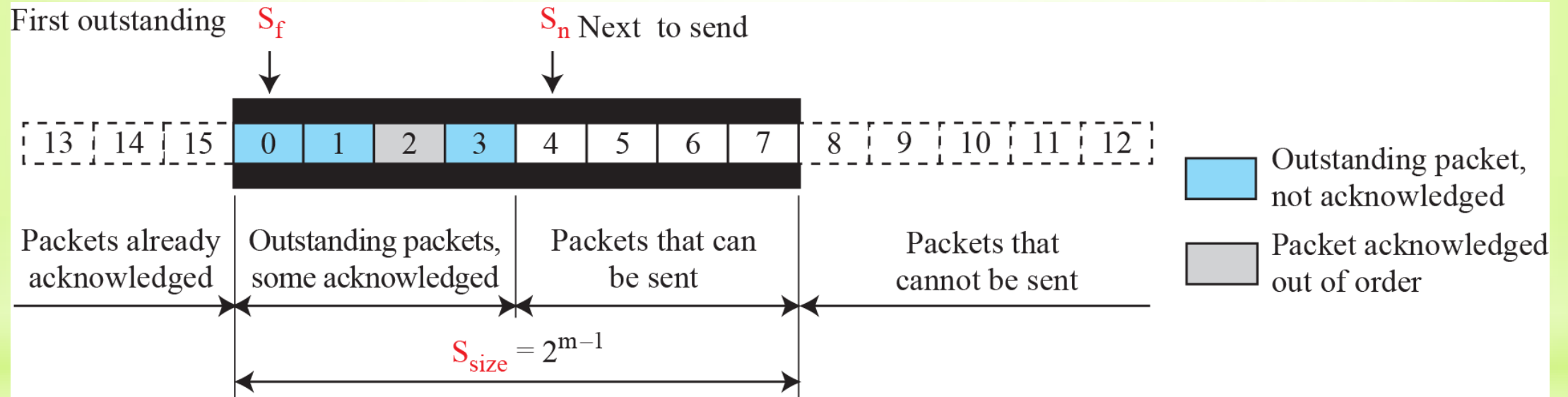
Another protocol, called the Selective-Repeat (SR) protocol, has been devised, which, as the name implies, resends only selective packets, those that are actually lost.

Outline of Selective-Repeat



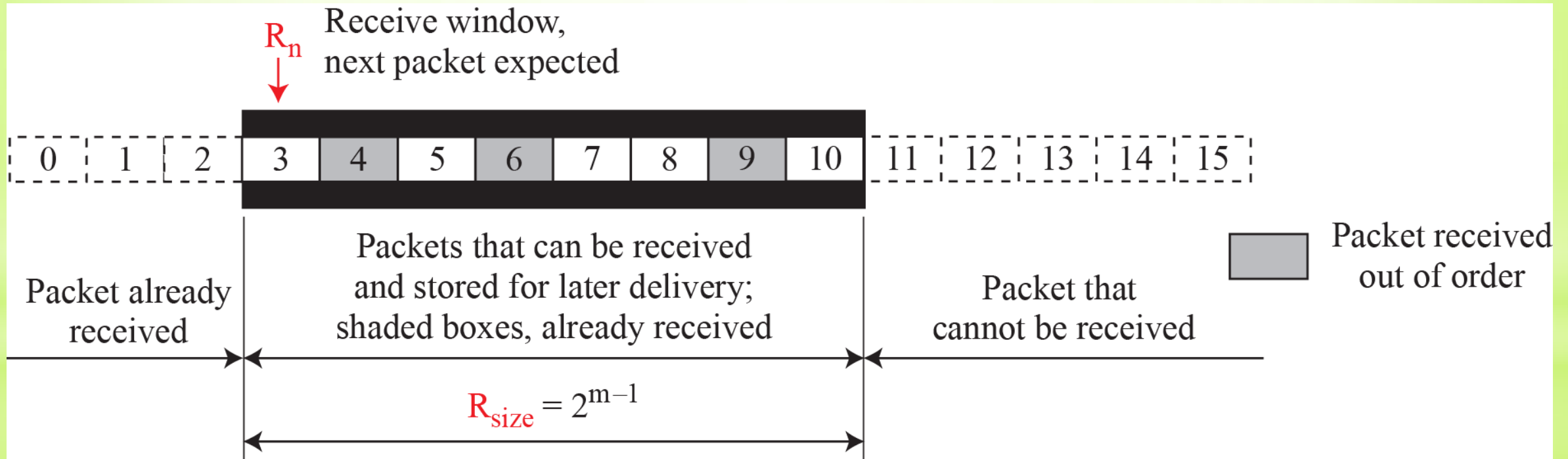
Send window for Selective-Repeat protocol

- The maximum **size of the send & receive window is 2^{m-1}** .



Receive window for Selective-Repeat protocol

- The Selective-Repeat protocol allows as many packets as the size of the receive window to arrive out of order and be kept until there is a set of consecutive packets to be delivered to the application layer.



- Selective-Repeat uses one timer for each outstanding packet.
- GBN treats outstanding packets as a group; SR treats them individually.
- In the Selective-Repeat protocol, an acknowledgment number defines the sequence number of the error-free packet received.

Flow diagram for selective repeat

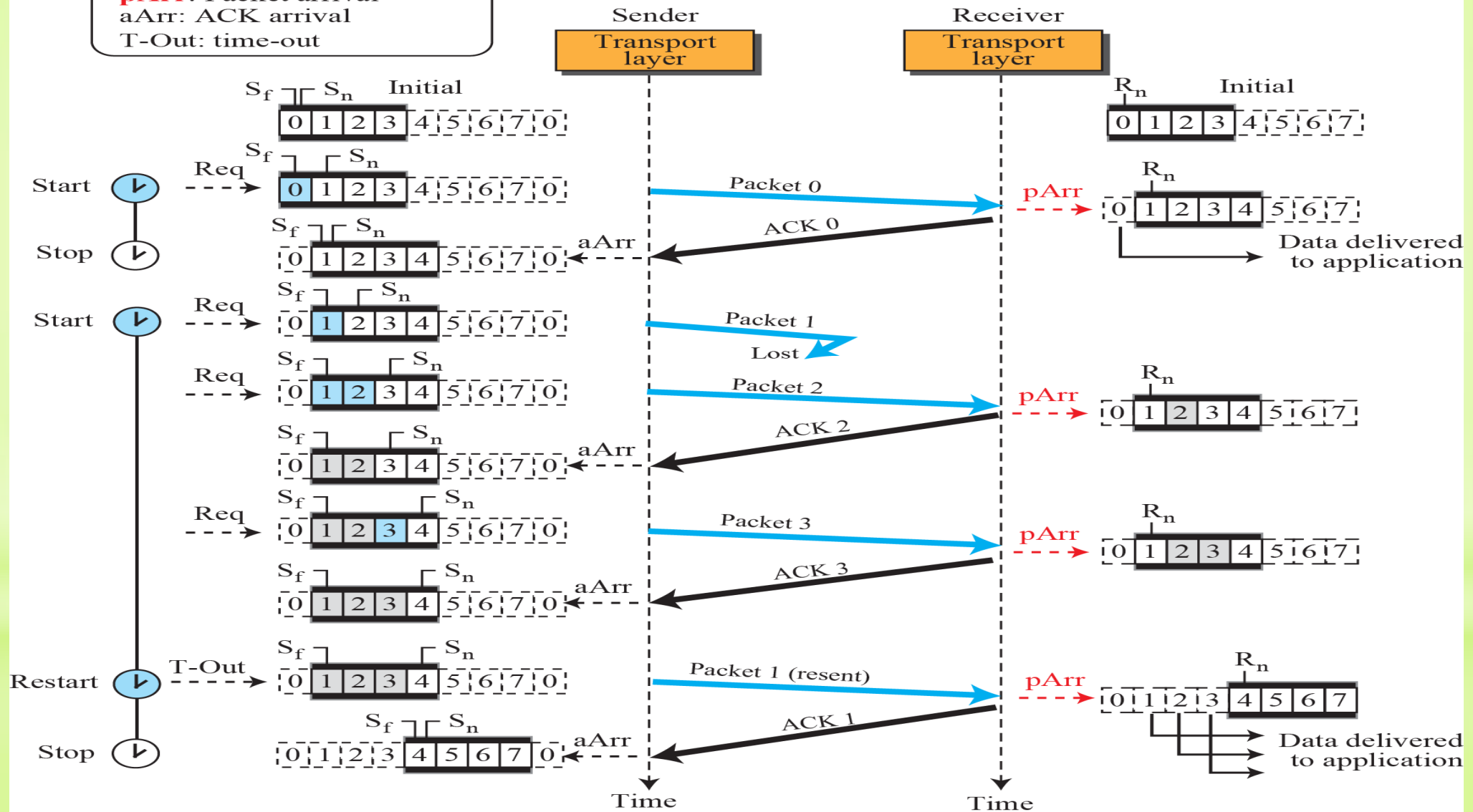
Events:

Req: Request from process

pArr: Packet arrival

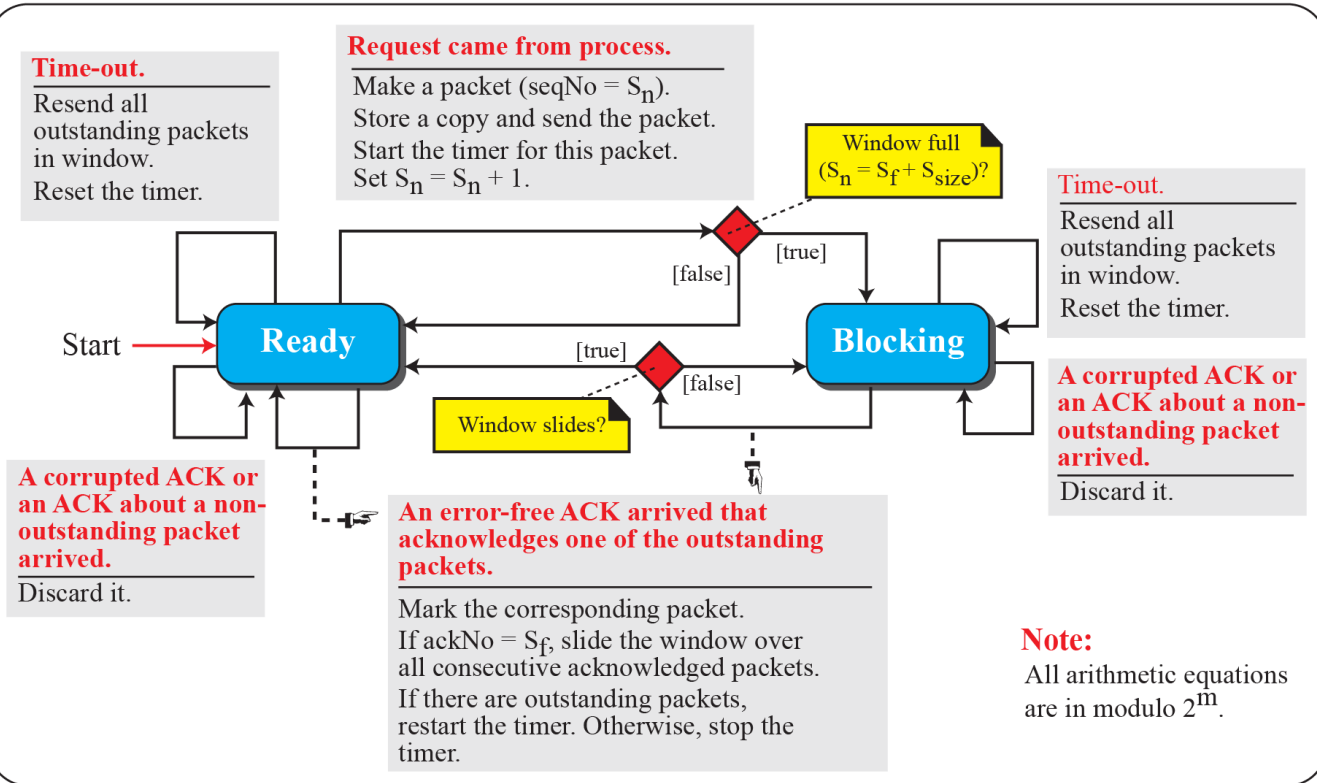
aArr: ACK arrival

T-Out: time-out

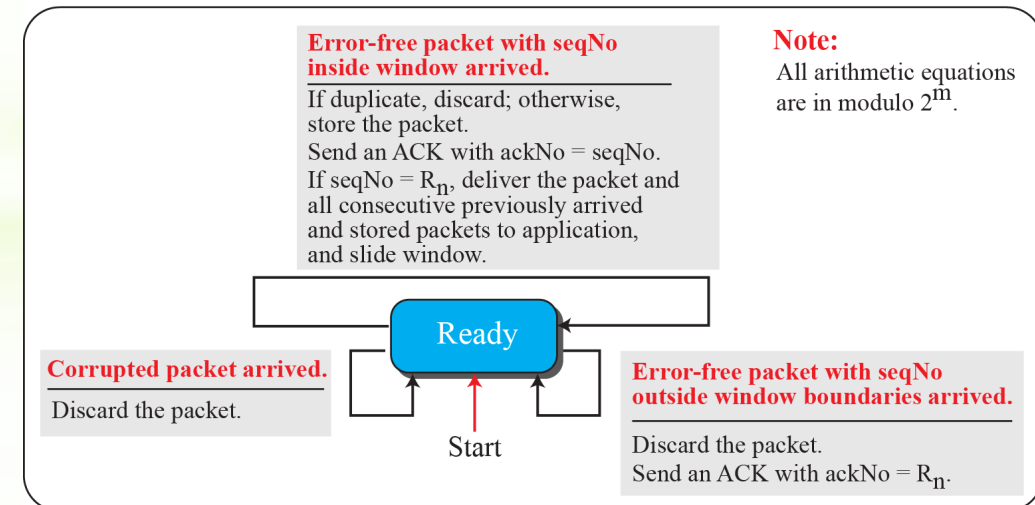


FSMs for SR protocol

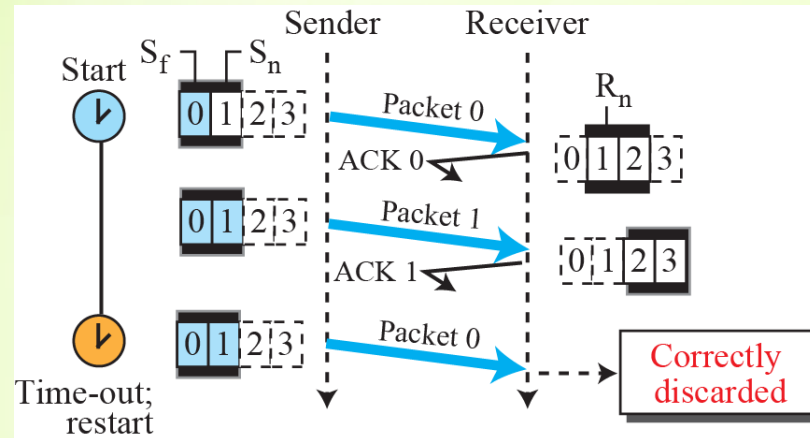
Sender



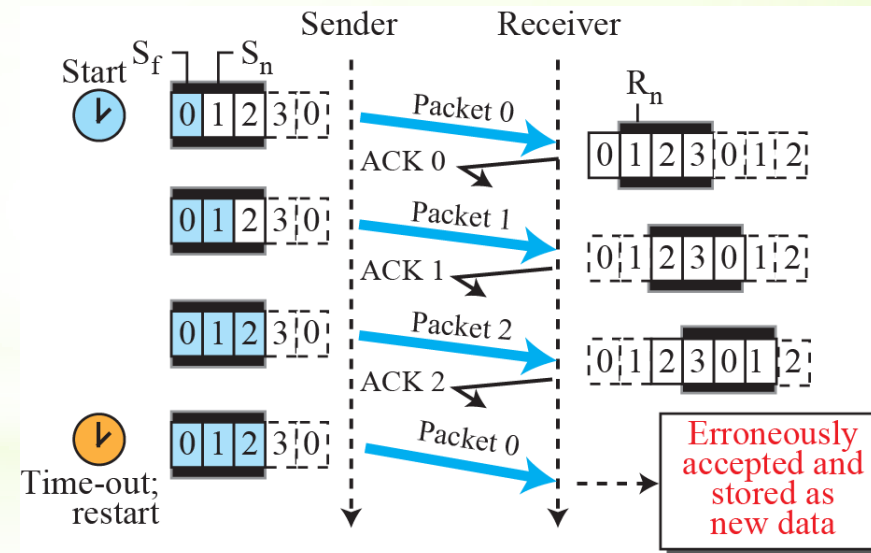
Receiver



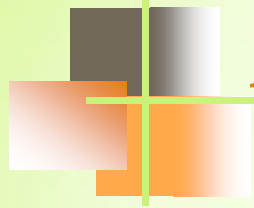
Selective-Repeat, window size



a. Send and receive windows of size $= 2^m - 1$



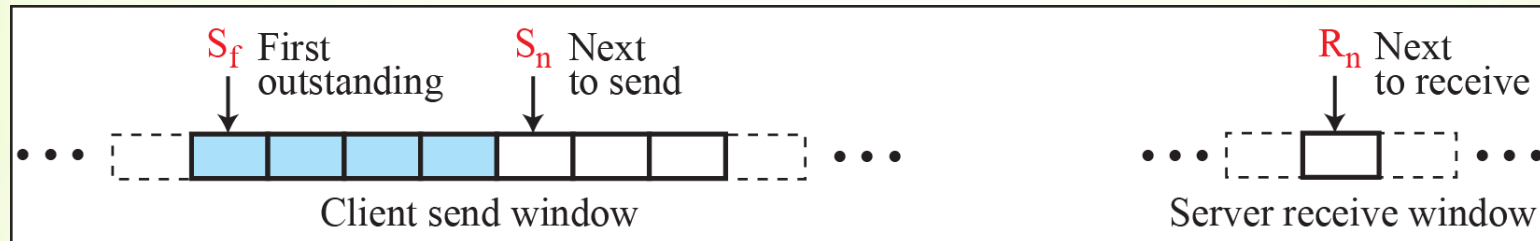
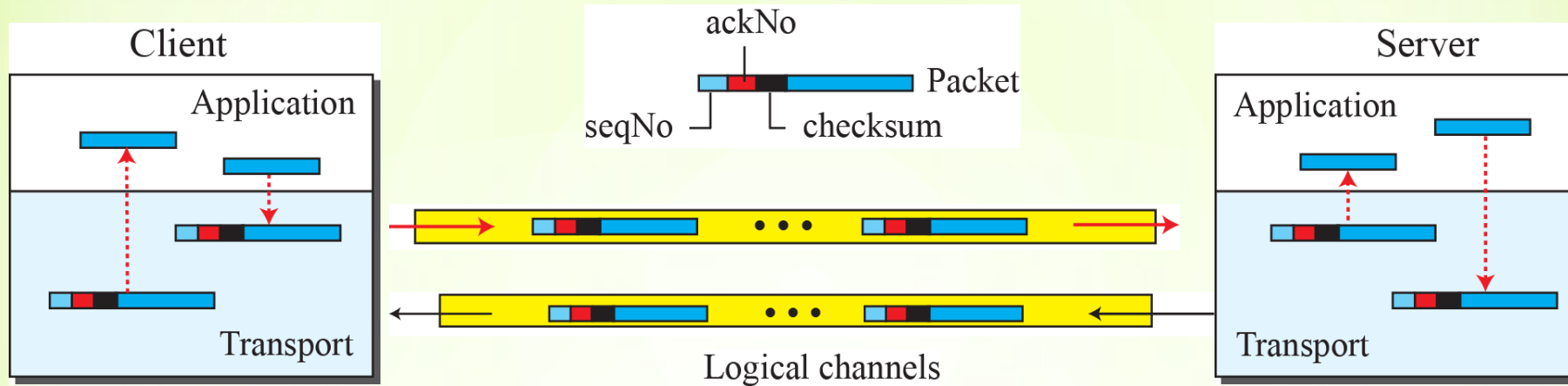
b. Send and receive windows of size $> 2^m - 1$



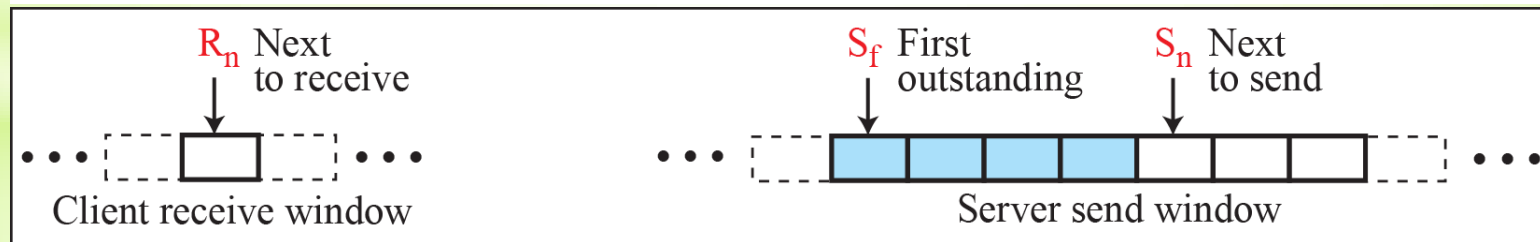
Bidirectional Protocols

The four protocols we discussed earlier in this section are all unidirectional: data packets flow in only one direction and acknowledgments travel in the other direction. In real life, data packets are normally flowing in both directions: from client to server and from server to client. This means that acknowledgments also need to flow in both directions. A technique called piggybacking is used to improve the efficiency of the bidirectional protocols.

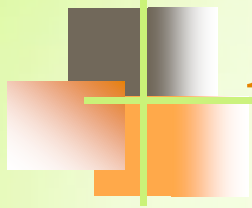
Design of piggybacking in Go-Back-N



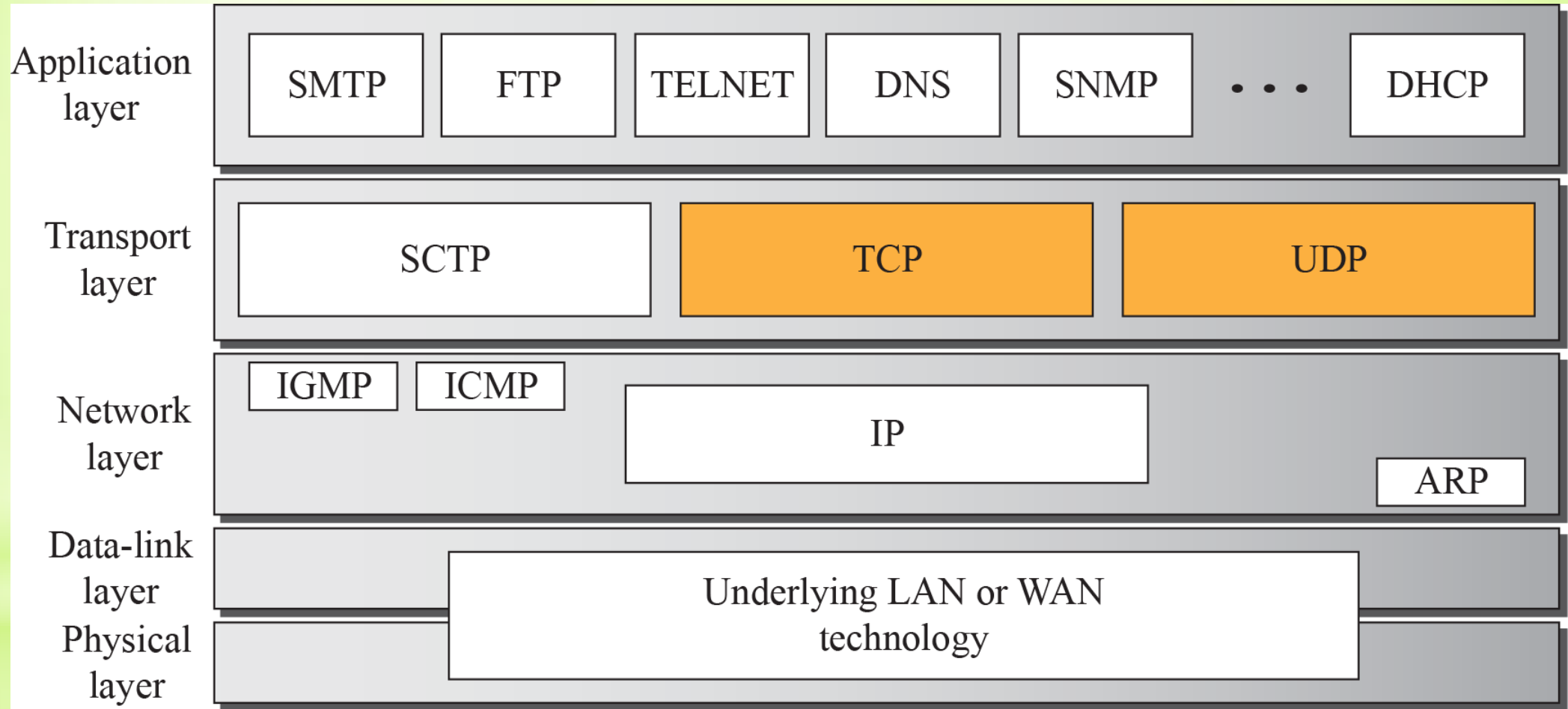
Windows for communication from client to server



Windows for communication from server to client



Internet Transport-Layer Protocols



UNIT - 4

Transport Layer

3.1 INTRODUCTION

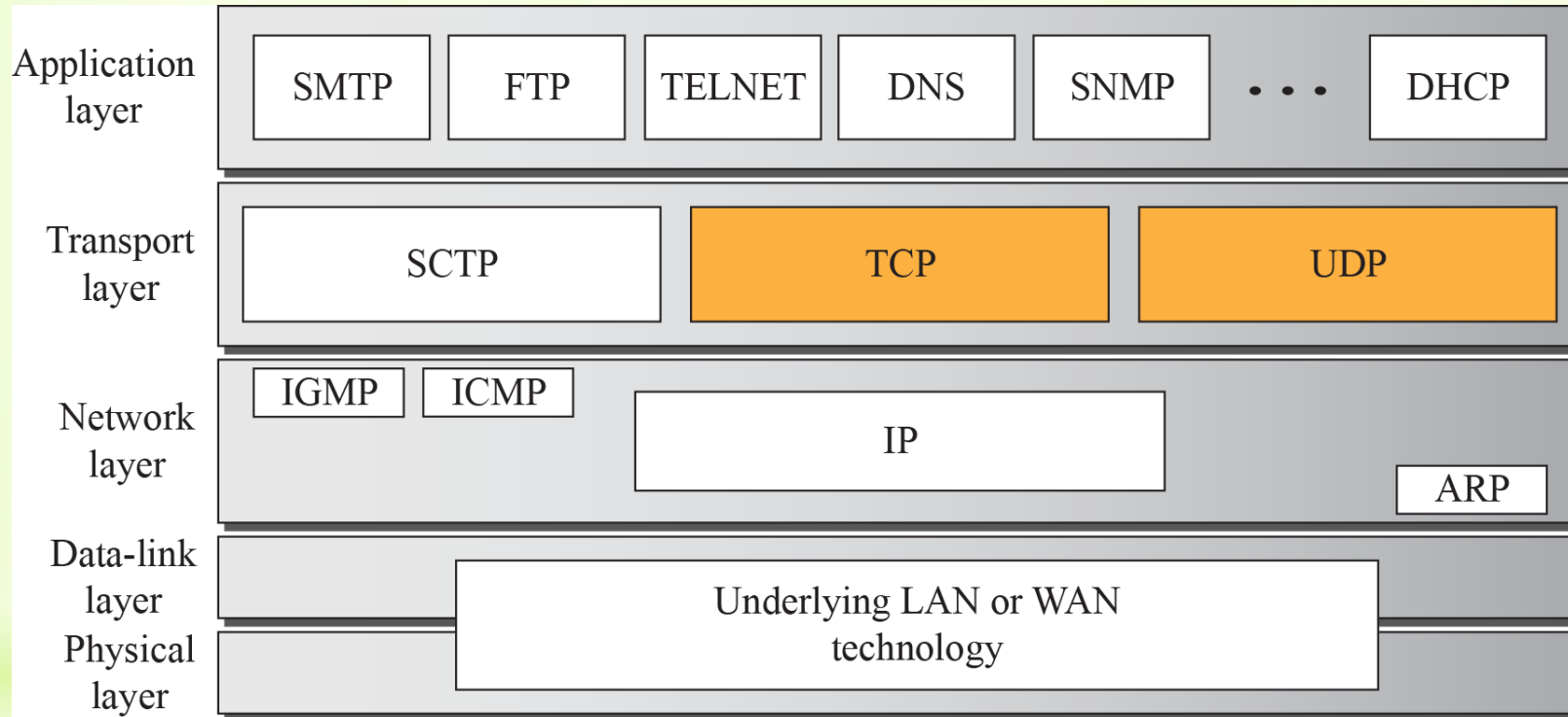
3.2 TRANSPORT-LAYER PROTOCOLS

- Simple Protocol
- Stop-and-Wait Protocol
- Go-Back-N Protocol (GBN)
- Selective-Repeat Protocol
- Bidirectional Protocols: Piggybacking

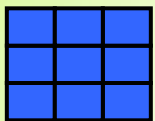
3.3 INTERNET TRANSPORT-LAYER PROTOCOLS

- USER DATAGRAM PROTOCOL (UDP)
- TRANSMISSION CONTROL PROTOCOL (TCP)

Internet Transport-Layer Protocols



***Position of transport-layer protocols in the TCP/IP
protocol suite***

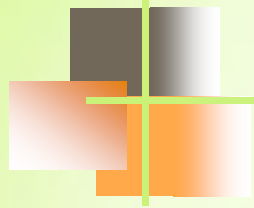


Some well-known ports used with UDP and TCP

<i>Port</i>	<i>Protocol</i>	<i>UDP</i>	<i>TCP</i>	<i>Description</i>
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, 21	FTP		√	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, 162	SNMP		√	Simple Network Management Protocol

USER DATAGRAM PROTOCOL (UDP)

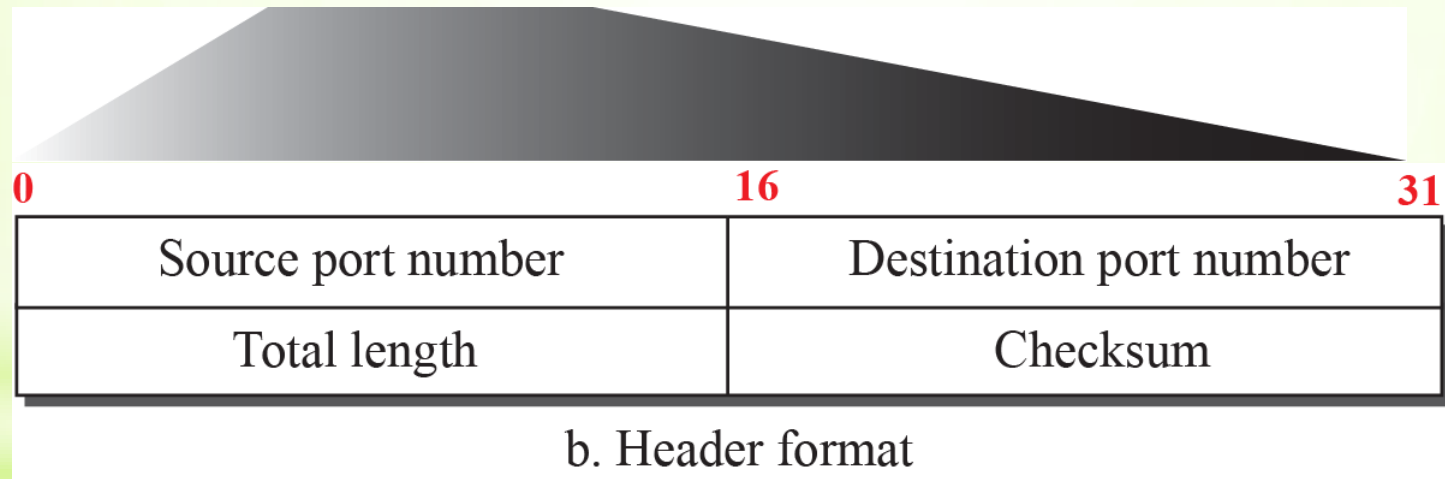
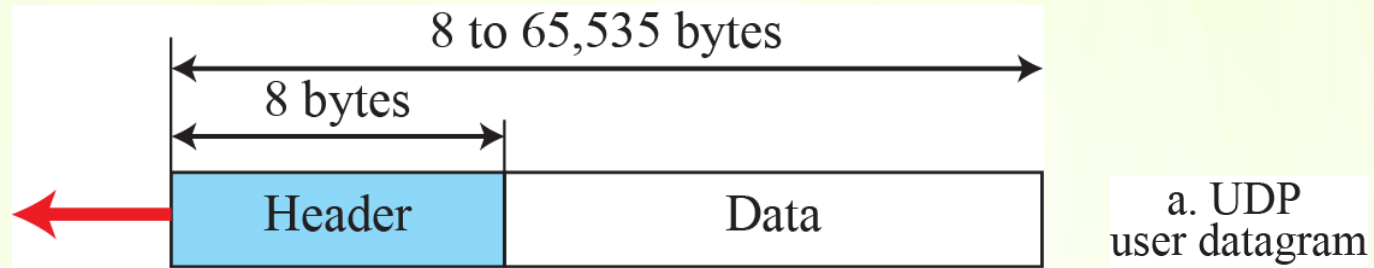
- ***The User Datagram Protocol (UDP) is a connectionless, unreliable transport protocol.***
- ***It does not add anything to the services of IP except for providing process-to-process instead of host-to-host communication.***
- ***UDP is a very simple protocol using a minimum of overhead.***



User Datagram

UDP packets, called user datagrams, have a fixed size header of 8 bytes made of four fields, each of 2 bytes (16 bits). the format of a user datagram. The first two fields define the source and destination port numbers. The third field defines the total length of the user datagram, header plus data. The 16 bits can define a total length of 0 to 65,535 bytes.

User datagram packet format



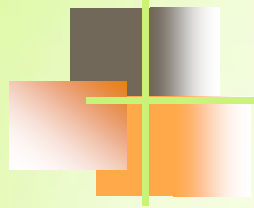
The following is the contents of a UDP header in hexadecimal format.

CB84000D001C001C

- a.** What is the source port number?
- b.** What is the destination port number?
- c.** What is the total length of the user datagram?
- d.** What is the length of the data?
- e.** Is the packet directed from a client to a server or vice versa?
- f.** What is the client process?

Solution

- a.** The source port number is the first four hexadecimal digits $(CB84)_{16}$ or 52100
- b.** The destination port number is the second four hexadecimal digits $(000D)_{16}$ or 13.
- c.** The third four hexadecimal digits $(001C)_{16}$ define the length of the whole UDP packet as 28 bytes.
- d.** The length of the data is the length of the whole packet minus the length of the header, or $28 - 8 = 20$ bytes.
- e.** Since the destination port number is 13 (well-known port), the packet is from the client to the server.
- f.** The client process is the Daytime (see Table 3.1).

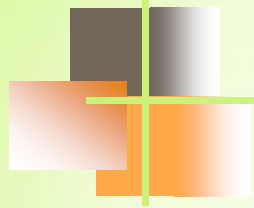


UDP Services

Earlier we discussed the general services provided by a transport-layer protocol. In this section, we discuss what portions of those general services are provided by UDP.

☐ Process-to-Process Communication

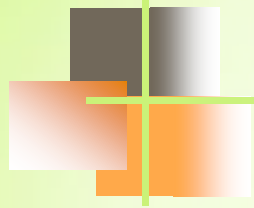
UDP provides process-to-process communication using socket addresses, a combination of IP addresses and port numbers.



UDP Services

□ Connectionless Services

- There is no connection establishment and no connection termination
- This means that each user datagram sent by UDP is an independent datagram. There is no relationship between the different user datagrams even if they are coming from the same source process and going to the same destination program.
- The user datagrams are not numbered.
- This means that each user datagram can travel on a different path.
- Only those processes sending short messages, messages less than 65,507 bytes (65,535 minus 8 bytes for the UDP header and minus 20 bytes for the IP header), can use UDP.



UDP Services

☐ Flow Control

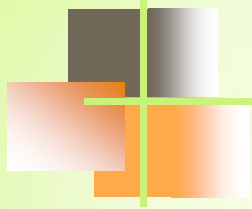
UDP is a very simple protocol. There is no flow control, and hence no window mechanism.

☐ Error Control

- There is no error control mechanism in UDP except for the checksum. This means that the sender does not know if a message has been lost or duplicated.
- When the receiver detects an error through the checksum, the user datagram is silently discarded.

☐ Congestion Control

- Since UDP is a connectionless protocol, it does not provide congestion control.
- UDP assumes that the packets sent are small and sporadic and cannot create congestion in the network.



❏ Checksum

UDP checksum calculation includes three sections: a pseudoheader, the UDP header, and the data coming from the application layer.

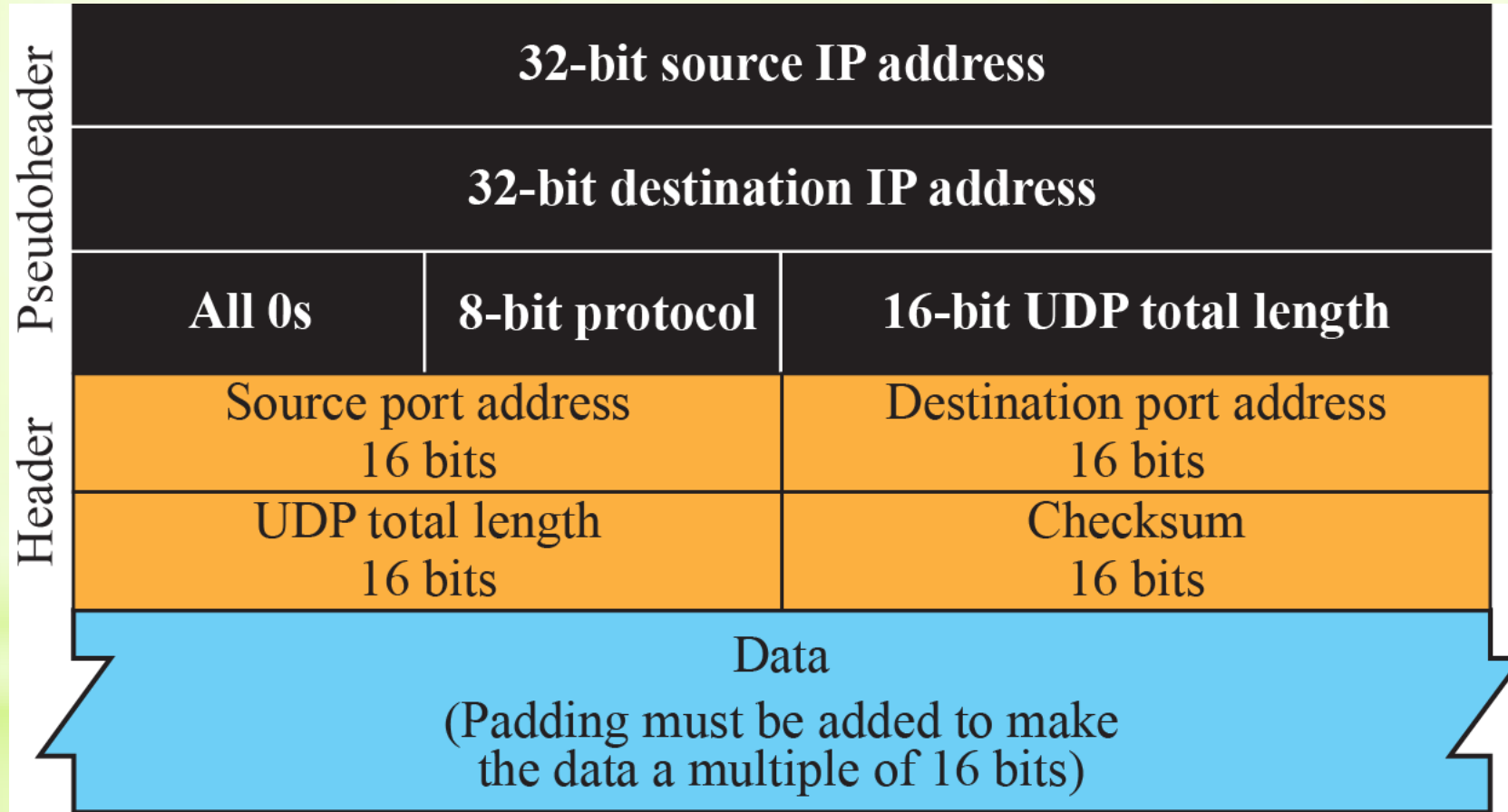
The pseudoheader is the part of the header of the IP packet in which the user datagram is to be encapsulated with some fields filled with 0s.

If the checksum does not include the pseudoheader, a user datagram may arrive safe and sound. However, if the IP header is corrupted, it may be delivered to the wrong host.

The protocol field is added to ensure that the packet belongs to UDP, and not to TCP. The value of the protocol field for UDP is 17.

If this value is changed during transmission, the checksum calculation at the receiver will detect it and UDP drops the packet. It is not delivered to the wrong protocol.

Pseudoheader for checksum calculation





❑ Encapsulation and Decapsulation

To send a message from one process to another, the UDP protocol encapsulates and decapsulates messages

❑ Queuing

At the client site, when a process starts, it requests a port number from the operating system. Some implementations create both an incoming and an outgoing queue associated with each process.

❑ Multiplexing and Demultiplexing

In a host running a TCP/IP protocol suite, there is only one UDP but possibly several processes that may want to use the services of UDP. To handle this situation, UDP multiplexes and demultiplexes

❑ Comparison : UDP and Simple Protocol

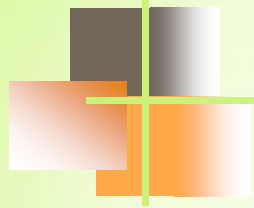
The only difference is that UDP provides an optional checksum to detect corrupted packets at the receiver site.

What value is sent for the checksum in one of the following hypothetical situations?

- a.** The sender decides not to include the checksum.
- b.** The sender decides to include the checksum, but the value of the sum is all 1s.
- c.** The sender decides to include the checksum, but the value of the sum is all 0s.

Solution

- a.** The value sent for the checksum field is all 0s to show that the checksum is not calculated.
- b.** When the sender complements the sum, the result is all 0s; the sender complements the result again before sending. The value sent for the checksum is all 1s. The second complement operation is needed to avoid confusion with the case in part a.
- c.** This situation never happens because it implies that the value of every term included in the calculation of the sum is all 0s, which is impossible; some fields in the pseudoheader have nonzero values.



UDP Applications

Although UDP meets almost none of the criteria we mentioned earlier for a reliable transport-layer protocol, UDP is preferable for some applications. The reason is that some services may have some side effects that are either unacceptable or not preferable. An application designer sometimes needs to compromise to get the optimum.

❑ UDP Features

- ❖ Connectionless Service
- ❖ Lack of Error Control
- ❖ Lack of Congestion Control

A client-server application such as DNS uses the services of UDP because a client needs to send a short request to a server and to receive a quick response from it. The request and response can each fit in one user datagram. Since only one message is exchanged in each direction, the connectionless feature is not an issue; the client or server does not worry that messages are delivered out of order.

A client-server application such as SMTP, which is used in electronic mail, cannot use the services of UDP because a user can send a long e-mail message, which may include multimedia (images, audio, or video). If the application uses UDP and the message does not fit in one single user datagram, the message must be split by the application into different user datagrams. Here the connectionless service may create problems. The user datagrams may arrive and be delivered to the receiver application out of order. The receiver application may not be able to reorder the pieces. This means the connectionless service has a disadvantage for an application program that sends long messages.

Assume we are downloading a very large text file from the Internet. We definitely need to use a transport layer that provides reliable service. We don't want part of the file to be missing or corrupted when we open the file. The delay created between the deliveries of the parts is not an overriding concern for us; we wait until the whole file is composed before looking at it. In this case, UDP is not a suitable transport layer.

Assume we are using a real-time interactive application, such as Skype. Audio and video are divided into frames and sent one after another. If the transport layer is supposed to resend a corrupted or lost frame, the synchronizing of the whole transmission may be lost. The viewer suddenly sees a blank screen and needs to wait until the second transmission arrives. This is not tolerable. However, if each small part of the screen is sent using one single user datagram, the receiving UDP can easily ignore the corrupted or lost packet and deliver the rest to the application program. That part of the screen is blank for a very short period of time, which most viewers do not even notice.

Typical Applications of UDP

- UDP is suitable for a process that requires simple request-response communication with little concern for flow and error control. It is not usually used for a process such as FTP that needs to send bulk data
- UDP is suitable for a process with internal flow- and error-control mechanisms. For example, the Trivial File Transfer Protocol (TFTP) process includes flow and error control. It can easily use UDP.
- UDP is a suitable transport protocol for multicasting. Multicasting capability is embedded in the UDP software but not in the TCP software.
- UDP is used for management processes such as SNMP
- UDP is used for some route updating protocols such as Routing Information Protocol (RIP)
- UDP is normally used for interactive real-time applications that cannot tolerate uneven delay between sections of a received message

UNIT - 4

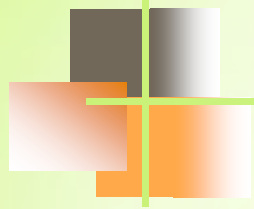
Transport Layer

- 
- INTRODUCTION
 - TRANSPORT-LAYER PROTOCOLS
 - INTERNET TRANSPORT-LAYER PROTOCOLS
 - ✓ USER DATAGRAM PROTOCOL (UDP)
 - ✓ **TRANSMISSION CONTROL PROTOCOL (TCP)**

TRANSMISSION CONTROL PROTOCOL (TCP)

Transmission Control Protocol (TCP) is a connection-oriented, reliable protocol. TCP explicitly defines connection establishment, data transfer, and connection teardown phases to provide a connection-oriented service.

TCP uses a combination of GBN and SR protocols to provide reliability.



TCP Services

- ☐ Process-to-Process Communication
- ☐ Stream Delivery Service
 - ❖ Sending and Receiving Buffers
 - ❖ Segments
- ☐ Full-Duplex Communication
- ☐ Multiplexing and Demultiplexing
- ☐ Connection-Oriented Service
- ☐ Reliable Service

TCP Services

❑ Process-to-Process Communication

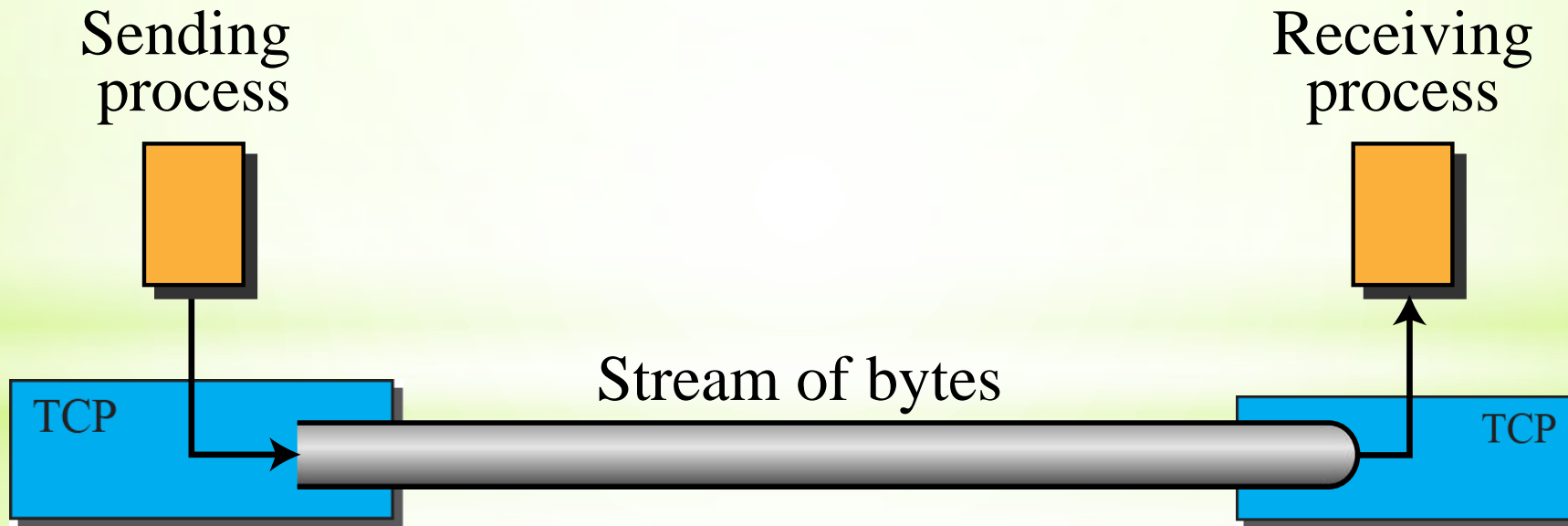
- TCP provides process-to-process communication using port numbers.
- Senders port number & Destination Port number are placed in header of TCP Segment

❑ Stream Delivery Service

- TCP groups a number of bytes together into a packet called a segment
- TCP adds a header to each segment (for control purposes) and delivers the segment to the network layer for transmission. The segments are encapsulated in an IP datagram and transmitted. This entire operation is transparent to the receiving process.
- The segments are not necessarily all the same size.

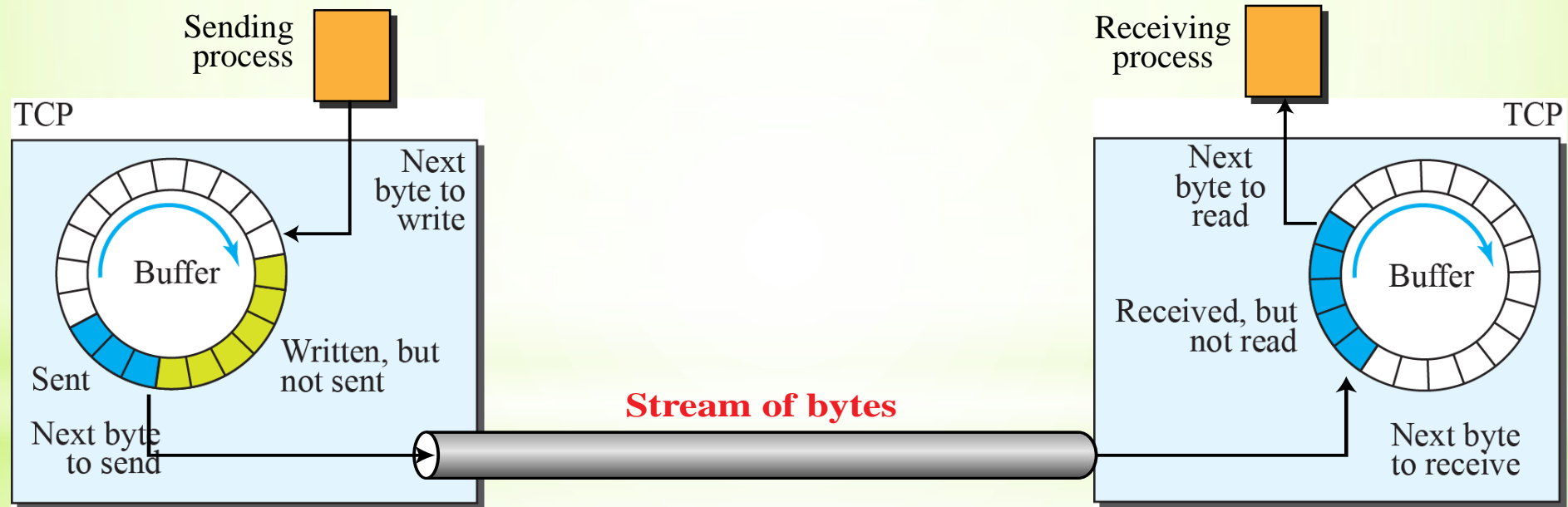
Stream delivery

- TCP, unlike UDP, is a stream-oriented protocol.
- TCP, allows the sending process to deliver data as a stream of bytes and allows the receiving process to obtain data as a stream of bytes.



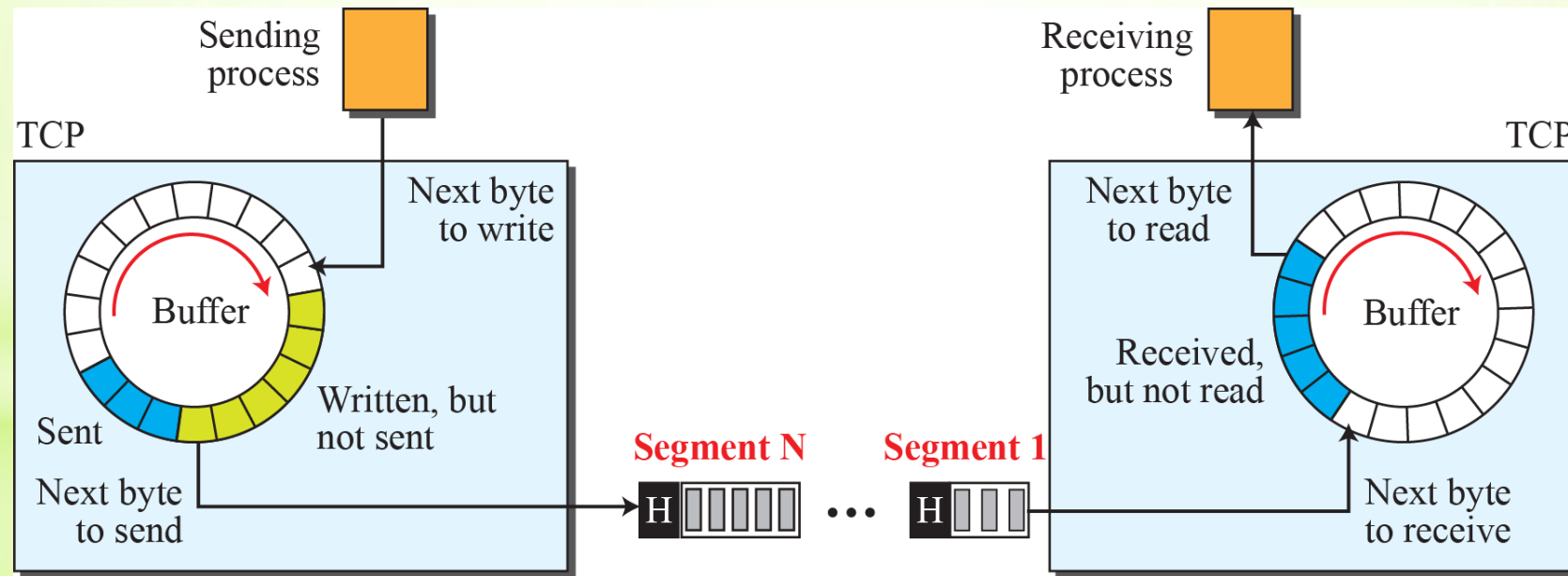
Sending and receiving buffers

- There are two buffers, the sending buffer and the receiving buffer, one for each direction
- The TCP sender keeps these bytes in the buffer until it receives an acknowledgment.



TCP segments

- TCP groups a number of bytes together into a packet called a segment
- TCP adds a header to each segment (for control purposes) and delivers the segment to the network layer for transmission. The segments are encapsulated in an IP datagram and transmitted. This entire operation is transparent to the receiving process.
- The segments are not necessarily all the same size.



❑ Full-Duplex Communication

TCP offers full-duplex service, where data can flow in both directions at the same time. Each TCP endpoint then has its own sending and receiving buffer, and segments move in both directions.

❑ Multiplexing and Demultiplexing

TCP performs multiplexing at the sender and demultiplexing at the receiver. However, since TCP is a connection-oriented protocol, a connection needs to be established for each pair of processes.

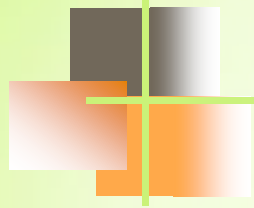
❑ Connection-Oriented Service

TCP is a connection-oriented protocol. With Three Phases:

1. The two TCP's establish a logical connection between them.
2. Data are exchanged in both directions.
3. The connection is terminated.

❑ Reliable Service

TCP is a reliable transport protocol. It uses an acknowledgment mechanism to check the safe and sound arrival of data.



TCP Features

To provide the services mentioned in the previous section, TCP has several features:

□ Numbering System

❖ Byte Number

- The bytes of data being transferred in each connection are numbered by TCP.
- The numbering starts with an arbitrarily generated number.
- Numbering is independent in each direction.
- TCP chooses an arbitrary number between 0 and $2^{32} - 1$ for the number of the first byte

❖ Sequence Number

- After the bytes have been numbered, TCP assigns a sequence number to each segment that is being sent. The sequence number, in each direction, is defined as follows:
 1. The sequence number of the first segment is the ISN (initial sequence number), which is a random number.
 2. The sequence number of any other segment is the sequence number of the previous segment plus the number of bytes (real or imaginary) carried by the previous segment.

❖ Acknowledgment Number

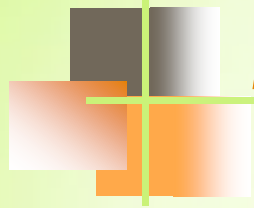
- TCP is full duplex; when a connection is established, both parties can send and receive data at the same time
- The value of the acknowledgment field in a segment defines the number of the next byte a party expects to receive. The acknowledgment number is cumulative.

Suppose a TCP connection is transferring a file of 5,000 bytes. The first byte is numbered 10,001. What are the sequence numbers for each segment if data are sent in five segments, each carrying 1,000 bytes?

Solution

The following shows the sequence number for each segment:

Segment 1	→	Sequence Number:	10,001	Range:	10,001	to	11,000
Segment 2	→	Sequence Number:	11,001	Range:	11,001	to	12,000
Segment 3	→	Sequence Number:	12,001	Range:	12,001	to	13,000
Segment 4	→	Sequence Number:	13,001	Range:	13,001	to	14,000
Segment 5	→	Sequence Number:	14,001	Range:	14,001	to	15,000



Segment

A packet in TCP is called a **segment**.

❑ Format

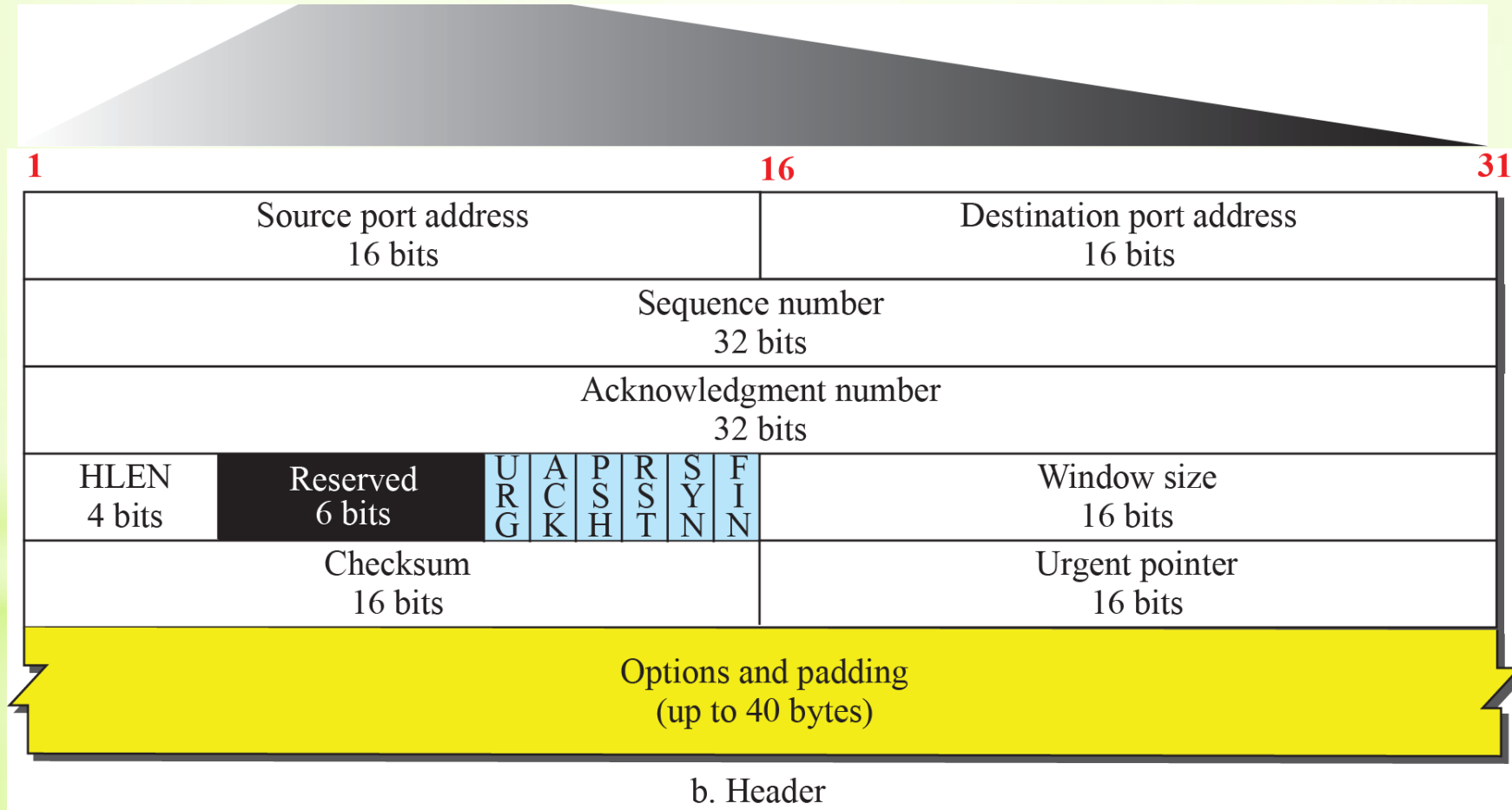
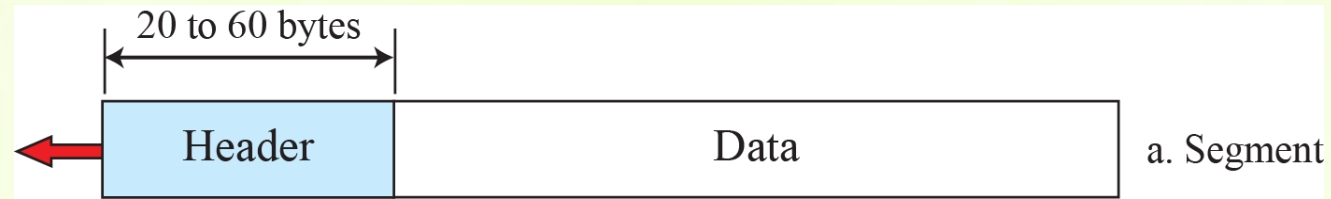
The segment consists of a header of 20 to 60 bytes, followed by data from the application program.

The header is 20 bytes if there are no options and up to 60 bytes if it contains options.

❑ Encapsulation

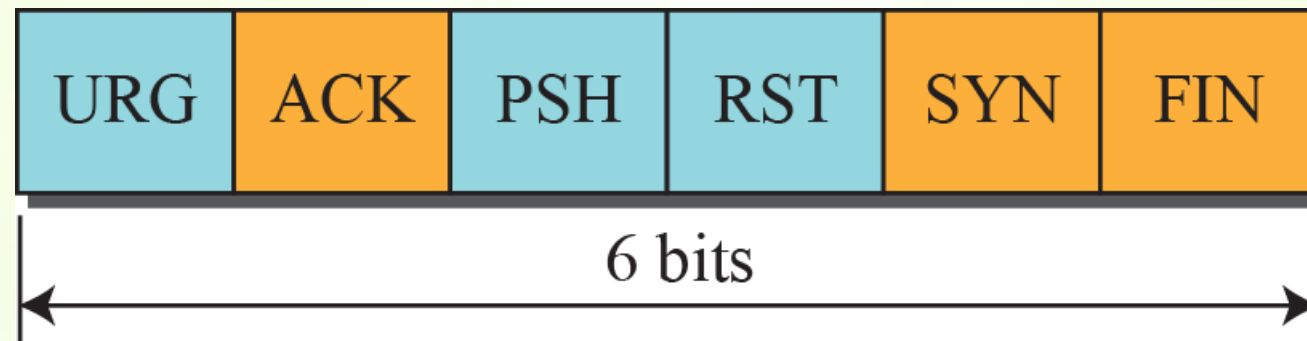
A TCP segment encapsulates the data received from the application layer. The TCP segment is encapsulated in an IP datagram, which in turn is encapsulated in a frame at the data-link layer

TCP segment format



- Source port address (16 bits): defines the port number of the application program in the host that is sending the segment
- Destination port address (16 bits): defines the port number of the application program in the host that is receiving the segment
- Sequence number(32 bits): defines the number assigned to the first byte of data contained in this segment
- Acknowledgment number(32 bits): defines the byte number that the receiver of the segment is expecting to receive from the other party
- Header length(4 bits): indicates the number of 4-byte words in the TCP header.

Control Field(6 bits): defines 6 different control bits or flags. One or more of these bits can be set at a time. These bits enable flow control, connection establishment and termination, connection abortion, and the mode of data transfer in TCP.



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

Window size(16 bits): defines the window size of the sending TCP in bytes. The maximum size of the window is 65,535 bytes. This value is normally referred to as the receiving window (rwnd) and is determined by the receiver.

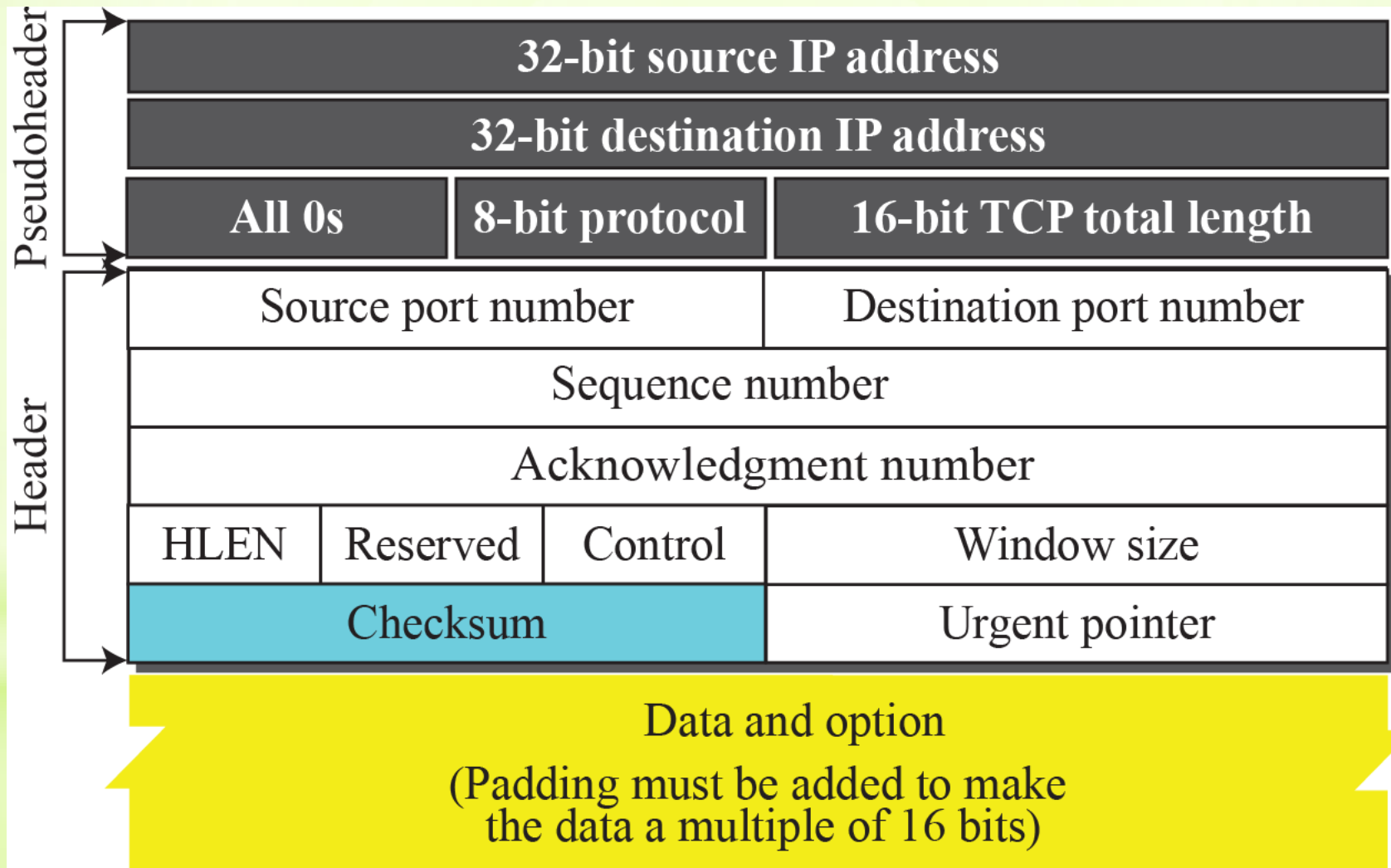
Checksum(16 bits): The calculation of the checksum for TCP follows the same procedure as the one described for UDP. However, the use of the checksum in the UDP datagram is optional, whereas the use of the checksum for TCP is mandatory.

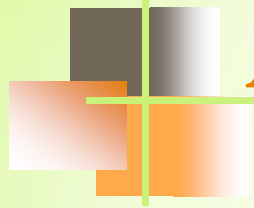
Urgent pointer(16 bits): which is valid only if the urgent flag is set, is used when the segment contains urgent data. It defines a value that must be added to the sequence number to obtain the number of the last urgent byte in the data section of the segment.

Options: There can be up to 40 bytes of optional information in the TCP header.

Pseudoheader added to the TCP datagram

TCP pseudoheader, the value for the protocol field is 6.





A TCP Connection

TCP is connection-oriented. As discussed before, a connection-oriented transport protocol establishes a logical path between the source and destination. All of the segments belonging to a message are then sent over this logical path. Using a single logical pathway for the entire message facilitates the acknowledgment process as well as retransmission of damaged or lost frames.



(continued)

- ☐ Connection Establishment

- ❖ Three-Way Handshaking
- ❖ SYN Flooding Attack

- ☐ Data Transfer

- ❖ Pushing Data
- ❖ Urgent Data

- ☐ Connection Termination

- ❖ Three-Way Handshaking
- ❖ Half-Close

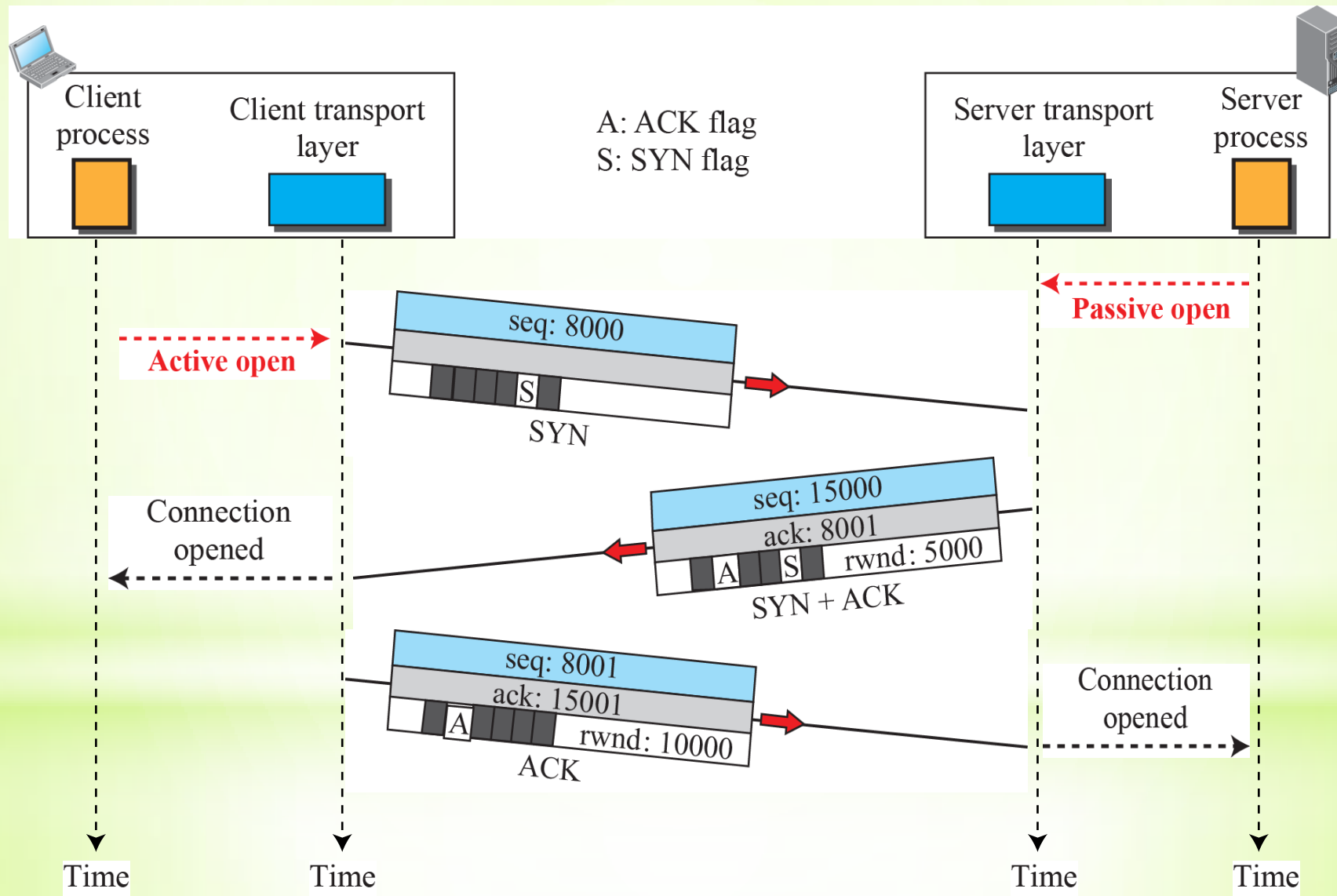
- ☐ Connection Reset

Connection Establishment

- TCP transmits data in full-duplex mode.
- The connection establishment in TCP is called *three-way handshaking*.
- The process starts with the server. The server program tells its TCP that it is ready to accept a connection. This request is called a *passive open*.
- The client program issues a request for an *active open*.
- The connection establishment is done by exchanging 3 messages between the two parties:
 - SYN segment
 - SYN + ACK segment
 - ACK segment

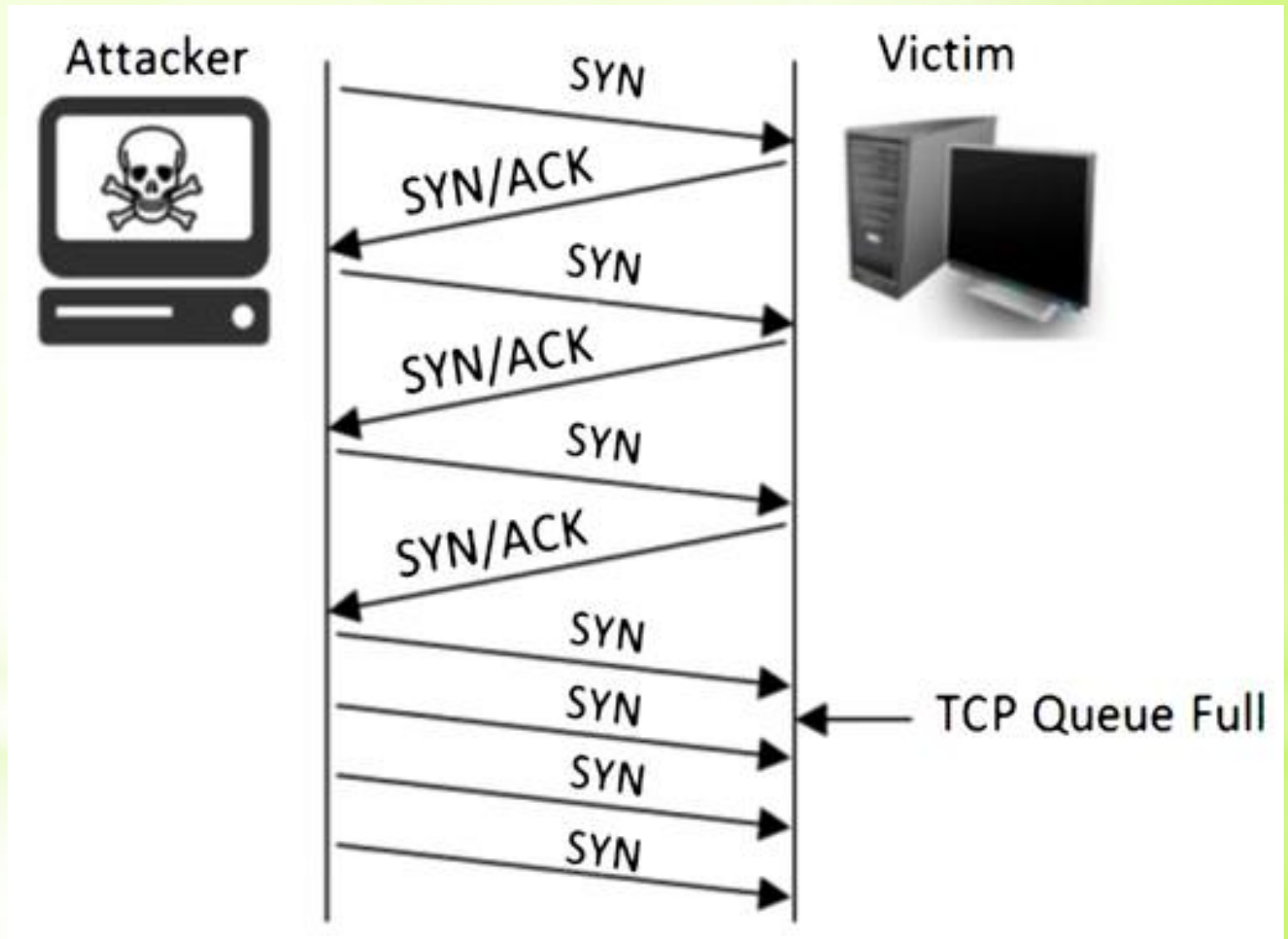
- A SYN segment cannot carry data, but it consumes one sequence number.
- A SYN + ACK segment cannot carry data, but it does consume one sequence number.
- An ACK segment, if carrying no data, consumes no sequence number.

Connection establishment using three-way handshaking

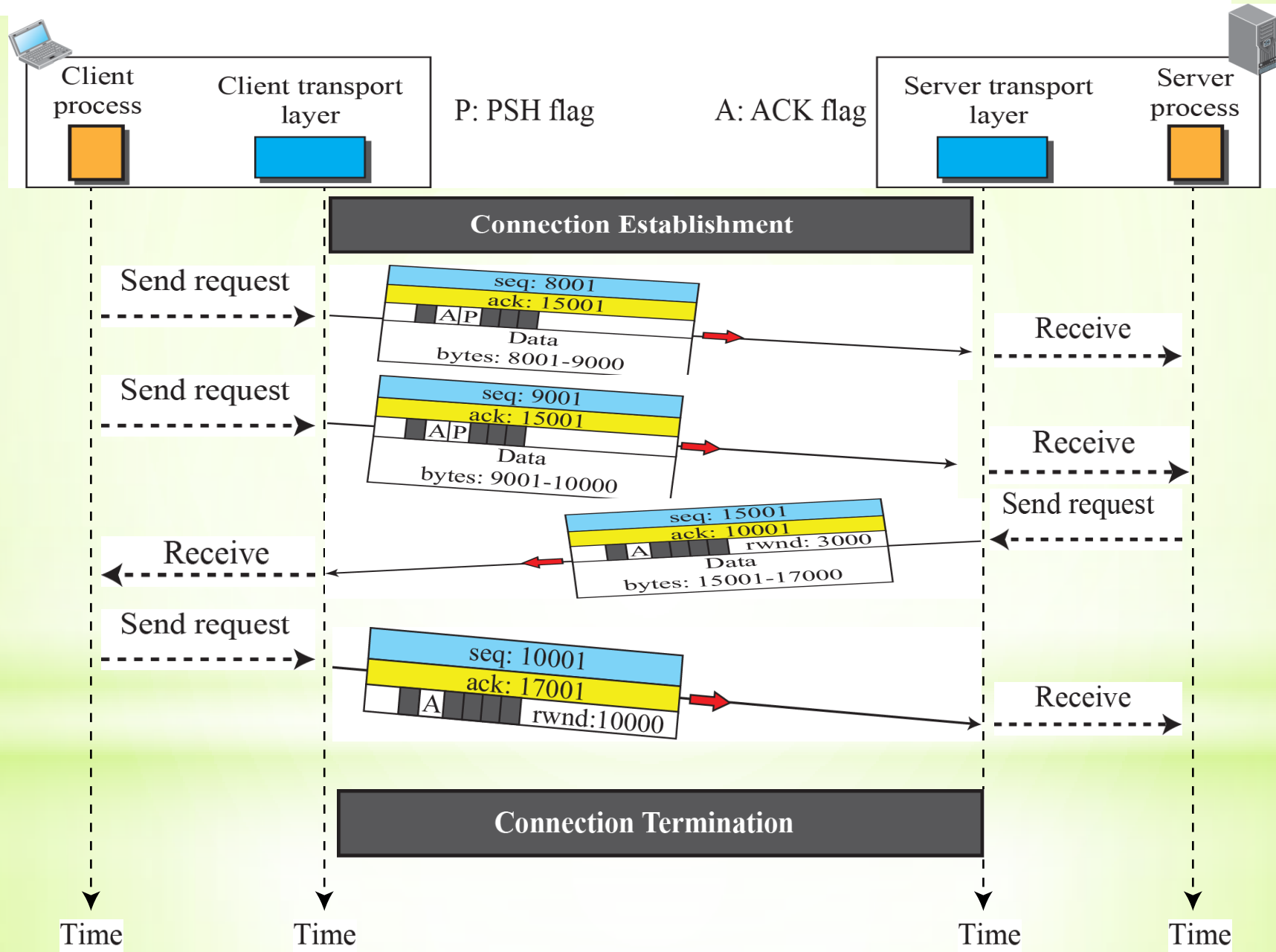


Denial of Service attack / SYN flooding attack

- The connection establishment procedure in TCP is susceptible to a serious security problem called **SYN flooding attack**.
- This SYN flooding attack belongs to a group of security attacks known as a **denial of service attack**
- The server, assuming that the clients are issuing an active open, allocates the necessary resources, such as creating transfer control block (TCB) tables and setting timers.
- One recent strategy is to postpone resource allocation until the server can verify that the connection request is coming from a valid IP address, by using what is called a **cookie**.



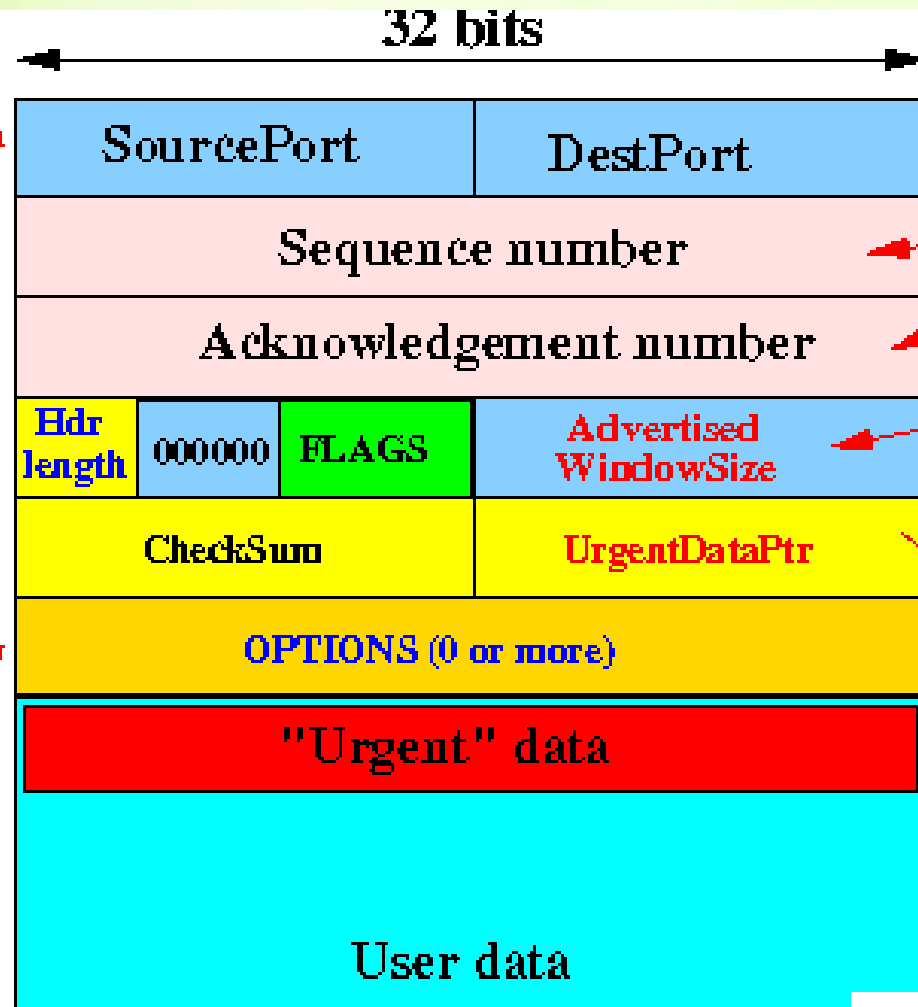
Data transfer



Urgent Pointer

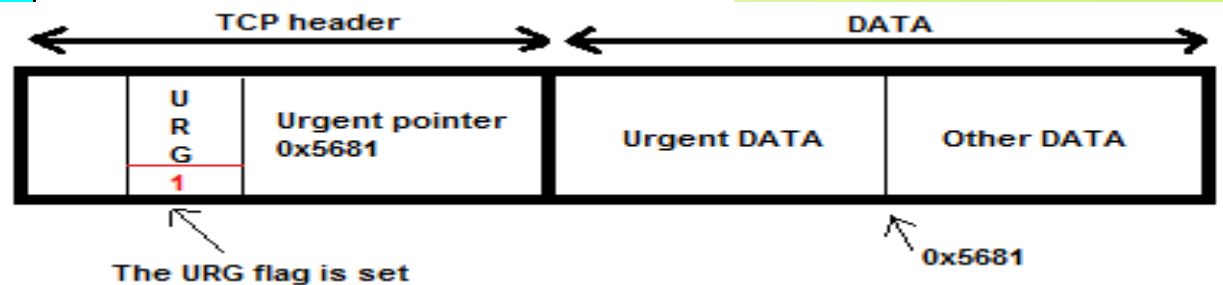
TCP header

Hdr
length
(# words)



Sliding window

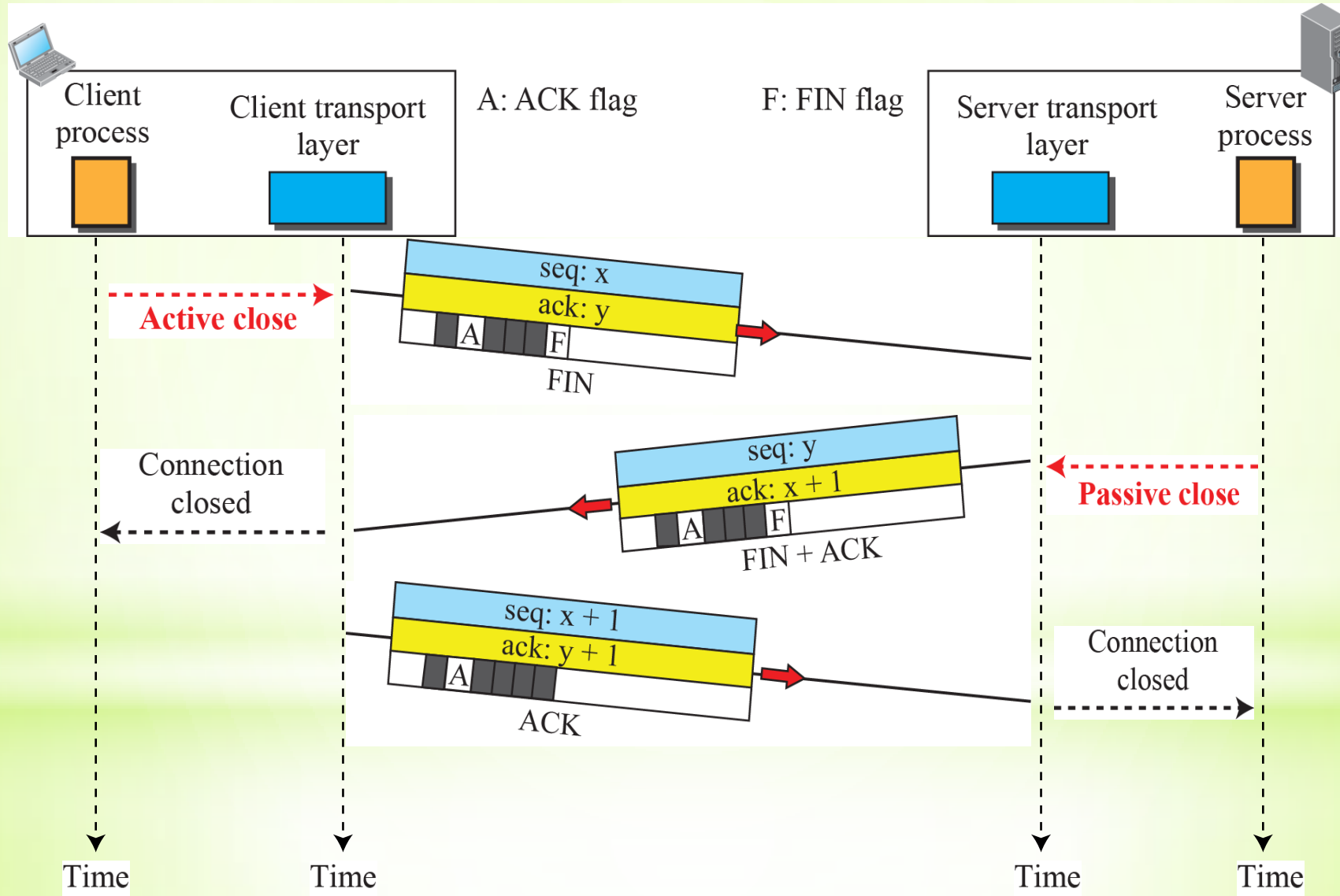
Flow control



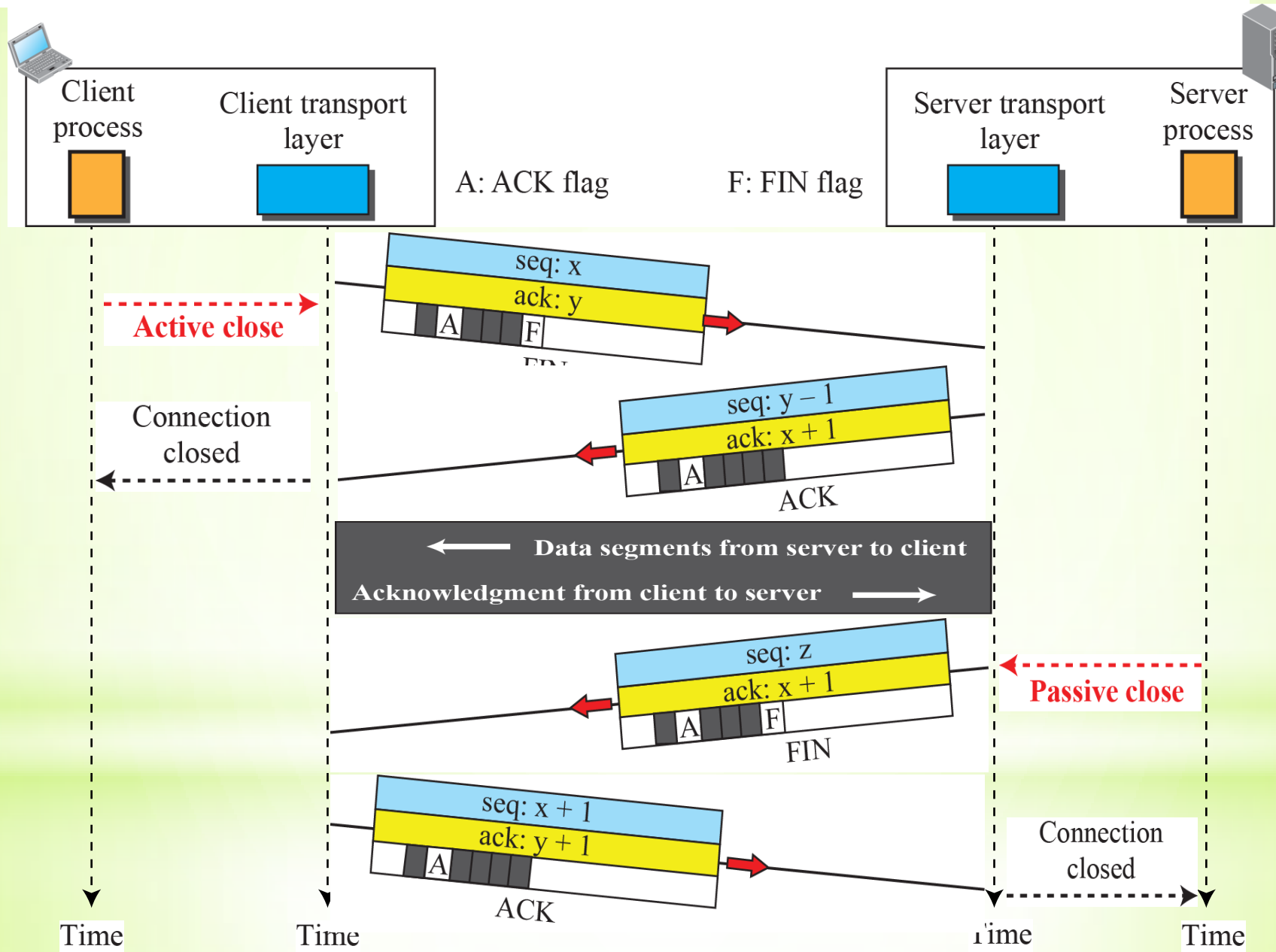
Connection Termination

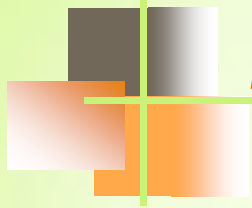
- Most implementations today allow two options for connection termination: three-way handshaking and four-way handshaking with a half-close option.
- *Three-way handshaking* for connection termination includes exchange of three messages between the two parties:
 - FIN segment
 - FIN + ACK segment
 - ACK segment
- The FIN segment consumes one sequence number if it does not carry data.
- The FIN + ACK segment consumes only one sequence number if it does not carry data.

Connection termination using three-way handshaking



Half-close



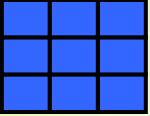


State Transmission Diagram

To keep track of all the different events happening during connection establishment, connection termination, and data transfer, TCP is specified as the finite state machine (FSM).

□ Scenarios

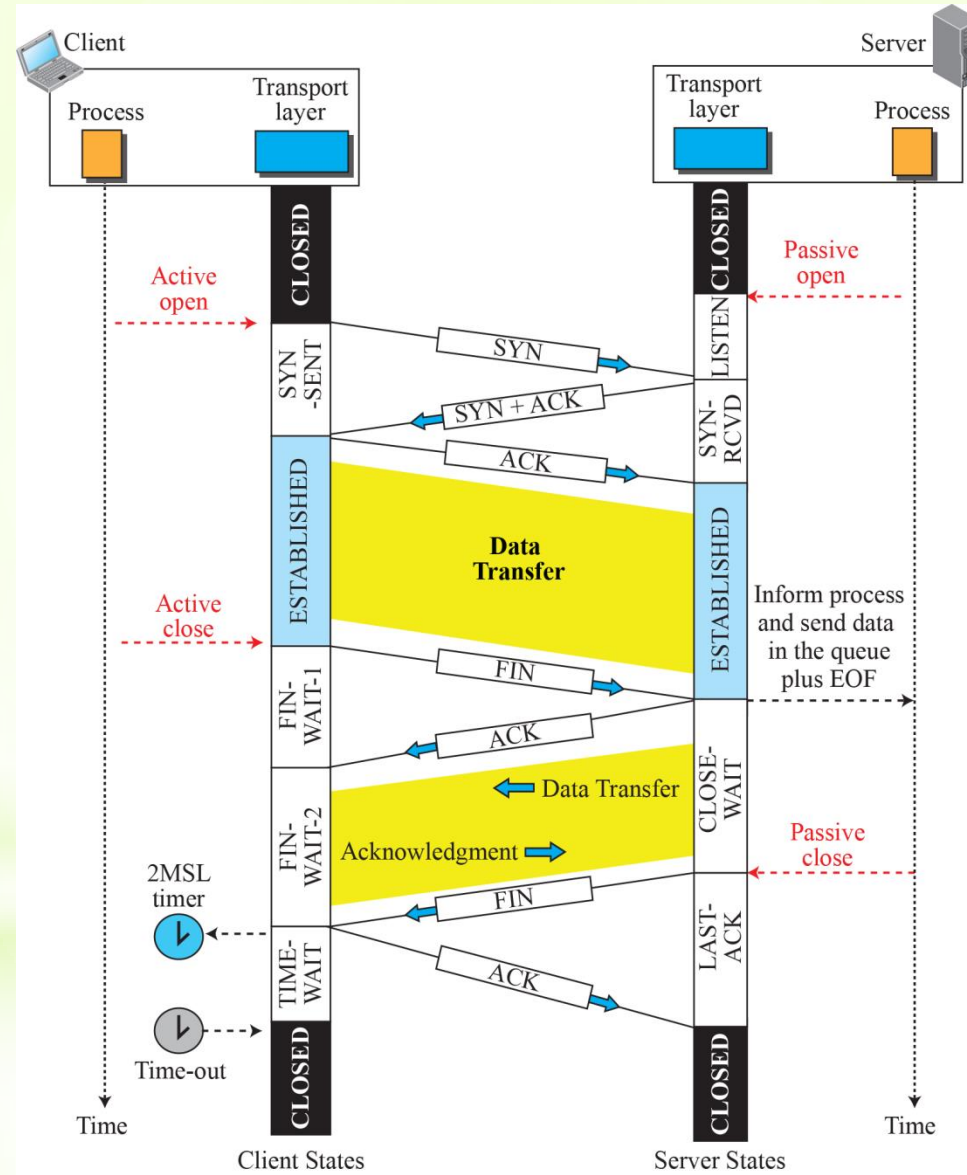
❖ A Half-Close Scenario



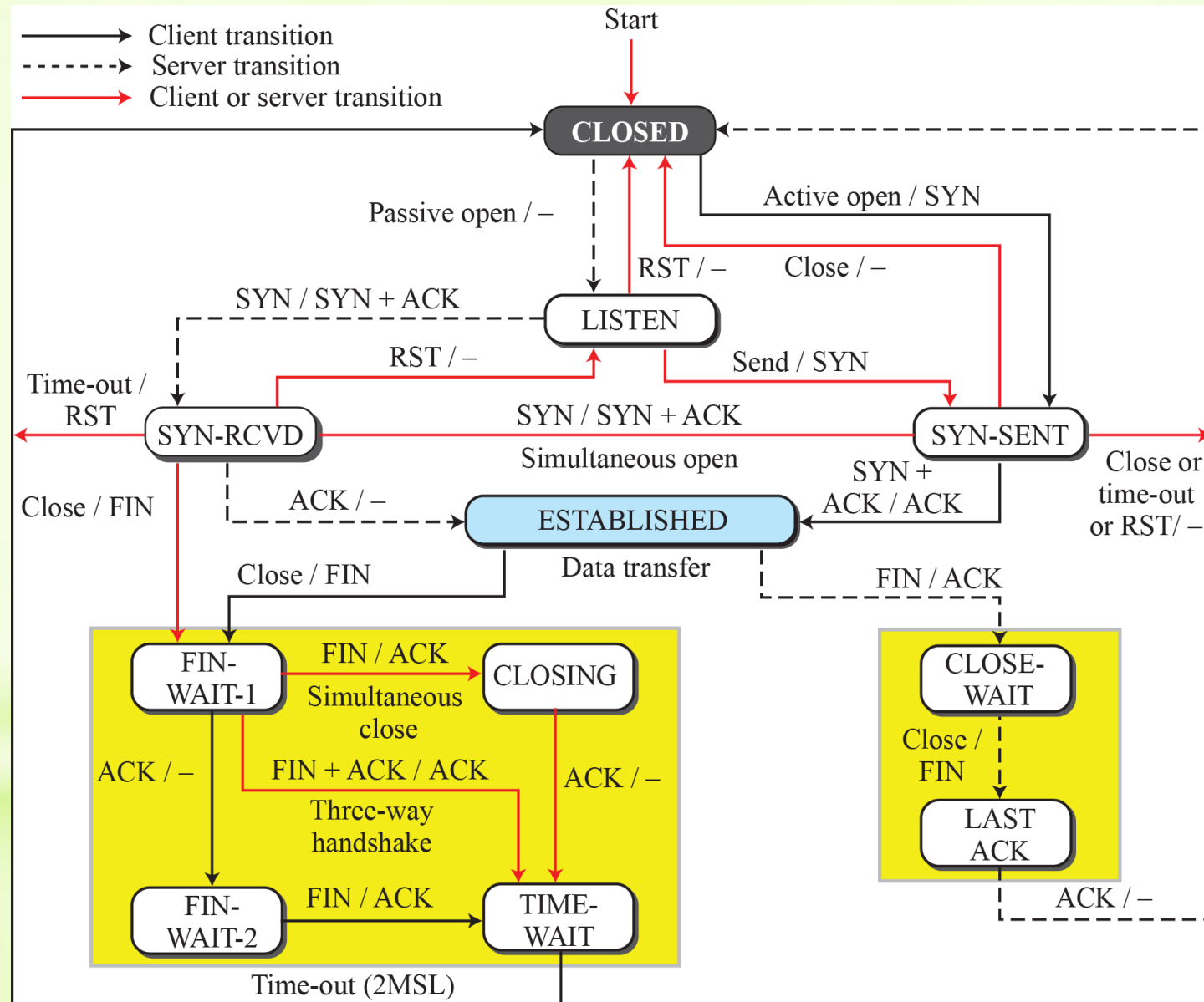
States for TCP

<i>State</i>	<i>Description</i>
CLOSED	No connection exists
LISTEN	Passive open received; waiting for SYN
SYN-SENT	SYN sent; waiting for ACK
SYN-RCVD	SYN+ACK sent; waiting for ACK
ESTABLISHED	Connection established; data transfer in progress
FIN-WAIT-1	First FIN sent; waiting for ACK
FIN-WAIT-2	ACK to first FIN received; waiting for second FIN
CLOSE-WAIT	First FIN received, ACK sent; waiting for application to close
TIME-WAIT	Second FIN received, ACK sent; waiting for 2MSL time-out
LAST-ACK	Second FIN sent; waiting for ACK
CLOSING	Both sides decided to close simultaneously

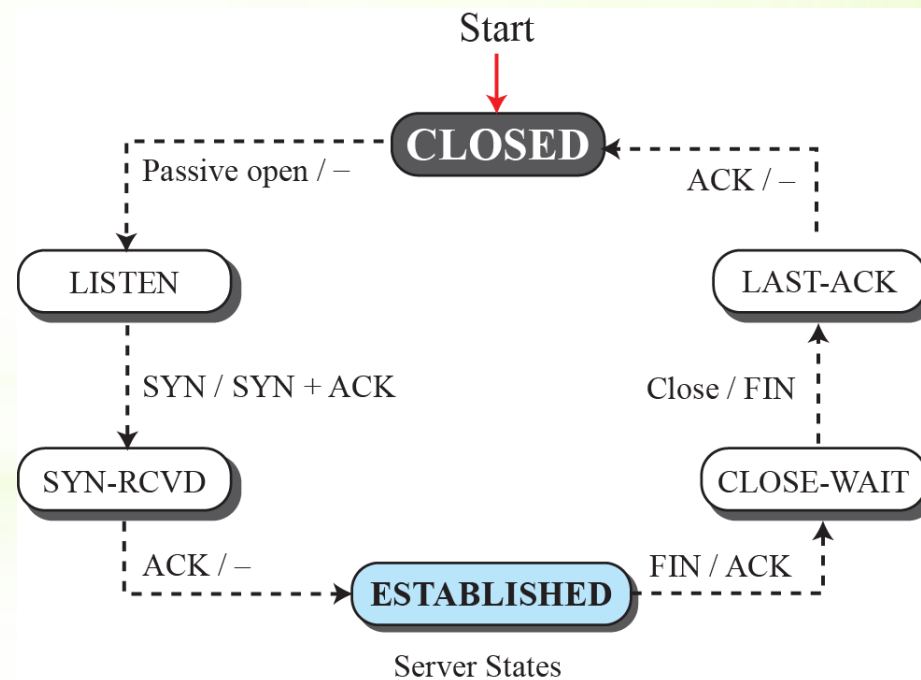
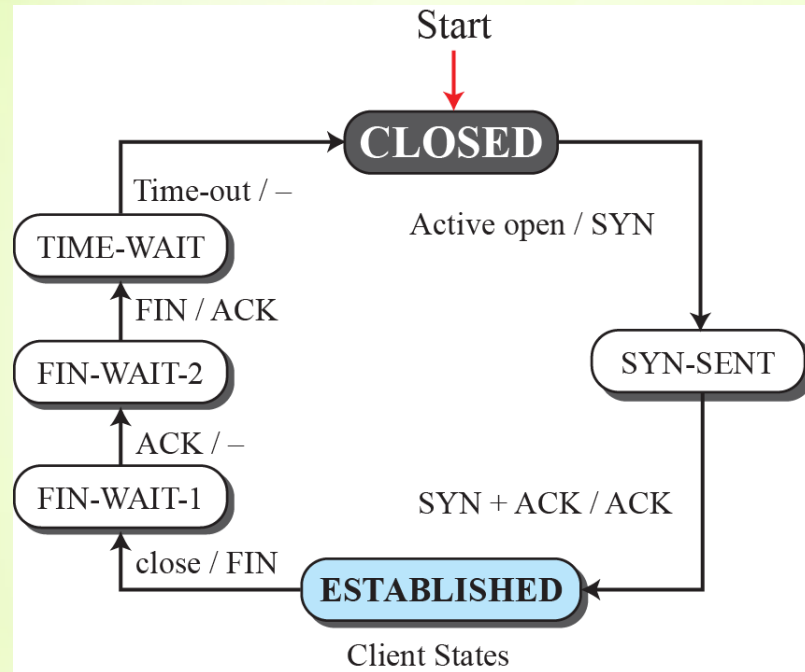
Time-line diagram for a common scenario



State transition diagram



Transition diagram with half-close connection termination





Windows in TCP

TCP uses two windows (send window and receive window) for each direction of data transfer, which means four windows for a bidirectional communication. To make the discussion simple, we make an unrealistic unidirectional; the bidirectional communication can be inferred using two unidirectional communications with piggybacking.

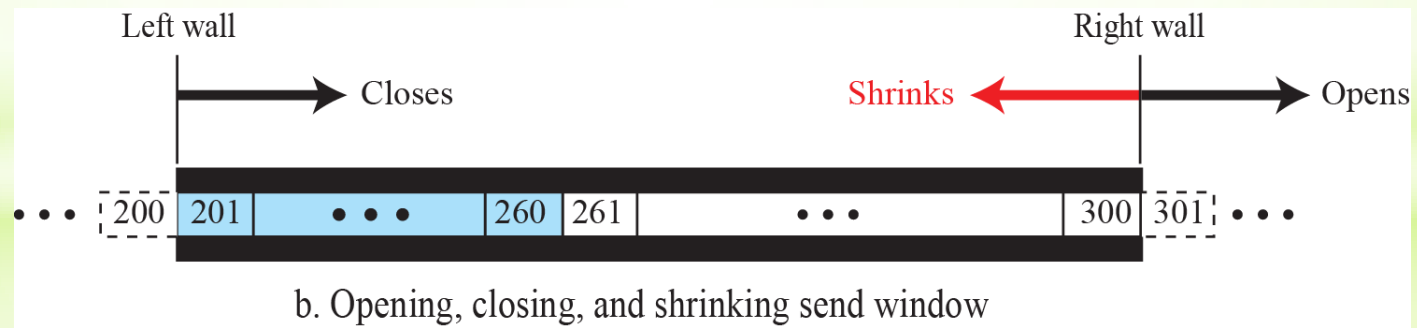
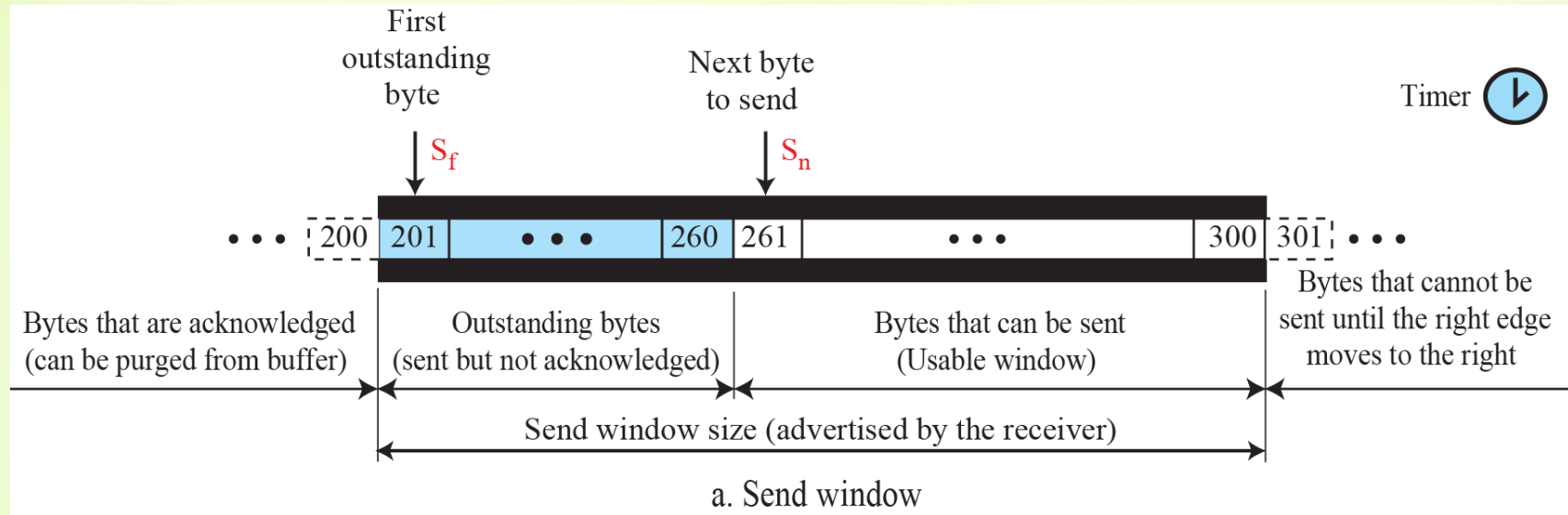
- Send Window

- Receive Window

Send window in TCP

- The send window size is dictated by the receiver (flow control) and the congestion in the underlying network (congestion control).
- A send window *opens, closes, or shrinks*
- The send window in TCP is similar to the one used with the Selective-Repeat protocol, but with some differences:
 - The window size in SR is the number of packets, but the window size in TCP is the number of bytes. Although actual transmission in TCP occurs segment by segment, the variables that control the window are expressed in bytes.
 - TCP can store data received from the process and send them later
 - The theoretical Selective-Repeat protocol may use several timers for each packet sent, but as mentioned before, the TCP protocol uses only one timer.

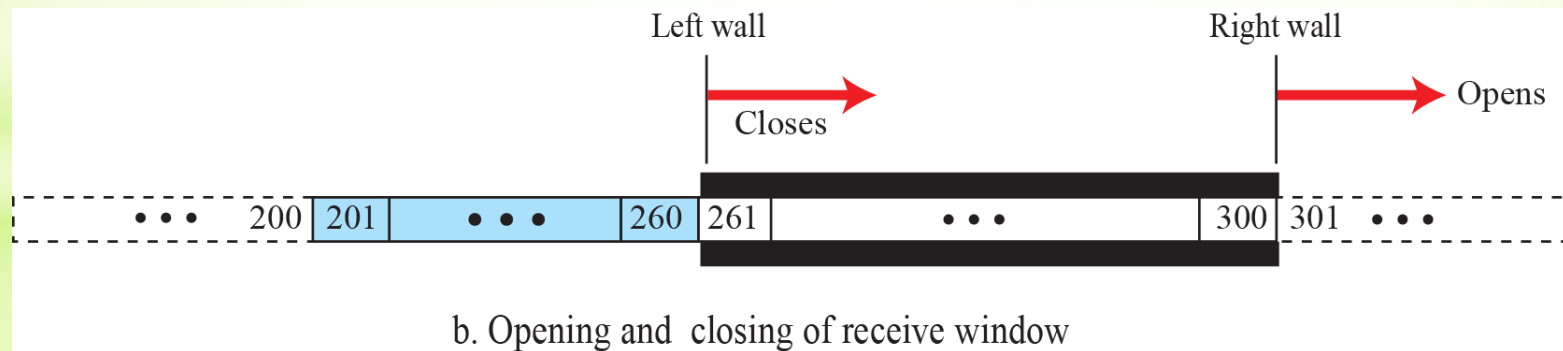
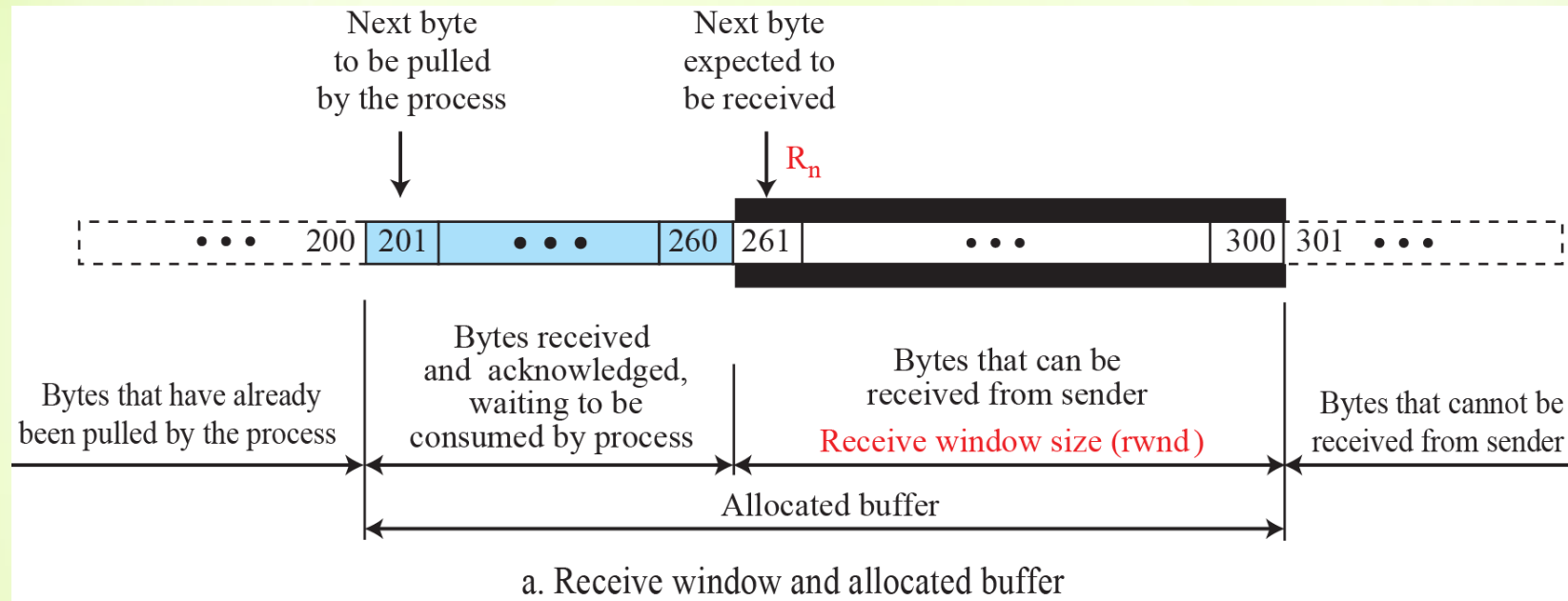
Send window in TCP

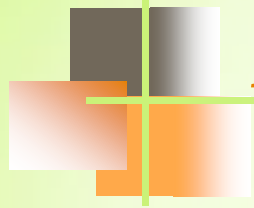


Receive window in TCP

- The receive window opens and closes; in practice, the window should never shrink.
- There are two differences between the receive window in TCP and the one we used for SR.
 - The first difference is that TCP allows the receiving process to pull data at its own pace. The receive window size is then always smaller than or equal to the buffer size
 $rwnd = \text{buffer size} - \text{number of waiting bytes to be pulled}$
 - Remember that an acknowledgement in SR is selective, defining the uncorrupted packets that have been received. The major acknowledgment mechanism in TCP is a cumulative acknowledgment announcing the next expected byte to receive (in this way TCP looks like GBN).

Receive window in TCP





Flow Control

As discussed before, flow control balances the rate a producer creates data with the rate a consumer can use the data. TCP separates flow control from error control. In this section we discuss flow control, ignoring error control. We assume that the logical channel between the sending and receiving TCP is error-free.



(continued)

- ❑ Opening and Closing Windows

- ❖ A Scenario

- ❑ Shrinking of Windows

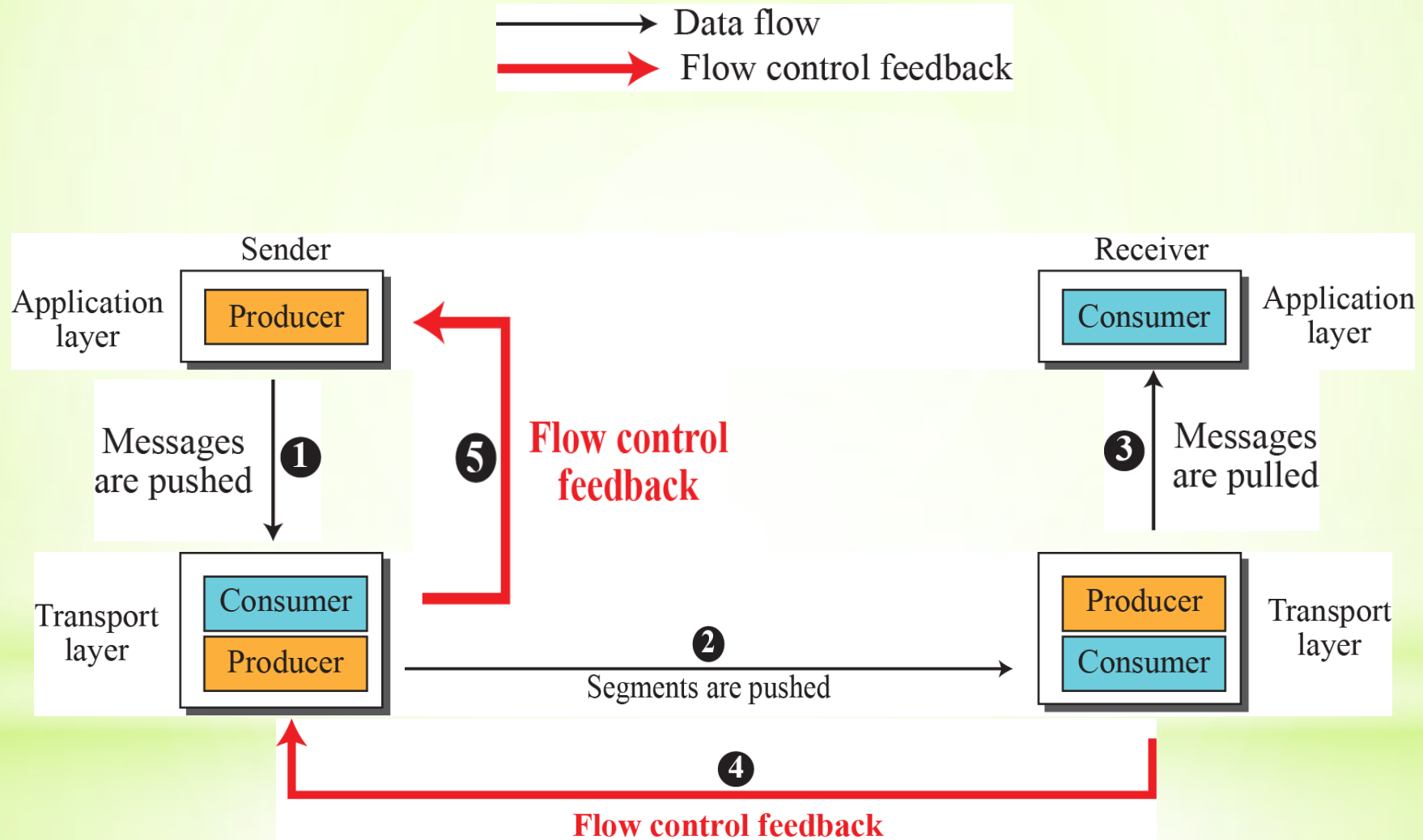
- ❖ Window Shutdown

- ❑ Silly Window Syndrome

- ❖ Syndrome Created by the Sender

- ❖ Syndrome Created by the Receiver

Data flow and flow control feedbacks in TCP



An example of flow control

Note: We assume only unidirectional communication from client to server. Therefore, only one window at each side is shown.

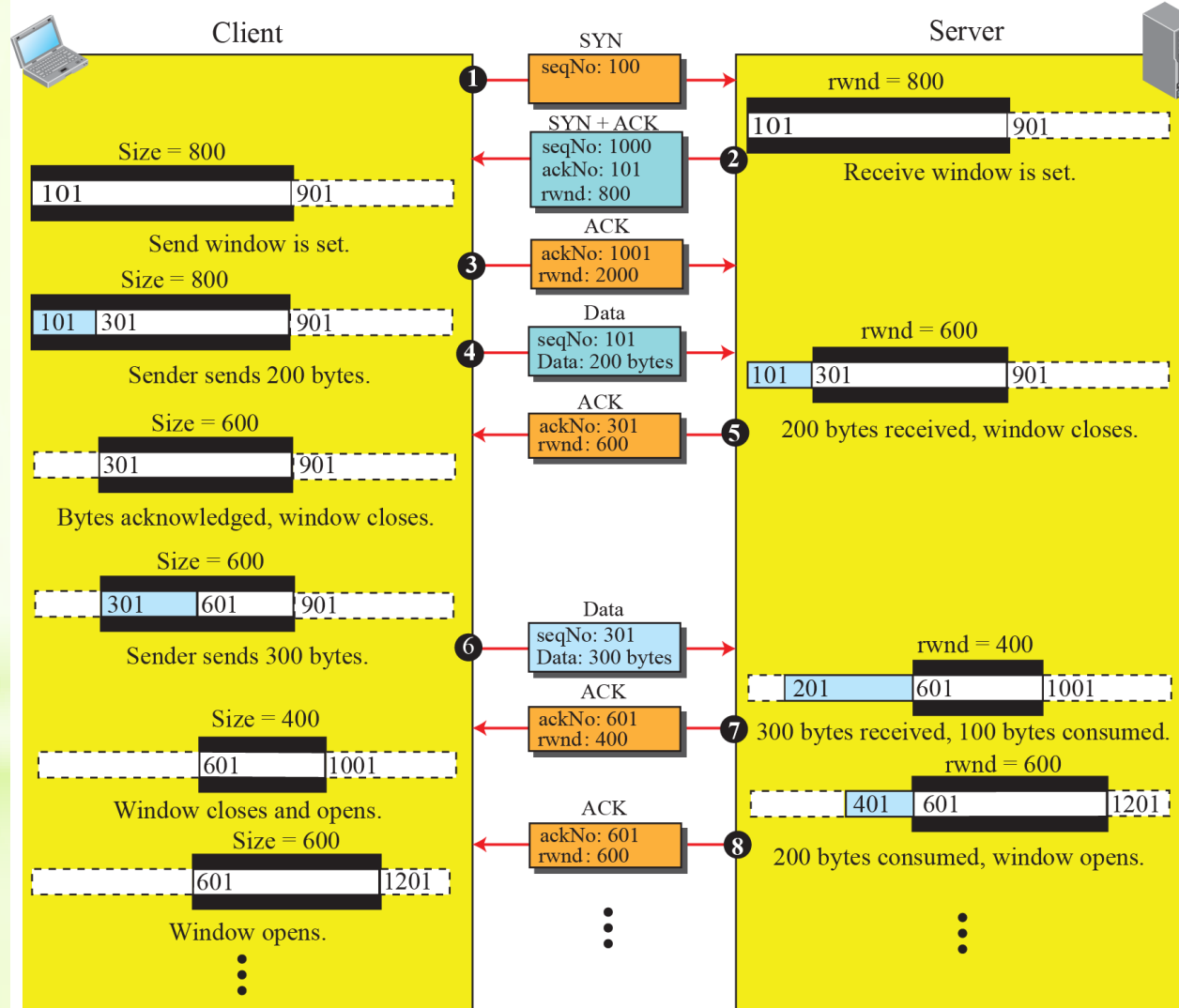


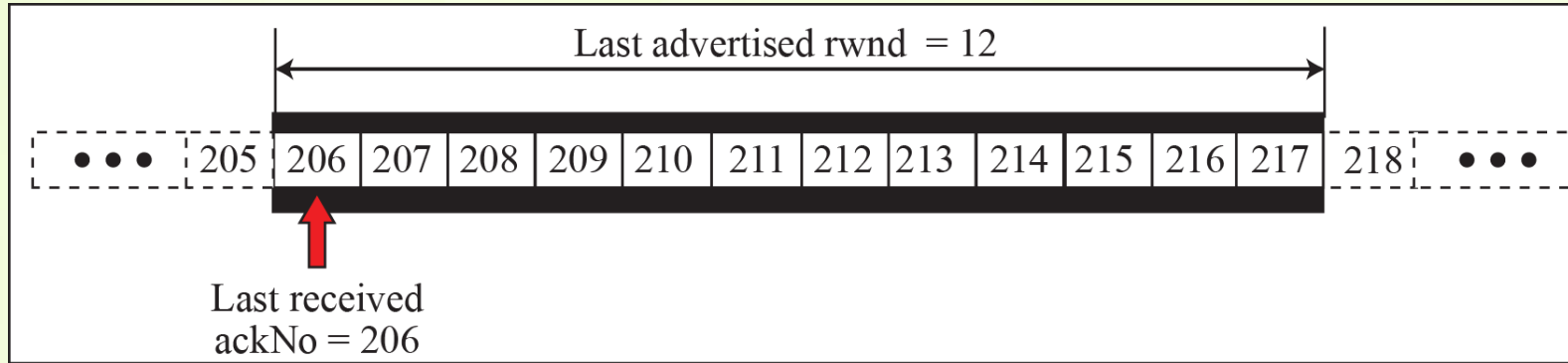
Figure shows the reason for this mandate.

Part a of the figure shows the values of the last acknowledgment and *rwnd*. Part b shows the situation in which the sender has sent bytes 206 to 214. Bytes 206 to 209 are acknowledged and purged. The new advertisement, however, defines the new value of *rwnd* as 4, in which $210 + 4 < 206 + 12$. When the send window shrinks, it creates a problem: byte 214, which has already been sent, is outside the window. The relation discussed before forces the receiver to maintain the right-hand wall of the window to be as shown in part a, because the receiver does not know which of the bytes 210 to 217 has already been sent. described above.

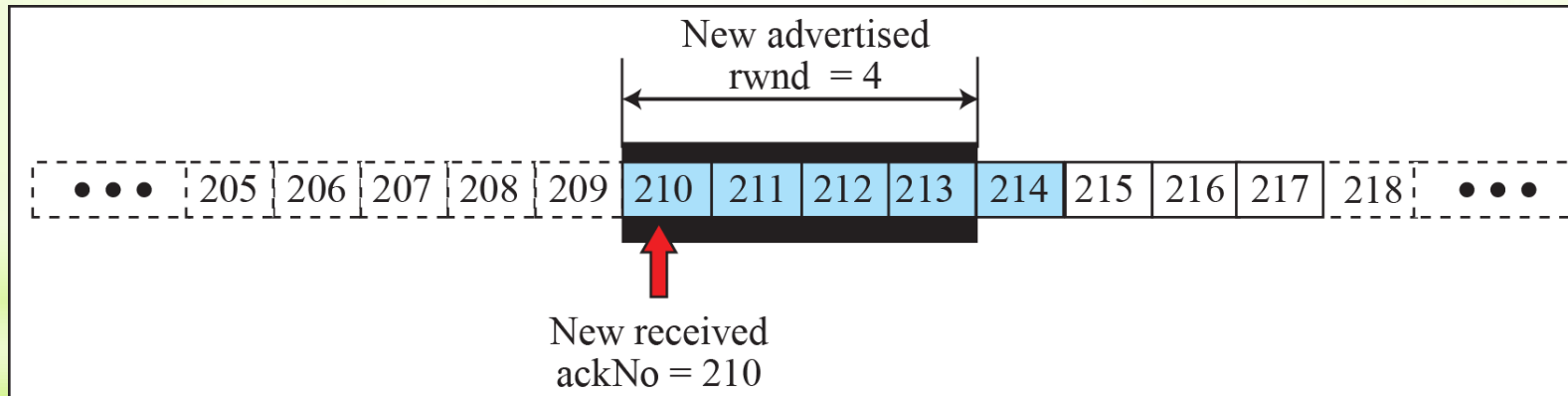
Shrinking of Windows

$$\text{new ackNo} + \text{new } rwnd \geq \text{last ackNo} + \text{last } rwnd$$

Example



a. The window after the last advertisement



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

Window Shutdown - the receiver can temporarily shut down the window by sending a *rwnd* of 0.

Silly Window Syndrome

A serious problem can arise in the sliding window operation when either the sending application program creates data slowly or the receiving application program consumes data slowly, or both. Any of these situations results in the sending of data in very small segments, which reduces the efficiency of the operation.

❖ Syndrome Created by the Sender

- TCP may create a silly window syndrome if it is serving an application program that creates data slowly, for example, 1 byte at a time.
- The result is a lot of 41-byte segments that are traveling through an internet.
- The solution is to prevent the sending TCP from sending the data byte by byte. The sending TCP must be forced to wait and collect data to send in a larger block. **Nagle** found an elegant solution.

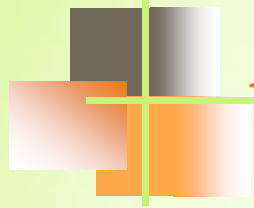
Nagle's algorithm is simple:

1. The sending TCP sends the first piece of data it receives from the sending application program even if it is only 1 byte.
 2. After sending the first segment, the sending TCP accumulates data in the output buffer and waits until either the receiving TCP sends an acknowledgment or until enough data have accumulated to fill a maximum-size segment. At this time, the sending TCP can send the segment.
 3. Step 2 is repeated for the rest of the transmission. Segment 3 is sent immediately if an acknowledgment is received for segment 2, or if enough data have accumulated to fill a maximum-size segment.
- ❖ If the application program is faster than the network, the segments are larger (maximum-size segments). If the application program is slower than the network, the segments are smaller (less than the maximum segment size).

❖ Syndrome Created by the Receiver

Suppose that the sending application program creates data in blocks of 1 kilobyte, but the receiving application program consumes data 1 byte at a time. Two solutions have been proposed:

- **Clark's solution** is to send an acknowledgment 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.
- The second solution is to delay sending the acknowledgment. This means that when a segment arrives, it is not acknowledged immediately. The receiver waits until there is a decent amount of space in its incoming buffer before acknowledging the arrived segments.
 - Delayed acknowledgment also has another advantage: it reduces traffic. The receiver does not have to acknowledge each segment.
 - However, there also is a disadvantage in that the delayed acknowledgment may result in the sender unnecessarily retransmitting the unacknowledged segments.



Error Control

TCP is a reliable transport-layer protocol. This means that an application program that delivers a stream of data to TCP relies on TCP to deliver the entire stream to the application program on the other end in order, without error, and without any part lost or duplicated.



- ❑ Checksum

- ❑ Acknowledgment

 - ❖ Cumulative Acknowledgment (ACK)

 - ❖ Selective Acknowledgment (SACK)

- ❑ Generating Acknowledgments

- ❑ Retransmission

 - ❖ Retransmission after RTO

 - ❖ Retransmission after Three Duplicate ACK

- ❑ Out-of-Order Segments



(continued)

□ FSMs for Data Transfer in TCP

- ❖ Sender-Side FSM
- ❖ Receiver-Side FSM

□ Some Scenarios

- ❖ Normal Operation
- ❖ Lost Segment
- ❖ Fast Retransmission
- ❖ Delayed Segment
- ❖ Duplicate Segment
- ❖ Automatically Corrected Lost ACK
- ❖ Correction by Resending a Segment
- ❖ Deadlock Created by Lost Acknowledgment

❑ Checksum

- Each segment includes a checksum field, which is used to check for a corrupted segment.
- If a segment is corrupted, as detected by an invalid checksum, the segment is discarded by the destination TCP and is considered as lost.

❑ Acknowledgment

ACK segments do not consume sequence numbers and are not acknowledged.

- ❖ **Cumulative Acknowledgment (ACK):** The receiver advertises the next byte it expects to receive, ignoring all segments received and stored out of order. This is sometimes referred to as *positive cumulative acknowledgment*, or ACK and is specified as 32-bit ACK field in the TCP header
- ❖ **Selective Acknowledgment (SACK):** A SACK reports a block of bytes that is out of order, and also a block of bytes that is duplicated, i.e., received more than once. SACK is implemented as an option at the end of the TCP header.

❑ Retransmission

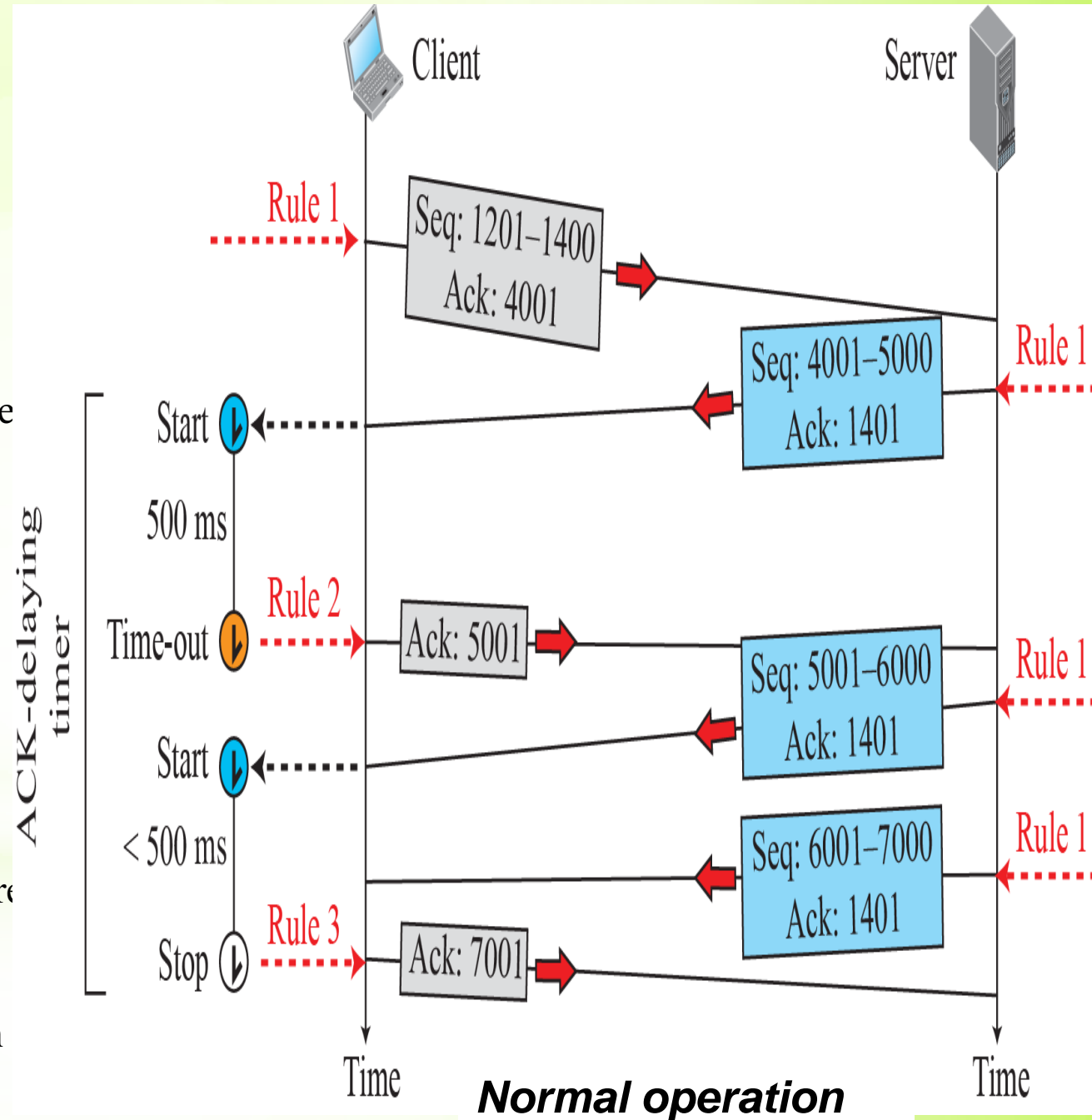
- **Retransmission after RTO:** The sending TCP maintains one **retransmission time-out (RTO)** for each connection. When the timer matures, i.e. times out, TCP resends the segment in the front of the queue (the segment with the smallest sequence number) and restarts the timer. RTO is dynamic in TCP and is updated based on the **round-trip time (RTT)** of segments.
- **Retransmission after Three Duplicate ACK:** To expedite service throughout the Internet by allowing senders to retransmit without waiting for a time out, most implementations today follow the three duplicate ACKs rule and retransmit the missing segment immediately. This feature is called ***fast retransmission***.

❑ Out-of-Order Segments

- TCP implementations today do not discard out-of-order segments. They store them temporarily and flag them as out-of-order segments until the missing segments arrive.
- **TCP guarantees that no out-of-order data are delivered to the process.**

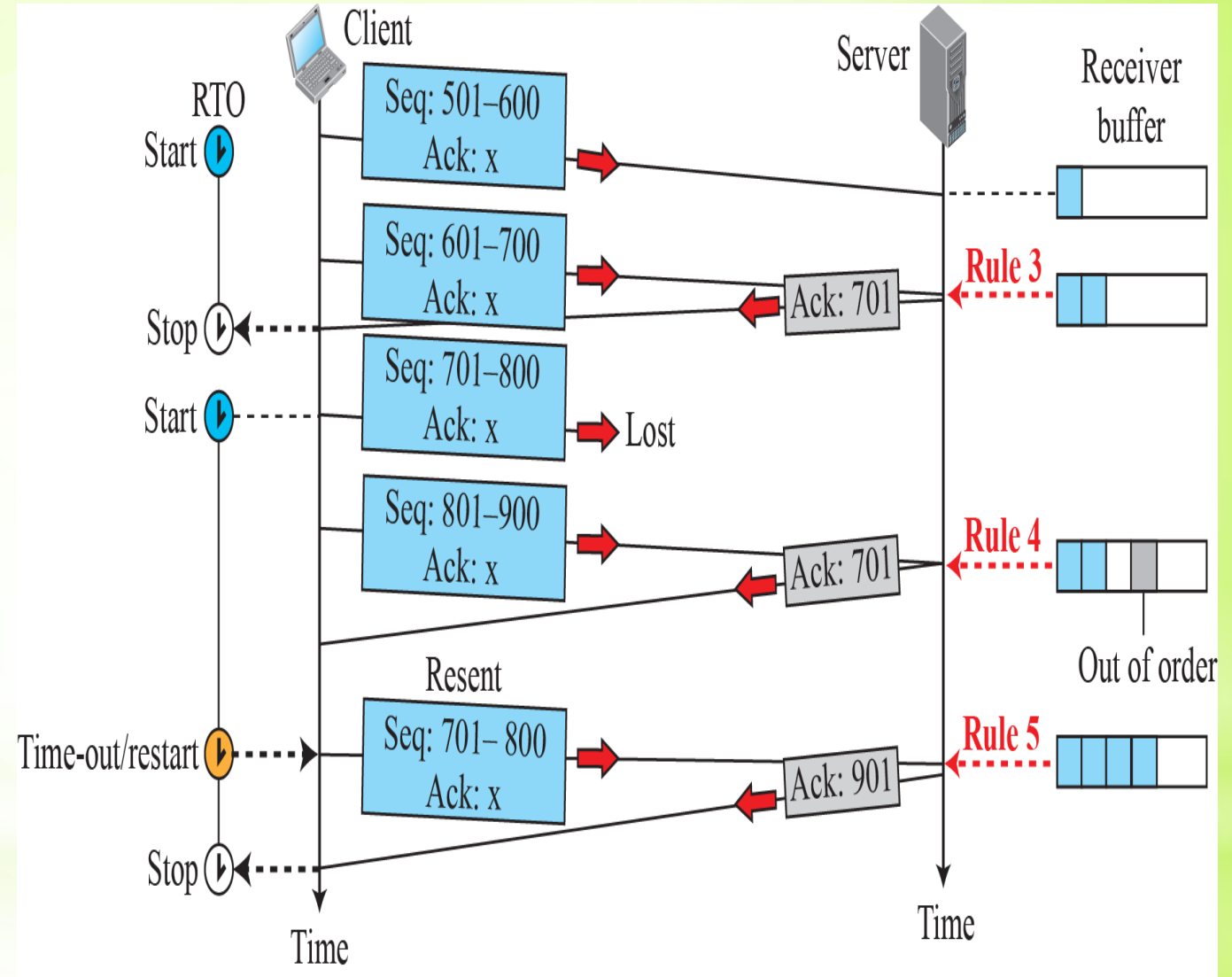
□ Generating Acknowledgments

1. When end A sends a data segment to end B, it must include (piggyback) an acknowledgment that gives the next sequence number it expects to receive. This rule decreases the number of segments needed and therefore reduces traffic.
2. When the receiver has no data to send and it receives an in-order segment (with expected sequence number) and the previous segment has already been acknowledged, the receiver delays sending an ACK segment until another segment arrives or until a period of time (normally 500 ms) has passed. In other words, the receiver needs to delay sending an ACK segment if there is only one outstanding in-order segment. This rule reduces ACK segments.
3. When a segment arrives with a sequence number that is expected by the receiver, and the previous in-order segment has not been acknowledged, the receiver immediately sends an ACK segment. In other words, there should not be more than two in-order unacknowledged segments at any time. This prevents the unnecessary retransmission of segments that may create congestion in the network.



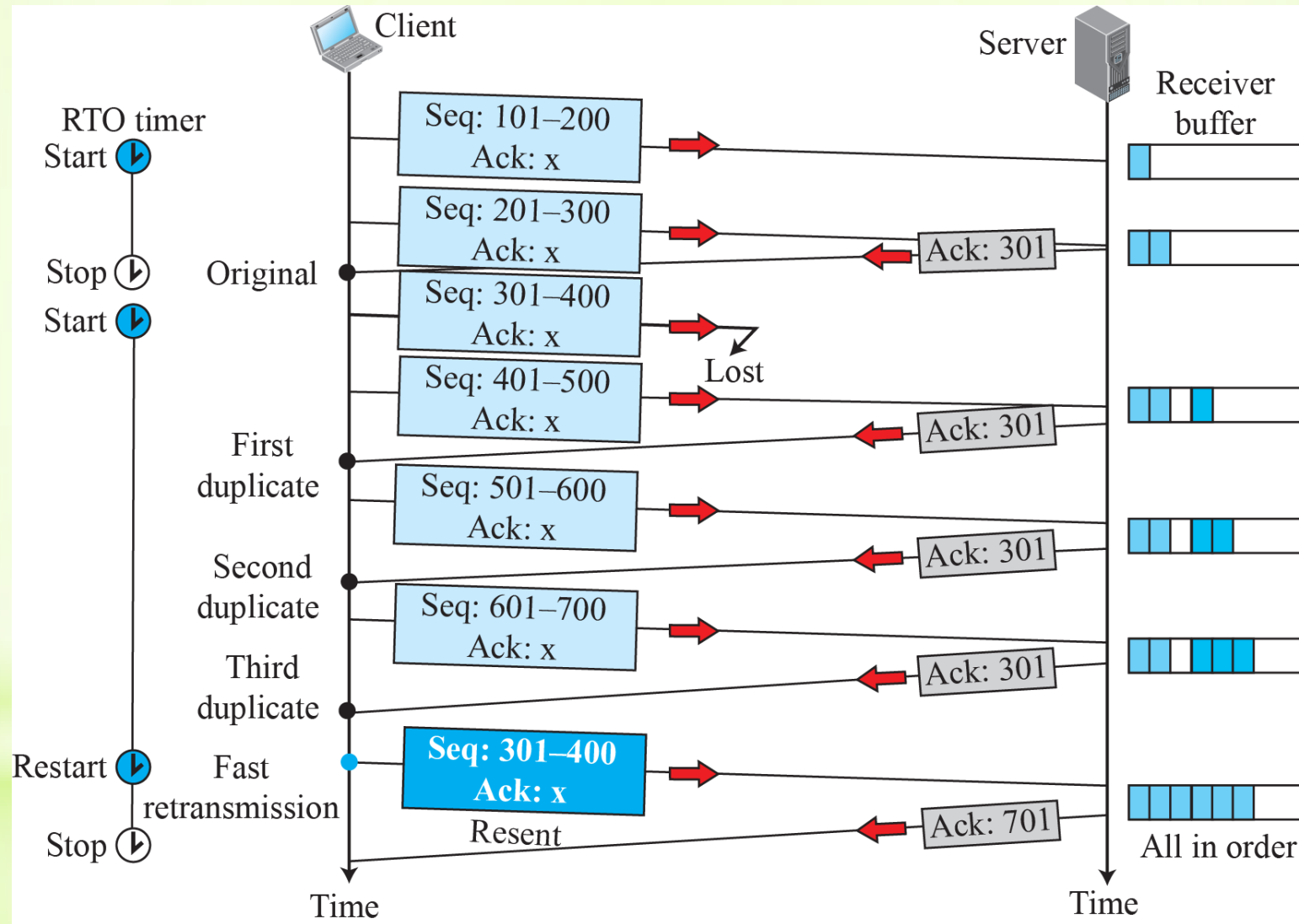
□ Generating Acknowledgments

4. When a segment arrives with an out-of-order sequence number that is higher than expected, the receiver immediately sends an ACK segment announcing the sequence number of the next expected segment. This leads to the **fast retransmission** of missing segments.
5. When a missing segment arrives, the receiver sends an ACK segment to announce the next sequence number expected. This informs the receiver that segments reported missing have been received.

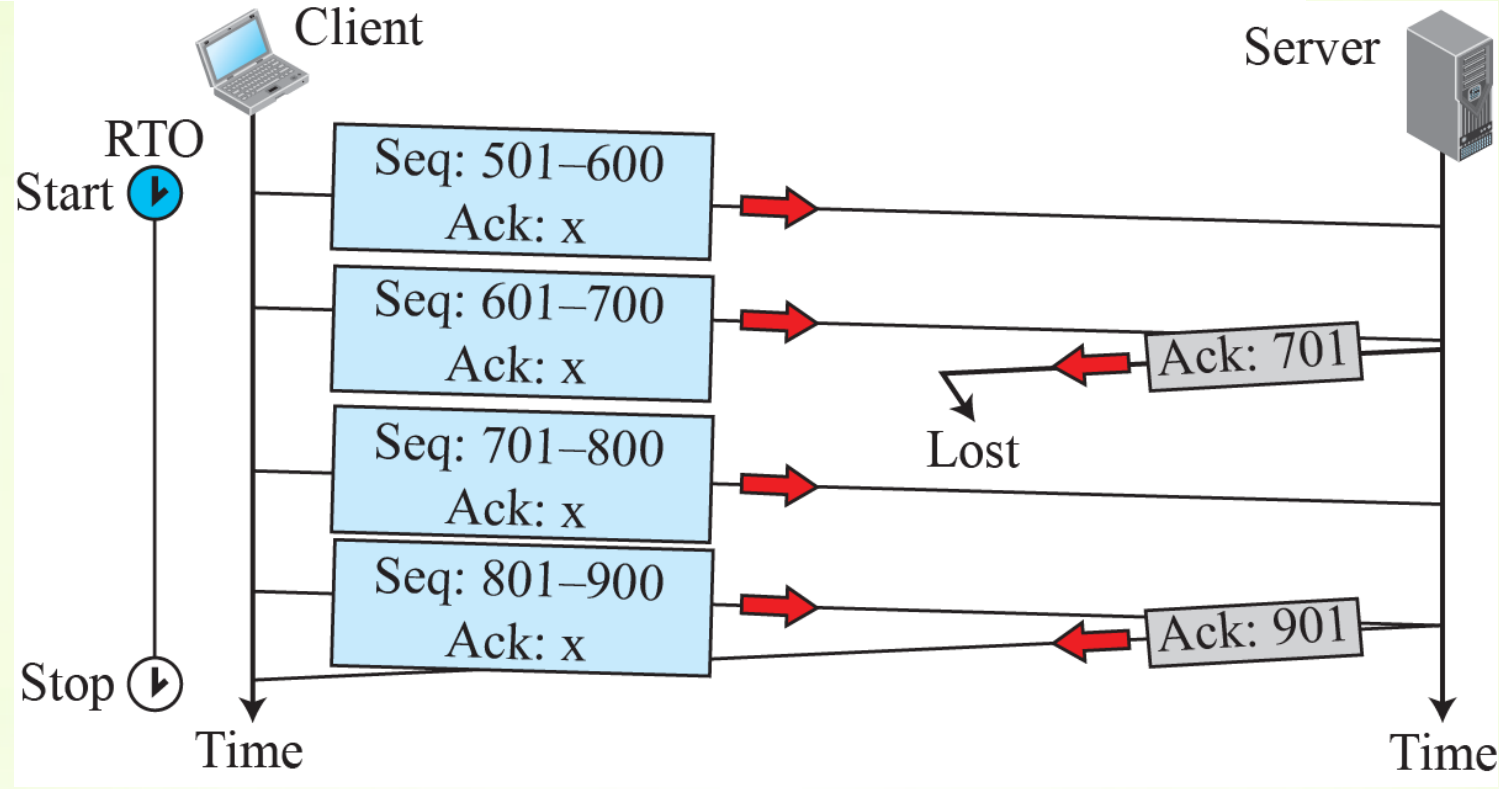


Lost segment

Fast retransmission (Rule-4)



Lost acknowledgment



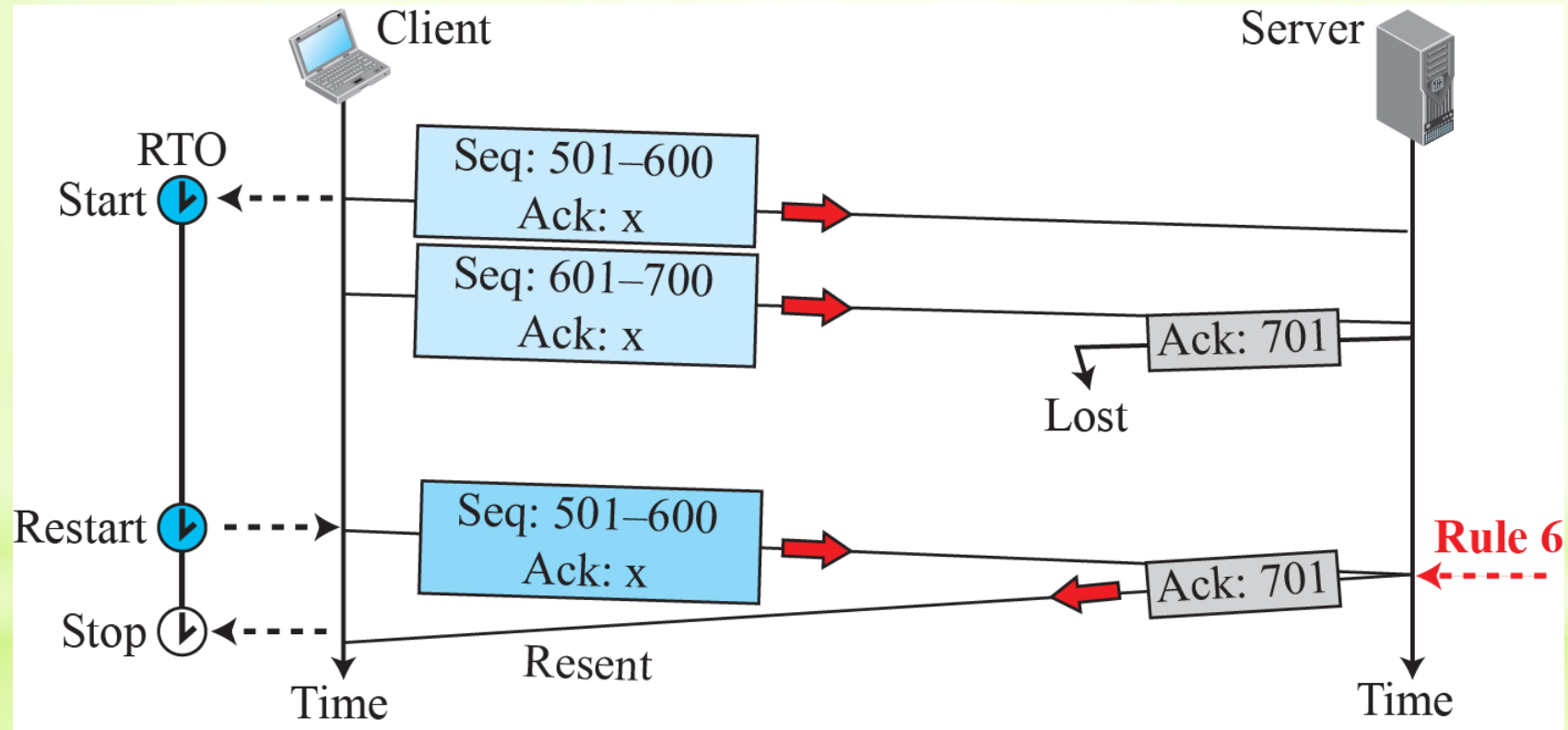
Deadlock Created by Lost Acknowledgment

Lost acknowledgments may create deadlock if they are not properly handled.

Eg: `ack(wsize=0)` → sender shutdown Window :: `ack(wsize!=0)` :: Ack Lost ::
Sender still in shutdown mode

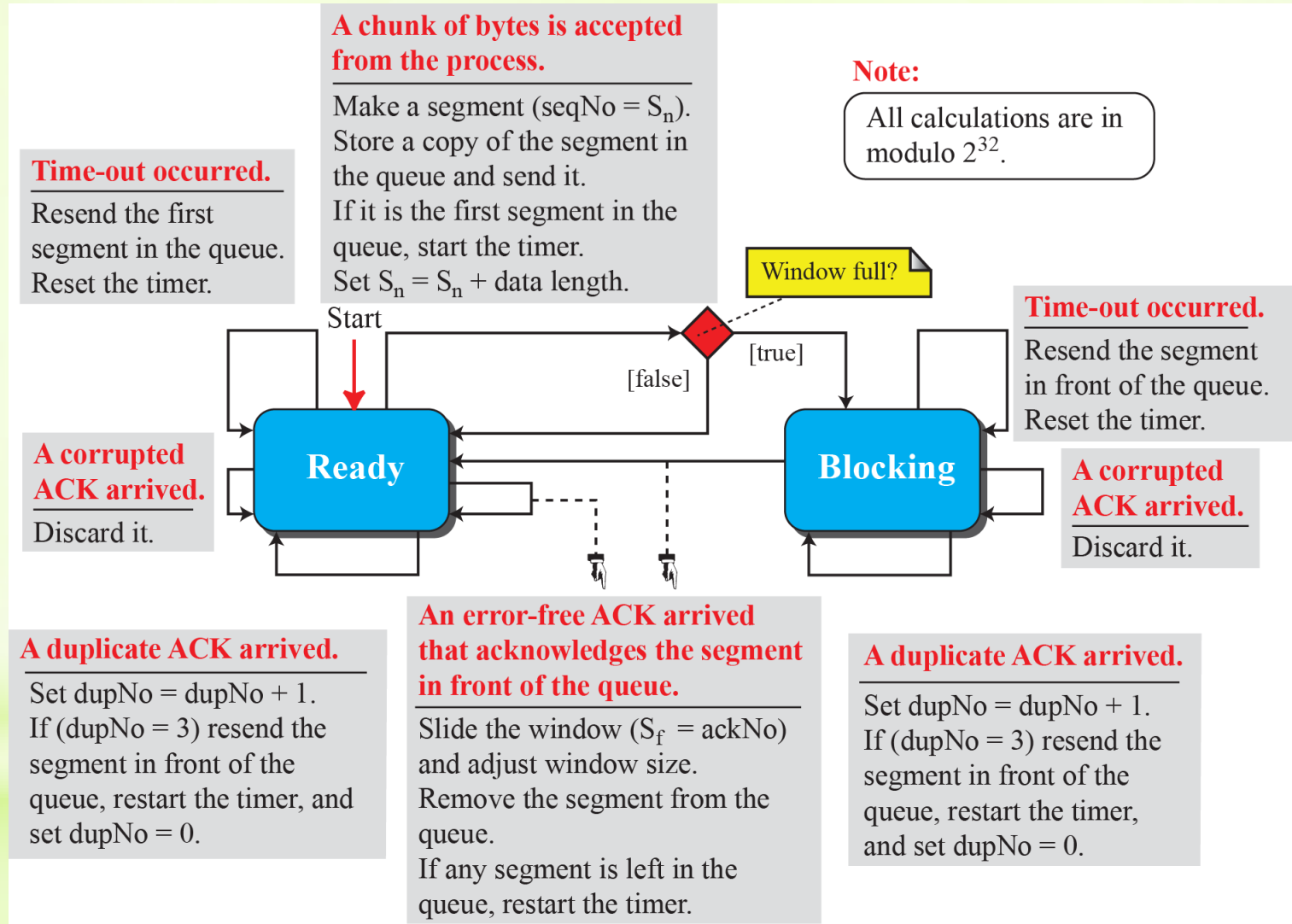
□ Generating Acknowledgments

6. If a duplicate segment arrives, the receiver discards the segment, but immediately sends an acknowledgment indicating the next in-order segment expected. This solves some problems when an ACK segment itself is lost.



Lost acknowledgment corrected by resending a segment

Simplified FSM for the TCP sender side



Simplified FSM for the TCP receiver side

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.

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).

An expected error-free segment arrived.

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

ACK-delaying timer expired.

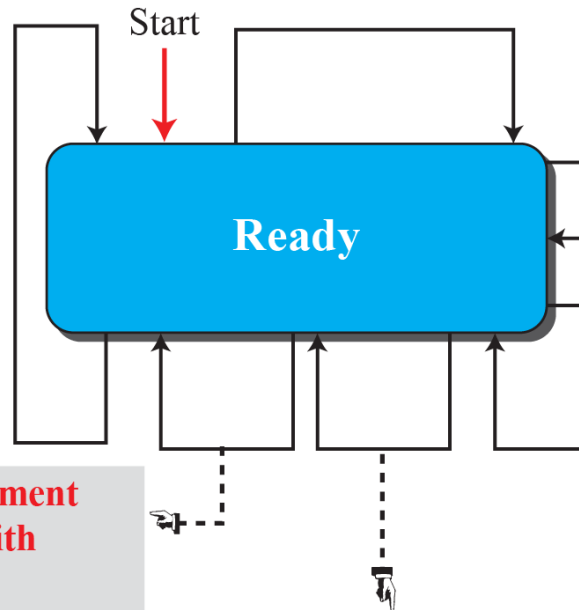
Send the delayed ACK.

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).

A corrupted segment arrived.

Discard the segment.





TCP Congestion Control

TCP uses different policies to handle the congestion in the network. We describe these policies in this section.

- ❑ Congestion Window & Receiver window
- ❑ Congestion Detection
 - ❖ Time out
 - ❖ Three Selective Acknowledgements
- ❑ Congestion Policies
 - ❖ Slow Start: Exponential Increase
 - ❖ Congestion Avoidance: Additive Increase
 - ❖ Fast Retransmission/Fast Recovery



(continued)

- ❑ Policy Transition

- ❖ Tahoe TCP

- ❖ Reno TCP

- ❖ NewReno TCP

- ❑ Additive Increase, Multiplicative Decrease

❑ Congestion Window & Receiver window

TCP is an end-to-end protocol that uses the service of IP. The congestion in the router is in the IP territory and should be taken care of by IP. IP is a simple protocol with no congestion control. TCP, itself, needs to be responsible for this problem. TCP cannot ignore the congestion in the network; it cannot aggressively send segments to the network. The result of such aggressiveness would hurt the TCP itself

TCP cannot be very conservative, either, sending a small number of segments in each time interval, because this means not utilizing the available bandwidth of the network. TCP needs to define policies that accelerate the data transmission when there is no congestion and decelerate the transmission when congestion is detected.

To control the number of segments to transmit, TCP uses another variable called a *congestion window*, *cwnd*, whose size is controlled by the congestion situation in the network (as we will explain shortly). The *cwnd* variable and the *rwnd* variable together define the size of the send window in TCP.

Actual window size = minimum (*rwnd*, *cwnd*)

❑ Congestion Detection

The TCP sender uses the occurrence of two events as signs of congestion in the network: time-out and receiving three duplicate ACKs.

- ❖ **Time out :** If a TCP sender does not receive an ACK for a segment or a group of segments before the time-out occurs, it assumes that the corresponding segment or segments are lost and the loss is due to congestion.
- ❖ **Three Selective Acknowledgements:** Recall that when a TCP receiver sends a duplicate ACK, it is the sign that a segment has been delayed, but sending three duplicate ACKs is the sign of a missing segment, which can be due to congestion in the network. When a receiver sends three duplicate ACKs, it means that one segment is missing, but three segments have been received. (slightly congested)

Maximum segment size (MSS)

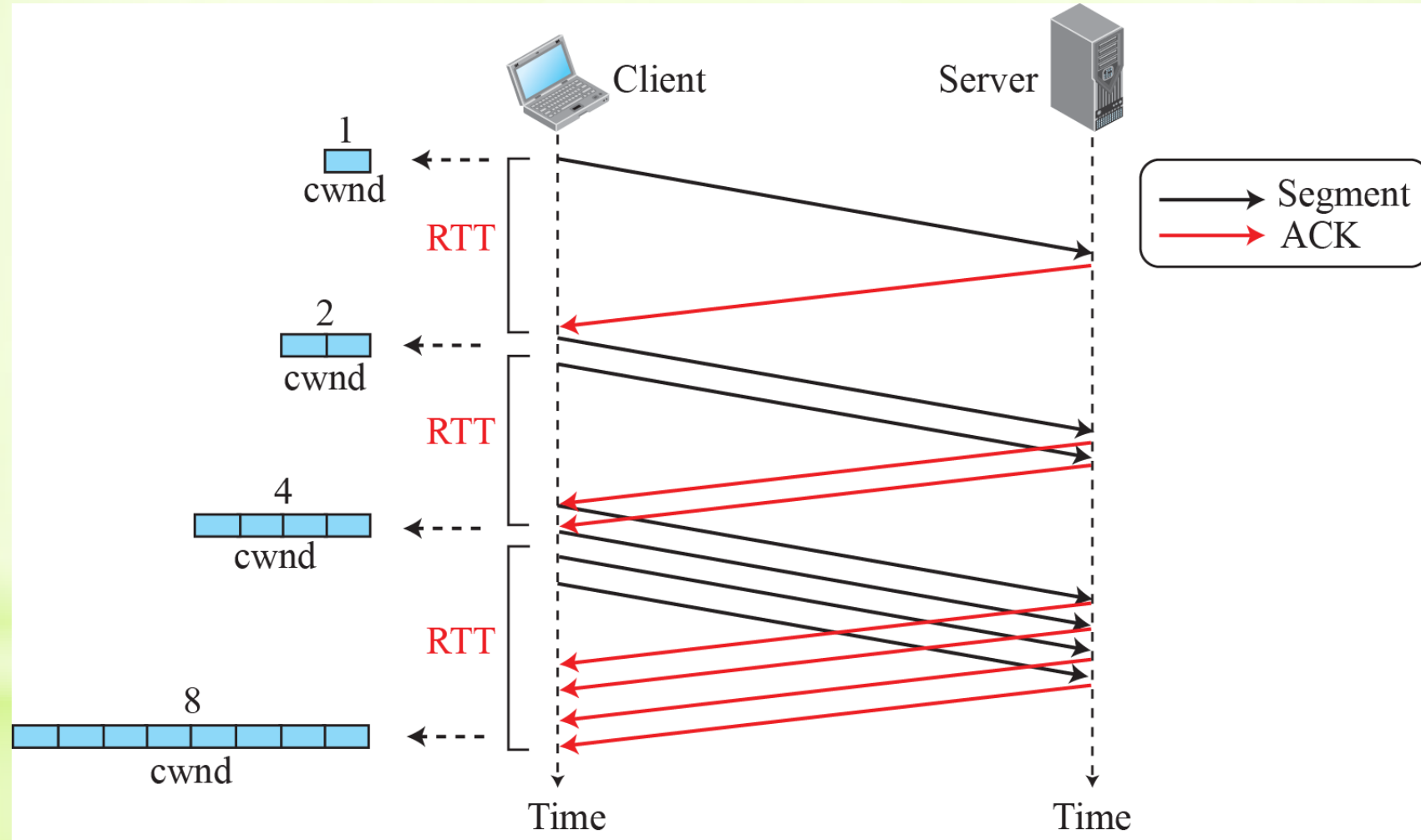
- The MSS is a value negotiated during the connection establishment, using an option of the same name
- Each segment is of the same size and carries MSS bytes

❑ Congestion Policies

❖ Slow Start: Exponential Increase

- The **slow-start algorithm** is based on the idea that the size of the congestion window (*cwnd*) starts with one maximum segment size (MSS), but it increases one MSS each time an acknowledgment arrives. i.e. window size grows exponentially.
- There must be a threshold to stop this phase. The sender keeps track of a variable named *ssthresh* (slow-start threshold). When the size of the window in bytes reaches this threshold, slow start stops and the next phase starts.

Slow start, exponential increase



❑ Congestion Policies

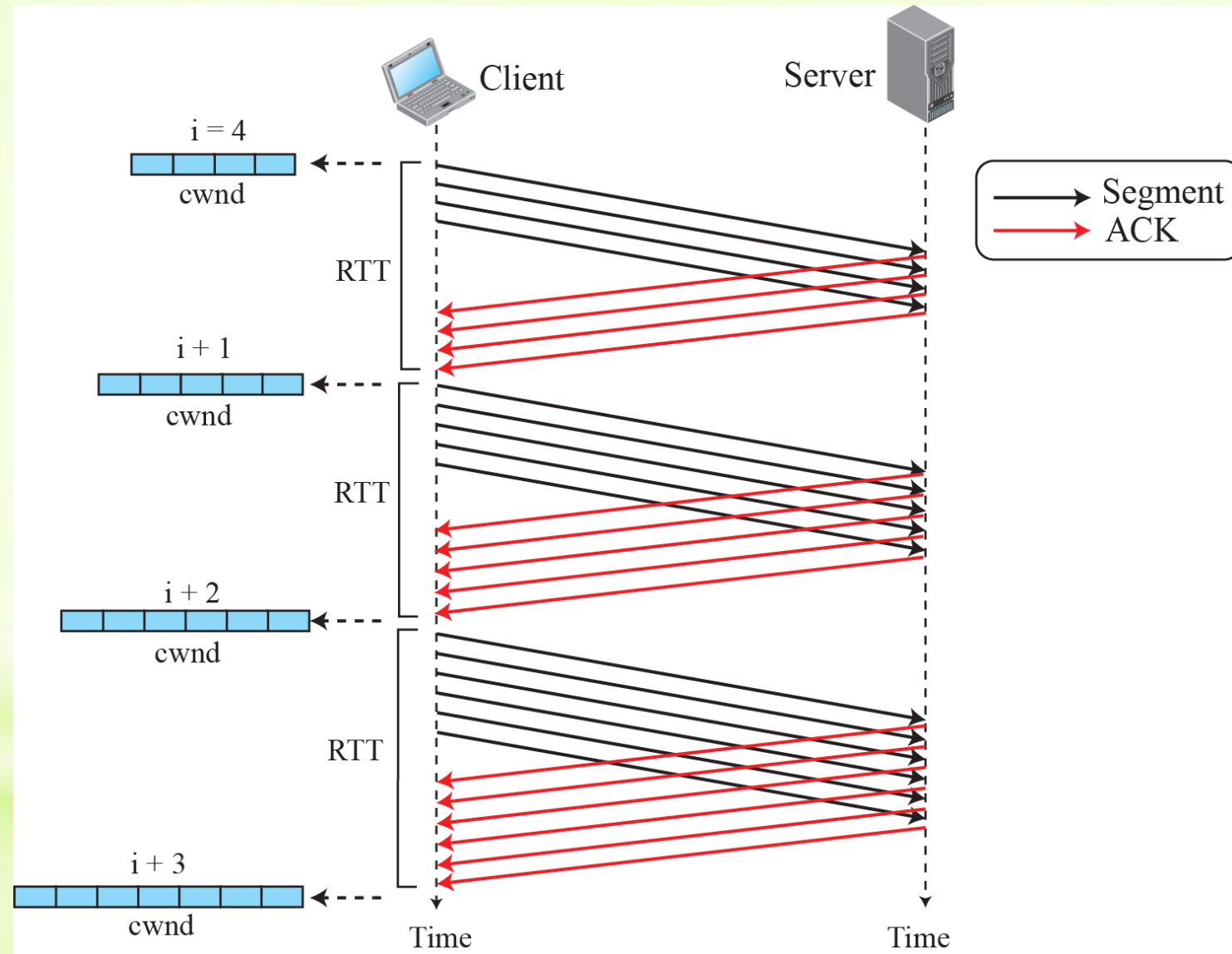
❖ Congestion Avoidance: Additive Increase

- When the size of the congestion window reaches the slow-start threshold, the slow-start phase stops and the additive phase begins. In this algorithm, each time the whole “window” of segments is acknowledged, the size of the congestion window is increased by one.
- **In the congestion-avoidance algorithm, the size of the congestion window increases additively until congestion is detected.**

❖ Fast Retransmission/Fast Recovery

- The **fast-recovery** algorithm is optional in TCP. The old version of TCP did not use it, but the new versions try to use it.

Congestion avoidance, additive increase



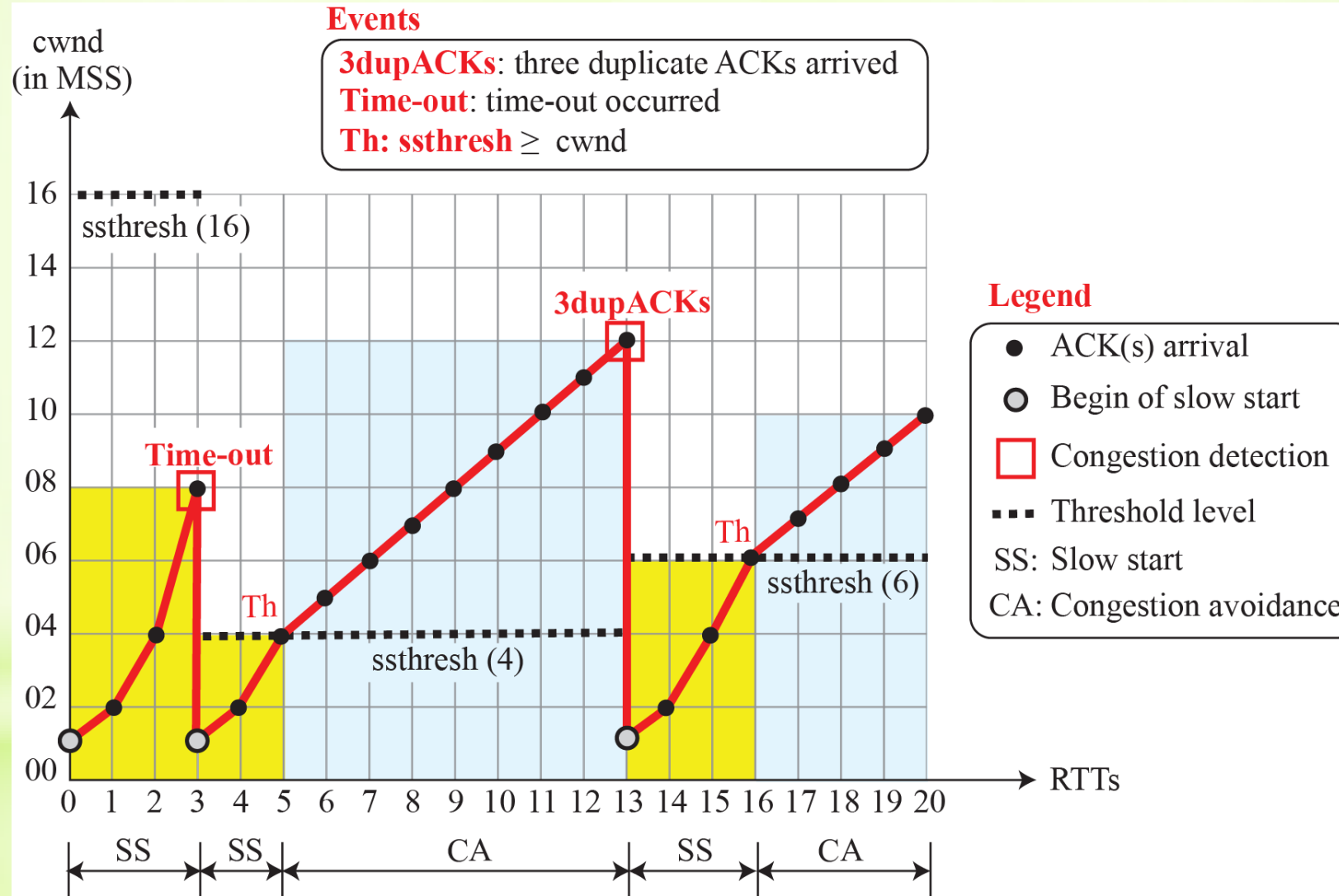
Three versions of TCP with different congestion policies

- ❖ Tahoe TCP
- ❖ Reno TCP
- ❖ NewReno TCP

Taho TCP

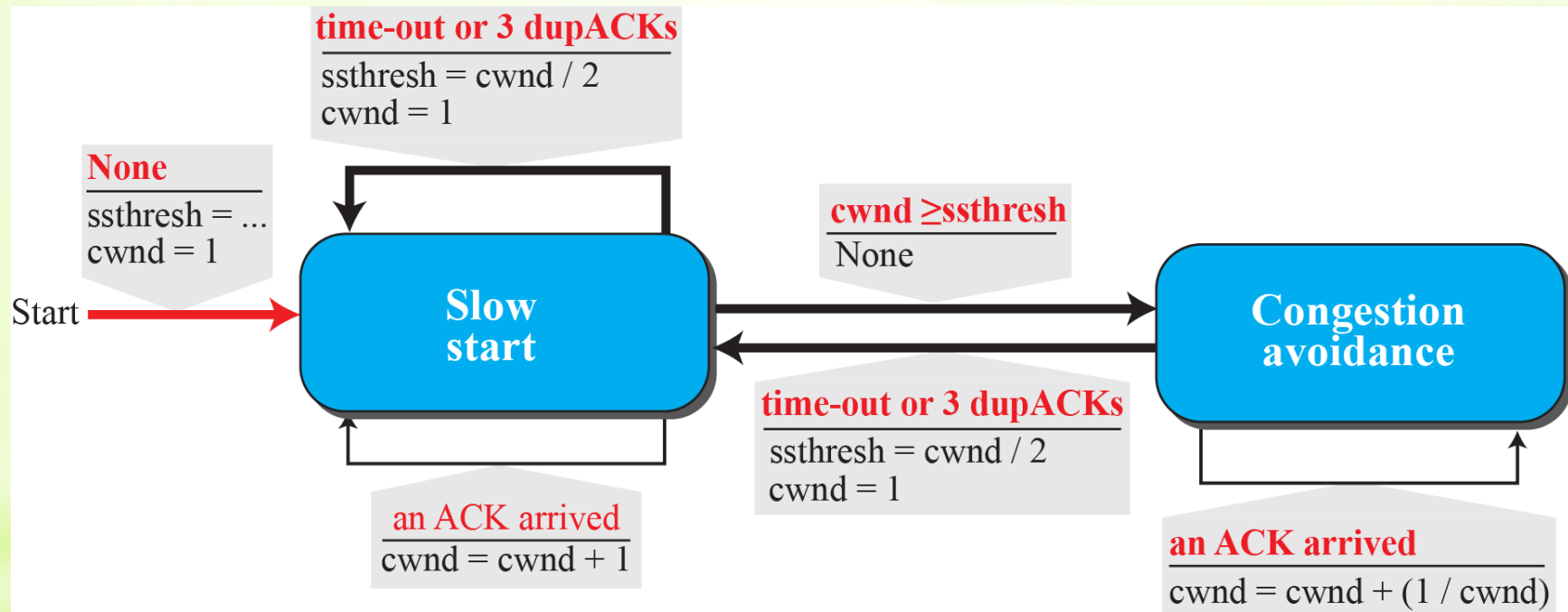
- The early TCP, known as *Taho TCP*, used only two different algorithms in their congestion policy: *slow start* and *congestion avoidance*
- Taho TCP treats the two signs used for congestion detection, time-out and three duplicate ACKs, in the same way.
- When the connection is established, TCP starts the slow-start algorithm and sets the *ssthresh* variable to a pre-agreed value (normally a multiple of MSS) and the *cwnd* to 1 MSS. Then it continues to congestion avoidance phase.
- If congestion is detected (occurrence of time-out or arrival of three duplicate ACKs), TCP immediately interrupts this aggressive growth and restarts a new slow start algorithm by **limiting the threshold to half of the current *cwnd*** and **resetting the congestion window to 1**.

Example of Tahoe TCP



Tahoe TCP

FSM for Tahoe TCP

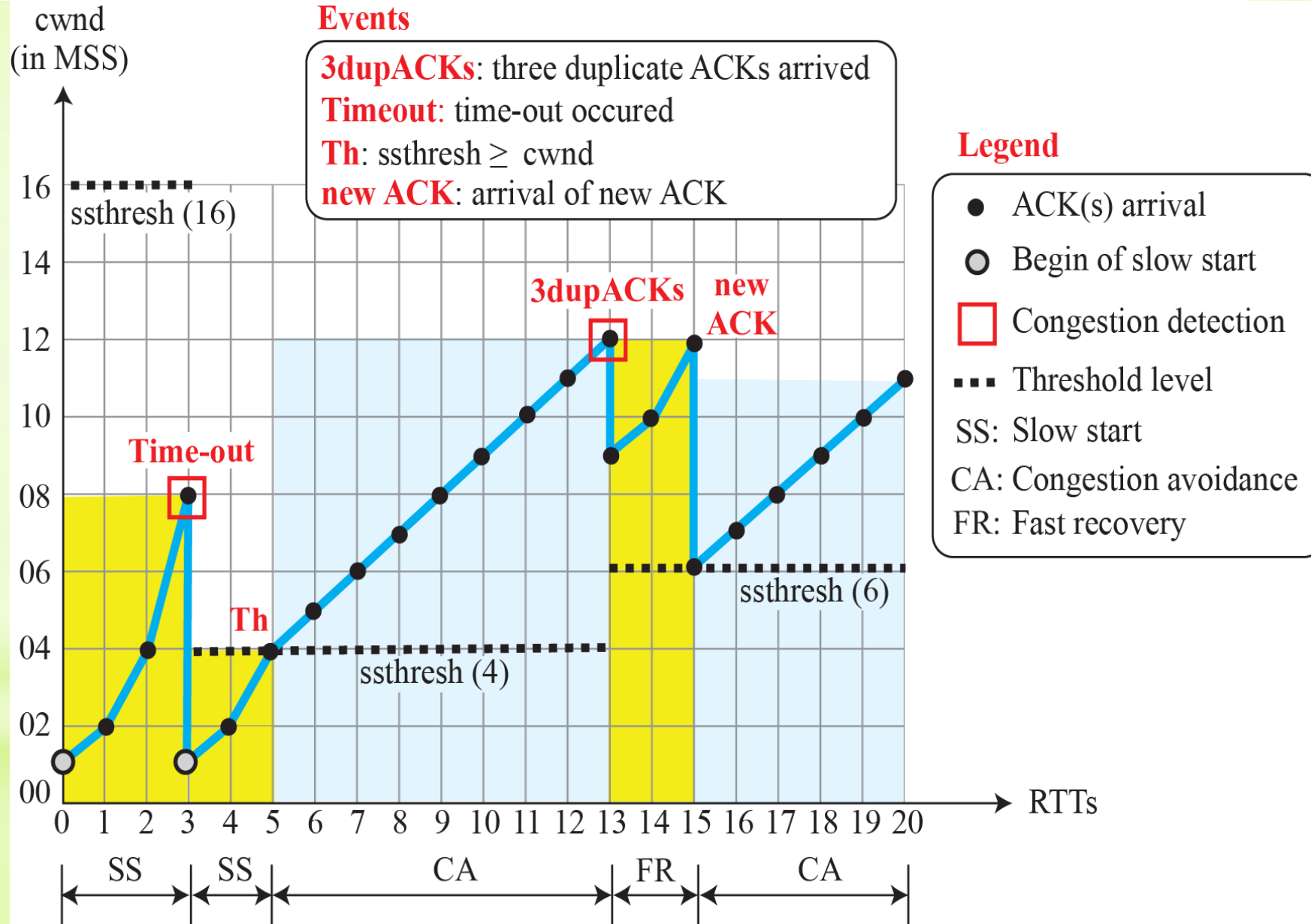


Reno TCP

- A newer version of TCP, called *Reno TCP*, added a new state to the congestion-control FSM, called the fast-recovery state.
- This version treated the two signals of congestion, time-out and the arrival of three duplicate ACKs, differently.
- In this version, **if a time-out occurs, TCP moves to the slow-start state** (or starts a new round if it is already in this state)
- **If three duplicate ACKs arrive, TCP moves to the fast-recovery state and remains there as long as more duplicate ACKs arrive.** The fast-recovery state is a state somewhere between the slow-start and the congestion-avoidance states.

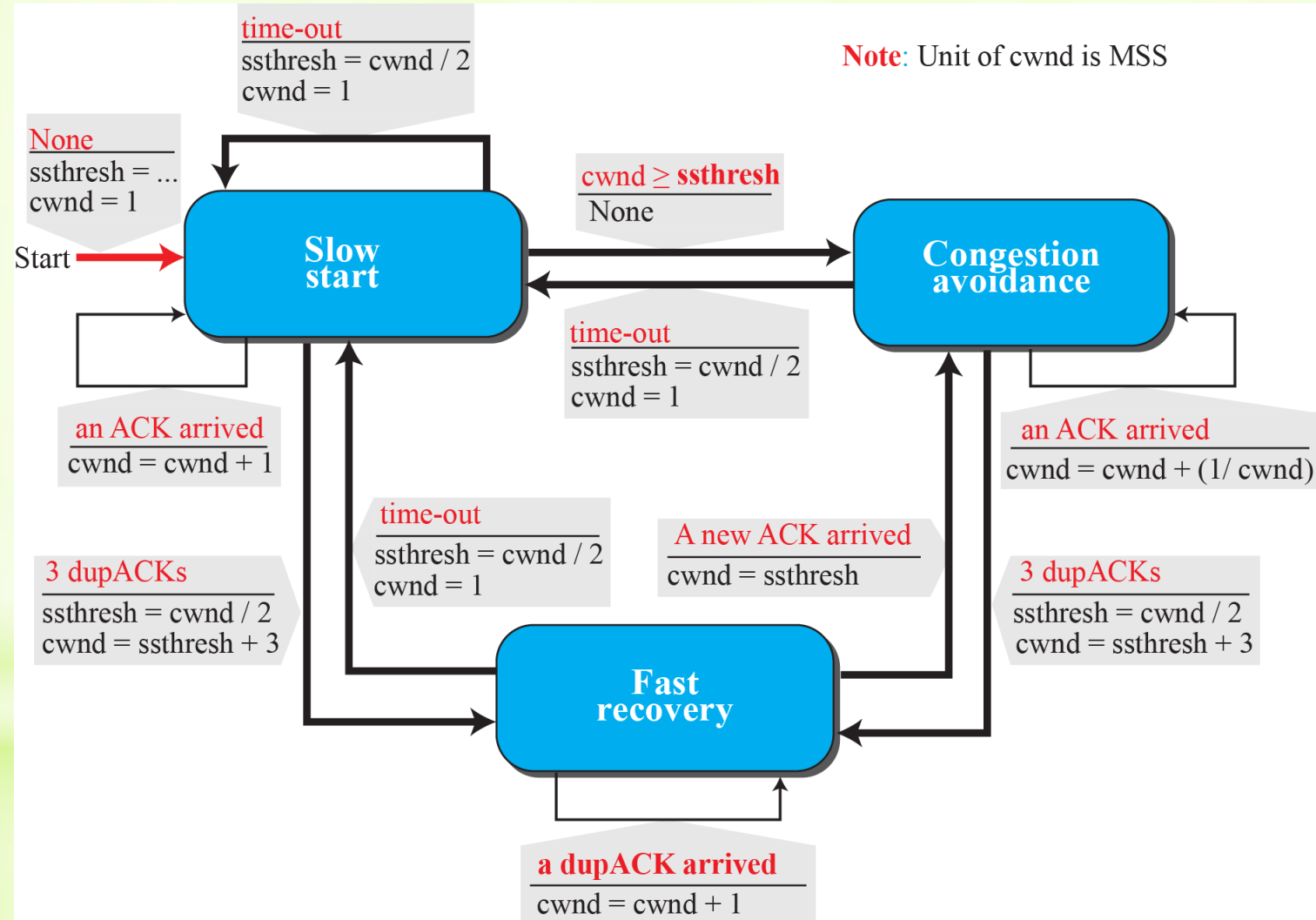
- In the fast-recovery state, it behaves like the slow start, in which the *cwnd* grows exponentially, but the *cwnd* starts with the value of *ssthresh* plus 3 MSS (instead of 1).
- When TCP enters the fast-recovery state, three major events may occur.
 - If duplicate ACKs continue to arrive, TCP stays in this state, but the *cwnd* grows exponentially.
 - If a time-out occurs, TCP assumes that there is real congestion in the network and moves to the slow-start state.
 - If a new (nonduplicate) ACK arrives, TCP moves to the congestion-avoidance state, but deflates the size of the *cwnd* to the *ssthresh* value, as though the three duplicate ACKs have not occurred, and transition is from the slow-start state to the congestion-avoidance state.

Example of a Reno TCP



Reno TCP

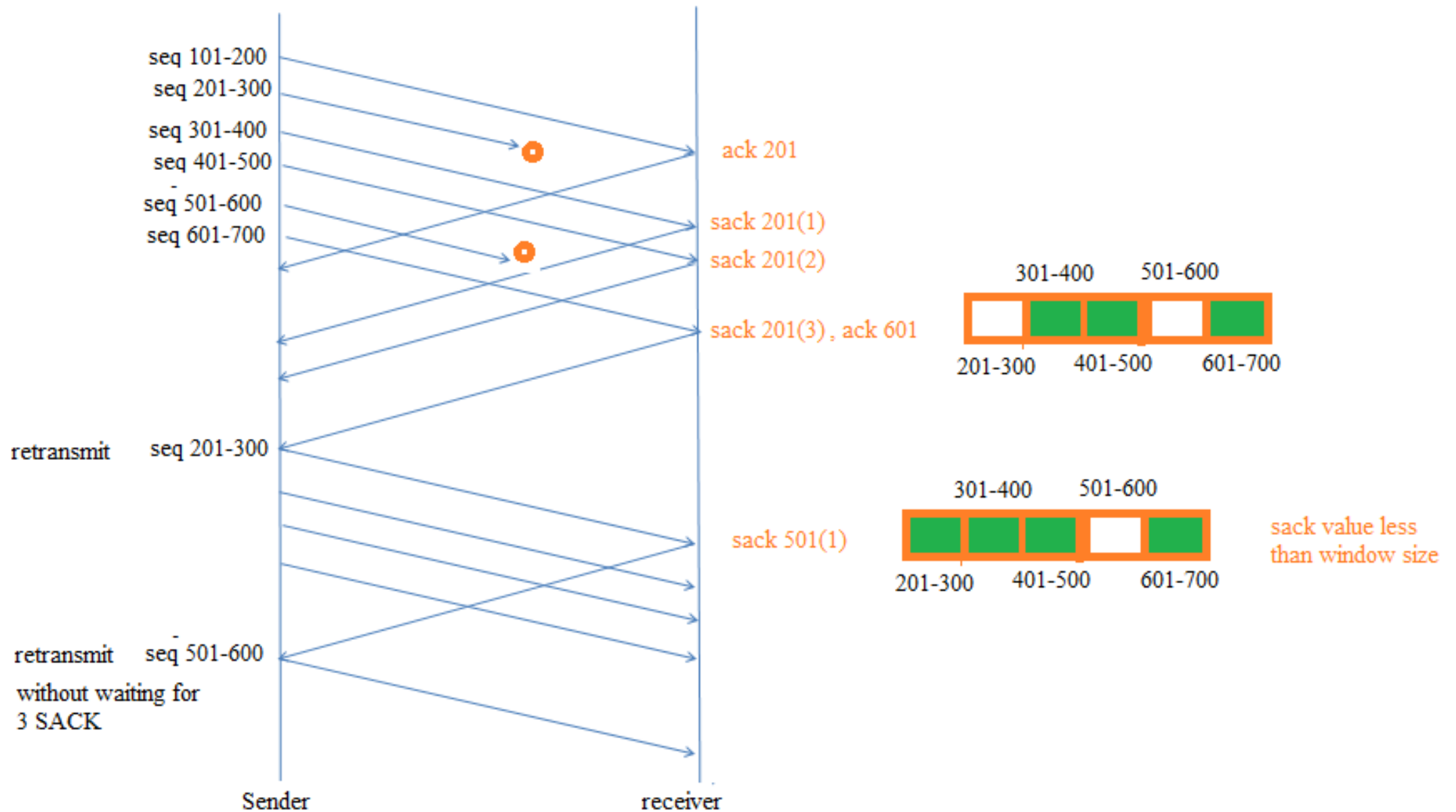
FSM for Reno TCP



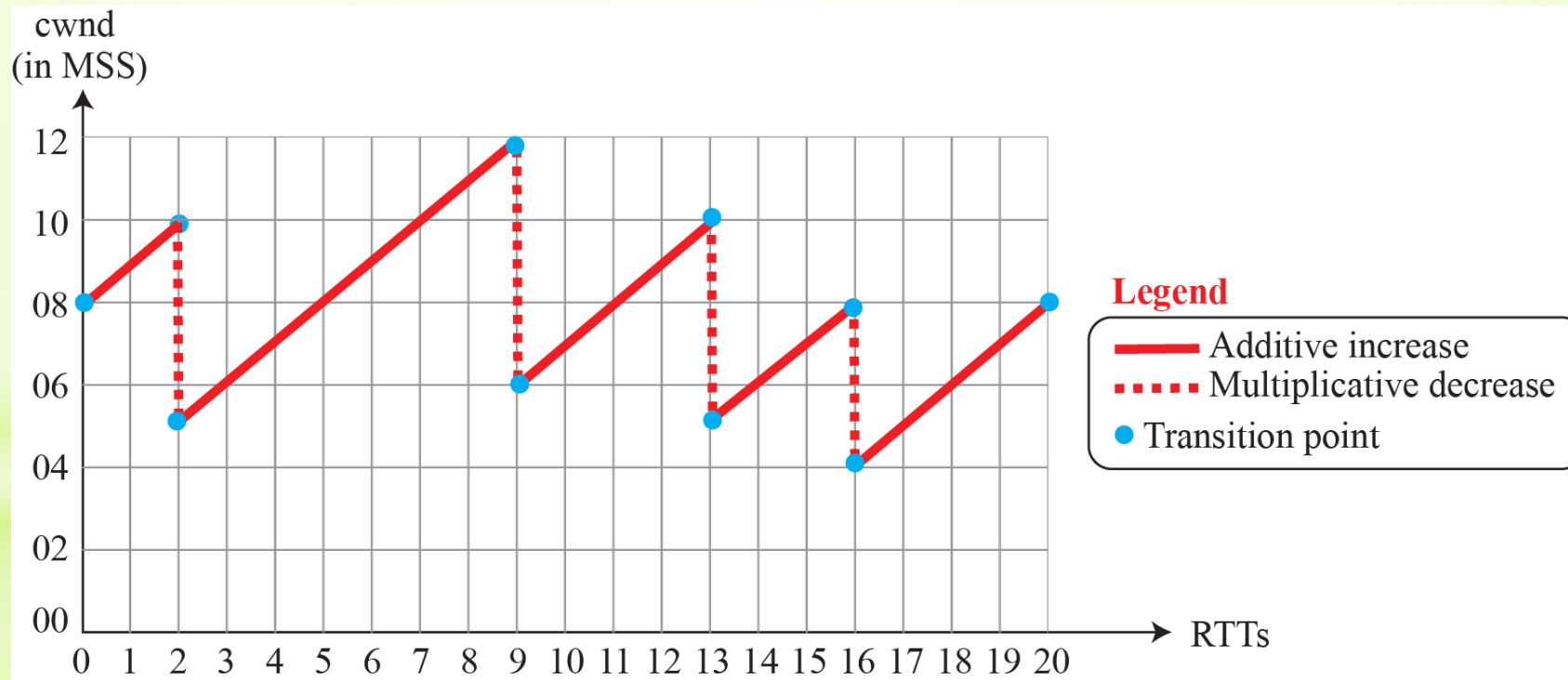
New Reno TCP

- A later version of TCP, called *NewReno TCP*, made an extra optimization on the Reno TCP.
- When TCP receives three duplicate ACKs, it retransmits the lost segment until a new ACK (not duplicate) arrives. (Fast Recovery Sate)
- If the new ACK defines the end of the window when the congestion was detected, TCP is certain that only one segment was lost. However, if the ACK number defines a position between the retransmitted segment and the end of the window, it is possible that the segment defined by the ACK is also lost.
- NewReno TCP retransmits this segment to avoid receiving more and more duplicate ACKs for it.

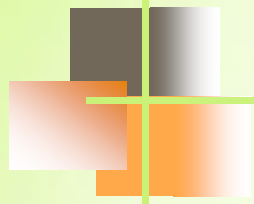
New Reno TCP



- Out of the three versions of TCP, the Reno version is most common today.
- If we ignore the slow-start states and short exponential growth during fast recovery, the TCP congestion window is $cwnd = cwnd + (1 / cwnd)$ when an ACK arrives (congestion avoidance), and $cwnd = cwnd / 2$ when congestion is detected. i.e. Additive increase, multiplicative decrease (AIMD)



Additive increase, multiplicative decrease (AIMD)



TCP Timers

To perform their operations smoothly, most TCP implementations use at least four timers.

- ❑ Retransmission Timer

- ❖ Round-Trip Time (RTT)
- ❖ Karn's Algorithm
- ❖ Exponential Backoff

- ❑ Persistence Timer

- ❑ Keepalive Timer

- ❑ TIME-WAIT Timer

Retransmission Timer

Round-Trip Time (RTT)

Measured RTT (RTT_M)

- How long it takes to send a segment and receive an acknowledgment for it. This is the measured RTT_M .
- In TCP the segments and their acknowledgments do not have a one-to-one relationship, several segments may be acknowledged together.
- In TCP, there can be only one RTT_M measurement in progress at any time.

Smoothed RTT (RTT_S)

- Most implementations use a smoothed RTT, called RTT_S , which is a weighted average of RTT_M and the previous RTT_S

Initially

→

No value

After first measurement

→

$RTT_S = RTT_M$

After each measurement

→

$RTT_S = (1 - \alpha) RTT_S + \alpha \times RTT_M$

α is implementation-dependent, but it is normally set to 1/8

❑ Retransmission Timer

Deviated RTT

Most implementations do not just use RTT_S ; they also calculate the RTT deviation, called RTT_D , based on the RTT_S and RTT_M

Initially	→	No value
After first measurement	→	$RTT_D = RTT_M / 2$
After each measurement	→	$RTT_D = (1 - \beta) RTT_D + \beta \times RTT_S - RTT_M $

β is also implementation-dependent, but is usually set to 1/4.

Retransmission Time Out (RTO)

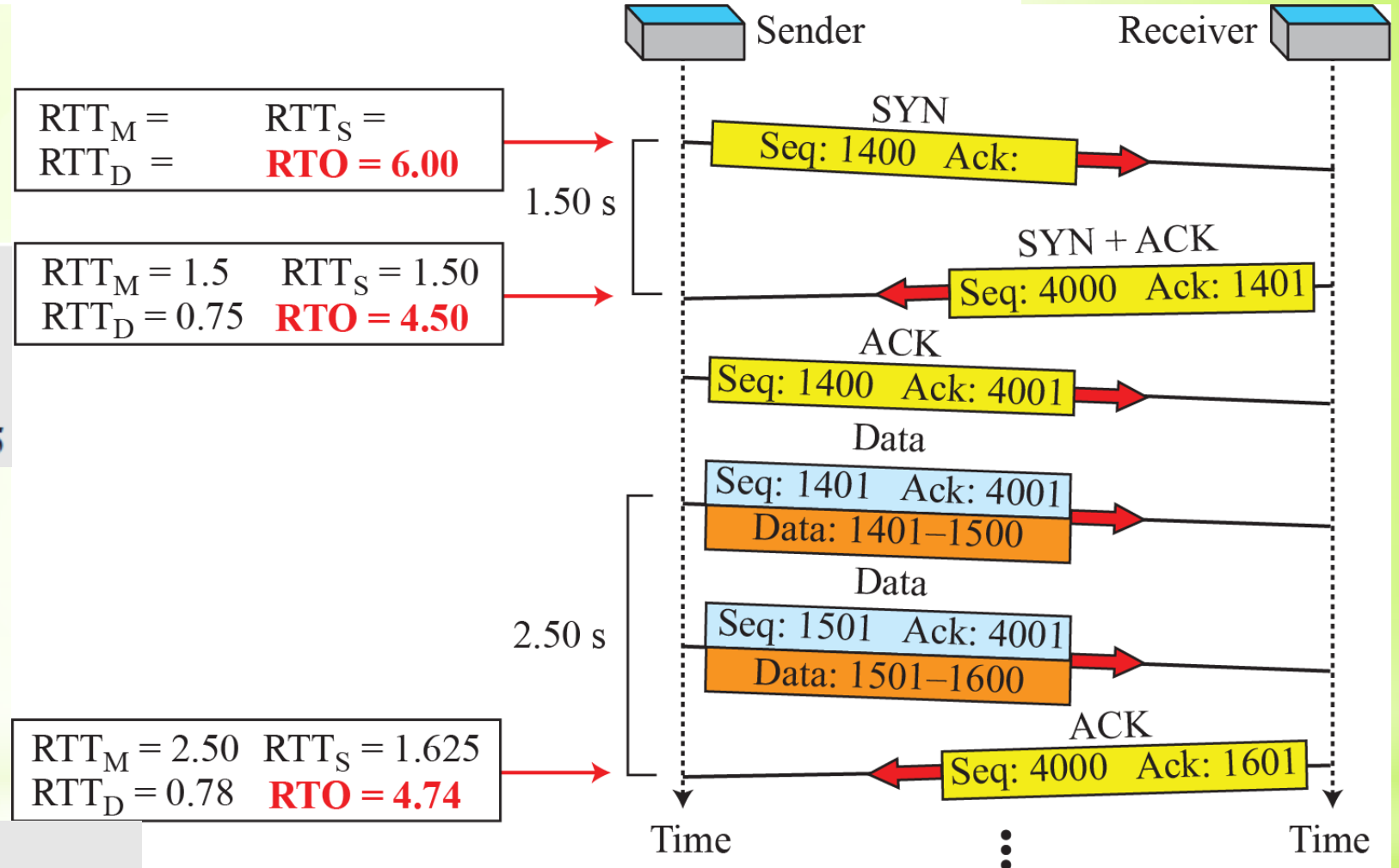
The value of RTO is based on the smoothed round trip time and its deviation.

Original	→	Initial value
After any measurement	→	$RTO = RTT_S + 4 \times RTT_D$

Example

$$\begin{aligned} \text{RTT}_M &= 1.5 \\ \text{RTT}_S &= 1.5 \\ \text{RTT}_D &= (1.5)/2 = 0.75 \\ \text{RTO} &= 1.5 + 4 \times 0.75 = 4.5 \end{aligned}$$

$$\begin{aligned} \text{RTT}_M &= 2.5 \\ \text{RTT}_S &= (7/8) \times (1.5) + (1/8) \times (2.5) = 1.625 \\ \text{RTT}_D &= (3/4) \times (0.75) + (1/4) \times |1.625 - 2.5| = 0.78 \\ \text{RTO} &= 1.625 + 4 \times (0.78) = 4.74 \end{aligned}$$



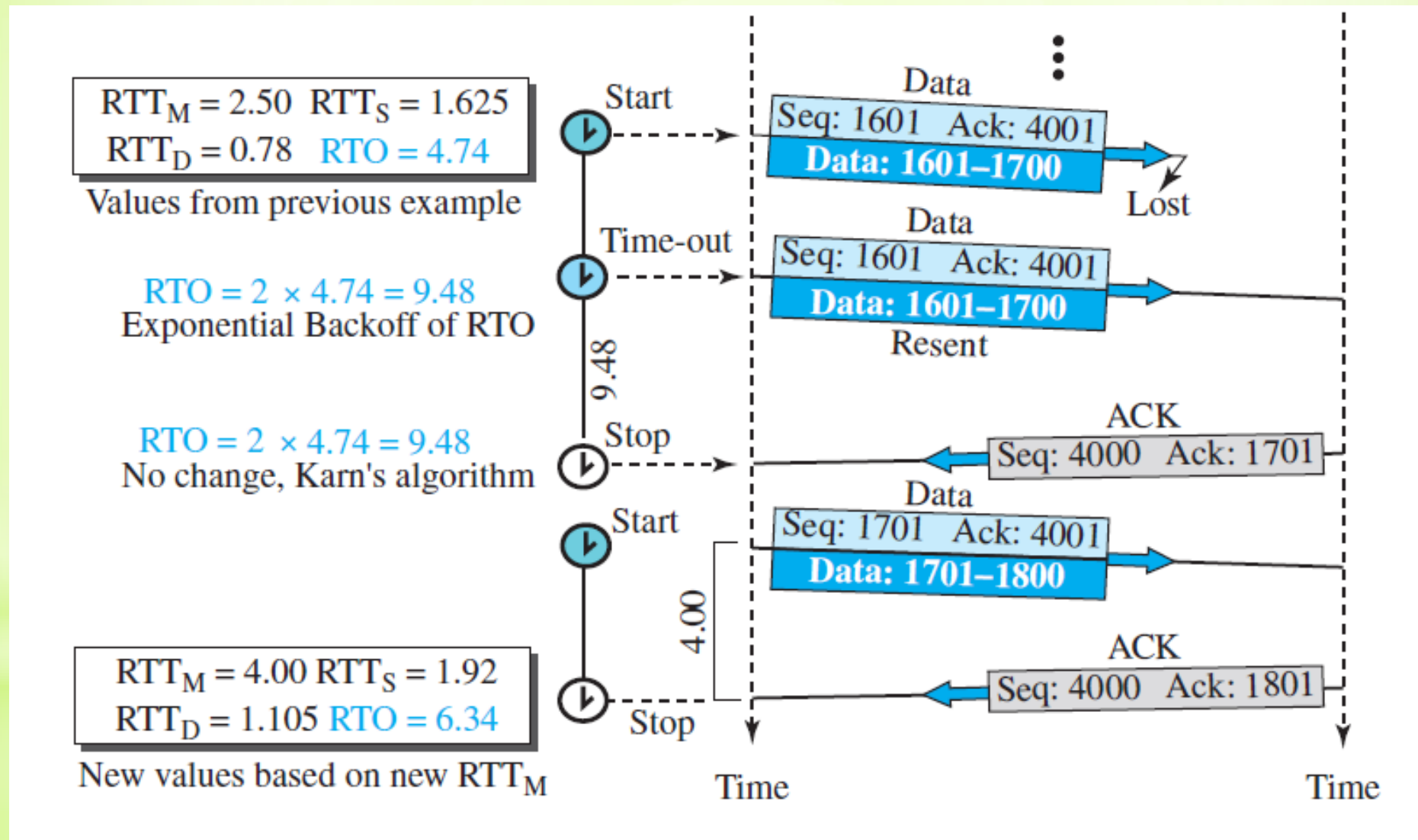
❖ Karn's Algorithm

- Suppose that a segment is not acknowledged during the retransmission time-out period and is therefore retransmitted. When the sending TCP receives an acknowledgment for this segment, it does not know if the acknowledgment is for the original segment or for the retransmitted one.
- Karn's algorithm is simple: **TCP does not consider the RTT of a retransmitted segment in its calculation of a new RTO.**

❖ Exponential Backoff

- Most TCP implementations use an exponential backoff strategy. The value of RTO is doubled for each retransmission. So if the segment is retransmitted once, the value is two times the RTO. If it transmitted twice, the value is four times the RTO, and so on.

Retransmission and Karn's algorithm is applied.



❑ Persistence Timer

- To deal with a zero-window-size advertisement, TCP needs another timer. If the receiving TCP announces a window size of zero, the sending TCP stops transmitting segments until the receiving TCP sends an ACK segment announcing a nonzero window size. This ACK segment can be lost. Both TCP's might continue to wait for each other forever (a deadlock).
- To correct this deadlock, TCP uses a persistence timer for each connection. When the sending TCP receives an acknowledgment with a window size of zero, it starts a persistence timer.
- When the persistence timer goes off, the sending TCP sends a special segment called a **probe**.
- This segment contains only 1 byte of new data. It has a sequence number, but its sequence number is never acknowledged.
- The value of the persistence timer is set to the value of the retransmission time. However, if a response is not received from the receiver, another probe segment is sent and the value of the persistence timer is doubled and reset until the value reaches a threshold (usually 60 s). After that the sender sends one probe segment every 60 seconds until the window is reopened.

❑ Keepalive Timer

- A keep alive timer is used in some implementations to prevent a long idle connection between two TCP's.
- The time-out is usually 2 hours. If the server does not hear from the client after 2 hours, it sends a probe segment. If there is no response after 10 probes, each of which is 75 seconds apart, it assumes that the client is down and terminates the connection.

❑ TIME-WAIT Timer

- The TIME-WAIT (2MSL) timer is used during connection termination. The 2MSL timer is used when TCP performs an active close and sends the final ACK. The connection must stay open for 2 MSL amount of time to allow TCP to resend the final ACK in case the ACK is lost.
- Common values are 30 seconds, 1 minute, or even 2 minutes.