

Last time

- Reliable Data Transfer with packet loss
 - ◆ rdt3.0

- Pipelining
 - ◆ Problems with stop-and-wait
 - ◆ Go-Back-N
 - ◆ Selective-Repeat

This time

- UDP socket programming
- TCP

Chapter 2: Application layer

- 2.1 Principles of network applications
- 2.2 Web and HTTP
- 2.3 FTP
- 2.4 Electronic Mail
 - ◆ SMTP, POP3, IMAP
- 2.5 DNS
- 2.6 P2P file sharing
- 2.7 Socket programming with TCP
- 2.8 Socket programming with UDP
- 2.9 Building a Web server

Socket programming

Goal: learn how to build client/server application that communicate using sockets

Socket API

- introduced in BSD4.1 UNIX, 1981
- explicitly created, used, released by apps
- client/server paradigm
- two types of transport service via socket API:
 - ◆ unreliable datagram
 - ◆ reliable, byte stream-oriented

socket

a *host-local*, *application-created*, *OS-controlled* interface (a “door”) into which application process can **both send and receive** messages to/from another application process

Socket programming *with UDP*

UDP: no “connection” between client and server

- no handshaking
- sender explicitly attaches IP address and port of destination to each packet
- server must extract IP address, port of sender from received packet

UDP: transmitted data may be received out of order, or lost

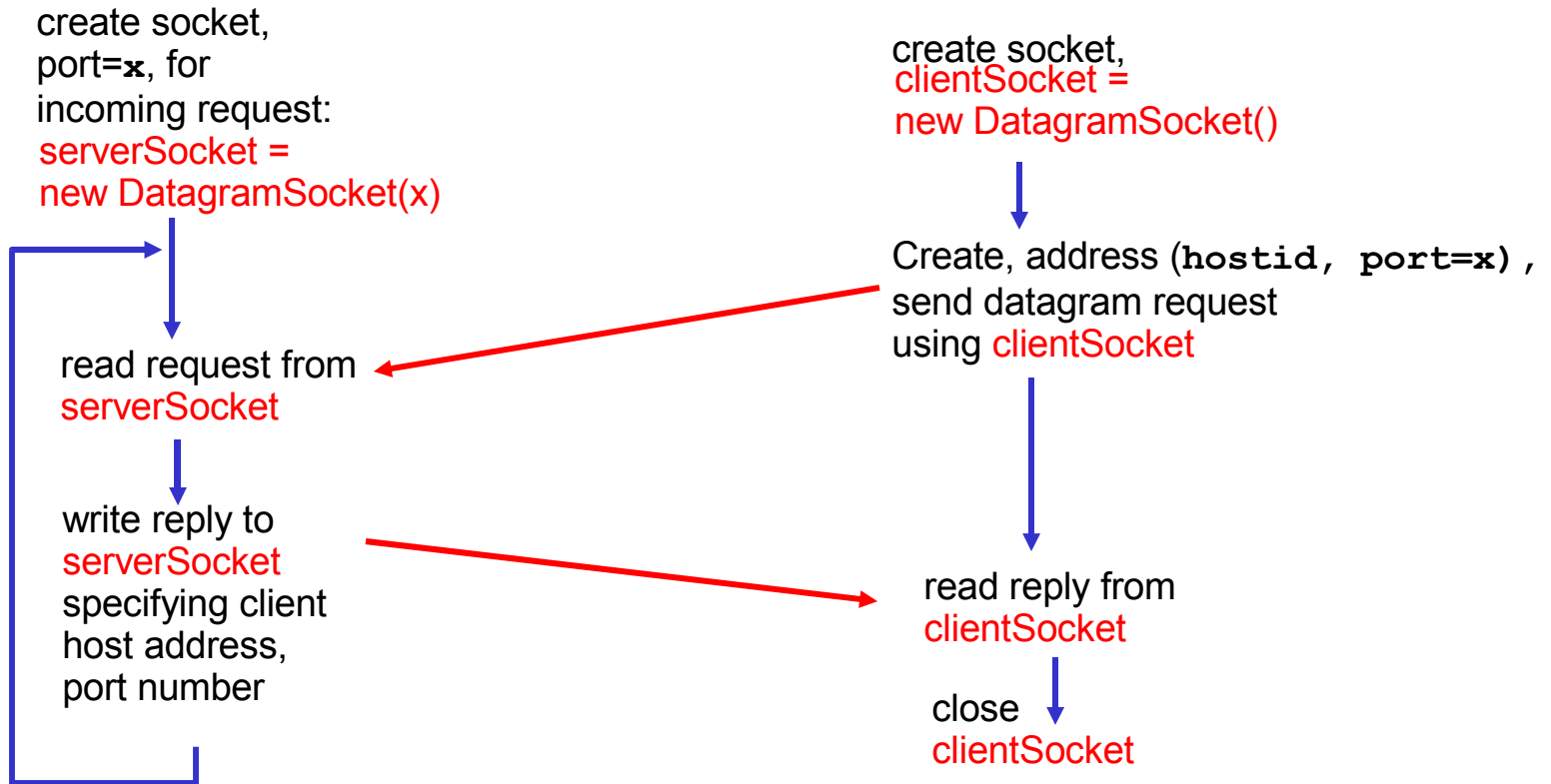
application viewpoint

UDP provides unreliable transfer of groups of bytes (“datagrams”) between client and server

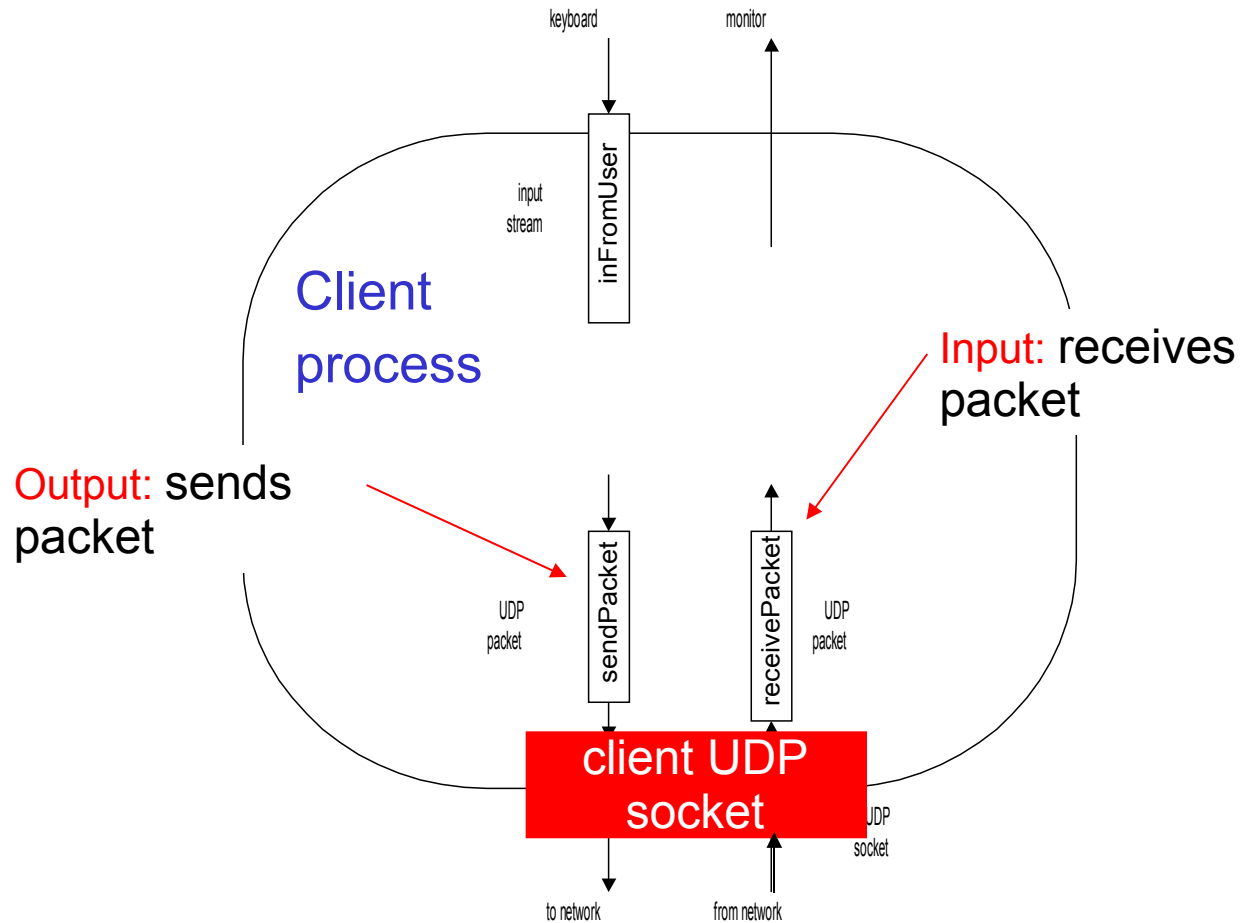
Client/server socket interaction: UDP

Server (running on `hostid`)

Client



Example: Java client (UDP)



Example: Java client (UDP)

```
import java.io.*;  
import java.net.*;
```

```
class UDPClient {  
    public static void main(String args[]) throws Exception  
    {
```

Create
input stream

```
        BufferedReader inFromUser =
```

Create
client socket

```
            new BufferedReader(new InputStreamReader(System.in));
```

```
        DatagramSocket clientSocket = new DatagramSocket();
```

Translate
hostname to IP
address using DNS

```
        InetAddress IPAddress = InetAddress.getByName("hostname");
```

```
        byte[] sendData = new byte[1024];
```

```
        byte[] receiveData = new byte[1024];
```

```
        String sentence = inFromUser.readLine();
```

```
        sendData = sentence.getBytes();
```


Example: Java client (UDP), cont.

Create datagram with
data-to-send,
length, IP addr, port

```
DatagramPacket sendPacket =  
    new DatagramPacket(sendData, sendData.length, IPAddress, 9876);
```

Send datagram
to server

```
clientSocket.send(sendPacket);
```

```
DatagramPacket receivePacket =  
    new DatagramPacket(receiveData, receiveData.length);
```

Read datagram
from server

```
clientSocket.receive(receivePacket);
```

```
String modifiedSentence =  
    new String(receivePacket.getData());
```

```
System.out.println("FROM SERVER:" + modifiedSentence);  
clientSocket.close();  
}
```

```
}
```

Example: Java server (UDP)

```
import java.io.*;  
import java.net.*;
```

```
class UDPServer {  
    public static void main(String args[]) throws Exception  
    {
```

Create
datagram socket
at port 9876



```
        DatagramSocket serverSocket = new DatagramSocket(9876);
```

```
        byte[] receiveData = new byte[1024];  
        byte[] sendData = new byte[1024];
```

```
        while(true)  
        {
```

Create space for
received datagram



```
            DatagramPacket receivePacket =  
                new DatagramPacket(receiveData, receiveData.length);
```

Receive
datagram



```
            serverSocket.receive(receivePacket);
```

Example: Java server (UDP), cont

String sentence = new String(receivePacket.getData());

Get IP addr
port #, of
sender

→ InetAddress IPAddress = receivePacket.getAddress();

→ int port = receivePacket.getPort();

String capitalizedSentence = sentence.toUpperCase();

sendData = capitalizedSentence.getBytes();

Create datagram
to send to client

→ DatagramPacket sendPacket =
new DatagramPacket(sendData, sendData.length, IPAddress,
port);

Write out
datagram
to socket

→ serverSocket.send(sendPacket);

}

}

}

End of while loop,
loop back and wait for
another datagram

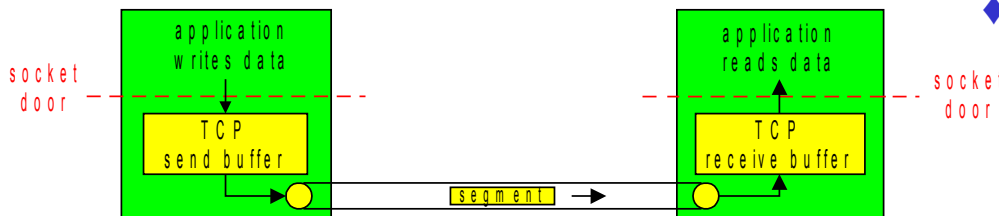
Chapter 3 outline

- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer
- 3.5 Connection-oriented transport: TCP
 - ◆ segment structure
 - ◆ reliable data transfer
 - ◆ flow control
 - ◆ connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

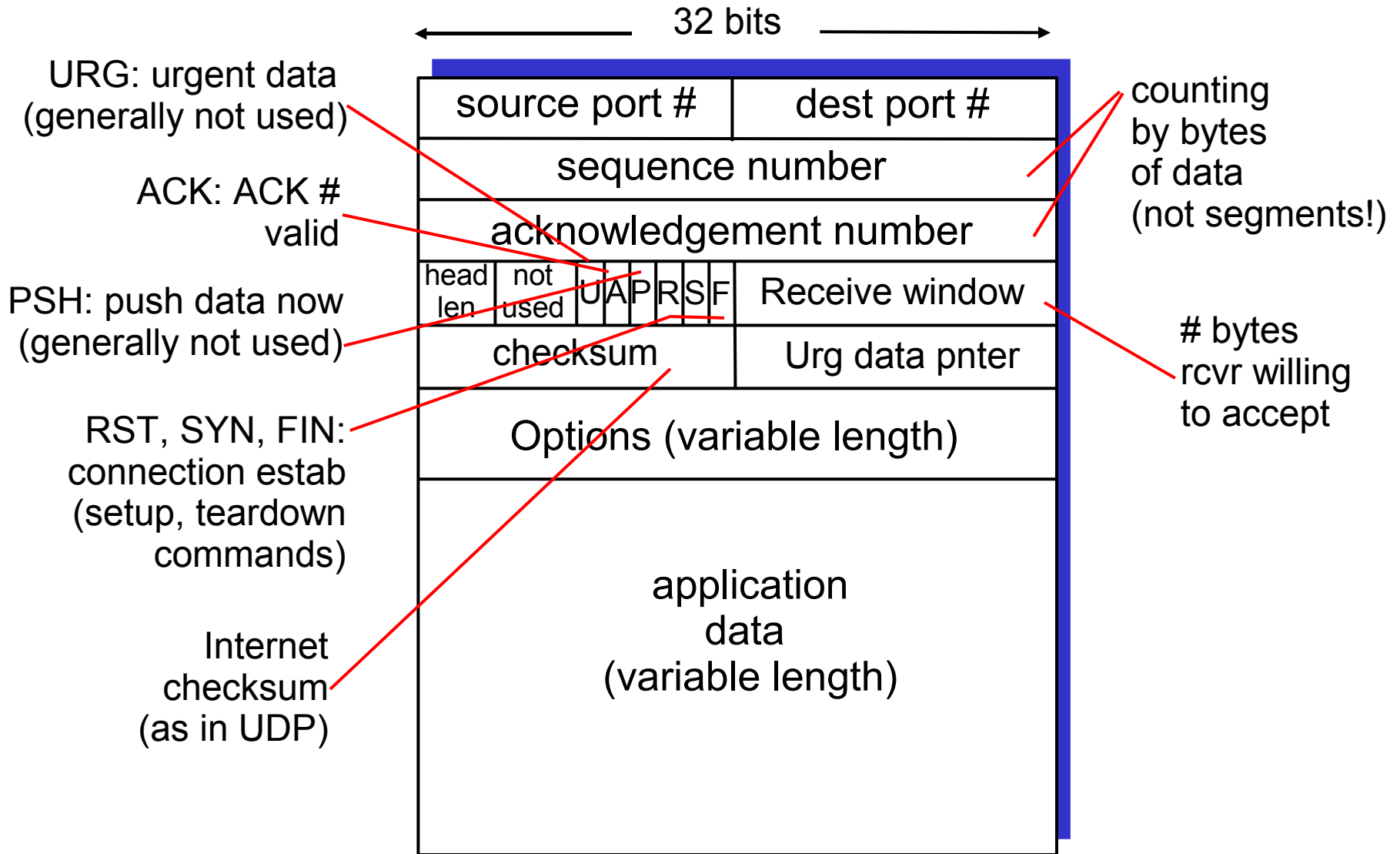
TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

- **Point-to-point:**
 - ◆ one sender, one receiver
- **Reliable, in-order *byte stream*:**
 - ◆ no “message boundaries”
- **Pipelined:**
 - ◆ TCP congestion and flow control set window size
- **Send & receive buffers**
- **Full duplex data:**
 - ◆ bi-directional data flow in same connection
 - ◆ MSS: maximum segment size
- **Connection-oriented:**
 - ◆ handshaking (exchange of control msgs) inits sender, receiver state before data exchange
- **Flow controlled:**
 - ◆ sender will not overwhelm receiver



TCP segment structure



TCP seq. #'s and ACKs

Seq. #'s:

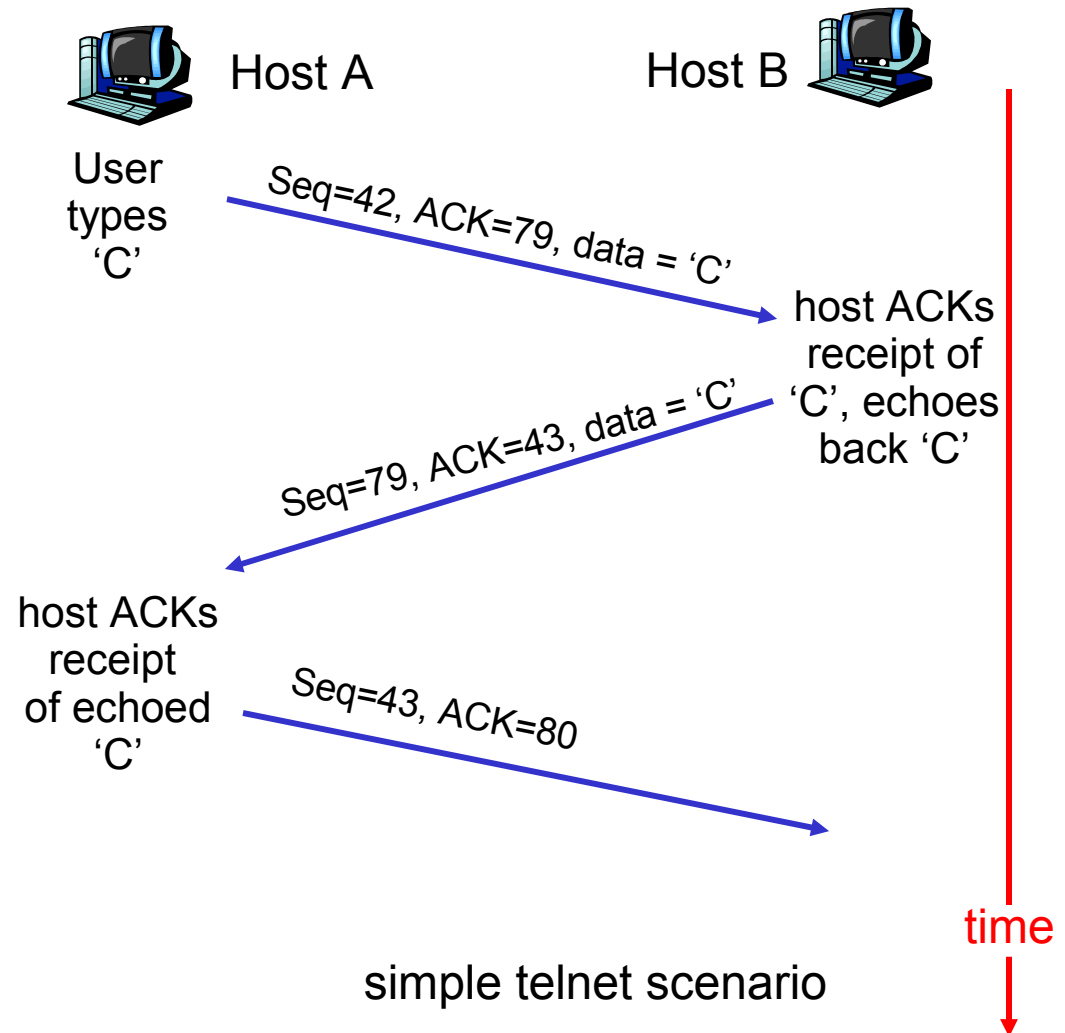
- ◆ byte stream
“number” of first byte in segment’s data

ACKs:

- ◆ seq # of next byte expected from other side
- ◆ cumulative ACK

Q: how receiver handles out-of-order segments

- ◆ A: TCP spec doesn’t say, - up to implementor



TCP Round Trip Time and Timeout

Q: how to set TCP timeout value?

- longer than RTT
 - ◆ but RTT varies
- too short: premature timeout
 - ◆ unnecessary retransmissions
- too long: slow reaction to segment loss

Q: how to estimate RTT?

- **SampleRTT**: measured time from segment transmission until ACK receipt
 - ◆ ignore retransmissions
- **SampleRTT** will vary, want estimated RTT “smoother”
 - ◆ average several recent measurements, not just current **SampleRTT**

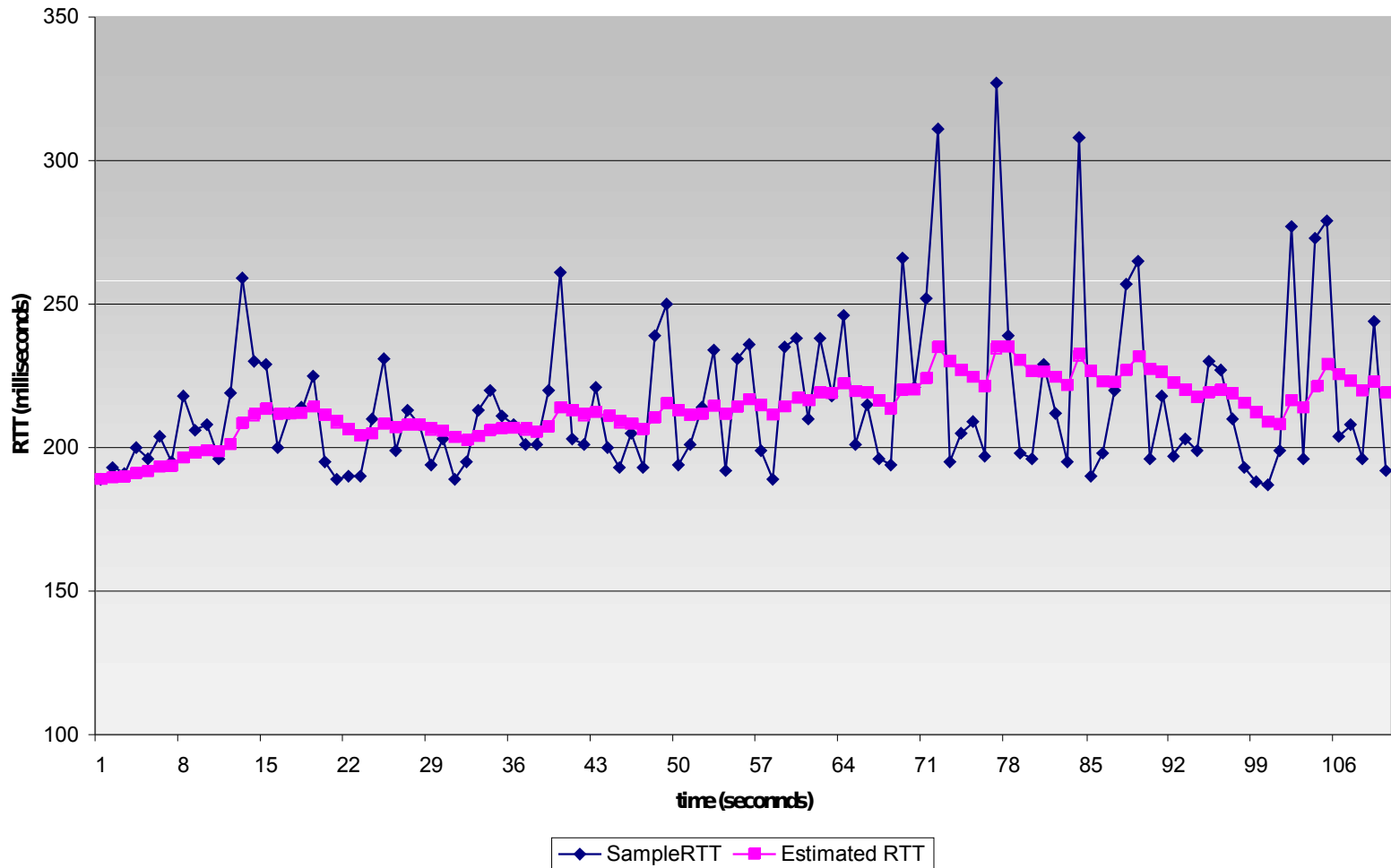
TCP Round Trip Time and Timeout

$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$$

- Exponentially weighted moving average
- influence of past sample decreases exponentially fast
- typical value: $\alpha = 0.125$

Example RTT estimation:

RTT: `gaiacs.umass.edu` to `fantasia.eurecom.fr`



TCP Round Trip Time and Timeout

Setting the timeout

- **EstimatedRTT** plus “safety margin”
 - ◆ large variation in **EstimatedRTT** → larger safety margin
- first estimate of how much **SampleRTT** deviates from **EstimatedRTT**:

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

(typically, $\beta = 0.25$)

Then set timeout interval:

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

Chapter 3 outline

- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer
- 3.5 Connection-oriented transport: TCP
 - ◆ segment structure
 - ◆ **reliable data transfer**
 - ◆ flow control
 - ◆ connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

TCP reliable data transfer

- TCP creates rdt service on top of IP's unreliable service
- Pipelined segments
- Cumulative acks
- TCP uses single retransmission timer
- Retransmissions are triggered by:
 - ◆ timeout events
 - ◆ duplicate acks
- Initially consider simplified TCP sender:
 - ◆ ignore duplicate acks
 - ◆ ignore flow control, congestion control

TCP sender events:

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 unacked segment)
- expiration interval:
TimeOutInterval

timeout:

- retransmit segment that caused timeout
- restart timer

Ack rcvd:

- If acknowledges previously unacked segments
 - ◆ update what is known to be acked
 - ◆ start timer if there are outstanding segments

NextSeqNum = InitialSeqNum

SendBase = InitialSeqNum

```
loop (forever) {  
    switch(event)
```

```
    event: data received from application above  
        create TCP segment with sequence number NextSeqNum  
        if (timer currently not running)  
            start timer  
        pass segment to IP  
        NextSeqNum = NextSeqNum + length(data)
```

```
    event: timer timeout  
        retransmit not-yet-acknowledged segment with  
            smallest sequence number  
        start timer
```

```
    event: ACK received, with ACK field value of y  
        if (y > SendBase) {  
            SendBase = y  
            if (there are currently not-yet-acknowledged segments)  
                start timer  
        }  
}
```

```
} /* end of loop forever */
```

TCP sender (simplified)

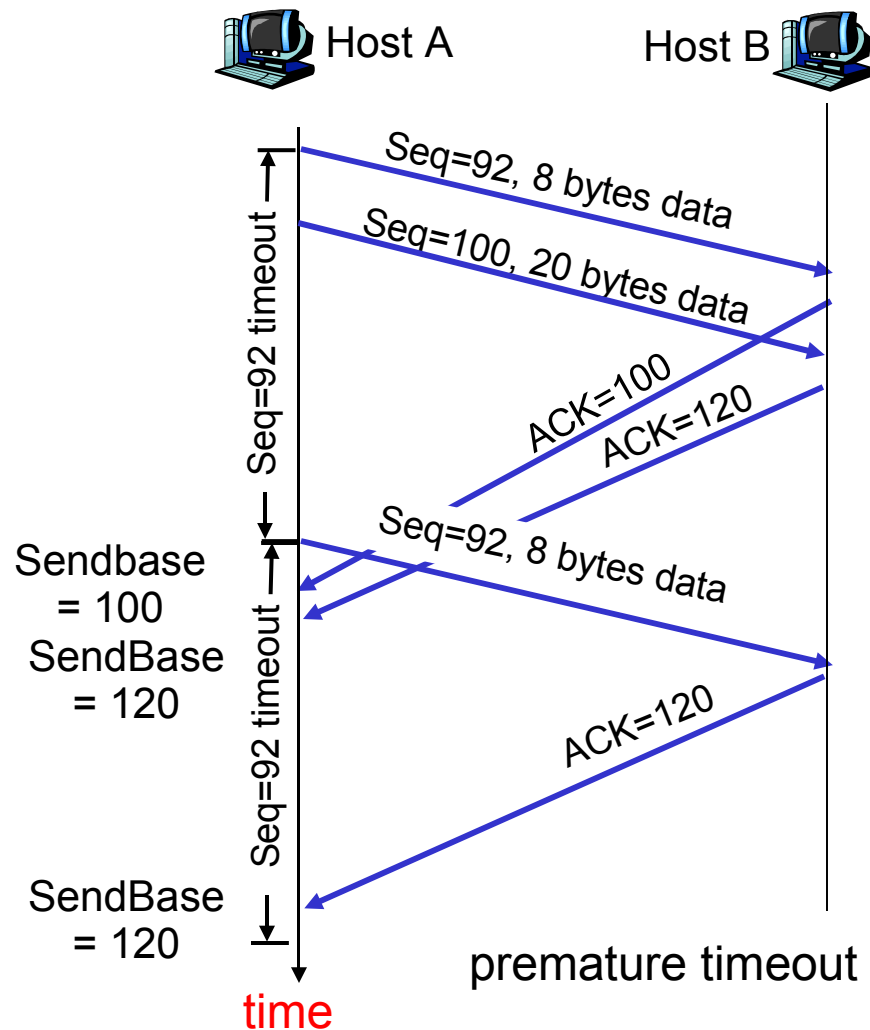
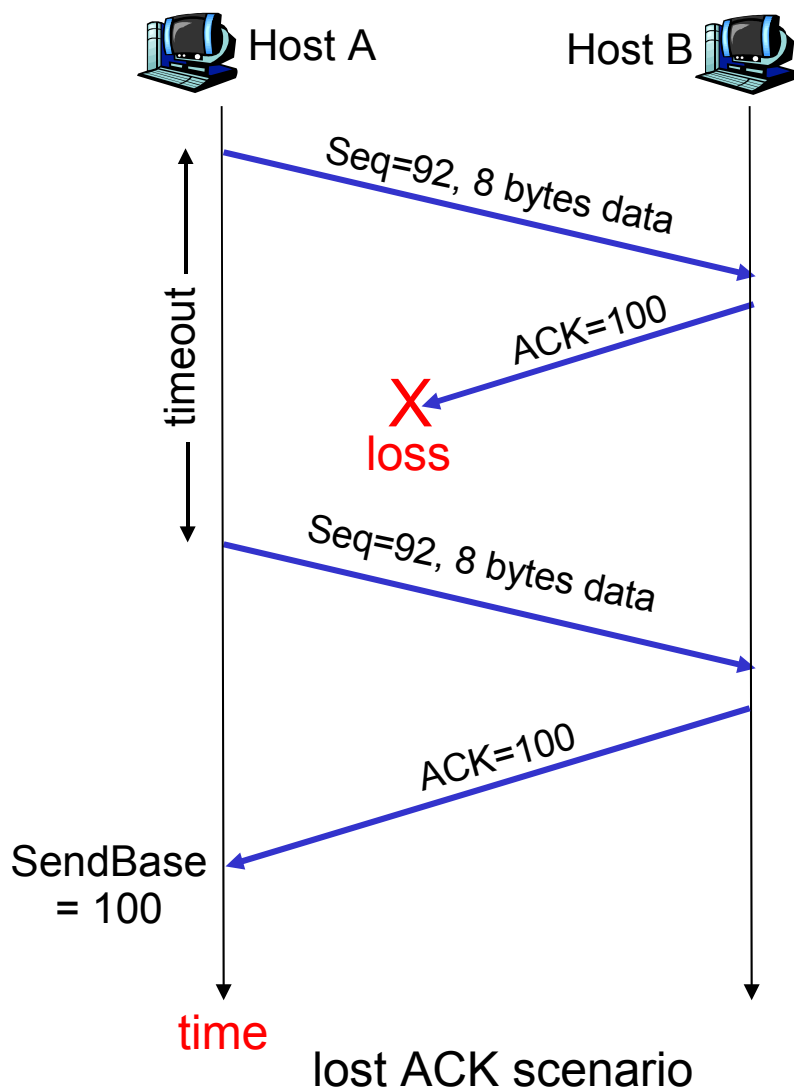
Comment:

- SendBase-1: last cumulatively ack'ed byte

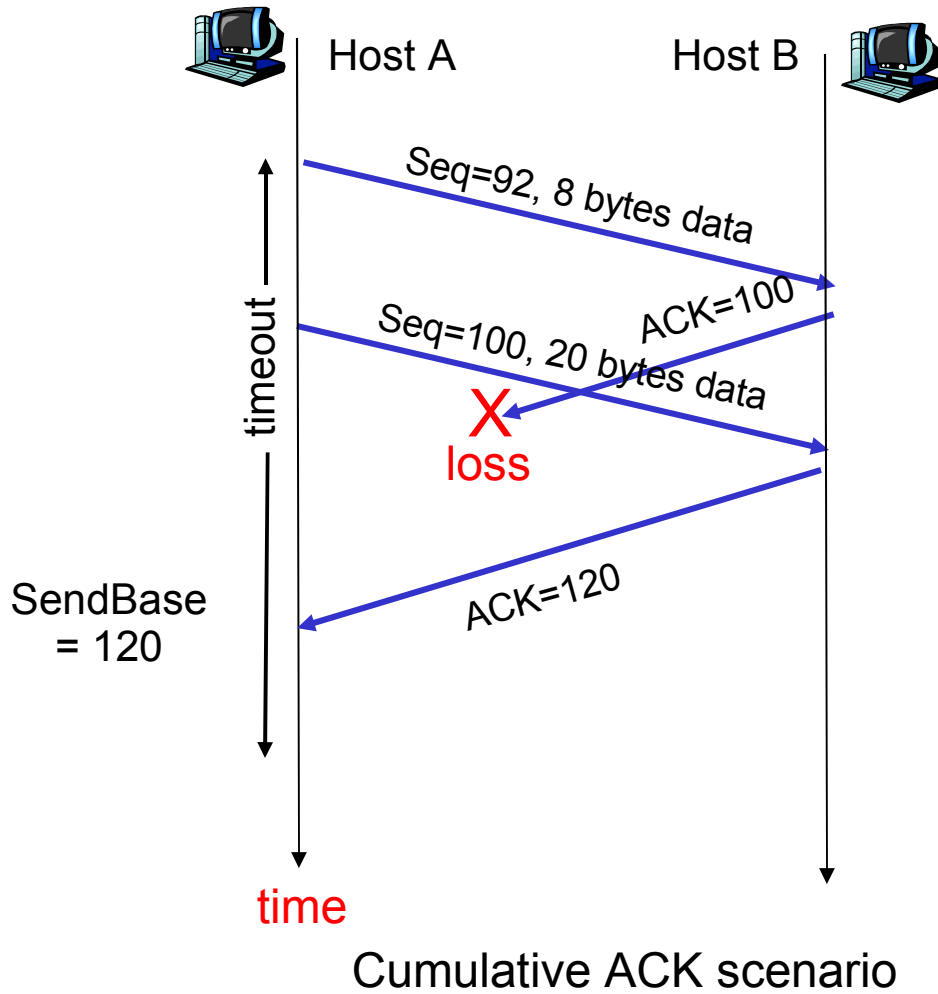
Example:

- SendBase-1 = 71;
y = 73, so the rcvr wants 73+ ;
y > SendBase, so that new data is acked

TCP: retransmission scenarios



TCP retransmission scenarios (more)



TCP ACK generation [RFC 1122, RFC 2581]

Event at Receiver

TCP Receiver action

Arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed

Delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK

Arrival of in-order segment with expected seq #. One other segment has ACK pending

Immediately send single cumulative ACK, ACKing both in-order segments

Arrival of out-of-order segment higher-than-expected seq. # . Gap detected

Immediately send *duplicate ACK*, indicating seq. # of next expected byte

Arrival of segment that partially or completely fills gap

Immediate send ACK, provided that segment starts at lower end of gap

Recap

- UDP socket programming
 - ◆ DatagramSocket, DatagramPacket

- TCP
 - ◆ Sequence numbers, ACKs
 - ◆ RTT, DevRTT, timeout calculations
 - ◆ Reliable data transfer algorithm

Next time

□ TCP

- ◆ Fast retransmit
- ◆ Flow control
- ◆ Connection management
- ◆ Congestion control