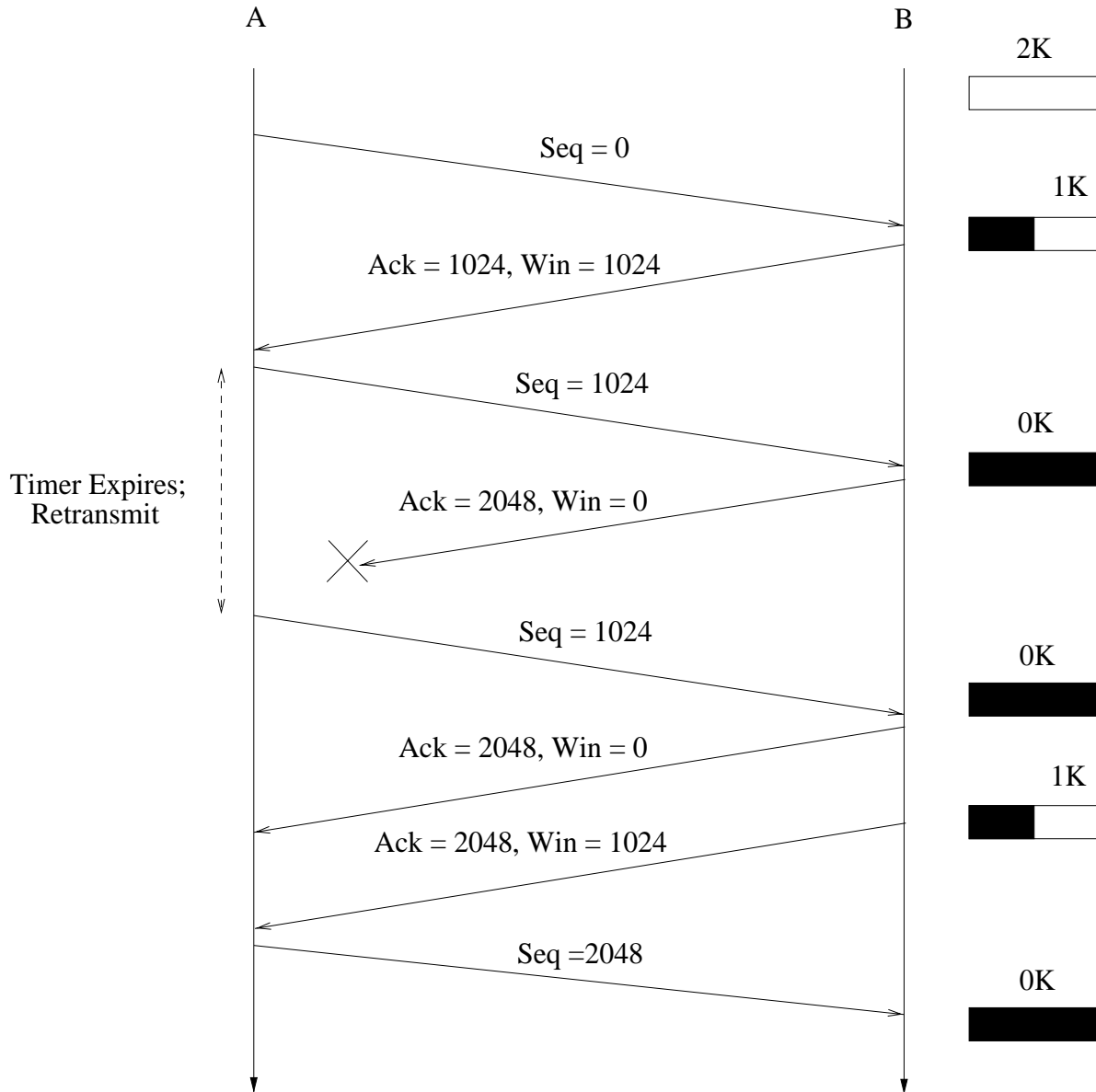
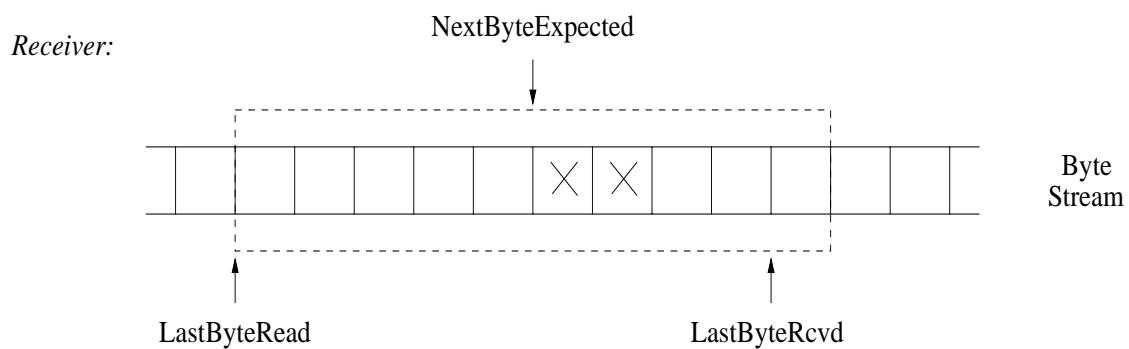
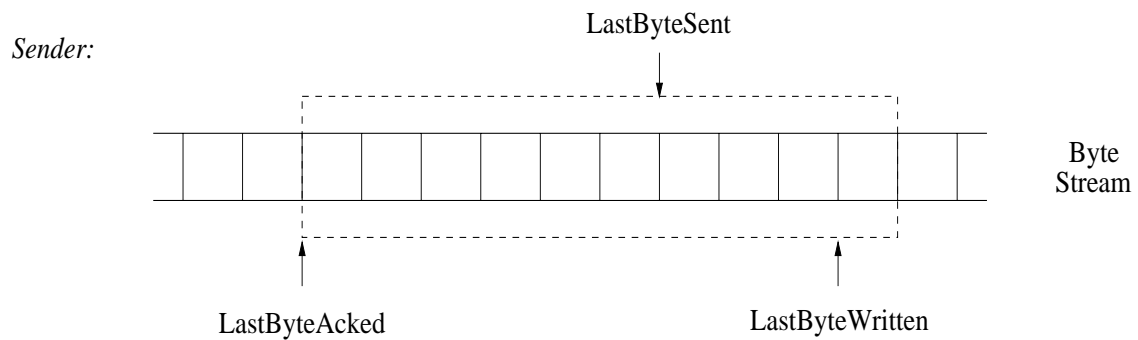


Basic TCP data transfer:



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

Even more importantly, “keeping-the-pipe-full” consideration.

bandwidth	delay-bandwidth product †
T1 (1.5 Mbps)	18 kB
Ethernet (10 Mbps)	122 kB
T3 (45 Mbps)	549 kB
FDDI (100 Mbps)	1.2 MB
OC-3 (155 Mbps)	1.8 MB
OC-12 (622 Mbps)	7.4 MB
OC-24 (1.2 Gbps)	14.8 MB

→ 100 ms latency

Also, throughput limitation imposed by TCP receiver window size.

→ e.g., high-performance grid apps

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

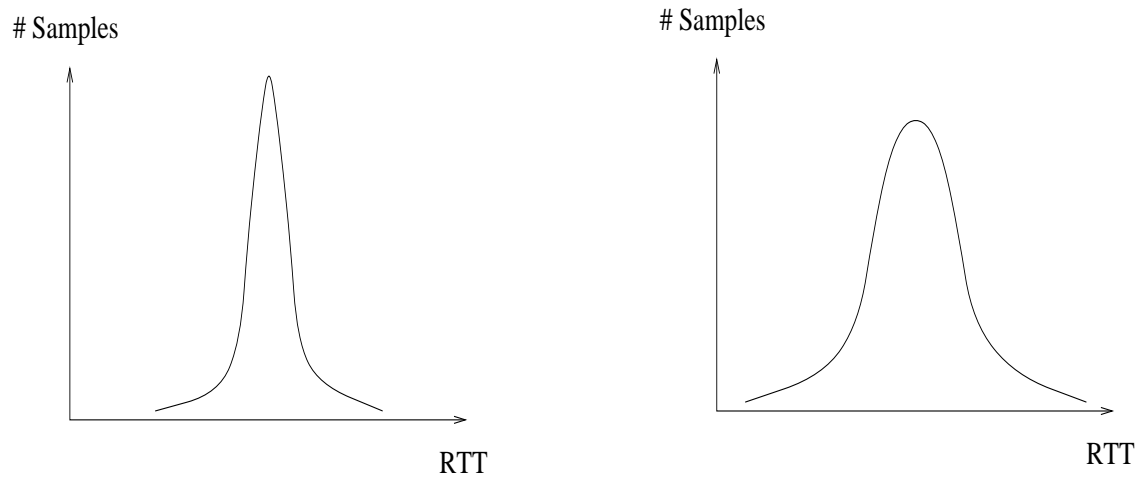
TimeOut $\leftarrow 2 \cdot \text{TimeOut}$ (if retransmit)

→ 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