

Reliable Byte-Stream (TCP)

- Outline
 - Connection Establishment/Termination
 - Sequence number selection
 - Connection tear-down
 - Round-trip estimation
 - Window flow control
 - Sliding Window Revisited
 - Adaptive Timeout

Slides courtesy:
 Ramesh Govindan @ USC
 Larry Peterson @ Princeton
 Jeffrey A. Six @ Delaware

Feb-18-03 4/598N: Computer Networks 1

End-to-End Protocols

- Underlying best-effort network
 - drop 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 flow control the sender
 - support multiple application processes on each host

Feb-18-03 4/598N: Computer Networks 2

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
- Header format
- Optional checksum
 - psuedo header + UDP header + data

0		16		31	
	SrcPort			DstPort	
	Checksum			Length	
Data					

Feb-18-03 4/598N: Computer Networks 3

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

Transmit segments

Feb-18-03 4/598N: Computer Networks 4

Data Link Versus Transport

- 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 node capacity
- Potentially different network capacity
 - need to be prepared for network congestion

Feb-18-03 4/598N: Computer Networks 5

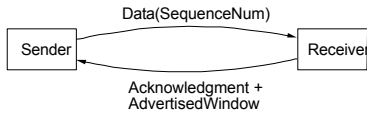
Segment Format

0		4		10	
	SrcPort			DstPort	
SequenceNum					
Acknowledgment					
HdrLen	0	Flags	AdvertisedWindow		
Checksum			UrgPtr		
Options (variable)					
Data					

Feb-18-03 4/598N: Computer Networks 6

Segment Format (cont)

- Each connection identified with 4-tuple:
 - (SrcPort, SrcIPAddr, DsrPort, DstIPAddr)
- Sliding window + flow control
 - acknowledgment, SequenceNum, AdvertisedWinow



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



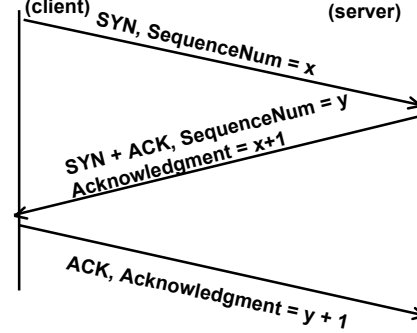
Feb-18-03

4/598N: Computer Networks

7

Connection Establishment and Termination

Active participant (client) Passive participant (server)



Feb-18-03

4/598N: Computer Networks

8

Sequence Number Selection

- Initial sequence number (ISN) selection
 - Why not simply chose 0?
 - Must avoid overlap with earlier incarnation
- Requirements for ISN selection
 - Must operate correctly
 - Without synchronized clocks
 - Despite node failures



Feb-18-03

4/598N: Computer Networks

9

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



Feb-18-03

4/598N: Computer Networks

10

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?

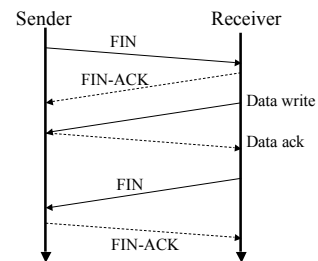


Feb-18-03

4/598N: Computer Networks

11

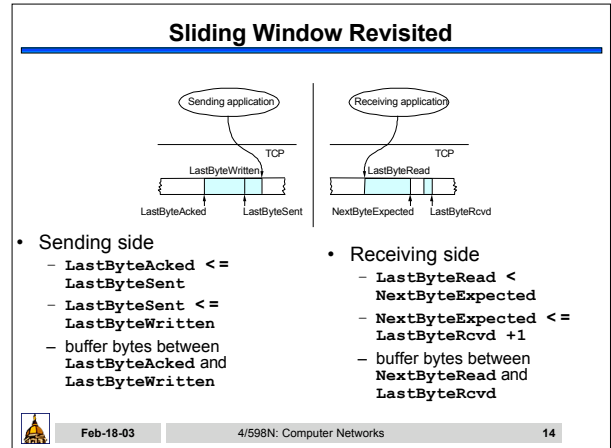
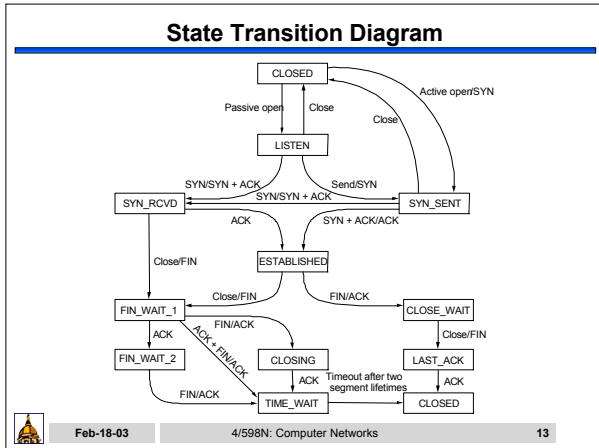
Tear-down Packet Exchange



Feb-18-03

4/598N: Computer Networks

12

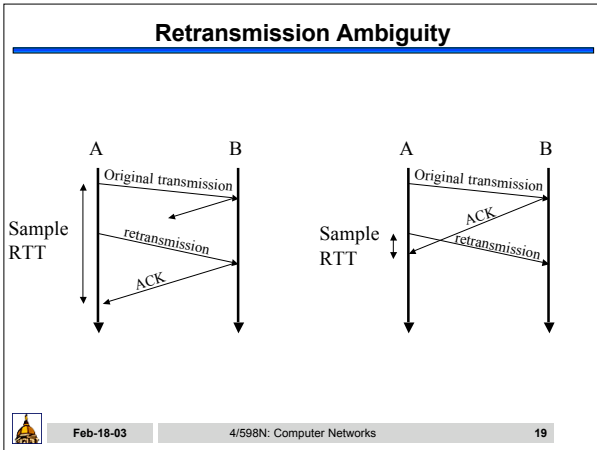


- ### Flow Control
- Fast sender can overrun receiver:
 - Packet loss, unnecessary retransmissions
 - Possible solutions:
 - Sender transmits at pre-negotiated rate
 - Sender limited to a window's worth of unacknowledged data
 - Flow control different from congestion control
- Feb-18-03 4/598N: Computer Networks 15

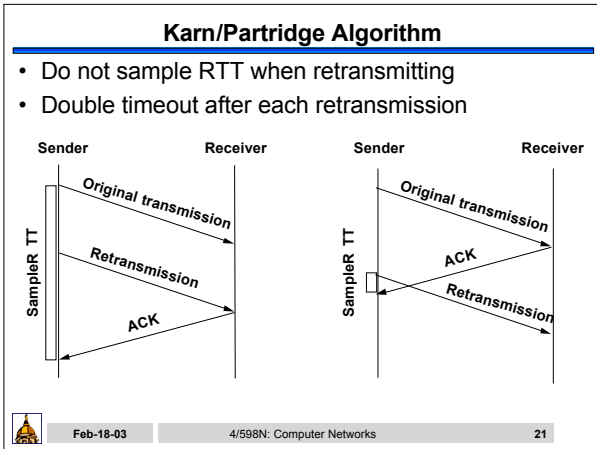
- ### Flow Control
- Send buffer size: MaxSendBuffer
 - Receive buffer size: MaxRcvBuffer
 - Receiving side
 - $LastByteRcvd - LastByteRead \leq MaxRcvBuffer$
 - $AdvertisedWindow = MaxRcvBuffer - (NextByteExpected - NextByteRead)$
 - Sending side
 - $LastByteSent - LastByteAcked \leq AdvertisedWindow$
 - $EffectiveWindow = AdvertisedWindow - (LastByteSent - LastByteAcked)$
 - $LastByteWritten - LastByteAcked \leq MaxSendBuffer$
 - block sender if $(LastByteWritten - LastByteAcked) + y > MaxSenderBuffer$
 - Always send ACK in response to arriving data segment
 - Persist when $AdvertisedWindow = 0$
- Feb-18-03 4/598N: Computer Networks 16

- ### 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!
- Feb-18-03 4/598N: Computer Networks 17

- ### Initial Round-trip Estimator
- Round trip times exponentially averaged:
- $New\ RTT = \alpha (old\ RTT) + (1 - \alpha) (new\ sample)$
 - Recommended value for α : 0.8 - 0.9
 - Retransmit timer set to αRTT , where $\alpha = 2$
 - Every time timer expires, RTO exponentially backed-off
- Feb-18-03 4/598N: Computer Networks 18



- ### 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
- Feb-18-03 4/598N: Computer Networks 20



- ### Jacobson's Retransmission Timeout Estimator
- Key observation:
 - Using \square RTT for timeout doesn't work
 - At high loads round trip variance is high
 - Solution:
 - If D denotes mean variation
 - Timeout = RTT + 4D
- Feb-18-03 4/598N: Computer Networks 22

- ### Jacobson/ Karels Algorithm
- New Calculations for average RTT
 - $Diff = SampleRTT - EstRTT$
 - $EstRTT = EstRTT + (d \times Diff)$
 - $Dev = Dev + d(|Diff| - Dev)$
 - where d is a factor between 0 and 1
 - Consider variance when setting timeout value
 - $TimeOut = m \times EstRTT + f \times Dev$
 - 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)
- Feb-18-03 4/598N: Computer Networks 23