

TCP

Kai Shen

10/15/2014 CSC 257/457 - Fall 2014 1

TCP: Overview

- **Connection-oriented:**
 - handshaking (exchange of control msgs) to initialize sender, receiver state before data exchange
- **Reliable data transfer:**
 - guaranteed arrival, no error, in order
- **Flow controlled:**
 - sender does not overwhelm receiver
- **Congestion controlled:**
 - sender does not overwhelm the network
- **No delay or bandwidth guarantee.**

10/15/2014 CSC 257/457 - Fall 2014 2

TCP Segment Structure

32 bits

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

counting by bytes of data (not segments!)

bytes willing to receive

10/15/2014 CSC 257/457 - Fall 2014 3

Maximum Segment Size (MSS)

- MSS is the maximum TCP segment that'd fit into the link layer frame
- Local MSS
- Path MSS
 - Try and error probing

10/15/2014 CSC 257/457 - Fall 2014 4

TCP Reliable Data Transfer

- TCP provides reliable data transfer service on top of IP's unreliable service
 - Pipelined transmissions
 - Cumulative ACKs
 - When the receiver receives out-of-order segments, it buffers them and re-ACKs the last in-order data
 - Retransmit a single segment at each timeout
 - The sender retransmits at timeout or when receiving duplicate ACKs
- Somewhere between Go-back-N and Selective Repeat, closer to which?

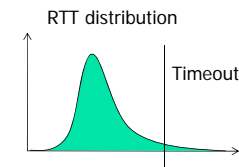
10/15/2014

CSC 257/457 - Fall 2014

5

TCP Timeout

- Q:** principles for setting transmission timeout value?
- too short: premature timeout and unnecessary retransmissions
 - too long: slow reaction to segment loss
 - longer than normal RTT (round trip time)
 - but RTT varies



10/15/2014

CSC 257/457 - Fall 2014

6

Expected Round Trip Time

Derive expected RTT from past RTT measurement.

- Basic measurement: measured time from segment transmission until ACK receipt
- **Stability:** RTT fluctuates, we want to avoid instability (pre-mature reaction to short-term spikes)
 - average several recent measurements, not just current RTT
- **Agility:** in case things do change, we want to adjust accordingly
 - give more recent measurements higher weight

10/15/2014

CSC 257/457 - Fall 2014

7

EWMA – Exponentially Weighted Moving Average

- influence of past samples decreases exponentially fast

$$\text{ExpectedRTT} = \frac{\text{SampleRTT}_1 + \alpha \cdot \text{SampleRTT}_2 + \alpha^2 \cdot \text{SampleRTT}_3 + \dots}{1 + \alpha + \alpha^2 + \dots}$$

SampleRTT_1 is RTT for the most recent data segment,
 SampleRTT_2 is RTT for the next recent data segment, etc.

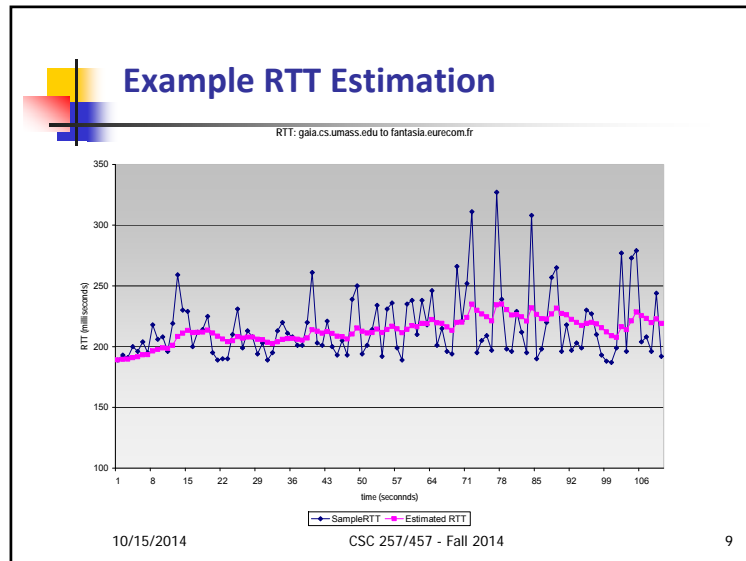
$$\Rightarrow \text{ExpectedRTT} = \alpha \cdot \text{ExpectedRTT}_{\text{last}} + (1 - \alpha) \cdot \text{SampleRTT}_1$$

- typical value in TCP: $\alpha = 0.875$

10/15/2014

CSC 257/457 - Fall 2014

8



TCP Timeout

Setting the timeout:

- ExpectedRTT plus “safety margin”
 - larger variation in RTT → larger safety margin
 - $\text{TimeoutInterval} = \text{ExpectedRTT} + 4 * \text{DevRTT}$
- DevRTT also estimated using EWMA of past measurements

10/15/2014 CSC 257/457 - Fall 2014 10

TCP Sender Events and Processing

Data ready to send:

- create segment with seq #
- seq # is byte-stream number of first data byte in segment
- start timer
- timeout value: we just decided it!!

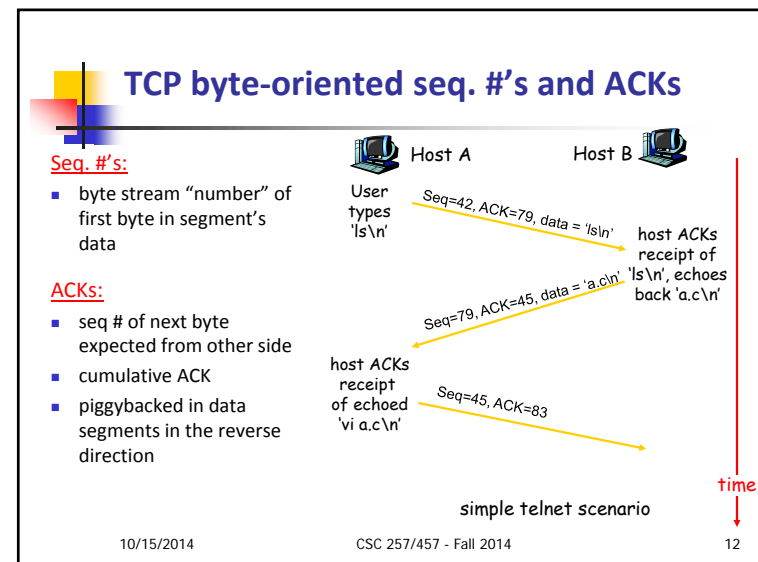
Timeout:

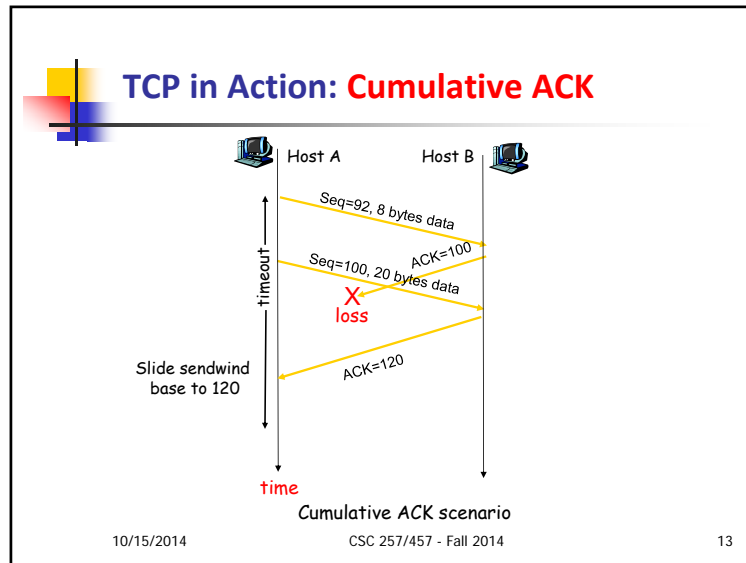
- retransmit segment that caused timeout
- restart timer

ACK rcvd:

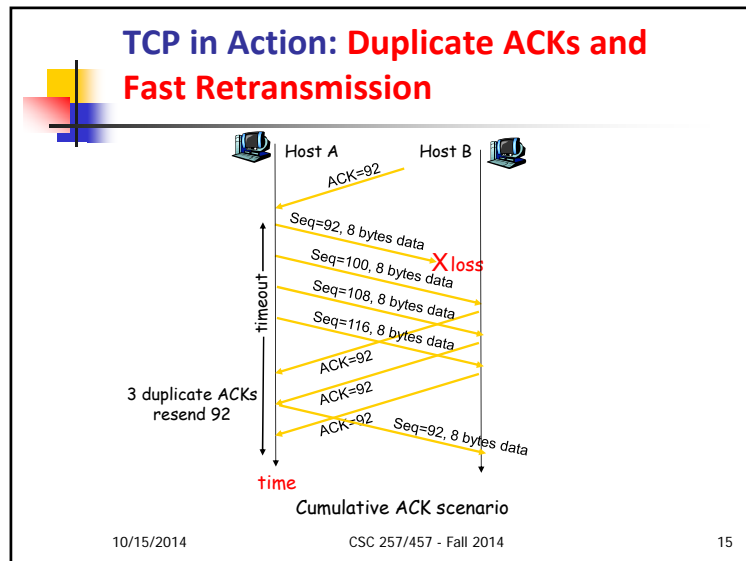
- slide sender window if acknowledges previously unacked segments
- retransmit if 3 duplicate ACKs

10/15/2014 CSC 257/457 - Fall 2014 11





- ### Fast Retransmission
- Time-out period often relatively long:
 - long delay before resending lost packet
 - When receiver receives out-of-order segments, it re-ACKs the last in-order byte
 - If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
 - **fast retransmission**: resend segment before timer expires, restart timer
- 10/15/2014 CSC 257/457 - Fall 2014 14



- ### Outline
- Segment structure
 - Reliable data transfer
 - **Flow control**
 - Connection management
- 10/15/2014 CSC 257/457 - Fall 2014 16

TCP Flow Control

- Receive side of TCP connection has a receive buffer:
- App process may be slow at reading from buffer.

Flow Control:

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

10/15/2014
CSC 257/457 - Fall 2014
17

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**)
- Teardown:**
 - Free up resources after mutually close

10/15/2014
CSC 257/457 - Fall 2014
18

TCP Connection Establishment

Three-way handshake:

- Step 1:** client (active open) 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

10/15/2014
CSC 257/457 - Fall 2014
19

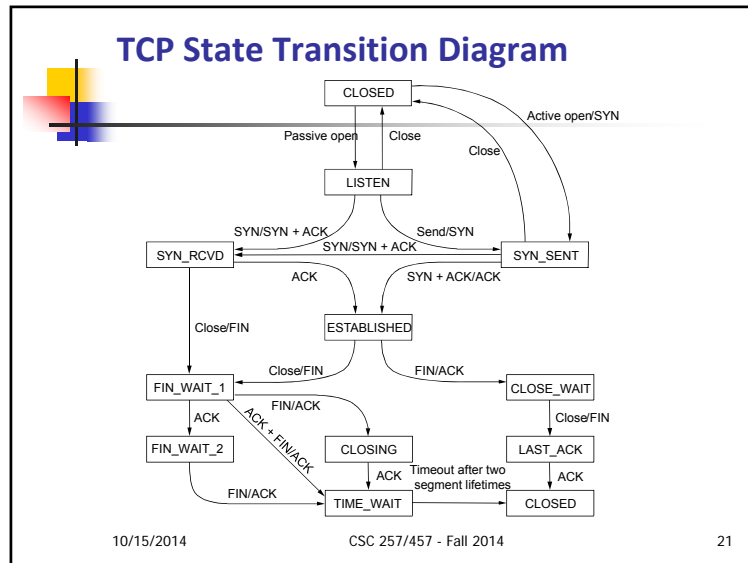
TCP Connection Teardown

Closing a connection:

close socket: `close(sockfd);`

- Step 1:** A (active closing host) sends TCP FIN control segment.
- 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.

10/15/2014
CSC 257/457 - Fall 2014
20



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).

10/15/2014 CSC 257/457 - Fall 2014 22