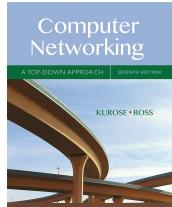


Chapter 3: Transport Layer

Our goals:

- ❑ understand **principles** behind transport layer services:
 - Multiplexing / demultiplexing
 - reliable data transfer
 - flow control
 - congestion control



- ❑ learn about transport layer protocols **implemented** in the Internet:
 - UDP: connectionless transport
 - TCP: connection-oriented transport
 - TCP congestion control

Computer Networking: A Top Down Approach
7th edition
Jim Kurose, Keith Ross
Pearson/Addison Wesley
April 2016

ELEC 331 1

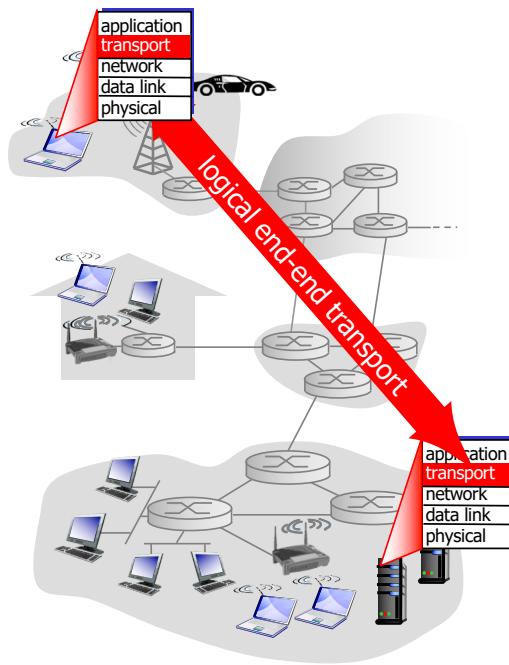
Chapter 3 outline

- ❑ 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

ELEC 331 2

Transport services and protocols

- ❑ provide *logical communication* between application processes running on different hosts
- ❑ transport protocols run in end systems
 - send side: breaks application messages into **segments**, passes to network layer
 - receive side: reassembles segments into messages, passes to application layer
- ❑ more than one transport protocol available to applications
 - Internet: TCP and UDP



ELEC 331 3

Transport vs. network layer

- ❑ *transport layer*: logical communication between **processes** running on different hosts
 - relies on, enhances, network layer services
- ❑ *network layer*: logical communication between **hosts**

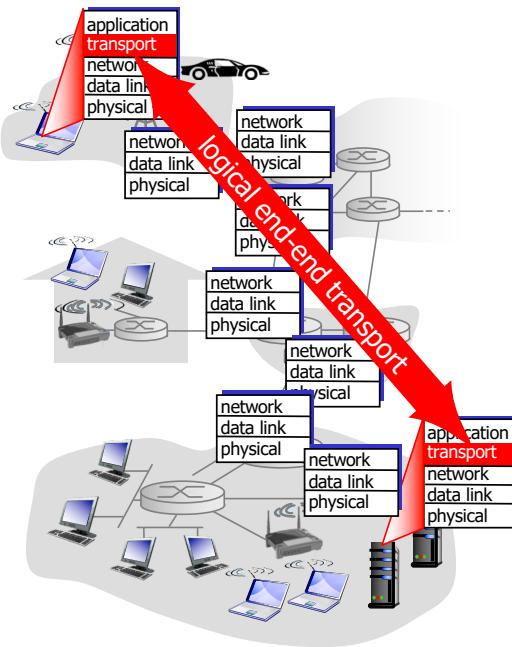
Household analogy:

- 12 kids sending letters to 12 kids
- ❑ processes = kids
- ❑ app messages = letters in envelopes
- ❑ hosts = houses
- ❑ transport layer protocol = Ann and Bill
- ❑ network layer protocol = postal service (including mail carriers)

ELEC 331 4

Internet transport layer protocols

- ❑ UDP: unreliable, unordered delivery
 - no-frills extension of “best-effort” IP
- ❑ TCP: reliable, in-order delivery
 - congestion control
 - flow control
 - connection setup
- ❑ services not available:
 - delay and bandwidth guarantee
- ❑ Other service: Encryption for confidentiality, SSL (Secure Sockets Layer)



ELEC 331 5

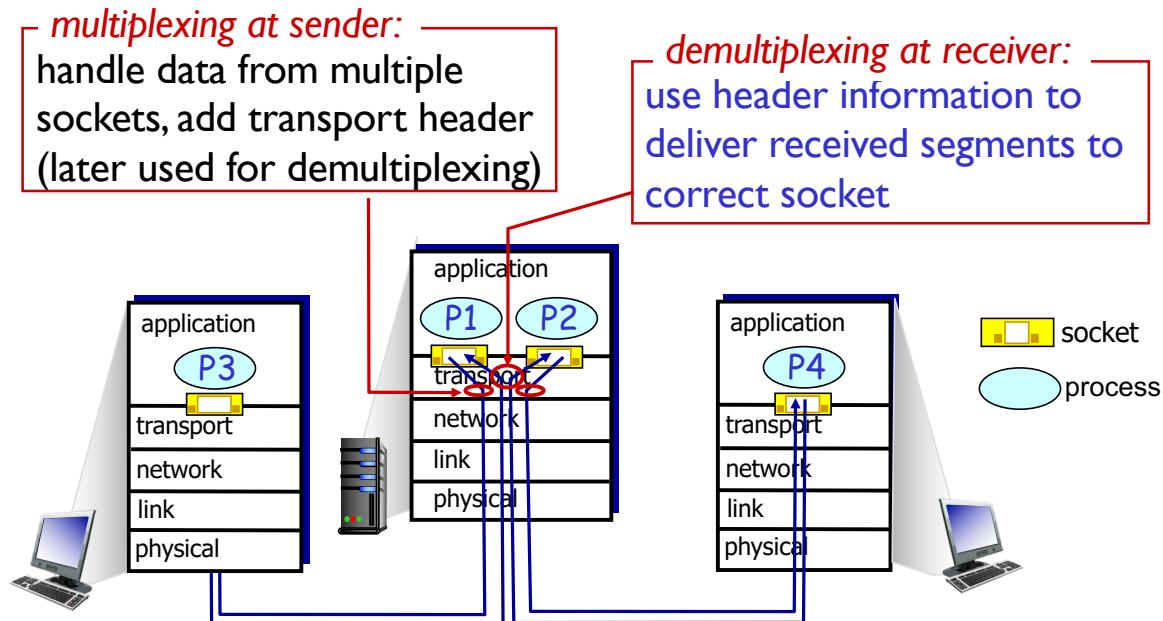
Chapter 3 outline

- ❑ 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

ELEC 331 6

Multiplexing/demultiplexing

- Extending host-to-host delivery to process-to-process delivery.



ELEC 331 7

Port Number

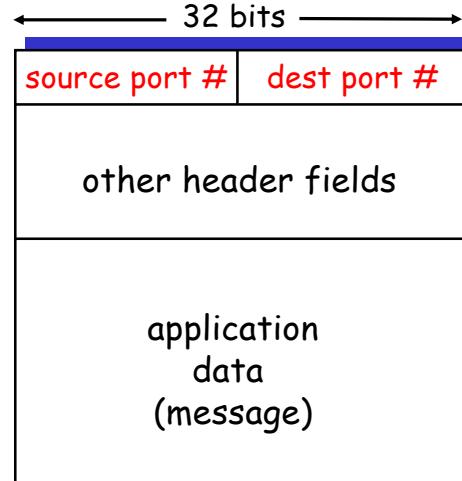
- IANA: Internet Assigned Number Authority.
<http://www.iana.org/assignments/service-names-port-numbers/service-names-port-numbers.xhtml>
- Well-known ports (0 – 1023): Assigned and controlled by IANA.
 - HTTP: 80, FTP: 21, SMTP: 25, DNS: 53
- Registered ports (1024 – 49,151): They are not assigned or controlled by IANA. They can only be registered with IANA to prevent duplication.
 - Microsoft 2000 SQL server: 1434
- Dynamic ports (49,152 – 65,535): Neither controlled nor registered. They can be used by any process.

ELEC 331 8

How demultiplexing works

❑ host receives IP datagrams

- each datagram has source IP address, destination IP address
- each datagram carries one transport layer segment
- each segment has source, destination port numbers (recall: well-known port numbers for specific applications)



❑ host uses IP addresses & port numbers to direct segment to appropriate socket

TCP/UDP segment format

ELEC 331 9

Connectionless Demultiplexing (UDP)

❑ recall: created socket has host-local port #:

```
clientSocket =  
socket(AF_INET, SOCK_DGRAM)  
clientSocket.bind(('', 19157))
```

❑ recall: when creating datagram to send into UDP socket, must specify

- destination IP address
- destination port #

❑ when host receives UDP segment:

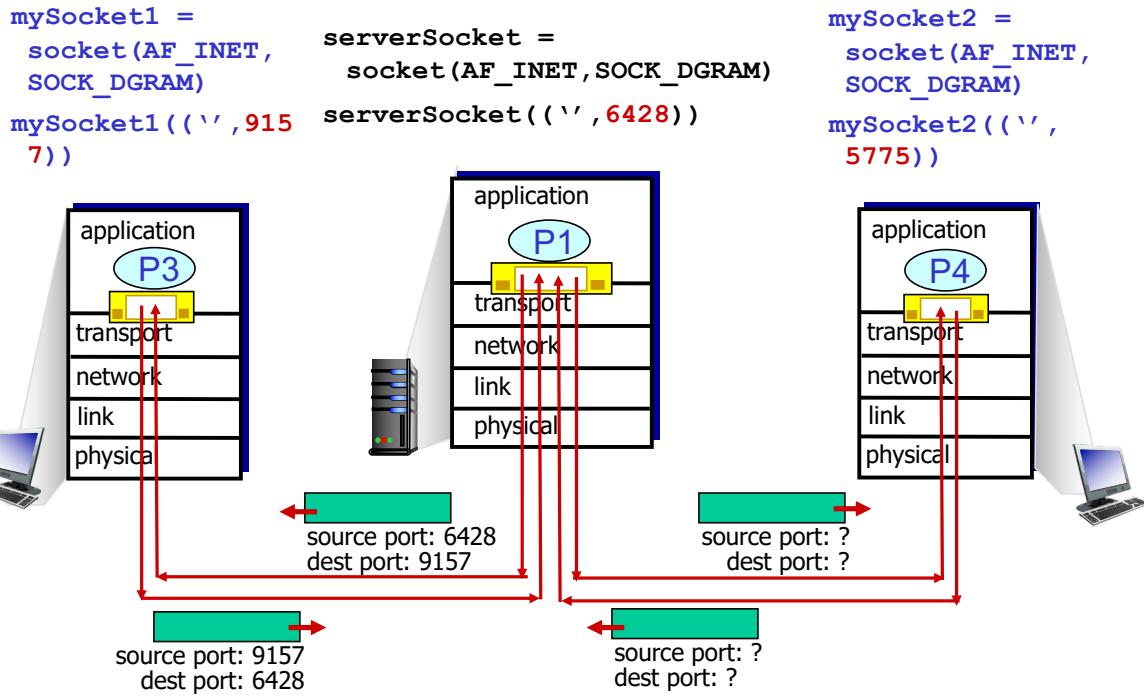
- checks destination port # in segment
- directs UDP segment to socket with that port #



IP datagrams with *same destination port #*, but different source IP addresses and/or source port numbers will be directed to *same socket* at destination

ELEC 331 10

Connectionless demux (cont.)



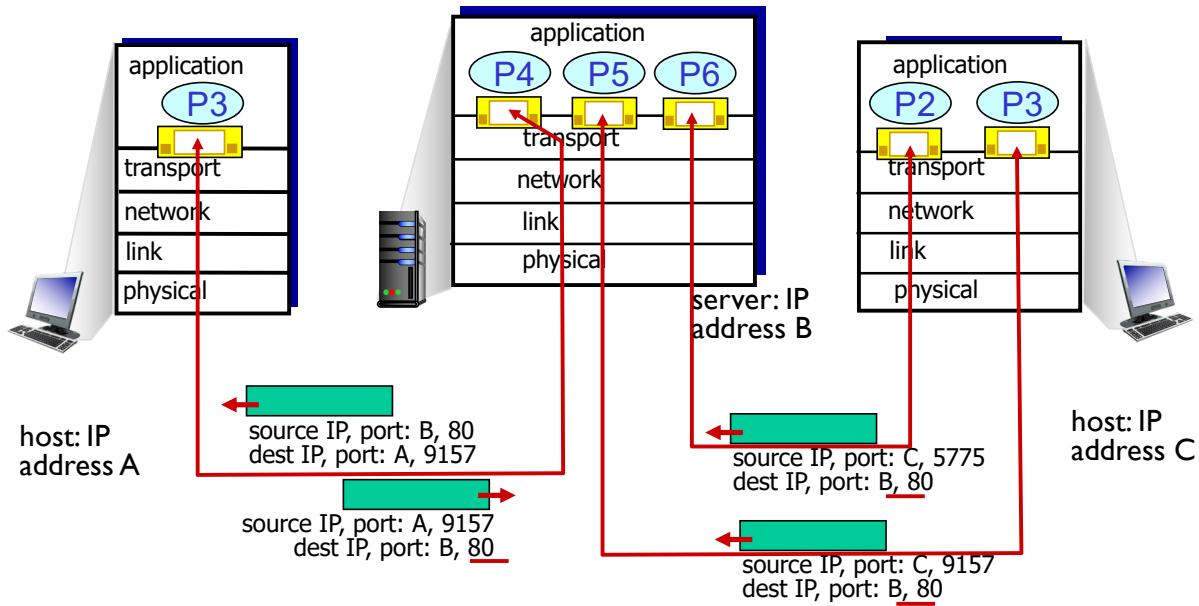
ELEC 331 11

Connection-oriented demux (TCP)

- ❑ TCP socket identified by 4-tuple:
 - source IP address
 - source port number
 - dest IP address
 - dest port number
- ❑ receiver host uses all four values to direct segment to appropriate socket
- ❑ Server host may support many simultaneous TCP sockets:
 - each connection socket is 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

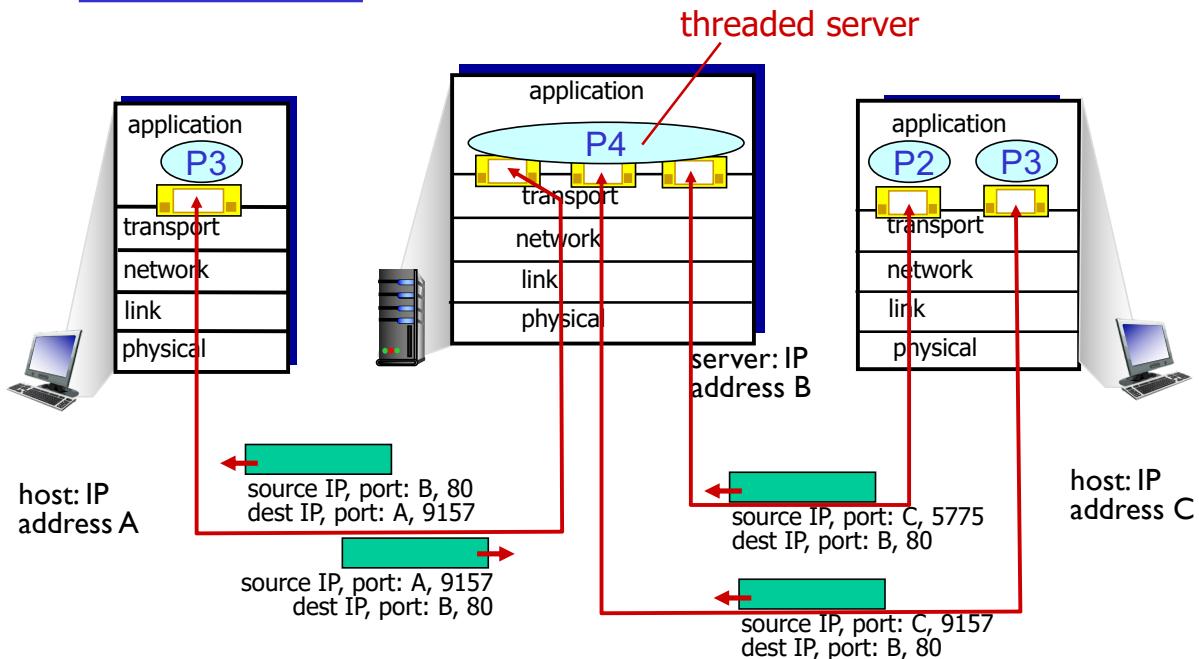
ELEC 331 12

Connection-oriented demux (cont)



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

Connection-oriented demux: Threaded Web Server



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

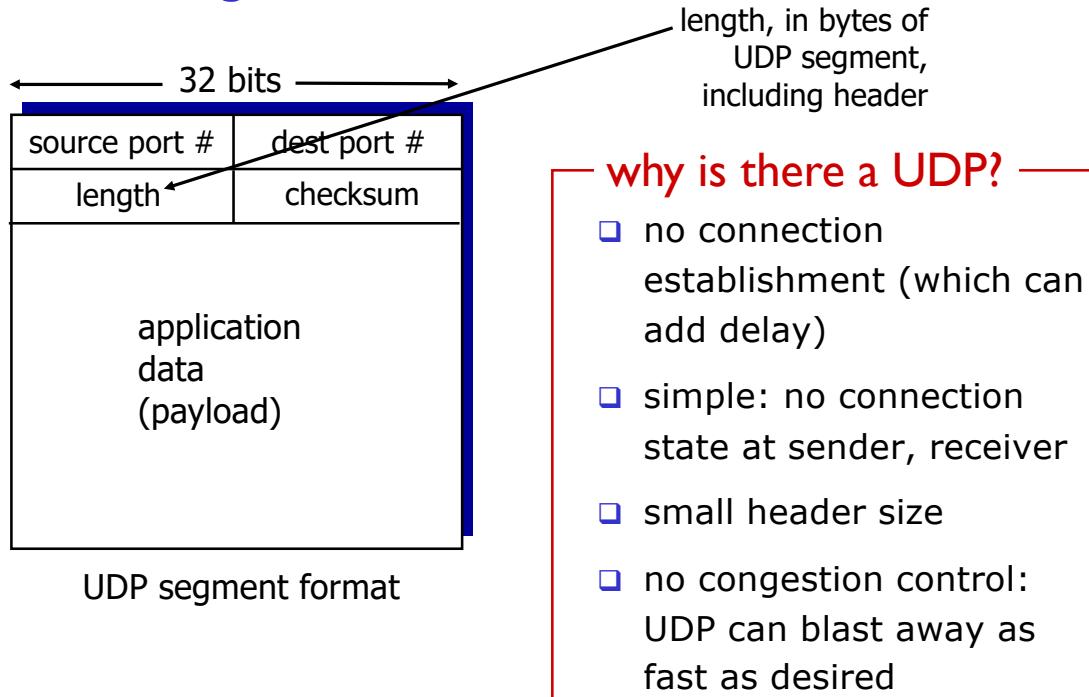
ELEC 331 15

UDP: User Datagram Protocol [RFC 768]

- ❑ “no frills,” “bare bones” Internet transport protocol
- ❑ “best effort” service, UDP segments may be:
 - lost
 - delivered out of order to application
- ❑ *connectionless*:
 - no handshaking between UDP sender, receiver
 - each UDP segment handled independently of others
- ❑ UDP use:
 - streaming multimedia apps (loss tolerant, rate sensitive)
 - DNS
 - SNMP
- ❑ Reliable transfer over UDP
 - add reliability at application layer
 - application-specific error recovery!

ELEC 331 16

UDP: segment header



ELEC 331 17

UDP checksum

Goal: detect “errors” (e.g., flipped bits) in transmitted segment

Sender:

- treat segment contents as sequence of **16-bit integers**
- checksum: addition (**1s complement of the sum**) of segment contents
- sender puts checksum value into UDP checksum field

Receiver:

- compute checksum of received segment
- check if computed checksum equals checksum field value:
 - NO - error detected
 - YES - no error detected.
But maybe errors nonetheless? More later

ELEC 331 18

Internet Checksum Example

- ❑ When adding numbers, an overflow from the most significant bit needs to be wrapped around.
- ❑ Given the three 16-bit words: 0110 0110 01100000,
0101010101010101, 1000111100001100
- ❑ The sum of the first two words is

$$\begin{array}{r} 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \\ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \\ \hline 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 1 \end{array}$$

- ❑ Adding the third word to the sum

$$\begin{array}{r} 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 1 \\ 1 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \\ \hline \end{array}$$

wraparound



Sum

Checksum

If packet has no errors, the sum at rcvr side is:

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

ELEC 331 19

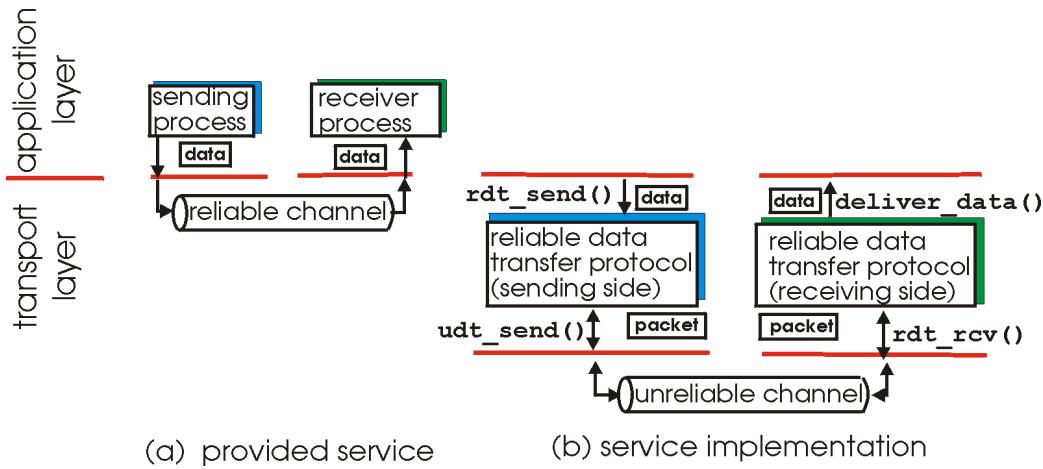
Chapter 3 outline

- ❑ 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

ELEC 331 20

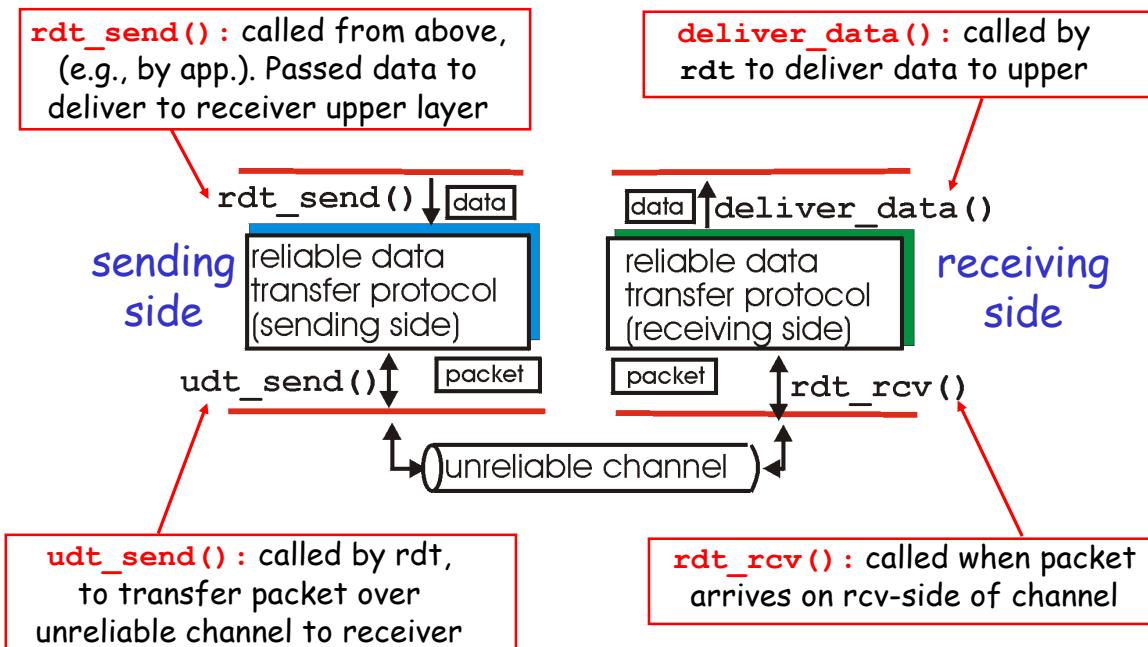
Principles of reliable data transfer

- ❑ important in application, transport, link layers



- ❑ characteristics of unreliable channel will determine complexity of **reliable data transfer protocol (rdt)**
- ❑ Consider only the case of **unidirectional data transfer**
- ❑ **Control packets** are sent from rcvng side to sending side.

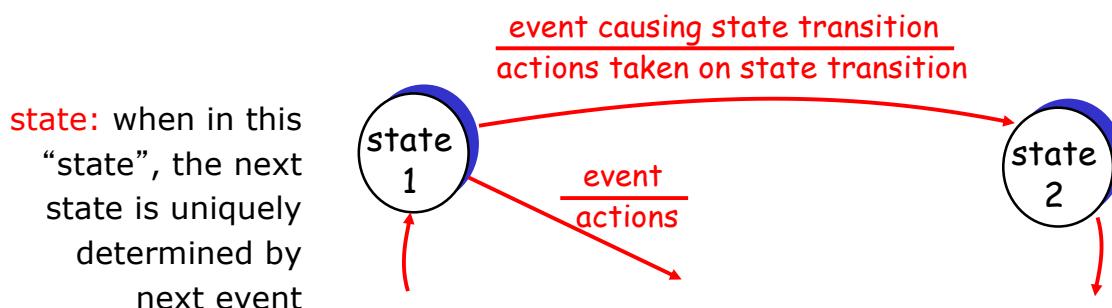
Reliable data transfer: getting started



Reliable data transfer: Getting started

We'll:

- ❑ incrementally develop sender, receiver sides of reliable data transfer (rdt) protocol.
- ❑ consider only unidirectional data transfer
 - but control info will flow on both directions!
- ❑ use **finite state machines (FSM)** to specify sender, receiver



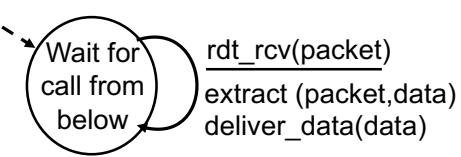
ELEC 331 23

rdt1.0: Data transfer over a perfectly reliable channel

- ❑ underlying channel is perfectly reliable
 - no bit errors
 - no loss of packets
- ❑ separate FSMs for sender, receiver:
 - sender sends data into underlying channel
 - receiver read data from underlying channel



rdt1.0: sender side



rdt1.0: receiver side

No feedback from receiver to sender

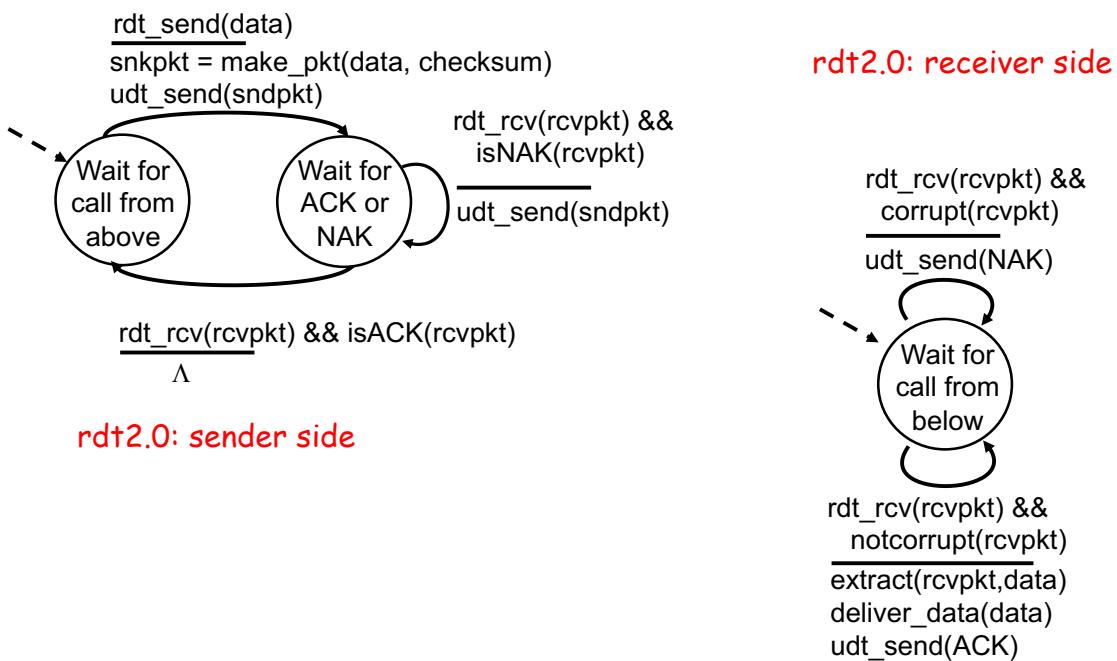
ELEC 331 24

rdt2.0: Channel with bit errors

- ❑ underlying channel may flip bits in packet
 - checksum to detect bit errors
- ❑ the question: how to recover from errors:
 - *Positive acknowledgements (ACKs)*: receiver explicitly tells sender that packet received is OK
 - *Negative acknowledgements (NAKs)*: receiver explicitly tells sender that packet had errors
 - sender **retransmits** pkt on receipt of NAK
- ❑ new mechanisms in **rdt2.0** (beyond **rdt1.0**):
 - error detection
 - receiver feedback: control msgs (ACK, NAK) rcvr -> sender
 - retransmission

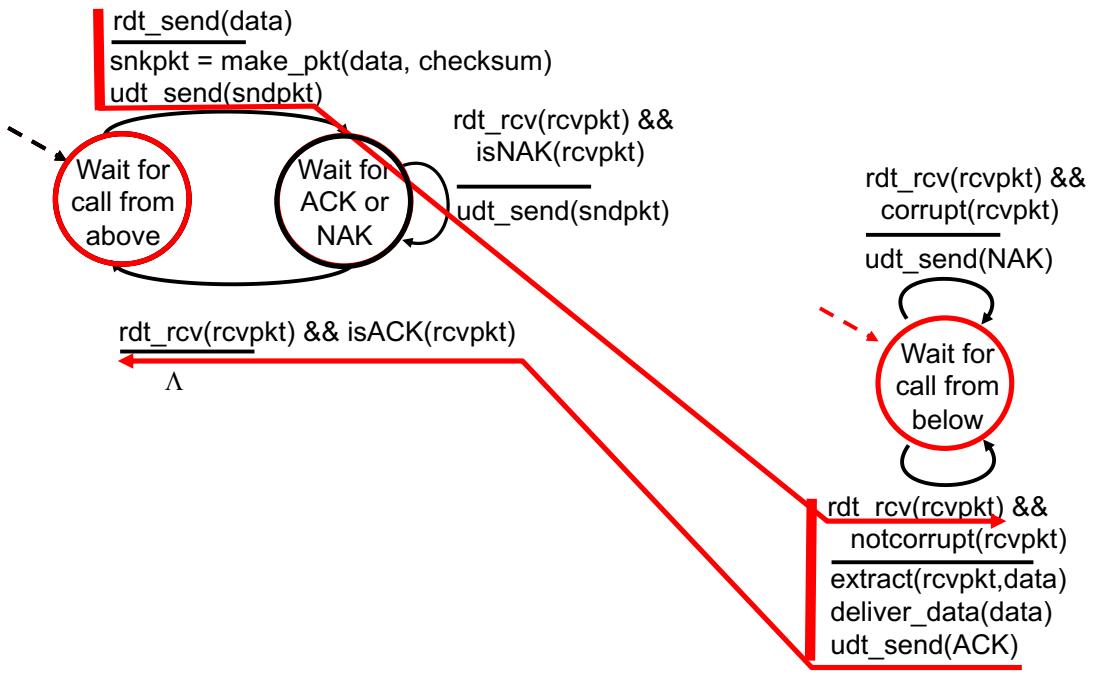
ELEC 331 25

rdt2.0: FSM specification



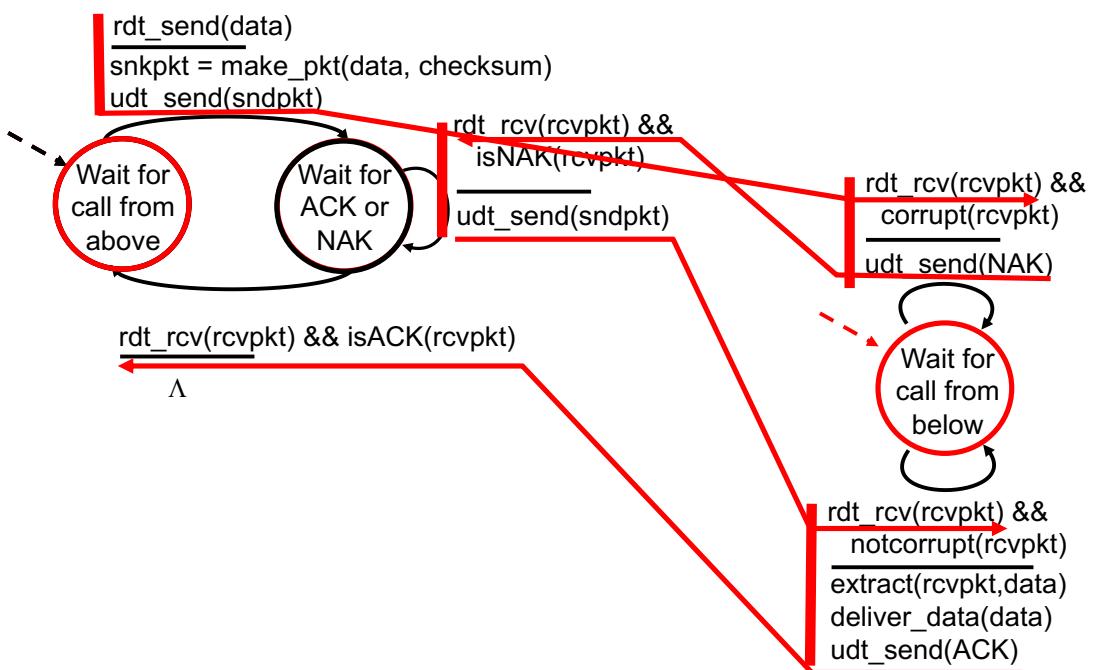
ELEC 331 26

rdt2.0: operation with no errors



ELEC 331 27

rdt2.0: error scenario



ELEC 331 28

rdt2.0 has a fatal flaw!

What happens if ACK/NAK
is corrupted?

- ❑ sender doesn't know what happened at receiver!
- ❑ can't just retransmit: possible **duplicate packets**

Handling duplicates:

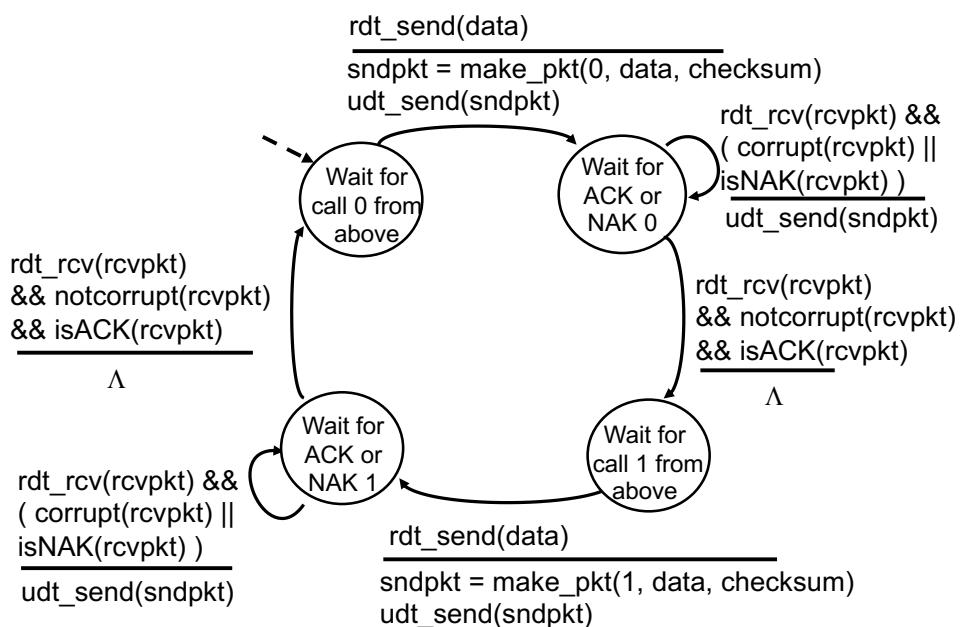
- ❑ sender adds *sequence number* to each pkt
- ❑ sender retransmits current pkt if ACK/NAK garbled
- ❑ receiver discards (doesn't deliver up) duplicate pkt

stop and wait

Sender sends one packet,
then waits for receiver
response

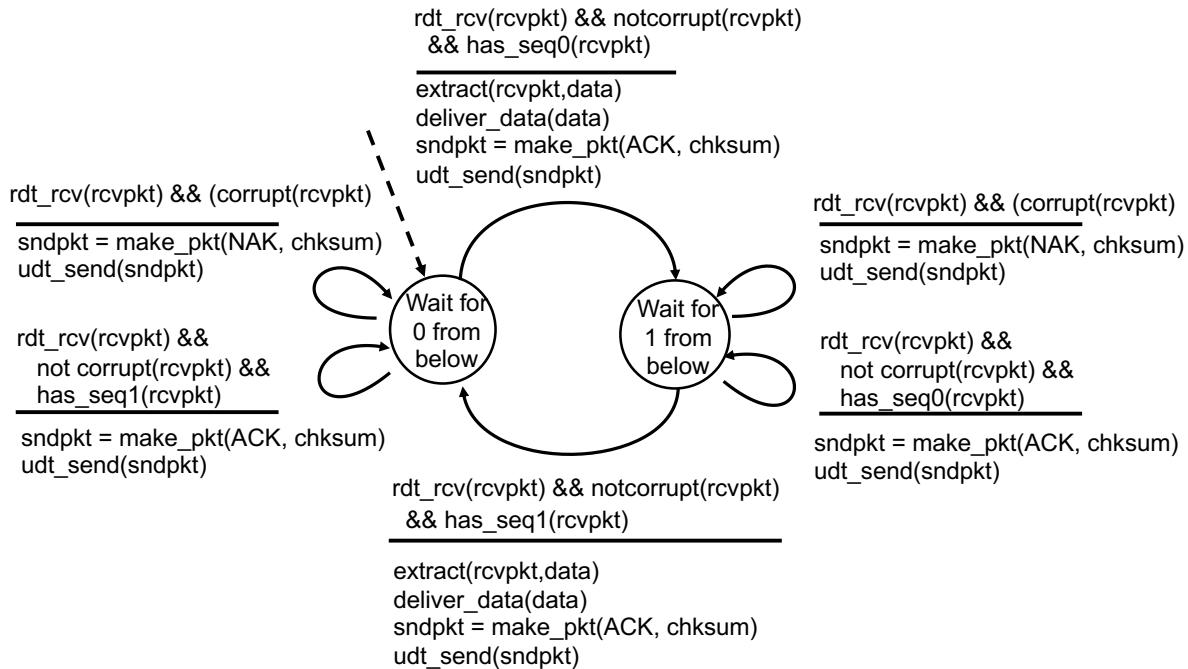
ELEC 331 29

rdt2.1: sender, handles garbled ACK/NAKs



ELEC 331 30

rdt2.1: receiver, handles garbled ACK/NAKs



ELEC 331 31

rdt2.1: discussion

Sender:

- ❑ seq # added to packet
- ❑ two seq. #'s (0,1) will suffice. Why?
- ❑ must check if received ACK/NAK corrupted
- ❑ twice as many states
 - state must “remember” whether “current” packet has 0 or 1 seq. #

Receiver:

- ❑ must check if received packet is duplicate
 - state indicates whether 0 or 1 is expected packet seq #
- ❑ note: receiver *cannot* know if its last ACK/NAK received OK at sender

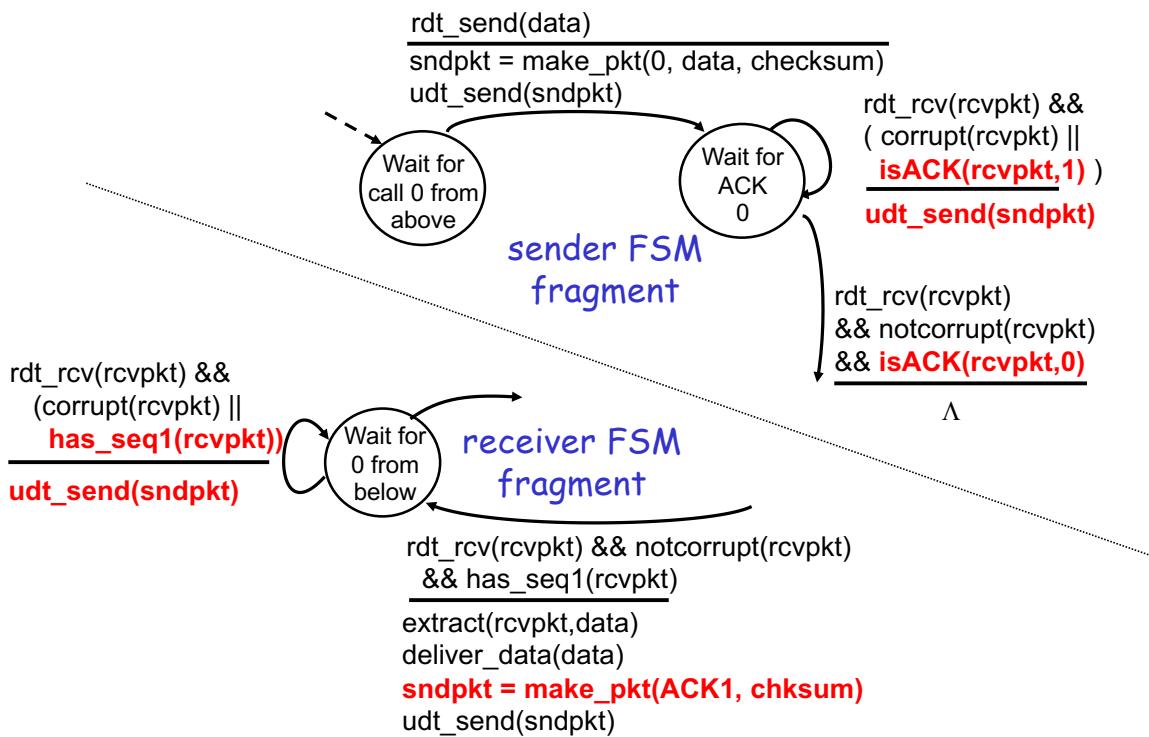
ELEC 331 32

rdt2.2: a NAK-free protocol

- ❑ same functionality as rdt2.1, using ACKs only
- ❑ instead of NAK, receiver sends **an ACK for the last correctly received packet.**
 - receiver must *explicitly include seq # of pkt being ACKed*
- ❑ duplicate ACK at sender results in same action as **NAK:** *retransmit current pkt*

ELEC 331 33

rdt2.2: sender, receiver fragments



ELEC 331 34

rdt3.0: channels with errors and loss

New assumption:

underlying channel can also lose packets (data packets or ACKs)

- ❑ How to detect packet loss?
- ❑ What to do when packet loss occurs?
- ❑ checksum, seq. #, ACKs, retransmissions will be of help, but not enough

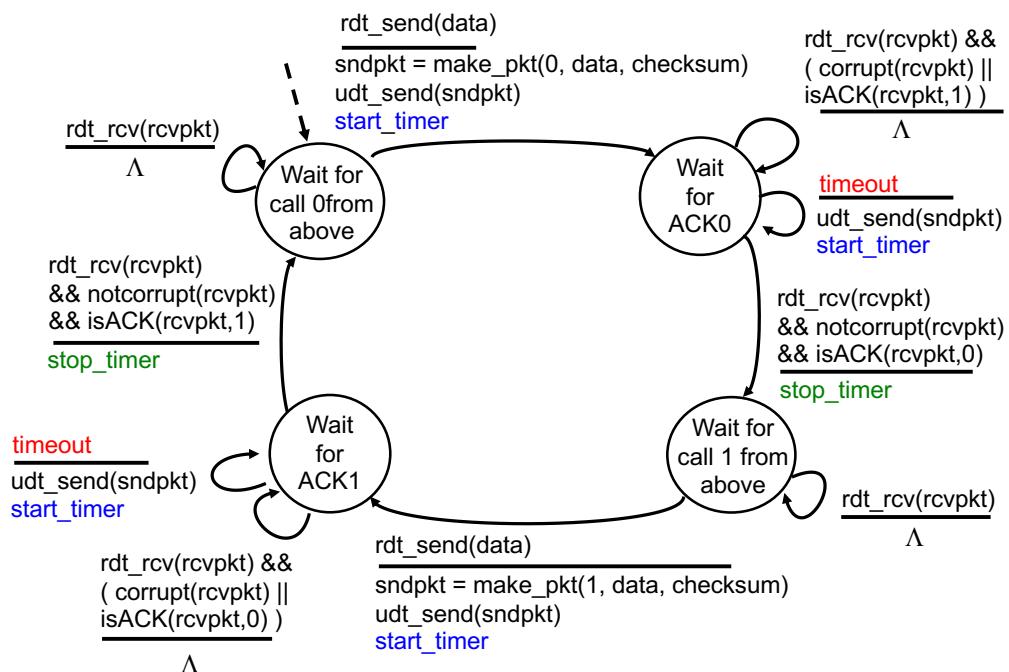
Approach: sender waits

“reasonable” amount of time for ACK

- ❑ requires countdown timer
- ❑ retransmits if no ACK is received in this time
- ❑ if packet (or ACK) just delayed (not lost):
 - retransmission will be duplicate, but use of seq. #'s already handles this
 - receiver must specify seq # of pkt being ACKed

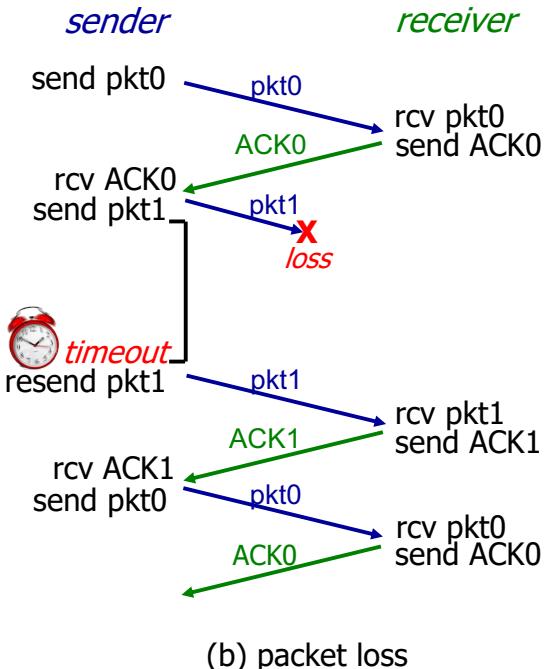
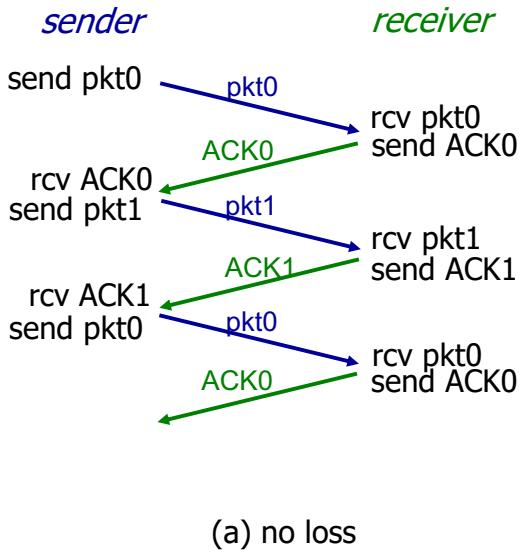
ELEC 331 35

rdt3.0 sender



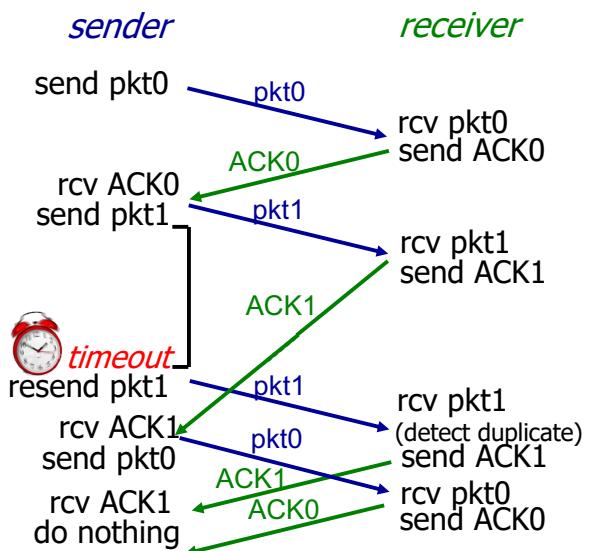
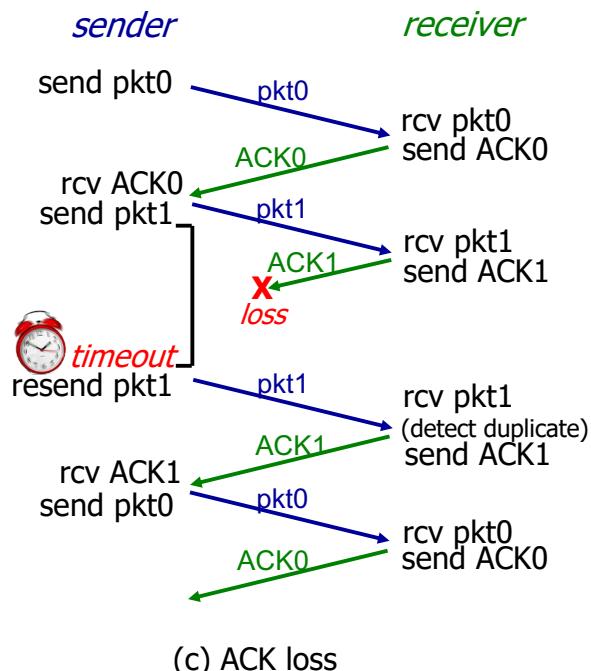
ELEC 331 36

rdt3.0 in action



ELEC 331 37

rdt3.0 in action



(d) premature timeout/ delayed ACK

Performance of rdt3.0

- ❑ rdt3.0 works, but performance stinks
- ❑ e.g., 1 Gbps link, 30 ms round-trip prop. delay, 1 KB packet

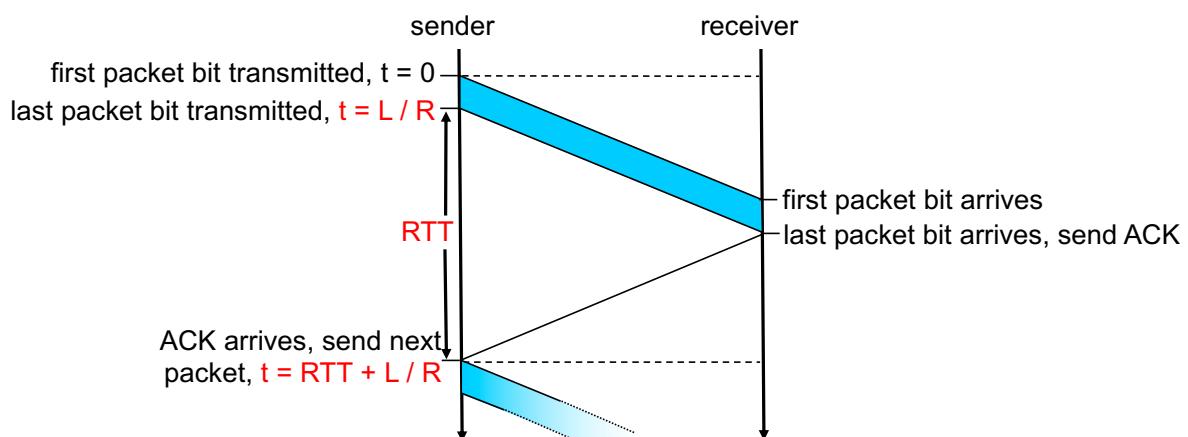
$$d_{trans} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8 \text{ kbytes/pkt}}{10^9 \text{ bits/sec}} = 8 \text{ microsec}$$

$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

- ❑ U_{sender} : Utilization (i.e., fraction of time sender busy sending bits)
- ❑ 1 KB packet every 30 msec -> 267 kbps effective throughput over 1 Gbps link
- ❑ network protocol limits use of physical resources!

ELEC 331 39

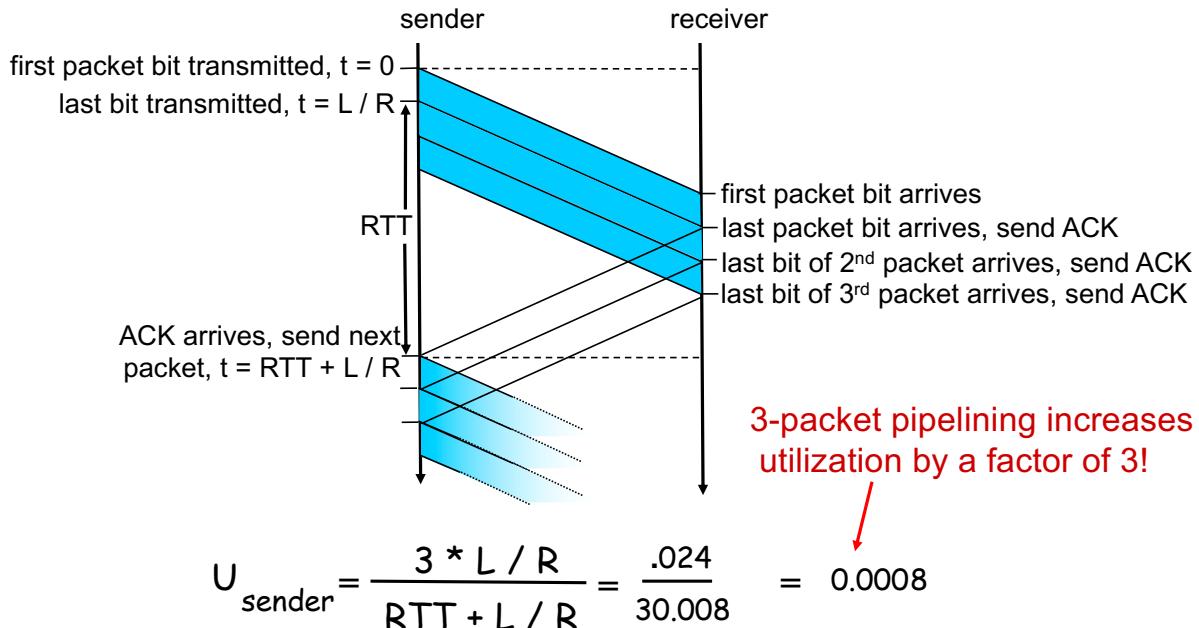
rdt3.0: stop-and-wait operation



$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

ELEC 331 40

Pipelining: increased utilization

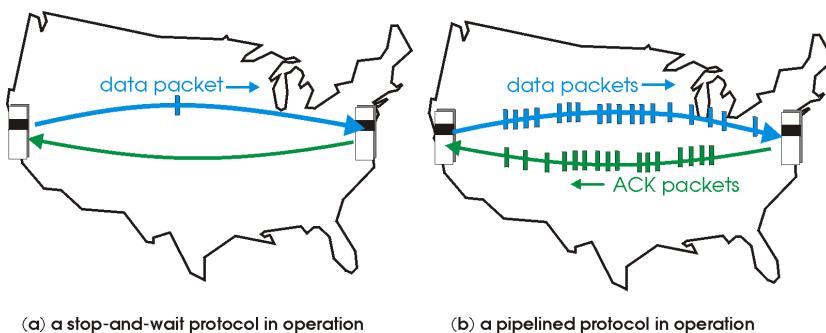


ELEC 331 41

Pipelined protocols

Pipelining: sender allows multiple, “in-flight”, yet-to-be-acknowledged pkts

- range of sequence numbers must be increased
- buffering at sender and/or receiver



- Two generic forms of pipelined protocols: *go-Back-N, selective repeat*

ELEC 331 42

Pipelining Protocols

Go-back-N: big picture

- ❑ Sender can have up to N unACKed packets in pipeline
- ❑ Receiver only sends **cumulative ACKs**
 - ❑ 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: big picture

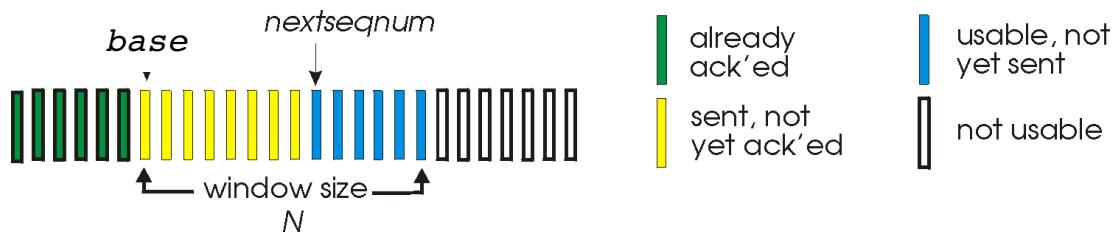
- ❑ Sender can have up to N unACKed packets in pipeline
- ❑ Receiver sends **individual ACK** for each packet
- ❑ sender maintains timer for each unACKed packet
 - ❑ When timer expires, retransmit only unACK packet

ELEC 331 43

Go-Back-N

Sender:

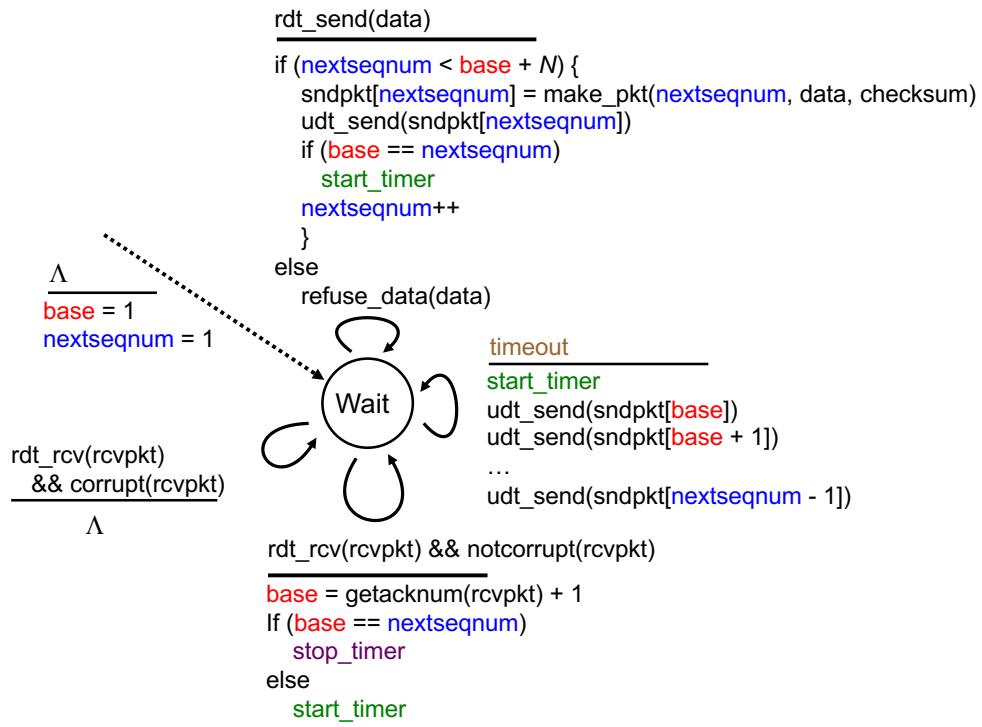
- ❑ k -bit seq # in pkt header. Range $[0, 2^k - 1]$
- ❑ “window” of up to N , consecutive unack’ed pkts allowed



- ❑ **ACK(j)**: ACKs all pkts up to, including seq # j - “**cumulative ACK**”
 - May receive duplicate ACKs (see receiver)
- ❑ **Single timer** for oldest unACKed packet
- ❑ **timeout(j)**: retransmit pkt j and all higher seq # pkts in window

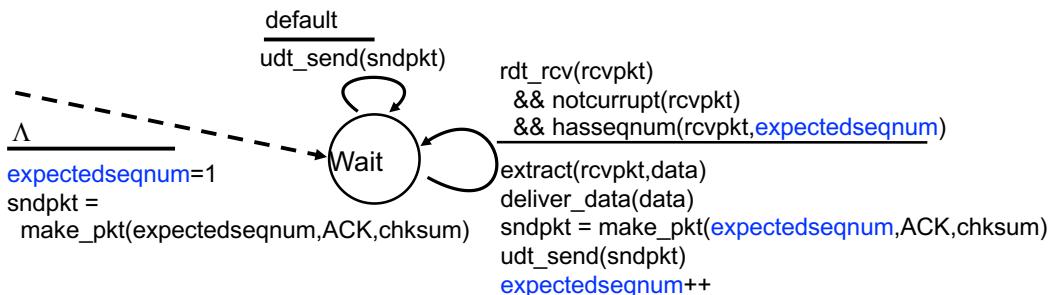
ELEC 331 44

GBN: sender extended FSM



ELEC 331 45

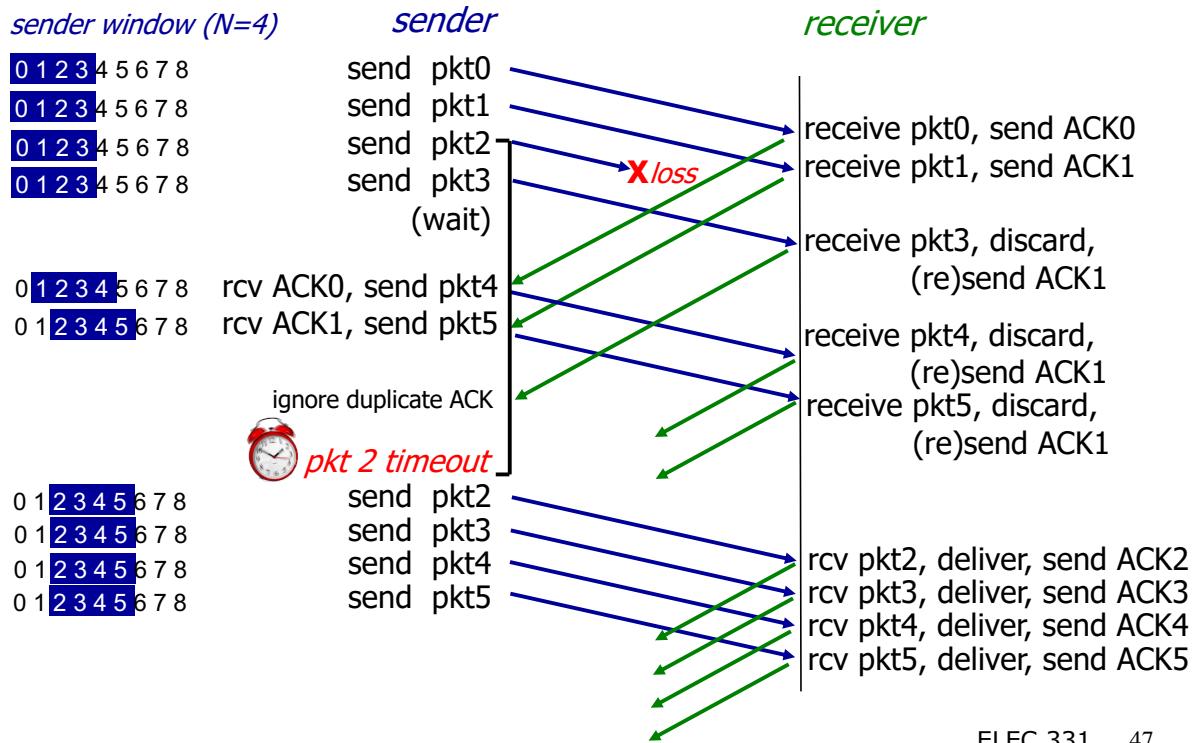
GBN: receiver extended FSM



- ❑ ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #
 - may generate duplicate ACKs
 - need only remember **expectedseqnum**
- ❑ out-of-order pkt:
 - discard (don't buffer): *no receiver buffering!*
 - re-ACK pkt with highest in-order seq #

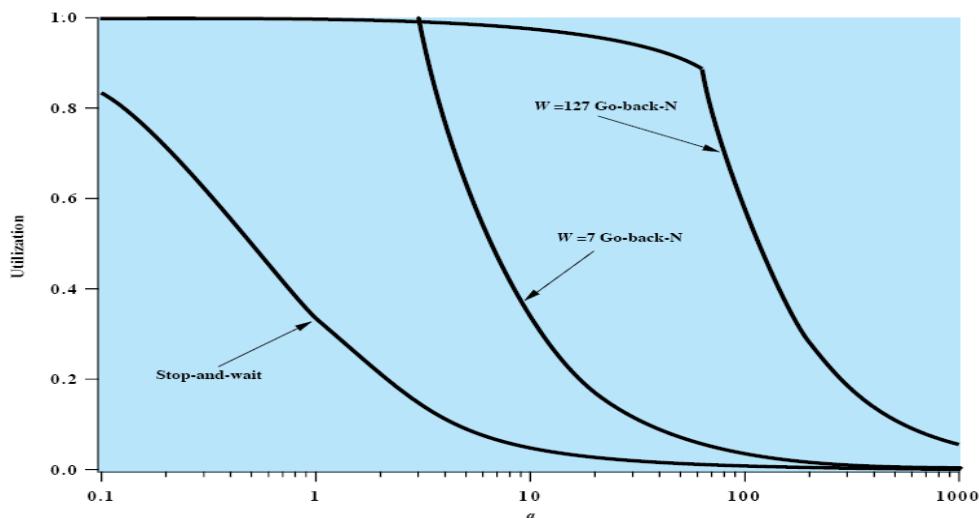
ELEC 331 46

GBN in action



ELEC 331 47

Performance of Go Back N



- $W = N$: Window Size
- $a = (\text{Propagation Delay}) / (\text{Transmission Time})$
- Packet error probability = 10^{-3} . Ignored errors in ACKs.

Source: Stallings, 7th edition

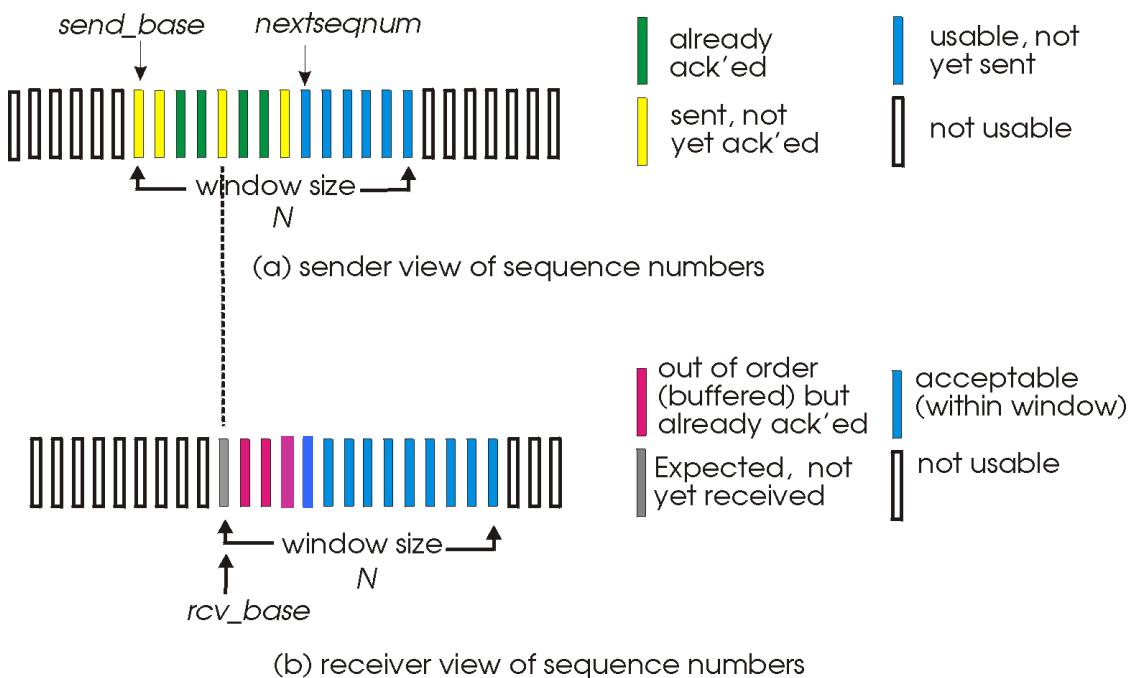
ELEC 331 48

Selective Repeat

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

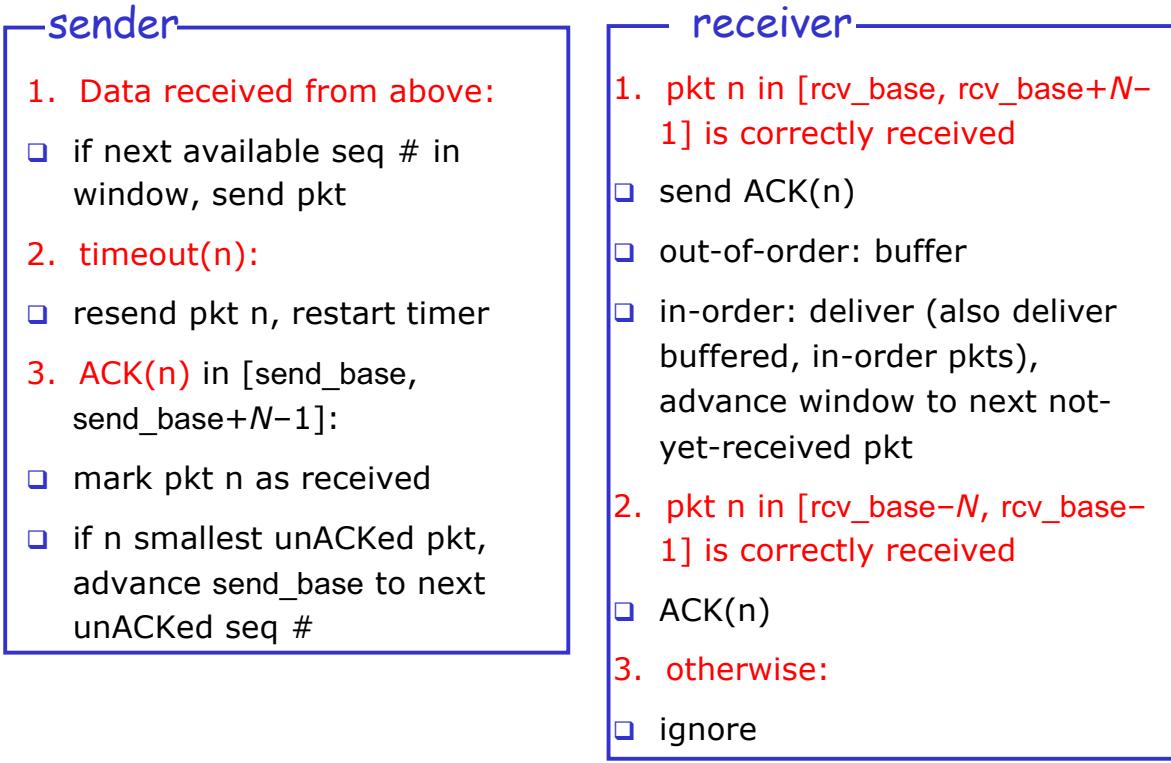
ELEC 331 49

Selective repeat: sender, receiver windows



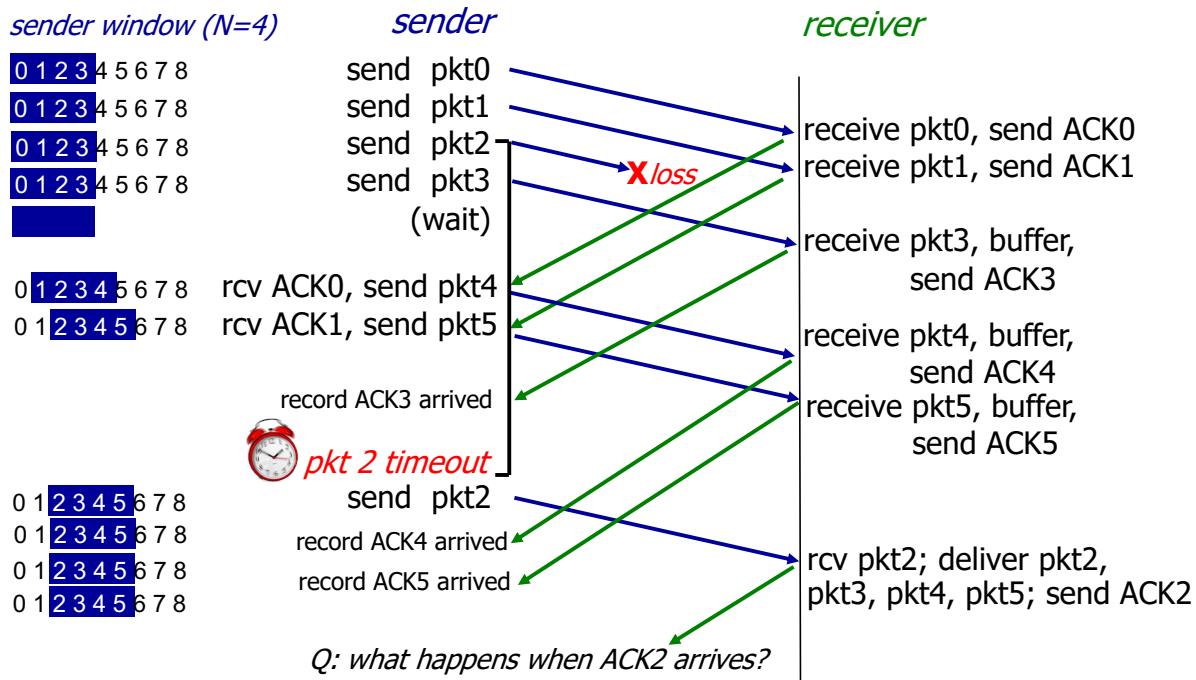
ELEC 331 50

Selective repeat



ELEC 331 51

Selective repeat in action



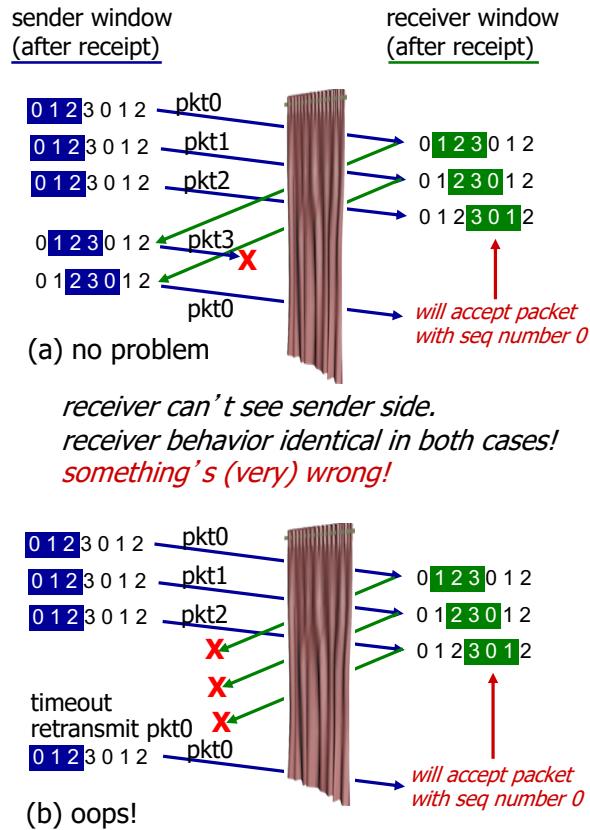
ELEC 331 52

Selective repeat: dilemma

Example:

- seq #'s: 0, 1, 2, 3
- window size = 3
- receiver sees no difference in two scenarios!
- incorrectly passes duplicate data as new in
 - (a)

Q: what relationship between seq # size and window size?



ELEC 331 53

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

ELEC 331 54

TCP: Overview (RFCs: 793, 1122, 1323, 2018, 2581)

- point-to-point:
 - one sender, one receiver
- reliable, in-order byte stream:
 - no “message boundaries”
- Use pipelining:
 - TCP congestion and flow control set window size
- send & receive buffers

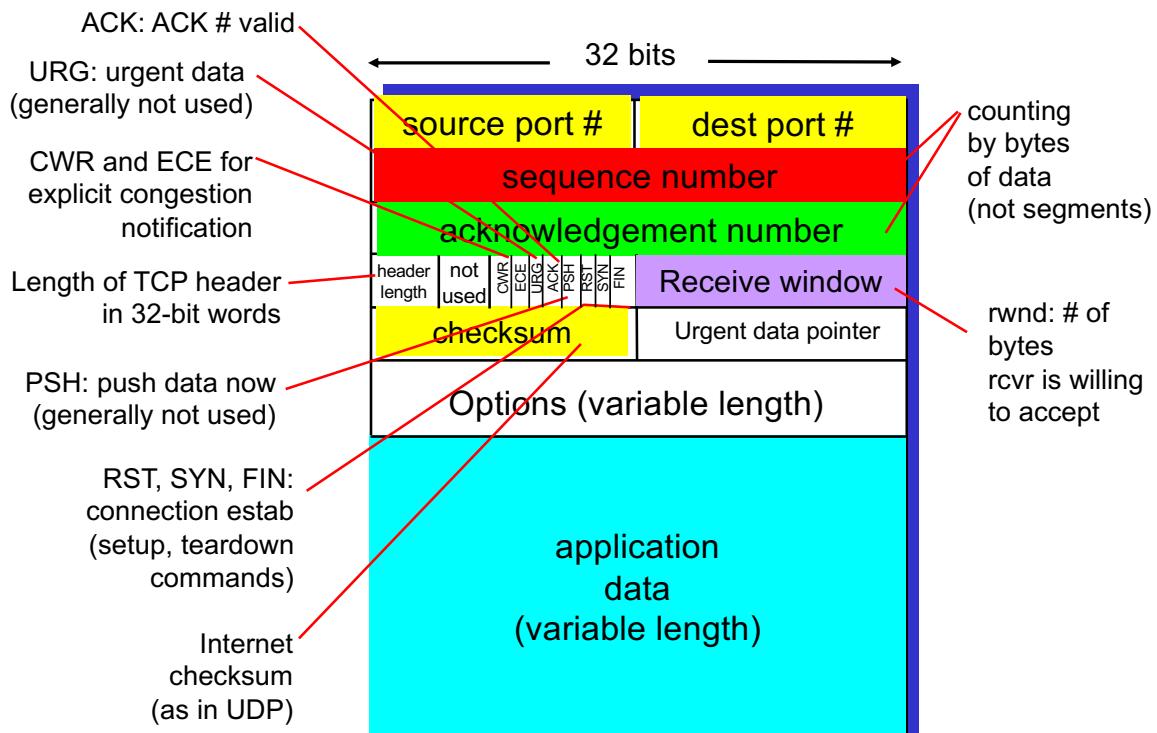


□ flow control:

- sender will not overwhelm receiver

ELEC 331 55

TCP segment structure



ELEC 331 56

TCP sequence number and ACKs

Seq. # of a segment:

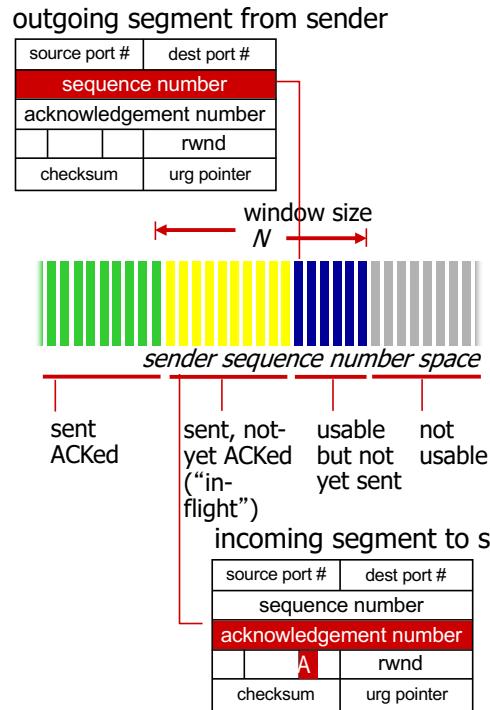
- o byte stream “number” of first byte in segment’s data

ACKs:

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

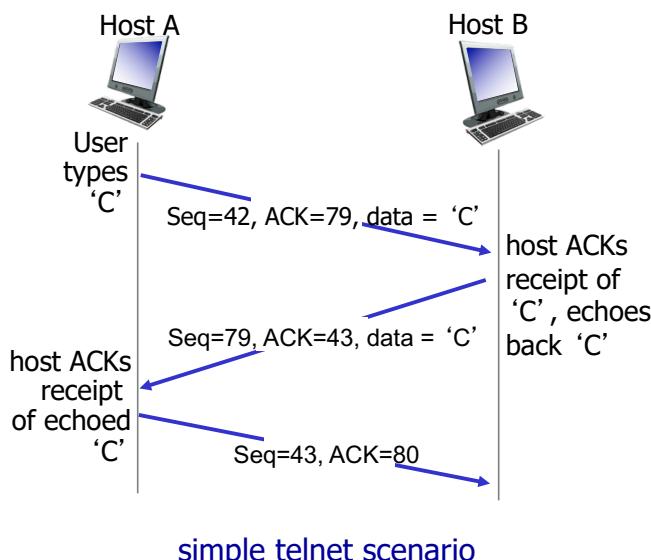
Q: how receiver handles out-of-order segments

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



ELEC 331 57

TCP sequence number and ACKs



ELEC 331 58

TCP round trip time estimation and timeout

Q: how to set TCP timeout value?

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

Q: how to estimate RTT?

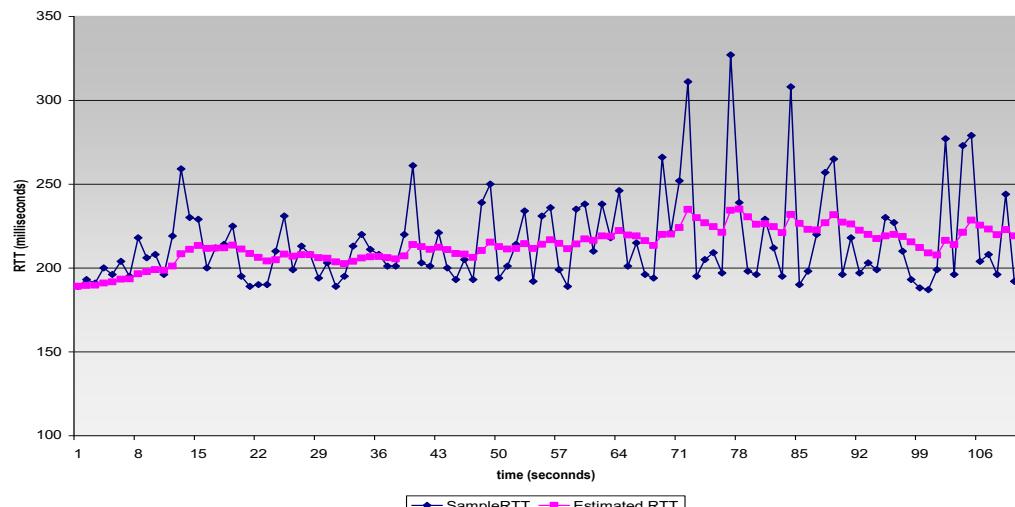
- SampleRTT**: measured time when segment is sent (passed to IP) until ACK receipt
 - ignore retransmissions
- SampleRTT** will vary, want estimated RTT “smoother”
 - average several recent measurements, not just current **SampleRTT**

ELEC 331 59

TCP round trip time estimation and timeout

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

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



ELEC 331 60

TCP round trip time estimation and timeout

Setting the timeout interval

- ❑ 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 the **timeout interval**:

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



estimated RTT

“safety margin”

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 (rdt) service

- ❑ 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

ELEC 331 63

TCP sender events:

Data rcvd from app above:

- ❑ 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

Timer timeout:

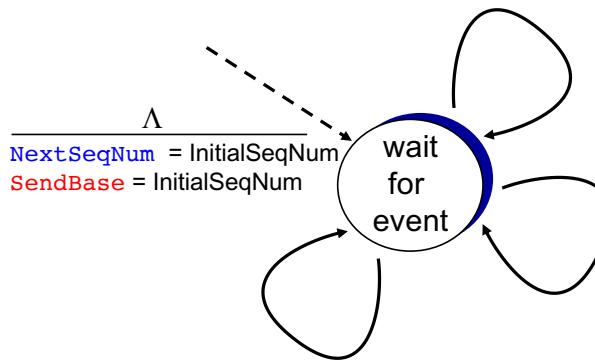
- ❑ retransmit segment that caused timeout

ACK rcvd:

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

ELEC 331 64

TCP sender (simplified)



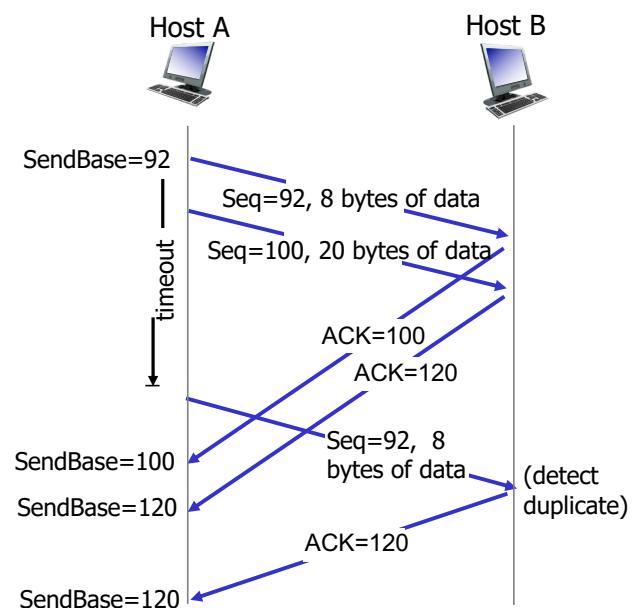
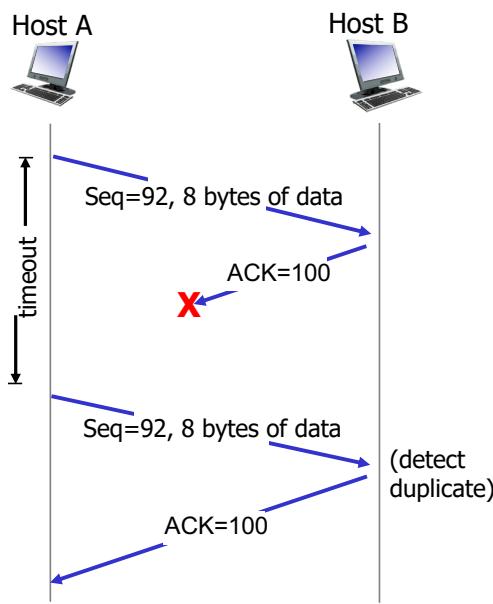
data received from application above
 create segment, seq. #: **NextSeqNum**
 pass segment to IP (i.e., “send”)
 $\text{NextSeqNum} = \text{NextSeqNum} + \text{length(data)}$
 if (timer currently not running)
 start timer

timer timeout
 retransmit not-yet-acked segment
 with smallest seq. #
 start timer

ACK received, with ACK field value y
 if ($y > \text{SendBase}$) {
 $\text{SendBase} = y$
 /* $\text{SendBase} - 1$: last cumulatively ACKed byte */
 if (there are currently any not-yet-acked segments)
 start timer
 else stop timer
 }

ELEC 331 65

TCP: retransmission scenarios

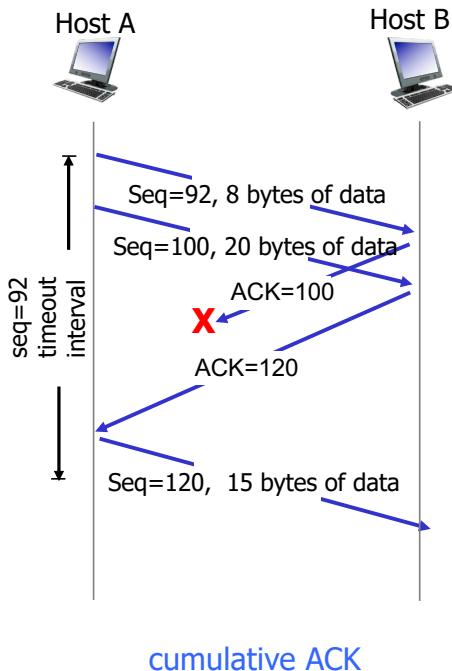


lost ACK scenario

premature timeout

ELEC 331 66

TCP retransmission scenarios (cont.)



ELEC 331 67

Doubling the Timeout Interval

- Each time TCP re-transmits, it sets the next timeout interval to **twice of previous value**,
 - rather than deriving it from the last **EstimatedRTT** and **DevRTT**
- E.g., **TimeoutInterval** = 0.75 sec before retransmission
 - After 1st retransmission, **TimeoutInterval** = 1.5 sec
 - After 2nd retransmission, **TimeoutInterval** = 3 sec
- After an ACK is received, **TimeoutInterval** is derived from the most recent values of **EstimatedRTT** and **DevRTT**.
- Limited form of congestion control.

ELEC 331 68

TCP ACK Generation Recommendation [RFC 5681]

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 500 ms 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 (which is lower end of gap)
Arrival of segment that partially or completely fills gap	Immediately send ACK, provided that segment starts at lower end of gap

ELEC 331 69

Fast Retransmit

- ❑ Timeout period often relatively long:
 - long delay before resending lost packet
- ❑ Detect lost segments via duplicate ACKs.
 - Sender often sends many segments back-to-back
 - If segment is lost, there will likely be many duplicate ACKs.

TCP fast retransmit

if sender receives 3 duplicate ACKs for same data (“**triple duplicate ACKs**”), resend unacked segment with smallest seq #

- likely that unacked segment lost, so don’t wait for timeout

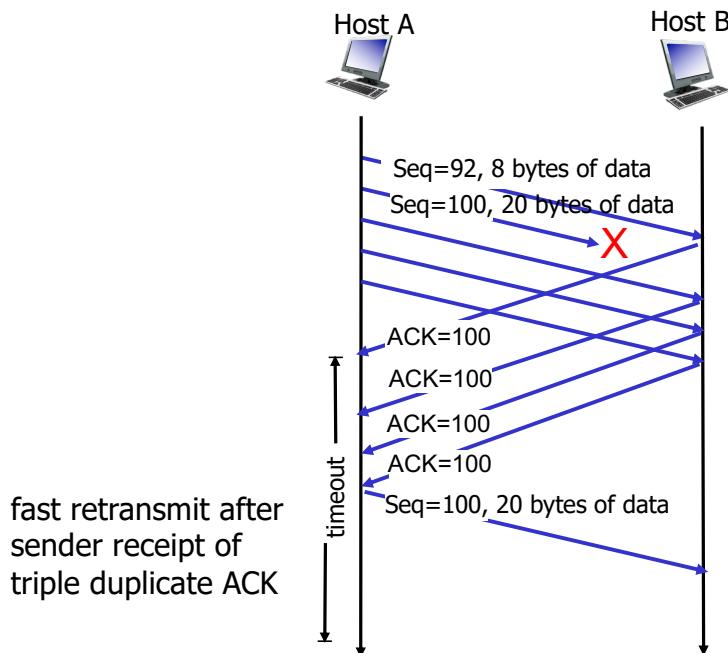
ELEC 331 70

Fast Retransmit Algorithm

```
event: ACK received, with ACK field value of y
    if (y > SendBase) {
        SendBase = y
        if (there are currently any not-yet-acknowledged segments)
            start timer
        else stop timer
    }
    else {
        increment number of duplicate ACKs received for y
        if (number of duplicate ACKs received for y == 3)
            resend segment with sequence number y
    }
}
a duplicate ACK for already ACKed segment
TCP Fast Retransmit
```

ELEC 331 71

Fast Retransmit based on Duplicate ACKs



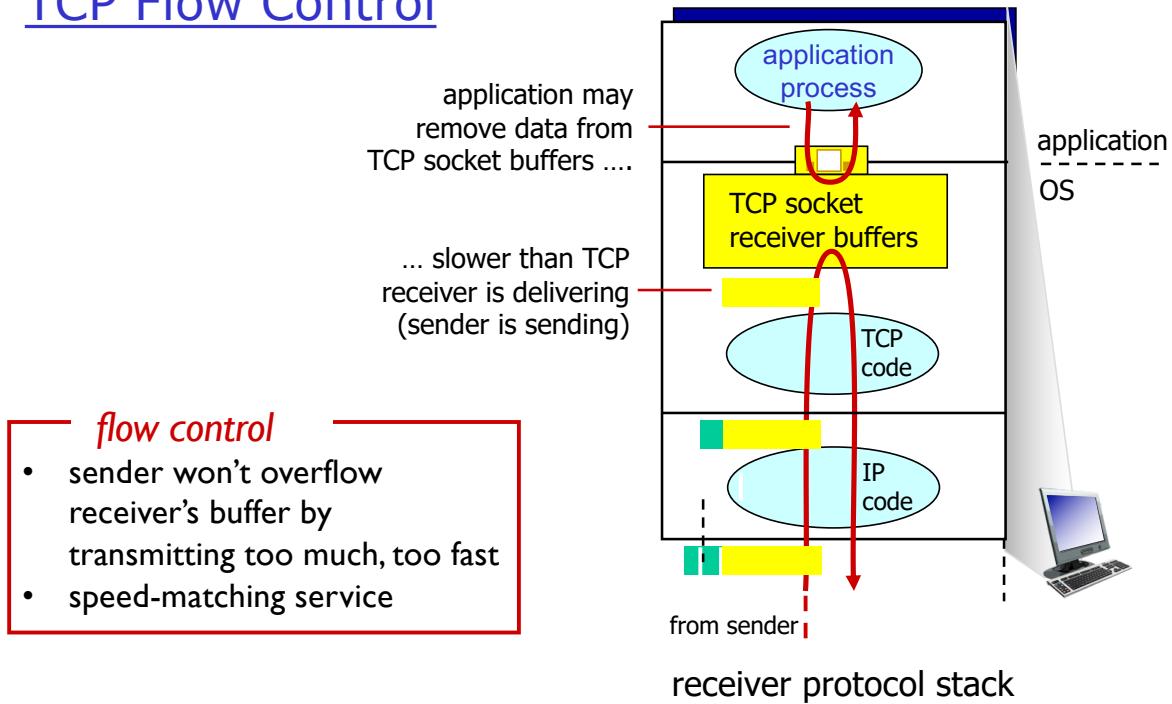
ELEC 331 72

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

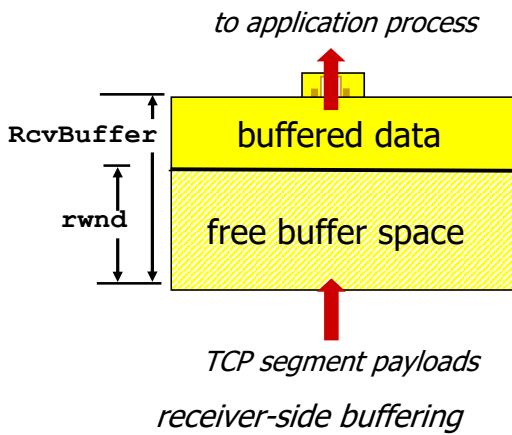
ELEC 331 73

TCP Flow Control



ELEC 331 74

TCP flow control



- ❑ spare room in buffer

$rwnd = RcvBuffer -$
 $[LastByteRcvd - LastByteRead]$

- ❑ Receiver advertises spare room by including value of **rwnd** in TCP header

- **RcvBuffer** size set via socket options (typical default is 4096 bytes)
 - many operating systems autoadjust **RcvBuffer**

- ❑ Sender limits unACKed data to **rwnd**

- guarantees receive buffer doesn't overflow

$$\boxed{\begin{aligned} & \text{LastByteSent} - \text{LastByteAcked} \\ & \leq rwnd \end{aligned}}$$

ELEC 331 75

Chapter 3 outline

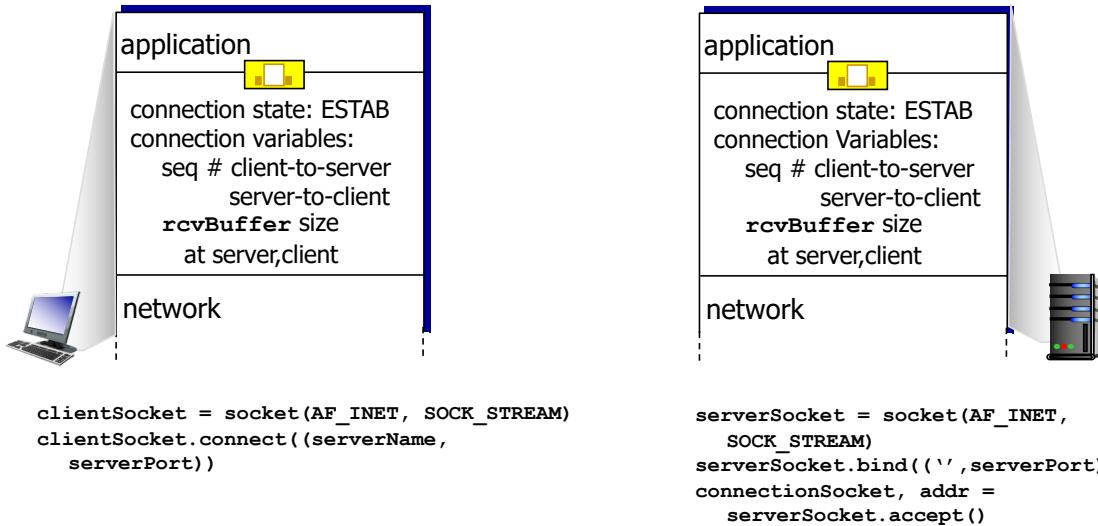
- ❑ 3.1 Transport-layer services
- ❑ 3.2 Multiplexing and demultiplexing
- ❑ 3.3 Connectionless transport: UDP
- ❑ 3.4 Principles of reliable data transfer
- ❑ 3.5 Connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - **connection management**
- ❑ 3.6 Principles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

ELEC 331 76

TCP Connection Management

before exchanging data, sender/receiver “handshake”:

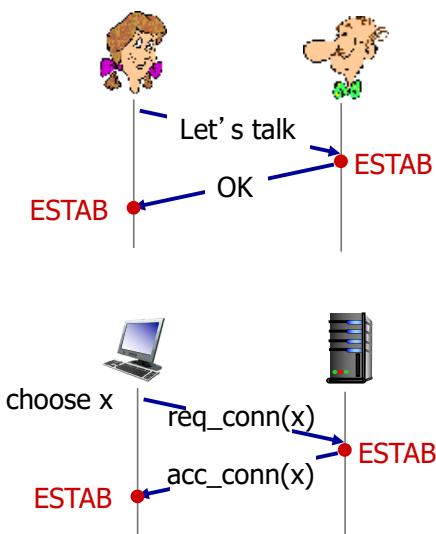
- ❑ agree to establish connection (each knowing the other willing to establish connection)
- ❑ agree on connection parameters



ELEC 331 77

Agreeing to establish a connection

2-way handshake:



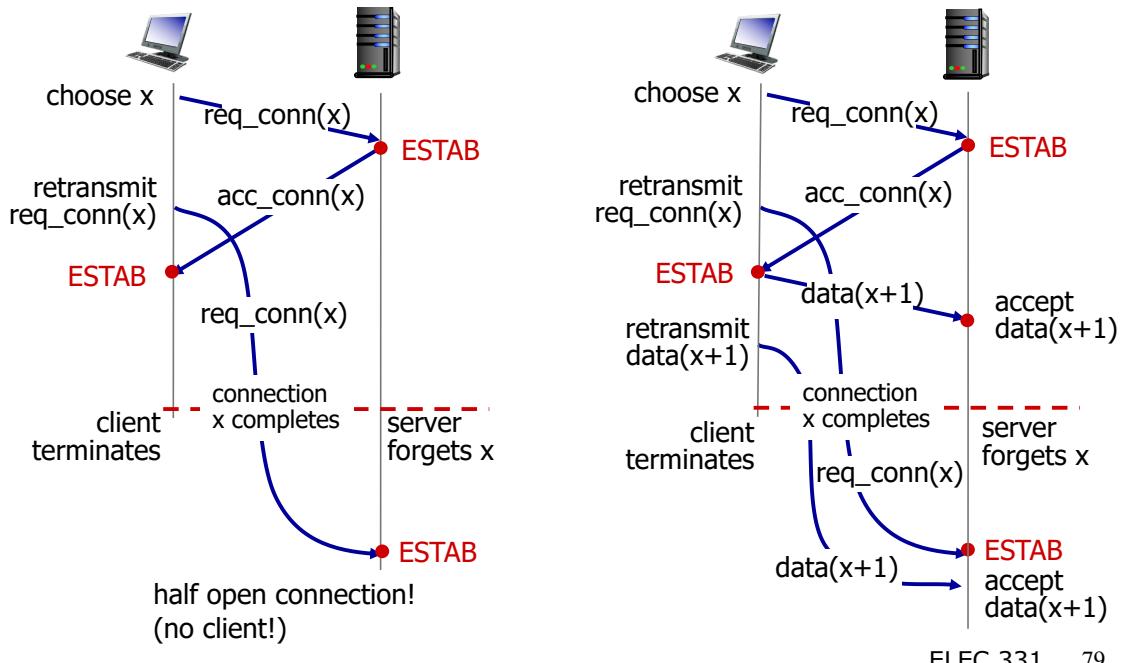
Q: Will 2-way handshake always work in network?

- ❑ variable delays
- ❑ retransmitted messages (e.g. req_conn(x)) due to message loss
- ❑ message reordering
- ❑ can't "see" other side

ELEC 331 78

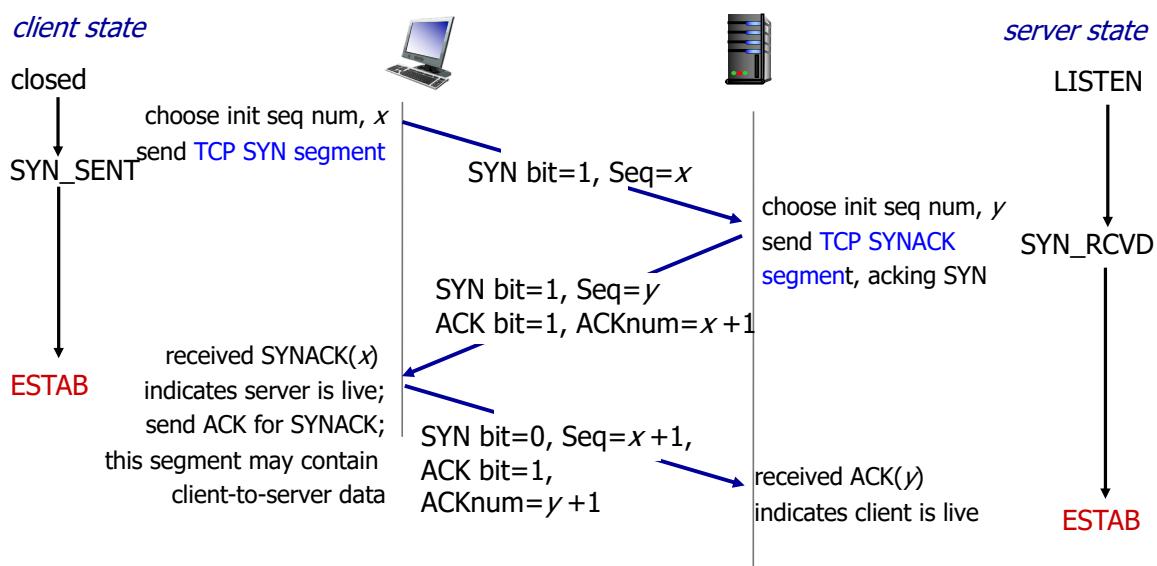
Agreeing to establish a connection

2-way handshake failure scenarios:



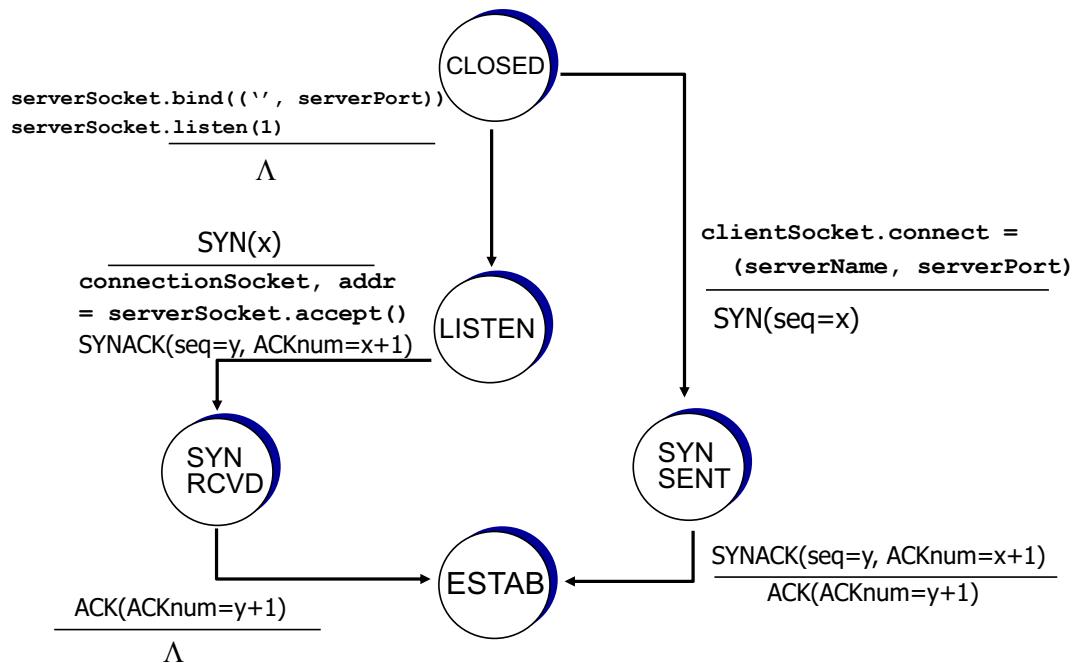
ELEC 331 79

TCP 3-way handshake



ELEC 331 80

TCP 3-way handshake: FSM



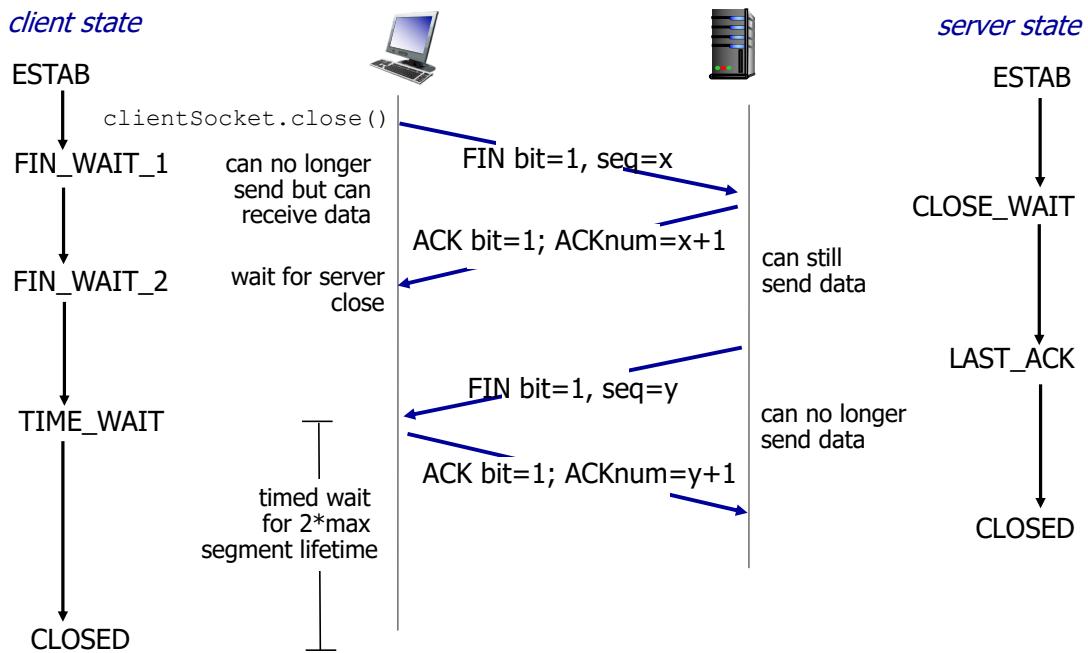
ELEC 331 81

TCP: closing a connection

- ❑ client, server each close their side of connection
 - send TCP segment with FIN bit = 1
- ❑ respond to received FIN with ACK
 - on receiving FIN, ACK can be combined with own FIN
- ❑ simultaneous FIN exchanges can be handled

ELEC 331 82

TCP: closing a connection



ELEC 331 83

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

ELEC 331 84

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)
- ❑ a top-10 problem!

ELEC 331 85

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

ELEC 331 86

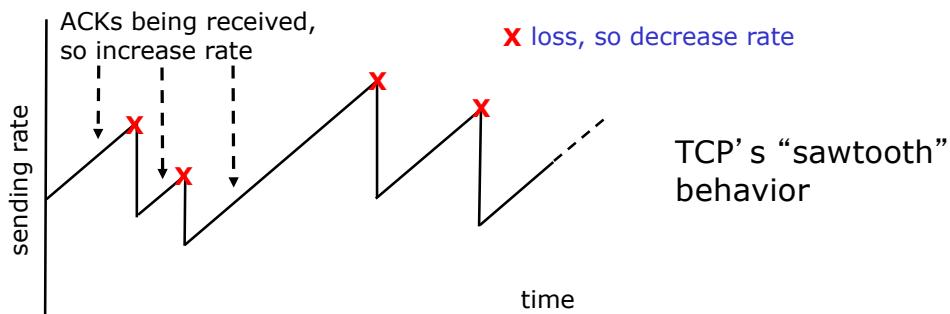
TCP Congestion Control

- ❑ **Goal:** TCP sender should transmit as fast as possible, but without congesting network
 - **Q:** How to find the rate just below congestion level?
- ❑ Decentralized: each TCP sender sets its own rate, based on **implicit feedback**:
 - **ACK:** segment received (a good thing!), network not congested, so **increase** sending rate
 - **lost segment:** assume loss due to congestion, so **decrease** sending rate

ELEC 331 87

TCP congestion control: Bandwidth probing

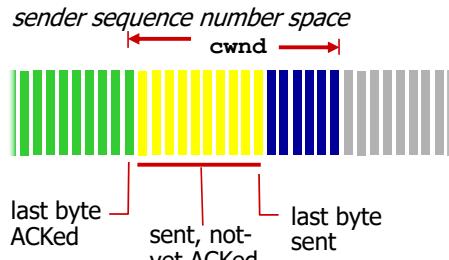
- ❑ “**probing for bandwidth**”: increase transmission rate on receipt of ACK, until eventually loss occurs, then decrease transmission rate
 - continue to increase on ACK, decrease on loss (since available bandwidth is changing, depending on other connections in network)



- ❑ **Q:** how fast to increase/decrease?
 - details to follow

ELEC 331 88

TCP Congestion Control: Details

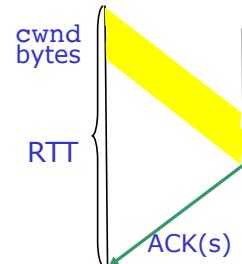


- ❑ Congestion window: **cwnd**
- ❑ sender limits rate by limiting number of unACKed bytes “in pipeline”:

$$\begin{aligned} & \text{LastByteSent} - \text{LastByteAcked} \\ & \leq \min\{\text{rwnd}, \text{cwnd}\} \end{aligned}$$

- ❑ Roughly,

$$\text{rate} \sim \frac{\text{cwnd}}{\text{RTT}} \text{ Bytes/sec}$$



- ❑ **cwnd** is dynamic, function of perceived network congestion

TCP Congestion Control: More Details

Loss event: reduce **cwnd**

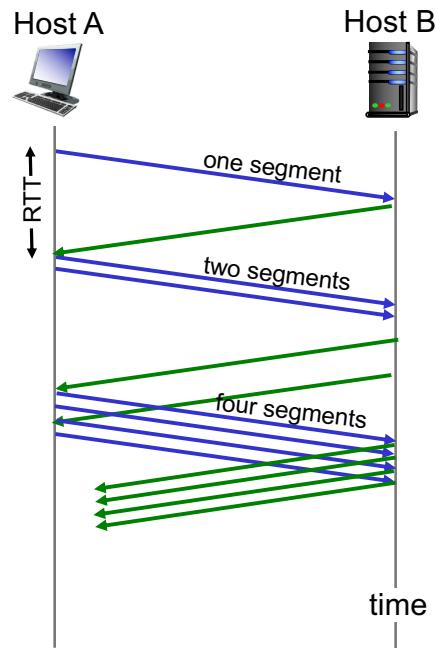
- ❑ **timeout**: no response from receiver
 - cut **cwnd** to 1 MSS
- ❑ **3 duplicate ACKs**: at least some segments getting through (recall fast retransmit)
 - cut **cwnd** in half, less aggressively than on timeout

ACK received: increase **cwnd**

- ❑ **Slow Start** phase:
 - increase exponentially fast (despite name) at connection start, or following timeout
- ❑ **Congestion Avoidance**:
 - increase linearly

TCP Slow Start

- ❑ when connection begins, increase rate exponentially until first loss event:
 - initially $cwnd = 1$ MSS
 - e.g., if $MSS = 500$ bytes, $RTT = 200$ ms, initial rate = 20 kbps
 - increment $cwnd$ by one MSS for each new ACK received
 - result in double $cwnd$ every RTT
- ❑ summary: initial rate is slow but ramps up exponentially fast



ELEC 331 91

TCP: detecting, reacting to loss

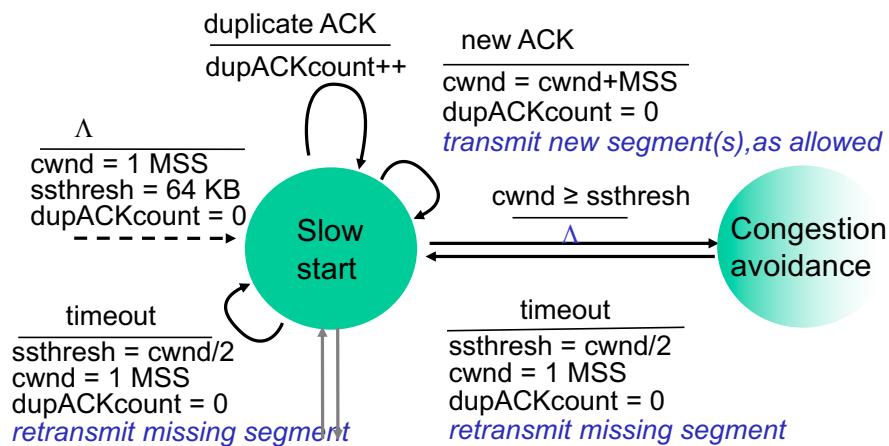
- ❑ loss indicated by timeout:
 - $cwnd$ set to 1 MSS;
 - window then grows exponentially (as in slow start) to threshold, then grows linearly
- ❑ loss indicated by 3 duplicate ACKs: TCP Reno
 - dup ACKs indicate network capable of delivering some segments
 - $cwnd$ is cut in half window then grows linearly
- ❑ TCP Tahoe always sets $cwnd$ to 1 (timeout or 3 duplicate acks)

ELEC 331 92

Transferring into/out of Slow Start

ssthresh: cwnd threshold maintained by TCP

- ❑ on loss event: set **ssthresh** to **cwnd/2**
 - remember (half of) TCP rate when congestion last occurred
- ❑ when **cwnd ≥ ssthresh**: transition from slow start to congestion avoidance phase



ELEC 331 93

TCP: Congestion avoidance

- ❑ when **cwnd > ssthresh**, increase **cwnd** linearly
 - increase **cwnd** by 1 MSS per RTT
- ❑ Implementation:
 $cwnd = cwnd + MSS(MSS/cwnd)$
 for each new ACK received
 - e.g., **MSS** = 1460 bytes, **cwnd** = 14600 bytes

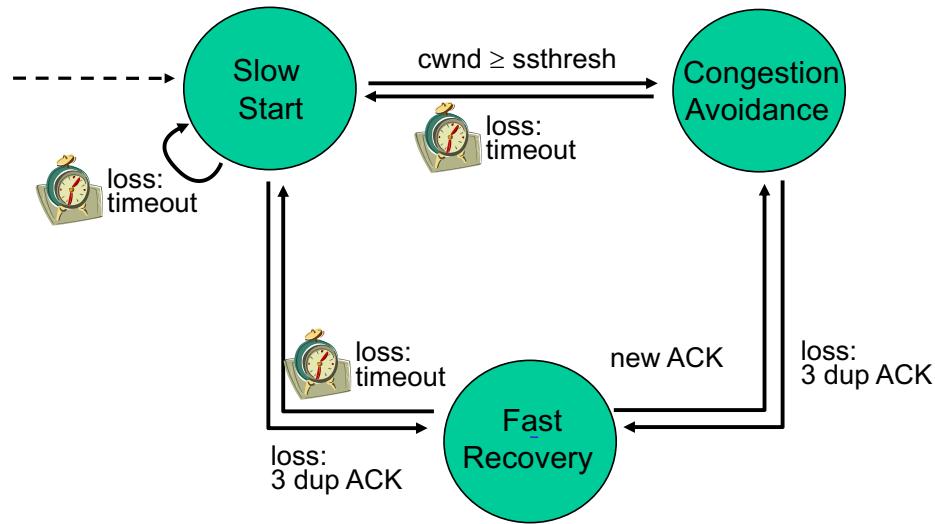
	cwnd (bytes)
	14600
1 st new ACK	14746
2 nd new ACK	14891
...	...
10 th new ACK	16000

AIMD

- ❑ ACKs: increase **cwnd** by 1 MSS per RTT: **additive increase**
 - ❑ loss: cut **cwnd** in half (non-timeout-detected loss): **multiplicative decrease**
 - ❑ **AIMD: Additive Increase Multiplicative Decrease**

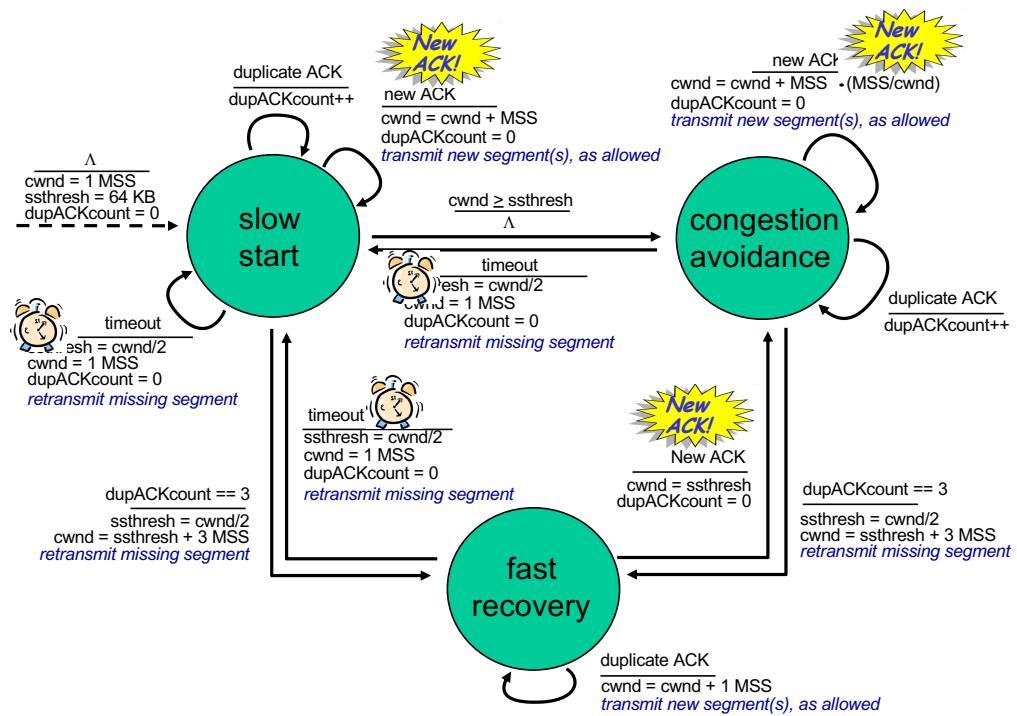
ELEC 331 94

TCP congestion control FSM: Overview



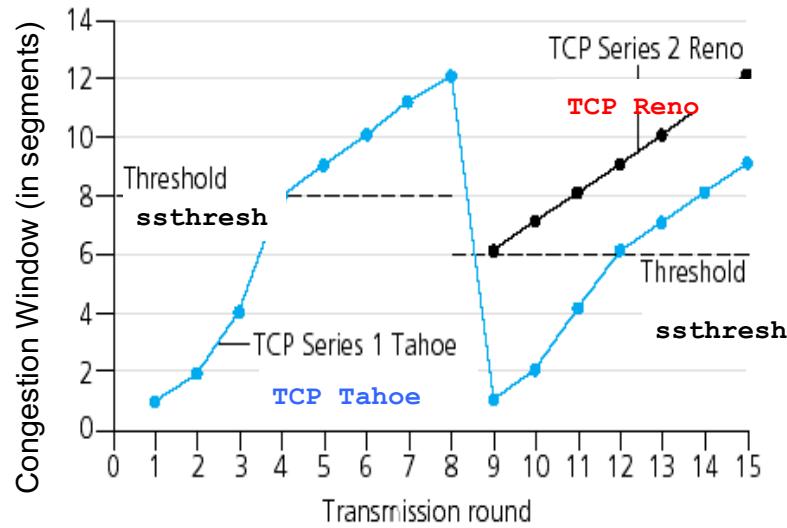
ELEC351 95

TCP congestion control FSM: details



ELEC 331 96

Popular Flavors of TCP



ELEC 331 97

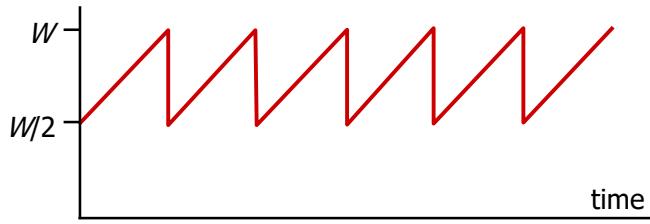
Summary: TCP Reno Congestion Control

- ❑ When **cwnd** is below **ssthresh**, sender in **slow-start** phase, window grows exponentially.
- ❑ When **cwnd** is above **ssthresh**, sender is in **congestion-avoidance** phase, window grows linearly.
- ❑ When **triple duplicate ACK** event occurs, **ssthresh** is set to **cwnd/2** and **cwnd** is set to **ssthresh**.
- ❑ When **timeout** occurs, **ssthresh** is set to **cwnd/2** and **cwnd** is set to 1 MSS.

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

TCP throughput

- ❑ What's the average throughput of TCP as a function of congestion window size and RTT?
 - Ignore slow start
- ❑ Let W be the congestion window size when loss occurs.
- ❑ When window is W , throughput is W/RTT
- ❑ Just after loss, window drops to $W/2$, throughput to $W/2RTT$.
- ❑ Assumption: RTT and W are approximately constant over the duration of the connection.
- ❑ Average throughout: $0.75 W/RTT$



ELEC 331 99

TCP Future

- ❑ Example: 1500-byte segments, 100 ms RTT, want 10 Gbps throughput
- ❑ Requires average window size = 83,333 in-flight segments
- ❑ Throughput in terms of loss rate L :

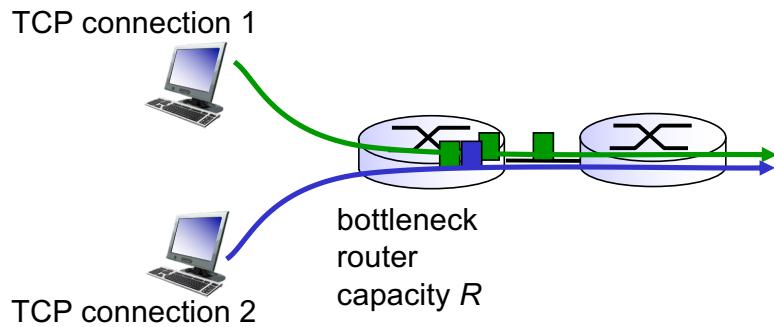
$$\frac{1.22 \text{ MSS}}{RTT\sqrt{L}}$$

- ❑ $\rightarrow L = 2 \times 10^{-10}$ (i.e., one loss event every 5×10^9 segments)
- ❑ New versions of TCP for high-speed are needed!

ELEC 331 100

TCP Fairness

Fairness goal: if K TCP sessions share same bottleneck link with transmission rate of R bits/sec, each should have average rate of R/K

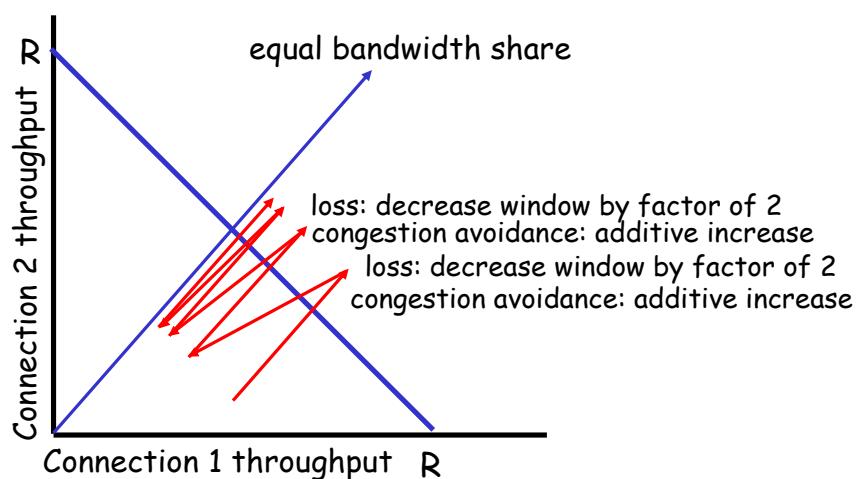


ELEC 331 101

Why is TCP fair?

Two competing sessions sharing a single link with rate R :

- ❑ Two connections have the same MSS and RTT.
- ❑ Ignore slow start. Connections operate in congestion avoidance.
- ❑ “Additive increase” gives slope of 1, as throughput increases.
- ❑ “Multiplicative decrease” decreases throughput proportionally.



ELEC 331 102

Fairness (more)

Fairness and UDP

- ❑ Multimedia apps often do not use TCP
 - ❑ do not want rate throttled by congestion control
- ❑ Instead use UDP:
 - ❑ pump audio/video at constant rate, tolerate packet loss
- ❑ Research area: TCP friendly

Fairness and parallel TCP connections

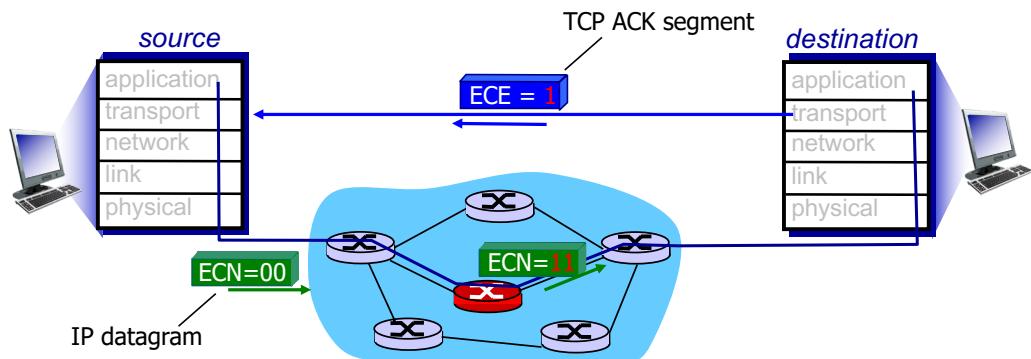
- ❑ nothing prevents app from opening parallel connections between 2 hosts.
- ❑ Web browsers do this
- ❑ Example: link of rate R supporting 9 connections;
 - ❑ new app asks for 1 TCP, gets rate $R/10$
 - ❑ new app asks for 11 TCPs, gets $R/2$

ELEC 331 103

Explicit Congestion Notification (ECN)

network-assisted congestion control:

- ❑ two bits in IP header (ToS field) marked by network router to indicate congestion
- ❑ congestion indication carried to receiving host
- ❑ receiver (seeing congestion indication in IP datagram) sets ECE (Explicit Congestion Notification Echo) bit on receiver-to-sender ACK segment to notify sender of congestion



ELEC 331 104

Chapter 3: Summary

- ❑ principles behind transport layer services:
 - multiplexing, demultiplexing
 - reliable data transfer
 - flow control
 - congestion control
- ❑ instantiation and implementation in the Internet
 - UDP
 - TCP

Next:

- ❑ leaving the network “edge” (application, transport layers)
- ❑ into the network “core”
- ❑ two network layer chapters:
 - data plane
 - control plane