

End-to-End Protocols

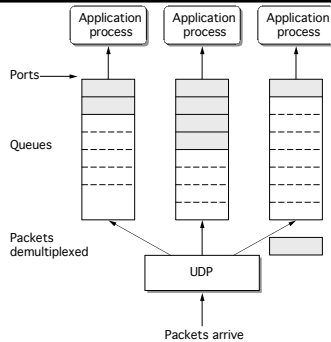
- Underlying best-effort network
 - drop messages
 - re-orders messages
 - delivers duplicate copies of a given message
 - limits messages to some finite size
 - delivers messages after an arbitrarily long delay
- Common end-to-end services
 - guarantee message delivery
 - deliver messages in the same 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



Mar-18-04

4/598N: Computer Networks

UDP message queue

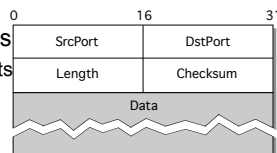


Mar-18-04

4/598N: Computer Networks

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

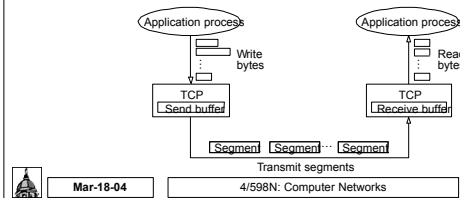


Mar-18-04

4/598N: Computer Networks

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



Mar-18-04

4/598N: Computer Networks

Data Link Versus Transport

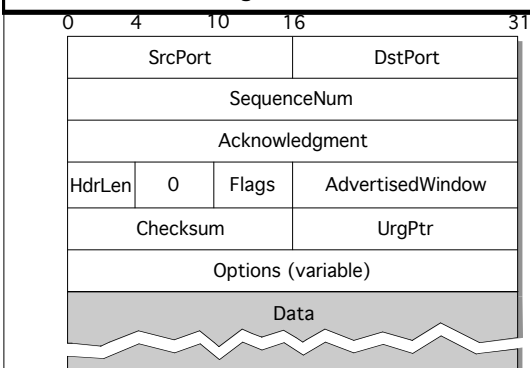
- Potentially connects many different hosts
 - need explicit connection establishment and termination
- Potentially different RTT
 - need adaptive timeout mechanism
- Potentially long delay in network
 - need to be prepared for arrival of very old packets
- Potentially different capacity at destination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



Mar-18-04

4/598N: Computer Networks

Segment Format

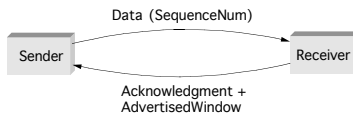


Mar-18-04

4/598N: Computer Networks

Segment Format (cont)

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



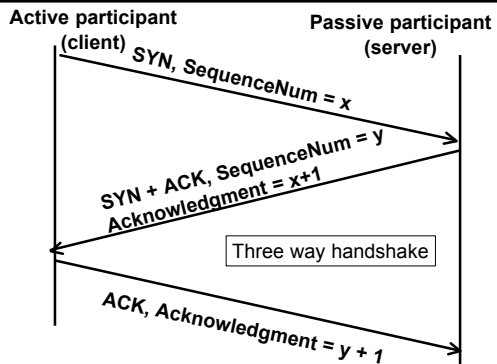
- Flags
 - SYN, FIN, RESET, PUSH, URG, ACK
- Checksum
 - pseudo header + TCP header + data



Mar-18-04

4/598N: Computer Networks

Connection Establishment and Termination



Mar-18-04

4/598N: Computer Networks

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



Mar-18-04

4/598N: Computer Networks

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



Mar-18-04

4/598N: Computer Networks

Connection Tear-down

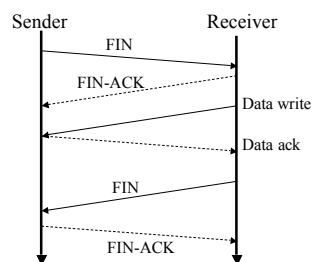
- Normal termination
 - Allow unilateral close
 - Avoid sequence number overlap
- TCP must continue to receive data even after closing
 - Cannot close connection immediately: what if a new connection restarts and uses same sequence number and receives retransmitted FIN from the current session?



Mar-18-04

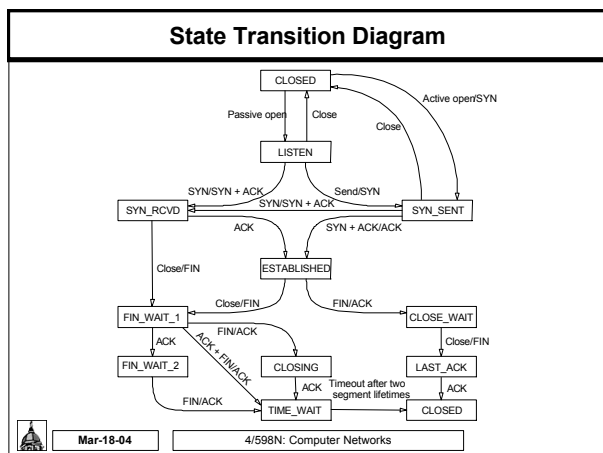
4/598N: Computer Networks

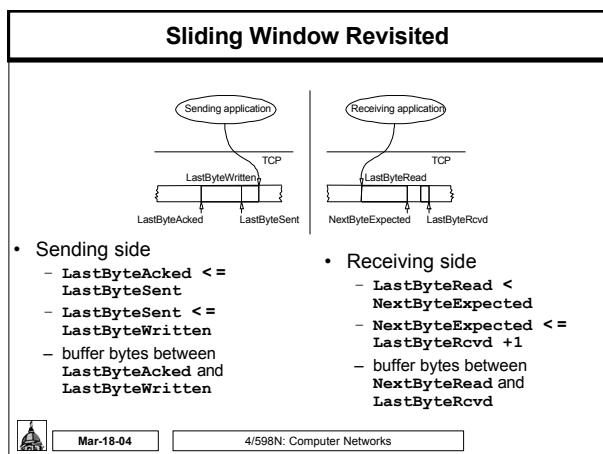
Tear-down Packet Exchange

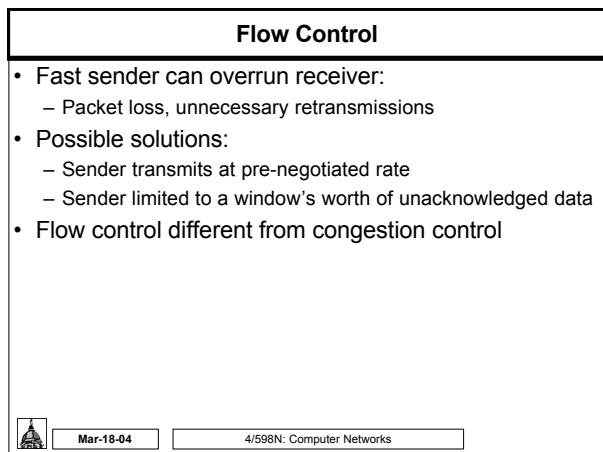


Mar-18-04

4/598N: Computer Networks







Flow Control

- Send buffer size: MaxSendBuffer
- Receive buffer size: MaxRcvBuffer
- Receiving side
 - $\text{LastByteRcvd} - \text{LastByteRead} \leq \text{MaxRcvBuffer}$
 - $\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{NextByteExpected} - \text{NextByteRead})$
- 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 arriving data segment
- Persist when $\text{AdvertisedWindow} = 0$



Mar-18-04

4/598N: Computer Networks

Round-trip Time Estimation

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



Mar-18-04

4/598N: Computer Networks

Initial Round-trip Estimator

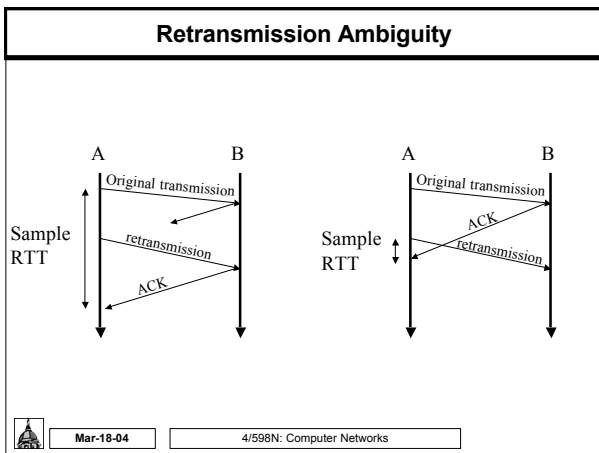
Round trip times exponentially averaged:

- New RTT = α (old RTT) + $(1 - \alpha)$ (new sample)
- Recommended value for α : 0.8 - 0.9
- Retransmit timer set to β RTT, where $\beta = 2$
- Every time timer expires, RTO exponentially backed-off



Mar-18-04

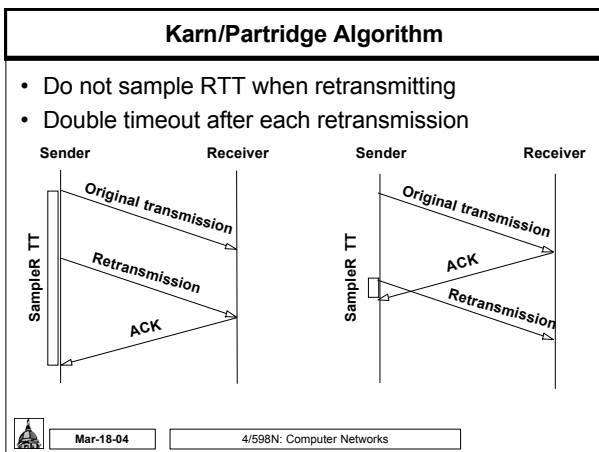
4/598N: Computer Networks



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

Mar-18-04 4/598N: Computer Networks



Jacobson's Retransmission Timeout Estimator

- Key observation:
 - Using β RTT for timeout doesn't work
 - At high loads round trip variance is high
- Solution:
 - If D denotes mean variation
 - Timeout = RTT + 4D



Mar-18-04

4/598N: Computer Networks

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)



Mar-18-04

4/598N: Computer Networks
