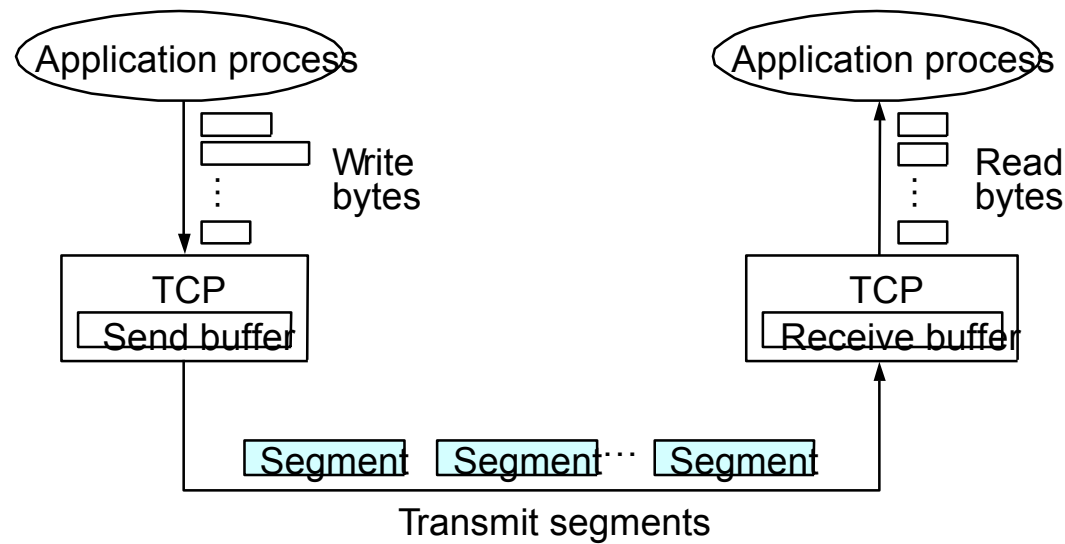# TCP Overview

▸ Connection-oriented

▸ Byte-stream
- app writes bytes
- TCP sends segments
- app reads bytes

▸ Full duplex

▸ Flow control: keep sender from overrunning receiver

▸ Congestion control: keep sender from overrunning network

Application process → Write bytes → TCP / Send buffer → Segment Segment ⋯ Segment → Transmit segments → TCP / Receive buffer → Read bytes → Application process

# Remember

- TCP: Segments, IP: packets, Ethernet: Frames
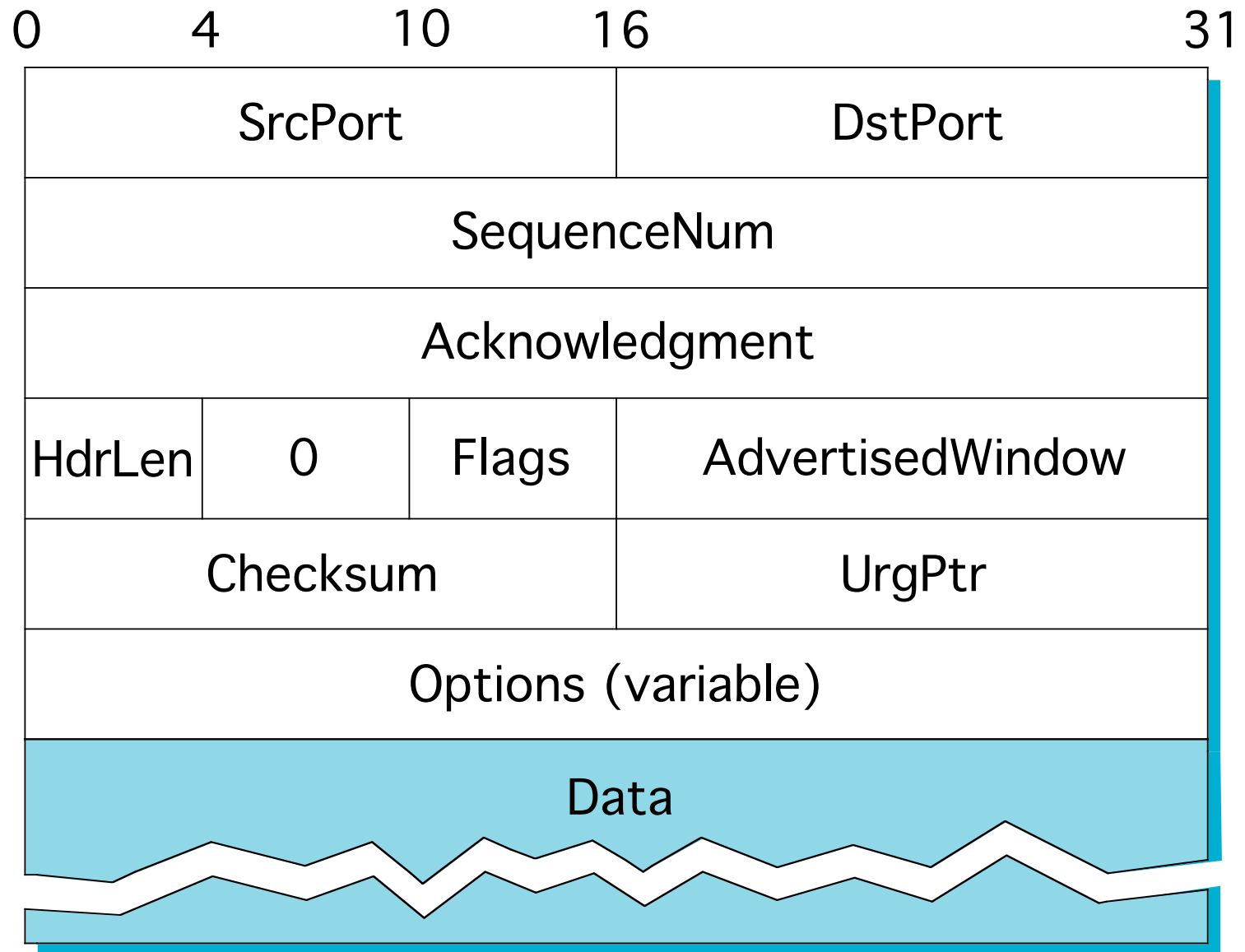- TCP segments are encapsulated as IP data which are encapsulated as Ethernet frame data

# Data Link Versus Transport (TCP)

- ▶ **Potentially connects many different hosts**
  - ■ need explicit connection establishment and termination
    - ● TCP three way hand-shake
- ▶ **Potentially different RTT**
  - ■ need adaptive timeout mechanism
- ▶ **Potentially long delay in network**
  - ■ need to be prepared for arrival of very old packets
    - ● Sequence number space should be well thought out
- ▶ **Potentially different capacity at destination**
  - ■ need to accommodate different node capacity
    - ● Flow control
- ▶ **Potentially different network capacity**
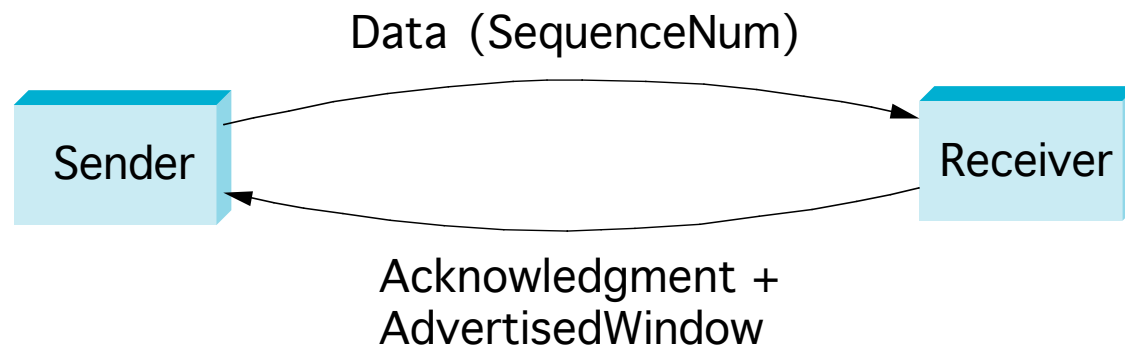  - ■ need to be prepared for network congestion
    - ● Congestion control

# Segment Format

| 0 | 4 | 10 | 16 | 31 |
|---|---|----|----|----|

| SrcPort | DstPort |
|---------|---------|

| SequenceNum |
|-------------|

| Acknowledgment |
|----------------|

| HdrLen | 0 | Flags | AdvertisedWindow |
|--------|---|-------|------------------|

| Checksum | UrgPtr |
|----------|--------|

| Options (variable) |
|--------------------|

| Data |
|------|

# Segment Format (cont)

▶ Each connection identified with 4-tuple:
  ■ (SrcPort, SrcIPAddr, DstPort, DstIPAddr)
▶ Sliding window + flow control
  ■ acknowledgment, SequenceNum, AdvertisedWindow

Data  (SequenceNum)

Sender → Receiver

Acknowledgment +
AdvertisedWindow

▶ Flags
  ■ SYN, FIN, RESET, PUSH, URG, ACK
▶ Checksum
  ■ pseudo header + TCP header + data

# Sequence Number Selection

▶ Initial sequence number (ISN) selection

■ Why not simply chose 0?

● Must avoid overlap with earlier incarnation.

● New sequence number should be larger than previous number.

■ Why can't the system remember the previous number used?

▶ Requirements for ISN selection

■ Must operate correctly

● Without synchronized clocks
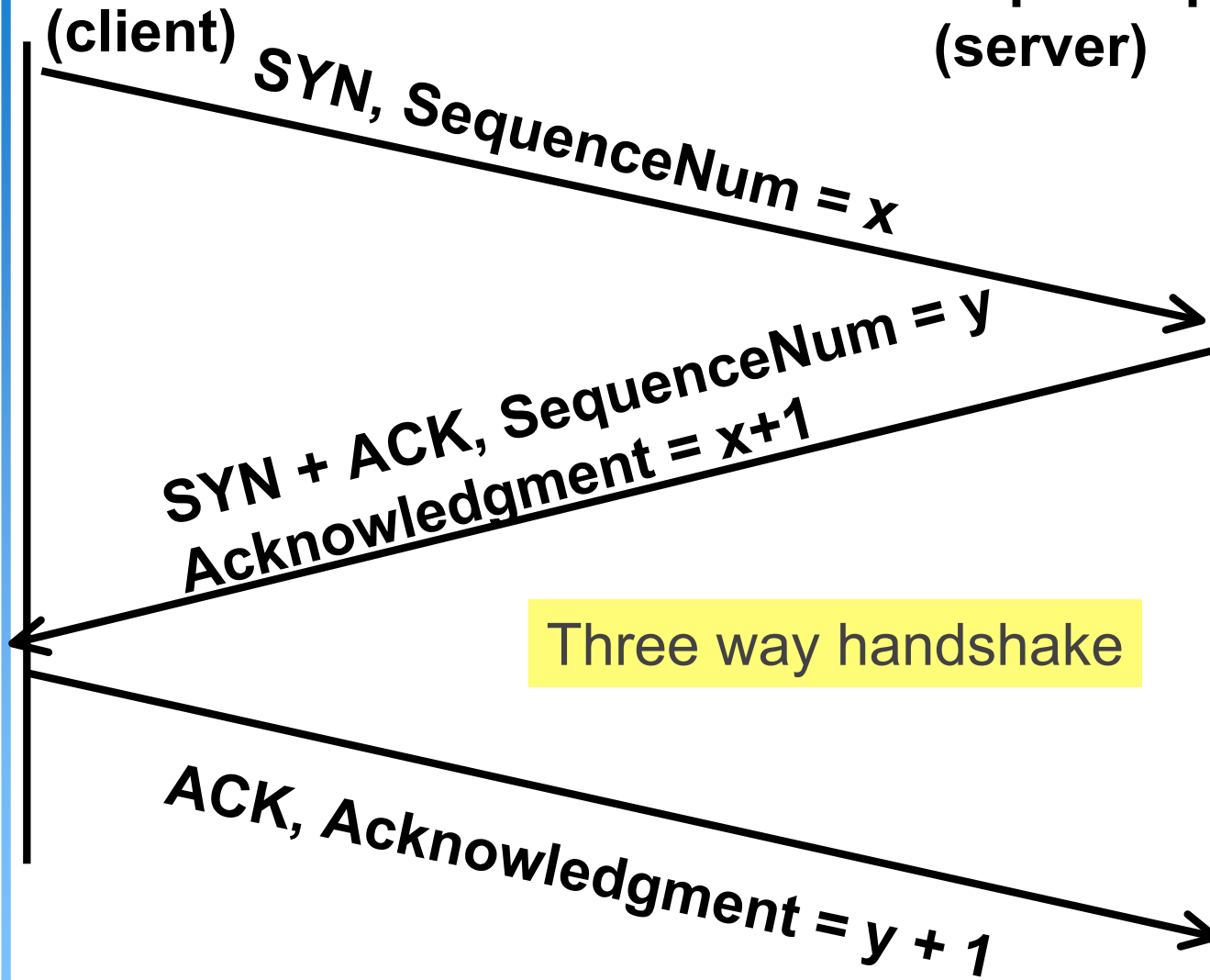
● Despite node failures

# ISN and Quiet Time

- ▶ Use local clock to select ISN
  - ■ Clock wraparound must be greater than max segment lifetime (MSL)
- ▶ Upon startup, cannot assign sequence numbers for MSL seconds
- ▶ Can still have sequence number overlap
  - ■ If sequence number space not large enough for high-bandwidth connections

# Connection Establishment

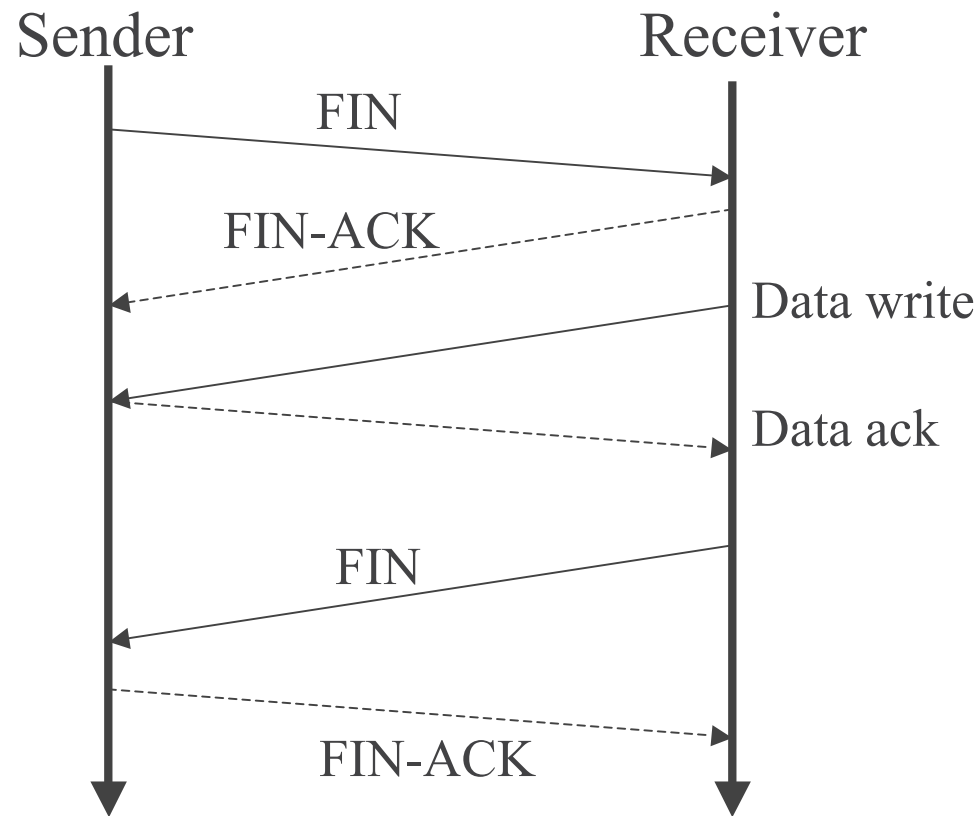**Active participant (client)**

**Passive participant (server)**

SYN, SequenceNum = x

SYN + ACK, SequenceNum = y Acknowledgment = x+1

Three way handshake

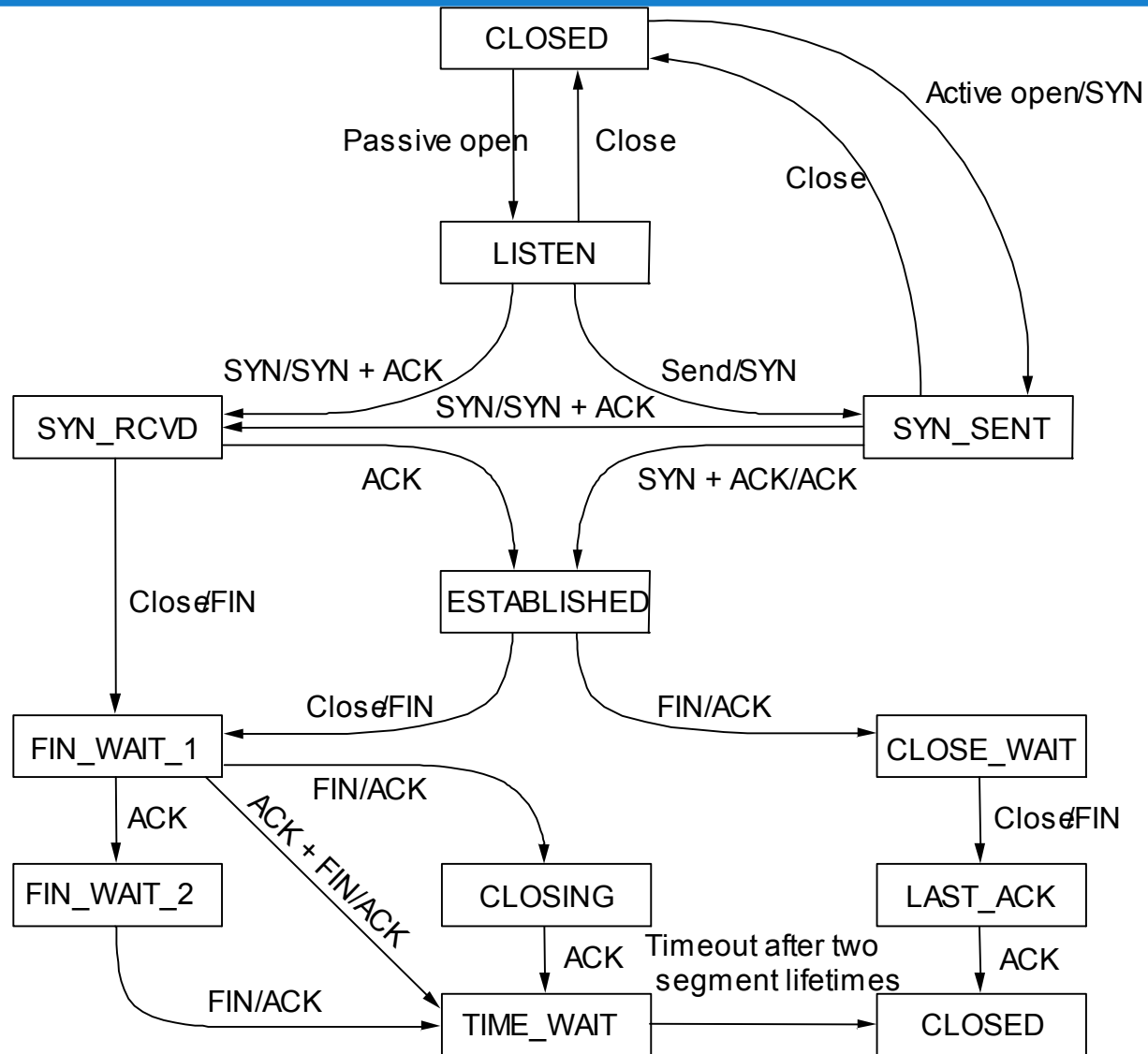ACK, Acknowledgment = y + 1

# Connection Tear-down

▸ Normal termination

  ■ Allow unilateral close

  ■ Avoid sequence number overlap

▸ TCP must continue to receive data even after closing

  ■ Cannot close connection immediately: what if a new connection restarts and uses same sequence number and receives retransmitted FIN from the current session?
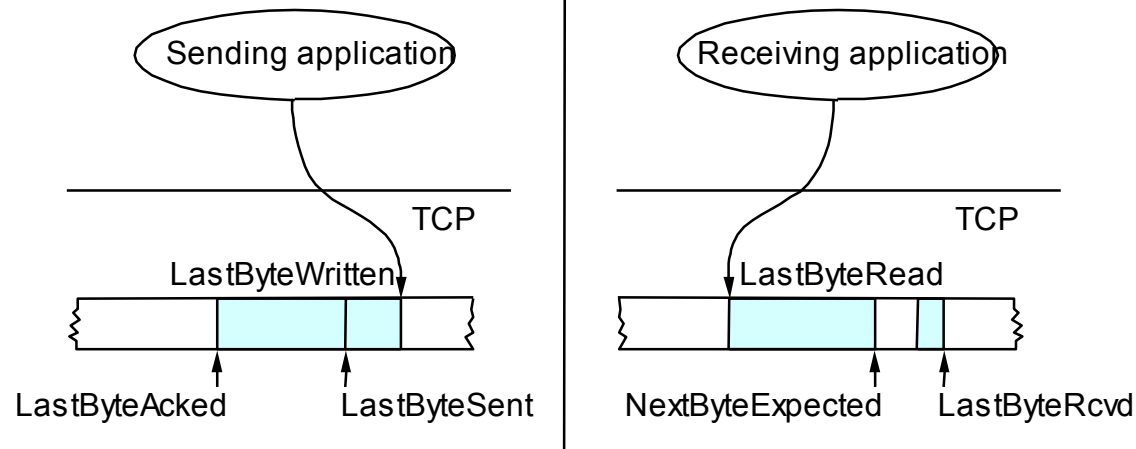
# Tear-down Packet Exchange

Sender                    Receiver

FIN →

← FIN-ACK

Data write →

← Data ack

← FIN

FIN-ACK →

# State Transition Diagram

# Sliding Window Revisited



- ▸ Sending side
  - ■ `LastByteAcked <= LastByteSent`
  - ■ `LastByteSent <= LastByteWritten`
  - ■ buffer bytes between `LastByteAcked` and `LastByteWritten`

- ▸ Receiving side
  - ■ `LastByteRead < NextByteExpected`
  - ■ `NextByteExpected <= LastByteRcvd +1`
  - ■ buffer bytes between `NextByteRead` and `LastByteRcvd`

# Flow Control

▸ Fast sender can overrun receiver:

- Packet loss, unnecessary retransmissions

▸ Possible solutions:

- Sender transmits at pre-negotiated rate
- Sender limited to a window's worth of unacknowledged data

▸ Flow control different from congestion control

# Flow Control

- Send buffer size: MaxSendBuffer
- Receive buffer size: MaxRcvBuffer
- Receiving side
  - LastByteRcvd - LastByteRead < = MaxRcvBuffer
  - AdvertisedWindow = MaxRcvBuffer - (NextByteExpected - NextByteRead)
- Sending side
  - LastByteSent - LastByteAcked < = AdvertisedWindow
  - EffectiveWindow = AdvertisedWindow - (LastByteSent - LastByteAcked)
  - LastByteWritten - LastByteAcked < = MaxSendBuffer
  - block sender if (LastByteWritten - LastByteAcked) + y > MaxSenderBuffer
- Always send ACK in response to arriving data segment
- Persist when AdvertisedWindow = 0

# Round-trip Time Estimation

- Wait at least one RTT before retransmitting
- Importance of accurate RTT estimators:
  - Low RTT -> unneeded retransmissions
  - High RTT -> poor throughput
- RTT estimator must adapt to change in RTT
  - But not too fast, or too slow!

- Problem: If the instantaneously calculated RTT is 10, 20, 5, 12, 3 , 5, 6; what RTT should we use for calculations?

# Initial Round-trip Estimator
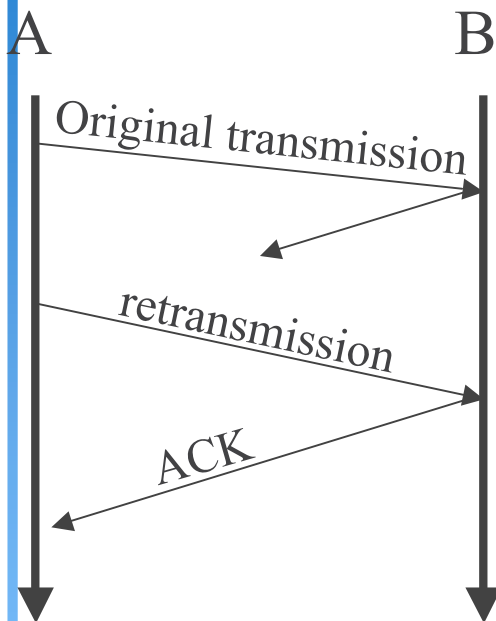
Round trip times exponentially averaged:

▸ New RTT = $\alpha$ (old RTT) + (1 - $\alpha$) (new sample)

▸ Recommended value for $\alpha$: 0.8 - 0.9

▸ Retransmit timer set to $\beta$ RTT, where $\beta$ = 2
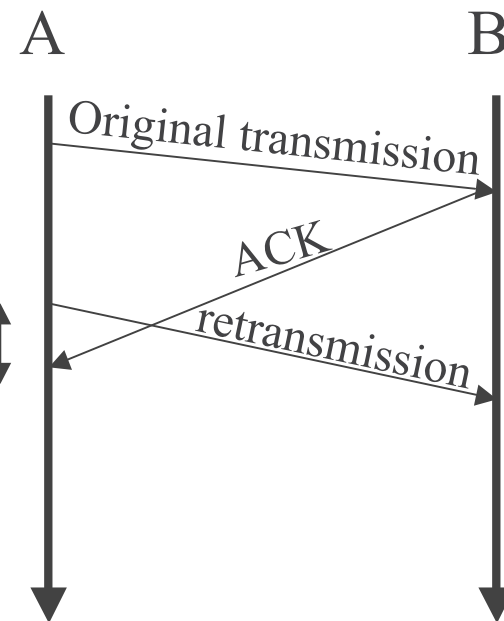
▸ Every time timer expires, RTO exponentially backed-off

# Retransmission Ambiguity



Sample RTT

A  B

Original transmission

retransmission

ACK

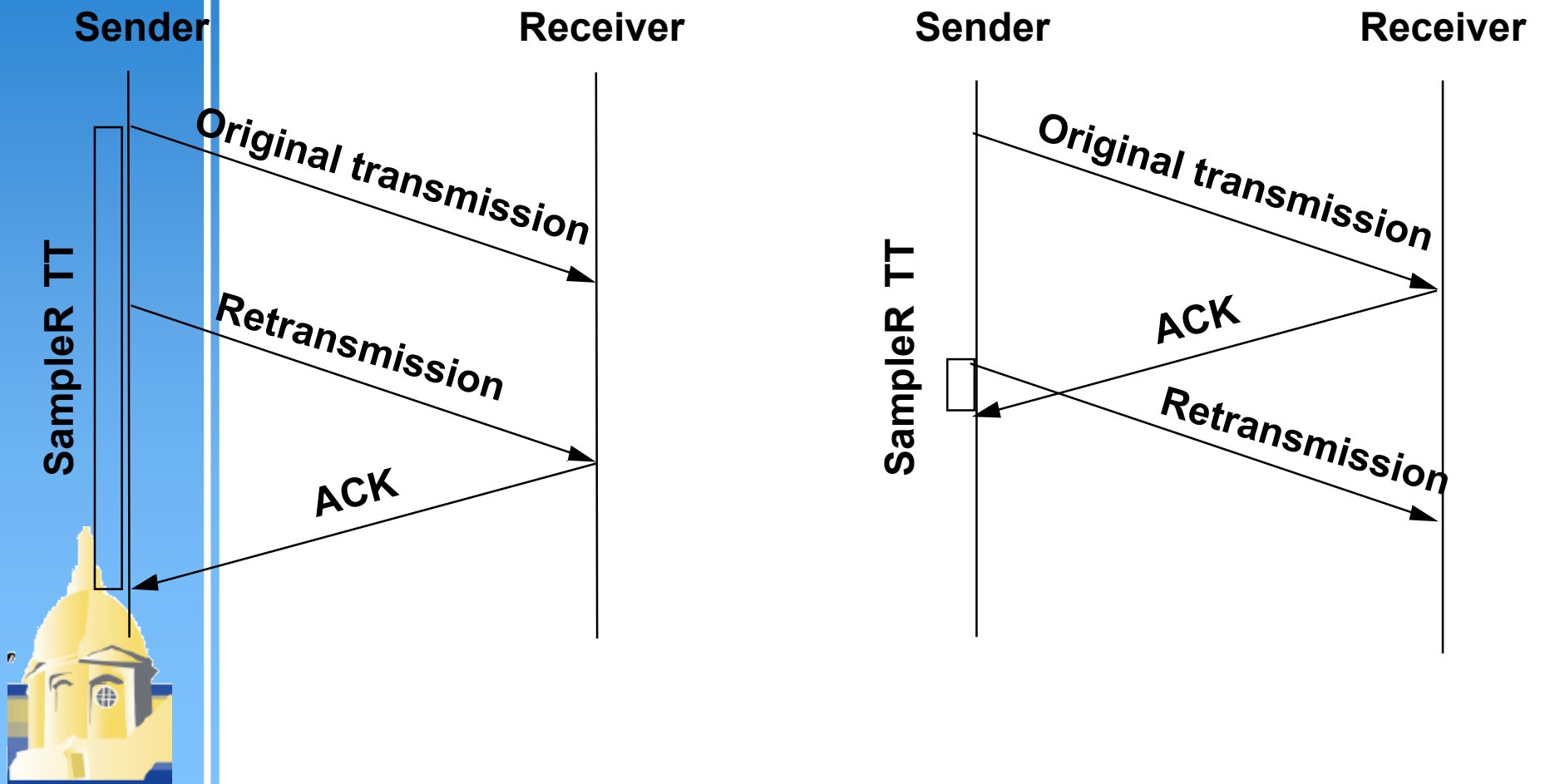Sample RTT

A  B

Original transmission

ACK

retransmission

# Karn's Retransmission Timeout Estimator

▸ Accounts for retransmission ambiguity

▸ If a segment has been retransmitted:

  ■ Don't count RTT sample on ACKs for this segment

  ■ Keep backed off time-out for next packet

  ■ Reuse RTT estimate only after one successful transmission

# Karn/Partridge Algorithm

▶ Do not sample RTT when retransmitting

▶ Double timeout after each retransmission

| Sender | Receiver | | Sender | Receiver |

**SampleR TT**

Original transmission

Retransmission

ACK

**SampleR TT**

Original transmission

ACK

Retransmission

# Jacobson's Retransmission Timeout Estimator

- Key observation:
  - Using $\beta$ RTT for timeout doesn't work
  - At high loads round trip variance is high
- Solution:
  - If D denotes mean variation
  - Timeout = RTT + 4D
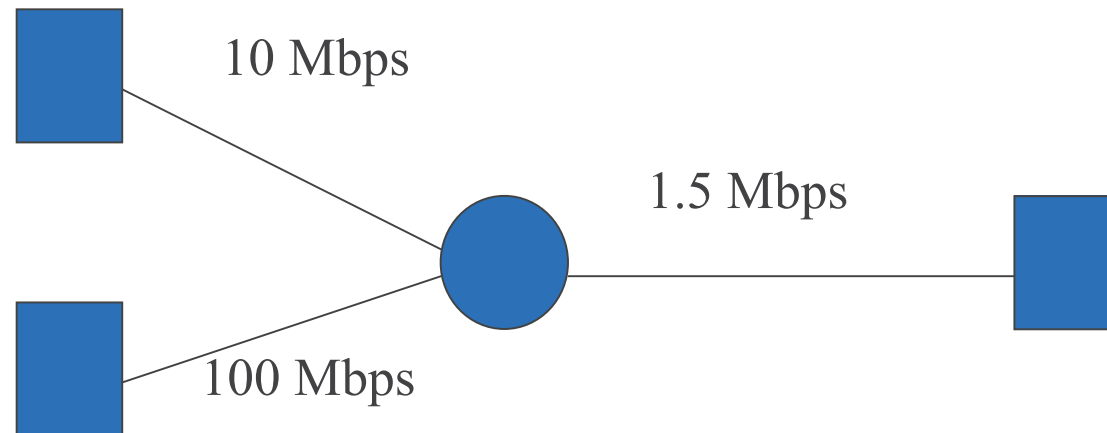
# Jacobson/ Karels Algorithm

- ▶ New Calculations for average RTT
- ▶ Diff = SampleRTT - EstRTT
- ▶ EstRTT = EstRTT + (d x Diff)
- ▶ Dev = Dev + d( |Diff| - Dev)
    - where d is a factor between 0 and 1
- ▶ Consider variance when setting timeout value
- ▶ TimeOut = m x EstRTT + f x Dev
    - where m = 1 and f = 4
- ▶ Notes
    - algorithm only as good as granularity of clock (500ms on Unix)
    - accurate timeout mechanism important to congestion control (later)
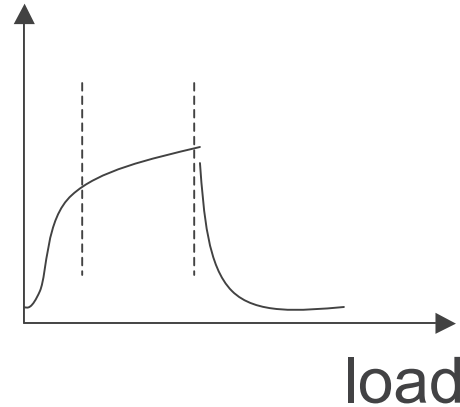
# Congestion

10 Mbps

1.5 Mbps

100 Mbps

▶ If both sources send full windows, we may get congestion collapse

▶ Other forms of congestion collapse:
  - Retransmissions of large packets after loss of a single fragment
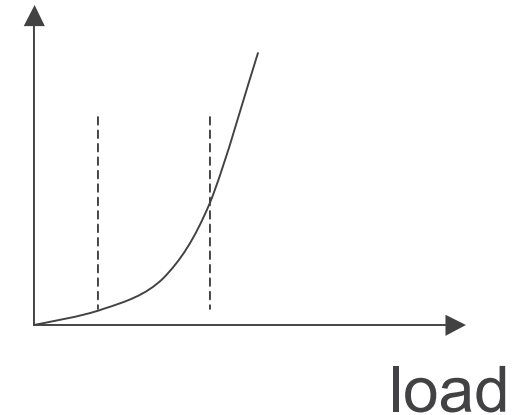  - Non-feedback controlled sources

# Congestion Response

throughput

delay

load

load

Avoidance keeps the system performing at the *knee*

Control kicks in once the system has reached a congested state

# Separation of Functionality

▶ Sending host must adjust amount of data it puts in the network based on detected congestion

▶ Routers can help by:

  ◼ Sending accurate congestion signals

  ◼ Isolating well-behaved from ill-behaved sources