

Connection-oriented Transport: TCP

Outline:

- segment structure
- reliable data transfer
- flow control
- connection management

TCP: Overview

- point-to-point:
 - one sender, one receiver
- reliable data transfer:
 - guaranteed arrival, no error, in order
- pipelined:
 - multiple in-flight segments
- full duplex data:
 - bi-directional data flow in same connection
- connection-oriented:
 - handshaking (exchange of control msgs) to initialize sender, receiver state before data exchange
- flow controlled:
 - sender does not overwhelm receiver
- congestion controlled:
 - sender does not overwhelm the network
- no delay or bandwidth guarantee.

10/21/2002 CSC 257/457 - Fall 2002 2

TCP Segment Structure

The diagram shows a 32-bit TCP segment structure. The header fields are: source port #, dest port #, sequence number, acknowledgement number, Receive window, Urg data pointer, checksum, and Options (variable length). The application data is also variable length. Annotations include: URG: urgent data (generally not used); ACK: ACK # valid; PSH: push data now (generally not used); RST, SYN, FIN: connection estab (setup, teardown commands); Internet checksum (as in UDP); counting by bytes of data (not segments!); and # bytes rcvr willing to accept.

10/21/2002 CSC 257/457 - Fall 2002 3

Outline

- segment structure
- reliable data transfer
- flow control
- connection management

10/21/2002 CSC 257/457 - Fall 2002 4

TCP Reliable Data Transfer

- TCP provides reliable data transfer service on top of IP's unreliable service
- Pipelined transmissions & cumulative ACKs
- Looks like Go-back-N, however the receiver could buffer out-of-order segments
- TCP uses single retransmission timer
- Retransmissions are triggered by:
 - timeout events
 - duplicate ACKs

10/21/2002 CSC 257/457 - Fall 2002 5

TCP Timeout

Q: how to set TCP timeout value?

- longer than RTT (round trip time)
 - but RTT varies
- too short: premature timeout and unnecessary retransmissions
- too long: slow reaction to segment loss

10/21/2002 CSC 257/457 - Fall 2002 6

Estimating TCP Round Trip Time

Q: how to estimate RTT?

- **SampleRTT**: measured time from segment transmission until ACK receipt
- **SampleRTT** fluctuates, we want estimated RTT "smoother" to avoid short-term spikes
 - average several recent measurements, not just current **SampleRTT**
- we also want to give more recent measurements higher weight in case things do change

10/21/2002 CSC 257/457 - Fall 2002 7

EWMA - Exponentially Weighted Moving Average

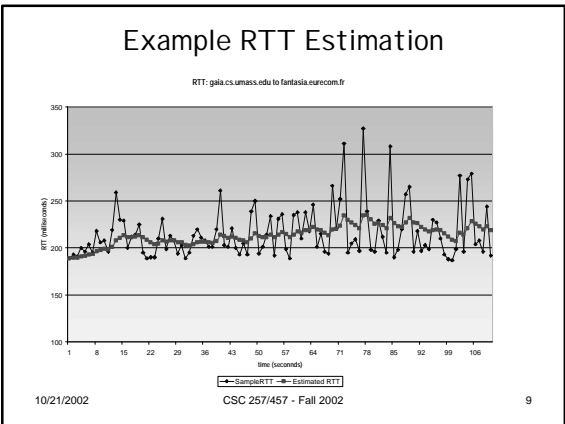
$$\text{EstimatedRTT} = a * \text{SampleRTT}_1 + a * (1-a) * \text{SampleRTT}_2 + a * (1-a)^2 * \text{SampleRTT}_3 + a * (1-a)^3 * \text{SampleRTT}_4 + \dots$$

SampleRTT₁ is the most recent sample TTL, **SampleRTT₂** is the next recent sample TTL, etc.

$$\text{EstimatedRTT} = (1-a) * \text{EstimatedRTT}_{\text{last}} + a * \text{SampleRTT}_1$$

- influence of past sample decreases exponentially fast
- typical value: a = 0.125

10/21/2002 CSC 257/457 - Fall 2002 8



TCP Timeout

Setting the timeout:

- **EstimatedRTT** plus "safety margin"
 - large variation in **EstimatedRTT** -> larger safety margin
- we need to estimate of how much **SampleRTT** deviates from **EstimatedRTT** (EWMA):

$$\text{DevRTT} = (1-b) * \text{DevRTT}_{\text{last}} + b * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically, b = 0.25)

Then set timeout interval:

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$

10/21/2002 CSC 257/457 - Fall 2002 10

TCP Sender Events and Processing

data rcvd from app:

- Create segment with seq #
- seq # is byte-stream number of first data byte in segment
- start timer if not already running (think of timer as for oldest un-ACKed segment)
- timeout vale: we just decided it!!

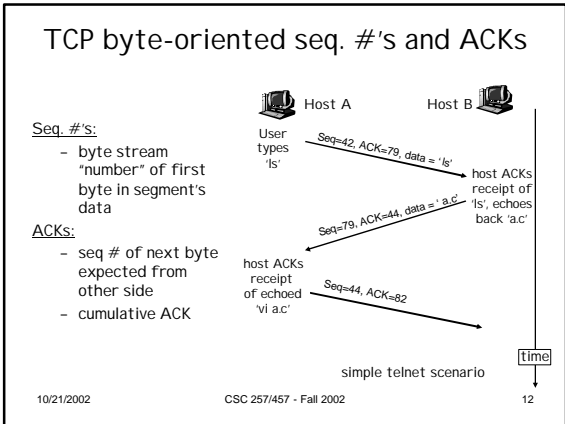
timeout:

- retransmit segment that caused timeout
- restart timer

ACK rcvd:

- If acknowledges previously unacked segments
 - slide sender window
 - start timer if there are outstanding segments

10/21/2002 CSC 257/457 - Fall 2002 11



TCP Flow Control: how it works

(Suppose TCP receiver discards out-of-order segments)

- spare room in buffer = **RcvWindow**
- = **RcvBuffer - [LastByteRcvd - LastByteRead]**

- Rcvr advertises spare room by including value of **RcvWindow** in segments
- Sender limits unACKed data to **RcvWindow**
 - guarantees receive buffer doesn't overflow

10/21/2002 CSC 257/457 - Fall 2002 19

Where is **RcvWindow** in each segment?

10/21/2002 CSC 257/457 - Fall 2002 20

Outline

- segment structure
- reliable data transfer
- flow control
- connection management

10/21/2002 CSC 257/457 - Fall 2002 21

TCP Connection Management

- Establishment:
 - TCP sender, receiver establish "connection" before exchanging data segments
 - initialize TCP variables: starting seq. #s, MSS, buffers, flow control info (e.g. **RcvWindow**)
- MSS is the maximum TCP segment size each side is willing to accept
 - typically the largest segment size fit into a link-layer frame
 - What is MSS for a IPv4 host connected with Ethernet?
- Teardown:
 - freeing up resources after mutually close

10/21/2002 CSC 257/457 - Fall 2002 22

TCP Connection Establishment

Three way handshake:

Step 1: client (active open) host sends TCP SYN segment to server

- specifies initial seq #
- no data

Step 2: server (passive open) host receives SYN, replies with SYNACK segment

- server allocates buffers
- specifies server initial seq. #

Step 3: client receives SYNACK, replies with ACK segment, which may contain data

10/21/2002 CSC 257/457 - Fall 2002 23

TCP Connection Teardown

Closing a connection:

close socket: `close(sockfd);`

Step 1: A (active closing host) sends TCP FIN control segment to server

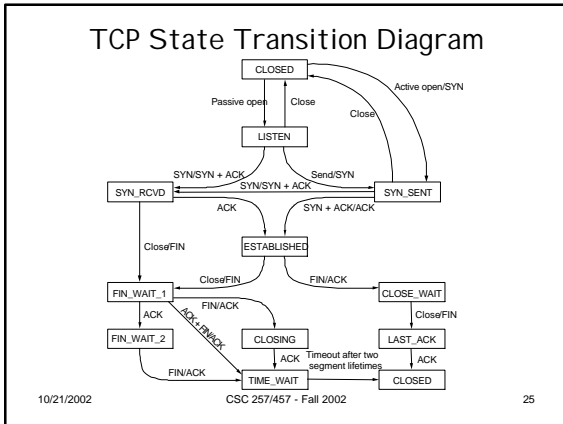
Step 2: B (passive closing host) receives FIN, replies with ACK. Closes connection, sends FIN.

Step 3: A receives FIN, replies with ACK.

- Enters "timed wait" - resend ACK in case it is lost

Step 4: B receives ACK. Connection closed.

10/21/2002 CSC 257/457 - Fall 2002 24



Disclaimer

- Parts of the lecture slides contain original work of James Kurose, Larry Peterson, and Keith Ross. The slides are intended for the sole purpose of instruction of computer networks at the University of Rochester. All copyrighted materials belong to their original owner(s).

10/21/2002 CSC 257/457 - Fall 2002 26