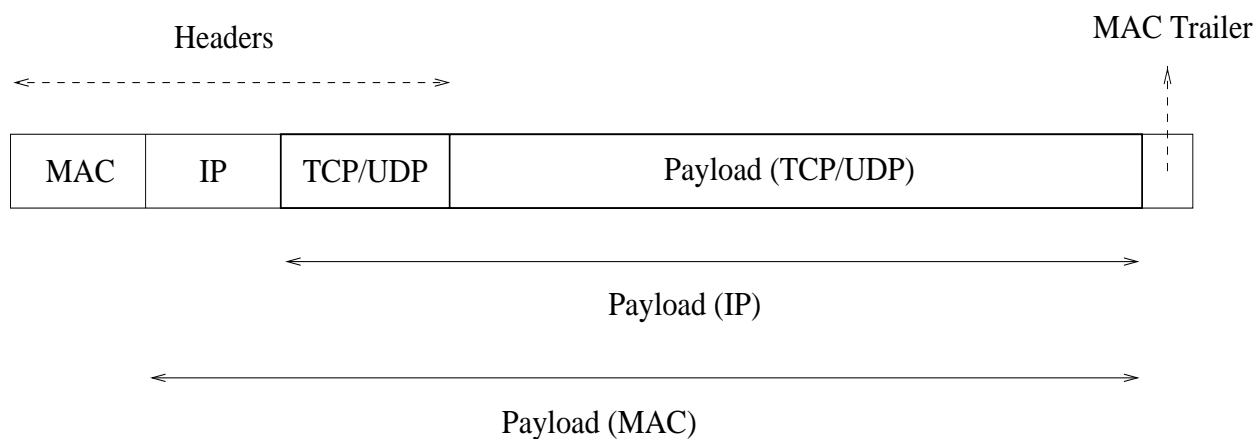


## Transport Protocols: TCP and UDP

- end-to-end protocol
- runs on top of network layer protocols
- treat network layer & below as black box

Three-level encapsulation:



- meaning of protocol “stack”: push/pop headers
- common TCP payload: HTTP

Network layer (IP) assumptions:

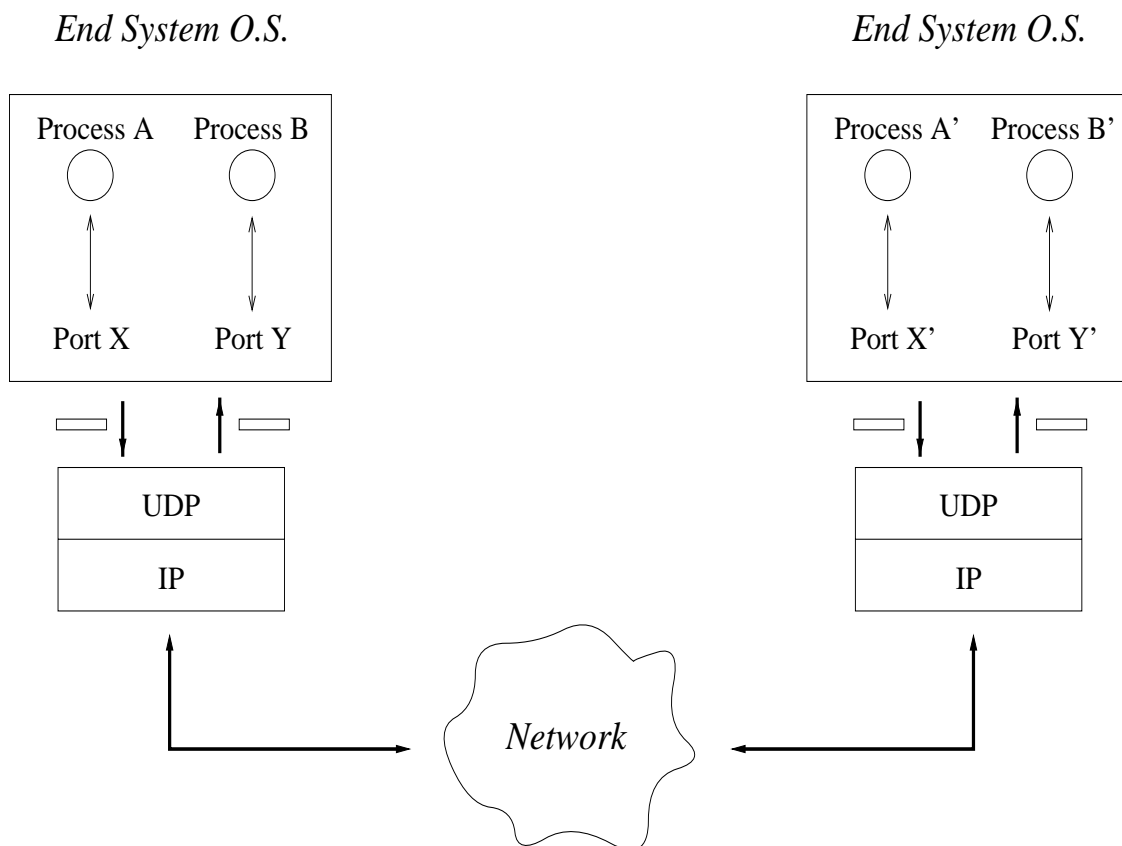
- unreliable
- out-of-order delivery
- absence of QoS guarantees (delay, throughput, etc.)
- insecure (IPv4)
  - IPsec

Additional performance properties:

- Works “ok”
- Can break down under high load conditions
  - Atlanta Olympics
  - DoS and worm attack
- Wide behavioral range
  - sometimes good, so so, or bad

Goal of UDP (User Datagram Protocol):

- process identification
- port number as demux key
- minimal support beyond IP



UDP packet format:

2		2
Source Port		Destination Port
Length		Checksum
Payload		

Checksum calculation (pseudo header):

4		
Source Address		
Destination Address		
00 ... 0	Protocol	UDP Length

→ pseudo header, UDP header and payload

UDP usage:

- multimedia streaming
  - lean and nimble
  - at minimum requires process identification
  - congestion control carried out above UDP
- stateless client/server applications
  - persistent state a hinderance
  - lightweight

Goals of TCP (Transmission Control Protocol):

- process identification
- reliable communication: ARQ
- speedy communication: congestion control
- segmentation
  - connection-oriented, i.e., stateful
  - complex mixture of functionalities

Segmentation task: provide “stream” interface to higher level protocols

—→ exported semantics: contiguous byte stream

—→ recall ARQ

- segment stream of bytes into blocks of fixed size
- segment size determined by TCP MTU (Maximum Transmission Unit)
- actual unit of transmission in ARQ

TCP packet format:

2

2

Source Port								Destination Port	
Sequence Number									
Acknowledgement Number									
Header Length	////	U	A	P	R	S	F	Window Size	
		R	C	S	S	Y	I		
		G	K	H	T	N	N		
Checksum								Urgent Pointer	
Options (if any)									
DATA (if any)									



- Sequence Number: position of first byte of payload
- Acknowledgement: next byte of data expected (receiver)
- Header Length (4 bits): 4 B units
- URG: urgent pointer flag
- ACK: ACK packet flag
- PSH: override TCP buffering
- RST: reset connection
- SYN: establish connection
- FIN: close connection
- Window Size: receiver's advertised window size
- Checksum: prepend pseudo-header
- Urgent Pointer: byte offset in current payload where urgent data begins
- Options: MTU; take min of sender & receiver (default 556 B)

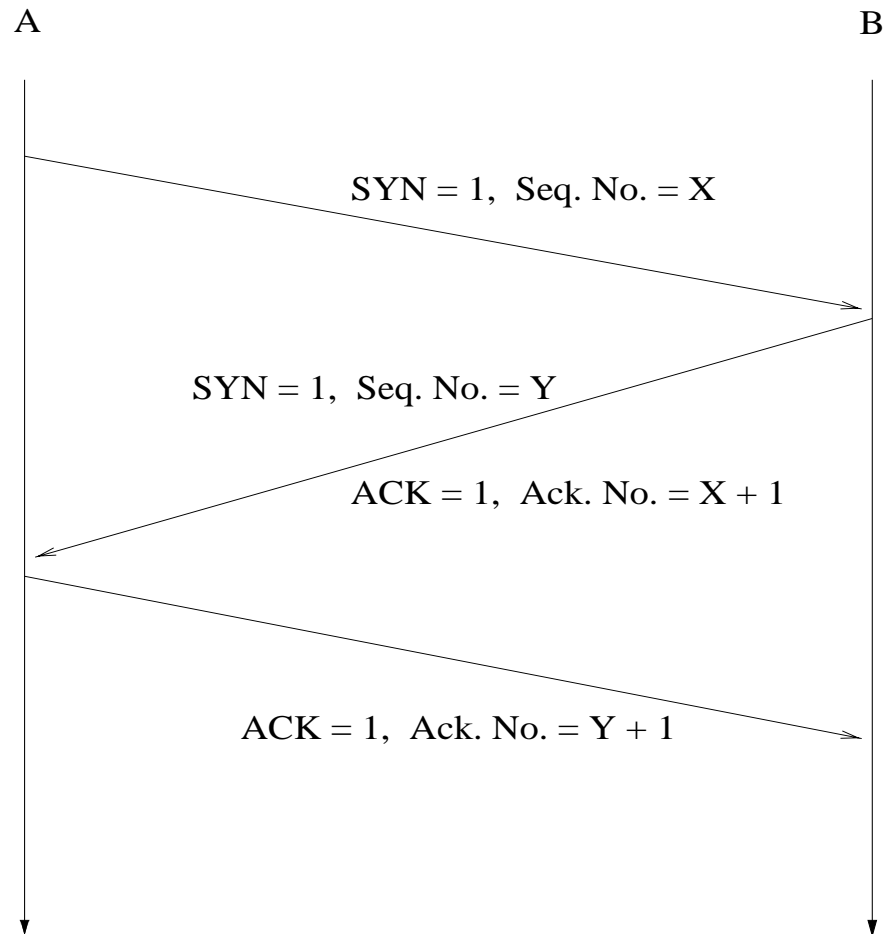
Checksum calculation (pseudo header):

4

Source Address		
Destination Address		
00 ... 0	Protocol	TCP Segment Length

→ pseudo header, TCP header and payload

TCP connection establishment (3-way handshake):

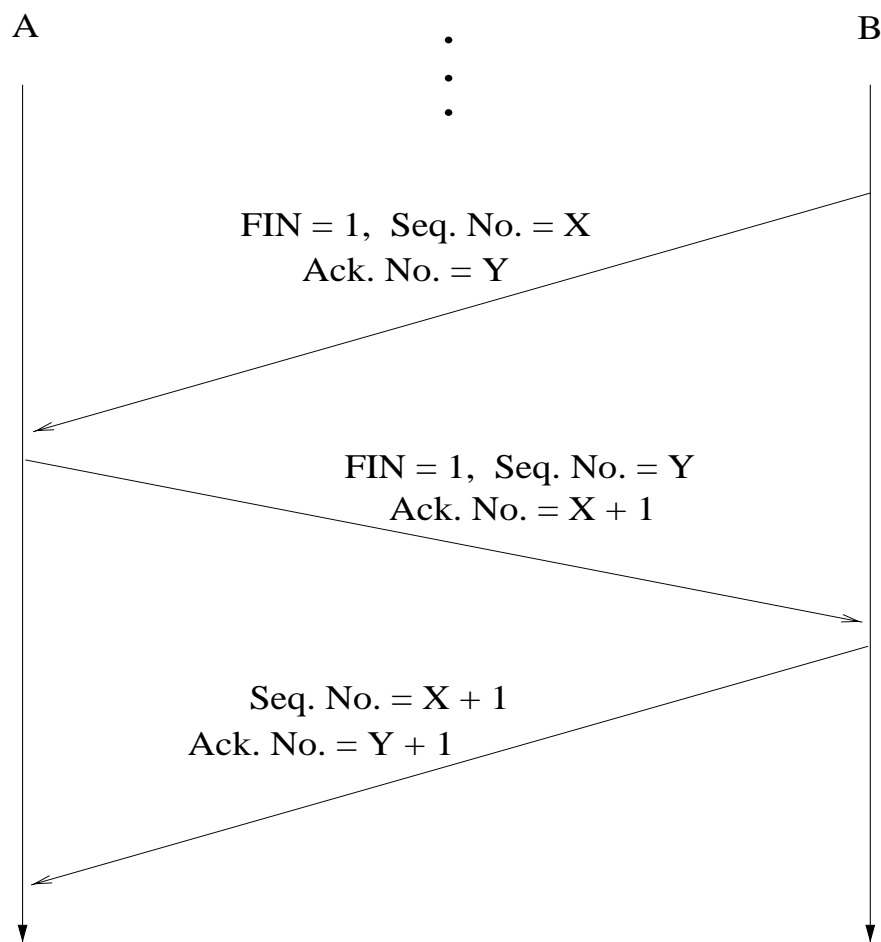


- $X, Y$  are chosen randomly  
→ sequence number prediction
- piggybacking

2-person consensus problem: are  $A$  and  $B$  in agreement about the state of affairs after 3-way handshake?

- in general: impossible
- can be proven
- “acknowledging the ACK problem”
- also TCP session ending
- lunch date problem

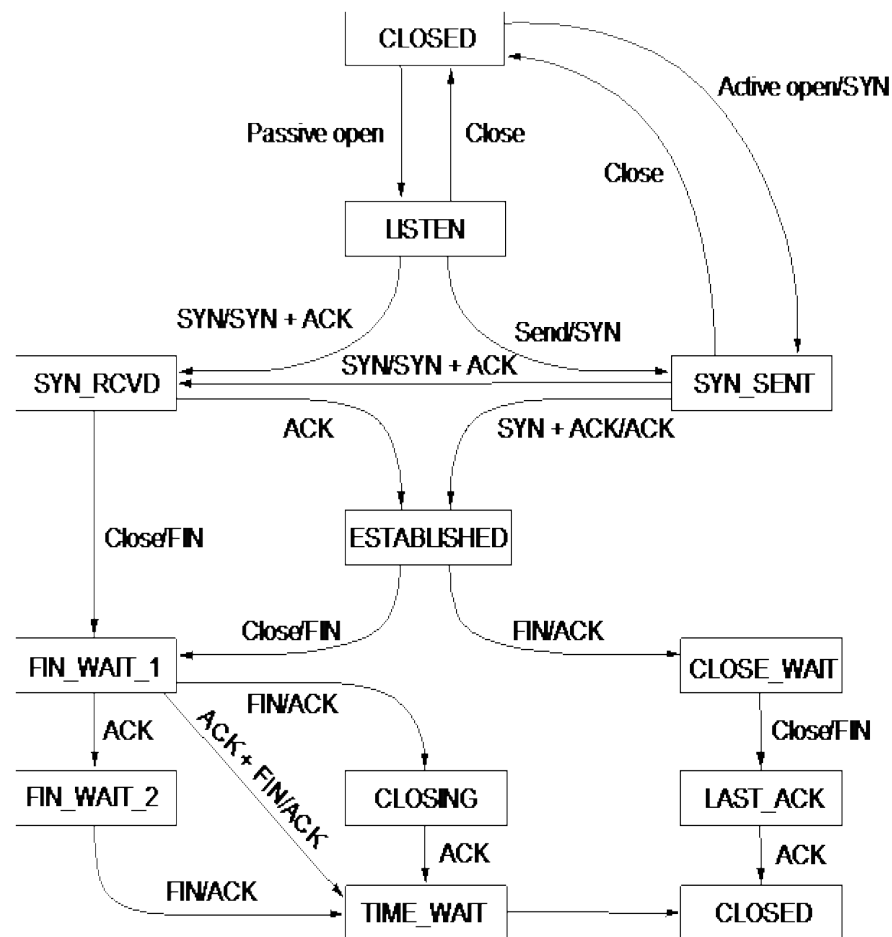
TCP connection termination:



- full duplex
- half duplex

More generally, finite state machine representation of TCP's control mechanism:

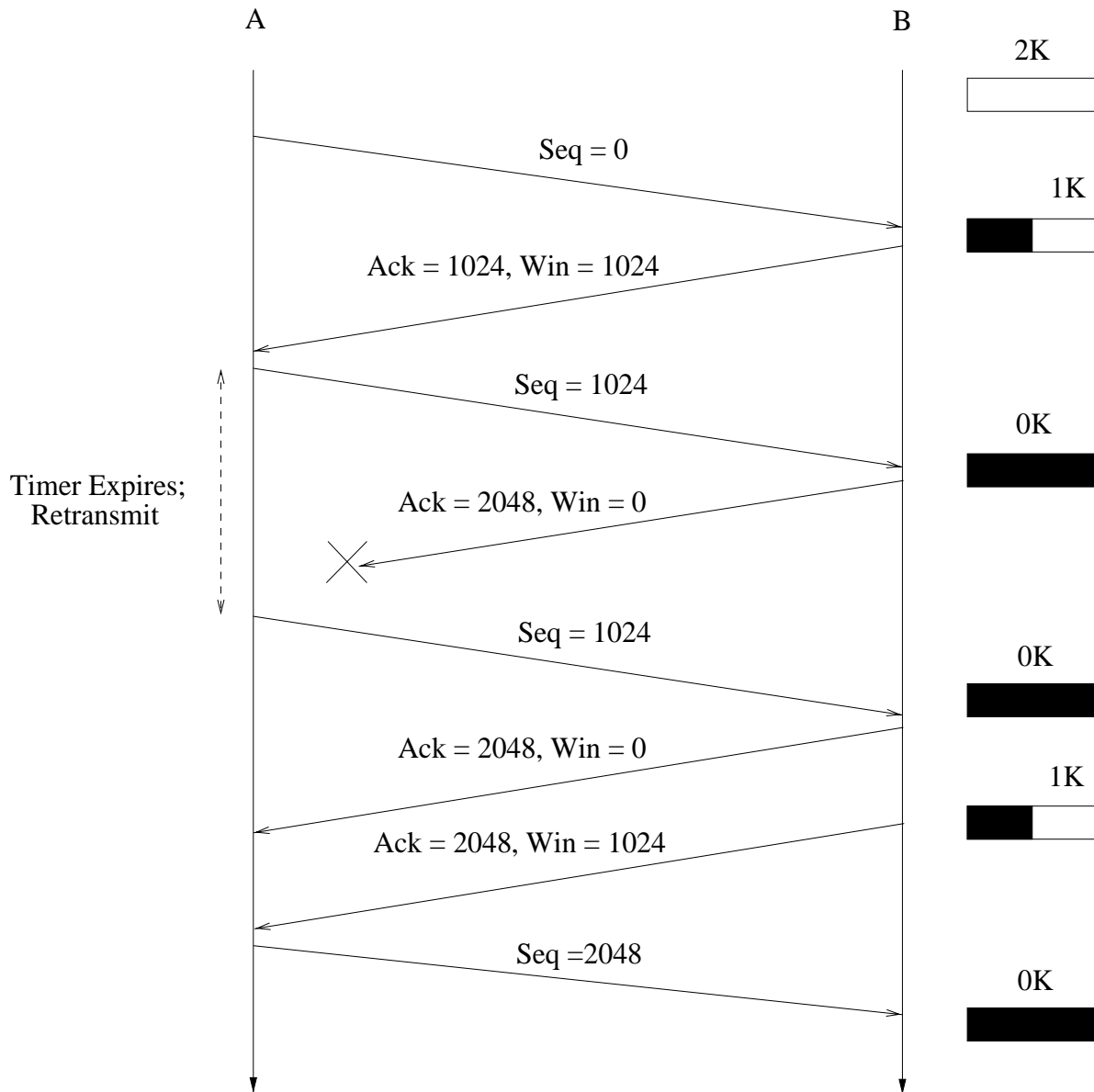
→ state transition diagram



Features to notice:

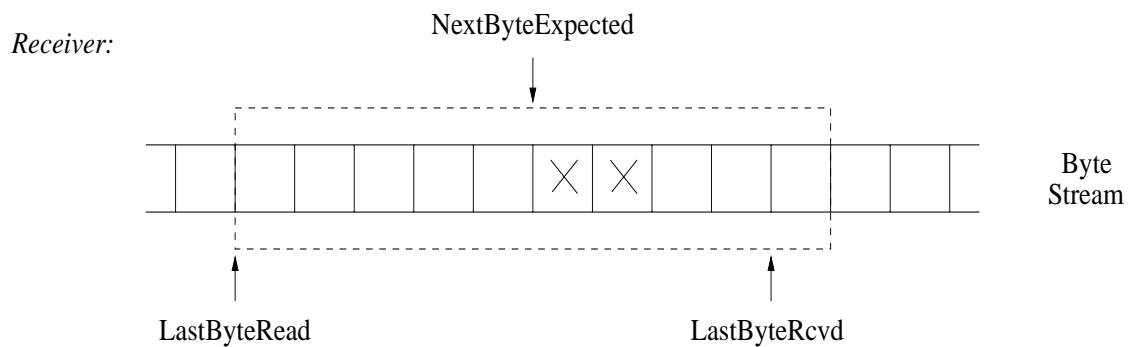
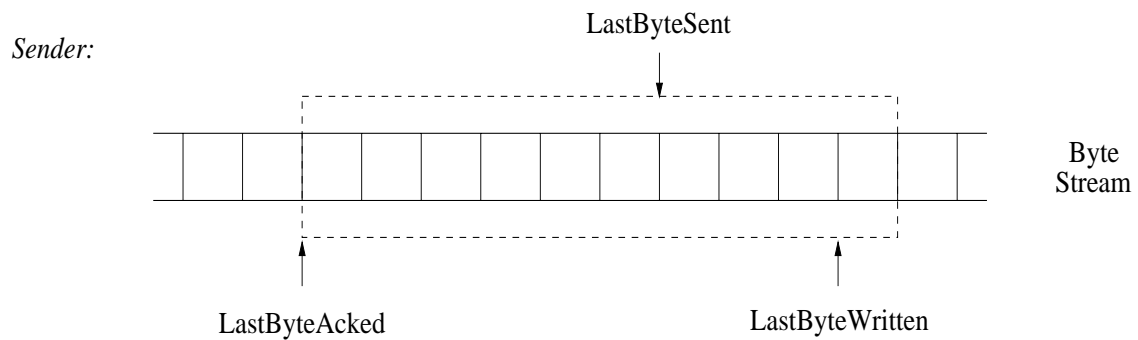
- Connection set-up:
  - client's transition to **ESTABLISHED** state without **ACK**
  - how is server to reach **ESTABLISHED** if client **ACK** is lost?
  - **ESTABLISHED** is macrostate (partial diagram)
- Connection tear-down:
  - three normal cases
  - special issue with **TIME WAIT** state
  - employs hack

Basic TCP data transfer:





## TCP's sliding window protocol



- sender, receiver maintain buffers `MaxSendBuffer`, `MaxRcvBuffer`

Note asynchrony between TCP module and application.

Sender side: maintain invariants

- $\text{LastByteAcked} \leq \text{LastByteSent} \leq \text{LastByteWritten}$
- $\text{LastByteWritten} - \text{LastByteAcked} < \text{MaxSendBuffer}$ 
  - buffer flushing (advance window)
  - application blocking
- $\text{LastByteSent} - \text{LastByteAcked} \leq \text{AdvertisedWindow}$

Thus,

$$\text{EffectiveWindow} = \text{AdvertisedWindow} - (\text{LastByteSent} - \text{LastByteAcked})$$

→ upper bound on new send volume

Actually, one additional refinement:

→ `CongestionWindow`

`EffectiveWindow` update procedure:

$$\text{EffectiveWindow} = \text{MaxWindow} - (\text{LastByteSent} - \text{LastByteAcked})$$

where

$$\text{MaxWindow} = \min\{\text{AdvertisedWindow}, \text{CongestionWindow}\}$$

How to set `CongestionWindow`.

→ domain of TCP congestion control

Receiver side: maintain invariants

- $\text{LastByteRead} < \text{NextByteExpected} \leq \text{LastByteRcvd} + 1$
- $\text{LastByteRcvd} - \text{NextByteRead} < \text{MaxRcvBuffer}$ 
  - buffer flushing (advance window)
  - application blocking

Thus,

$$\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$

Issues:

How to let sender know of change in receiver window size after `AdvertisedWindow` becomes 0?

- trigger ACK event on receiver side when `AdvertisedWindow` becomes positive
- sender periodically sends 1-byte probing packet
  - design choice: smart sender/dumb receiver
  - same situation for congestion control

Silly window syndrome: Assuming receiver buffer is full, what if application reads one byte at a time with long pauses?

- can cause excessive 1-byte traffic
- if `AdvertisedWindow < MSS` then set  
`AdvertisedWindow ← 0`

Do not want to send too many 1 B payload packets.

Nagle's algorithm:

- rule: connection can have only one such unacknowledged packet outstanding
- while waiting for ACK, incoming bytes are accumulated (i.e., buffered)

... compromise between real-time constraints and efficiency.

→ useful for **telnet**-type applications

Sequence number wrap-around problem: recall sufficient condition

$$\text{SenderWindowSize} < (\text{MaxSeqNum} + 1)/2$$

→ 32-bit sequence space/16-bit window space

However, more importantly, time until wrap-around important due to possibility of roaming packets.

bandwidth	time until wrap-around †
T1 (1.5 Mbps)	6.4 hrs
Ethernet (10 Mbps)	57 min
T3 (45 Mbps)	13 min
F/E (100 Mbps)	6 min
OC-3 (155 Mbps)	4 min
OC-12 (622 Mbps)	55 sec
OC-24 (1.2 Gbps)	28 sec



## RTT estimation

... important to not underestimate nor overestimate.

Karn/Partridge: Maintain running average with precautions

$$\text{EstimateRTT} \leftarrow \alpha \cdot \text{EstimateRTT} + \beta \cdot \text{SampleRTT}$$

- **SampleRTT** computed by sender using timer

- $\alpha + \beta = 1$ ;  $0.8 \leq \alpha \leq 0.9$ ,  $0.1 \leq \beta \leq 0.2$

- **TimeOut**  $\leftarrow 2 \cdot \text{EstimateRTT}$  or

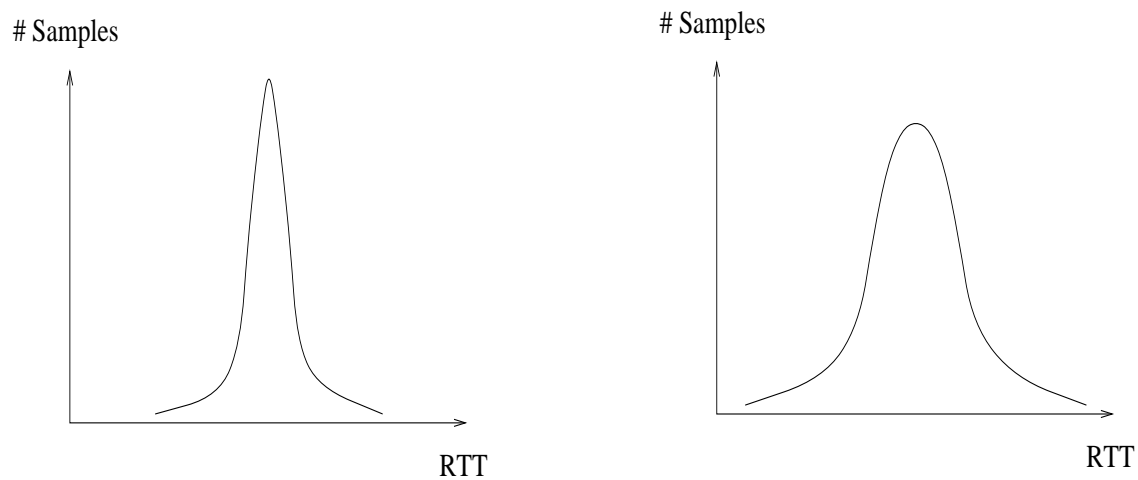
**TimeOut**  $\leftarrow 2 \cdot \text{TimeOut}$  (if retransmit)

→ need to be careful when taking **SampleRTT**

→ infusion of complexity

→ still remaining problems

Hypothetical RTT distribution:



→ need to account for variance

→ not nearly as nice

Jacobson/Karels:

- $\text{Difference} = \text{SampleRTT} - \text{EstimatedRTT}$
- $\text{EstimatedRTT} = \text{EstimatedRTT} + \delta \cdot \text{Difference}$
- $\text{Deviation} = \text{Deviation} + \delta(|\text{Difference}| - \text{Deviation})$

Here  $0 < \delta < 1$ .

Finally,

- $\text{TimeOut} = \mu \cdot \text{EstimatedRTT} + \phi \cdot \text{Deviation}$

where  $\mu = 1$ ,  $\phi = 4$ .

→ persistence timer

→ how to keep multiple timers in UNIX