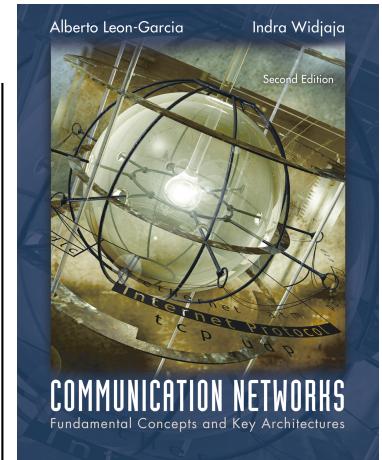
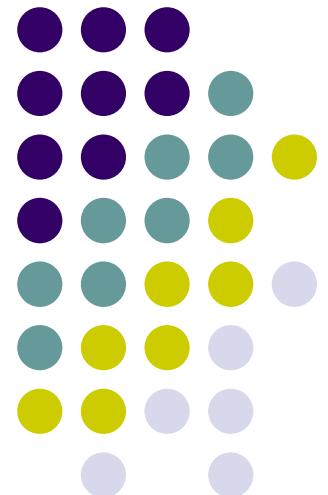


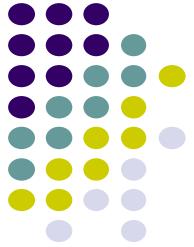
# *ARQ Protocols and Reliable Data Transfer*



**Stop-and-Wait** (pp 282-300)

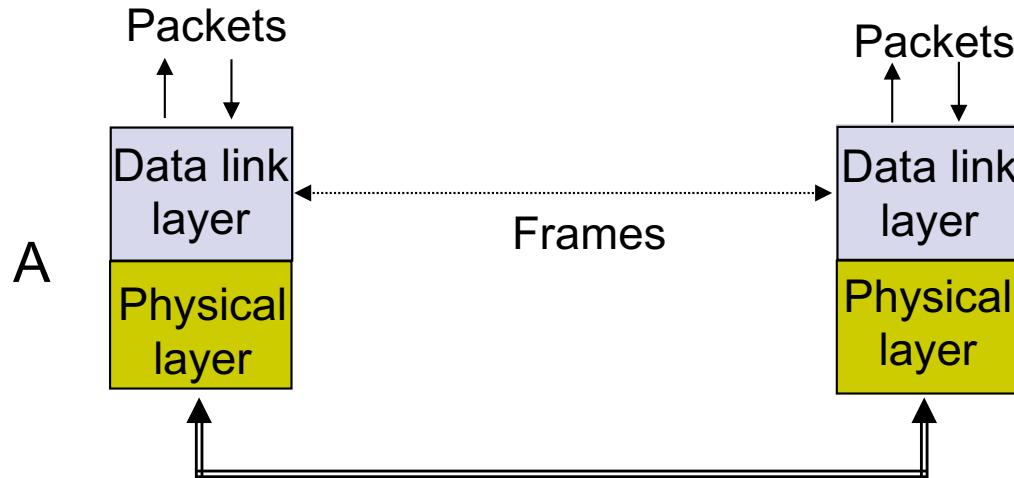
**Selective Repeat** (pp 300-315)



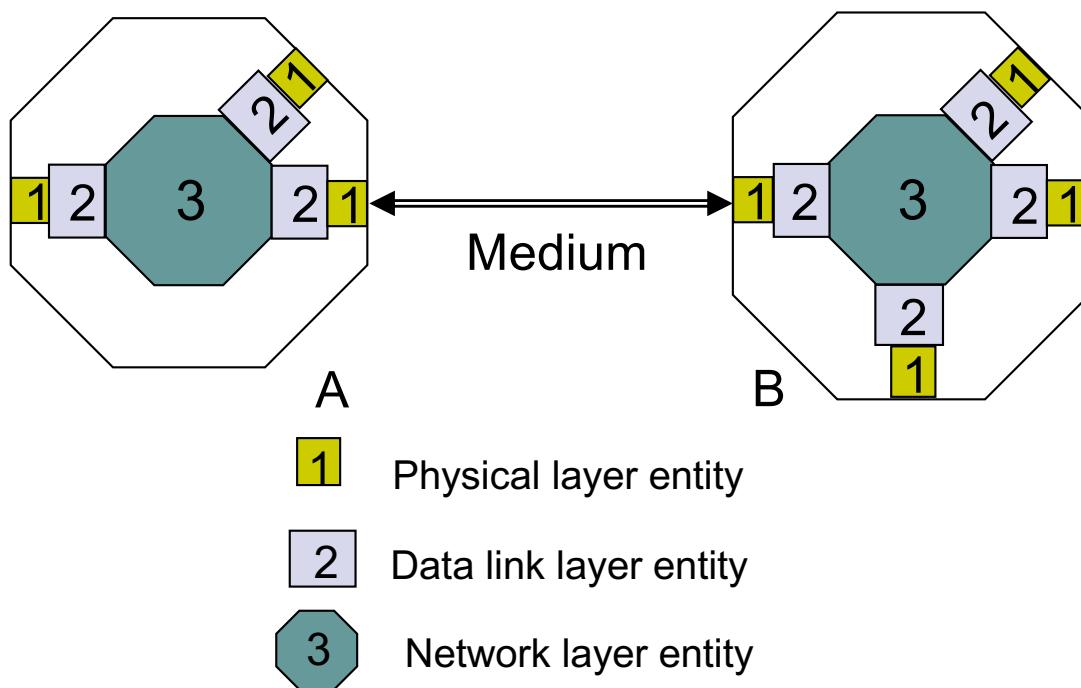


# Error control in Data Link Layer

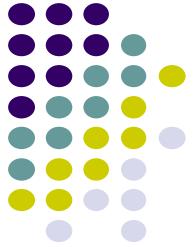
(a)



(b)

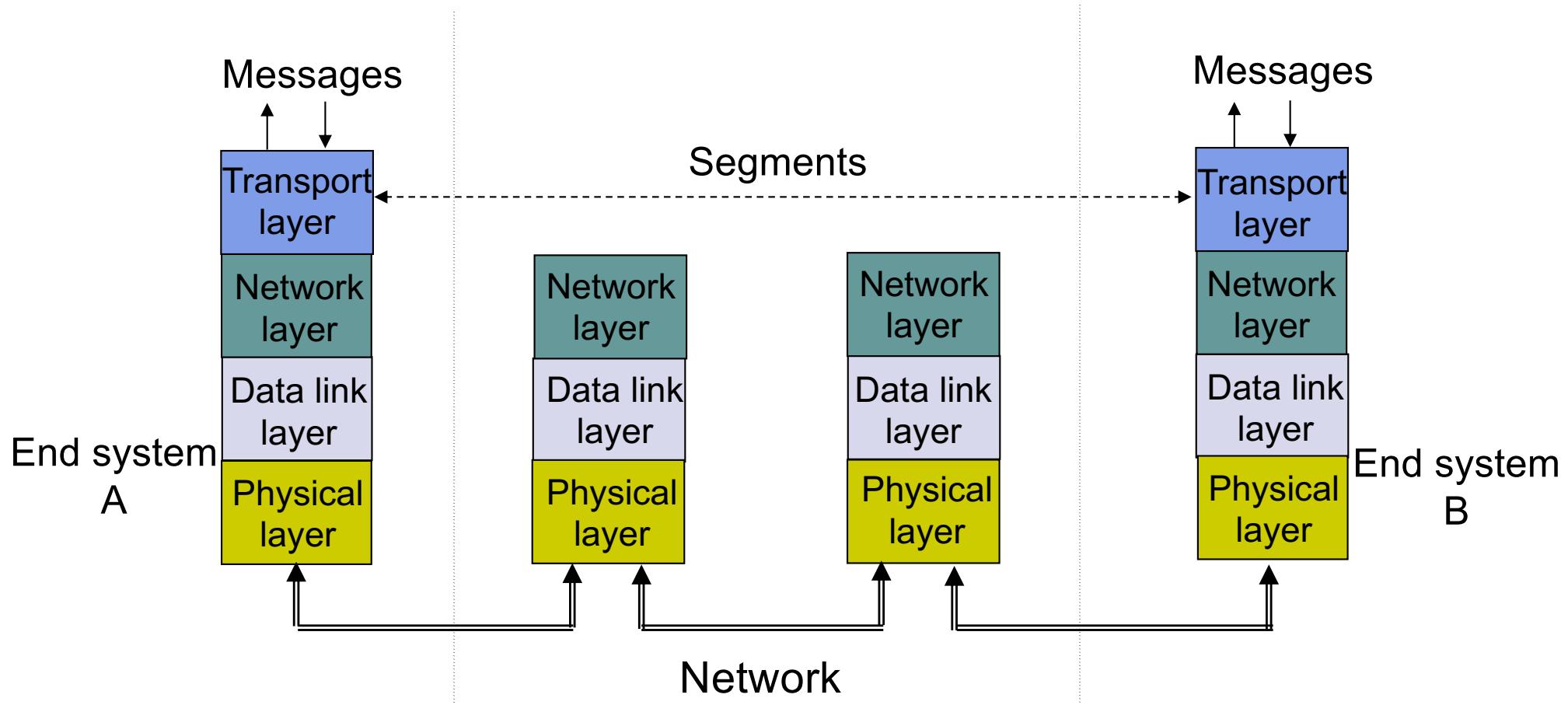


- Data Link operates over wire-like, directly-connected systems
- Frames can be corrupted or lost, but ***arrive in order***
- Data link performs error-checking & retransmission
- Ensures error-free packet transfer between two systems

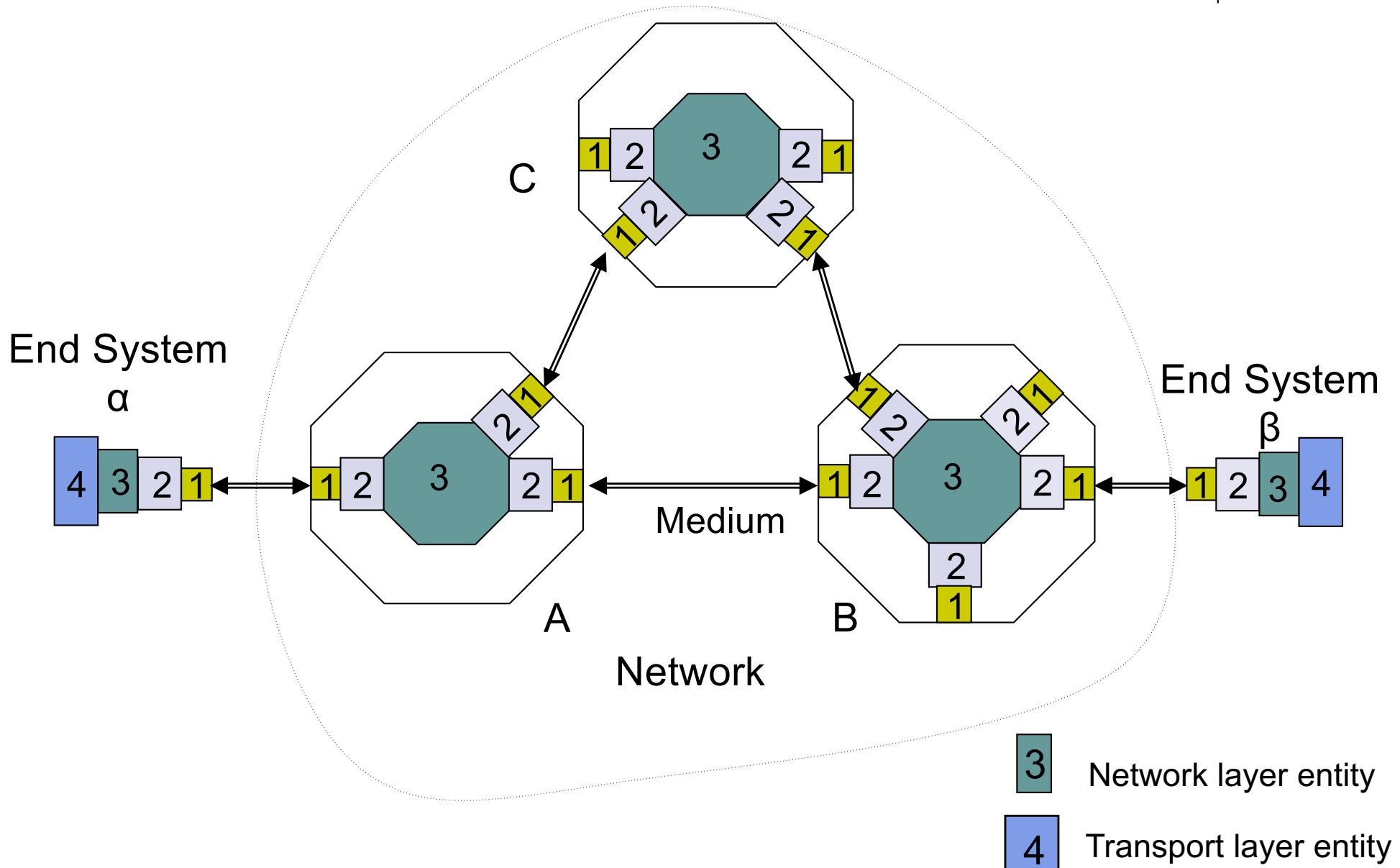
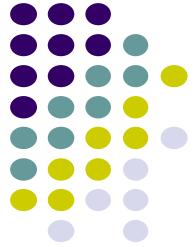


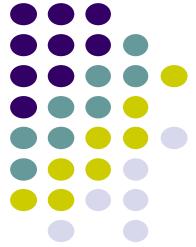
# Error Control in Transport Layer

- Transport layer protocol (e.g. TCP) sends segments across network and performs end-to-end error checking & retransmission
- Underlying network is assumed to be unreliable



- Segments can experience long delays, be lost, or **arrive out-of-order** because packets can follow different paths across network
- End-to-end error control protocol more difficult

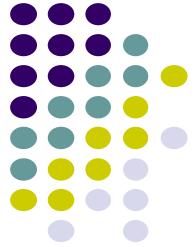




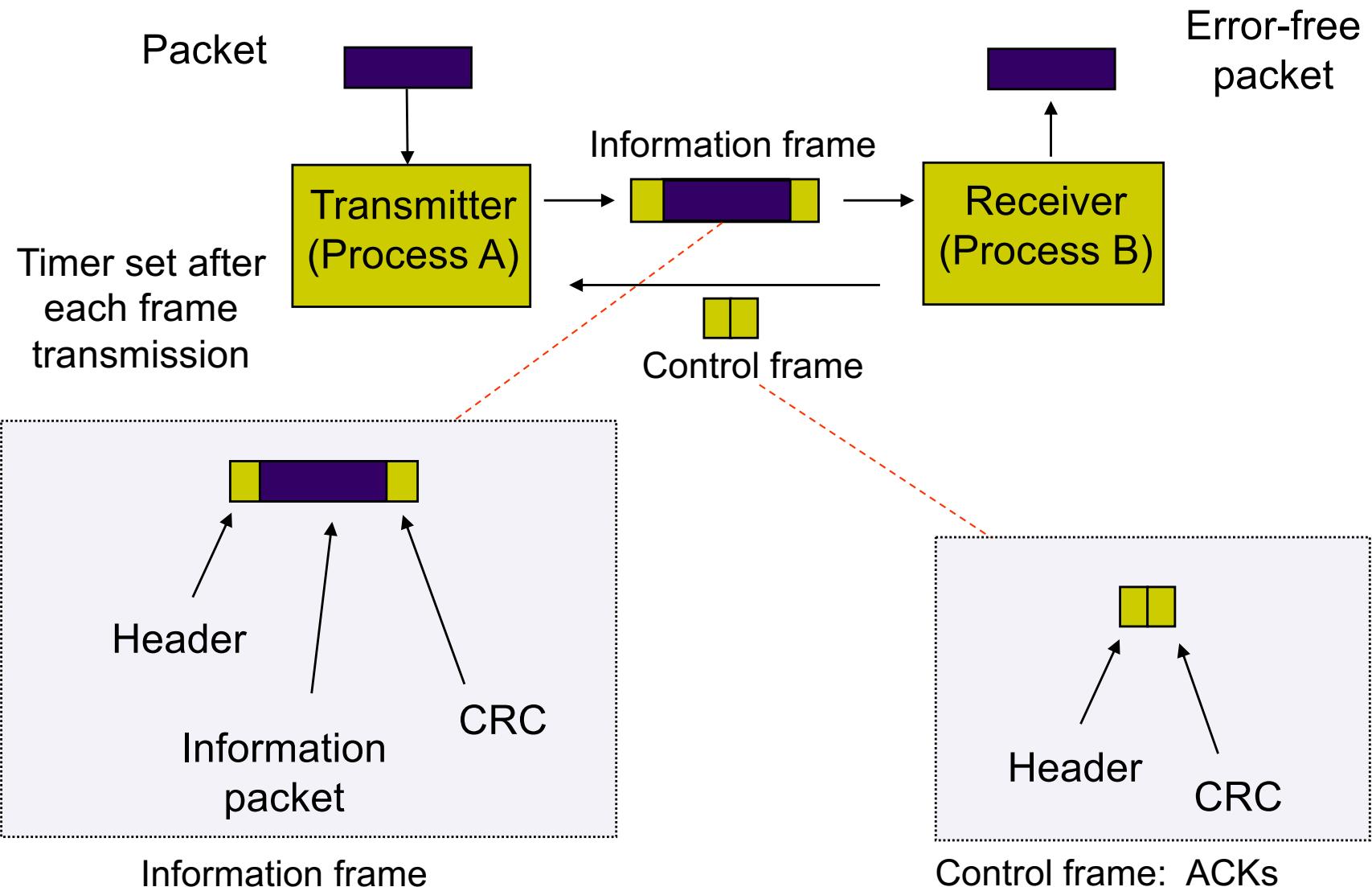
# Automatic Repeat Request (ARQ)

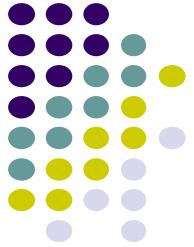
- *Purpose:* to ensure a sequence of information blocks (frames) is delivered in order and without errors or duplications despite transmission errors & losses
- We will look at:
  - Stop-and-Wait ARQ
  - Selective Repeat ARQ
- Basic elements of ARQ:
  - *Error-detecting code* with high error coverage
  - ACKs (positive acknowledgments) used extensively
  - *Timeout mechanism* used extensively

# Stop-and-Wait ARQ



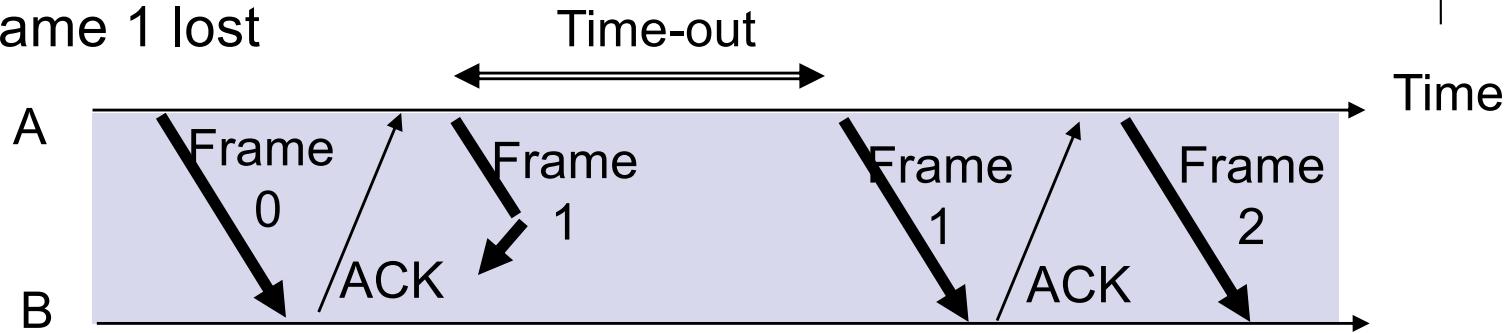
Transmit a frame, wait for ACK



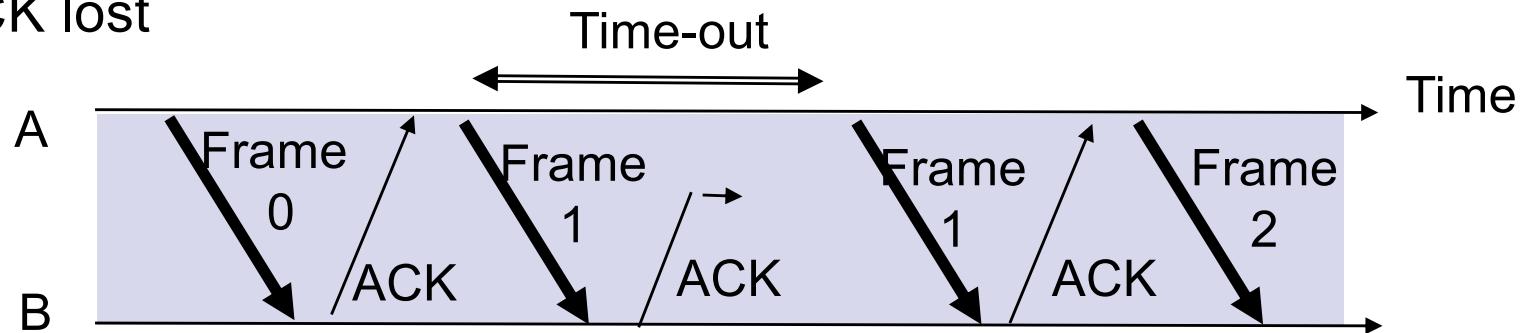


# Need for Sequence Numbers

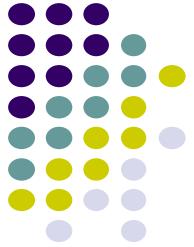
(a) Frame 1 lost



(b) ACK lost

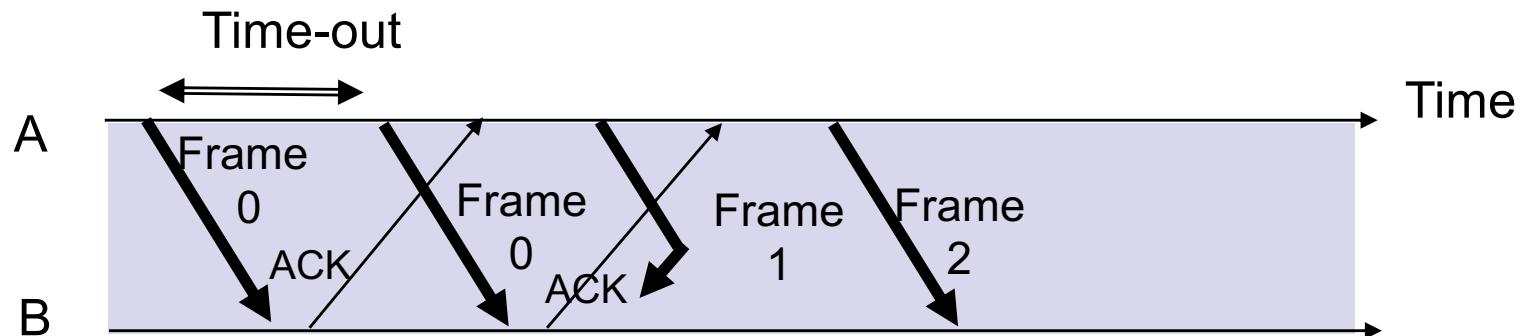


- In cases (a) & (b) the transmitting station A acts the same way
- But in case (b) the receiving station B accepts frame 1 twice
- Question: How is the receiver to know the second frame is also frame 1?
- Answer: **Add frame sequence number in header**
- $S_{last}$  is sequence number of most recent transmitted frame



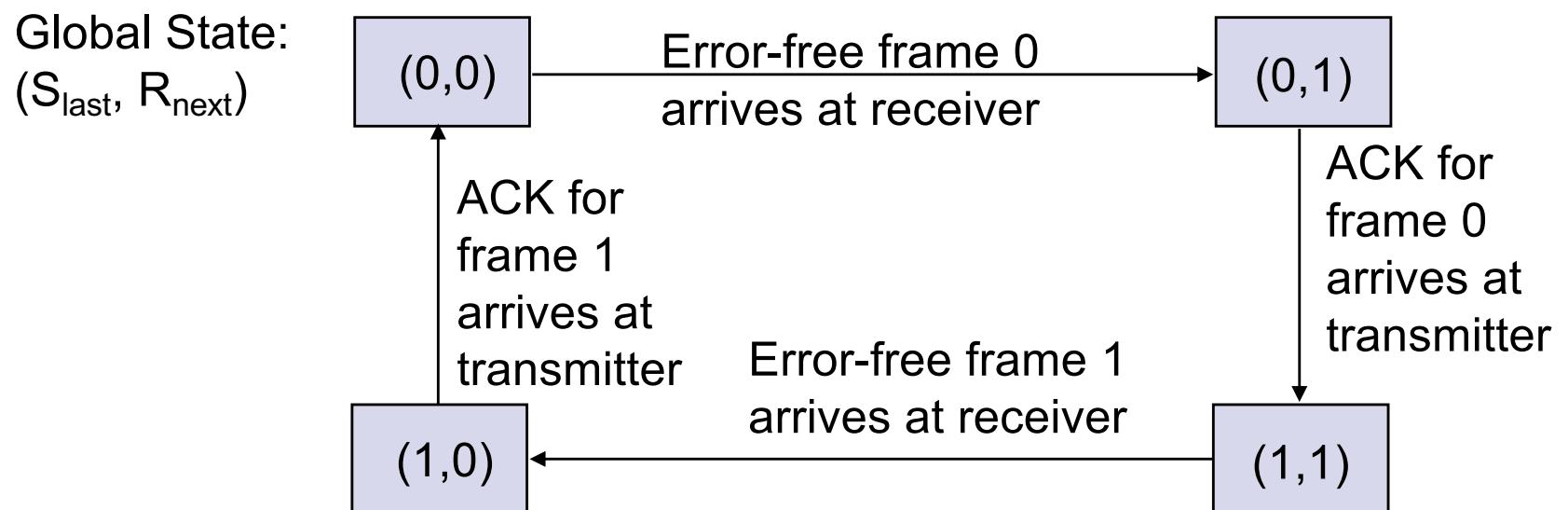
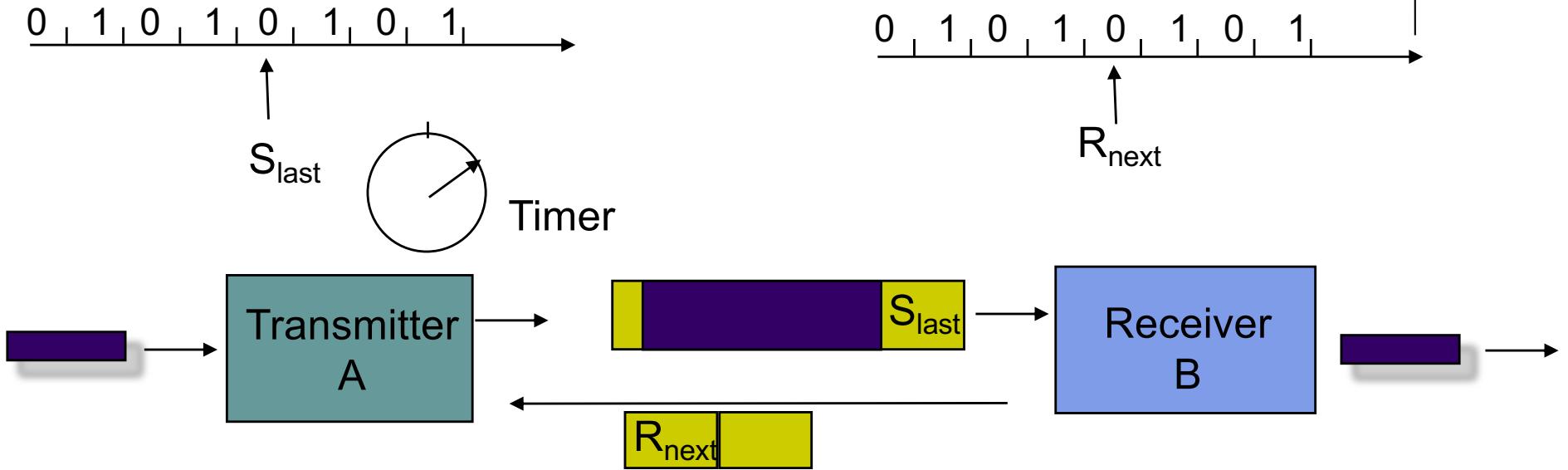
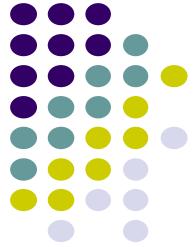
# Sequence Numbers

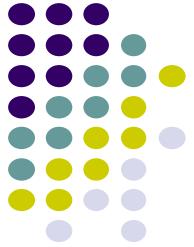
(c) Premature Time-out



- The transmitting station A misinterprets duplicate ACKs
- Incorrectly assumes second ACK acknowledges Frame 1
- Question: How is the receiver to know second ACK is for frame 0?
- Answer: **Add frame sequence number in ACK header**
- $R_{next}$  is sequence number of next frame expected by the receiver
- Implicitly acknowledges receipt of all prior frames

# 1-Bit Sequence Numbering Suffixes





# Stop-and-Wait ARQ

## Transmitter

### Ready state

- Await request from higher layer for packet transfer
- When request arrives, transmit frame with updated  $S_{last}$  and CRC
- Go to Wait State

### Wait state

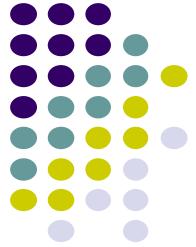
- Wait for ACK or timer to expire; block requests from higher layer
- If timeout expires
  - retransmit frame and reset timer
- If ACK received:
  - If sequence number is incorrect or if errors detected: ignore ACK
  - If sequence number is correct ( $R_{next} = S_{last} + 1$ ): accept frame, go to Ready state

## Receiver

### Always in Ready State

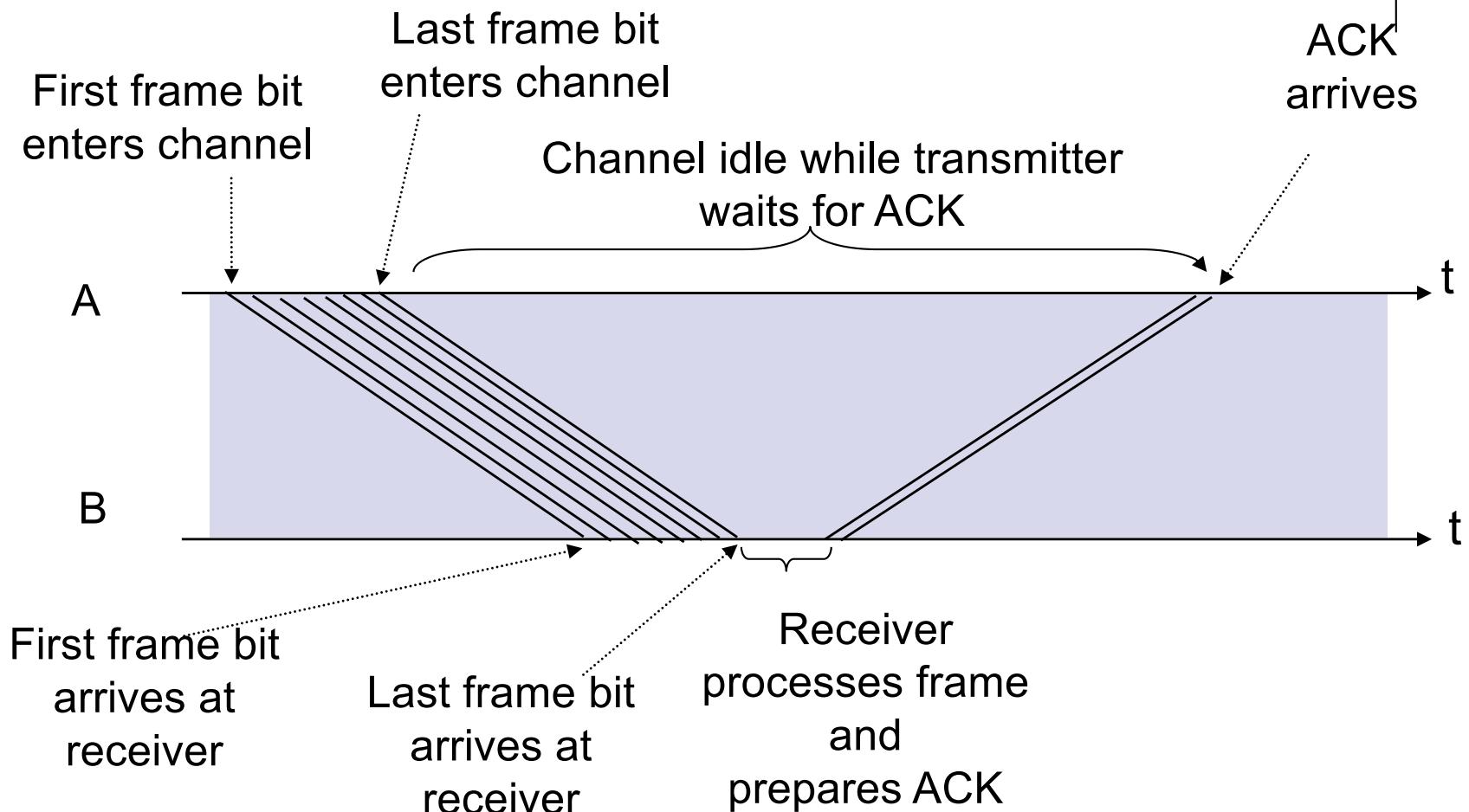
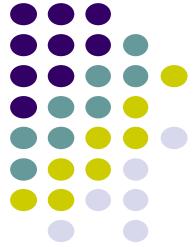
- Wait for arrival of new frame
- When frame arrives, check for errors
- If no errors detected and sequence number is correct ( $S_{last} = R_{next}$ ), then
  - accept frame,
  - update  $R_{next}$ ,
  - send ACK frame with  $R_{next}$ ,
  - deliver packet to higher layer
- If no errors detected and wrong sequence number
  - discard frame
  - send ACK frame with  $R_{next}$
- If errors detected
  - discard frame

# Applications of Stop-and-Wait ARQ



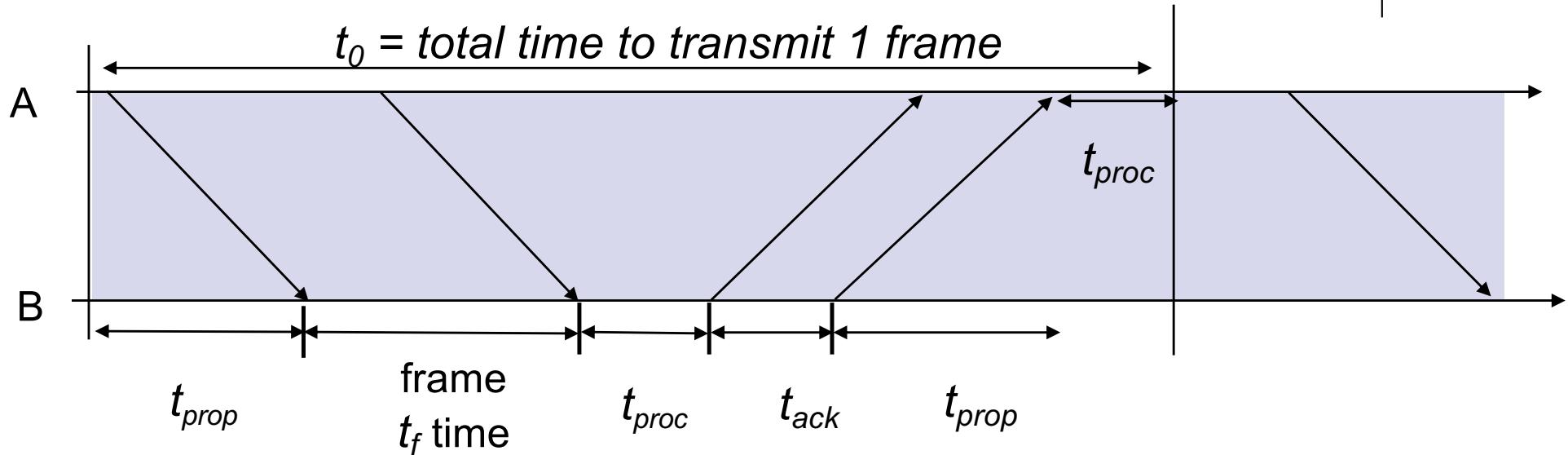
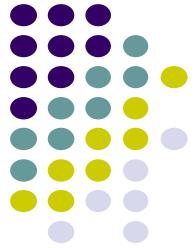
- IBM *Binary Synchronous Communications protocol* (Bisync): character-oriented data link control
- *Xmodem*: modem file transfer protocol
- *Trivial File Transfer Protocol* (RFC 1350): simple protocol for file transfer over UDP

# Stop-and-Wait Efficiency



- 10000 bit frame @ 1 Mbps takes 10 ms to transmit
- If wait for ACK = 1 ms, then efficiency =  $10/11 = 91\%$
- If wait for ACK = 20 ms, then efficiency =  $10/30 = 33\%$

# Stop-and-Wait Model

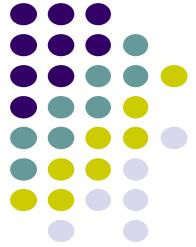


$$t_0 = 2t_{prop} + 2t_{proc} + t_f + t_{ack} \quad \text{bits/info frame}$$

$$= 2t_{prop} + 2t_{proc} + \frac{n_f}{R} + \frac{n_a}{R} \quad \text{bits/ACK frame}$$

channel transmission rate

# S&W Efficiency on Error-free channel



**Effective transmission rate:**

$$R_{eff}^0 = \frac{\text{number of information bits delivered to destination}}{\text{total time required to deliver the information bits}} = \frac{n_f - n_o}{t_0},$$

bits for header & CRC

**Transmission efficiency:**

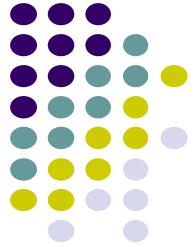
$$\eta_0 = \frac{R_{eff}}{R} = \frac{\frac{n_f - n_o}{t_0}}{R} = \frac{\frac{n_f - n_o}{n_f}}{\frac{t_0}{n_f}} = \frac{1 - \frac{n_o}{n_f}}{1 + \frac{n_a}{n_f} + \frac{2(t_{prop} + t_{proc})R}{n_f}}$$

Effect of frame overhead

Effect of ACK frame

Effect of **Delay-Bandwidth Product**

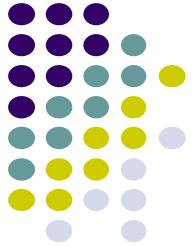
# Example: Impact of Delay-Bandwidth Product



$n_f=1250$  bytes = 10000 bits,  $n_a=n_o=25$  bytes = 200 bits

2xDelayxBW Efficiency	1 ms 200 km	10 ms 2000 km	100 ms 20000 km	1 sec 200000 km
1 Mbps	$10^3$ 88%	$10^4$ 49%	$10^5$ 9%	$10^6$ 1%
1 Gbps	$10^6$ 1%	$10^7$ 0.1%	$10^8$ 0.01%	$10^9$ 0.001%

*Stop-and-Wait does not work well for very high speeds or long propagation delays*



# Geometric RV

$X = \#$  transmissions till  
first success

- $X=1, (1-P_f)$
- $X=2, (1-P_f)P_f$
- $X=3, (1-P_f)P_f^2$
- ...
- $X=k, (1-P_f)P_f^{k-1}$
- ...

Mean of X

$$E[X]$$

$$\begin{aligned} &= (1-P_f) + 2(1-P_f)P_f + \\ &\quad + \dots + k(1-P_f)P_f^{k-1} + \dots \end{aligned}$$

$$E[X]/(1-P_f)$$

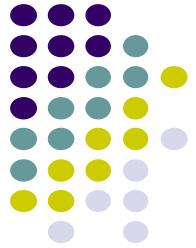
$$= 1 + 2P_f + \dots + kP_f^{k-1} + \dots$$

$$= d/dx \{P_f + P_f^2 + \dots + P_f^k + \dots\}$$

$$= 1/(1-P_f)^2$$

$$E[X] = 1/(1-P_f)$$

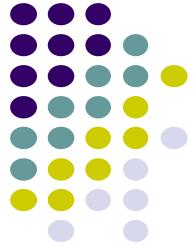
# S&W Efficiency in Channel with Errors



- Let  $1 - P_f$  = probability frame arrives w/o errors
- Avg. # of transmissions to first correct arrival is then  $1 / (1 - P_f)$
- “If 1-in-10 get through without error, then avg. 10 tries to success”
- Avg. Total Time per frame is then  $t_0 / (1 - P_f)$

$$\eta_{SW} = \frac{R_{eff}}{R} = \frac{\frac{n_f - n_o}{t_0 / (1 - P_f)}}{R} = \frac{1 - \frac{n_o}{n_f}}{1 + \frac{n_a}{n_f} + \frac{2(t_{prop} + t_{proc})R}{n_f}} (1 - P_f)$$

Effect of  
frame loss



## Example: Impact Bit Error Rate

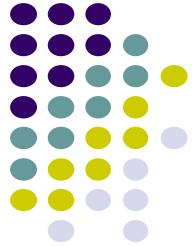
$n_f = 1250$  bytes = 10000 bits,  $n_a = n_o = 25$  bytes = 200 bits

Find efficiency for random bit errors with  $p=0, 10^{-6}, 10^{-5}, 10^{-4}$

$$1 - P_f = (1 - p)^{n_f} \approx e^{-n_f p} \text{ for large } n_f \text{ and small } p$$

$1 - P_f$ <b>Efficiency</b>	0	$10^{-6}$	$10^{-5}$	$10^{-4}$
1 Mbps	1	0.99	0.905	0.368
& 1 ms	88%	86.6%	79.2%	32.2%

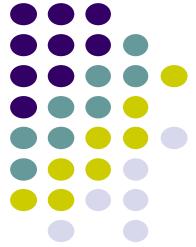
*Bit errors impact performance as  $n_f p$  approach 1*



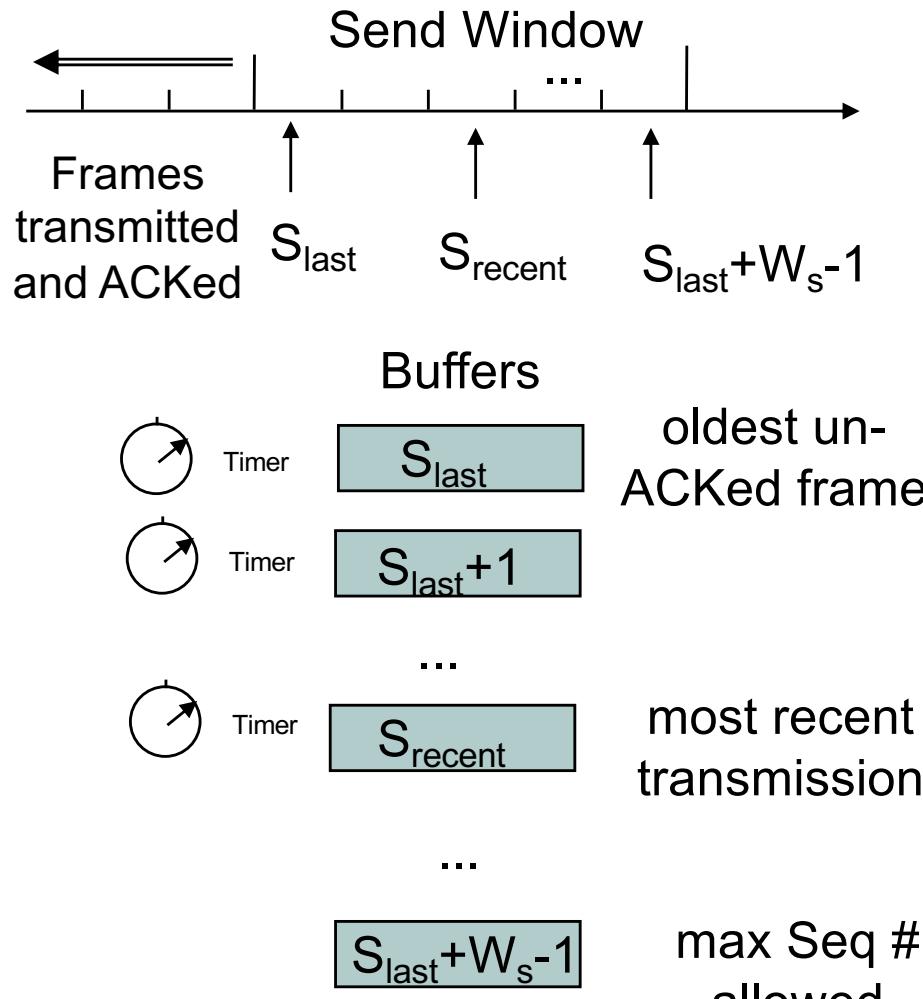
# Improvement: Don't Wait!

- Keep channel busy by continuing to send frames
- Allow a window of up to  $W_s$  outstanding frames
- Use  $m$ -bit sequence numbering
- If ACK for oldest frame arrives before window is exhausted, we can continue transmitting
- Many design options:
  - Use of timeout mechanism
  - What to ACK?
  - How to deal with out-of-order frames?
  - When and What to Retransmit?

# Send Many, Receive One

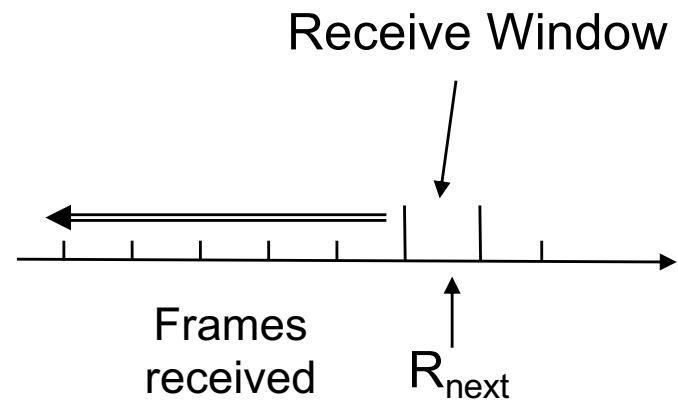


Transmitter



Maximum Send Window Size is  $W_s = 2^m - 1$

Receiver

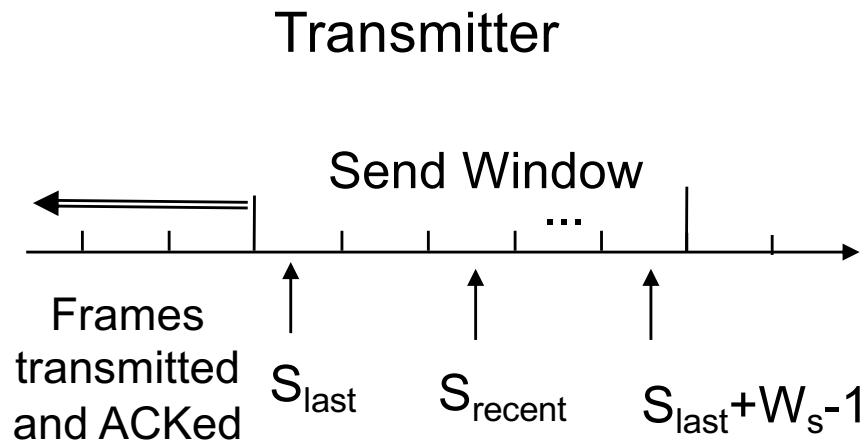
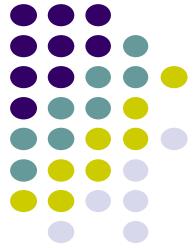


Receiver will only accept a frame that is error-free and that has sequence number  $R_{next}$

When such frame arrives  $R_{next}$  is incremented by one, so the *receive window slides forward* by one

If error-free, but out-of-sequence frame arrives, ACK with  $R_{next}$  is sent

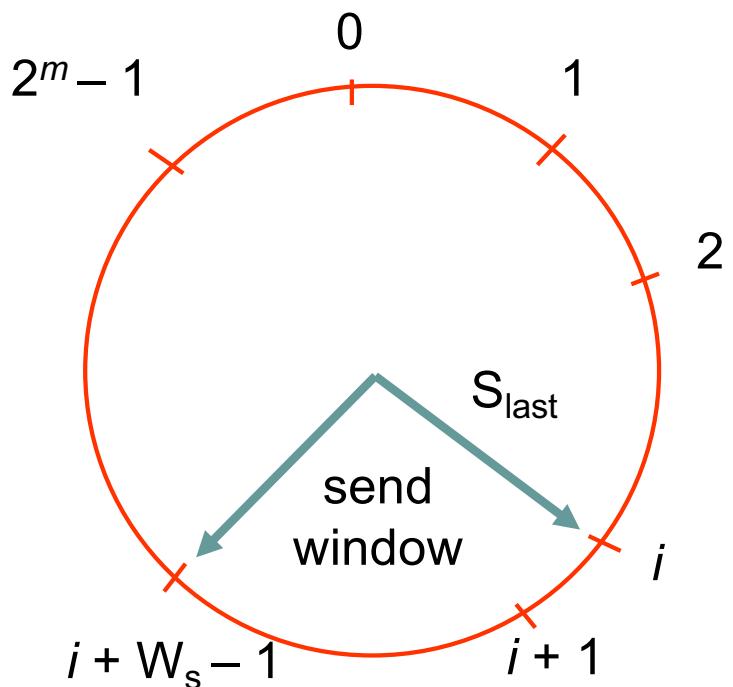
# Sliding Window Operation



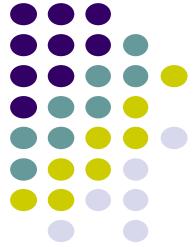
When error-free ACK frame arrives with  $R_{next}$  in interval between  $S_{last}$  and  $S_{recent}$ , then  $S_{last}$  is moved to  $R_{next}$ , since all frames prior to  $R_{next}$  have been received correctly.

Note that send window can move forward by more than one.

$m$ -bit Sequence Numbering

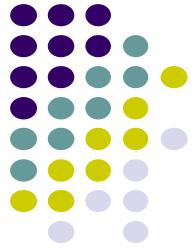


# Selective Repeat ARQ: Send Many, Receive Many

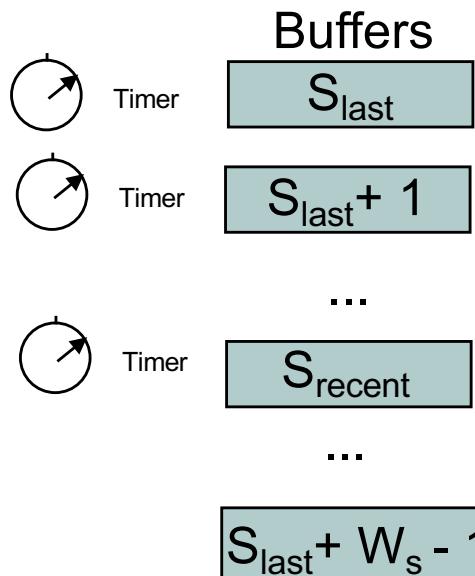
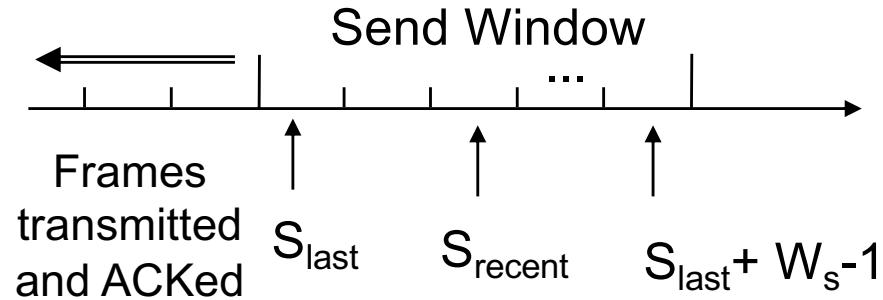


- Selective Repeat retransmits *only missing frames*
  - Timeout causes individual corresponding frame to be resent
- Receiver maintains a *receive window*  $W_r$  of sequence numbers that can be accepted
  - Error-free, but out-of-sequence frames with sequence numbers within the receive window are buffered
  - Arrival of frame with  $R_{next}$  causes window to slide forward by 1 or more
- $W_r + W_s \leq 2^m$  for unambiguous delivery

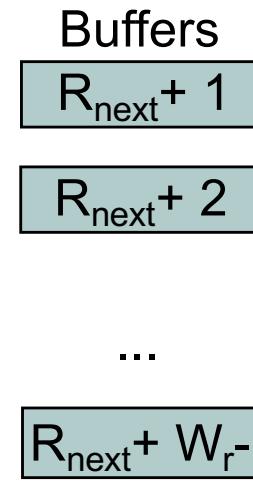
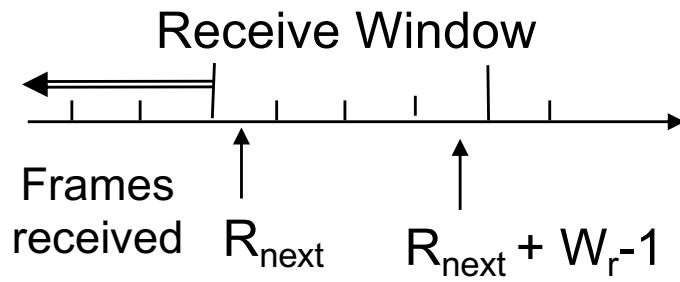
# Selective Repeat ARQ



Transmitter



Receiver

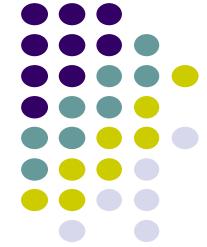


max Seq #  
accepted

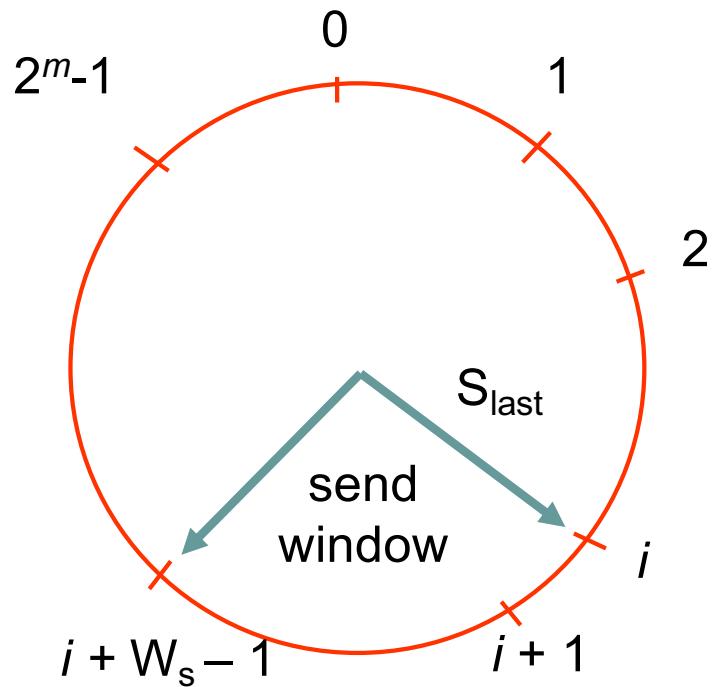
Accept error-free frames in receive window; Send ACK with new  $R_{next}$   
And list of missing frames

If error-free ACK frame arrives &  $R_{next}$  in send window, then move  $S_{last}$  to  $R_{next}$ .  
Update status of all frames in send window

# Send & Receive Windows

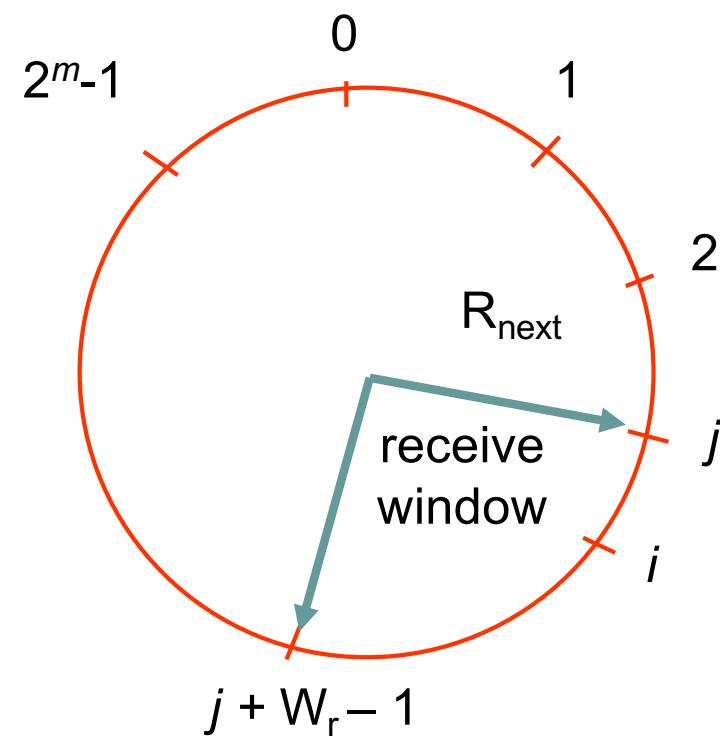


Transmitter



Moves  $k$  forward when ACK arrives with  $R_{\text{next}} = S_{\text{last}} + k$   
 $k = 1, \dots, W_s - 1$

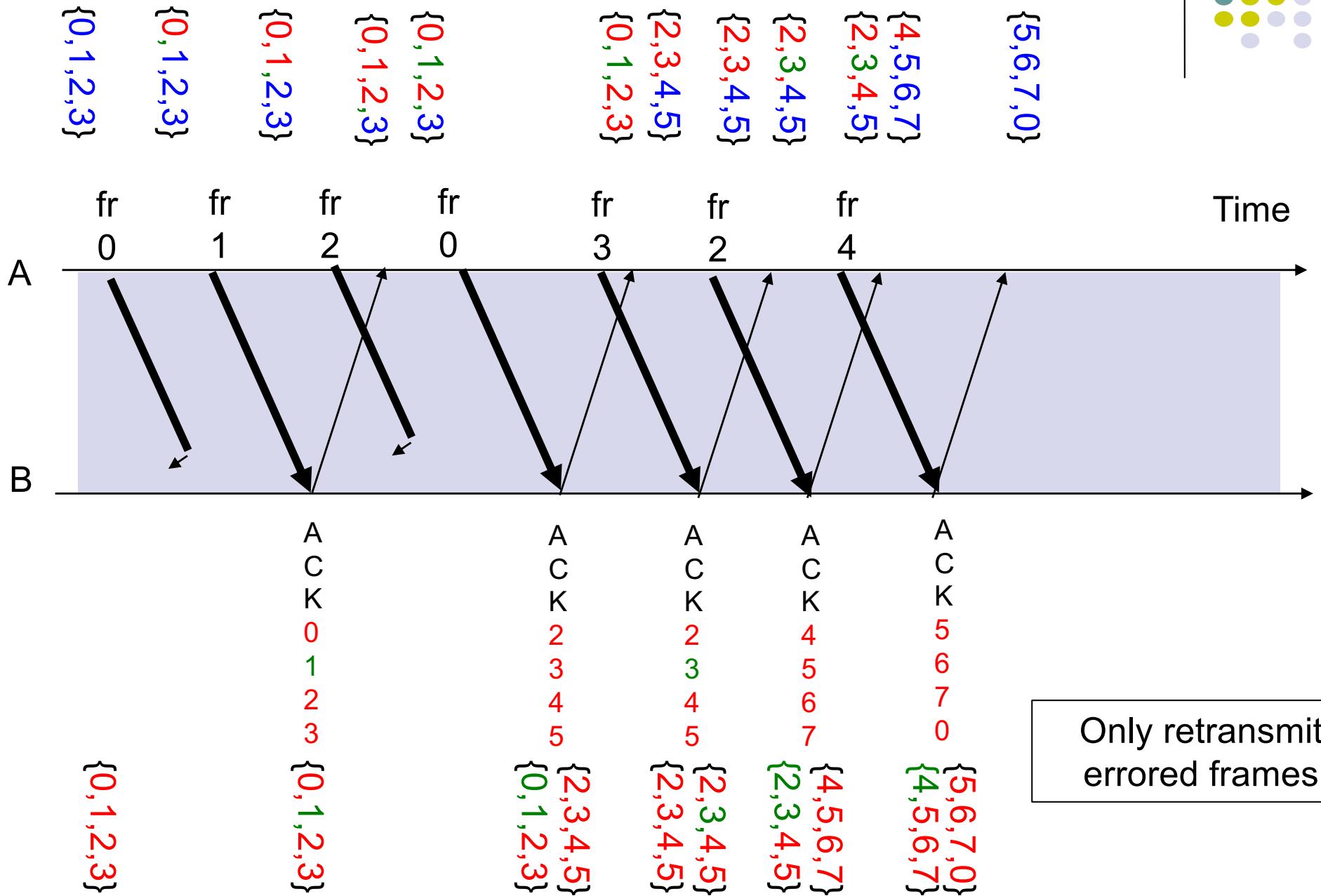
Receiver

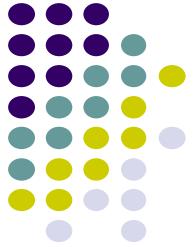


Moves forward by 1 or more when frame arrives with Seq. # =  $R_{\text{next}}$

# Selective Repeat ARQ

$M = 2^3 = 8$ ,  $Ws=4$ ,  $Wr=4$



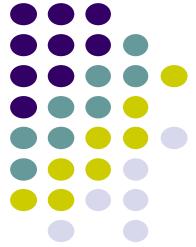


# Efficiency of Selective Repeat

- Assume  $P_f$  frame loss probability, then number of transmissions required to deliver a frame is:
  - $t_f / (1-P_f)$

$$\eta_{SR} = \frac{\frac{n_f - n_o}{t_f / (1 - P_f)}}{R} = \left(1 - \frac{n_o}{n_f}\right)(1 - P_f)$$

# Example: Impact Bit Error Rate on Selective Repeat

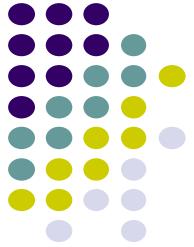


$n_f=1250$  bytes = 10000 bits,  $n_a=n_o=25$  bytes = 200 bits

Compare S&W & SR efficiency for random bit errors with  $p=0$ ,  $10^{-6}$ ,  $10^{-5}$ ,  $10^{-4}$  and  $R= 1$  Mbps & 100 ms

Efficiency	0	$10^{-6}$	$10^{-5}$	$10^{-4}$
S&W	8.9%	8.8%	8.0%	3.3%
SR	98%	97%	89%	36%

- Selective Repeat outperforms S&W, but efficiency drops as error rate increases



# Comparison of ARQ Efficiencies

Assume

$$n_a \text{ & } n_o \text{ are negligible relative to } n_f,$$
$$L = 2(t_{prop} + t_{proc})R/n_f = (W_s - 1), \text{ then}$$

Selective-Repeat:

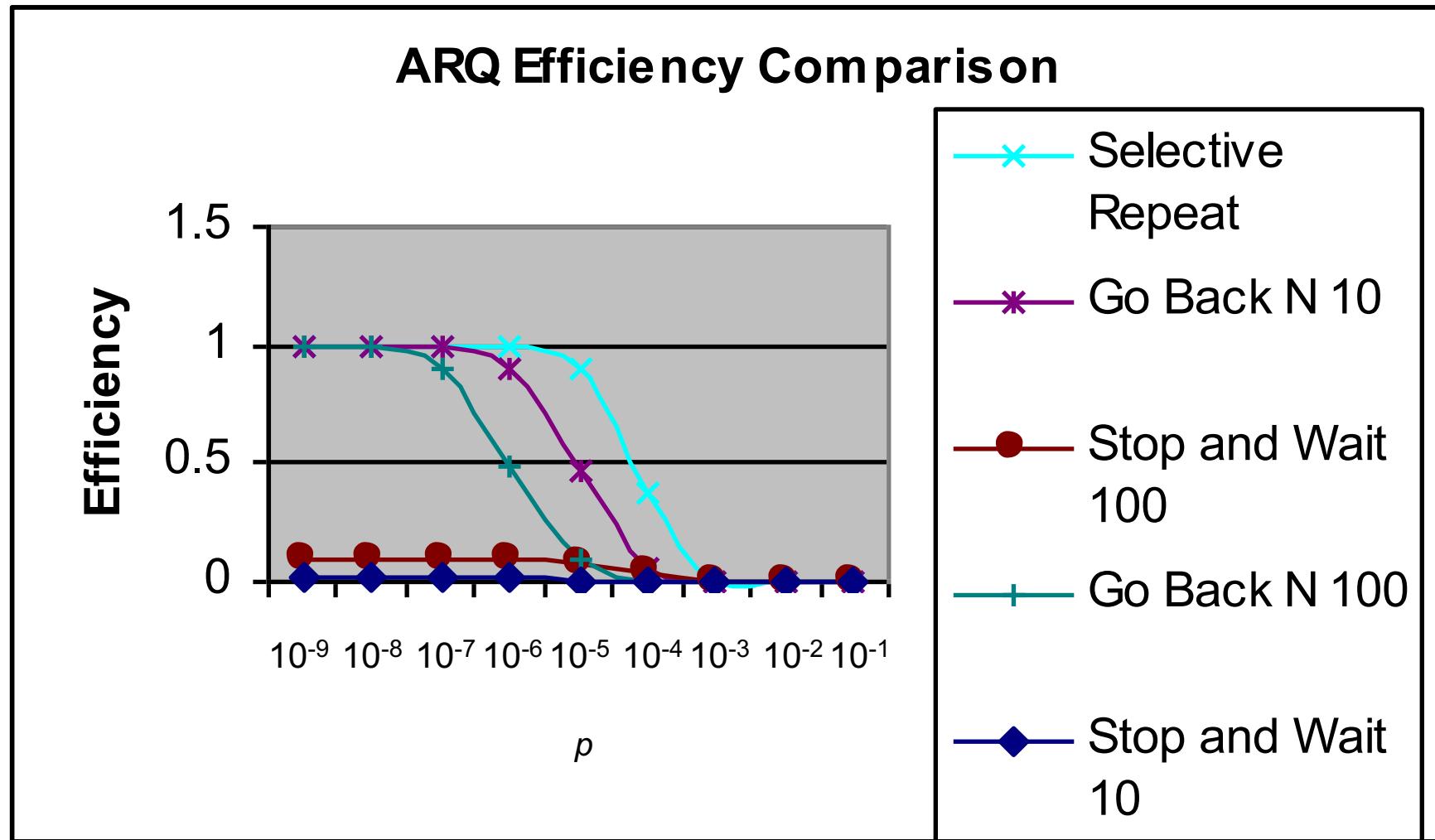
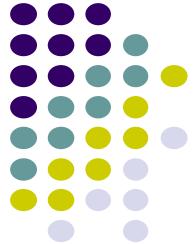
$$\eta_{SR} = (1 - P_f) \left(1 - \frac{n_o}{n_f}\right) \approx (1 - P_f)$$

For  $P_f \rightarrow 1$ , GBN & SW same

Stop-and-Wait:

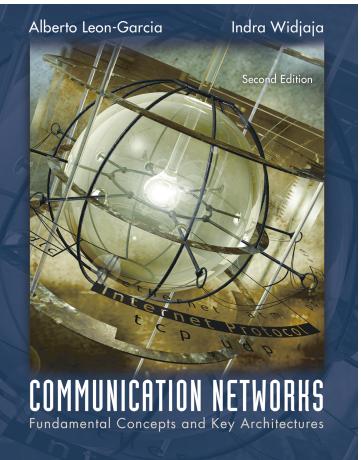
$$\eta_{SW} = \frac{(1 - P_f)}{1 + \frac{n_a}{n_f} + \frac{2(t_{prop} + t_{proc})R}{n_f}} \approx \frac{1 - P_f}{1 + L}$$

# ARQ Efficiencies

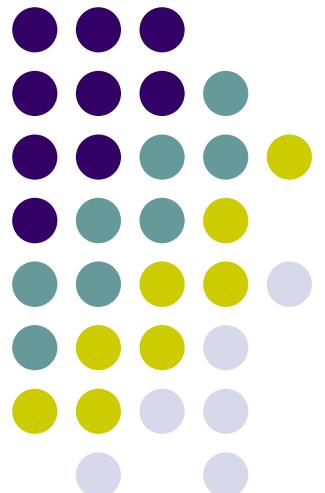


Delay-Bandwidth product = 10, 100

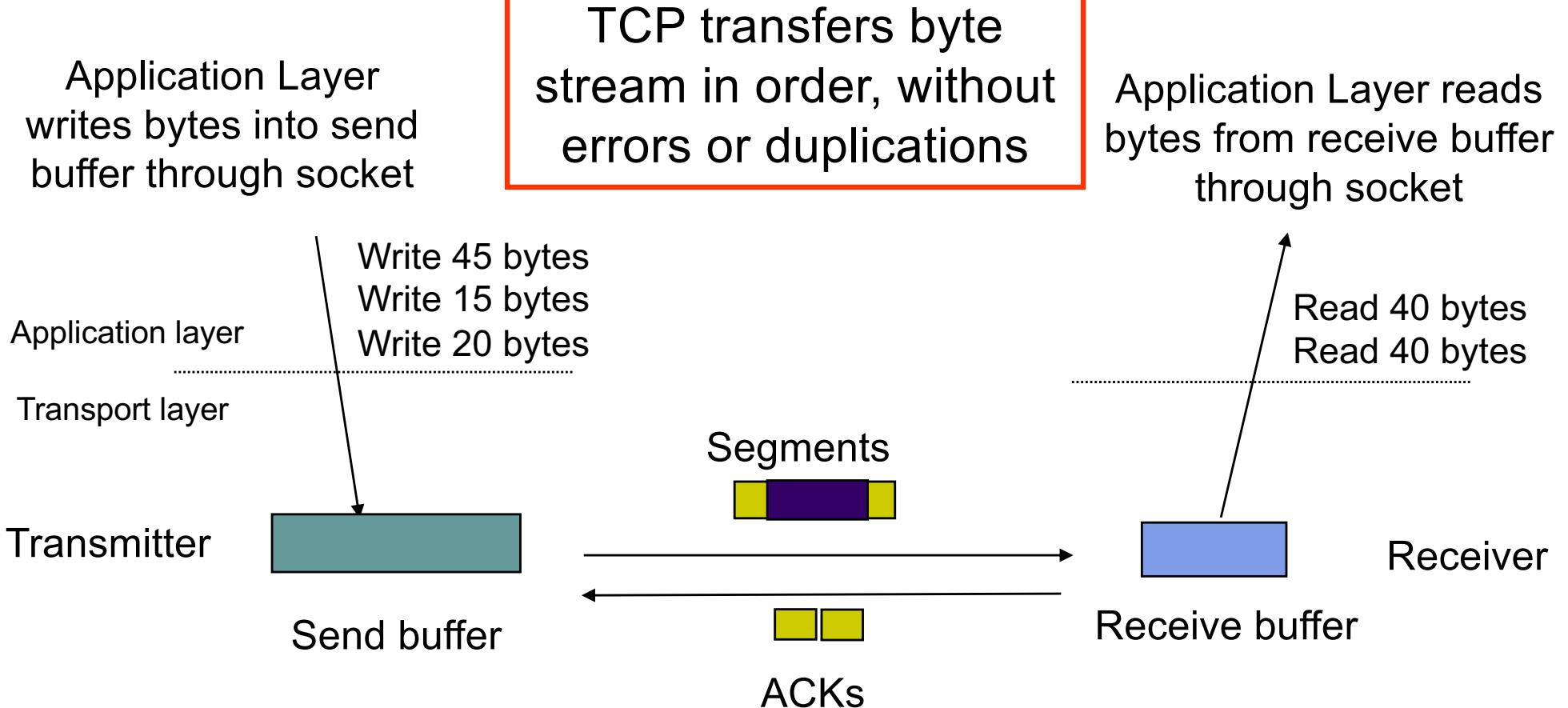
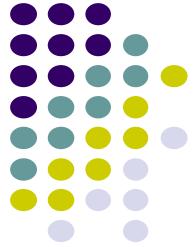
# TCP Protocol

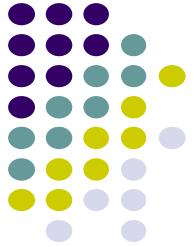


***Reliable Stream Service  
Flow Control  
Congestion control***  
pp 320-324, 602-620



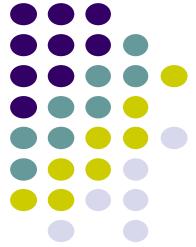
# TCP Reliable Stream Service



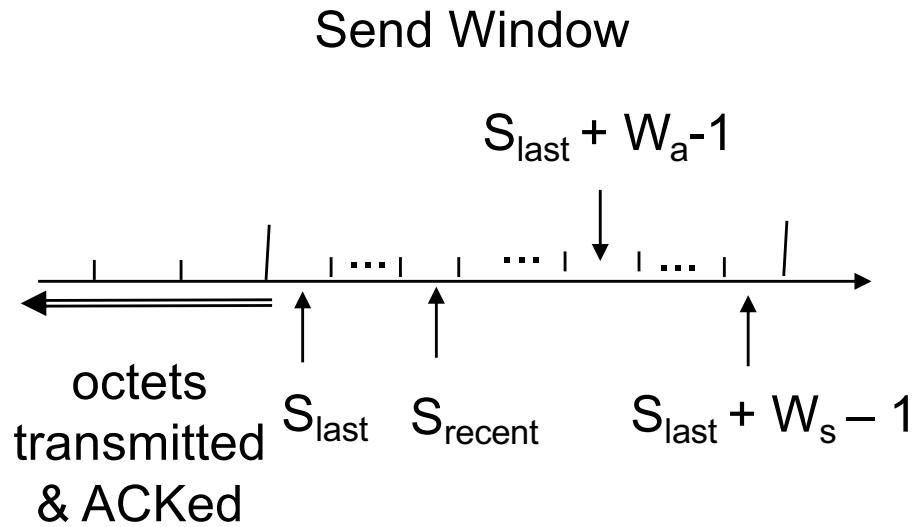


# TCP ARQ Method

- TCP uses variation of *Selective Repeat ARQ*
  - Transfers **byte stream** without preserving boundaries
  - Accepts segments within a receive window
  - Uses  $R_{next}$  for acknowledgements
- Operates over best effort service of IP
  - Packets can arrive with errors or be lost
  - Packets can arrive out-of-order
  - Packets can arrive after very long delays
  - Duplicate segments must be detected & discarded
  - Must protect against segments from previous connections
- Sequence Numbers
  - Seq. # is number of first byte in segment payload
  - Very long Seq. #s (32 bits) to deal with long delays
  - Initial sequence numbers negotiated during connection setup (to deal with very old duplicates)

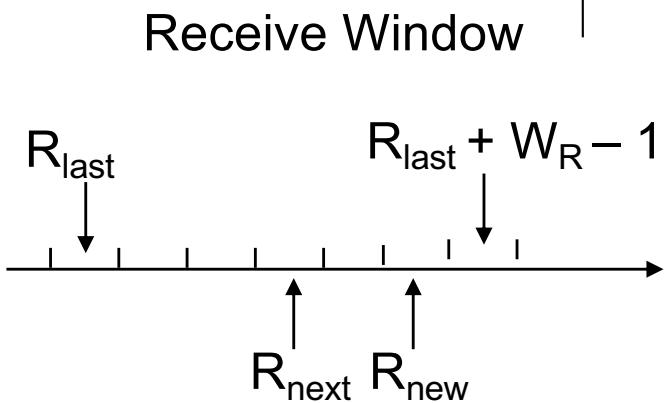


## TCP Sender

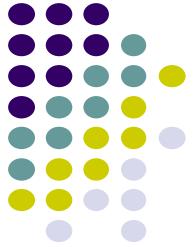


$S_{last}$  oldest unacknowledged byte  
 $S_{recent}$  highest-numbered transmitted byte  
 $S_{last} + W_a - 1$  highest-numbered byte that can be transmitted  
 $S_{last} + W_s - 1$  highest-numbered byte that can be accepted from the application

## TCP Receiver

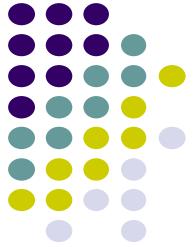


$R_{last}$  highest-numbered byte not yet read by the application  
 $R_{next}$  next expected byte  
 $R_{new}$  highest numbered byte received correctly  
 $R_{last} + W_R - 1$  highest-numbered byte that can be accommodated in receive buffer

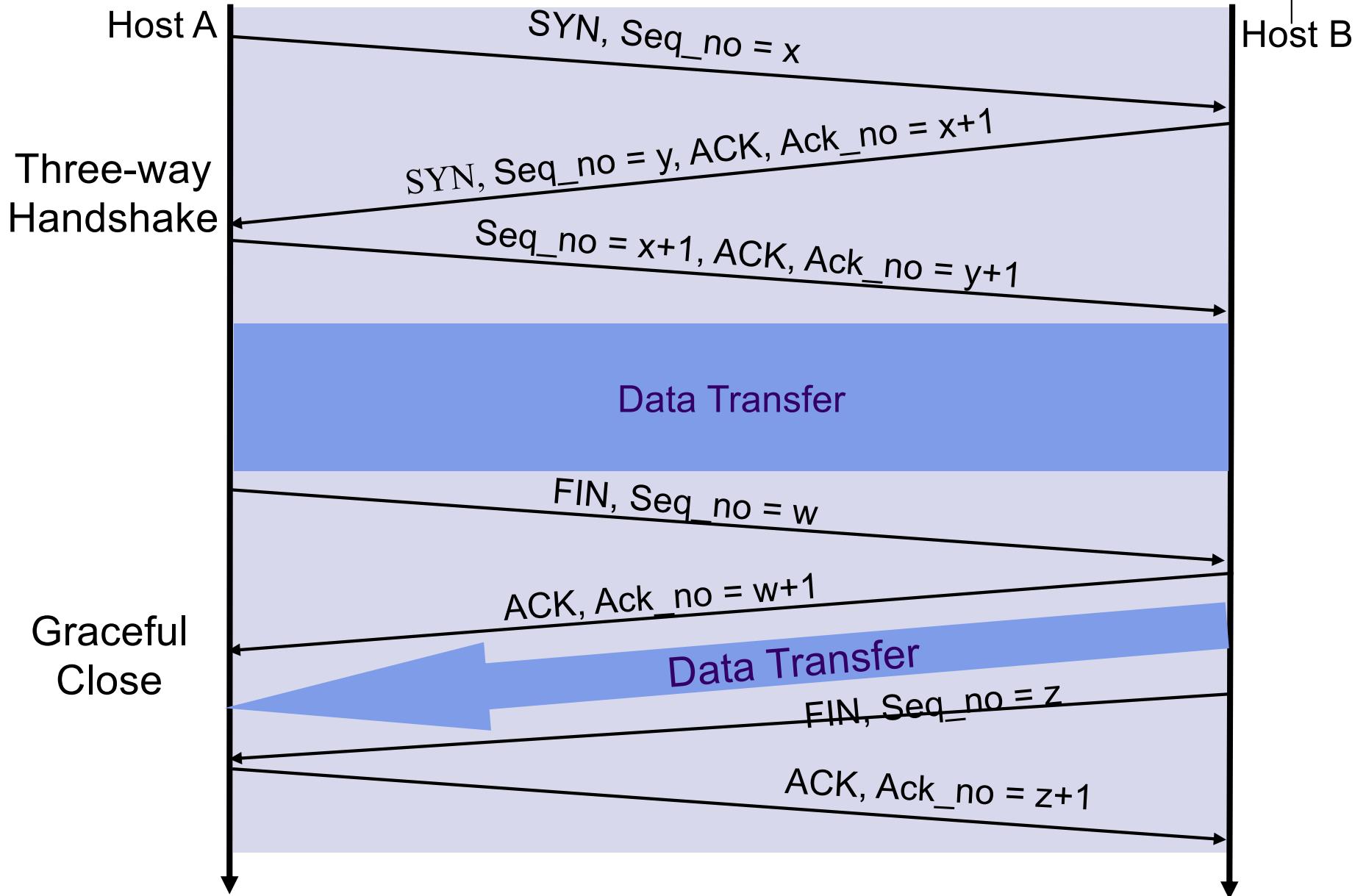


# TCP Connections

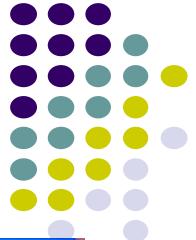
- TCP Connection
  - One connection each way
  - Each connection is identified uniquely by:  
*Send IP Address, Send TCP Port #, Receive IP Address, Receive TCP Port #*
- Connection Setup with Three-Way Handshake
  - Three-way exchange to negotiate initial Seq. #'s for connections in each direction
- Data Transfer
  - Exchange segments carrying data
- Graceful Close
  - Close each direction separately



# Three Phases of TCP Connection



# 1st Handshake: Client-Server Connection Request



→

**TCP Telnet Capture - Ethereal**

File Edit Capture Display Tools Help

No.	Time	Source	Destination	Protocol	Info
1	0.000000	65.95.113.77	128.113.26.22	TCP	2743 > telnet [SYN] Seq=1839733355 Ack=0 win=31988 Len=0
2	0.144934	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [SYN, ACK] Seq=1877388864 Ack=1839733356 wi
3	0.145270	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733356 Ack=1877388865 win=319
4	0.322432	128.113.26.22	65.95.113.77	TELNET	Telnet Data ...
5	0.323617	65.95.113.77	128.113.26.22	TELNET	Telnet Data ...
6	21.606250	65.95.113.77	128.113.26.22	TCP	2743 > telnet [FIN, ACK] Seq=1839733427 Ack=1877389120 wi
7	21.751944	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [ACK] Seq=1877389120 Ack=1839733428 win=491
8	21.757136	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [FIN, ACK] Seq=1877389120 Ack=1839733428 wi
9	21.757468	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733428 Ack=1877389121 win=317

Internet Protocol, Src Addr: 65.95.113.77 (65.95.113.77), Dst Addr: 128.113.26.22 (128.113.26.22)  
Transmission Control Protocol, Src Port: 2743 (2743), Dst Port: telnet (23), Seq: 1839733355, Ack: 0, Len: 0  
Source port: 2743 (2743)  
Destination port: telnet (23)  
Sequence number: 1839733355  
Header length: 28 bytes  
Flags: 0x0002 (SYN)  
0.... .... = Congestion Window Reduced (CWR): Not set  
.0..... = ECN-Echo: Not set  
.00.... = Urgent: Not set  
.000.... = Acknowledgment: Not set  
.000.. = Push: Not set  
.000..0 = Reset: Not set  
.000..1. = Syn: Set  
.000..0 = Fin: Not set  
window size: 31988  
Checksum: 0x2644 (correct)  
options: (8 bytes)

Initial Seq. # from client to server

SYN bit set indicates request to establish connection from client to server

0000 00 90 1a 40 1d 17 00 80 c6 e9 fe 08 88 64 11 00 ...@.... d...  
0010 15 97 00 32 00 21 45 00 00 30 4e 2f 40 00 80 06 ...2.!E. .ON/@...  
0020 5f 65 41 5f 71 4d 80 71 1a 16 0a b7 00 17 6d a8 \_eA\_qM.q .....m.  
0030 1a 6b 00 00 00 00 70 02 7c f4 26 44 00 00 02 04 .k....p. |.&....  
0040 05 86 01 01 04 02 .....  
Filter:    Reset  Apply File: TCP Telnet Capture

# 2<sup>nd</sup> Handshake: ACK from Server



**E TCP Telnet Capture - Ethereal**

No. Time Source Destination Protocol Info

1	0.000000	65.95.113.77	128.113.26.22	TCP	2743 > telnet [SYN] Seq=1839733355 Ack=0 win=31988 Len=0
2	0.144934	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [SYN, ACK] Seq=1877388864 Ack=1839733356 Win=31988
3	0.145270	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733356 Ack=1877388865 Win=31988
4	0.322432	128.113.26.22	65.95.113.77	TELNET	Telnet Data ...
5	0.323617	65.95.113.77	128.113.26.22	TELNET	Telnet Data ...
6	21.606250	65.95.113.77	128.113.26.22	TCP	2743 > telnet [FIN, ACK] Seq=1839733427 Ack=1877389120 Win=317
7	21.751944	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [ACK] Seq=1877389120 Ack=1839733428 Win=491
8	21.757136	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [FIN, ACK] Seq=1877389120 Ack=1839733428 Win=491
9	21.757468	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733428 Ack=1877389121 Win=317

Internet Protocol, src Addr: 128.113.26.22 (128.113.26.22), Dst Addr: 65.95.113.77 (65.95.113.77)  
Transmission Control Protocol, Src Port: telnet (23), Dst Port: 2743 (2743), Seq: 1877388864, Ack: 1839733356  
Source port: telnet (23)  
Destination port: 2743 (2743)  
Sequence number: 1877388864  
Acknowledgement number: 1839733356  
Header length: 24 bytes  
Flags: 0x0012 (SYN, ACK)  
0.... .... = Congestion Window Reduced (CWR): Not set  
.0.... .... = ECN-Echo: Not set  
.0.... .... = Urgent: Not set  
.1.... .... = Acknowledgment: Set  
.... 0... = Push: Not set  
.... .0.. = Reset: Not set  
.... ..1. = Syn: Set  
.... ...0 = Fin: Not set  
window size: 49152  
Checksum: 0xd9d8 (correct)  
Options: (4 bytes)

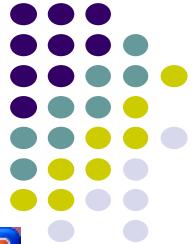
ACK Seq. # = Init.  
Seq. # + 1

ACK bit set acknowledges connection request; Client-to-Server connection established

0000 00 80 c6 e9 fe 08 00 90 1a 40 1d 17 88 64 11 00  
0010 15 97 00 2e 00 21 45 00 00 2c c9 1c 40 00 33 06  
0020 31 7c 80 71 1a 16 41 5f 71 4d 00 17 0a b7 6f e6  
0030 ae 40 6d a8 1a 6c 60 12 c0 00 d9 d8 00 00 02 04  
0040 05 b4

Filter:     File: TCP Telnet Capture

# 2nd Handshake: Server-Client Connection Request



**TCP Telnet Capture - Ethereal**

No. Time Source Destination Protocol Info

1	0.000000	65.95.113.77	128.113.26.22	TCP	2743 > telnet [SYN] Seq=1839733355 Ack=0 win=31988 Len=0
2	0.144934	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [SYN, ACK] Seq=1877388864 Ack=1839733356 wi
3	0.145270	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733356 Ack=1877388865 win=319
4	0.322432	128.113.26.22	65.95.113.77	TELNET	Telnet Data ...
5	0.323617	65.95.113.77	128.113.26.22	TELNET	Telnet Data ...
6	21.606250	65.95.113.77	128.113.26.22	TCP	2743 > telnet [FIN, ACK] Seq=1839733427 Ack=1877389120 wi
7	21.751944	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [ACK] Seq=1877389120 Ack=1839733428 win=491
8	21.757136	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [FIN, ACK] Seq=1877389120 Ack=1839733428 wi
9	21.757468	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733428 Ack=1877389121 win=317

Internet Protocol, Src Addr: 128.113.26.22 (128.113.26.22), Dst Addr: 65.95.113.77 (65.95.113.77)  
Transmission Control Protocol, Src Port: telnet (23) Dst Port: 2743 (2743), Seq: 1877388864, Ack: 1839733356  
Source port: telnet (23)  
Destination port: 2743 (2743)  
Sequence number: 1877388864  
Acknowledgement number: 1839733356  
Header length: 24 bytes  
Flags: 0x0012 (SYN, ACK)  
.... .... = Congestion Window Reduced (CWR): Not set  
.0. .... = ECN-Echo: Not set  
.0. .... = Urgent: Not set  
.1 .... = Acknowledgment: Set  
.0... = Push: Not set  
.0... = Reset: Not set  
.1. = Syn: Set  
.0 = Fin: Not set  
window size: 49152  
Checksum: 0xd9d8 (correct)  
Options: (4 bytes)

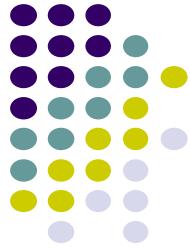
Initial Seq. # from server to client

SYN bit set indicates request to establish connection from server to client

0000 00 80 c6 e9 fe 08 00 90 1a 40 1d 17 88 64 11 00 ..... @..d..  
0010 15 97 00 2e 00 21 45 00 00 2c c9 1c 40 00 33 06 .....!E. ....@.3.  
0020 31 7c 80 71 1a 16 41 5f 71 4d 00 17 0a b7 6f e6 1|.q..A\_ qM....o.  
0030 ae 40 6d a8 1a 6c 60 12 c0 00 d9 d8 00 00 02 04 .@m..1' ..  
0040 05 b4 ..

Filter:    File: TCP Telnet Capture

# 3<sup>rd</sup> Handshake: ACK from Client



→

**TCP Telnet Capture - Ethereal**

File Edit Capture Display Tools Help

No.	Time	Source	Destination	Protocol	Info
1	0.000000	65.95.113.77	128.113.26.22	TCP	2743 > telnet [SYN] Seq=1839733355 Ack=0 win=31988 Len=0
2	0.144934	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [SYN, ACK] Seq=1877388864 Ack=1839733356 win=491
3	0.145270	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733356 Ack=1877388865 win=31988 Len=0
4	0.322432	128.113.26.22	65.95.113.77	TELNET	Telnet Data ...
5	0.323617	65.95.113.77	128.113.26.22	TELNET	Telnet Data ...
6	21.606250	65.95.113.77	128.113.26.22	TCP	2743 > telnet [FIN, ACK] Seq=1839733427 Ack=1877389120 win=317
7	21.751944	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [ACK] Seq=1877389120 Ack=1839733428 win=49152 Len=0
8	21.757136	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [FIN, ACK] Seq=1877389120 Ack=1839733428 win=491
9	21.757468	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733428 Ack=1877389121 win=31733 Len=0

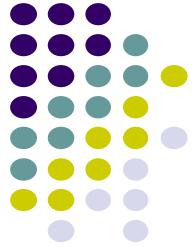
Internet Protocol, Src Addr: 65.95.113.77 (65.95.113.77), Dst Addr: 128.113.26.22 (128.113.26.22)  
Transmission Control Protocol, Src Port: 2743 (2743), Dst Port: telnet (23), seq: 1839733356, Ack: 1877388865, Len: 0  
Source port: 2743 (2743)  
Destination port: telnet (23)  
Sequence number: 1839733356  
Acknowledgement number: 1877388865  
Header length: 20 bytes  
Flags: 0x0010 (ACK)  
0... .... = Congestion Window Reduced (CWR): Not set  
.0. .... = ECN-Echo: Not set  
.0. .... = Urgent: Not set  
.1 .... = Acknowledgment: Set  
.... 0.. = Push: Not set  
.... .0.. = Reset: Not set  
.... ..0. = Syn: Not set  
.... ...0 = Fin: Not set  
window size: 31988  
Checksum: 0x34a2 (correct)

ACK Seq. # = Init.  
Seq. # + 1

ACK bit set acknowledges connection request; Connections in both directions established

0000 00 90 1a 40 1d 17 00 80 c6 e9 fe 08 88 64 11 00  
0010 15 97 00 2a 00 21 45 00 00 28 4e 30 40 00 80 06  
0020 5f 6c 41 5f 71 4d 80 71 1a 16 0a b7 00 17 6d a8  
0030 1a 6c 6f e6 ae 41 50 10 7c f4 34 a2 00 00

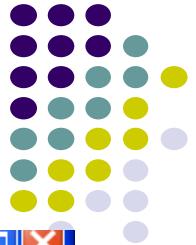
Filter:  Reset Apply File: TCP Telnet Capture



# TCP Data Exchange

- Application Layers write bytes into buffers
- TCP sender forms segments
  - When bytes exceed threshold or timer expires
  - Upon PUSH command from applications
  - Consecutive bytes from buffer inserted in payload
  - Sequence # & ACK # inserted in header
  - Checksum calculated and included in header
- TCP receiver
  - Performs selective repeat ARQ functions,  $R_{next}$
  - Writes error-free, in-sequence bytes to receive buffer

# Data Transfer: Server-to-Client Segment



→

**TCP Telnet Capture - Ethereal**

File Edit Capture Display Tools Help

No.	Time	Source	Destination	Protocol	Info
1	0.000000	65.95.113.77	128.113.26.22	TCP	2743 > telnet [SYN] Seq=1839733355 Ack=0 Win=31988 Len=0
2	0.144934	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [SYN, ACK] Seq=1877388864 Ack=1839733356 Win=49152 L
3	0.145270	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733356 Ack=1877388865 Win=31988 Len=0
4	0.322432	128.113.26.22	65.95.113.77	TELNET	Telnet Data ...
5	0.323617	65.95.113.77	128.113.26.22	TELNET	Telnet Data ...
6	21.606250	65.95.113.77	128.113.26.22	TCP	2743 > telnet [FIN, ACK] Seq=1839733427 Ack=1877389120 Win=31733 L
7	21.751944	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [ACK] Seq=1877389120 Ack=1839733428 Win=49152 Len=0
8	21.757136	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [FIN, ACK] Seq=1877389120 Ack=1839733428 Win=49152 L
9	21.757468	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733428 Ack=1877389121 Win=31733 Len=0

Internet Protocol, Src Addr: 128.113.26.22 (128.113.26.22), Dst Addr: 65.95.113.77 (65.95.113.77)  
Transmission Control Protocol, Src Port: telnet (23), Dst Port: 2743 (2743), Seq: 1877388865, Ack: 1839733356, Len: 12  
Source port: telnet (23)  
destination port: 2743 (2743)

Sequence number: 1877388865  
Next sequence number: 1877388877  
Acknowledgement number: 1839733356  
Header Length: 20 bytes  
Flags: 0x0018 (PSH, ACK)  
0.... .... = Congestion Window Reduced (CWR): Not set  
.0.... .... = ECN-Echo: Not set  
.0.... .... = Urgent: Not set  
.1.... .... = Acknowledgment: Set  
.... 1.... = Push: Set  
.... .0... = Reset: Not set  
.... .0... = Syn: Not set  
.... .0... = Fin: Not set  
Window size: 49152  
Checksum: 0xba41 (correct)

**Push set**

12 bytes of payload

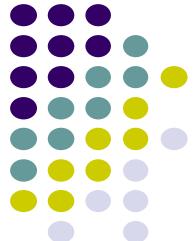
Telnet

Command: Do Terminal Type  
Command: Do Terminal Speed  
Command: Do X Display Location  
Command: Do Environment Option

12 bytes of payload carries telnet option negotiation

0010 15 97 00 36 00 21 45 00 00 34 c9 2b 40 00 33 06 ...6.!E. .4.+@.3.  
0020 21 65 90 71 1> 16 41 5f 71 4d 00 17 02 b7 6f 06 1e.q..A\_ qM....0.  
0030 ae 41 6d a8 1a 6c 50 18 c0 00 ba 41 00 00 ff fd .Am.1P. ....A....  
0040 18 ff fd 20 ff fd 23 ff fd 24 ...#..\$.  
Filter:  Reset  Apply Telnet (telnet), 12 bytes

# Graceful Close: Client-to-Server Connection



**TCP Telnet Capture - Ethereal**

File Edit Capture Display Tools Help

No.	Time	Source	Destination	Protocol	Info
1	0.000000	65.95.113.77	128.113.26.22	TCP	2743 > telnet [SYN] Seq=1839733355 Ack=0 Win=31988 Len=0
2	0.144934	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [SYN, ACK] Seq=1877388864 Ack=1839733356 win=49152 Len=0
3	0.145270	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733356 Ack=1877388865 win=31988 Len=0
4	0.322432	128.113.26.22	65.95.113.77	TELNET	Telnet Data ...
5	0.323617	65.95.113.77	128.113.26.22	TELNET	Telnet Data ...
6	21.606250	65.95.113.77	128.113.26.22	TCP	2743 > telnet [FIN, ACK] Seq=1839733427 Ack=1877389120 win=31733 Len=0
7	21.751944	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [ACK] Seq=1877389120 Ack=1839733428 win=49152 Len=0
8	21.757136	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [FIN, ACK] Seq=1877389120 Ack=1839733428 win=49152 Len=0
9	21.757468	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733428 Ack=1877389121 win=31733 Len=0

Internet Protocol, Src Addr: 65.95.113.77 (65.95.113.77), Dst Addr: 128.113.26.22 (128.113.26.22)  
Transmission Control Protocol, src Port: 2743 (2743), Dst Port: telnet (23), Seq: 1839733427, Ack: 1877389120, Len: 0

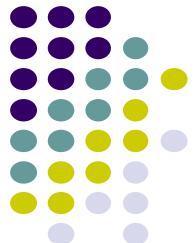
Source port: 2743 (2743)  
Destination port: telnet (23)  
Sequence number: 1839733427  
Acknowledgement number: 1877389120  
Header length: 20 bytes  
Flags: 0x0011 (FIN, ACK)  
0... .... = Congestion Window Reduced (CWR): Not set  
.0... .... = ECN-Echo: Not set  
.0. .... = Urgent: Not set  
.1 .... = Acknowledgment: Set  
.0... = Push: Not set  
.0.. = Reset: Not set  
.0. = Syn: Not set  
.1 = Fin: Set  
Window size: 31733  
Checksum: 0x345a (correct)

Client initiates closing of its connection to server

0000 00 90 1a 40 1d 17 00 80 c6 e9 fe 08 88 64 11 00 ...@.... ....d..  
0010 15 97 00 2a 00 21 45 00 00 28 4e 55 40 00 80 06 ...\*!.E. .(NU@...  
0020 5f 47 41 5f 71 4d 80 71 1a 16 0a b7 00 17 6d a8 \_GA\_qM.q .....m.  
0030 1a b3 6f e6 af 40 50 11 7b f5 34 5a 00 00 ..o..@P. {.4Z..

Filter: |   File: TCP Telnet Capture

# Graceful Close: Client-to-Server Connection



**TCP Telnet Capture - Ethereal**

File Edit Capture Display Tools Help

No.	Time	Source	Destination	Protocol	Info
1	0.000000	65.95.113.77	128.113.26.22	TCP	2743 > telnet [SYN] Seq=1839733355 Ack=0 Win=31988 Len=0
2	0.144934	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [SYN, ACK] Seq=1877388864 Ack=1839733356 Win=49152 Len=0
3	0.145270	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733356 Ack=1877388865 Win=31988 Len=0
4	0.322432	128.113.26.22	65.95.113.77	TELNET	Telnet Data ...
5	0.323617	65.95.113.77	128.113.26.22	TELNET	Telnet Data ...
6	21.606250	65.95.113.77	128.113.26.22	TCP	2743 > telnet [FIN, ACK] Seq=1839733427 Ack=1877389120 Win=31733 Len=0
7	21.751944	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [ACK] Seq=1877389120 Ack=1839733428 Win=49152 Len=0
8	21.757136	128.113.26.22	65.95.113.77	TCP	telnet > 2743 [FIN, ACK] Seq=1877389120 Ack=1839733428 Win=49152 Len=0
9	21.757468	65.95.113.77	128.113.26.22	TCP	2743 > telnet [ACK] Seq=1839733428 Ack=1877389121 Win=31733 Len=0

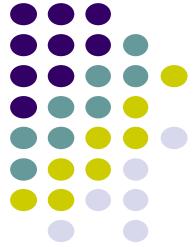
Internet Protocol, Src Addr: 128.113.26.22 (128.113.26.22), Dst Addr: 65.95.113.77 (65.95.113.77)  
Transmission Control Protocol, Src Port: telnet (23), Dst Port: 2743 (2743), seq: 1877389120, Ack: 1839733428, Len: 0  
source port: telnet (23)  
Destination port: 2743 (2743)  
Sequence number: 1877389120  
Acknowledgement number: 1839733428  
Header length: 20 bytes  
Flags: 0x0010 (ACK)  
0... .... = Congestion Window Reduced (CWR): Not set  
.0... .... = ECN-Echo: Not set  
.0. .... = Urgent: Not set  
.1 .... = Acknowledgment: Set  
.0... = Push: Not set  
.0.. = Reset: Not set  
.0.0. = Syn: Not set  
.0..0 = Fin: Not set  
window size: 49152  
Checksum: 0xf04e (correct)

ACK Seq. # = Previous Seq. # + 1

Server ACKs request; client-to-server connection closed

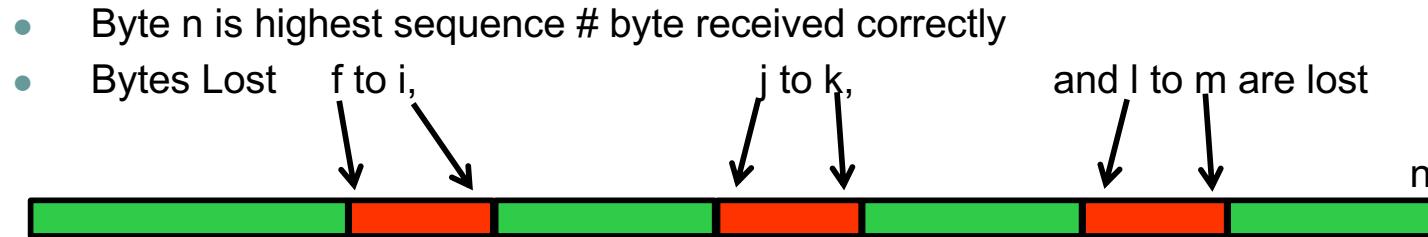
0000 00 80 c6 e9 fe 08 00 90 1a 40 1d 17 88 64 11 00 .@...d..  
0010 15 97 00 2a 00 21 45 00 00 28 c9 81 40 00 33 06 ..!\*!E. @.3.  
0020 31 1b 80 71 1a 16 41 5f 71 4d 00 17 0a b7 6f e6 1..q..A\_ qM....o.  
0030 af 40 6d a8 1a b4 50 10 c0 00 f0 4e 00 00 .@m...P. ....N..

Filter:  Reset  Apply File: TCP Telnet Capture



# TCP Selective ACK (RFC 2018)

- SACK option allows selective retransmission
- Assume

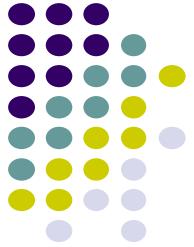


Receiver sends back:

- CACK=f, cumulative ACK “all bytes up to f-1 received correctly”
- SACK=(i+1 to j-1, k+1 to l-1, m+1 to n)

Sender then knows that:

- blocks f to i, j to k, and l to m are lost and retransmits them
- RFC 2018 uses 2 32-bit words to identify boundaries of up to 4 blocks



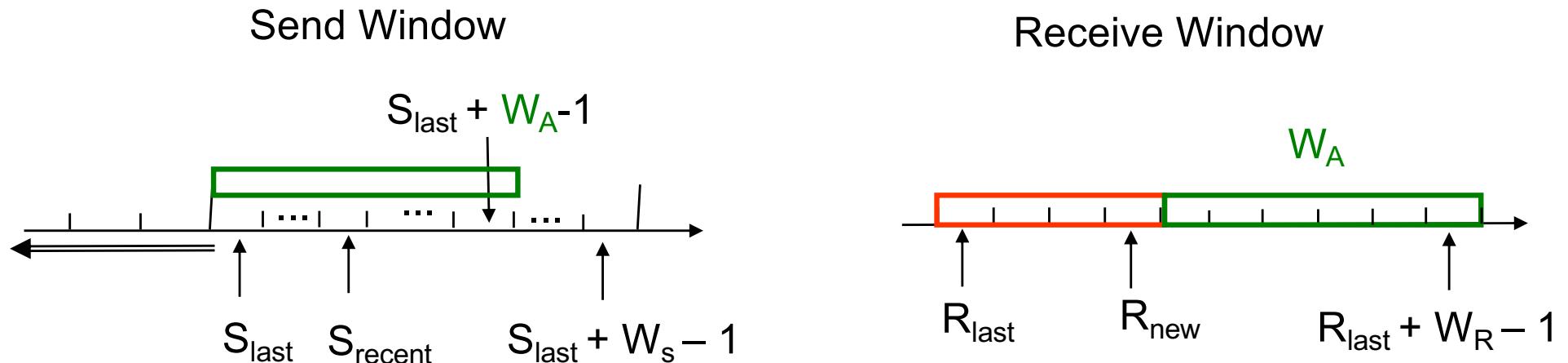
# Flow Control

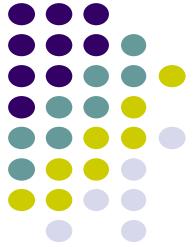
- TCP receiver controls rate at which sender transmits to prevent buffer overflow in the receiver buffer
- TCP receiver advertises a window size specifying number of bytes that can be accommodated by receiver

$$W_A = W_R - (R_{new} - R_{last})$$

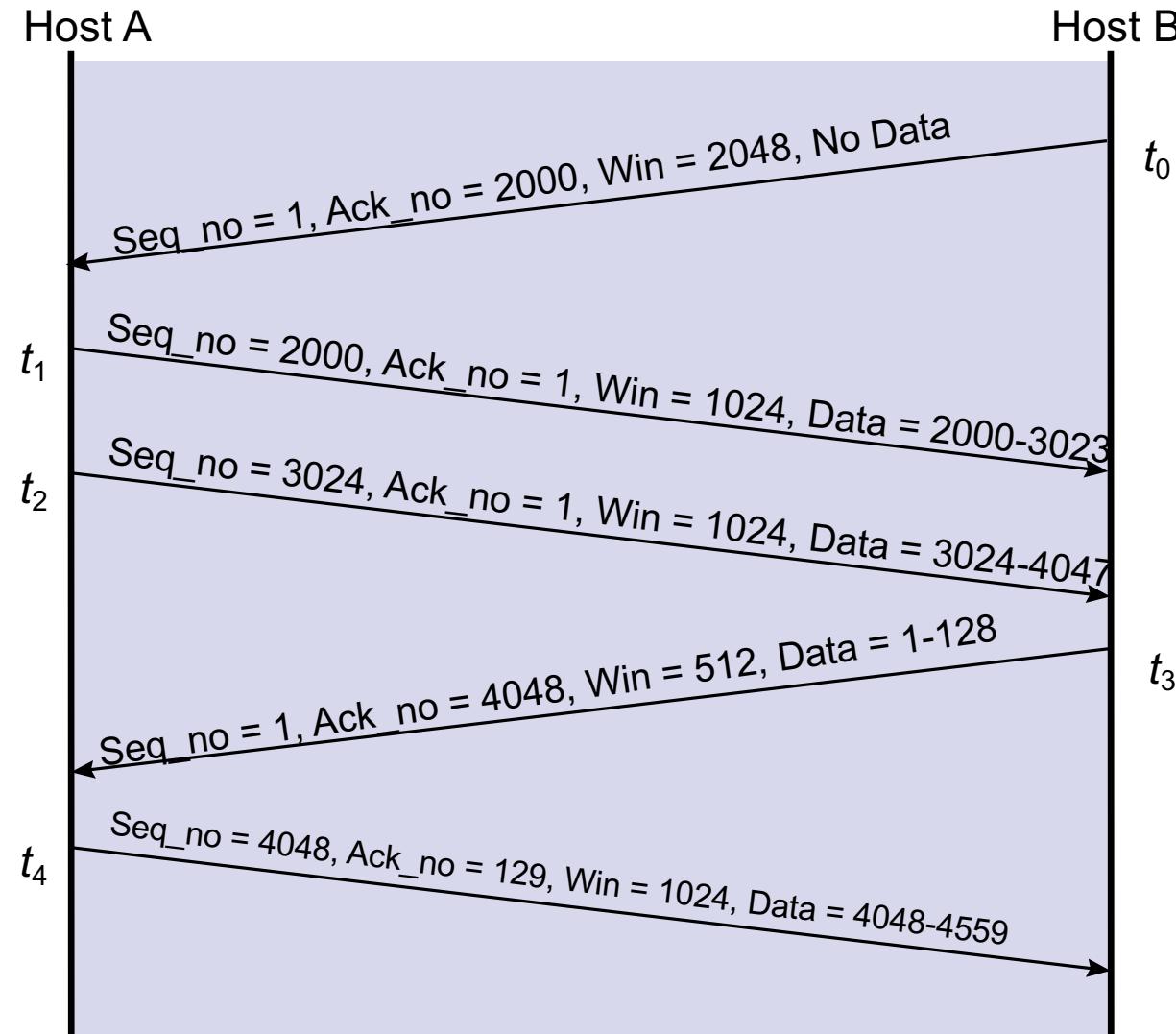
- TCP sender obliged to keep # outstanding bytes below  $W_A$

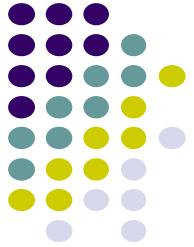
$$(S_{recent} - S_{last}) \leq W_A$$





# TCP window flow control



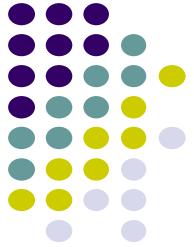


# TCP Retransmission Timeout

- TCP retransmits a segment after timeout period
  - If timeout too short: excessive number of retransmissions
  - If timeout too long: recovery too slow
  - Timeout depends on RTT: time from when segment is sent to when ACK is received
- Round trip time (RTT) in Internet is highly variable
  - Routes vary and can change in mid-connection
  - Traffic fluctuates
- TCP uses adaptive estimation of RTT
  - Measure RTT each time ACK received:  $\tau_n$

$$t_{RTT}(\text{new}) = \alpha t_{RTT}(\text{old}) + (1 - \alpha) \tau_n$$

- $\alpha = 7/8$  typical



# RTT Variability

- Estimate variance  $\sigma^2$  of RTT variation
- Estimate for timeout:

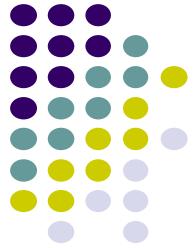
$$t_{out} = t_{RTT} + k \sigma_{RTT}$$

- If RTT highly variable, timeout increase accordingly
- If RTT nearly constant, timeout close to RTT estimate
- Approximate estimation of deviation

$$d_{RTT}(new) = \beta d_{RTT}(old) + (1-\beta) | \tau_n - t_{RTT} |$$

$$t_{out} = t_{RTT} + 4 d_{RTT}$$

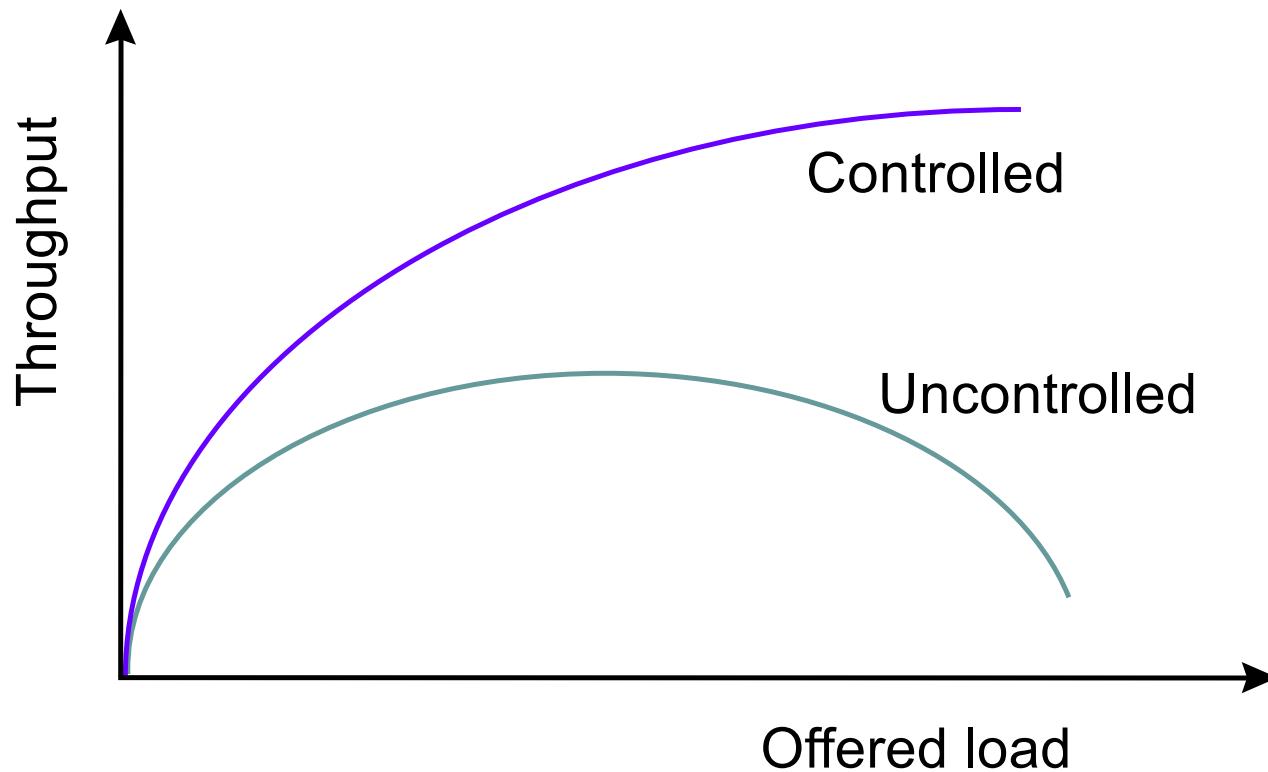
# Congestion Control



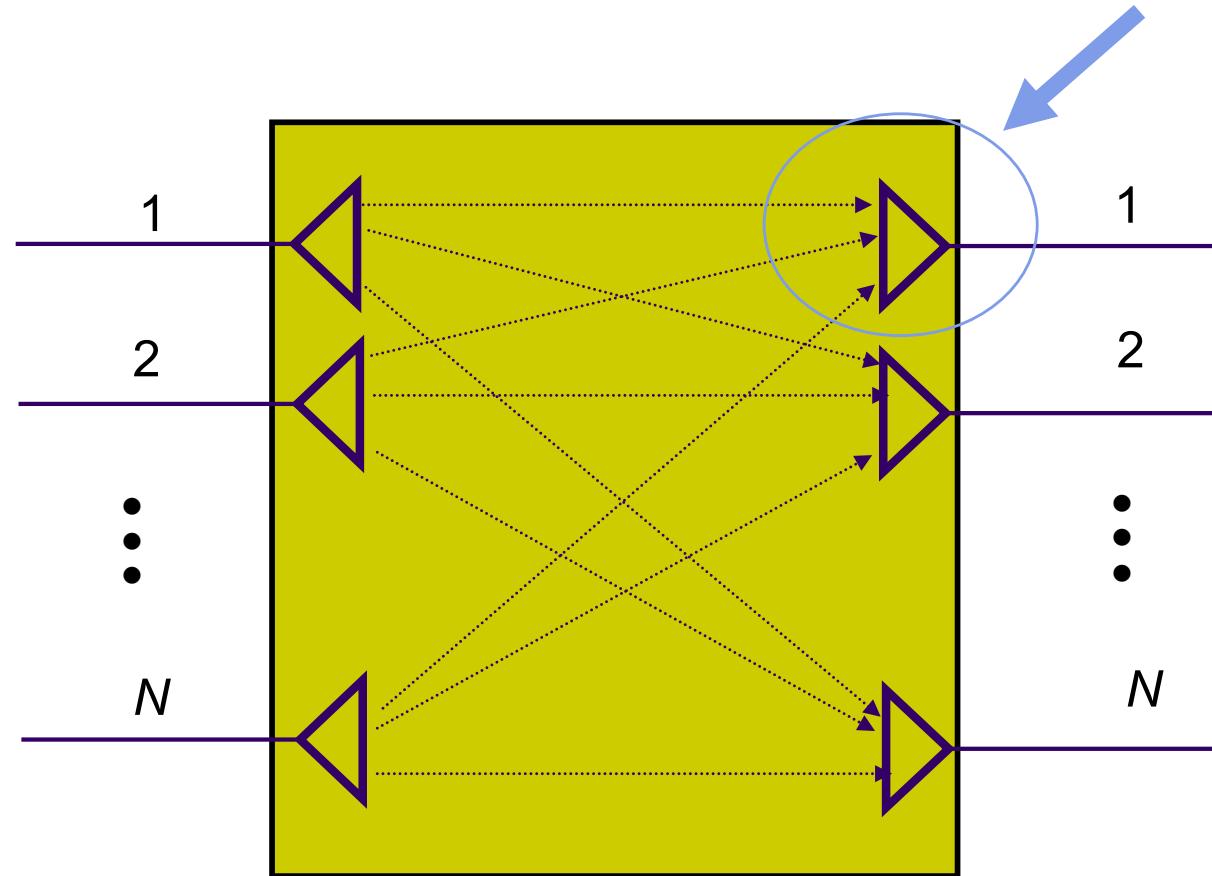
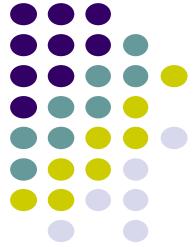
Ideal effect of congestion control:

Resources used efficiently up to capacity available

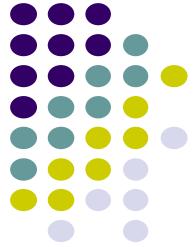
Throughput (Bytes/second) increases steadily without drop



# Multiplexers inherent in Packet Switches

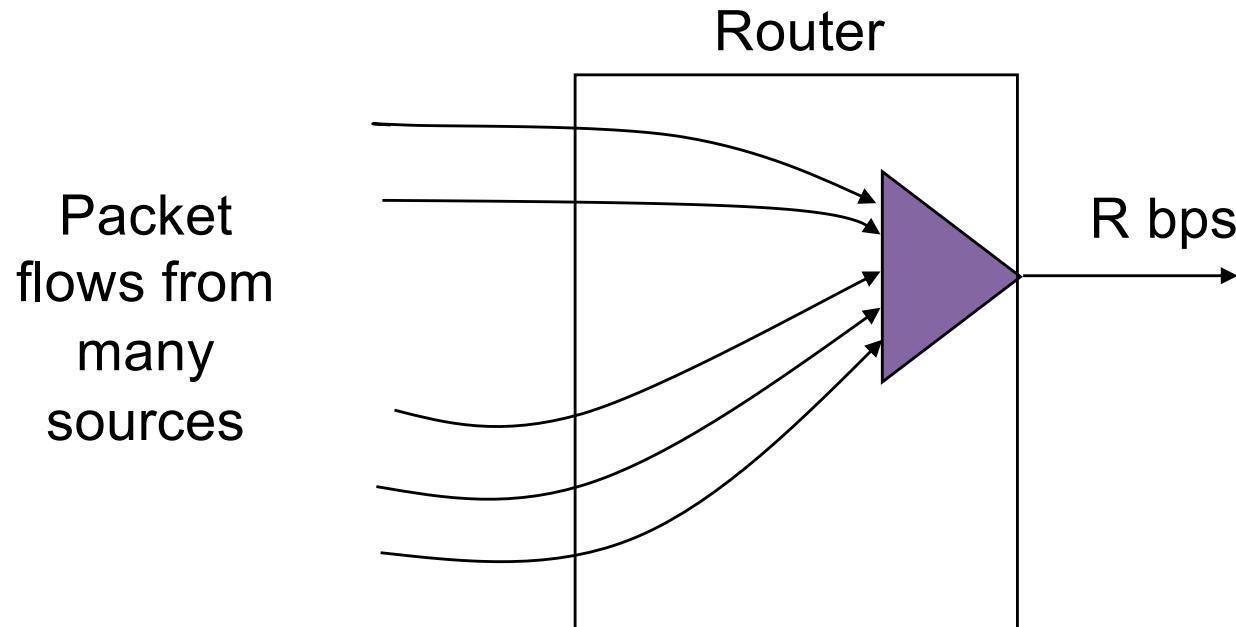


- Packets/frames forwarded to buffer prior to transmission from switch
- Multiplexing occurs in these buffers

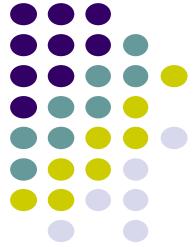


# TCP Congestion Control

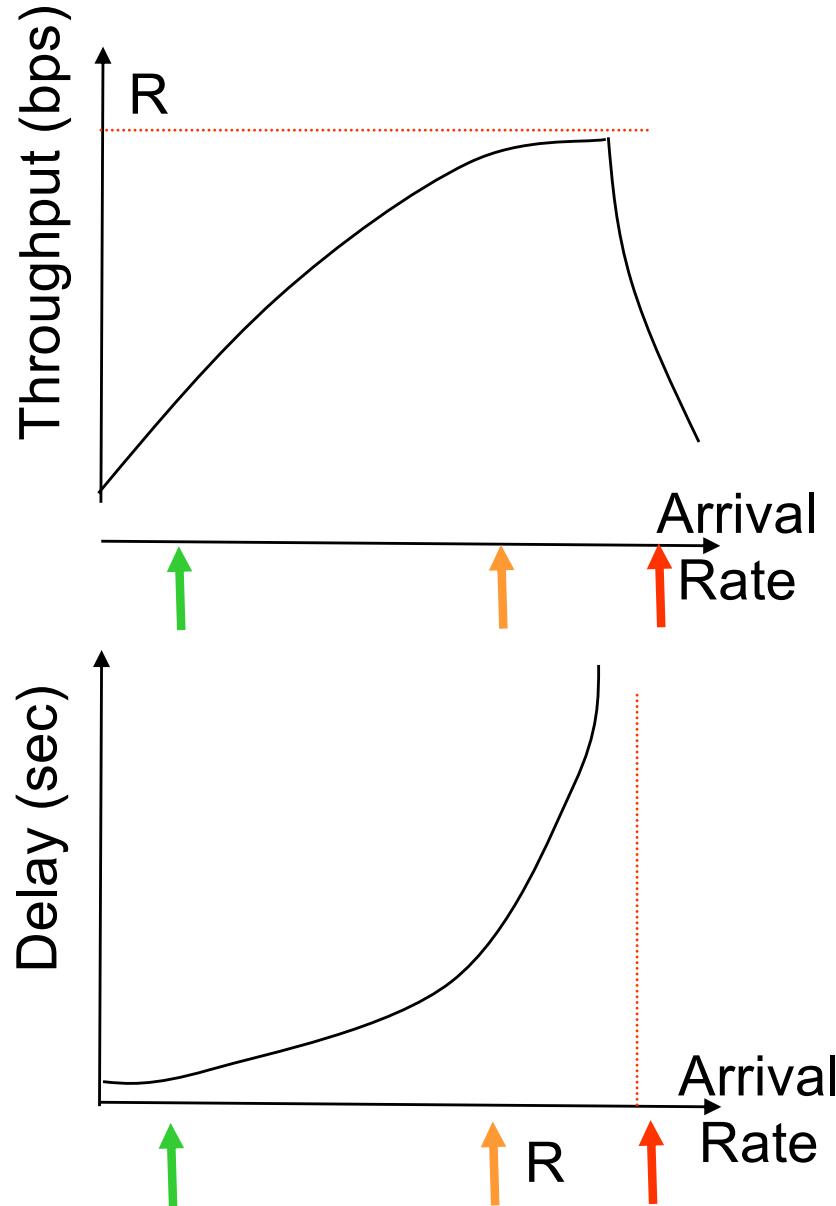
- *Advertised window size* is used to ensure that receiver's buffer will not overflow
- However, buffers at intermediate routers between source and destination may overflow



- Congestion occurs when total arrival rate from all packet flows exceeds  $R$  over a sustained period of time
- Buffers at multiplexer will fill and packets will be lost



# Phases of Congestion Behavior



## Light traffic

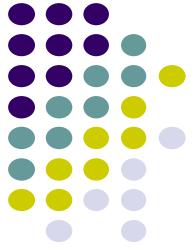
- Arrival Rate  $<< R$
- Low delay
- Can accommodate more

## Knee (congestion onset)

- Arrival rate approaches  $R$
- Delay increases rapidly
- Throughput begins to saturate

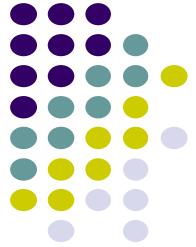
## Congestion Collapse

- Arrival rate  $> R$
- Large delays, packet loss
- Useful application throughput drops



# Window Congestion Control

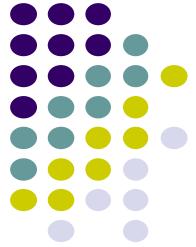
- Desired operating point: just before knee
  - Sources must control their sending rates so that aggregate arrival rate is just before knee
- TCP sender maintains a *congestion window* cwnd to control congestion at intermediate routers
- Effective window is minimum of congestion window and advertised window
- Problem: source does not know what its “fair” share of available bandwidth should be
- Solution: adapt dynamically to available BW
  - Sources probe the network by increasing cwnd
  - When congestion detected, sources reduce rate
  - Ideally, sources sending rate stabilizes near ideal point



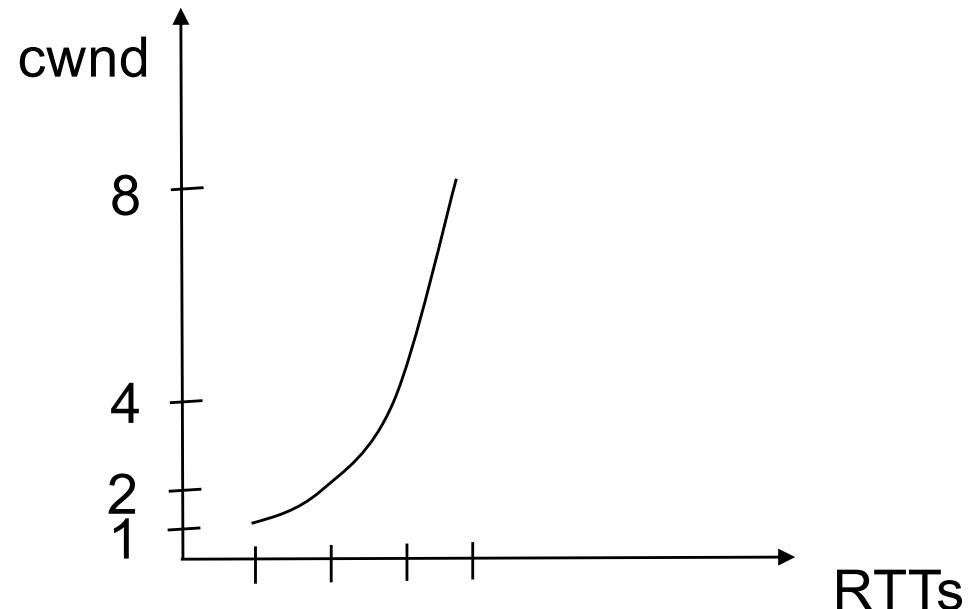
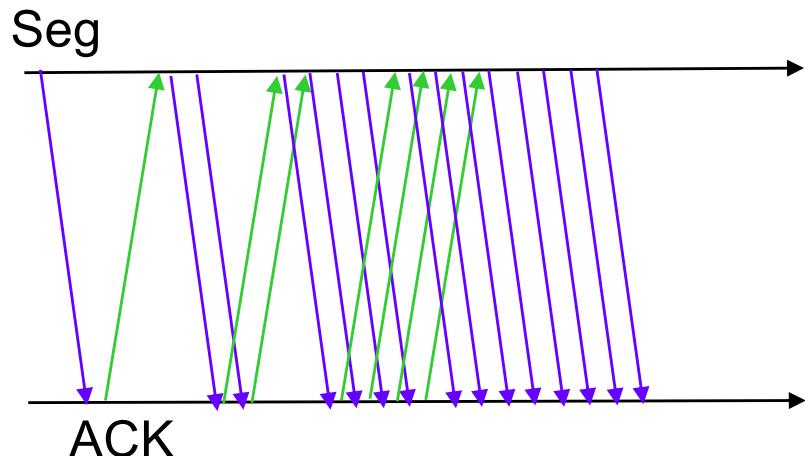
# Congestion Window

- How does the TCP congestion algorithm change congestion window dynamically according to the most up-to-date state of the network?
- At light traffic: each segment is ACKed quickly
  - Increase cwnd aggressively
- At knee: segment ACKs arrive, but more slowly
  - Slow down increase in cwnd
- At congestion: segments encounter large delays (so retransmission timeouts occur); segments are dropped in router buffers (resulting in duplicate ACKs)
  - Reduce transmission rate, then probe again

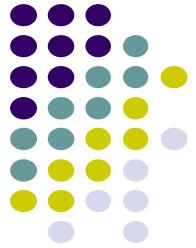
# TCP Congestion Control: Slow Start



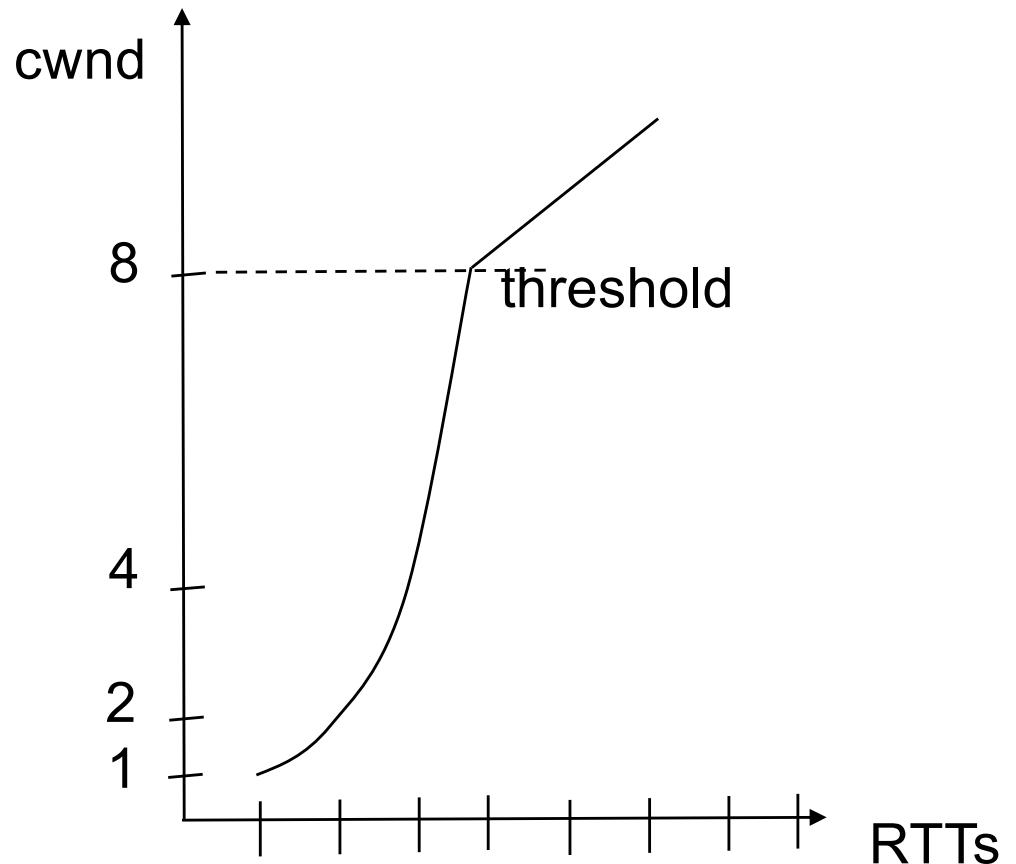
- **Slow start:** increase congestion window size by one segment upon receiving an ACK from receiver
  - initialized at 1 segment
  - used at (re)start of data transfer
  - congestion window increases exponentially



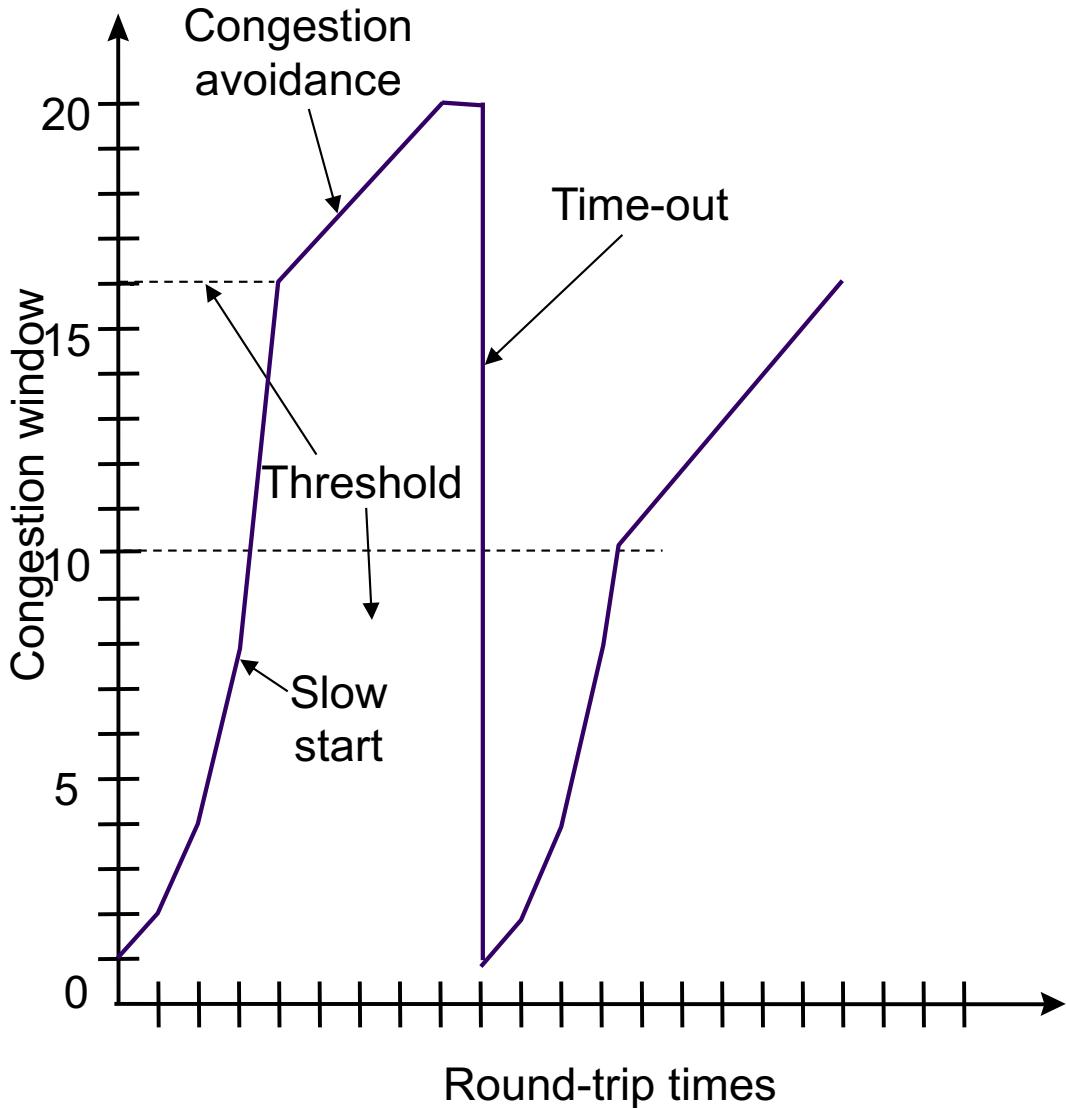
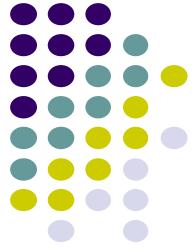
# TCP Congestion Control: Congestion Avoidance



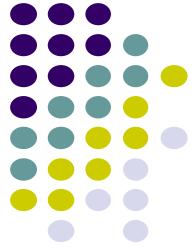
- Algorithm progressively sets a *congestion threshold*
  - When  $cwnd > \text{threshold}$ , slow down rate at which  $cwnd$  is increased
- Increase congestion window size by one segment per round-trip-time (RTT)
  - Each time an ACK arrives,  $cwnd$  is increased by  $1/cwnd$
  - In one RTT,  $cwnd$  segments are sent, so total increase in  $cwnd$  is  $cwnd \times 1/cwnd = 1$
  - $cwnd$  grows linearly with time



# TCP Congestion Control: Congestion

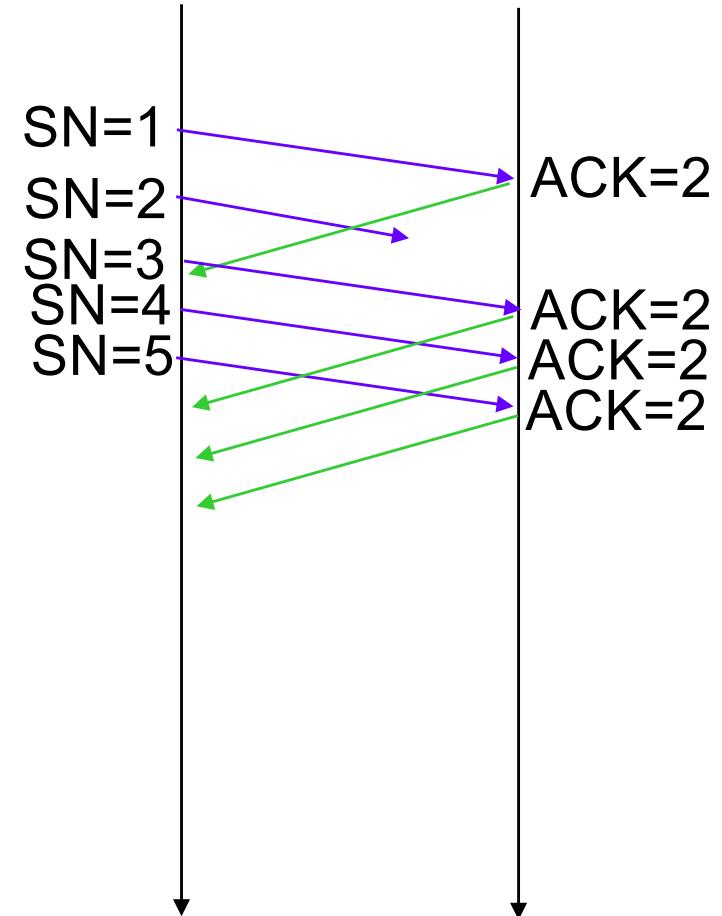


- Congestion is detected upon timeout or receipt of duplicate ACKs
- Assume current cwnd corresponds to available bandwidth
- Adjust congestion threshold =  $\frac{1}{2} \times$  current cwnd
- Reset cwnd to 1
- Go back to slow-start
- Over several cycles expect to converge to congestion threshold equal to about  $\frac{1}{2}$  the available bandwidth

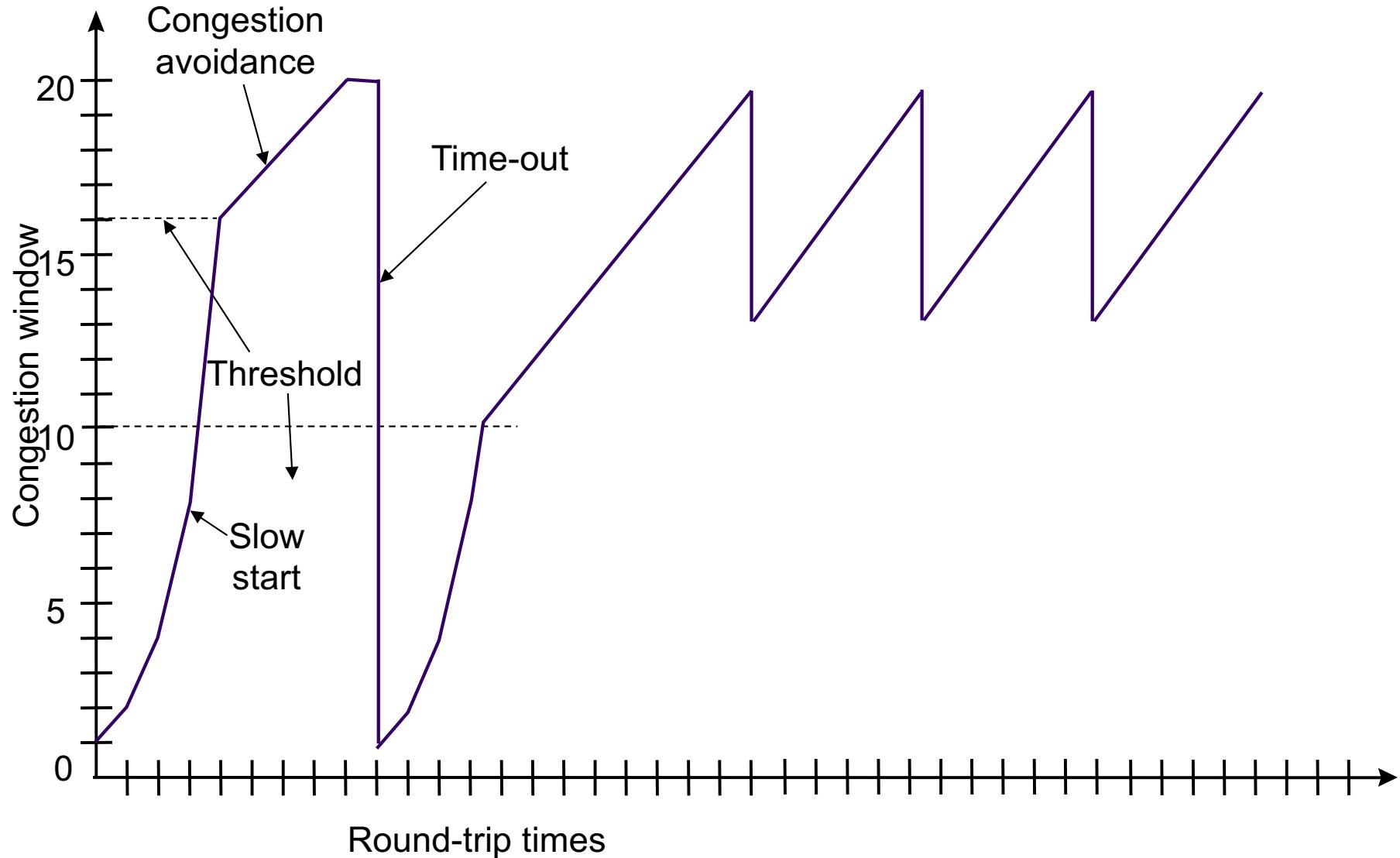
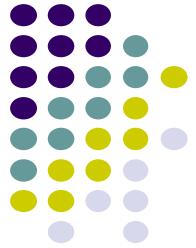


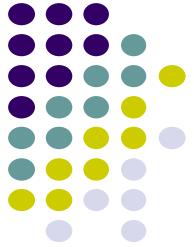
# Fast Retransmit & Fast Recovery

- Congestion causes many segments to be dropped
- If only a single segment is dropped, then subsequent segments trigger duplicate ACKs before timeout
- Can avoid large decrease in cwnd as follows:
  - When three duplicate ACKs arrive, retransmit lost segment immediately
  - Reset congestion threshold to  $\frac{1}{2}$  cwnd
  - Reset cwnd to congestion threshold + 3 to account for the three segments that triggered duplicate ACKs
  - Remain in congestion avoidance phase
  - However if timeout expires, reset cwnd to 1
  - In absence of timeouts, cwnd will oscillate around optimal value



# TCP Congestion Control: Fast Retransmit & Fast Recovery





# Random Early Detection (RED)

- Packets produced by TCP will reduce input rate in response to network congestion
- Early drop: discard packets before buffers are full
- Random drop causes some sources to reduce rate before others, causing gradual reduction in aggregate input rate

Algorithm:

- Maintain running average of queue length
- If  $Q_{avg} < \text{minthreshold}$ , do nothing
- If  $Q_{avg} > \text{maxthreshold}$ , drop packet
- If in between, drop packet according to probability
- Flows that send more packets are more likely to have packets dropped

# Packet Drop Profile in RED

