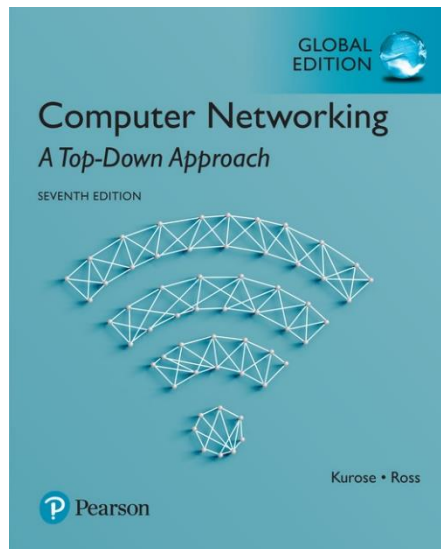


## Dr. Saad Aslam

Office: Room AE-327, 3<sup>rd</sup> Floor, NUB

Contact: +603 74918622  
(Ext: 7143)

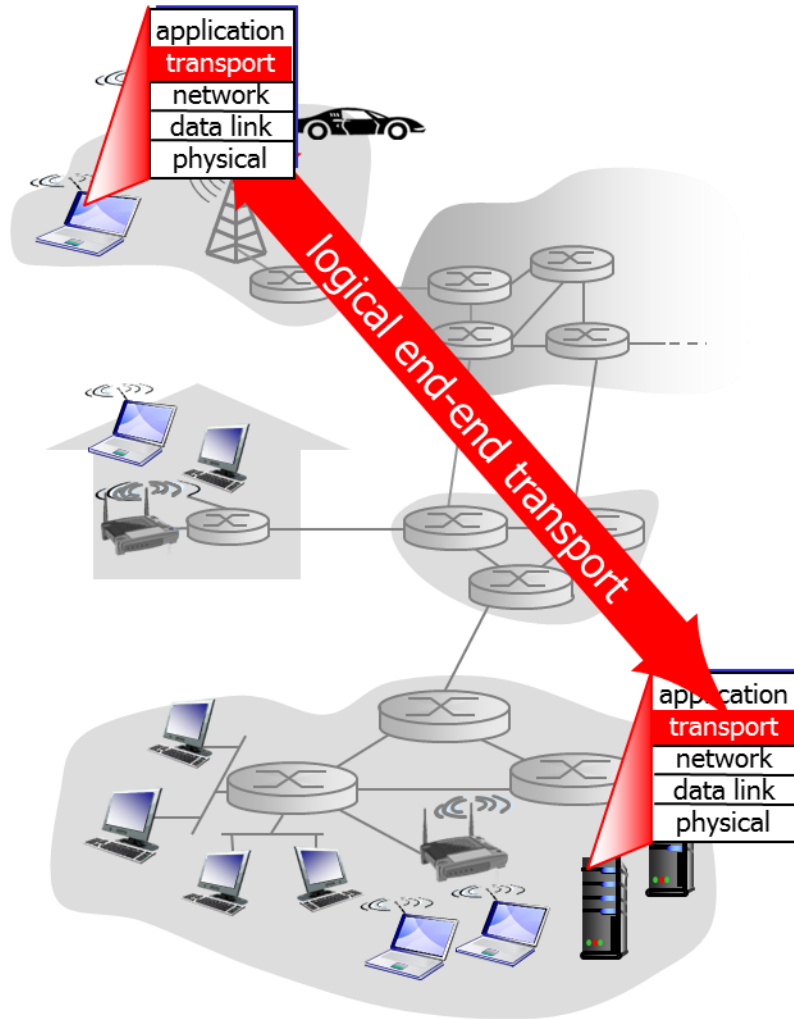
Email: [saada@sunway.edu.my](mailto:saada@sunway.edu.my)



*Computer Networking: A Top-down Approach, 7<sup>th</sup> edition.*  
Jim Kurose, Keith Ross  
Pearson

Section	Topic	Slides	
3.1	<b>Introduction and Transport-Layer Services</b>	2-3	
3.1 3.3 3.5	<ul style="list-style-type: none"><li>• Overview of the Transport Layer in the Internet</li></ul>		
3.4	<b>Principles of Reliable Data Transfer</b>		
	<ul style="list-style-type: none"><li>• Pipelined Reliable Data Transfer Protocols</li></ul>		4-9
	<ul style="list-style-type: none"><li>• Go-Back-N (GBN)</li></ul>	10-12	
	<ul style="list-style-type: none"><li>• Selective Repeat (SR)</li></ul>	13-15	
3.5	<b>Connection-Oriented Transport: TCP</b>		
	<ul style="list-style-type: none"><li>• Reliable Data Transfer</li></ul>	16-18	
	<ul style="list-style-type: none"><li>• Flow Control</li></ul>	19	
3.7	<b>TCP Congestion Control</b>	20-23	

## 3.1 Introduction and Transport-Layer Services



### • Transport Layer

- Provide logical communication between processes running on different hosts
  - Sender host: Break a message into segments
  - Destination host: Reassemble segments into the message, and pass to application layer
- Rely on network layer services
  - *network layer*
    - Provide logical communication between hosts
- Two type of protocol
  - **Transmission Control Protocol (TCP)**
  - **User Datagram Protocol (UDP)**

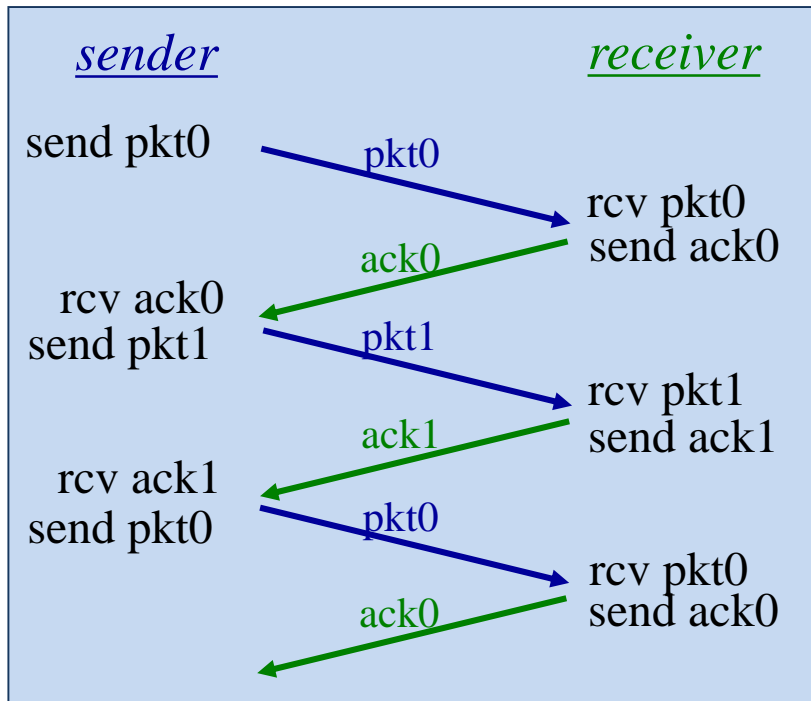
# 3.1, 3.3, 3.5 Introduction and Transport-Layer Services : A Summary

## Comparison between TCP and UDP

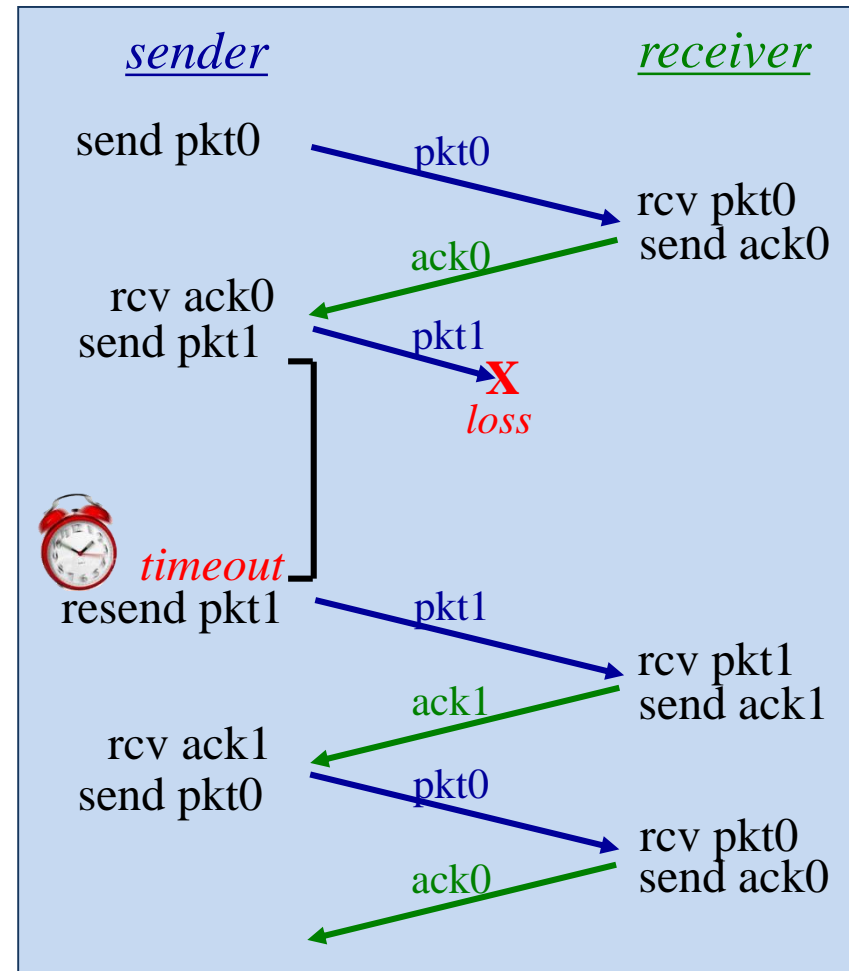
TCP	UDP
<ul style="list-style-type: none"> <li>• <b>Reliable</b> <ul style="list-style-type: none"> <li>• Destination host receive all data in-order</li> <li>• Use Go-Back-N</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• <b>Unreliable</b> <ul style="list-style-type: none"> <li>• Destination host may not receive some lost data and out-of-order</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>• <b>Pipelined</b> <ul style="list-style-type: none"> <li>• <b>Congestion control</b> <ul style="list-style-type: none"> <li>• Resend segment until successful transmission</li> <li>• Slower: Retransmission incurs delay</li> <li>• Control a source host's sending rate so that each connection traversing a congested link get an equal share of link bandwidth</li> </ul> </li> <li>• <b>Flow control</b> <ul style="list-style-type: none"> <li>• Control a source host's sending rate so that it does not overflow receive host's buffer</li> </ul> </li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• <b>No pipelined</b> <ul style="list-style-type: none"> <li>• <b>No congestion control</b> <ul style="list-style-type: none"> <li>• Do not resend segment until successful transmission</li> <li>• Faster</li> <li>• Do not control a source host's sending rate</li> </ul> </li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>• <b>Connection-oriented</b> <ul style="list-style-type: none"> <li>• Establish connection, then start transmission <ul style="list-style-type: none"> <li>• Establishing connection require handshaking (exchange of control messages to initialize sender and receiver state)</li> <li>• Each connection is full duplex (bi-directional data flow)</li> <li>• Slower: Establishing connection incurs delay</li> </ul> </li> <li>• <b>Point-to-point</b> <ul style="list-style-type: none"> <li>• One sender, one receiver</li> </ul> </li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• <b>Connectionless</b> <ul style="list-style-type: none"> <li>• Do not establish connection <ul style="list-style-type: none"> <li>• Faster</li> </ul> </li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>• <b>Stateful</b> <ul style="list-style-type: none"> <li>• Host maintain buffer utilization, congestion control parameter, acknowledgment number</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• <b>Stateless</b> <ul style="list-style-type: none"> <li>• Host do not maintain buffer utilization, congestion control parameter, acknowledgment number</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>• <b>Packet header</b> <ul style="list-style-type: none"> <li>• <b>Larger</b></li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• <b>Packet header</b> <ul style="list-style-type: none"> <li>• <b>Smaller</b></li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>• <b>Example of usage</b> <ul style="list-style-type: none"> <li>• Email</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• <b>Example of usage</b> <ul style="list-style-type: none"> <li>• Audio and video traffic</li> </ul> </li> </ul>

## 3.4 Reliable Data Transfer Protocol

- Consider FOUR scenarios
  - No packet loss
  - Packet loss
  - ACK loss
  - Timeout / Delayed ACK



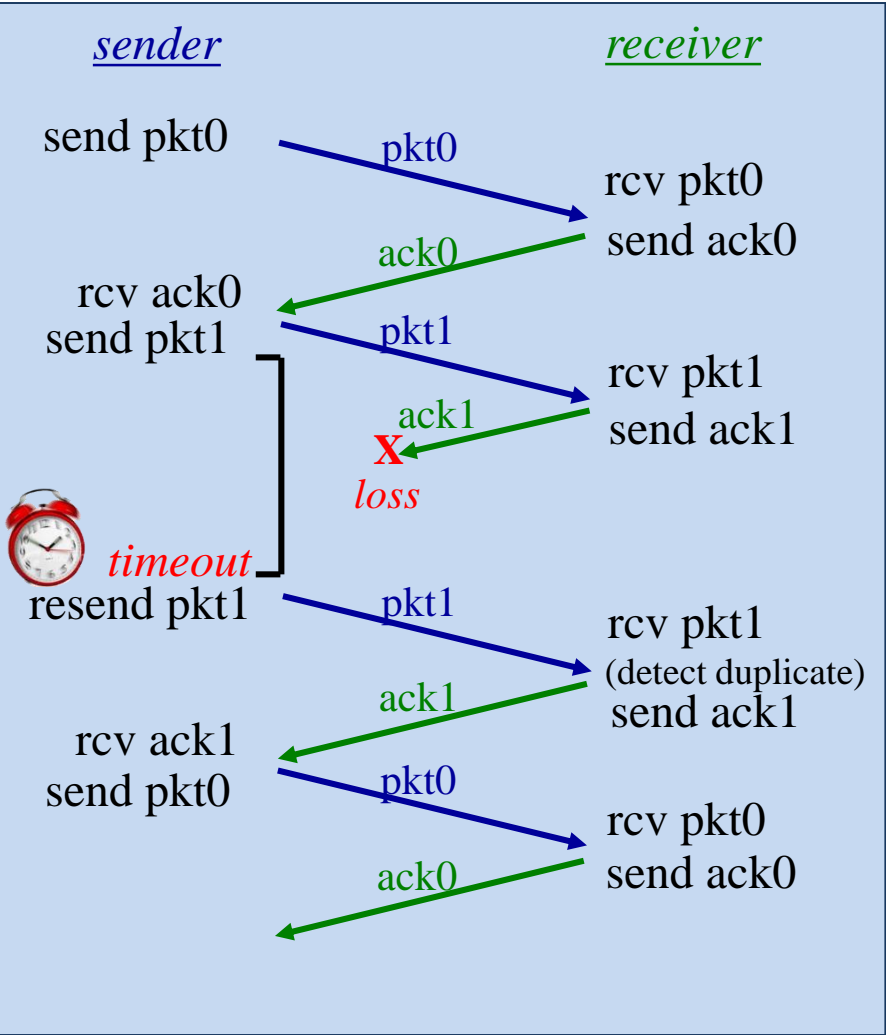
(a) No packet loss



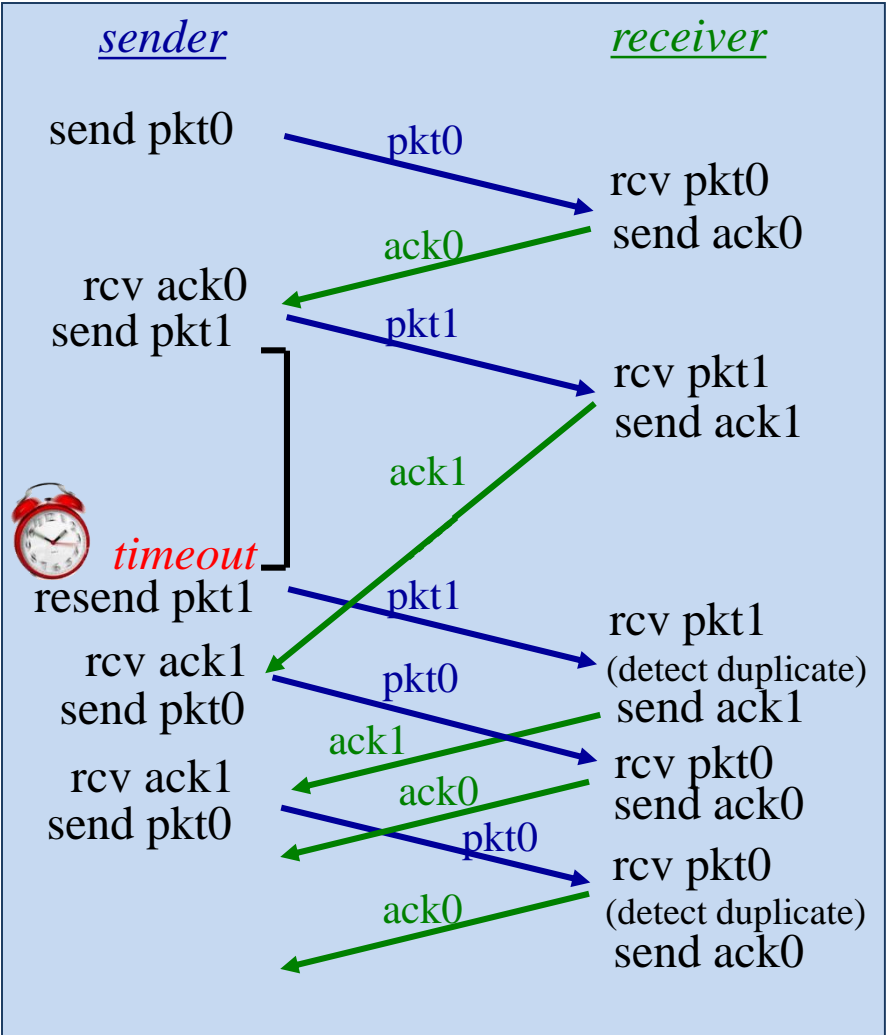
(b) Packet loss

Note: Timer is restarted when sender receive ack

# 3.4 Reliable Data Transfer Protocol

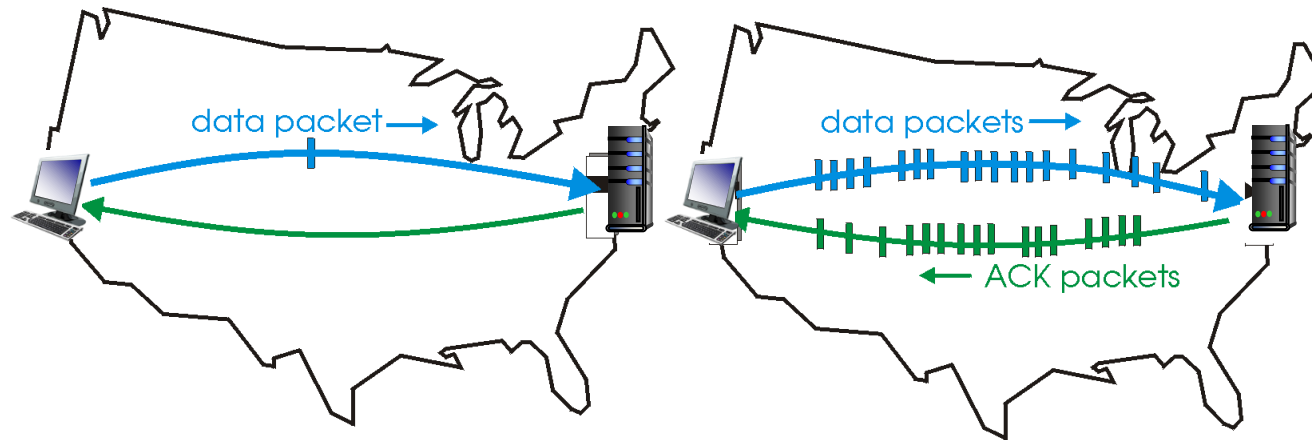


(c) ACK loss



(d) Timeout / Delayed ACK

### 3.4 Reliable Data Transfer Protocol: Comparison between stop-and-wait and pipelined protocols



(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

- **Stop-and-wait protocol**

- Sender send a single packet, receive acknowledgment; then only send the next packet

- **Pipelined protocol**

- Sender send multiple unacknowledged packets
  - Must increase the range of sequence number
  - Must have sufficient buffer size at sender and receiver
- Two types
  - **Go-Back-N (GBN)**
  - **Selective Repeat (SR)**

### 3.4 Reliable Data Transfer Protocol: Comparison between stop-and-wait and pipelined protocols

- **Stop-and-wait protocol**

- Low utilization
- Example:
  - Packet size,  $L$ : 8000 bit
  - Transmission rate,  $R$ : 1 Gbps
  - Propagation delay: 15 ms
  - Round trip time, RTT: 30 ms
- What is the utilization?

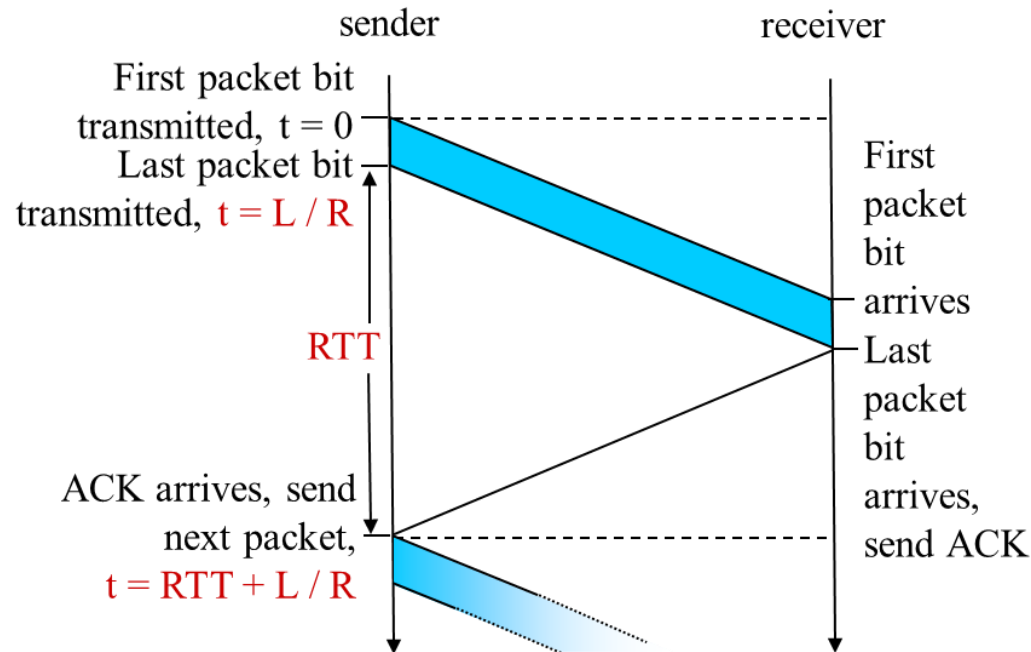
- Transmission delay
  - Time to transmit packet into the 1 Gbps link

$$d_{trans} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bits/sec}} = 8\mu s$$

- Utilization

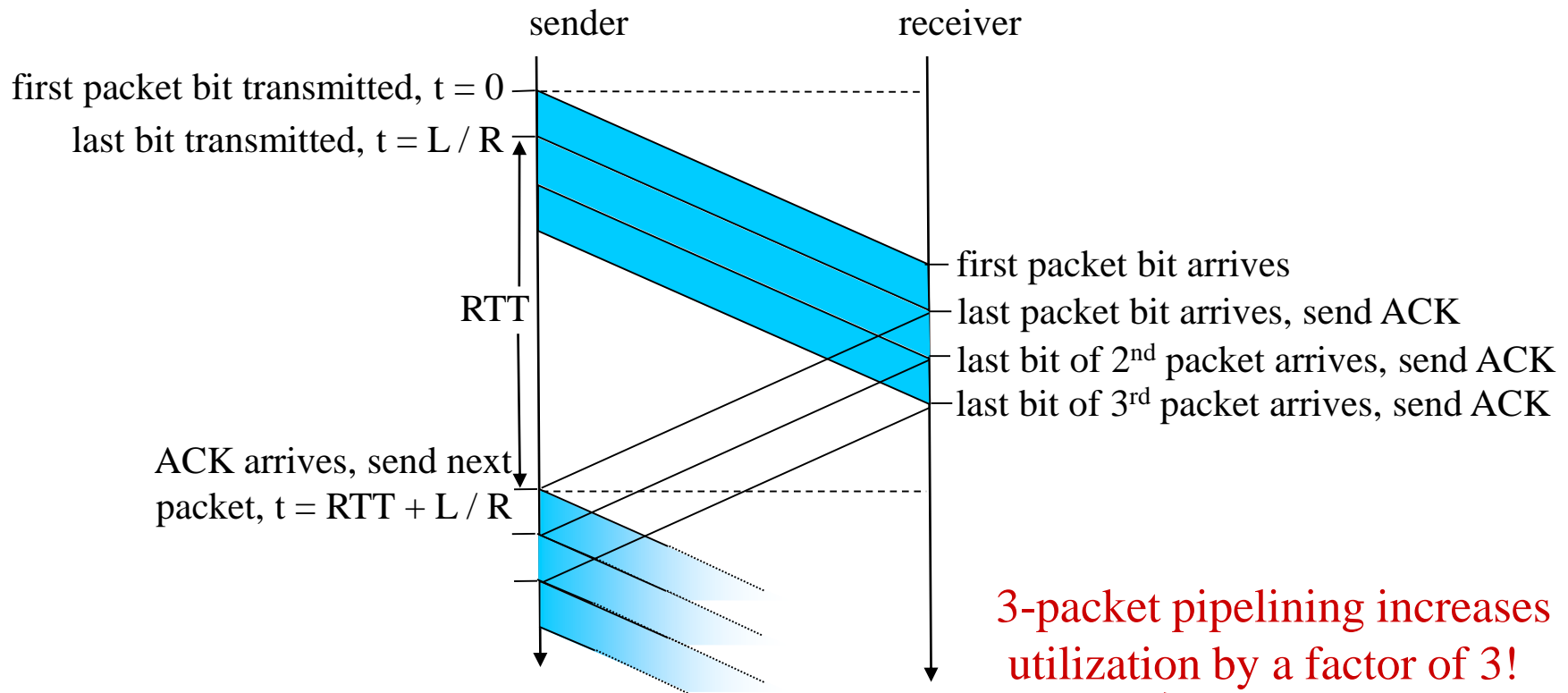
$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{0.008}{30.008} = 0.00027$$

- So, the effective throughput is 270 kbps only for a 1 Gbps link!



### 3.4 Reliable Data Transfer Protocol: Comparison between stop-and-wait and pipelined protocols

- **Pipelined protocol**
  - Higher utilization



$$U_{sender} = \frac{3L/R}{RTT+L/R} = \frac{0.0024}{30.008} = 0.00081$$



## 3.4 Reliable Data Transfer Protocol: Pipelined protocols

### • Go-back-N

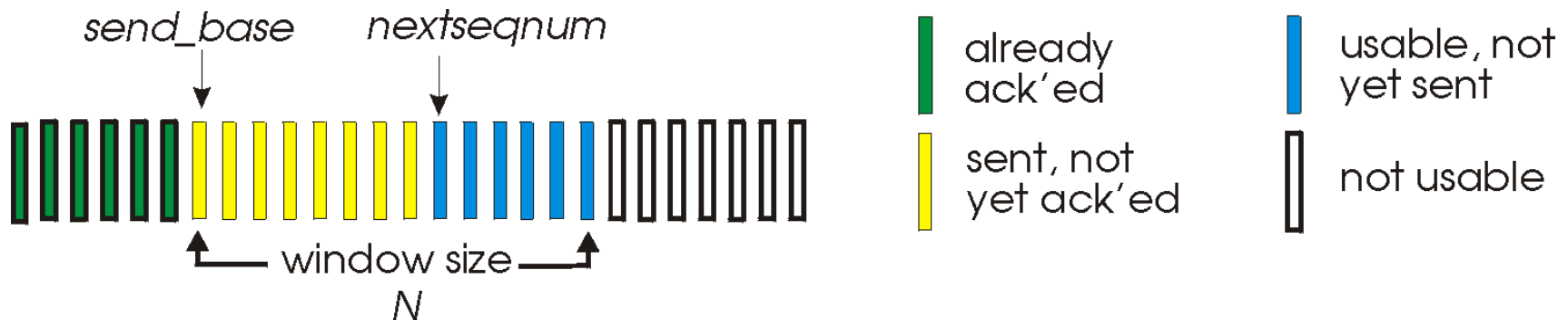
- Sender
  - Can send  $N$  unacknowledged packets in pipeline
- Receiver
  - Only send cumulative acknowledgment
    - This means an acknowledgment for packet  $n$  indicates all packets with sequence number up to and including  $n$  have been received correctly
  - No buffer
- Timeout
  - Sender has a single timer for the oldest unacknowledged packet
  - When timer expire, retransmit all sent but unacknowledged packets
    - If window size is large
      - There can be many unacknowledged packets in pipeline, so a single packet loss can cause sender to retransmit many packets

### • Selective Repeat

- Sender
  - Can send  $N$  unacknowledged packets in pipeline
- Receiver
  - Send acknowledgment for each packet
  - Has buffer
- Timeout
  - Sender has timer for each unacknowledged packet, so it has many timers
  - When timer expire, retransmit the unacknowledged packet only
    - Address the disadvantage of Go-back-N

## 3.4 Go-Back-N (GBN)

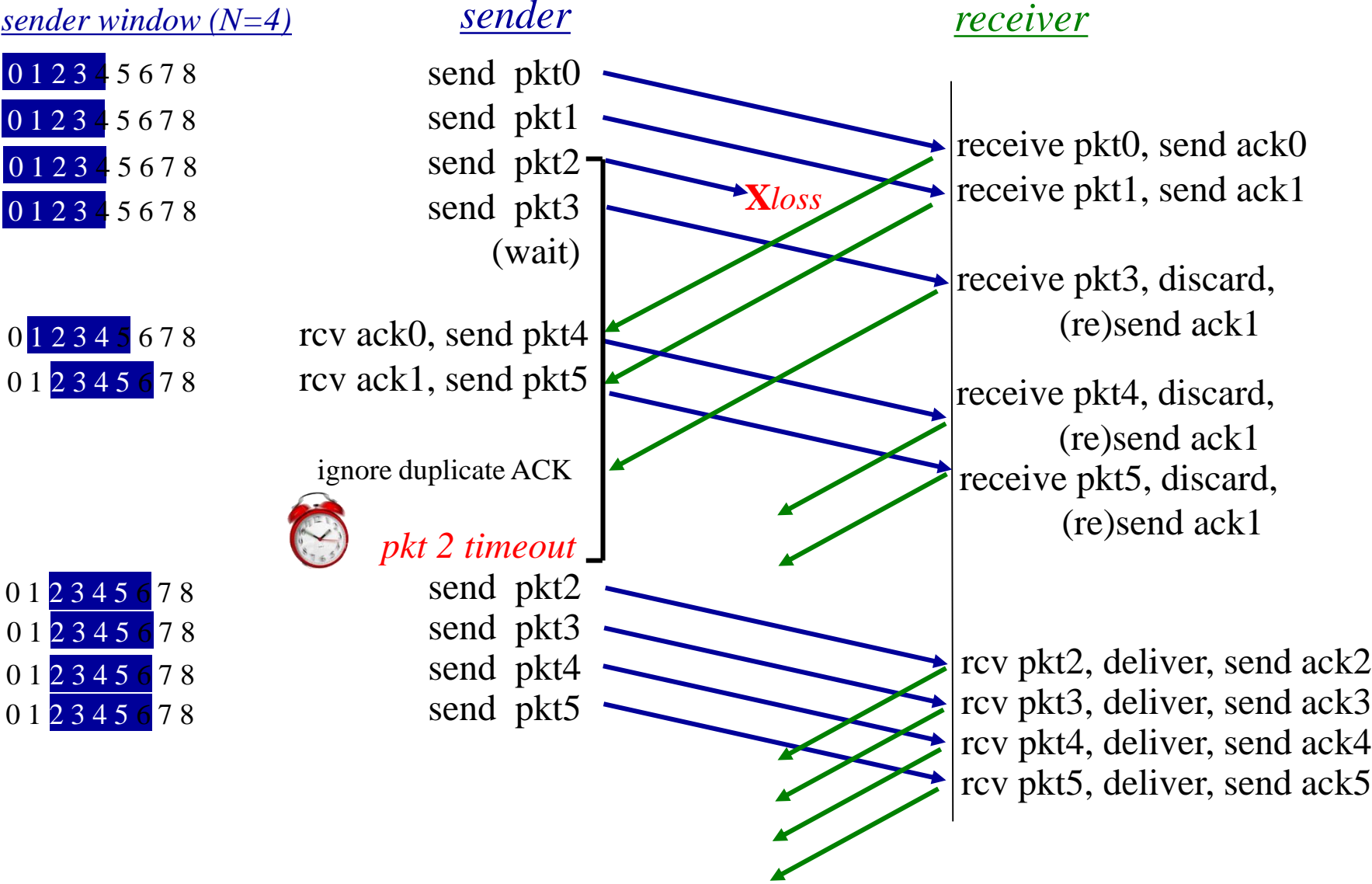
- Go-back-N
  - Maintain a window
    - Window size =  $N$ 
      - A sender can send  $N$  unacknowledged packets in pipeline
    - When sender receive an acknowledgment, it slide forward window so that a sequence number change from *not usable* to *usable, not yet sent*
  - Each packet has sequence number
    - Sequence number



## 3.4 Go-Back-N (GBN)

- Go-back-N
  - Receiver
    - Always send ACK for correctly-received packet with the highest in-order sequence number
    - Generate duplicate ACKs when receive out-of-order packet
      - Discard out-of-order packet because receiver does not have buffer
    - Re-acknowledge packet with the highest in-order sequence number
      - E.g.:
        - Suppose,
          - Receiver
            - Expect to receive packet  $n$
            - However, it receive packet  $n+1$ . This mean that packet  $n$  is lost
            - Discard packet  $n+1$
            - Resend acknowledgement for packet  $n-1$
          - Sender
            - Receive duplicated ACKs
            - Resend packet  $n$  and  $n+1$

# 3.4 Go-Back-N (GBN)



## 3.4 Selective Repeat (SR)

### sender

#### Data from upper layer

- If there is *usable, not yet sent* in window, send the packet

#### Timeout of packet $n$

- Resend packet  $n$ , restart timer of packet  $n$

#### Receive acknowledge for packet $n$

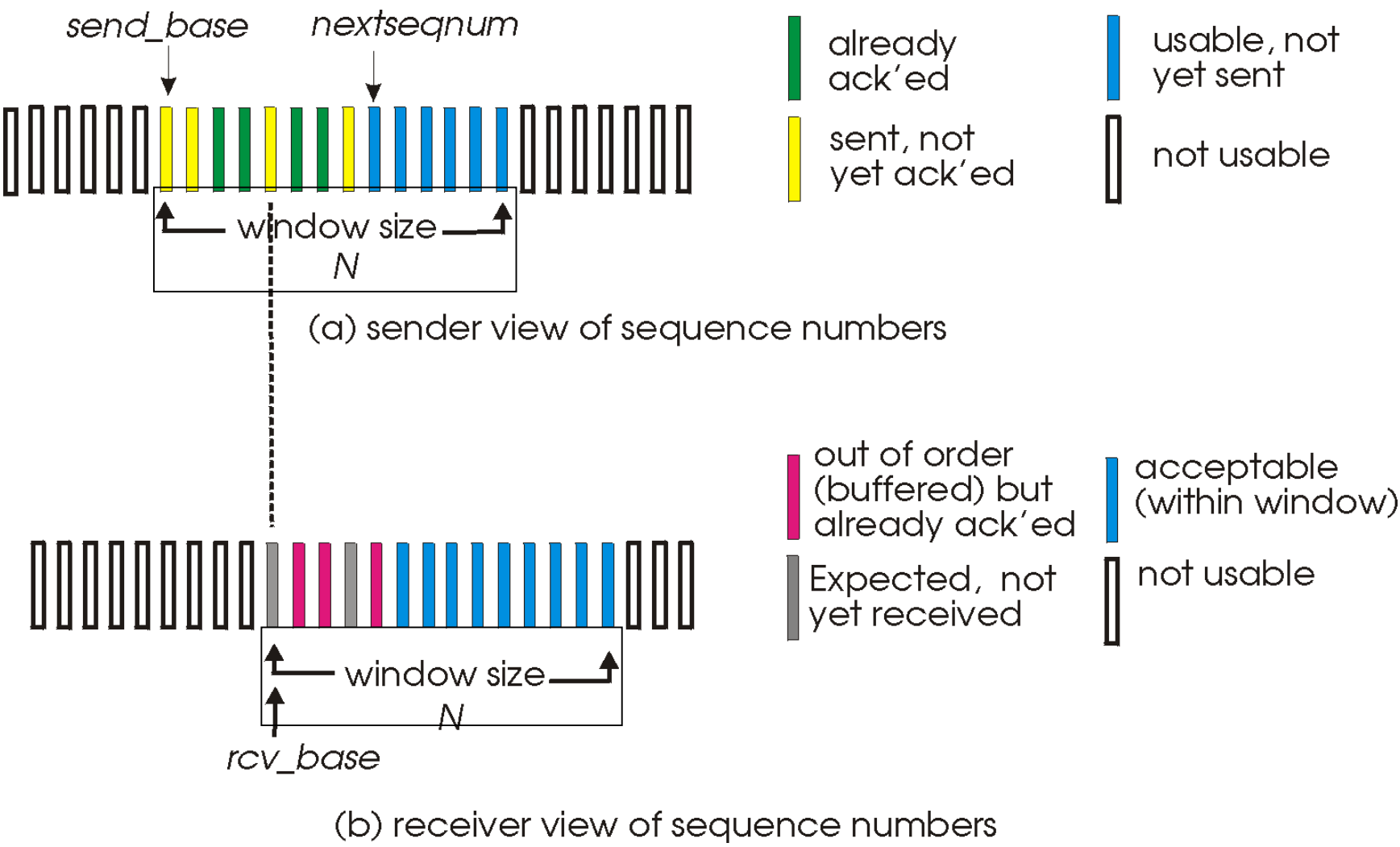
- When sender receive an acknowledgment for send\_base packet, it slide forward window so that a sequence number change from not usable to usable, not yet sent

### receiver

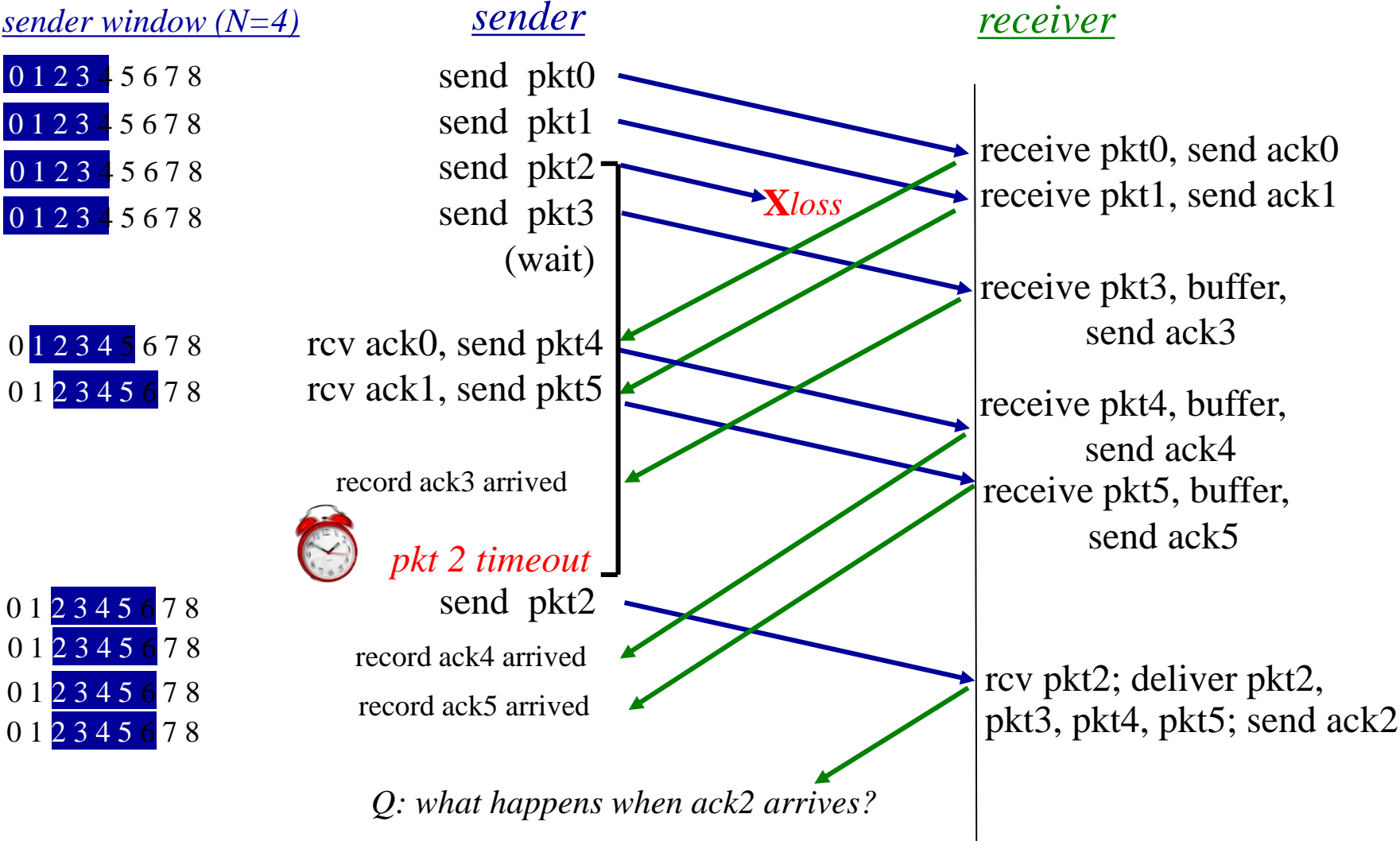
#### Receive packet $n$

- Send an acknowledgment for packet  $n$
- If
  - Out-of-order packet
    - Buffer the packet
  - In-order packet
    - Deliver (also deliver buffered, in-order packets) to upper layer

### 3.4 Selective Repeat (SR)



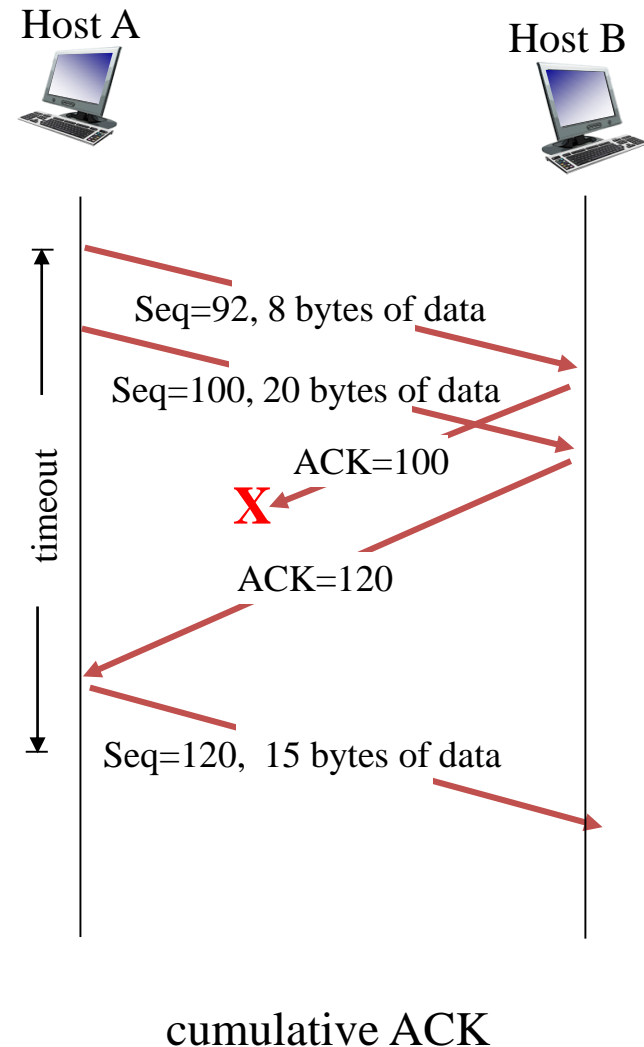
# 3.4 Selective Repeat (SR)



## 3.5 Connection-Oriented Transport: TCP

### Reliable Data Transfer: Use Go-Back-N

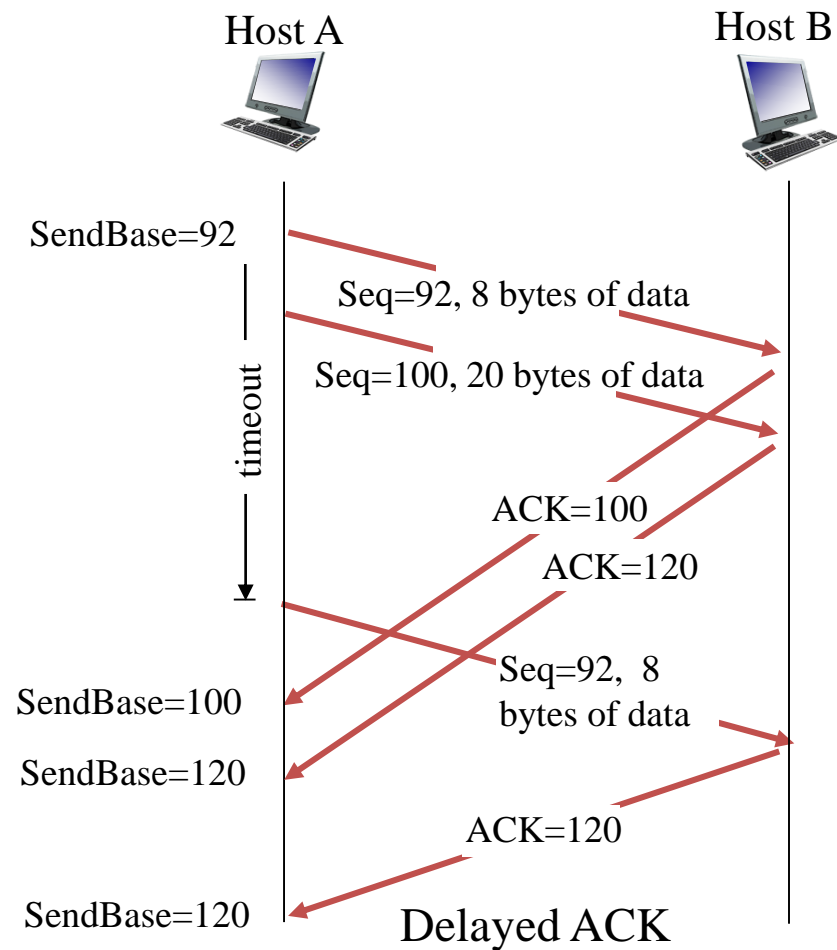
- Use Go-Back-N
  - Receiver
    - Only send cumulative acknowledgment



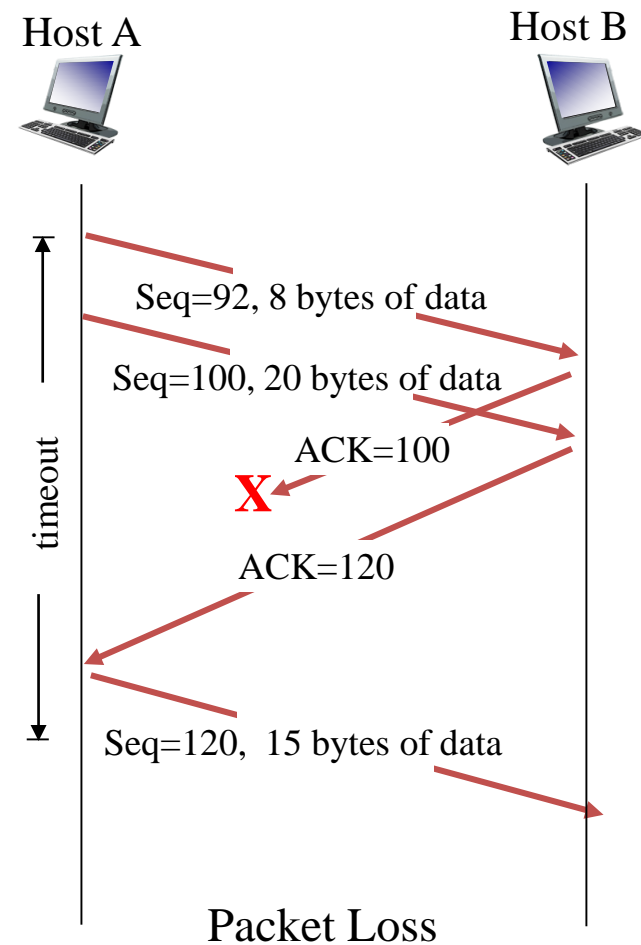


# 3.5 Connection-Oriented Transport: TCP

## Reliable Data Transfer: TCP Retransmission Scenario



Host A receive ACK = 120 before timer for the retransmission of packet Seq = 92 timeout. This indicate packet Seq = 100 has been received correctly, so packet Seq = 100 is not retransmitted

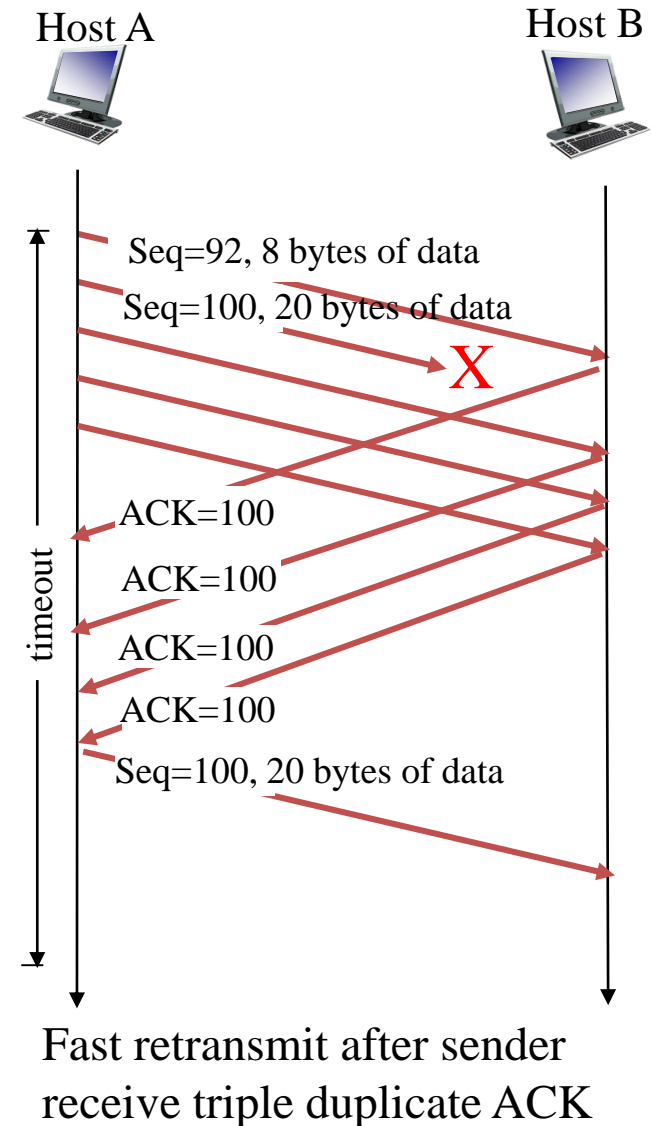


Host A receive ACK = 120. This indicate packet Seq = 92 and Seq = 100 have been received correctly. So, even though ACK = 100 is lost, packet Seq = 92 is not retransmitted

### 3.5 Connection-Oriented Transport: TCP

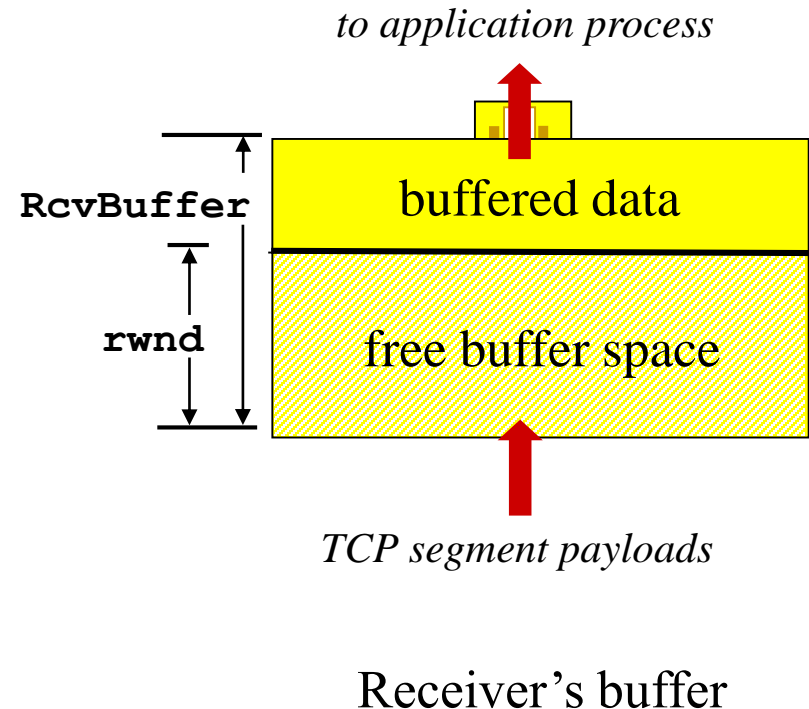
#### Reliable Data Transfer: Use Fast Retransmission

- How to detect packet loss as soon as possible?
  - Problem
    - Timeout period is too long
      - Sender wait long before retransmit lost packet
        - Increase delay
  - Solution
    - Fast retransmission
      - Sender receive triple duplicate ACK for packet  $n$ 
        - Retransmit packet  $n+1$



## 3.5 Flow Control

- Receiver
  - Control a source host's sending rate so that it does not overflow receive host's buffer
  - Inform sender the receive window value **rwnd**
    - **rwnd**
      - Unit: byte
      - Indicate the buffer space availability at receiver
- Sender
  - Send packets of size  $< \mathbf{rwnd}$  bytes to receiver

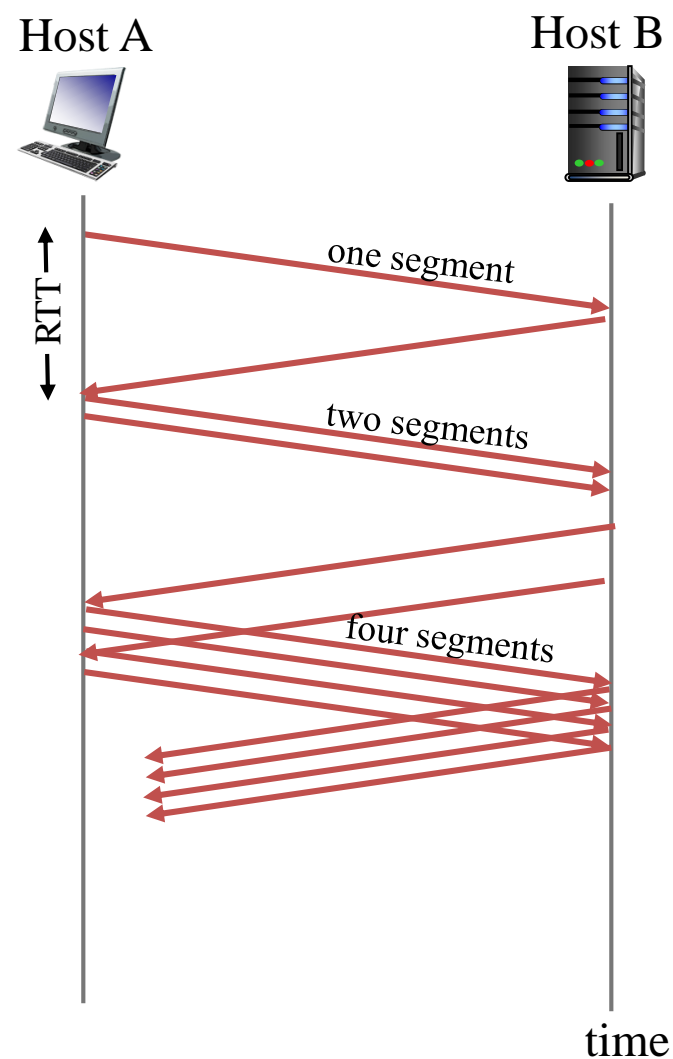


## 3.7 TCP Congestion Control

- Characteristics
  - End-to-end
    - Destination host infer congestion from observed packet loss and end-to-end delay
      - End-to-end delay =  $RTT / 2$  where RTT is Round Trip Time
- Control a source host's sending rate so that each connection traversing a congested link get an equal share of link bandwidth
  - Sender
    - Send packets of size  $< \min \{ \text{rwnd}, \text{cwnd} \}$  bytes to receiver
      - **rwnd** is receive window
      - **cwnd** is congestion window
      - So, either flow control or congestion control take effect
    - Sending rate =  $\text{cwnd} / RTT$  bytes/sec
      - This means, sender send **cwnd** bytes, wait RTT for ACK, then send more
- Two main components
  - **Slow start**
  - **Congestion avoidance**

# 3.7 TCP Congestion Control: Slow Start

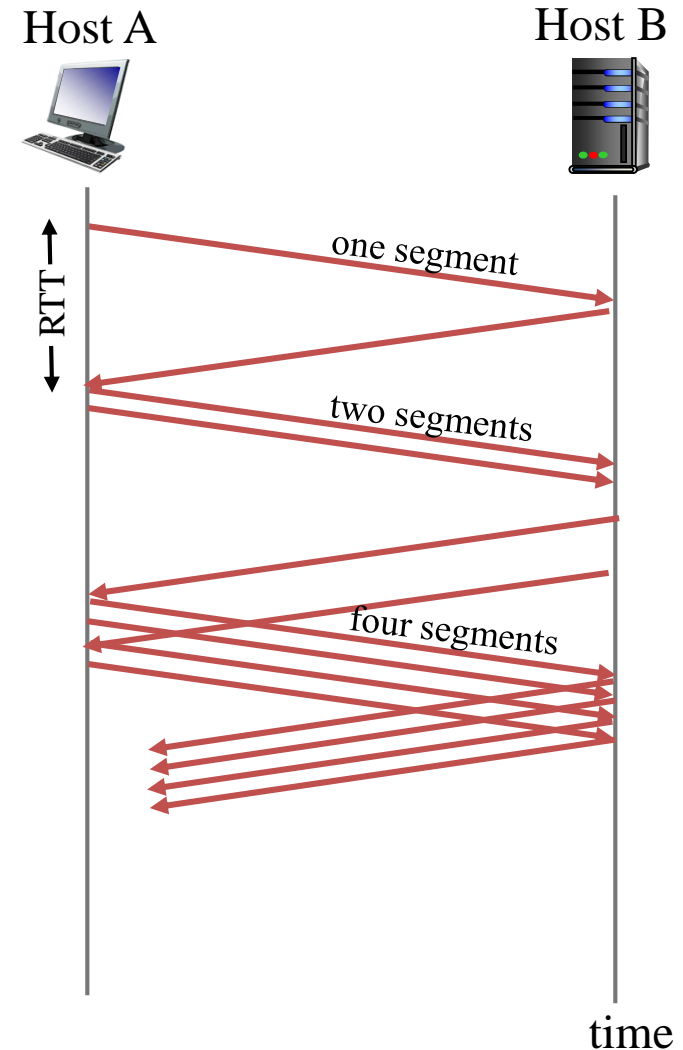
- How it work?
  - When connection begin, increase rate exponentially until
    - first packet loss or
    - **cwnd = ssthresh**
      - **ssthresh** threshold is reached
  - Increase **cwnd** for every ACK received
    - Initially, **cwnd** = 1 MSS
      - Maximum Segment Size (MSS) is the maximum amount of data that can be included in a segment
    - Then, double **cwnd** every RTT
      - Example:
        - First round: 1 segment
        - Second round: 2 segment
        - Third round: 4 segment
  - So, initial rate is slow but increase exponentially





## 3.7 TCP Congestion Control: Congestion Avoidance

- How it work?
  - When connection **cwnd** = **ssthresh**
  - Increase **cwnd** by 1 MSS for every ACK received (or 1 RTT)
    - Note: 1 RTT = 1 ACK received
  - So, initial rate is slow but increase exponentially
- During packet loss
  - Firstly, set **ssthresh** = **cwnd** / 2
  - Then, set **cwnd** = 1 MSS



## 3.7 TCP Congestion Control: Packet Loss

- During packet loss (Either timeout or triple duplicate ACK is received)
  - TCP Tahoe
    - Firstly, set **ssthresh** = **cwnd** / 2
    - Then, set **cwnd** = 1 MSS
  - TCP Reno
    - Firstly, set **ssthresh** = **cwnd** / 2
    - Then,
      - For triple duplicate ACK
        - set **cwnd** = **cwnd** / 2
      - For timeout
        - set **cwnd** = 1 MSS
- So, TCP congestion control applies **Additive Increase Multiplicative Decrease (AIMD)**
  - Additive increase
    - Increase **cwnd** by 1 MSS
  - Multiplicative decrease
    - set **ssthresh** = **cwnd** / 2

