

End-to-End (Transport) Protocols

- Underlying best-effort network
 - drops messages
 - re-orders messages
 - delivers duplicate copies of a given message
 - limits messages to some finite size
 - delivers messages after an arbitrarily long delay
- Common end-to-end services
 - guarantee message delivery
 - deliver messages in the same order they are sent
 - deliver at most one copy of each message
 - support arbitrarily large messages
 - support synchronization
 - allow the receiver to apply flow control to the sender
 - support multiple application processes on each host

Simple Demultiplexor (UDP)

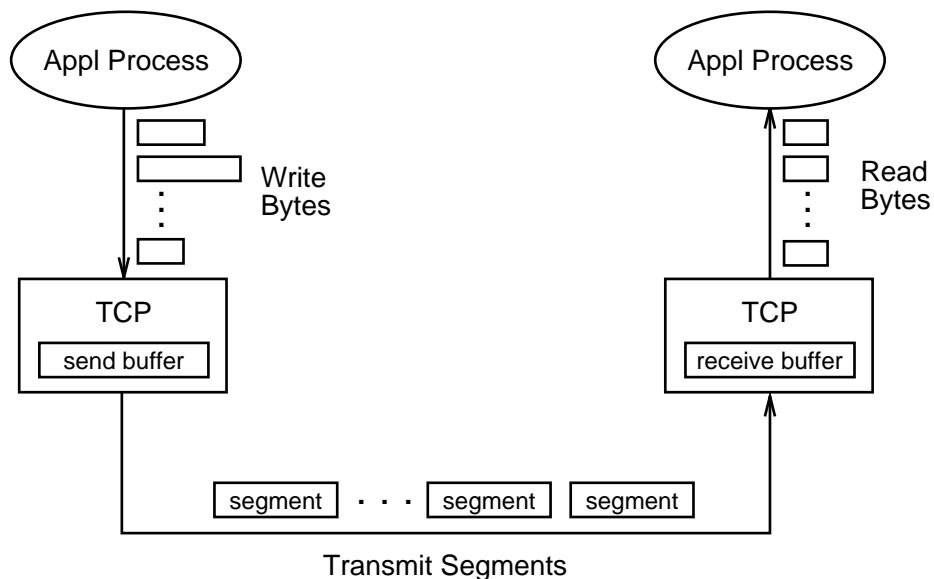
- Unreliable and unordered datagram service
- Adds multiplexing
- No flow control
- Endpoints identified by ports
 - servers have *well-known* ports
 - see `/etc/services` on Unix
- Optional checksum
 - pseudo header + udp header + data
- Header format

Src Port	Dest Port
Check Sum	Length

Reliable Byte-Stream (TCP)

Overview

- Connection-oriented
- Byte-stream
 - sending process writes some number of bytes
 - TCP breaks into *segments* and sends via IP
 - receiving process reads some number of bytes



- Full duplex
- Flow control: keep sender from overrunning receiver
- Congestion control: keep sender from overrunning network

End-to-End Issues

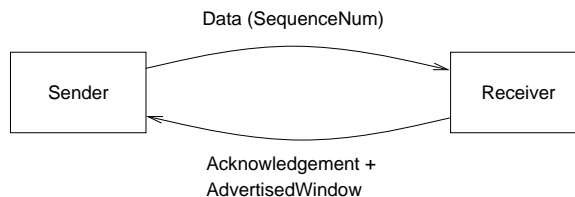
Based on sliding window protocol used at data link level, but the situation is very different.

- Potentially connects many different hosts
 - need explicit connection establishment and termination
- Potentially different RTT
 - need adaptive timeout mechanism
- Potentially long delay in network
 - need to be prepared for arrival of very old packets
- Potentially different capacity at destination
 - need to accommodate different amounts of buffering
- Potentially different network capacity
 - need to be prepared for network congestion

Segment Format

Src Port		Dest Port	
SequenceNum			
Acknowledgement			
HdrLen (4)	0 (6)	Flags (6)	Advertised Window
Checksum		UrgPtr	
options (variable)			
data			

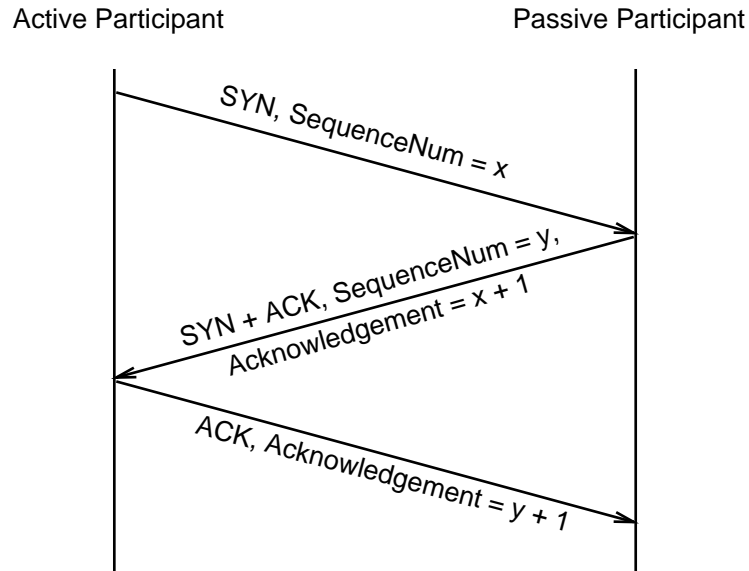
- Each connection identified with 4-tuple:
 - $\langle \text{SrcPort}, \text{SrcIPAddr}, \text{DstPort}, \text{DstIPAddr} \rangle$
- Sliding window + flow control
 - Acknowledgment, SequenceNum, AdvertisedWindow



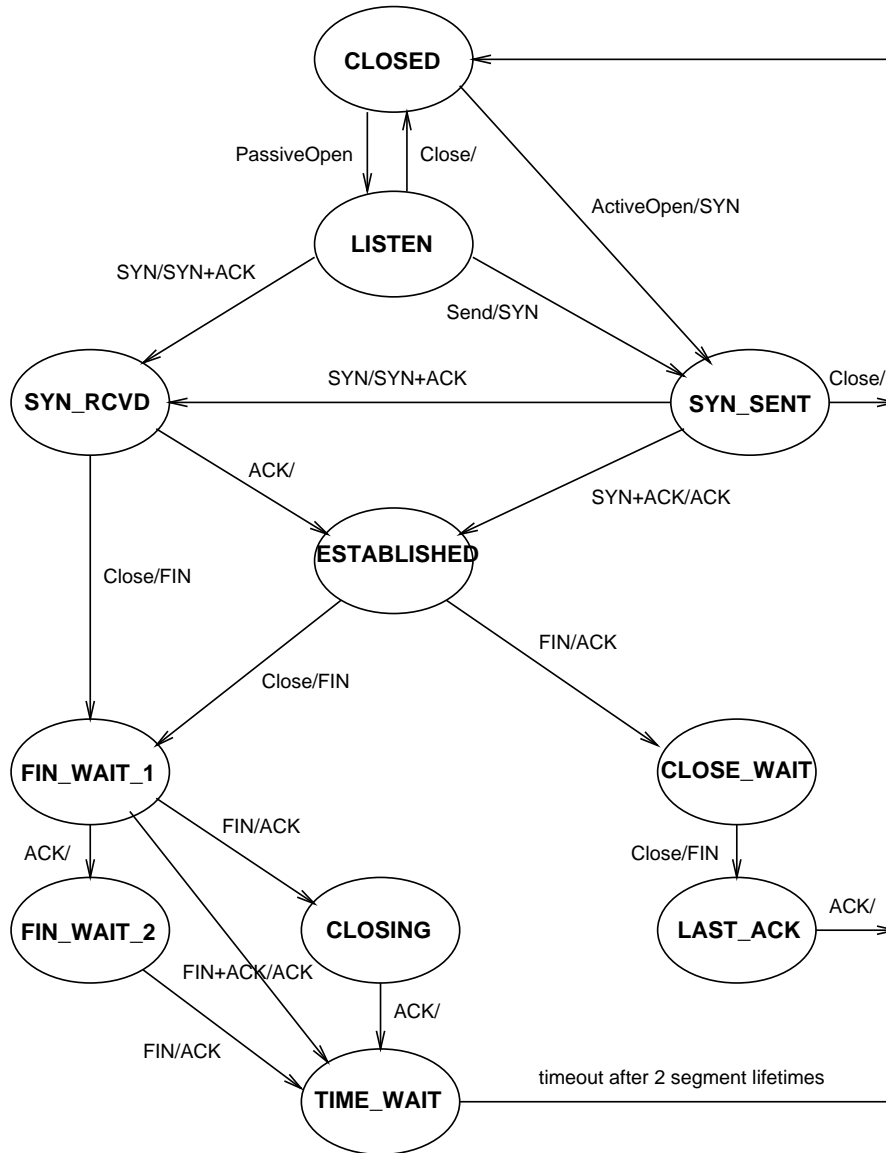
- Flags
 - SYN, FIN, RESET, PUSH, URG, ACK
- Checksum
 - pseudo header + tcp header + data

Connection Establishment and Termination

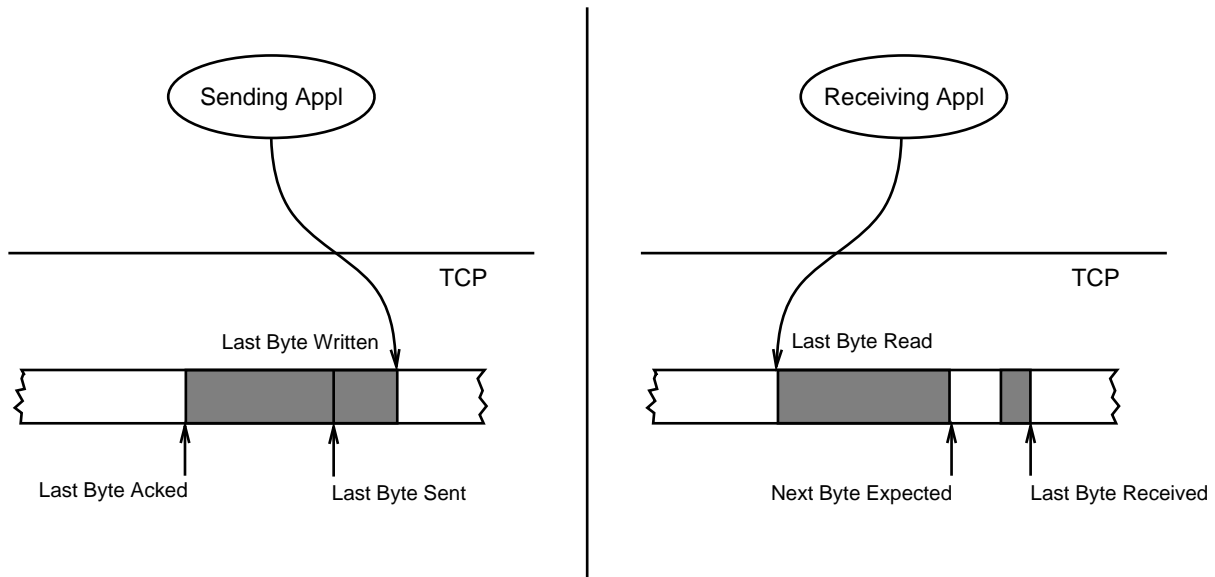
Three-Way Handshake



State Transition Diagram



Sliding Window Revisited



- Each byte has a sequence number
- ACKs are cumulative
- Sending side
 - $\text{LastByteAcked} \leq \text{LastByteSent}$
 - $\text{LastByteSent} \leq \text{LastByteWritten}$
 - bytes between **LastByteAcked** and **LastByteWritten** must be buffered
- Receiving side
 - $\text{LastByteRead} < \text{NextByteExpected}$
 - $\text{NextByteExpected} \leq \text{LastByteRcvd} + 1$
 - bytes between **NextByteRead** and **LastByteRcvd** must be buffered

Flow Control

- Sender buffer size: `MaxSendBuffer`
- Receive buffer size: `MaxRcvBuffer`
- Receiving side
 - $\text{LastByteRcvd} - \text{NextByteRead} \leq \text{MaxRcvBuffer}$
 - $\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{NextByteRead})$
- Sending side
 - $\text{LastByteSent} - \text{LastByteAcked} \leq \text{AdvertisedWindow}$
 - $\text{EffectiveWindow} = \text{AdvertisedWindow} - (\text{LastByteSent} - \text{LastByteAcked})$
 - $\text{LastByteWritten} - \text{LastByteAcked} \leq \text{MaxSendBuffer}$
 - block sender if $(\text{LastByteWritten} - \text{LastByteAcked}) + y > \text{MaxSendBuffer}$
- Always send ACK in response to an arriving data segment
- Persist when `AdvertisedWindow=0`

Keeping the Pipe Full

Wrap Around: 32-bit **SequenceNum**

Bandwidth	Time Until Wrap Around
T1 (1.5Mbps)	6.4 hours
Ethernet (10Mbps)	57 minutes
T3 (45Mbps)	13 minutes
FDDI (100Mbps)	6 minutes
STS-3 (155Mbps)	4 minutes
STS-12 (622Mbps)	55 seconds
STS-24 (1.2Gbps)	28 seconds

Bytes in Transit: 16-bit **AdvertisedWindow**

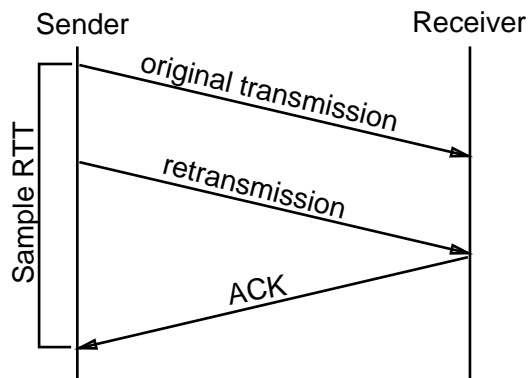
Bandwidth	Delay × Bandwidth Product
T1 (1.5Mbps)	18KB
Ethernet (10Mbps)	122KB
T3 (45Mbps)	549KB
FDDI (100Mbps)	1.2MB
STS-3 (155Mbps)	1.8MB
STS-12 (622Mbps)	7.4MB
STS-24 (1.2Gbps)	14.8MB

Adaptive Retransmission

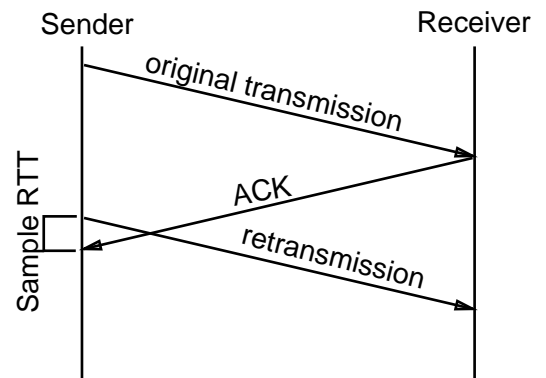
Original Algorithm

- Measure **SampleRTT** for each segment/ACK pair
- Compute weighted average of RTT
 - **EstimatedRTT** = $\alpha \times \text{EstimatedRTT} + \beta \times \text{SampleRTT}$
 - where $\alpha + \beta = 1$
 - α between 0.8 and 0.9
 - β between 0.1 and 0.2
- Set timeout based on **EstimatedRTT**
 - **TimeOut** = $2 \times \text{EstimatedRTT}$

Karn/Partridge Algorithm



(a) Sample RTT too long



(b) Sample RTT too short

- Do not sample RTT when retransmitting
- Double timeout after each retransmission

Jacobson/Karels Algorithm

- New calculation for average RTT

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

$\text{EstimatedRTT} = \text{EstimatedRTT} + (\delta \times \text{Difference})$

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

– where δ is a fraction between 0 and 1

- Consider variance when setting timeout value

$\text{Timeout} = \mu \times \text{EstimatedRTT} + \phi \times \text{Deviation}$

– where $\mu = 1$ and $\phi = 4$

- Notes

– algorithm only as good as granularity of clock (500ms on Unix)

– accurate timeout mechanism important to congestion control (later)

TCP Extensions

- Implemented as header options
- Store timestamp in outgoing segments
- Use 32-bit timestamp to extend sequence space (PAWS)
- Shift (scale) advertised window