# Transport Layer

Our goals:
- Understand principles behind transport layer services:
  - Multiplexing, demultiplexing
  - Reliable data transfer
  - Flow control
  - Congestion control
- Learn about Internet transport layer protocols:
  - UDP: connectionless transport
  - TCP: connection-oriented reliable transport
  - TCP congestion control

Roadmap
1. Transport-layer services
2. Multiplexing and demultiplexing
3. Connectionless transport: UDP
4. Principles of reliable data transfer
5. Connection-oriented transport: TCP
   a. Segment structure
   b. Reliable data transfer
   c. Flow control
   d. Connection management
6. Principles of congestion control
7. TCP congestion control
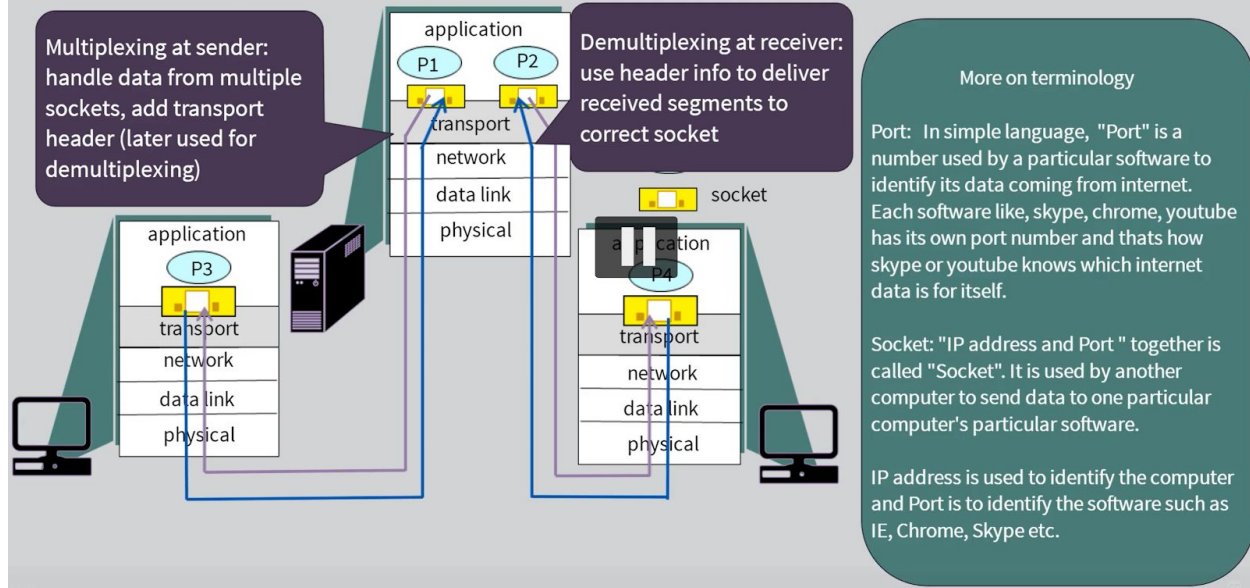
## Transport Services and Protocols
- Provide *logical communication* between app processes running on different hosts
- Transport protocols run in end systems
  - Send side: breaks app messages into segments, passes to network layer
  - Receiving side: reassembles segments into messages, passes to app layer
- More than one transport protocol available to apps
  - Internet: TCP and UDP

## Internet Transport-layer Protocols
- Reliable, in-order delivery (TCP)
  - Congestion control
  - Flow control
  - Connection setup
- Unreliable, unordered delivery: UDP
  - No-frills extension of "best effort" IP
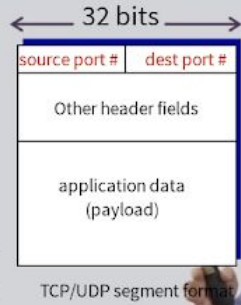- Services not available
  - Delay guarantees

○ Bandwidth guarantees

# Multiplexing and Demultiplexing

**Multiplexing at sender:** handle data from multiple sockets, add transport header (later used for demultiplexing)

**Demultiplexing at receiver:** use header info to deliver received segments to correct socket

application
P1    P2

transport
network
data link
physical

socket

application
P3

transport
network
data link
physical

application
P4

transport
network
data link
physical

**More on terminology**

Port:  In simple language,  "Port" is a number used by a particular software to identify its data coming from internet. Each software like, skype, chrome, youtube has its own port number and thats how skype or youtube knows which internet data is for itself.

Socket: "IP address and Port " together is called "Socket". It is used by another computer to send data to one particular computer's particular software.

IP address is used to identify the computer and Port is to identify the software such as IE, Chrome, Skype etc.
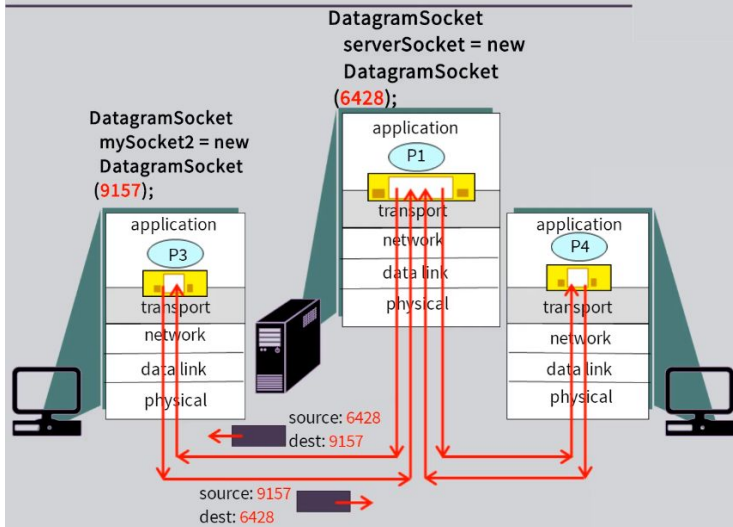
# How Demultiplexing Works

- Host receives IP datagrams
  - Each datagram has source IP address, destination IP address
  - Each datagram carries one transport-layer segment
- Host uses *IP addresses & port numbers* to direct segment to appropriate socket



← 32 bits →

| source port # | dest port # |

Other header fields

application data
(payload)

TCP/UDP segment format

# Connectionless Demultiplexing

- When host receives UDP segment:
  - Checks destination port # in segment
  - Directs UDP segment to socket with that port #

- IP datagrams with *same dest. port #,* but different source IP addresses and/or source port numbers will be directed to *same socket* at dest
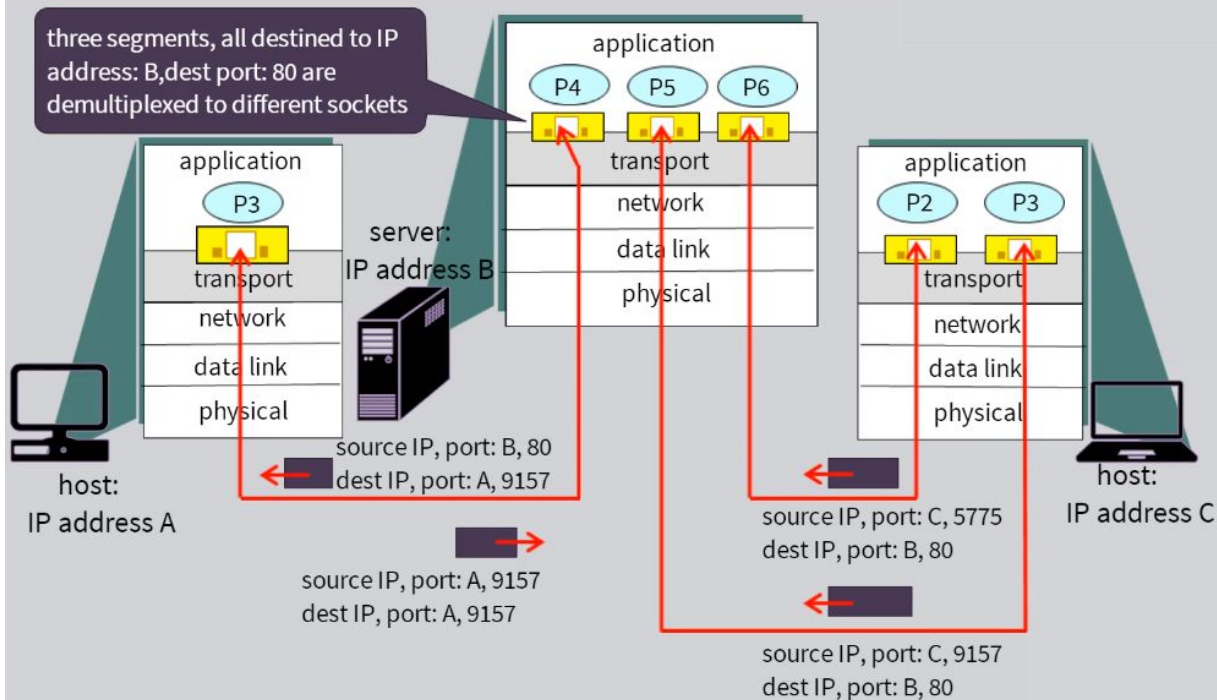
# Connectionless Demux: Example



**DatagramSocket serverSocket = new DatagramSocket (6428);**

**DatagramSocket mySocket2 = new DatagramSocket (9157);**

application
P1

transport
network
data link
physical

application
P3

transport
network
data link
physical

application
P4

transport
network
data link
physical

source: 6428
dest: 9157

source: 9157
dest: 6428

# Connection-Oriented Demux

- TCP socket identified by 4-tuple:
  - Source IP address
  - Source port number
  - Dest IP address
  - Dest port number
- Server host may support many simultaneous TCP sockets: each socket identified by its own 4-tuple
- Web servers have different sockets for each connecting client
  - Non-persistent HTTP will have different socket for each request
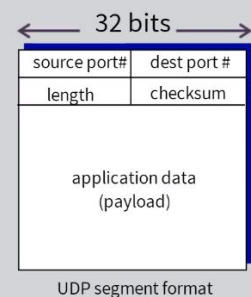
# Example

three segments, all destined to IP address: B, dest port: 80 are demultiplexed to different sockets

application
P4   P5   P6
transport
network
data link
physical

application
P3
transport
network
data link
physical

server:
IP address B

host:
IP address A

source IP, port: B, 80
dest IP, port: A, 9157

source IP, port: A, 9157
dest IP, port: A, 9157

application
P2   P3
transport
network
data link
physical

host:
IP address C

source IP, port: C, 5775
dest IP, port: B, 80

source IP, port: C, 9157
dest IP, port: B, 80

# Connectionless Transport: UDP

- *Connectionless:*
  - No handshaking between UDP sender, receiver
  - Each UDP segment handled independently of others

- UDP use:
  - Streaming multimedia apps (loss tolerant, rate sensitive)
  - DNS

- "Best effort" service
  UDP segments may be:
  - Lost
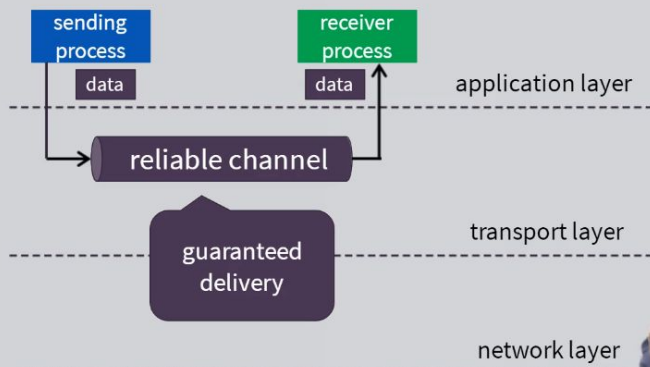  - Delivered out-of-order to app
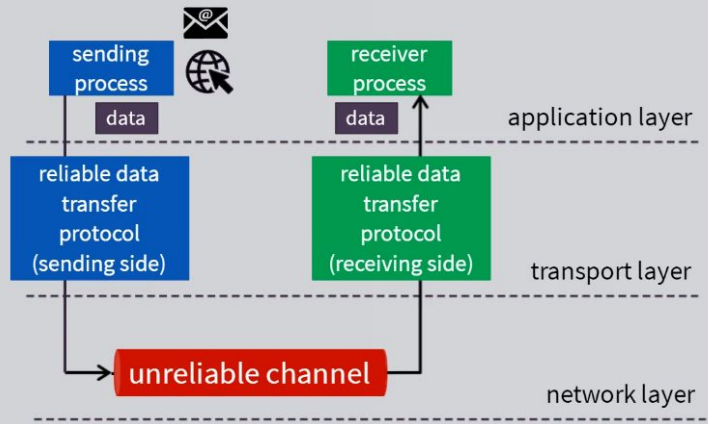
# UDP: Segment Header

Why is there a UDP?
- No connection establishment (which can add delay)
- Simple: no connection state at sender, receiver
- Small header size
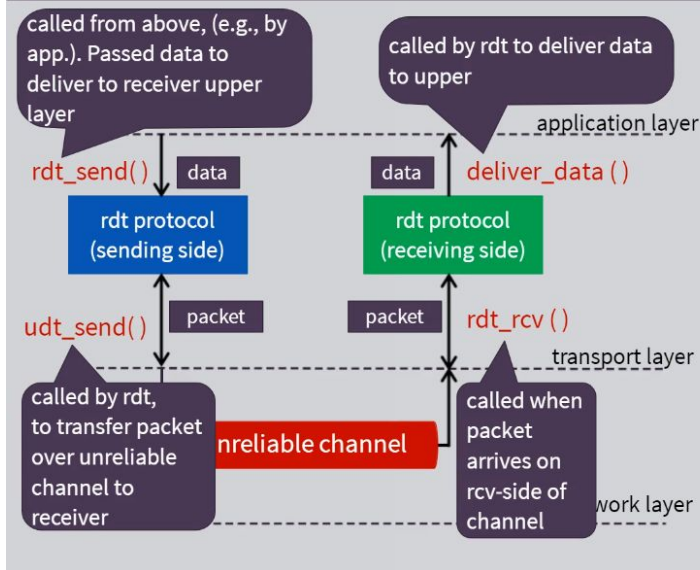- No congestion control: UDP can blast away as fast as desired

32 bits

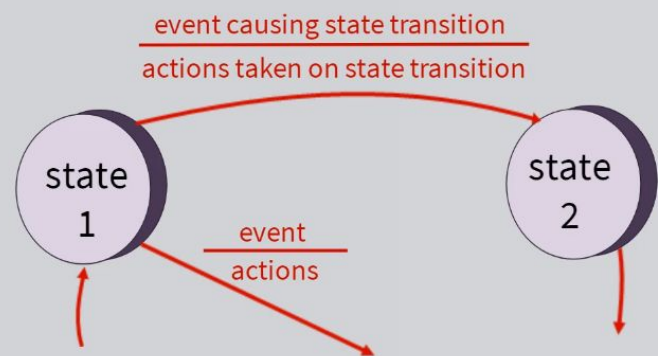| source port# | dest port # |
|---|---|
| length | checksum |

application data
(payload)

UDP segment format

## Principles of Reliable Data Transfer



sending process — data — application layer
receiver process — data

reliable channel — transport layer

guaranteed delivery

network layer

## Principles of Reliable Data Transfer



sending process — data — application layer
receiver process — data

reliable data transfer protocol (sending side)
reliable data transfer protocol (receiving side)
transport layer

unreliable channel
network layer

## RDT: Getting Started



called from above, (e.g., by app.). Passed data to deliver to receiver upper layer

called by rdt to deliver data to upper

application layer

rdt_send( ) — data
data — deliver_data ( )

rdt protocol (sending side)
rdt protocol (receiving side)

udt_send( ) — packet
packet — rdt_rcv ( )

transport layer

called by rdt, to transfer packet over unreliable channel to receiver

unreliable channel

called when packet arrives on rcv-side of channel

network layer

## RDT: Getting Started



event causing state transition
actions taken on state transition

state 1          state 2

event
actions

## RDT1.0

Reliable transfer over a reliable channel

- Underlying channel perfectly reliable
  - No bit errors
  - No loss of packets
- Separate FSMs for sender, receiver:
  - Sender sends data into underlying channel
  - Receiver reads data from underlying channel

## RDT2.0: Channel with Bit Errors

- Underlying channel may flip bits in packet
  - Checksum to detect bit
- *The* question: how to recover from errors?
  - *Acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK
  - *Negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors
  - Sender retransmits pkt on receipt of NAK

## RDT2.0 has a Fatal Flaw!

**What happens if ACK/NAK corrupted?**
- Sender doesn't know what happened at receiver!
- Can't just retransmit: possible duplicate

**Handling duplicates:**
- Sender retransmits current pkt if ACK/NAK corrupted
- Sender adds *sequence number* to each pkt
- Receiver discards (doesn't deliver up) duplicate pkt

**Stop and wait**
- Sender sends one packet, then waits for receiver respond

## RDT2.1: Discussion

**Sender:**
- Seq # added to pkt
- Two seq. #'s (0,1) will suffice.
- Must check if received ACK/NAK corrupted
- Twice as many states
  - State must "remember" whether "expected" pkt should have seq # of 0 or 1

**Receiver:**
- Must check if received packet is duplicate
  - State indicates whether 0 or 1 is expected pkt seq #
- Note: receiver can *not* know if its last ACK/NAK received OK at sender

## RDT2.2: a NAK-free Protocol

- Same functionality as rdt2.1, using ACKs only
- Instead of NAK, receiver sends ACK for last pkt received OK
  - Receiver must *explicitly* include seq # of pkt being ACKed
- Duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

## RDT3.0: Channels with Errors and Loss

**New assumption:**
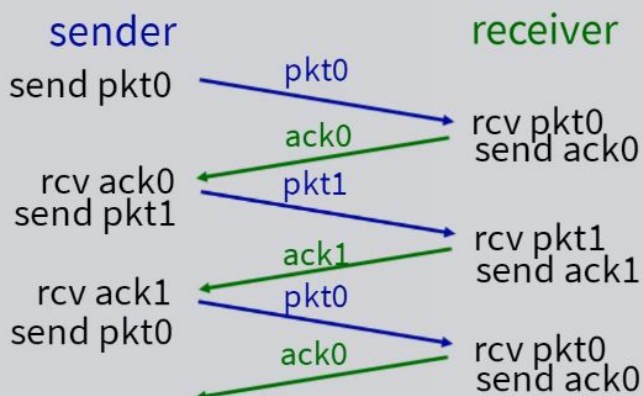Underlying channel can also lose packets (data, ACKs)
  - Checksum, seq. #, ACKs, retransmissions will be of help … but not enough

**Approach:** Sender waits "reasonable" amount of time for ACK
- Retransmits if no ACK received in this time
- If pkt (or ACK) just delayed (not lost):
  - Retransmission will be duplicate, but seq. #'s already handles this
  - Receiver must specify seq # of pkt being ACKed
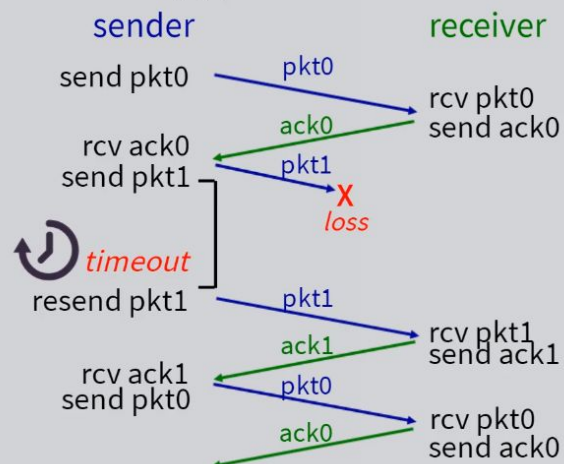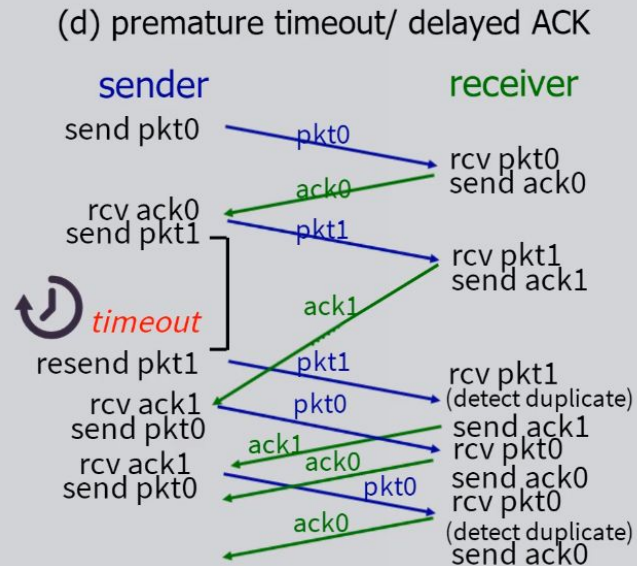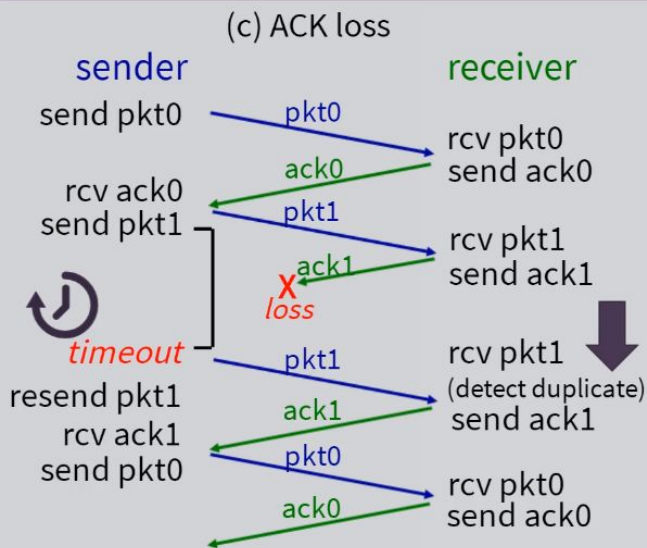- Requires countdown timer

## RDT3.0 in Action



(a) no loss

## RDT3.0 in Action



(b) packet loss

## RDT3.0 in Action

### (c) ACK loss

sender · receiver

send pkt0 — pkt0 → rcv pkt0
send ack0
rcv ack0 ← ack0
send pkt1 — pkt1 → rcv pkt1
send ack1
X ack1 loss
timeout
resend pkt1 — pkt1 → rcv pkt1
(detect duplicate)
send ack1
rcv ack1 ← ack1
send pkt0 — pkt0 → rcv pkt0
send ack0
← ack0

## RDT3.0 in Action

### (d) premature timeout/ delayed ACK

sender · receiver

send pkt0 — pkt0 → rcv pkt0
send ack0
rcv ack0 ← ack0
send pkt1 — pkt1 → rcv pkt1
send ack1
timeout ack1
resend pkt1 — pkt1 → rcv pkt1
(detect duplicate)
rcv ack1 ← pkt0 send ack1
send pkt0 rcv pkt0
rcv ack1 ← ack1 send ack0
send pkt0 ← ack0 rcv pkt0
pkt0 → (detect duplicate)
← ack0 send ack0

## Performance of RDT3.0

- Rdt3.0 is correct, but performance stinks
- e.g.: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$D_{trans} = \frac{L}{R} = \frac{8000 \ bits}{10^9 \ bits/sec} = 8 \ \mu s$$

$U_{sender}$: *utilization* – fraction of time sender busy sending

$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$
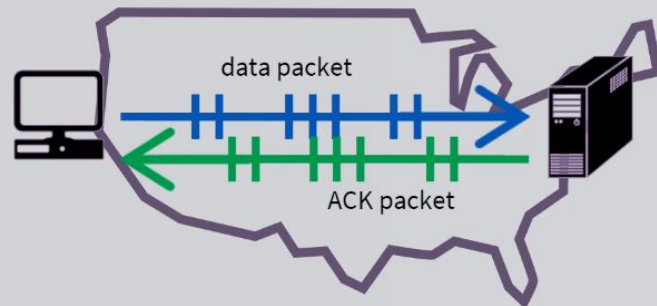
- If RTT=30 millisecond,
  1KB pkt every 30 millisecond:
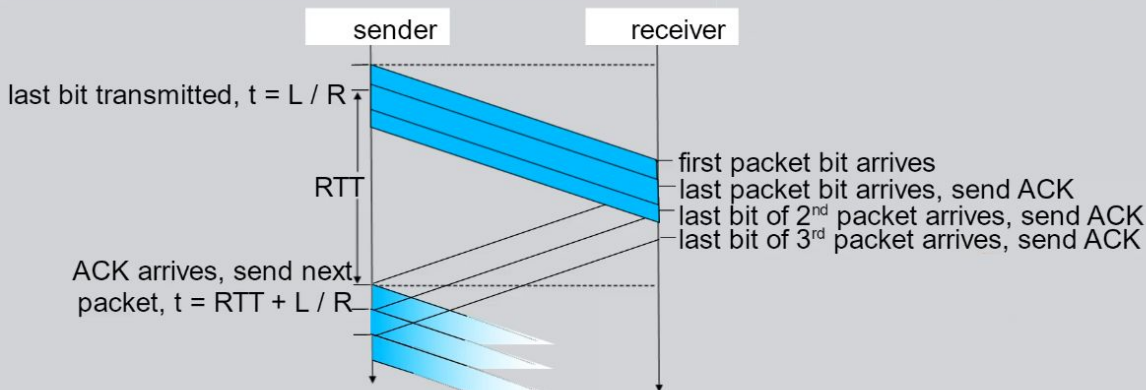  33kB/sec throughput over 1 Gbps link

## Pipelined Protocols

Two generic form:
- Go-Back-N
- Selective repeat

data packet

ACK packet

## Pipelining: Increased Utilization

sender · receiver

last bit transmitted, t = L / R

RTT

first packet bit arrives
last packet bit arrives, send ACK
last bit of 2nd packet arrives, send ACK
last bit of 3rd packet arrives, send ACK

ACK arrives, send next packet, t = RTT + L / R

$$U_{sender} = \frac{3L/R}{RTT + L/R} = \frac{.0024}{30.008} = 0.00081$$
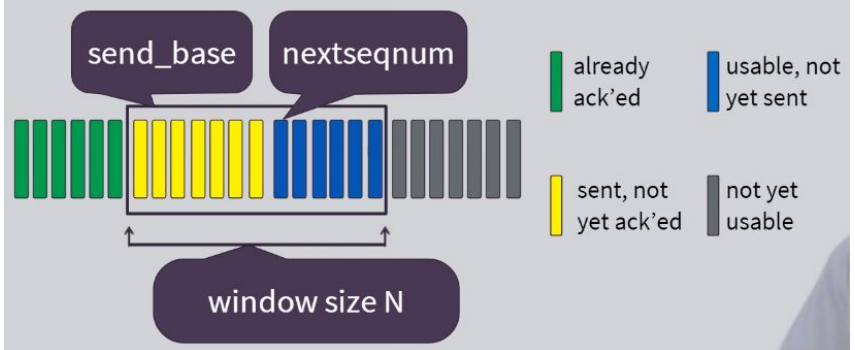
## Pipelined Protocols: Overview

**Go-back-N:**
- Sender can have up to N unacked packets in pipeline
- Receiver only sends *cumulative ack*
  - doesn't ack packet if there's a gap
- Sender has timer for oldest unacked packet
  - When timer expires, retransmit *all* unacked packets

**Selective Repeat:**
- Sender can have up to N unack'ed packets in pipeline
- rcvr sends *individual ack* for each packet
- Sender maintains timer for each unacked packet
  - When timer expires, retransmit only that unacked packet

## Go-Back-N: Sender



send_base   nextseqnum

window size N

- already ack'ed (green)
- usable, not yet sent (blue)
- sent, not yet ack'ed (yellow)
- not yet usable (gray)

## GBN in Action

NYU TANDON SCHOOL OF ENGINEERING



sender window (N=4)

| sender | receiver |
|--------|----------|
| send pkt0 | receive pkt0, send ack0 |
| send pkt1 | receive pkt1, send ack1 |
| send pkt2 ✗ loss | |
| send pkt3 | receive pkt3, discard, (re)send ack1 |
| (wait) | |
| rcv ack0, send pkt4 | receive pkt4, discard, (re)send ack1 |
| rcv ack1, send pkt5 | receive pkt5, discard, (re)send ack1 |
| ignore duplicate ACK | |
| pkt 2 timeout | |
| send pkt2 | rcv pkt2, deliver, send ack2 |
| send pkt3 | rcv pkt3, deliver, send ack3 |
| send pkt4 | rcv pkt4, deliver, send ack4 |
| send pkt5 | rcv pkt5, deliver, send ack5 |

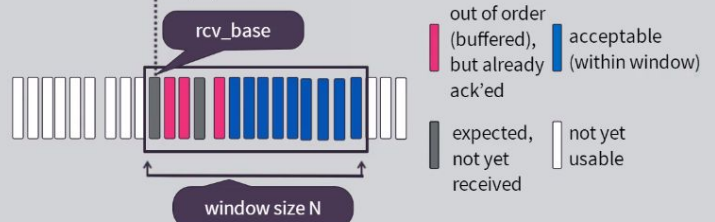0 1 2 3 4 5 6 7 8

## Selective Repeat

- receiver *individually* acknowledges all correctly received pkts
  - buffers pkts, as needed, for eventual in-order delivery to upper layer
- sender only resends pkts for which ACK not received
  - sender timer for each unACKed pkt
- sender window
  - *N* consecutive seq #'s
  - limits seq #s of sent, unACKed pkts

## Selective Repeat



(a) sender view

send_base   nextseqnum

window size N

- already ack'ed (green)
- usable, not yet sent (blue)
- sent, not yet ack'ed (yellow)
- not yet usable (gray)

(b) receiver view

rcv_base

window size N

- out of order (buffered), but already ack'ed (pink)
- acceptable (within window) (blue)
- expected, not yet received (gray)
- not yet usable (white)

# Selective Repeat in Action

**sender window (N=4)**    **sender**          **receiver**

```
0 1 2 3 4 5 6 7 8      send  pkt0
0 1 2 3 4 5 6 7 8      send  pkt1
0 1 2 3 4 5 6 7 8      send  pkt2                           receive pkt0, send ack0
0 1 2 3 4 5 6 7 8      send  pkt3        X loss             receive pkt1, send ack1
                       (wait)

                                                            receive pkt3, buffer,
  0 1 2 3 4 5 6 7 8    rcv ack0, send pkt4                      send ack3
  0 1 2 3 4 5 6 7 8    rcv ack1, send pkt5
                                                            receive pkt4, buffer,
                       record ack3 arrived                     send ack4
                       pkt 2 timeout                        receive pkt5, buffer,
                                                               send ack5
  0 1 2 3 4 5 6 7 8    send  pkt2
  0 1 2 3 4 5 6 7 8    record ack4 arrived                 rcv pkt2; deliver pkt2,
  0 1 2 3 4 5 6 7 8    record ack4 arrived
  0 1 2 3 4 5 6 7 8                                         pkt3, pkt4, pkt5; send ack2
```

*Q: what happens when ack2 arrives?*
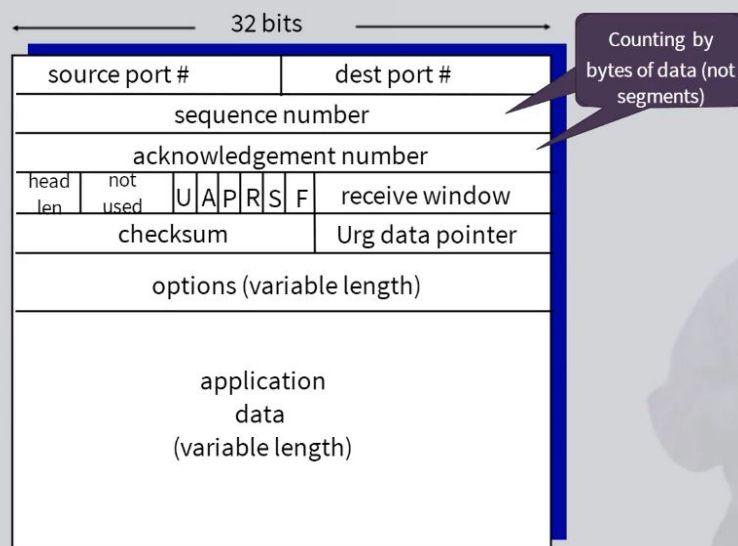
---

# TCP: Overview

## RFCs: 793,1122,1323, 2018, 2581

- Point-to-point:
  - One sender, one receiver
- Reliable, in-order *byte steam:*
  - No "message boundaries"
- Pipelined:
  - TCP congestion and flow control set window size

- 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

| 32 bits | |
|---|---|
| source port # | dest port # |
| sequence number | |
| acknowledgement number | |

Counting by bytes of data (not segments)

| head len | not used | U A P R S F | receive window |
|---|---|---|---|
| checksum | | | Urg data pointer |

options (variable length)

application
data
(variable length)

## TCP Seq. Numbers, ACKs
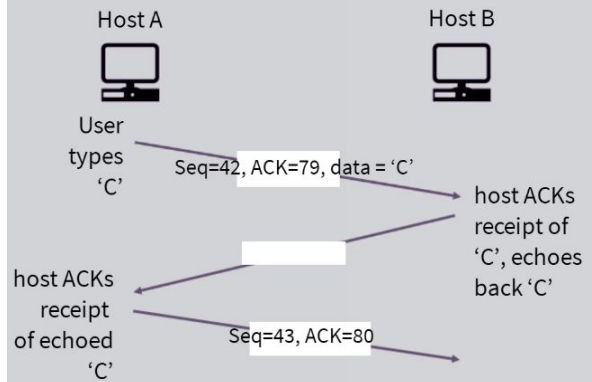
Sequence numbers:
- Byte stream "number" of first byte in segment's data

acknowledgements:
- Seq # of next byte expected from other side
- Cumulative ACK

## TCP Seq. Numbers, ACKs

simple telnet scenario

Host A                                    Host B

User types 'C'
Seq=42, ACK=79, data = 'C'

host ACKs receipt of 'C', echoes back 'C'

host ACKs receipt of echoed 'C'
Seq=43, ACK=80

## TCP Round Trip Time, 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"
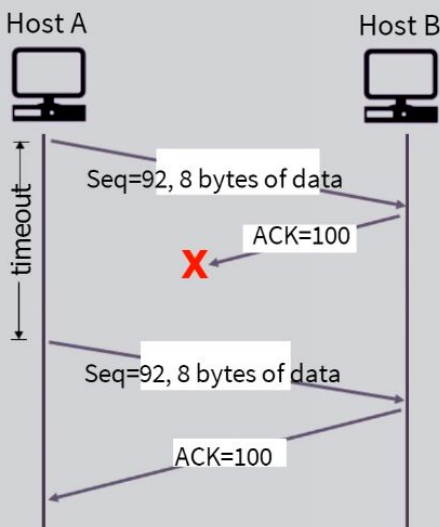  - average several *recent* measurements, not just current **SampleRTT**

## TCP Round Trip Time, Timeout

```
EstimatedRTT = (1- α)*EstimatedRTT + α*SampleRTT
typical value: a = 0.125
```
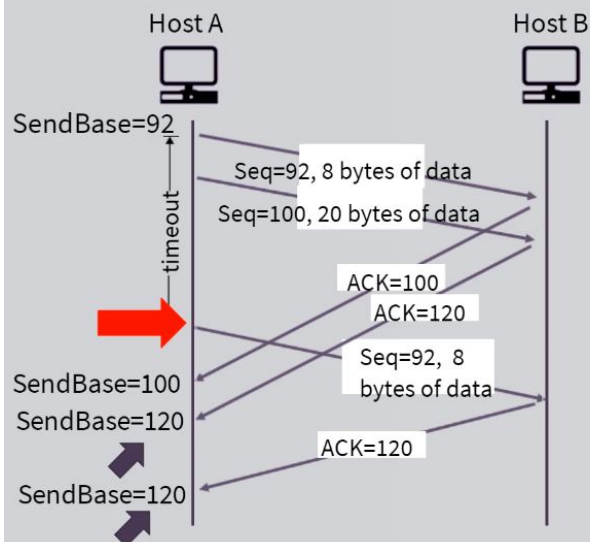
RTT (milliseconds)

time (seconds)

## TCP: Retransmission Scenarios

lost ACK scenario

Host A                          Host B

Seq=92, 8 bytes of data

ACK=100

X

timeout

Seq=92, 8 bytes of data

ACK=100

## TCP: Retransmission Scenarios

premature timeout

Host A                          Host B

SendBase=92

Seq=92, 8 bytes of data
Seq=100, 20 bytes of data

timeout

ACK=100
ACK=120

Seq=92, 8 bytes of data

SendBase=100
SendBase=120

ACK=120

SendBase=120

# TCP: Retransmission Scenarios

cumulative ACK



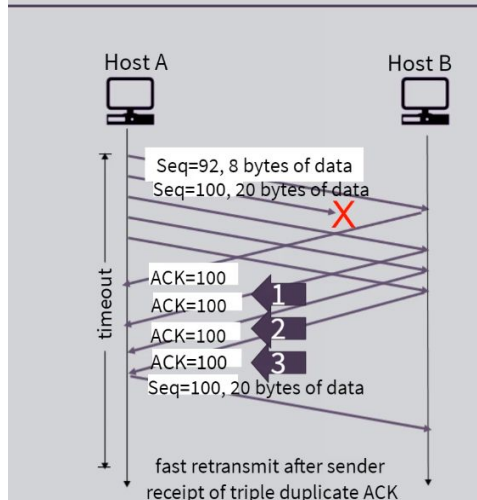# TCP ACK Generation [RFC 1122, RFC 2581]

| 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-expect 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 |

# TCP Fast Retransmit



fast retransmit after sender receipt of triple duplicate ACK

# TCP Fast Retransmit

If sender receives 3 ACKs for same data ("triple duplicate ACKs"), resend unacked segment with smallest seq #
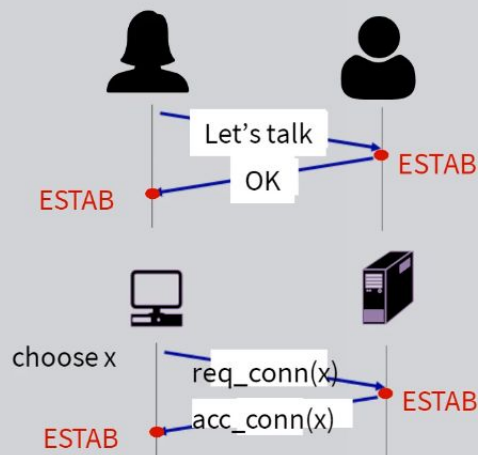- Likely that unacked segment lost, so don't wait for timeout

# TCP Flow Control

- Receiver "advertises" free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments
  - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
  - Many operating systems autoadjust **RcvBuffer**
- Sender limits amount of unacked ("in-flight") data to receiver's **rwnd** value
- Guarantees receive buffer will not overflow



to application process

buffered data

free buffer space

TCP segment payloads

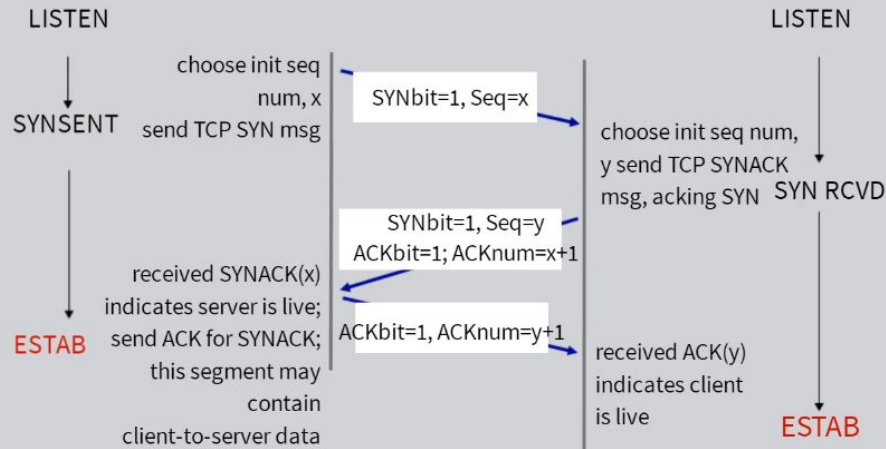receiver-side buffering

# Agreeing to Establish a Connection

2-way handshake:



Let's talk

OK

ESTAB

ESTAB

choose x

req_conn(x)

acc_conn(x)

ESTAB

ESTAB

# TCP 3-Way Handshake

client state

sever state

LISTEN

LISTEN

SYNSENT — choose init seq num, x send TCP SYN msg

**SYNbit=1, Seq=x**

choose init seq num, y send TCP SYNACK msg, acking SYN  SYN RCVD

**SYNbit=1, Seq=y ACKbit=1; ACKnum=x+1**

received SYNACK(x) indicates server is live; send ACK for SYNACK; this segment may contain client-to-server data

ESTAB

**ACKbit=1, ACKnum=y+1**

received ACK(y) indicates client is live

ESTAB

# TCP: Closing a Connection

client state

ESTAB

FIN_WAIT_1 — clientSocket.close()

**FINbit=1, seq=x**

can no longer send but can receive data wait for server close

FIN_WAIT_2

**ACKbit=1; ACKnum=x+1**

can still send data

server state

ESTAB

CLOSE_WAIT

TIMED_WAIT

**FINbit=1, seq=y**

can no longer send data

LAST_ACK

**ACKbit=1; ACKnum=y+1**

timed wait for 2*max segment lifetime

CLOSED

CLOSED

## Congestion: Scenario



$\lambda_{in}$ : original data
$\lambda'_{in}$ : original data, *plus* retransmitted data
$\lambda_{out}$

Host A
Host B
Host D
Host C

finite shared output link buffers
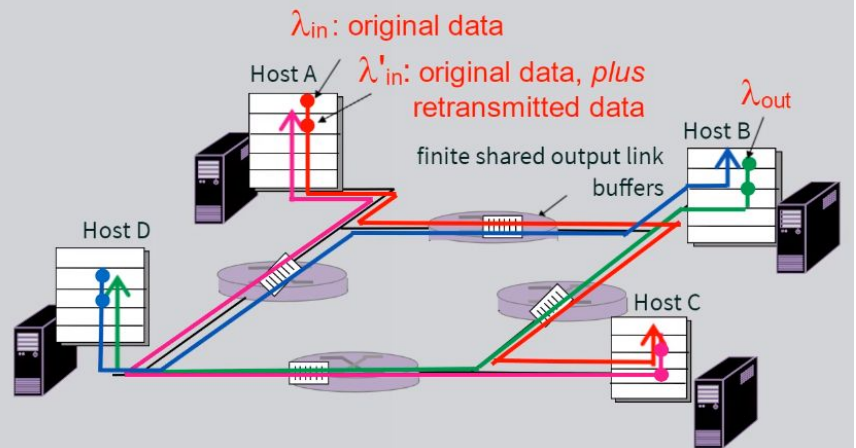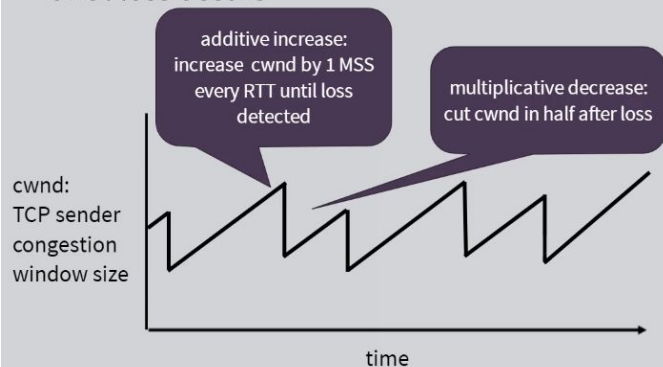
## Principles of Congestion Control

Congestion:
- Informally: "too many sources sending too much data too fast for network to handle"
- Different from flow control!
- Manifestations:
- Lost packets (buffer overflow at routers)
- Long delays (queueing in router buffers)

## Congestion: Scenario

Another "cost"of congestion: when packet dropped, any "upstream transmission capacity used for that packet was wasted!

## TCP Congestion Control

Additive Increase Multiplicative Decrease (AIMD)

- *Approach: S*ender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
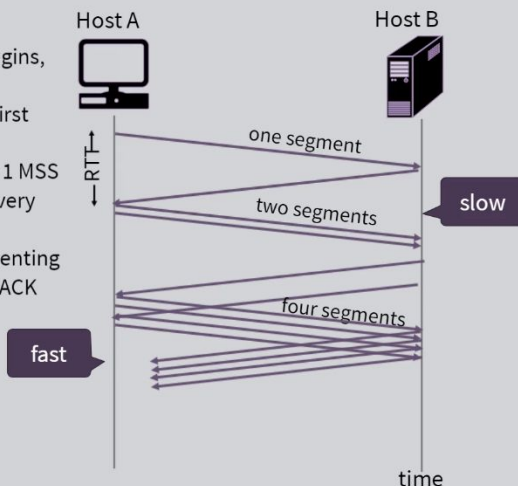


additive increase: increase cwnd by 1 MSS every RTT until loss detected

multiplicative decrease: cut cwnd in half after loss

cwnd: TCP sender congestion window size

time

## TCP Congestion Control

### TCP Slow Start

- When connection begins, increase rate exponentially until first loss event:
  - Initially **cwnd** = 1 MSS
  - Double **cwnd** every RTT
  - Done by incrementing **cwnd** for every ACK received

Host A
Host B
RTT
one segment
two segments
four segments
slow
fast
time

## Summary

- Principles behind transport layer services:
  - Multiplexing, demultiplexing
  - Reliable data transfer
  - Flow control
  - Congestion control
- Instantiation, implementation in the Internet
  - UDP
  - TCP