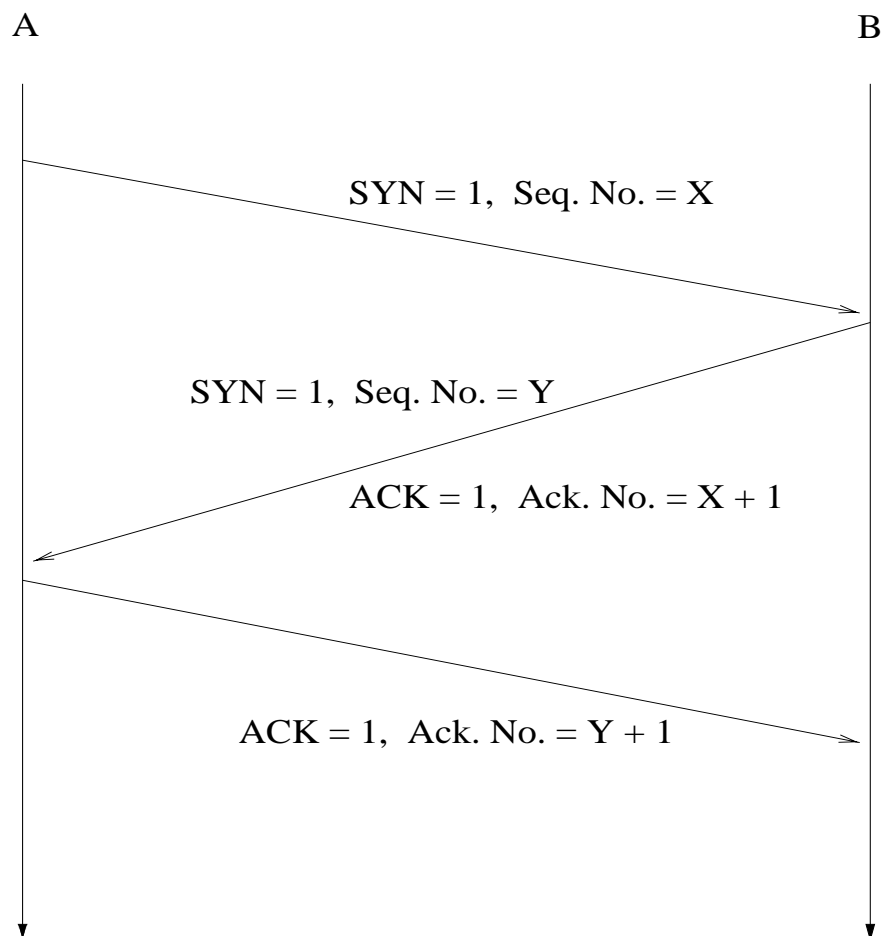TCP connection establishment (3-way handshake):
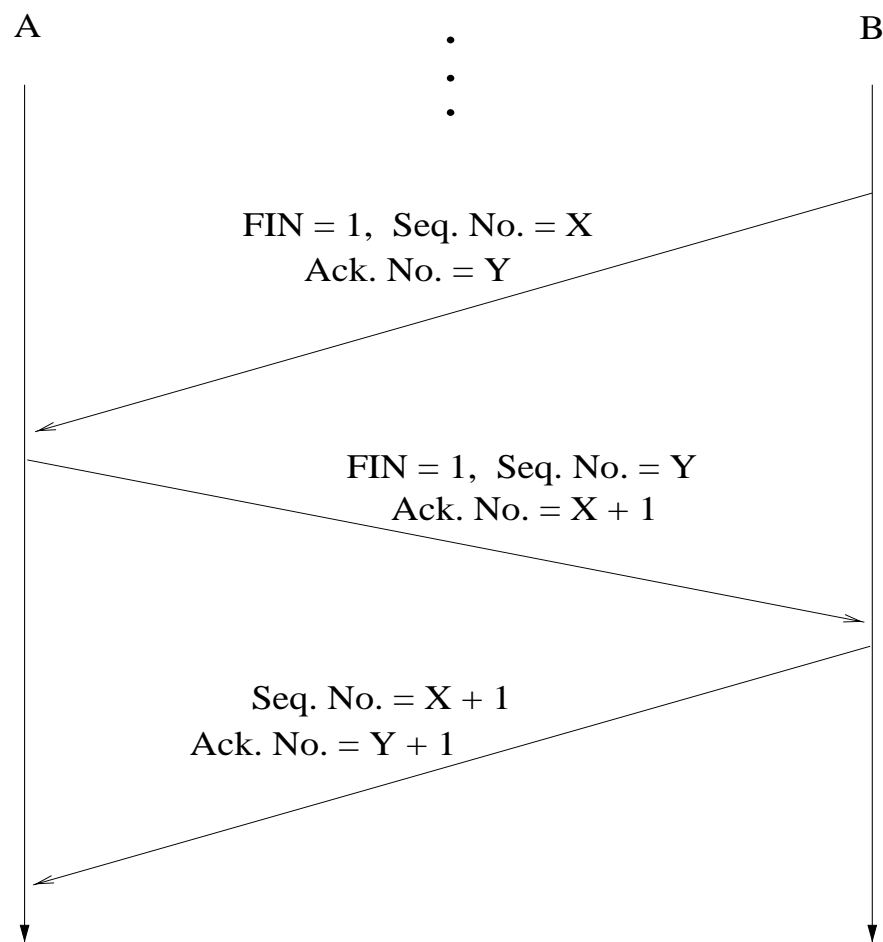


- $X, Y$ are chosen randomly

  $\rightarrow$ sequence number prediction

- piggybacking

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

   $\longrightarrow$   in general: impossible

   $\longrightarrow$   can be proven

   $\longrightarrow$   "acknowledging the ACK problem"

   $\longrightarrow$   also TCP session ending

   $\longrightarrow$   lunch date problem

# TCP connection termination:

A          B

FIN = 1, Seq. No. = X
Ack. No. = Y

FIN = 1, Seq. No. = Y
Ack. No. = X + 1

Seq. No. = X + 1
Ack. No. = Y + 1

- full duplex

- half duplex

More generally, finite state machine representation of TCP's control mechanism:
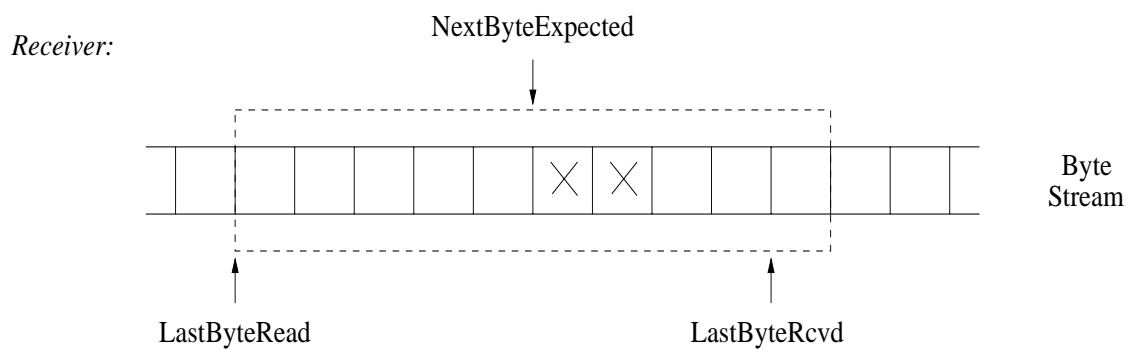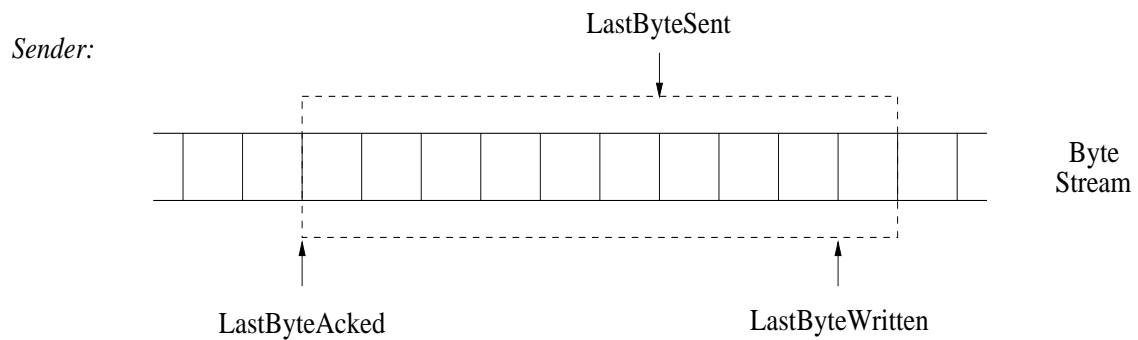
$\longrightarrow$   state transition diagram

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

# TCP's sliding window protocol

*Sender:*

LastByteSent

Byte
Stream

LastByteAcked

LastByteWritten

*Receiver:*

NextByteExpected

Byte
Stream

LastByteRead

LastByteRcvd

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

Same as generic sliding window

$\longrightarrow$   data unit: byte, not packet

Sender side: maintain invariants

- `LastByteAcked` $\leq$ `LastByteSent` $\leq$ `LastByteWritten`

- `LastByteWritten` $-$ `LastByteAcked` $<$ `MaxSendBuffer`

  $\longrightarrow$   buffer flushing (advance window)

  $\longrightarrow$   application blocking

- `LastByteSent` $-$ `LastByteAcked` $\leq$ `AdvertisedWindow`

  $\longrightarrow$   `AdvertisedWindow`: receiver side free space

  $\longrightarrow$   throttling effect

How much sender can still send:

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

$\longrightarrow$  upper bound

$\longrightarrow$  sender may choose to send less

$\longrightarrow$  self-throttling

Affected through sender side variable

$\longrightarrow$  $\texttt{CongestionWindow}$

`EffectiveWindow` update procedure:

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

where

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

How to set `CongestionWindow`.

$$\longrightarrow \quad \text{TCP congestion control}$$

Receiver side: maintain invariants

- $\texttt{LastByteRead} < \texttt{NextByteExpected} \leq$
  $\texttt{LastByteRcvd} + 1$

- $\texttt{LastByteRcvd} - \texttt{NextByteRead} < \texttt{MaxRcvBuffer}$

  $\longrightarrow$ buffer flushing (advance window)

  $\longrightarrow$ application blocking

Thus,

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

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

  $\longrightarrow$ design choice: smart sender/dumb receiver

  $\longrightarrow$ 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 $\texttt{AdvertisedWindow} < \text{MSS}$ then set
  $\texttt{AdvertisedWindow} \leftarrow 0$

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

Nagle's method:

- 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.

$\rightarrow$ useful for `telnet/ssh`-type interactive applications

## RTT estimation

. . . important to not underestimate nor overestimate.

Karn/Partridge: Maintain running average with precautions

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

- `SampleRTT` computed by sender using timer

- $\alpha + \beta = 1$; $\ 0.8 \leq \alpha \leq 0.9$, $0.1 \leq \beta \leq 0.2$
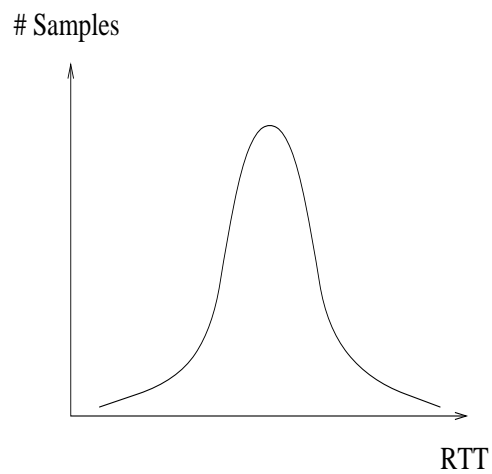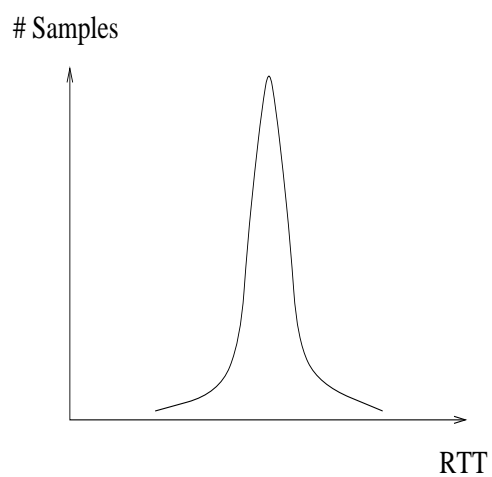
- `TimeOut` $\leftarrow 2 \cdot$ `EstimateRTT`   or

  `TimeOut` $\leftarrow 2 \cdot$ `TimeOut`   (if retransmit)

      $\longrightarrow$   need to be careful when taking `SampleRTT`

      $\longrightarrow$   infusion of complexity

      $\longrightarrow$   still remaining problems

# Hypothetical RTT distribution:

# Samples
RTT

# Samples
RTT

$\longrightarrow$   need to account for variance

$\longrightarrow$   not nearly as nice

Jacobson/Karels:

- $\texttt{Difference} = \texttt{SampleRTT} - \texttt{EstimatedRTT}$

- $\texttt{EstimatedRTT} = \texttt{EstimatedRTT} + \delta \cdot \texttt{Difference}$

- $\texttt{Deviation} = \texttt{Deviation} + \delta(|\texttt{Difference}| - \texttt{Deviation})$

Here $0 < \delta < 1$.

Finally,

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

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

$\longrightarrow$    persistence timer

$\longrightarrow$    how to keep multiple timers in UNIX