

CS162
Operating Systems and
Systems Programming
Lecture 17
TCP, Flow Control, Reliability

October 29, 2012
Ion Stoica
<http://inst.eecs.berkeley.edu/~cs162>

Goals for Today

- Reliable Transfer & flow control
- TCP
 - Open connection (3-way handshake)
 - Tear-down connection
 - Flow control

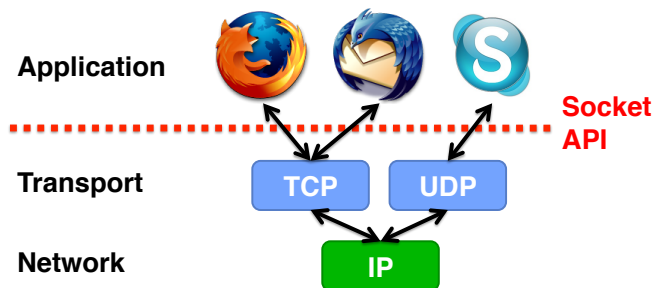
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.2

Socket API

- Socket API: Network programming interface



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.3

BSD Socket API

- Created at UC Berkeley (1980s)
- Most popular network API
- Ported to various OSes, various languages
 - Windows Winsock, BSD, OS X, Linux, Solaris, ...
 - Socket modules in Java, Python, Perl, ...
- Similar to Unix file I/O API
 - In the form of *file descriptor* (sort of handle).
 - Can share same `read()`/`write()`/`close()` system calls

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.4

TCP: Transport Control Protocol

- Reliable, in-order, and at most once delivery
- Stream oriented: messages can be of arbitrary length
- Provides multiplexing/demultiplexing to IP
- Provides congestion and flow control
- Application examples: file transfer, chat

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.5

TCP Service

- 1) Open connection: 3-way handshaking
- 2) Reliable byte stream transfer from (IPa, TCP_Port1) to (IPb, TCP_Port2)
 - Indication if connection fails: Reset
- 3) Close (tear-down) connection

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.6

Open Connection: 3-Way Handshaking

- Goal: agree on a set of parameters, i.e., the start sequence number for each side
 - Starting sequence number: sequence of first byte in stream
 - Starting sequence numbers are random

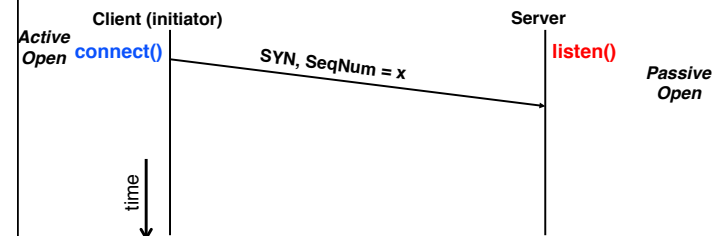
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.7

Open Connection: 3-Way Handshaking

- Server waits for new connection calling **listen()**
- Sender call **connect()** passing socket which contains server's IP address and port number
 - OS sends a special packet (SYN) containing a proposal for first sequence number, x



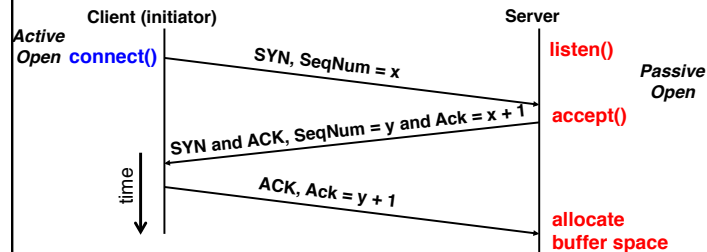
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.8

Open Connection: 3-Way Handshaking

- If it has enough resources, server calls `accept()` to accept connection, and sends back a SYN ACK packet containing
 - client's sequence number incremented by one, $(x + 1)$
 - » Why is this needed?
 - A sequence number proposal, y , for first byte server will send



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.9

3-Way Handshaking (cont' d)

- Three-way handshake adds 1 RTT delay
- Why?
 - Congestion control: SYN (40 byte) acts as cheap probe
 - Protects against delayed packets from other connection (would confuse receiver)

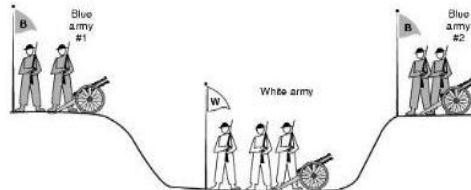
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.10

Close Connection (Two Generals Problem)

- Goal: both sides agree to close the connection
- Two-army problem:
 - “Two blue armies need to simultaneously attack the white army to win; otherwise they will be defeated. The blue army can communicate only across the area controlled by the white army which can intercept the messengers.”



- What is the solution?

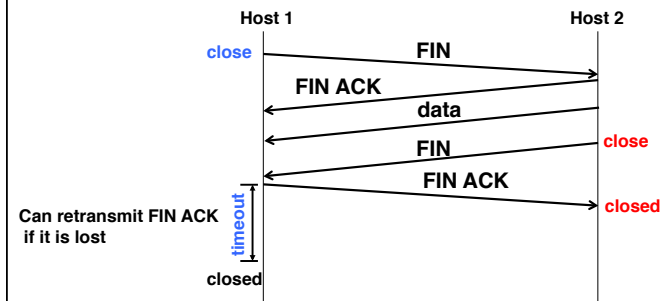
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.11

Close Connection

- 4-ways tear down connection



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.12

Reliable Transfer

- Retransmit missing packets
 - Numbering of packets and ACKs
- Do this efficiently
 - Keep transmitting whenever possible
 - Detect missing packets and retransmit quickly
- Two schemes
 - Stop & Wait
 - Sliding Window (Go-back-n and Selective Repeat)

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.13

Detecting Packet Loss?

- Timeouts
 - Sender timeouts on not receiving ACK
- Missing ACKs
 - Sender ACKs each packet
 - Receiver detects a missing packet when seeing a gap in the sequence of ACKs
 - Need to be careful! Packets and acks might be reordered
- NACK: Negative ACK
 - Receiver sends a NACK specifying a packet it is missing

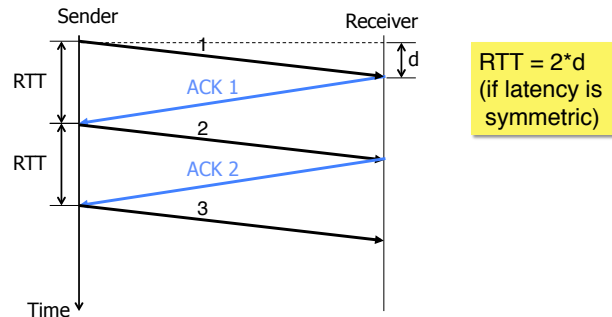
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.14

Stop & Wait w/o Errors

- Send; wait for ack; repeat
- RTT: Round Trip Time (RTT): time it takes a packet to travel from sender to receiver and back
 - One-way latency (d): one way delay from sender and receiver



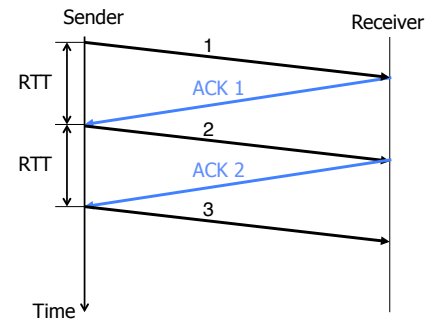
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.15

Stop & Wait w/o Errors

- How many packets can you send?
- 1 packet / RTT
- Throughput: number of bits delivered to receiver per sec



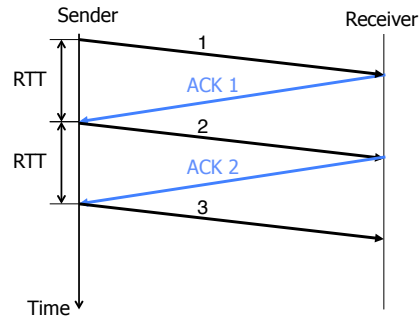
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.16

Stop & Wait w/o Errors

- Say, RTT = 100ms
- 1 packet = 1500 bytes
- Throughput = $1500 \cdot 8 \text{ bits} / 0.1 \text{ s} = 120 \text{ Kbps}$



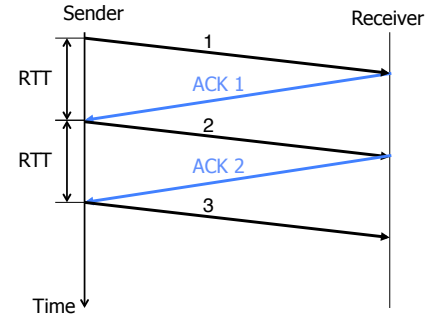
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.17

Stop & Wait w/o Errors

- Can be highly inefficient for high capacity links
- Throughput doesn't depend on the network capacity → even if capacity is 1 Gbps, we can only send 120 Kbps!



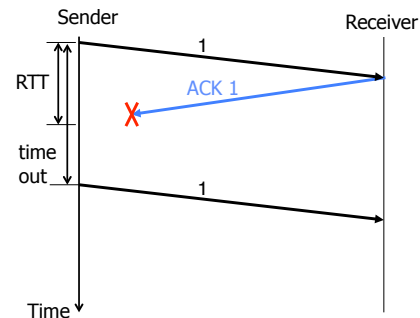
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.18

Stop & Wait with Errors

- If a loss wait for a retransmission timeout and retransmit
- How do you pick the timeout?



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.19

Sliding Window

- *window* = set of adjacent sequence numbers
- The size of the set is the *window size*
- Assume window size is n
- Let A be the last ack'd packet of sender without gap; then window of sender = $\{A+1, A+2, \dots, A+n\}$
- Sender can send packets in its window
- Let B be the last received packet without gap by receiver, then window of receiver = $\{B+1, \dots, B+n\}$
- Receiver can accept out of sequence, if in window

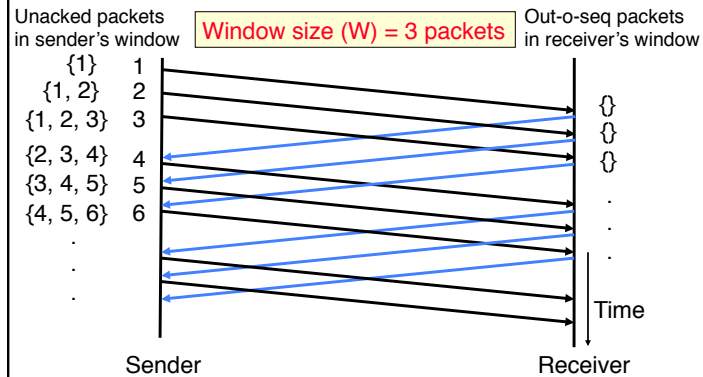
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.20

Sliding Window w/o Errors

- Throughput = $W \cdot \text{packet_size} / \text{RTT}$



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.21

Example: Sliding Window w/o Errors

- Assume
 - Link capacity, $C = 1\text{Gbps}$
 - Latency between end-hosts, $\text{RTT} = 80\text{ms}$
 - packet_length = 1000 bytes
- What is the window size W to match link's capacity, C ?
- Solution
 - We want Throughput = C
 - Throughput = $W \cdot \text{packet_size} / \text{RTT}$
 - $C = W \cdot \text{packet_size} / \text{RTT}$
 - $W = C \cdot \text{RTT} / \text{packet_size} = 10^9\text{bps} \cdot 80 \cdot 10^{-3}\text{s} / (8000\text{b}) = 10^4$ packets

Window size \sim Bandwidth (Capacity), delay ($\text{RTT}/2$)

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.22

Sliding Window with Errors

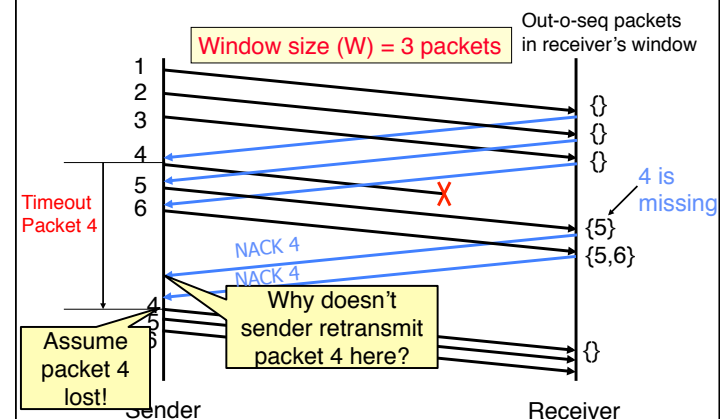
- Two approaches
 - Go-Back- n (GBN)
 - Selective Repeat (SR)
- In the absence of errors they behave identically
- Go-Back- n (GBN)
 - Transmit up to n unacknowledged packets
 - If timeout for $\text{ACK}(k)$, retransmit $k, k+1, \dots$
 - Typically uses NACKs instead of ACKs
 - Recall, NACK specifies first in-sequence packet missed by receiver

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.23

GBN Example with Errors



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.24

Selective Repeat (SR)

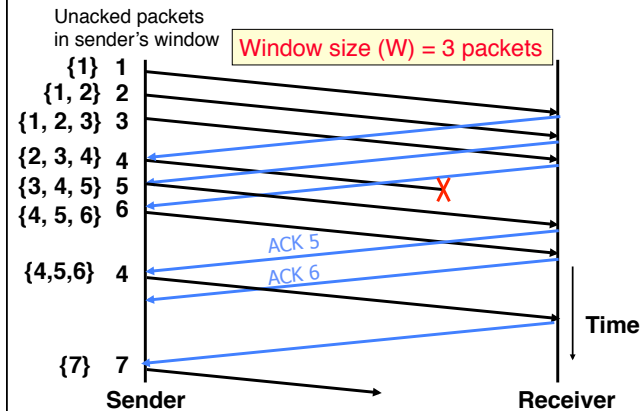
- Sender: transmit up to n unacknowledged packets;
- Assume packet k is lost
- Receiver: indicate packet k is missing (use ACKs)
- Sender: retransmit packet k

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.25

SR Example with Errors



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.26

Summary

- TCP: Reliable Byte Stream
 - Open connection (3-way handshaking)
 - Close connection: no perfect solution; no way for two parties to agree in the presence of arbitrary message losses (Two General problem)
- Reliable transmission
 - S&W not efficient for links with large capacity (bandwidth) delay product
 - Sliding window more efficient but more complex

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.27

Announcements

- I'll no longer be away this Wednesday, so I'll be teaching the lecture
- Project 3 initial design due: **Thursday, Nov 1**

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.28

5min Break

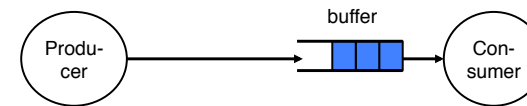
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.29

Flow Control

- Recall: Flow control ensures a fast sender does not overwhelm a slow receiver
- Example: Producer-consumer with bounded buffer (Lecture 5)
 - A buffer between producer and consumer
 - Producer puts items into buffer as long as buffer **not full**
 - Consumer consumes items from buffer



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.30

TCP Flow Control

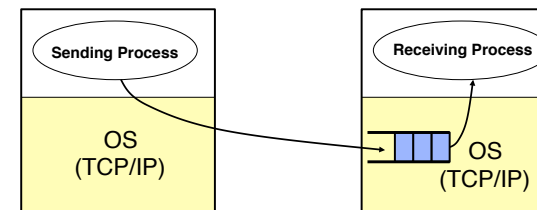
- TCP: sliding window protocol at byte (not packet) level
 - Go-back-N: TCP Tahoe, Reno, New Reno
 - Selective Repeat (SR): TCP Sack
- Receiver tells sender how many more bytes it can receive without overflowing its buffer (i.e., AdvertisedWindow)
- The ACK contains sequence number N of **next byte the receiver expects**, i.e., receiver has received all bytes in **sequence** up to and including N-1

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.31

TCP Flow Control

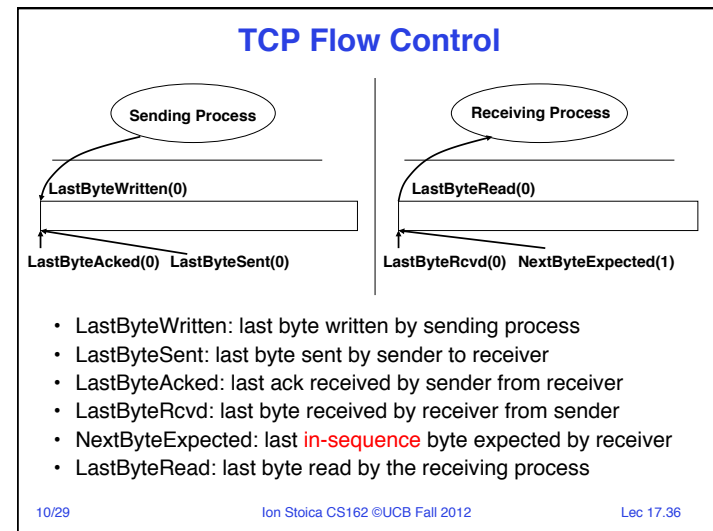
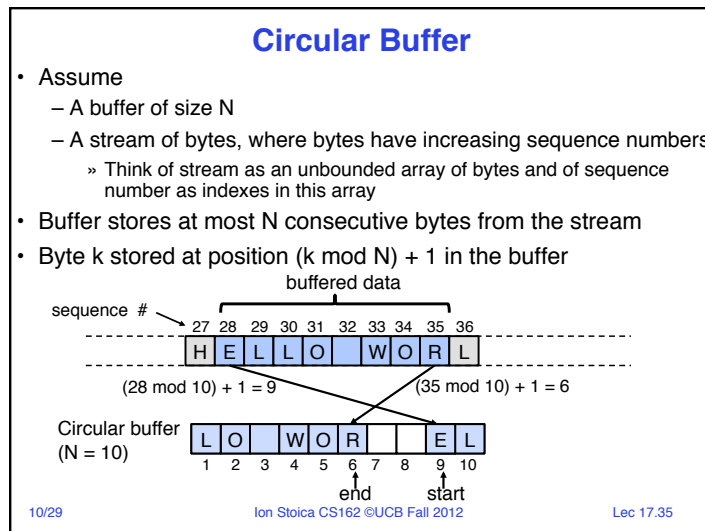
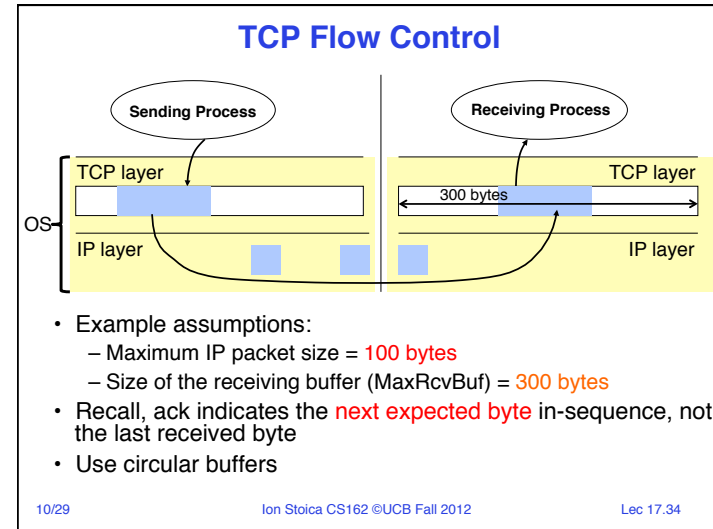
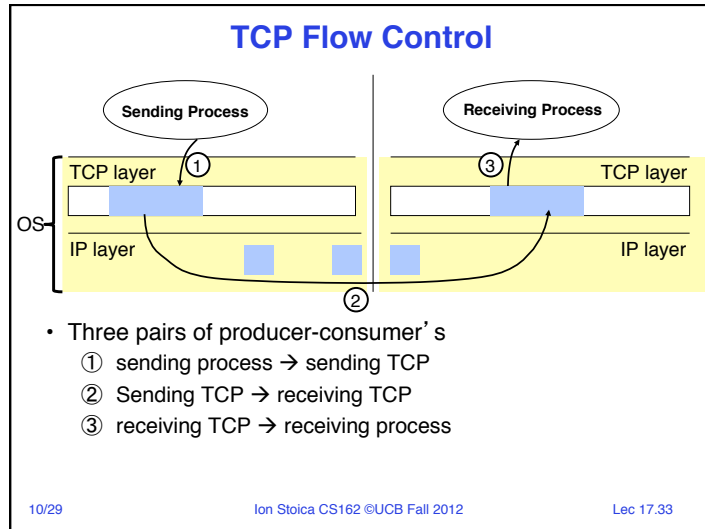


- TCP/IP implemented by OS (Kernel)
 - Cannot do context switching on sending/receiving every packet
 - » At 1Gbps, it takes 12 usec to send a 1500 bytes, and 0.8usec to send a 100 byte packet
- Need buffers to match ...
 - sending app with sending TCP
 - receiving TCP with receiving app

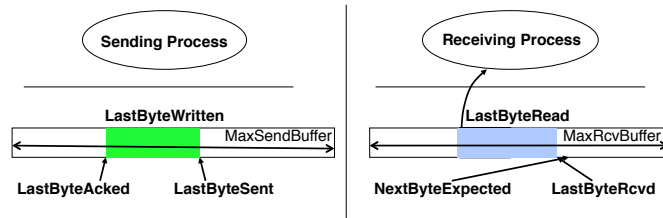
10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.32



TCP Flow Control



- AdvertisedWindow: number of bytes TCP receiver can receive

$$\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$

- SenderWindow: number of bytes TCP sender can send

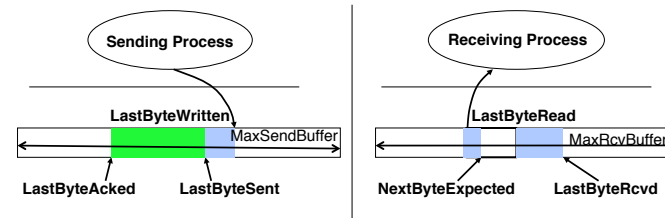
$$\text{SenderWindow} = \text{AdvertisedWindow} - (\text{LastByteSent} - \text{LastByteAcked})$$

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.37

TCP Flow Control



- Still true if receiver missed data....

$$\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$

- WriteWindow: number of bytes sending process can write

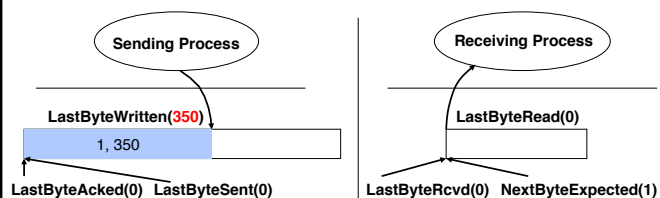
$$\text{WriteWindow} = \text{MaxSendBuffer} - (\text{LastByteWritten} - \text{LastByteAcked})$$

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.38

TCP Flow Control



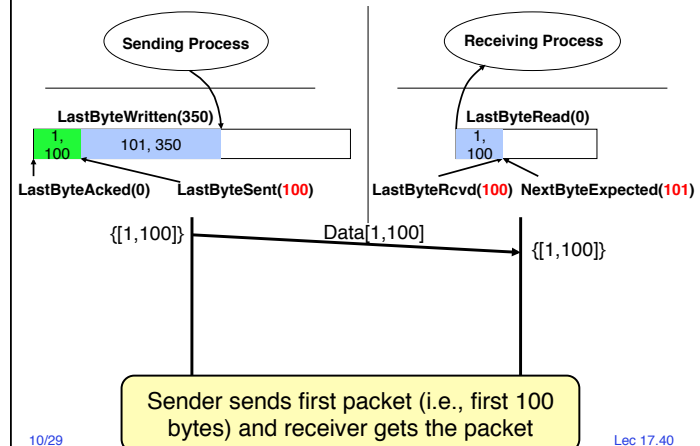
- Sending app sends 350 bytes
- Recall:
 - We assume IP only accepts packets no larger than 100 bytes
 - MaxRcvBuf = 300 bytes, so initial Advertised Window = 300 bytes

10/29

Ion Stoica CS162 ©UCB Fall 2012

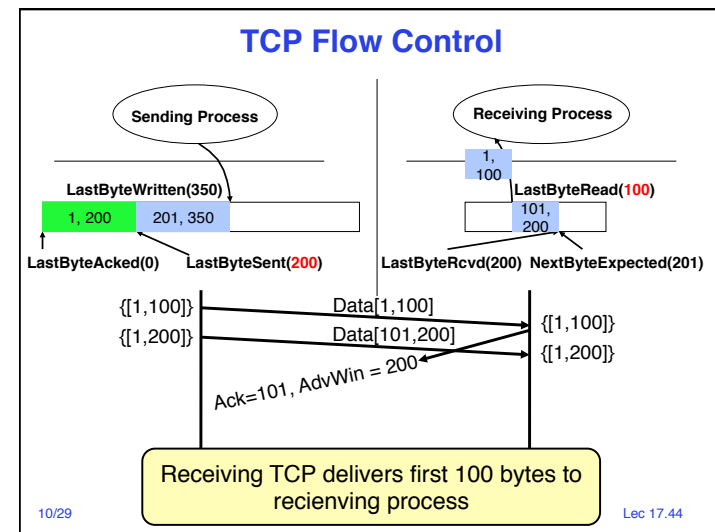
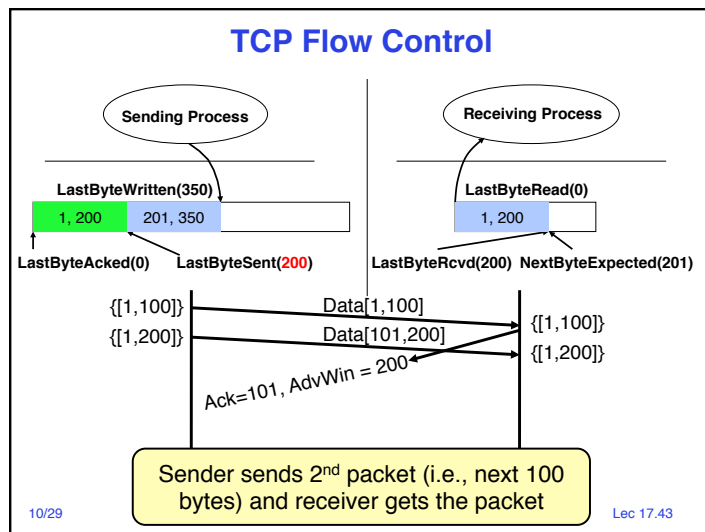
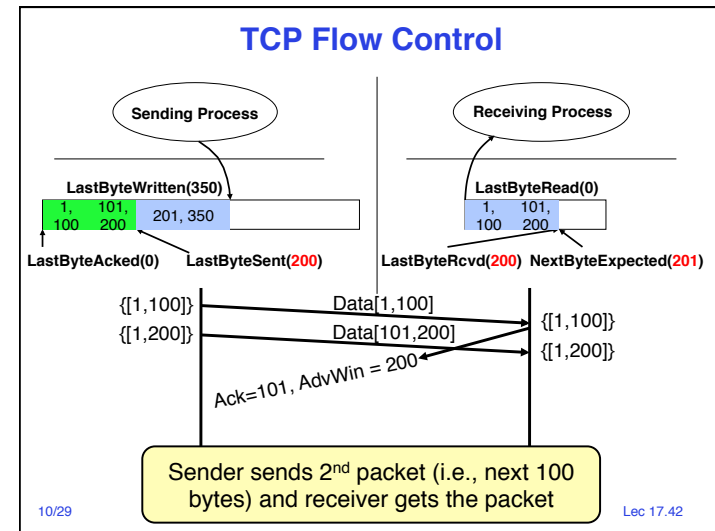
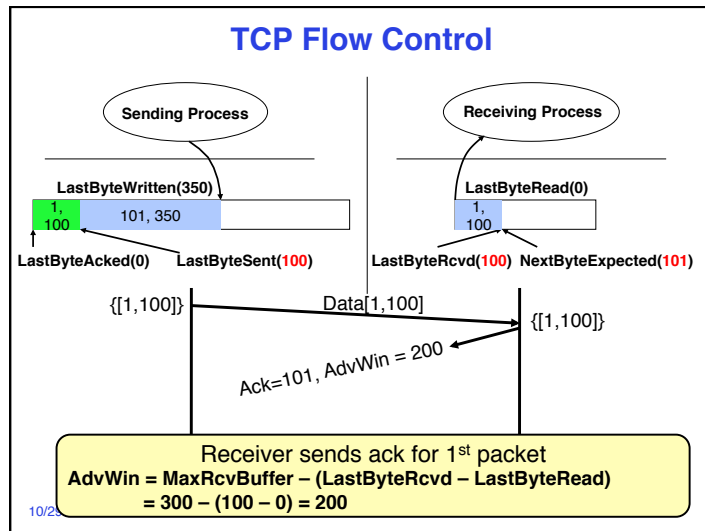
Lec 17.39

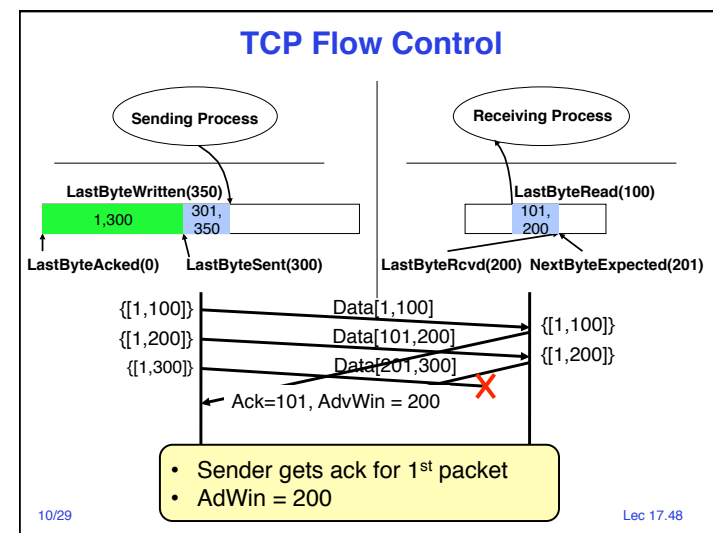
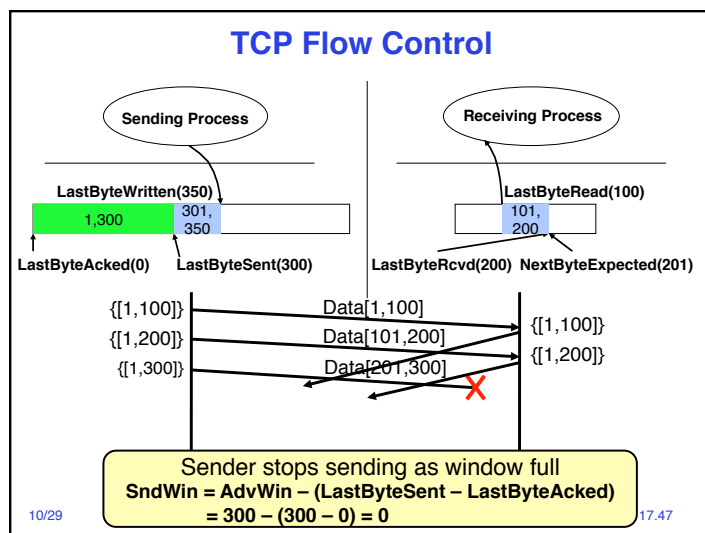
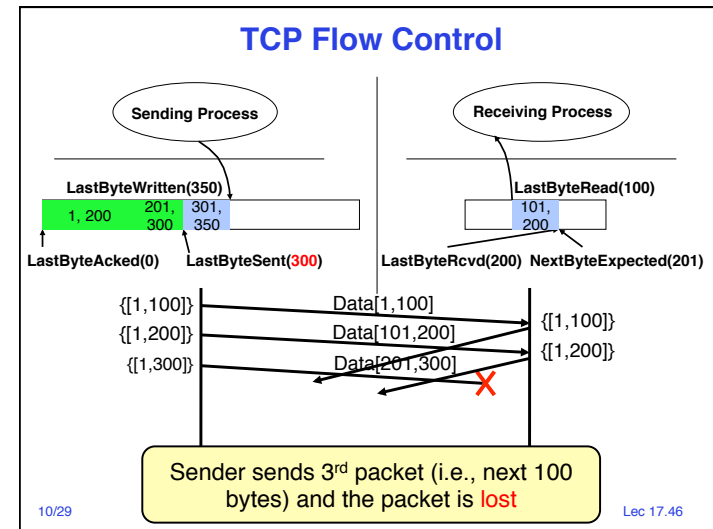
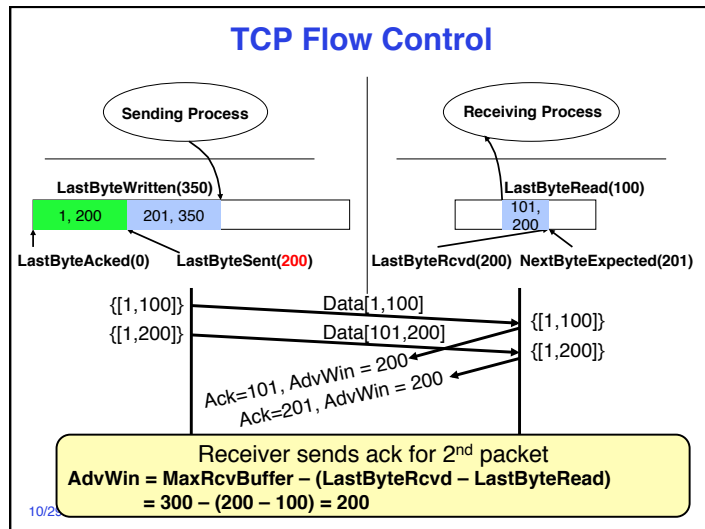
TCP Flow Control



10/29

Lec 17.40





TCP Flow Control



TCP Flow Control

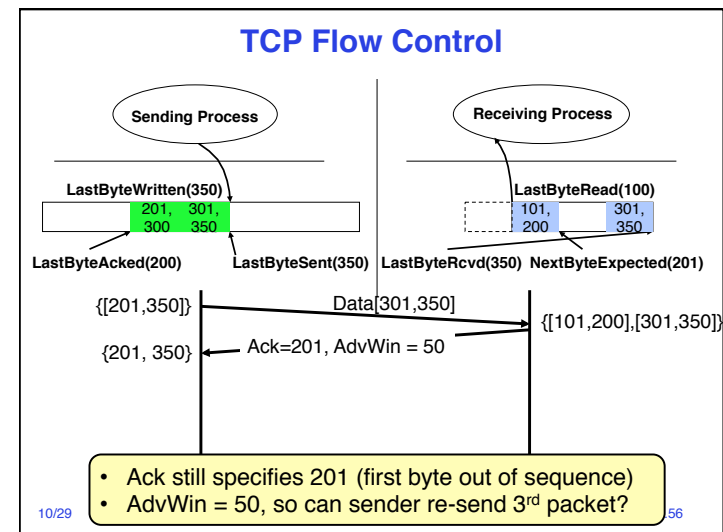
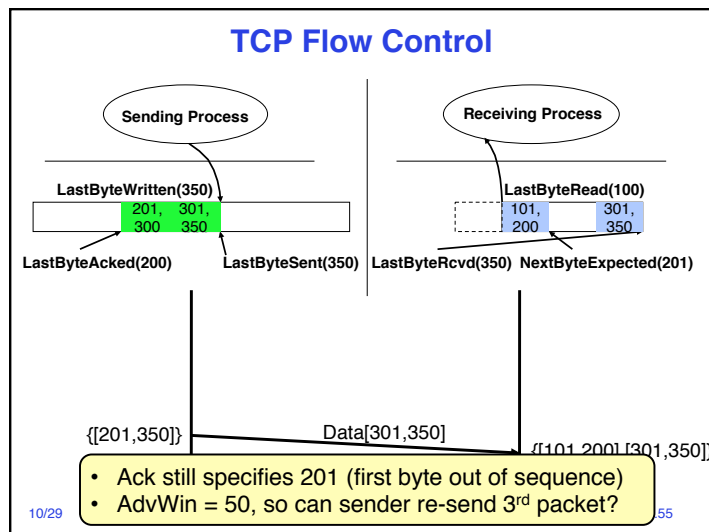
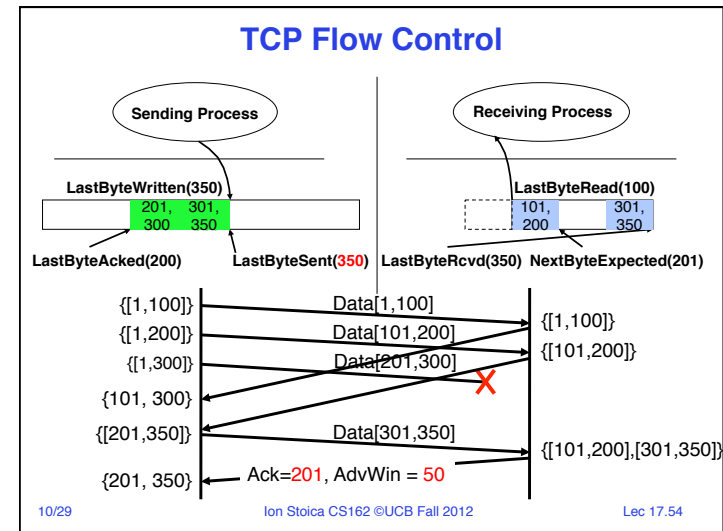
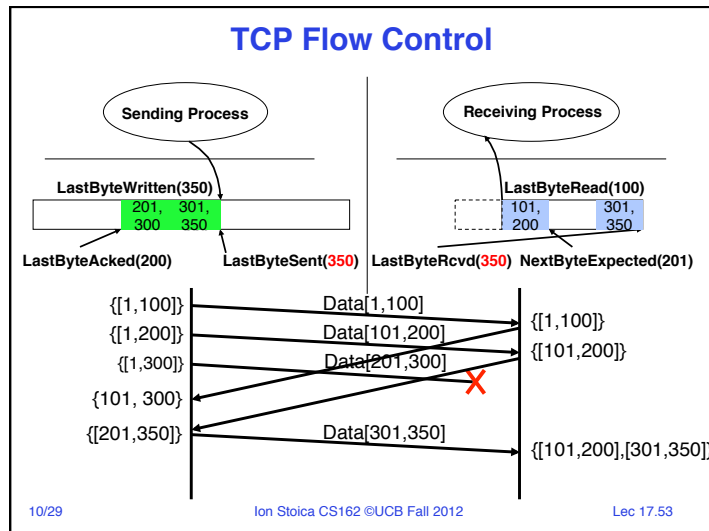


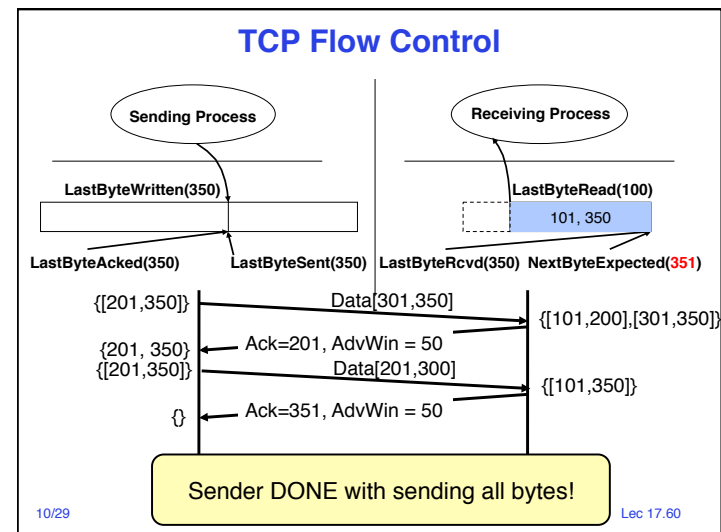
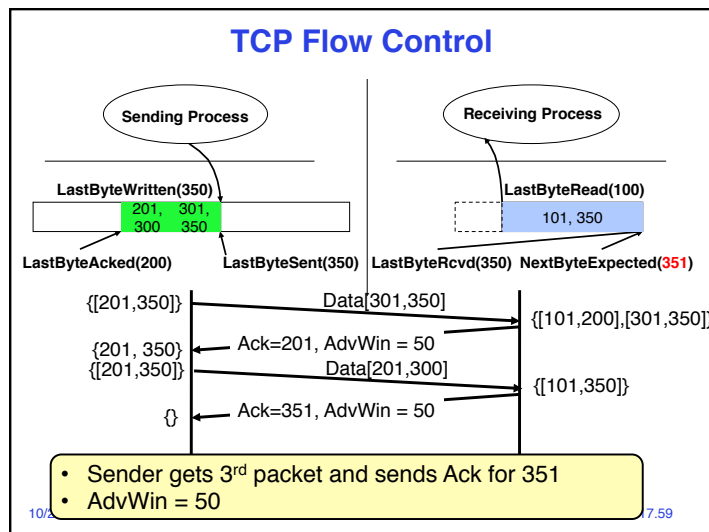
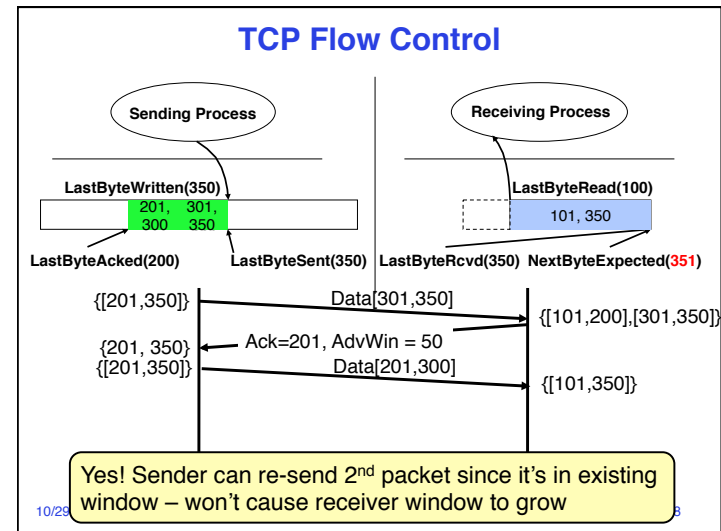
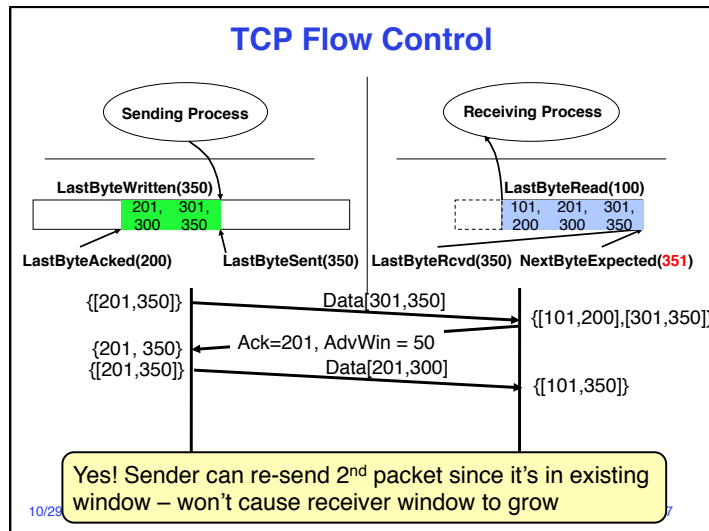
TCP Flow Control



TCP Flow Control







Discussion

- Why not have a huge buffer at the receiver (memory is cheap!)?
- Sending window (SndWnd) also depends on network congestion
 - **Congestion control**: ensure that a fast receiver doesn't overwhelm a router in the network (discussed in detail in ee122)
- In practice there is another set of buffers in the protocol stack, at the **link layer** (i.e., Network Interface Card)

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.61

Summary: Reliability & Flow Control

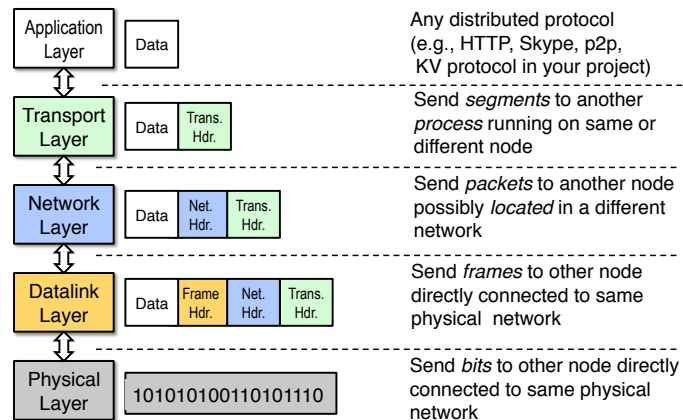
- Flow control: three pairs of producer consumers
 - Sending process → sending TCP
 - Sending TCP → receiving TCP
 - Receiving TCP → receiving process
- Reliable transmission
 - S&W not efficient for links with large capacity (bandwidth) delay product
 - Sliding window far more efficient
- TCP: Reliable Byte Stream
 - Open connection (3-way handshaking)
 - Close connection: no perfect solution; no way for two parties to agree in the presence of arbitrary message losses (Two General problem)

10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.62

Summary: Networking (Internet Layering)



10/29

Ion Stoica CS162 ©UCB Fall 2012

Lec 17.63