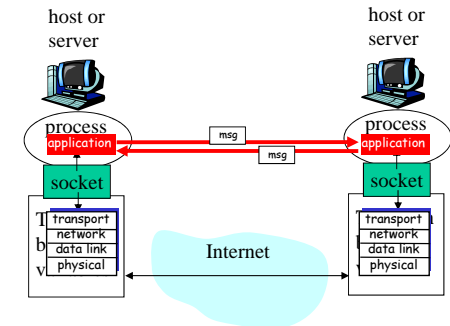


## Chapter 3 Transport Layer

Transport Layer 3-1

## Transport Layer

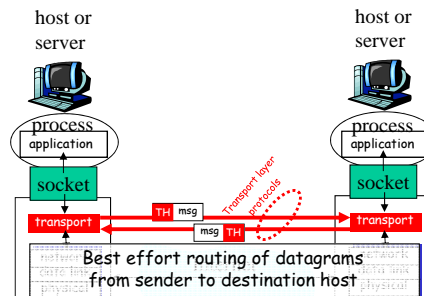
- Central piece of the layered network architecture
  - Provides communications services to applications
  - ...relying on the service of the network layer



Transport Layer 3-2

## Transport Layer

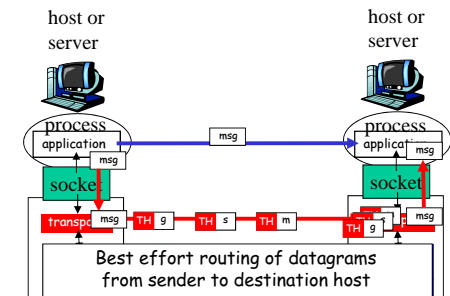
- Central piece of the layered network architecture
  - Provides communications services to applications
  - ...relying on the service of the network layer



Transport Layer 3-3

## Transport Layer

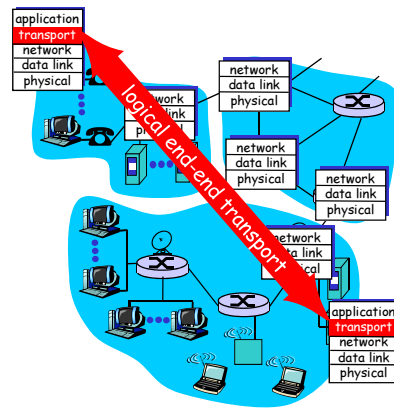
- Transport entities run in end systems
  - send side: breaks app messages into **segments**, passes to network layer
  - rcv side: reassembles segments into messages, passes to app layer
- Extends host to host datagram delivery to app to app message transfer



Transport Layer 3-4

## 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



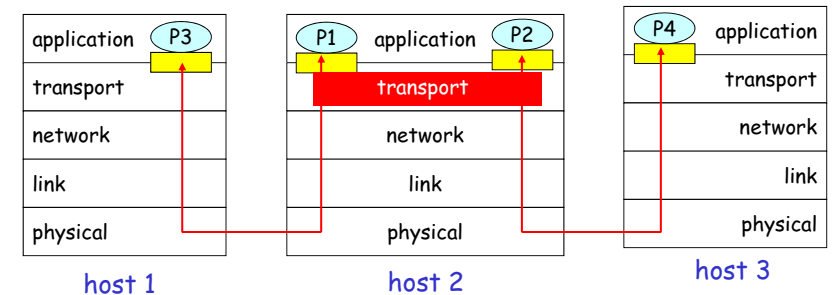
Transport Layer 3-5

## Multiplexing/demultiplexing

**Demultiplexing at rcv host:**  
delivering received segments to correct socket

**Multiplexing at send host:**  
gathering data from multiple sockets, enveloping data with header (later used for demultiplexing)

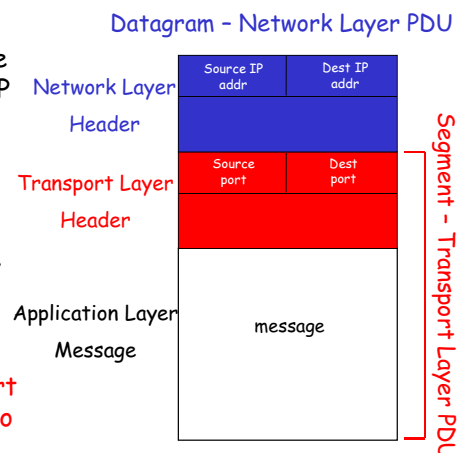
□ = socket ○ = process



Transport Layer 3-6

## How demultiplexing works

- host receives IP datagrams
  - each datagram has source IP address, destination IP address (in the network layer header)
  - each datagram carries 1 transport-layer segment
  - each segment has source, destination port number (in the transport layer header)
- host uses IP addresses & port numbers to direct segment to appropriate socket



Transport Layer 3-7

## Connectionless demultiplexing

- Create sockets with port numbers:

```
DatagramSocket mySocket1 = new
    DatagramSocket(9911);
DatagramSocket mySocket2 = new
    DatagramSocket(9922);
```

- UDP socket identified by two-tuple:

(IP address, port number)

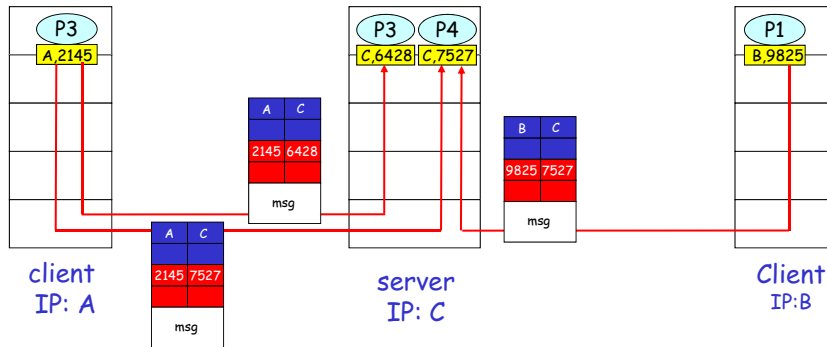
- When host receives UDP segment:

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

- IP datagrams with different source IP addresses and/or source port # but same destination port # are directed to same socket

Transport Layer 3-8

## Connectionless demux (cont)



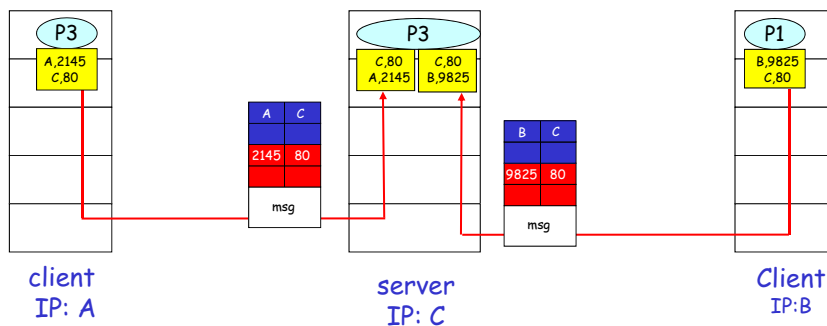
Transport Layer 3-9

## 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
- ❑ recv host uses all four values to direct segment to appropriate socket

Transport Layer 3-10

## Connection-oriented demux (cont)



Transport Layer 3-11

## 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 bare IP service, UDP adds
  - Mux/demux
  - checksum
- ❑ **connectionless**:
  - no handshaking between UDP sender, receiver
  - each UDP segment handled independently of others

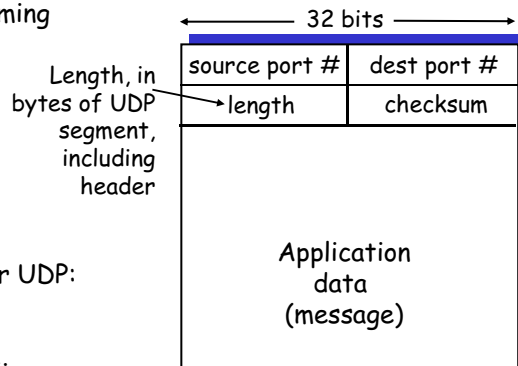
### Why is there a UDP?

- ❑ no connection establishment (which can add delay)
- ❑ simple: no connection state at sender, receiver
- ❑ small segment header
- ❑ no congestion control: UDP can blast away as fast as desired

Transport Layer 3-12

## UDP: more

- often used for streaming multimedia apps
  - loss tolerant
  - rate sensitive
- other UDP uses
  - DNS
  - SNMP
- reliable transfer over UDP: add reliability at application layer
  - application-specific error recovery!



UDP segment format

## UDP checksum

**Goal:** detect "errors" (e.g., flipped bits) in transmitted segment

### Sender:

- treat segment contents as sequence of 16-bit integers
- checksum: addition (1's complement 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*
- ....

## Internet checksum: example

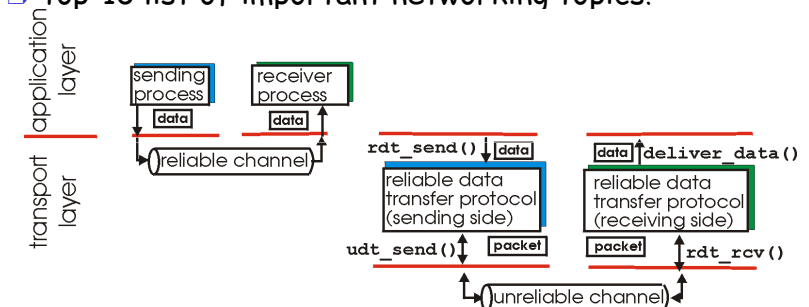
example: add two 16-bit integers

	1	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0
	1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
wraparound	1	1	0	1	1	1	0	1	1	1	0	1	1	1	0	1
sum	1	0	1	1	1	0	1	1	1	0	1	1	1	1	0	0
checksum	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1	1

**Note:** when adding numbers, a carryout from the most significant bit needs to be added to the result

## Principles of Reliable data transfer

- important in app., transport, link layers
- top-10 list of important networking topics!

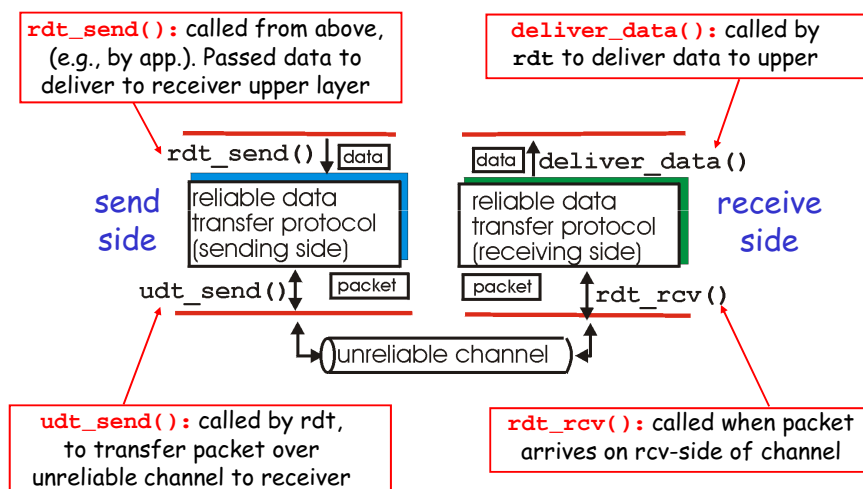


(a) provided service

(b) service implementation

- characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)

## Reliable data transfer: getting started

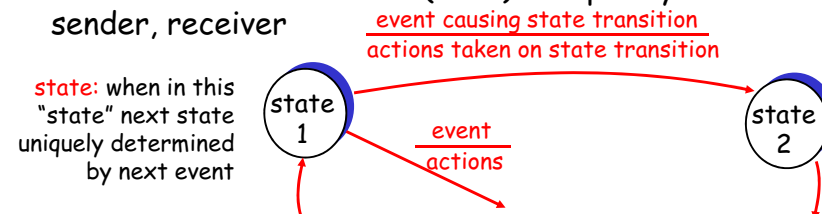


Transport Layer 3-17

## Reliable data transfer: getting started

We'll:

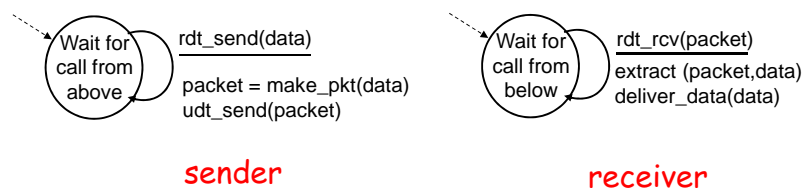
- incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)
- consider only unidirectional data transfer
  - but control info will flow on both directions!
- use finite state machines (FSM) to specify sender, receiver



Transport Layer 3-18

## 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 read data from underlying channel



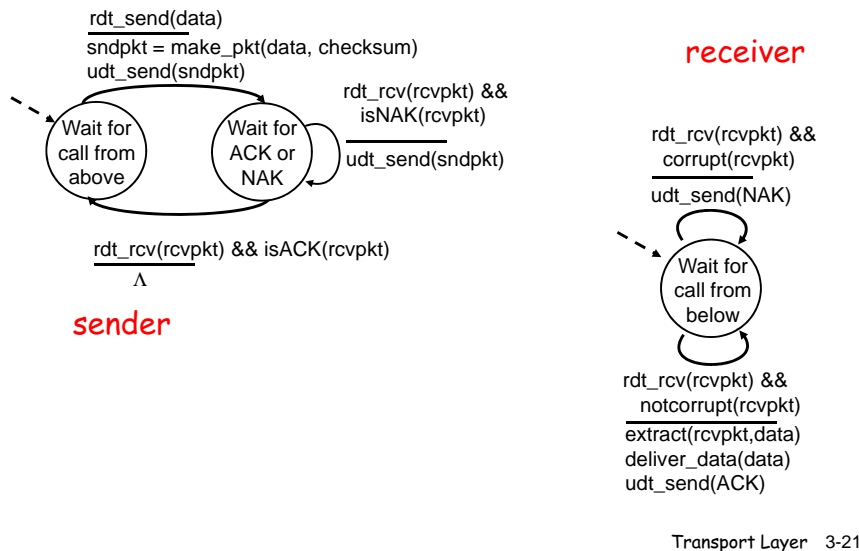
Transport Layer 3-19

## Rdt2.0: channel with bit errors

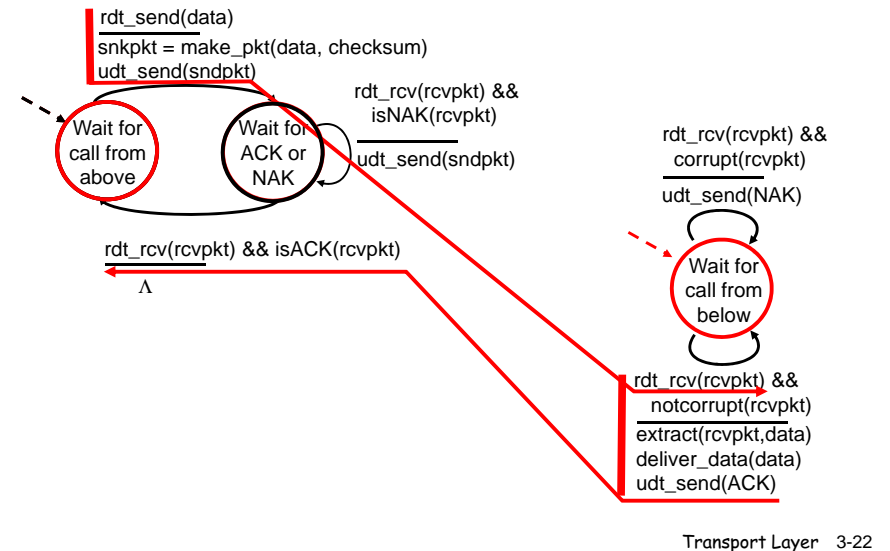
- underlying channel may flip bits in packet
  - recall: UDP checksum to detect bit errors
- the question: how to recover from errors:
  - acknowledgements (ACKs)**: receiver explicitly tells sender that pkt received OK
  - negative acknowledgements (NAKs)**: receiver explicitly tells sender that pkt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

Transport Layer 3-20

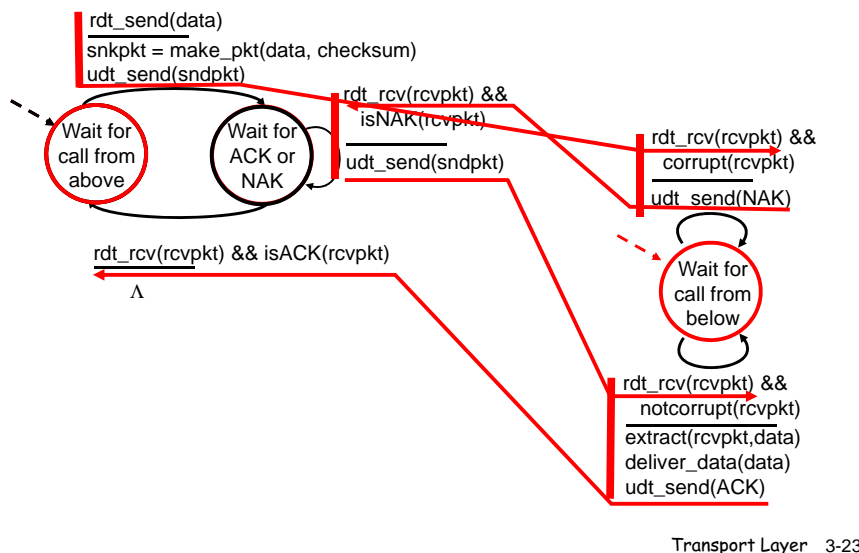
## rdt2.0: FSM specification



## rdt2.0: operation with no errors



## rdt2.0: error scenario



## rdt2.0 has a fatal flaw!

What happens if  
**ACK/NAK corrupted?**

- sender doesn't know what happened at receiver!

What to do?

- sender ACKs/NAKs receiver's ACK/NAK? What if sender ACK/NAK lost?
- retransmit, but this might cause retransmission of correctly received pkt!

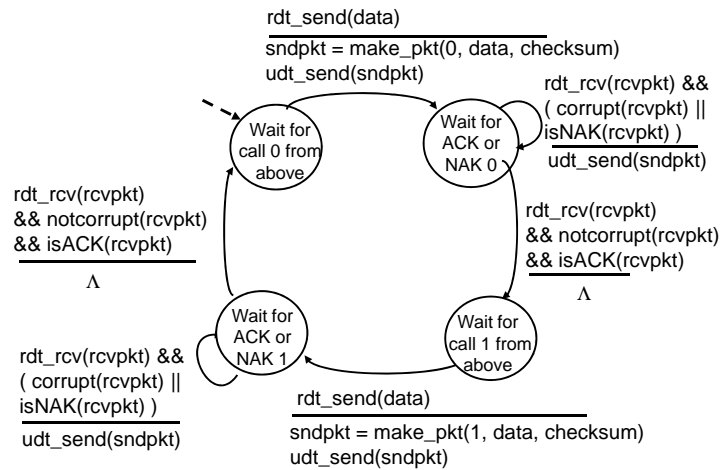
Handling duplicates:

- 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**

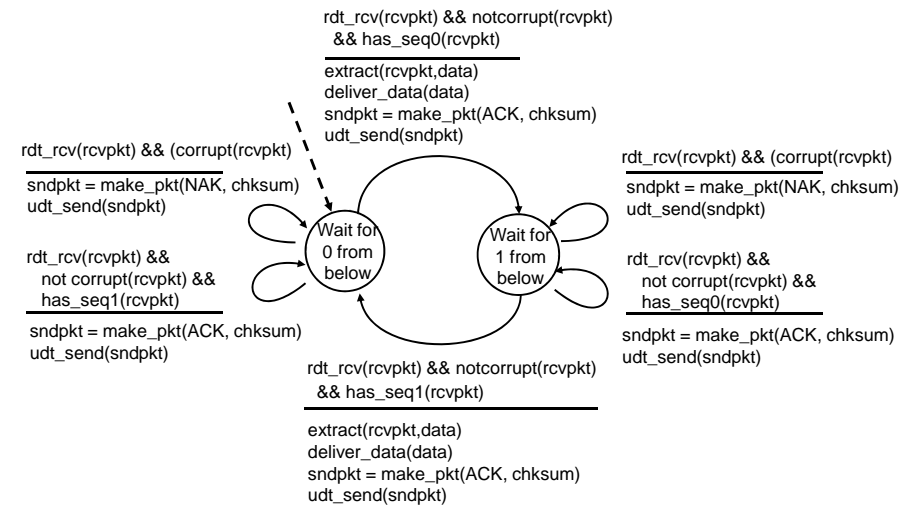
Sender sends one packet, then waits for receiver response

## rdt2.1: sender, handles garbled ACK/NAKs



Transport Layer 3-25

## rdt2.1: receiver, handles garbled ACK/NAKs



Transport Layer 3-26

## rdt2.1: discussion

### Sender:

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

### 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

Transport Layer 3-27

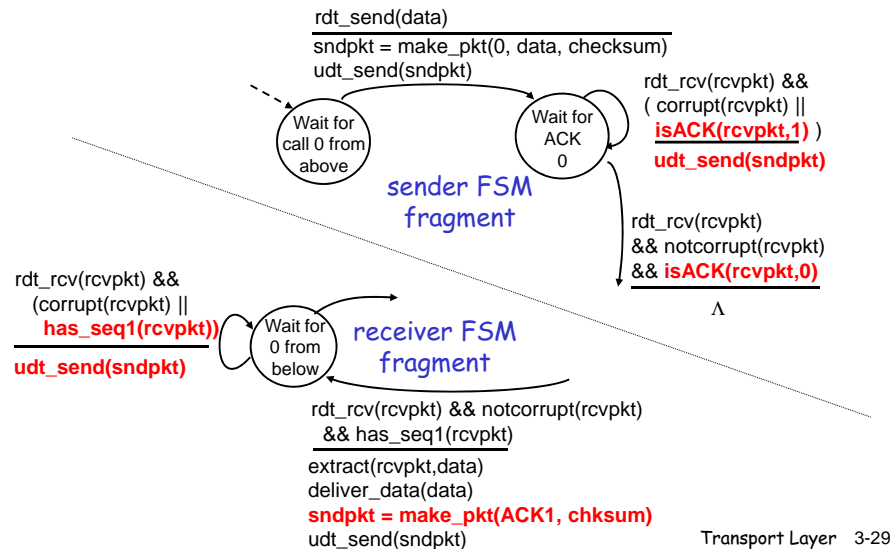
## rdt2.2: a NAK-free protocol

- ❑ 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*

Transport Layer 3-28



## rdt2.2: sender, receiver fragments



## rdt3.0: channels with errors and loss

### New assumption:

underlying channel can also lose packets (data or ACKs)

- checksum, seq. #, ACKs, retransmissions will be of help, but not enough

### Q: how to deal with loss?

- sender waits until certain data or ACK lost, then retransmits
- yuck: drawbacks?

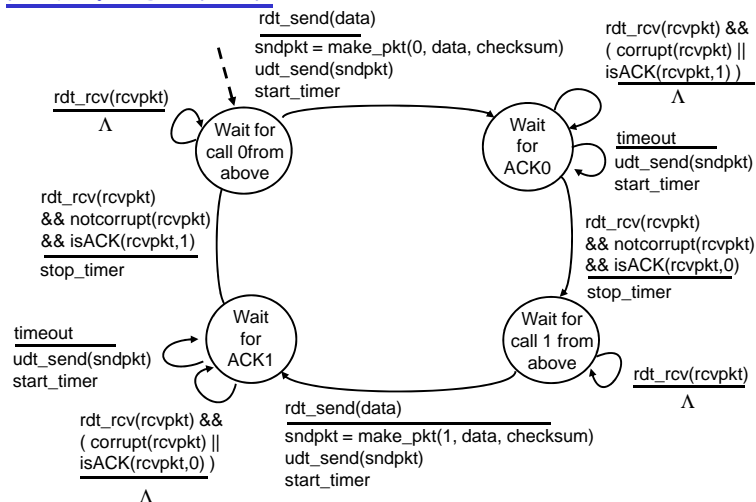
### Approach:

sender waits "reasonable" amount of time for ACK

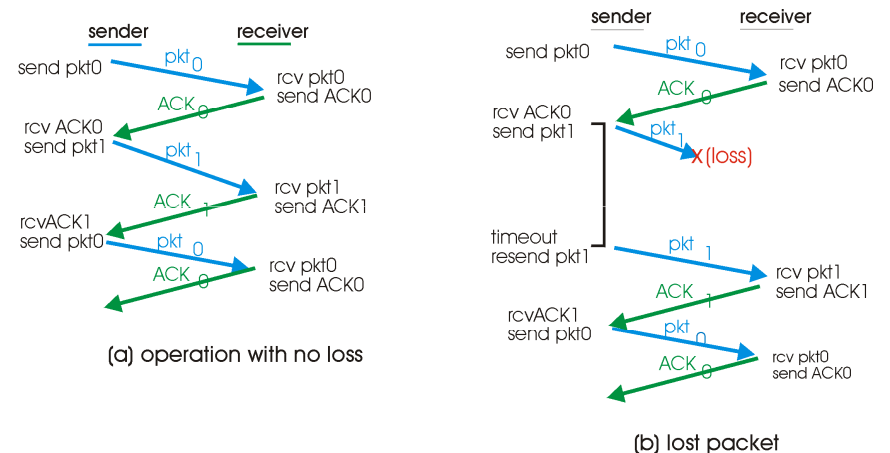
- retransmits if no ACK received in this time
- if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate, but use of seq. #'s already handles this
  - receiver must specify seq # of pkt being ACKed
- requires countdown timer

Transport Layer 3-30

## rdt3.0 sender

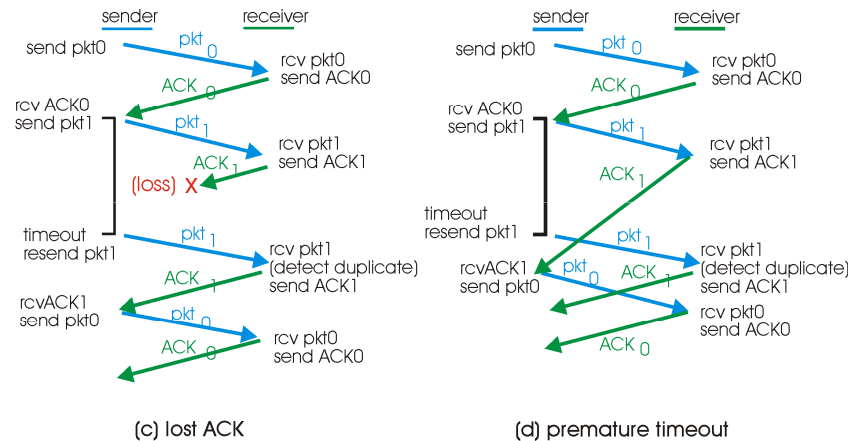


## rdt3.0 in action





## rdt3.0 in action



## Performance of rdt3.0

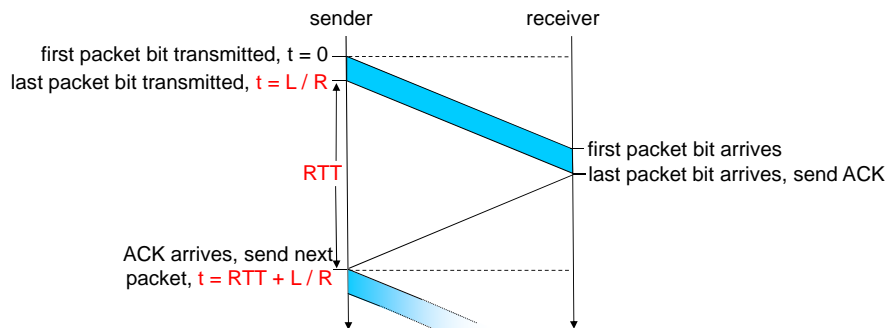
- rdt3.0 works, but performance stinks
- example: 1 Gbps link, 15 ms e-e prop. delay, 1KB packet:

$$T_{\text{transmit}} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8\text{kb/pkt}}{10^9 \text{ b/sec}} = 8 \text{ microsec}$$

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

- $U_{\text{sender}}$ : **utilization** - fraction of time sender busy sending
- 1KB pkt every 30 msec  $\rightarrow$  267kbps thrupt over 1 Gbps link
- network protocol limits use of physical resources!

## rdt3.0: stop-and-wait operation

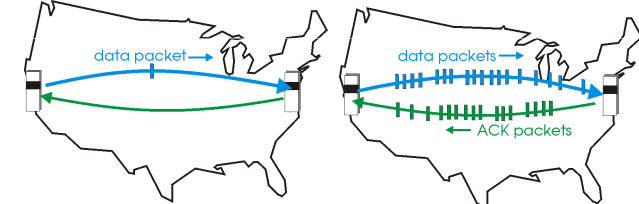


$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

## 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

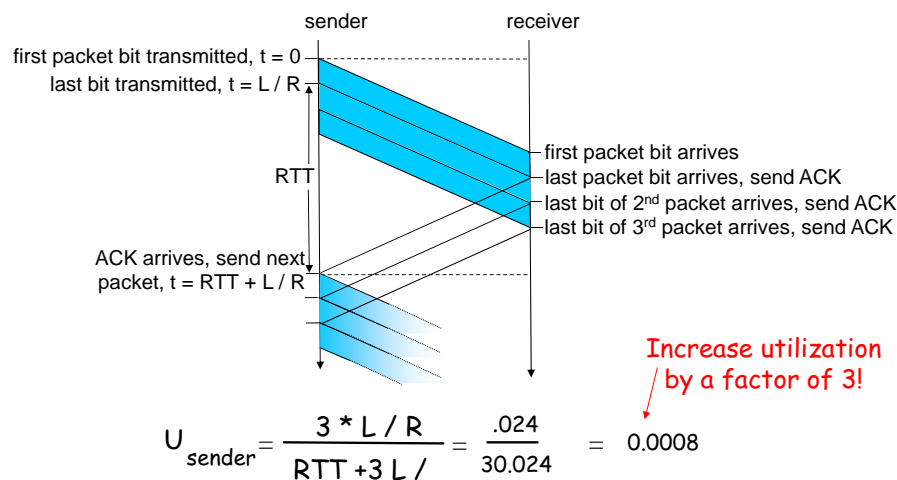


(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

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

## Pipelining: increased utilization



## Pipelining Protocols

### Go-back-N: big picture:

- ❑ Sender can have up to N unacked packets in pipeline
- ❑ Rcvr only sends cumulative acks
  - Doesn't ack packet if there's a gap
- ❑ Sender has timer for oldest unacked packet
  - If timer expires, retransmit all unacked packets

### Selective Repeat: big pic

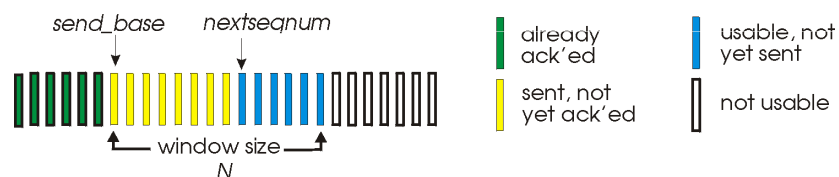
- ❑ Sender can have up to N unacked packets in pipeline
- ❑ Rcvr acks individual packets
- ❑ Sender maintains timer for each unacked packet
  - When timer expires, retransmit only unack packet

## Go-Back-N

- ❑ Transmit multiple packets (up to N) without waiting for ACK

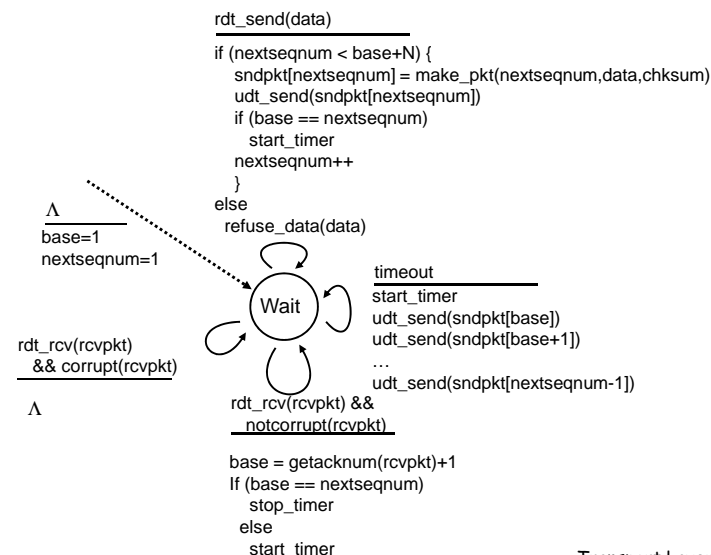
### Sender:

- ❑ k-bit seq # in pkt header
- ❑ "window" of up to N, consecutive unack'd pkts allowed

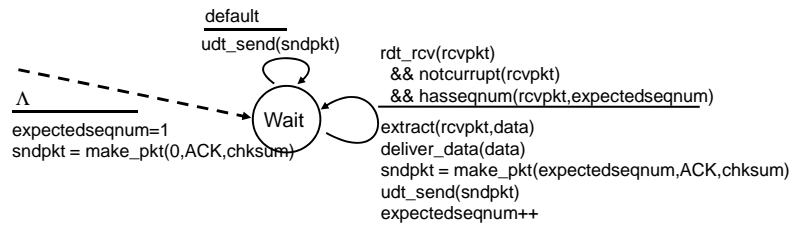


- ❑ ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  - may receive duplicate ACKs (see receiver)
- ❑ timeout(n): retransmit pkt n and all higher seq # pkts in window

## GBN: sender extended FSM (1 timer)



## GBN: receiver extended FSM

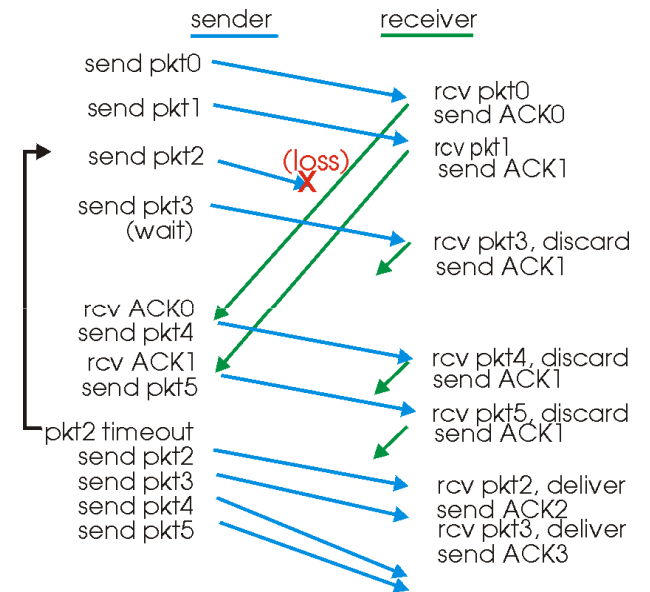


ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #

- 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 #

Transport Layer 3-41

## GBN in action



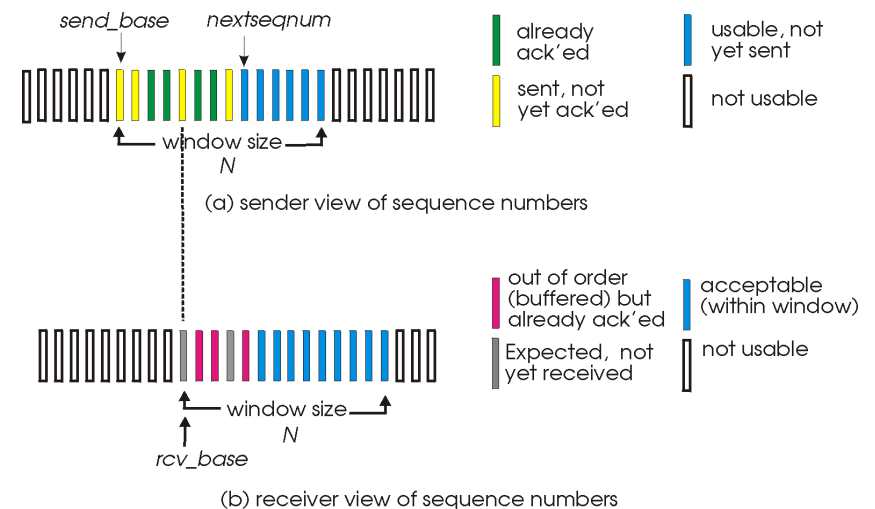
Transport Layer 3-42

## 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
  - again limits seq #'s of sent, unACKed pkts

Transport Layer 3-43

## Selective repeat: sender, receiver windows



Transport Layer 3-44

## Selective repeat

### sender

data from above :

- if next available seq # in window, send pkt

timeout(n):

- resend pkt n, restart timer

ACK(n) in [sendbase, sendbase+N]:

- mark pkt n as received
- if n smallest unACKed pkt, advance window base to next unACKed seq #

### receiver

pkt n in [rcvbase, rcvbase+N-1]

- send ACK(n)
- out-of-order: buffer
- in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

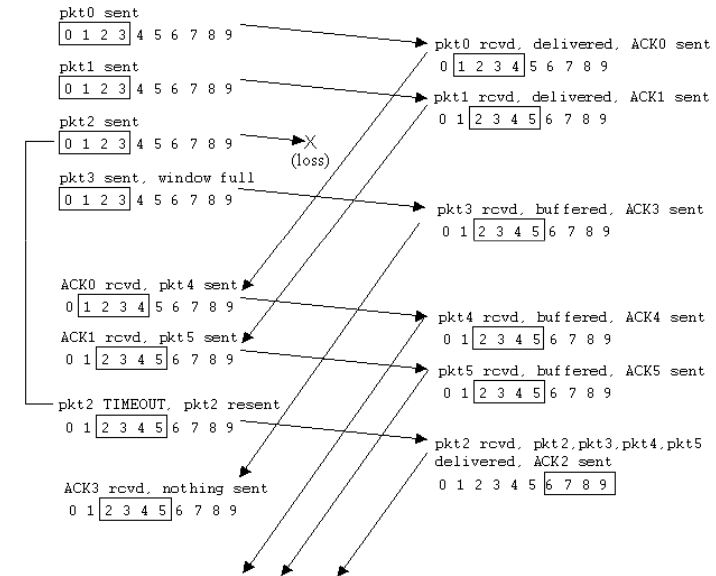
pkt n in [rcvbase-N, rcvbase-1]

- ACK(n)
- otherwise:

- ignore

Transport Layer 3-45

## Selective repeat in action



† Layer 3-46

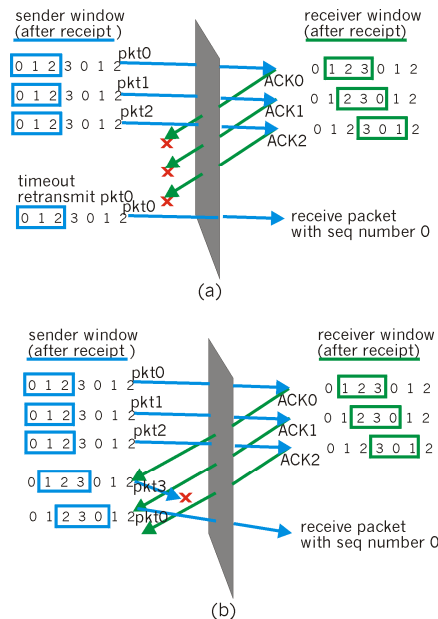
## Selective repeat: dilemma

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?

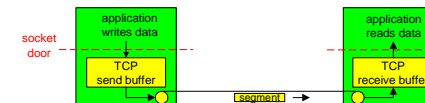


Transport Layer 3-47

## TCP: Overview

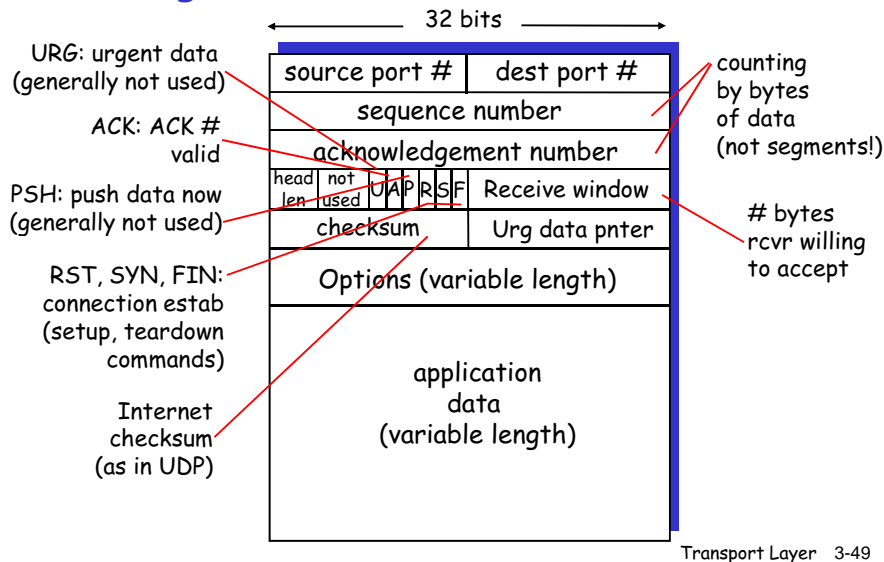
RFCs: 793, 1122, 1323, 2018, 2581

- 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) init's sender, receiver state before data exchange
- flow controlled:
  - sender will not overwhelm receiver



Transport Layer 3-48

## TCP segment structure



## Sequence and Acknowledgement Numbers

- TCP views data as unstructured, but ordered, data
- In a segment:
  - **Sequence number**: is the byte-stream number of the first byte in the segment
    - Initial sequence number is randomly chosen
  - **Ack number**: is the number of the next byte expected from the other side
    - TCP uses cumulative acknowledgements
- Q: how receiver handles out-of-order segments
  - A: TCP spec doesn't say - up to implementor

Transport Layer 3-50

## TCP seq. numbers, ACKs

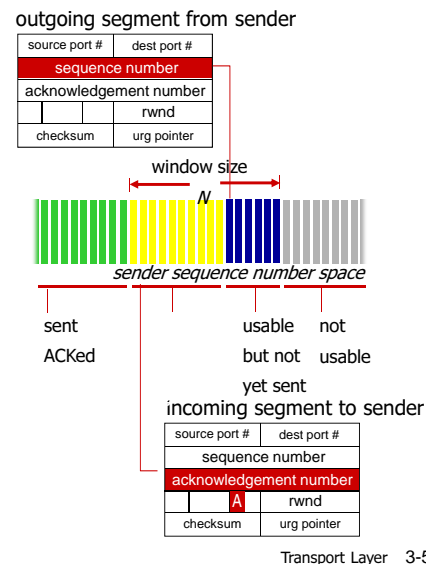
### sequence numbers:

- byte stream "number" of first byte in segment's data

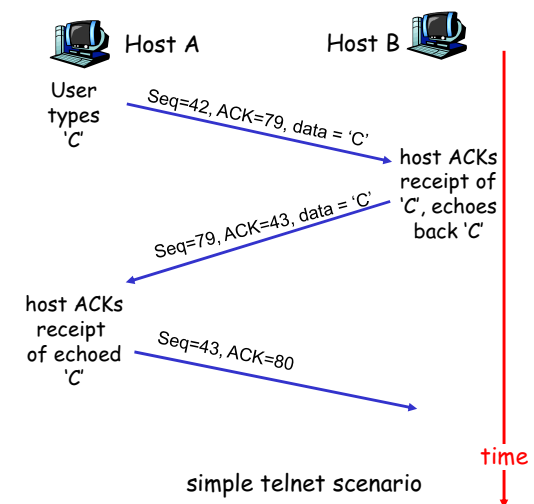
### acknowledgements:

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

- Q: how receiver handles out-of-order segments
- A: TCP spec doesn't say, - up to implementor



## TCP seq. #'s and ACKs



Transport Layer 3-52

## TCP reliable data transfer

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

Transport Layer 3-53

## TCP sender events:

### data rcvd from app:

- ❑ Create segment with seq #
- ❑ seq # is byte-stream number of first data byte in segment
- ❑ start timer for that 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

Transport Layer 3-54

```
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)
    break
```

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

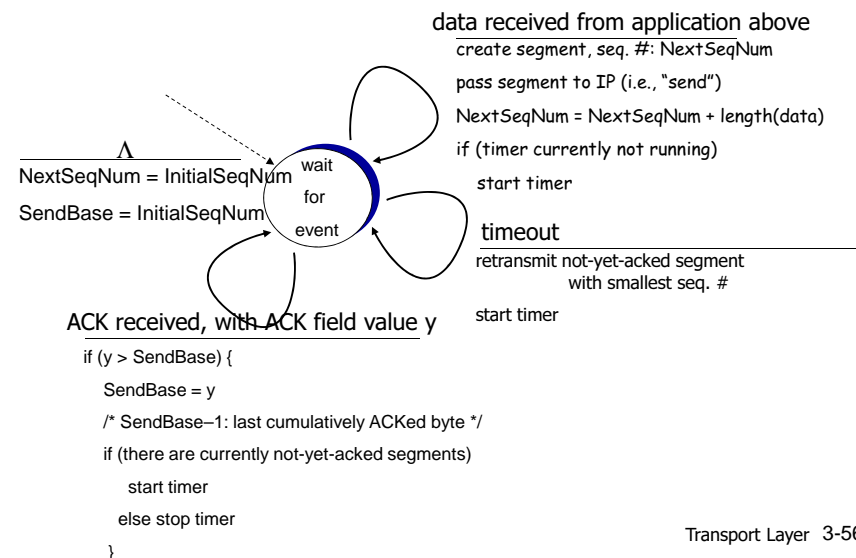
```
  event: ACK received, with ACK field value of y
    if (y > SendBase) {
      SendBase=y;
      if( there are currently any not yet acked segment)
        start timer
    }
    break
```

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

## TCP sender (simplified)

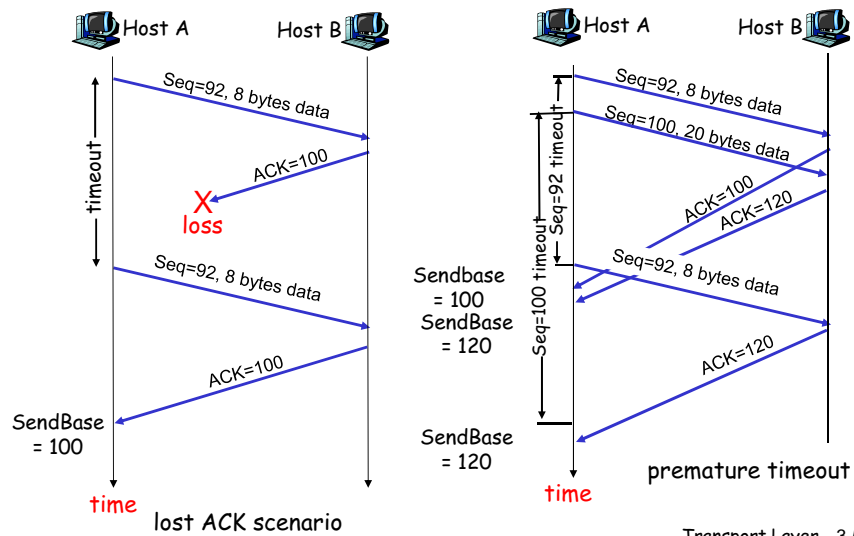
Transport Layer 3-55

## TCP sender (simplified)



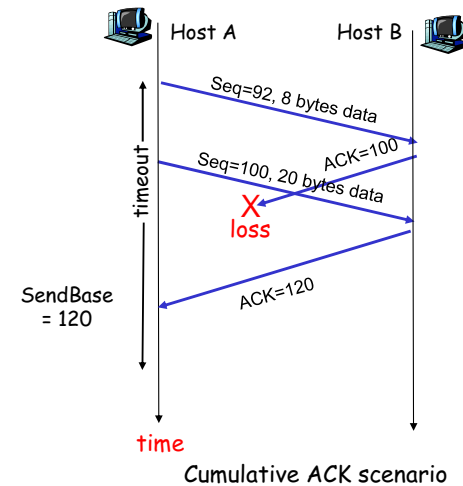
Transport Layer 3-56

## TCP: retransmission scenarios



Transport Layer 3-57

## TCP retransmission scenarios (more)



Transport Layer 3-58

## 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

Transport Layer 3-59

## Fast Retransmit

- Time-out 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.
- If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
  - fast retransmit: resend segment before timer expires

Transport Layer 3-60



### Fast retransmit algorithm:

```

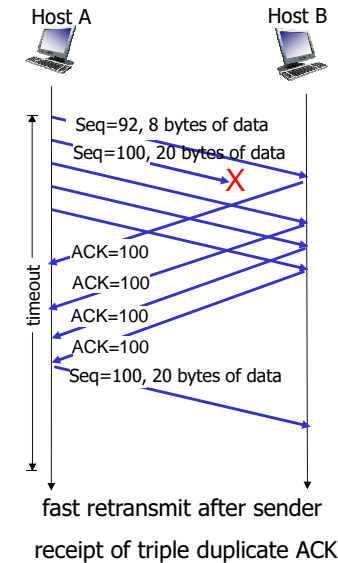
event: ACK received, with ACK field value of y
    if (y > SendBase) {
        Sendbase=y;
        if ( there are currently any not yet acked segment)
            start timer
    }
    else {
        increment count of dup ACKs received for y
        if (count of dup ACKs received for y = 3) {
            resend segment with sequence number y
        }
        break
    }

```

a duplicate ACK for  
already ACKed segment

fast retransmit

## TCP fast retransmit



## 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?

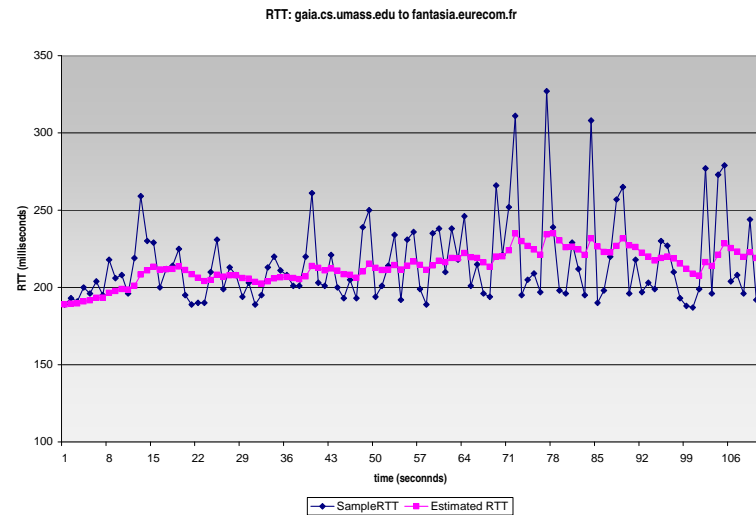
- ❑ **SampleRTT**: measured time from segment transmission until ACK receipt
  - ignore retransmitted segments
- ❑ **SampleRTT** will vary, want estimated RTT “smoother”
  - average several recent measurements, not just current **SampleRTT**

## TCP Round Trip Time and Timeout

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

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

## Example RTT estimation:



Transport Layer 3-65

## 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:

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

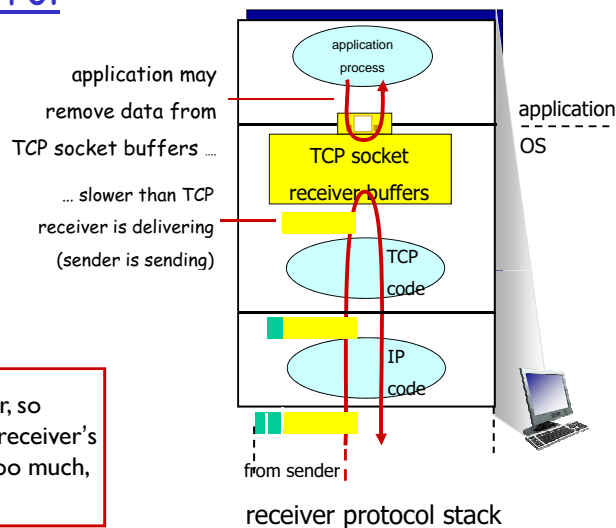
(typically,  $\beta = 0.25$ )

Then set timeout interval:

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

Transport Layer 3-66

## TCP flow control

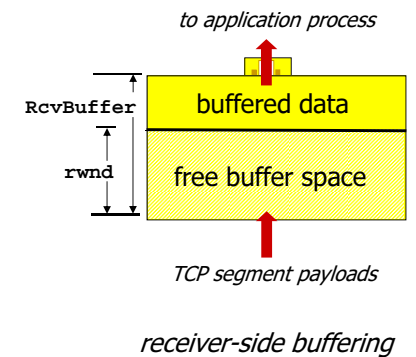


**flow control**  
receiver controls sender, so sender won't overflow receiver's buffer by transmitting too much, too fast

Transport Layer 3-67

## 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
- guaranteed receive



Transport Layer 3-68

## TCP Connection Management

**Recall:** TCP sender, receiver establish "connection" before exchanging data segments

- initialize TCP variables:

- seq. #s
- buffers, flow control info (e.g. RcvWindow)

- *client*: connection initiator

```
Socket clientSocket = new Socket("hostname", "port number");
```

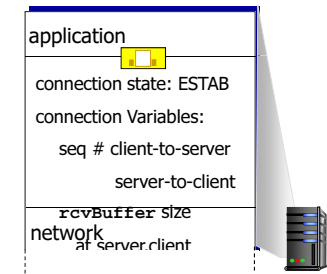
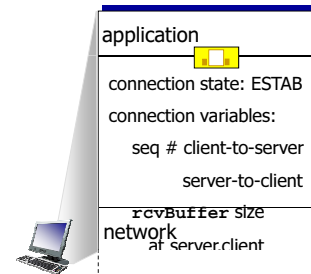
- *server*: contacted by client

```
Socket connectionSocket = welcomeSocket.accept();
```

## Connection Management

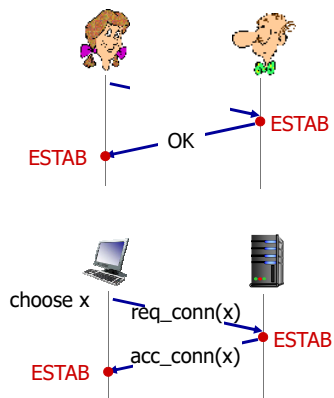
before exchanging data, sender/receiver "handshake":

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



## Agreeing to establish a connection

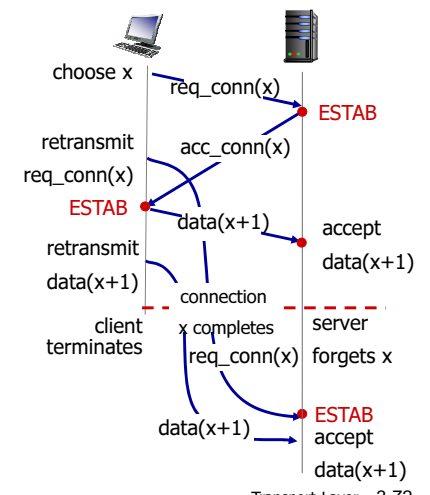
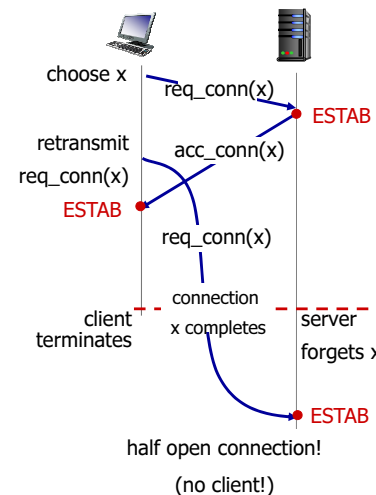
2-way handshake:



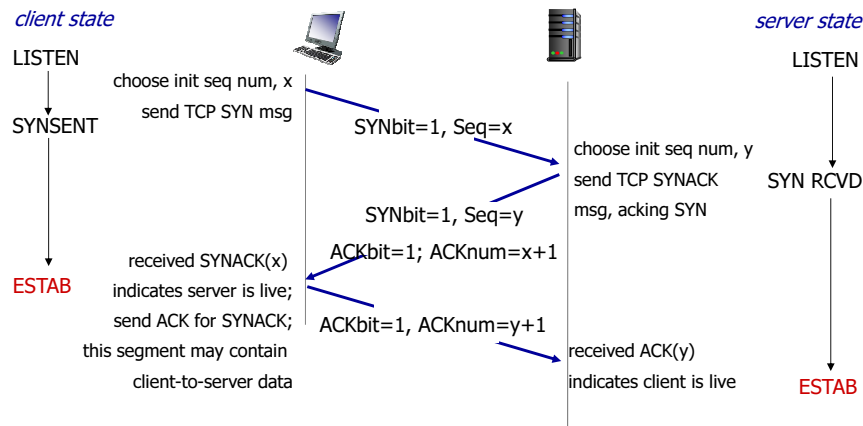
- Q:** will 2-way handshake always work in network?
- variable delays
  - retransmitted messages (e.g. req\_conn(x)) due to message loss
  - message reordering
  - can't "see" other side

## Agreeing to establish a connection

2-way handshake failure scenarios:

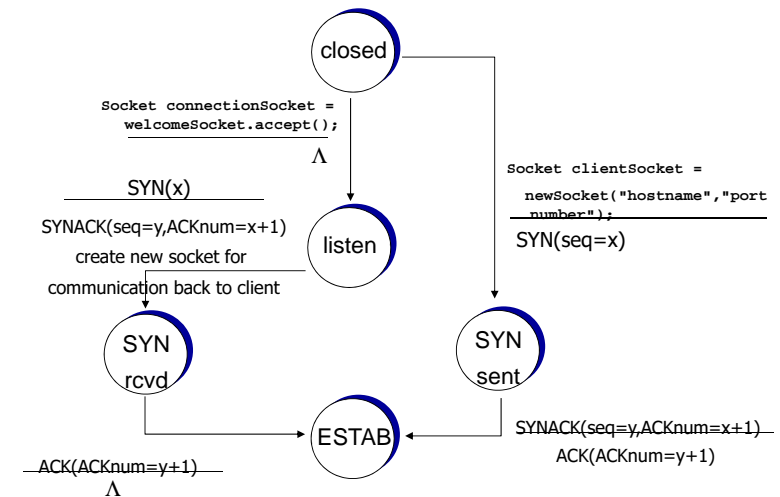


## TCP 3-way handshake



Transport Layer 3-73

## TCP 3-way handshake: FSM



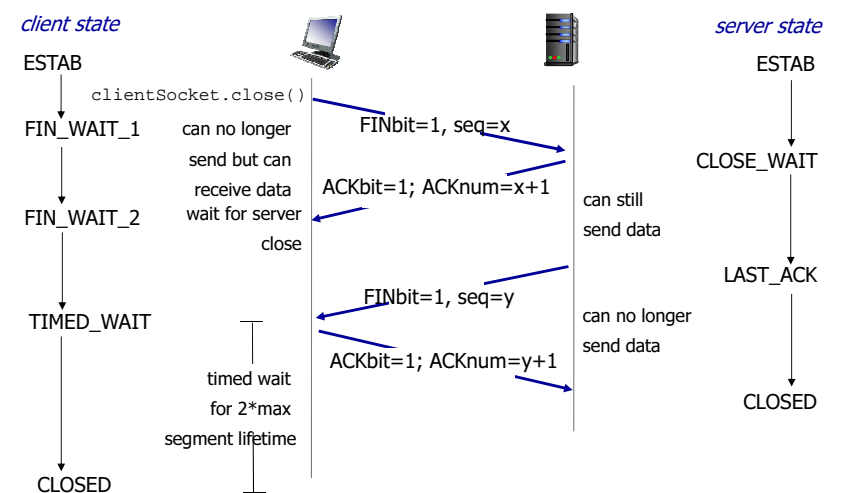
Transport Layer 3-74

## 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

Transport Layer 3-75

## TCP: closing a connection



Transport Layer 3-76

## Principles of Congestion Control

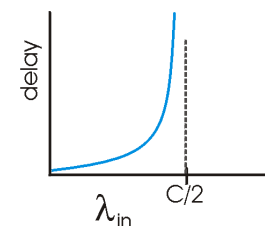
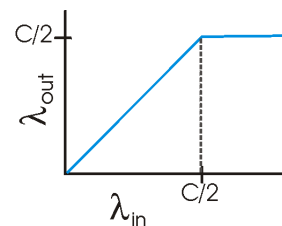
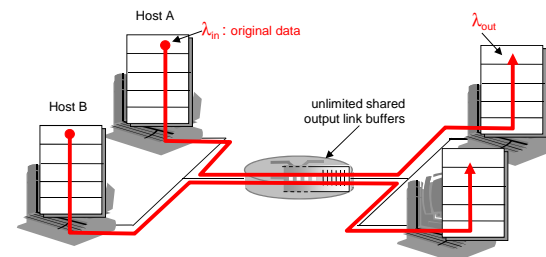
### Congestion:

- 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!

Transport Layer 3-77

## Causes/costs of congestion: scenario 1

- two senders, two receivers
- one router, infinite buffers
- no retransmission

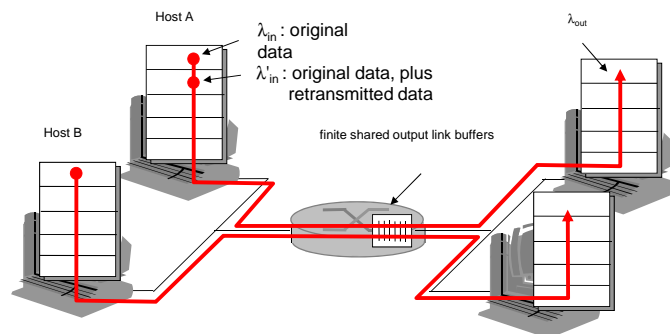


- large delays when congested
- maximum achievable throughput

Transport Layer 3-78

## Causes/costs of congestion: scenario 2

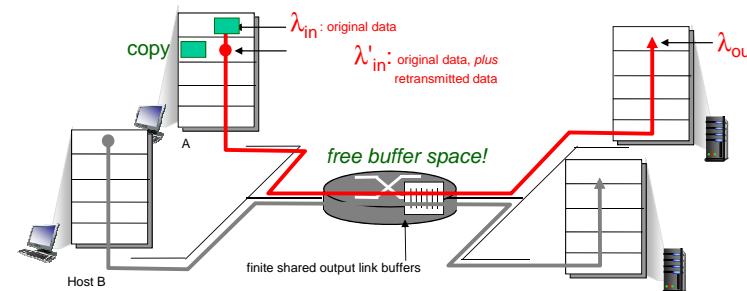
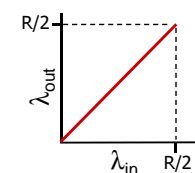
- one router, *finite* buffers
- sender retransmission of lost packet



Transport Layer 3-79

## Causes/costs of congestion: scenario 2

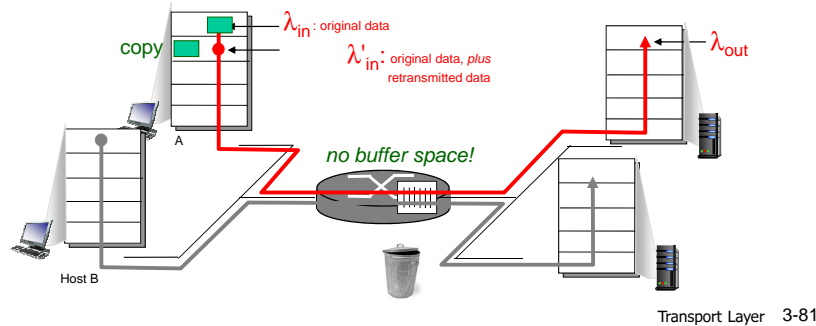
- idealization: perfect knowledge
- sender sends only when router buffers available



Transport Layer 3-80

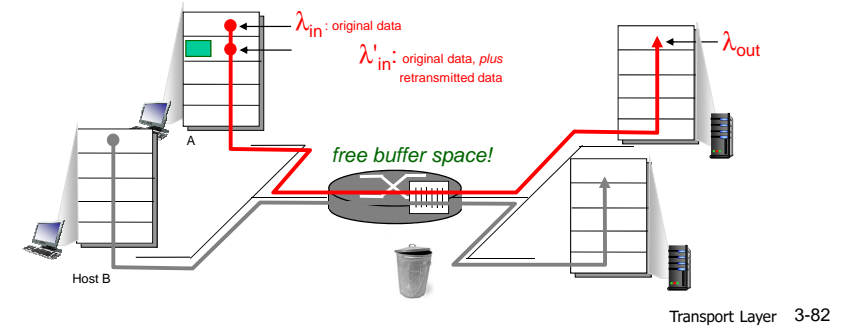
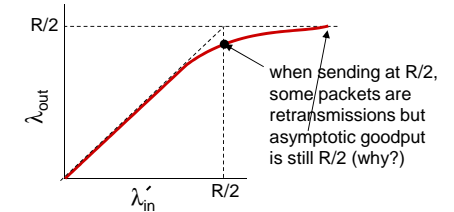
## Causes/costs of congestion: scenario 2

- **Idealization: known loss** packets can be lost, dropped at router due to full buffers
- sender only resends if packet *known* to be lost



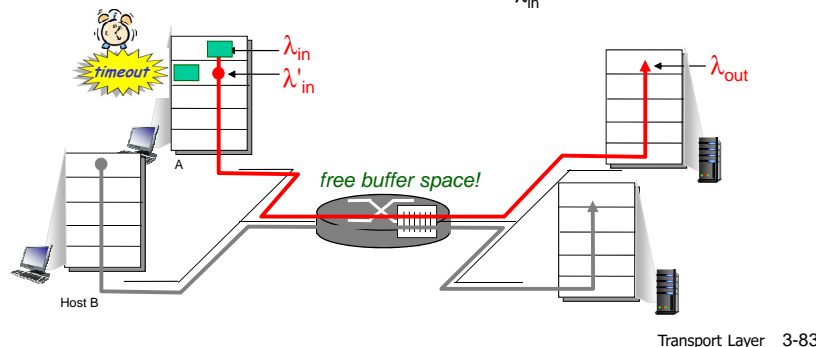
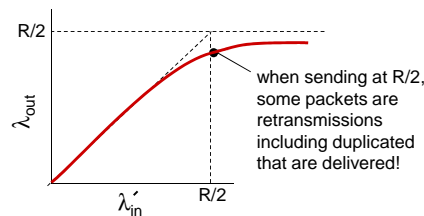
## Causes/costs of congestion: scenario 2

- **Idealization: known loss** packets can be lost, dropped at router due to full buffers
- sender only resends if packet *known* to be lost



## Causes/costs of congestion: scenario 2

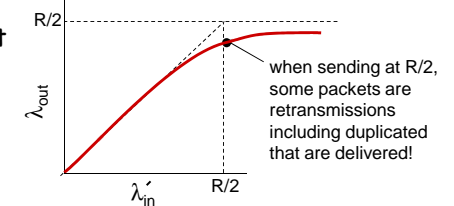
- **Realistic: duplicates**
- packets can be lost, dropped at router due to full buffers
- sender times out prematurely, sending two copies, both of which are delivered



## Causes/costs of congestion: scenario 2

### Realistic: duplicates

- packets can be lost, dropped at router due to full buffers
- sender times out prematurely, sending *two* copies, both of which are delivered

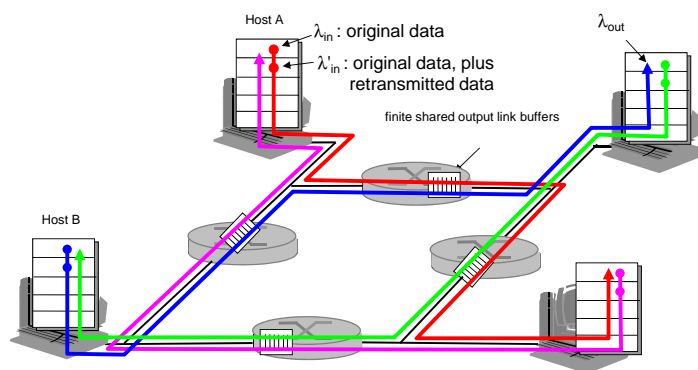


- “costs” of congestion:
- more work (retrans) for given “goodput”
- unneeded retransmissions: link carries multiple copies of pkt
  - decreasing goodput

## Causes/costs of congestion: scenario 3

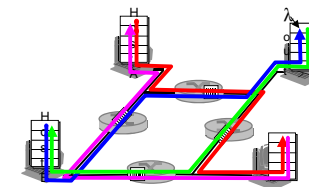
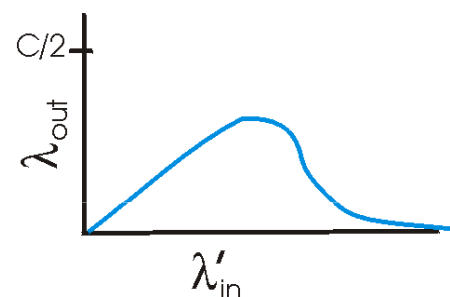
- four senders
- multihop paths
- timeout/retransmit

Q: what happens as  $\lambda_{in}$  and  $\lambda'_{in}$  increase?



Transport Layer 3-85

## Causes/costs of congestion: scenario 3



Another "cost" of congestion:

- when packet dropped, any "upstream transmission capacity used for that packet was wasted!"

Transport Layer 3-86

## Approaches towards congestion control

Two broad approaches towards congestion control:

### End-end congestion control:

- no explicit feedback from network
- congestion inferred from end-system observed loss, delay
- approach taken by TCP

### Network-assisted congestion control:

- routers provide feedback to end systems
  - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
  - explicit rate sender should send at

Transport Layer 3-87

## Case study: ATM ABR congestion control

### ABR: available bit rate:

- "elastic service"
- if sender's path "underloaded":
  - sender should use available bandwidth
- if sender's path congested:
  - sender throttled to minimum guaranteed rate

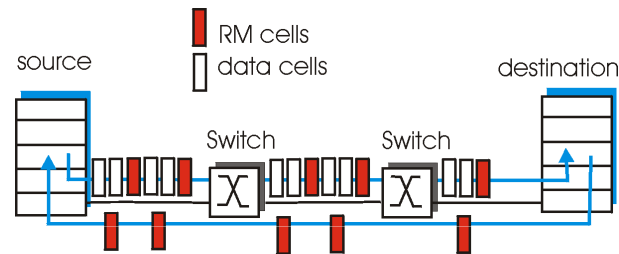
### RM (resource management) cells:

- sent by sender, interspersed with data cells
- bits in RM cell set by switches ("network-assisted")
  - NI bit: no increase in rate (mild congestion)
  - CI bit: congestion indication
- RM cells returned to sender by receiver, with bits intact

Transport Layer 3-88



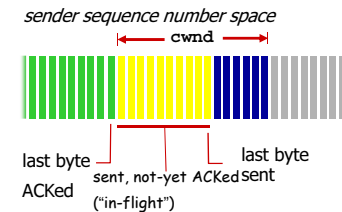
## Case study: ATM ABR congestion control



- two-byte ER (explicit rate) field in RM cell
  - congested switch may lower ER value in cell
  - sender's send rate thus maximum supportable rate on path
- EFCI bit in data cells: set to 1 in congested switch
  - if data cell preceding RM cell has EFCI set, sender sets CI bit in returned RM cell

Transport Layer 3-89

## TCP Congestion Control



- end-end control (no network assistance)
- sender limits transmission:
 
$$\text{LastByteSent} - \text{LastByteAcked} \leq \min\{\text{CongWin}, \text{RcvWindow}\}$$
- Roughly,

$$\text{rate} = \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$$

- CongWin is dynamic, function of perceived network congestion

### How does sender perceive congestion?

- loss event = timeout or 3 duplicate acks
- TCP sender reduces rate (CongWin) after loss event

### three mechanisms:

- slow start
- AIMD
- conservative after timeout events

Transport Layer 3-90

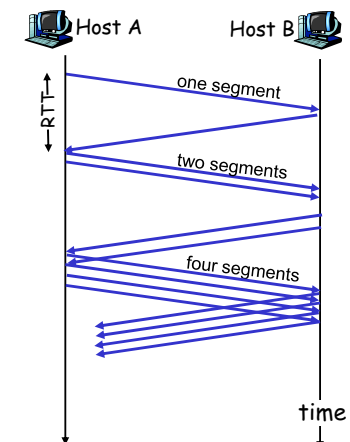
## TCP Slow Start

- When connection begins, CongWin = 1 MSS
  - Example: MSS = 500 bytes & RTT = 200 msec
  - initial rate = 20 kbps
- available bandwidth may be  $\gg \text{MSS}/\text{RTT}$ 
  - desirable to quickly ramp up to respectable rate
- When connection begins, increase rate exponentially fast until first loss event

Transport Layer 3-91

## TCP Slow Start (more)

- When connection begins, increase rate exponentially until first loss event:
  - double CongWin every RTT
  - done by incrementing CongWin for every ACK received
- Summary: initial rate is slow but ramps up exponentially fast



Transport Layer 3-92

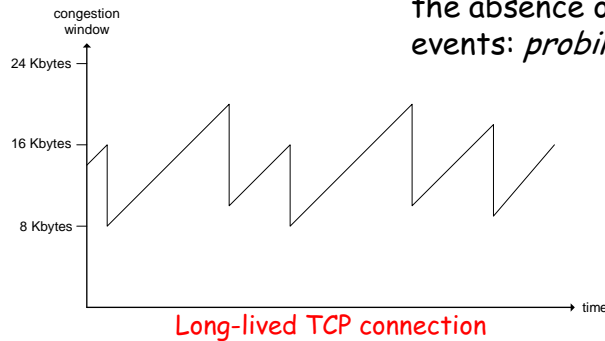
## TCP AIMD

### multiplicative decrease:

cut CongWin in half  
after loss event

### additive increase:

increase CongWin by  
1 MSS every RTT in  
the absence of loss  
events: *probing*



## Refinement: inferring loss

- After 3 dup ACKs:
  - CongWin is cut in half
  - window then grows linearly
- But after timeout event:
  - CongWin instead set to 1 MSS;
  - window then grows exponentially
  - to a threshold, then grows linearly

### Philosophy:

- 3 dup ACKs indicates network capable of delivering some segments
- timeout before 3 dup ACKs is "more alarming"

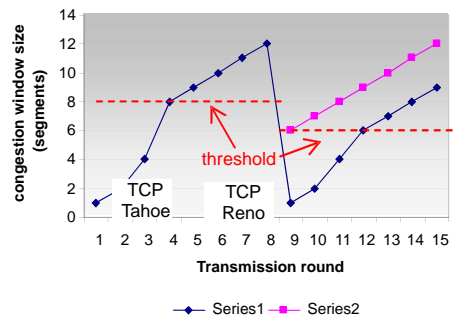
## Refinement (more)

**Q:** When should the exponential increase switch to linear?

**A:** When CongWin gets to 1/2 of its value before timeout.

### Implementation:

- Variable Threshold
- At loss event, Threshold is set to 1/2 of CongWin just before loss event



## Summary: TCP Congestion Control

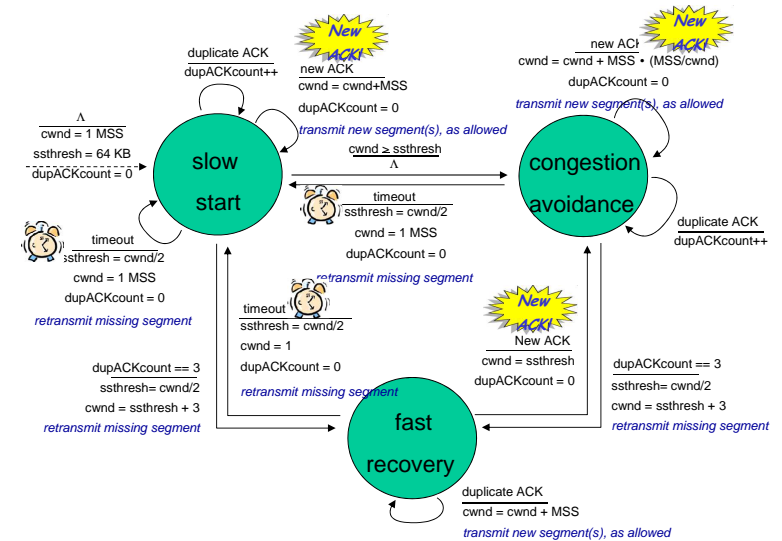
- When CongWin is below Threshold, sender in **slow-start** phase, window grows exponentially.
- When CongWin is above Threshold, sender is in **congestion-avoidance** phase, window grows linearly.
- When a **triple duplicate ACK** occurs, Threshold set to CongWin/2 and CongWin set to Threshold.
- When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

## TCP sender congestion control

State	Event	TCP Sender Action	Commentary
Slow Start (SS)	ACK receipt for previously unacked data	CongWin = CongWin + MSS, If (CongWin > Threshold) set state to "Congestion Avoidance"	Resulting in a doubling of CongWin every RTT
Congestion Avoidance (CA)	ACK receipt for previously unacked data	CongWin = CongWin + MSS * (MSS/CongWin)	Additive increase, resulting in increase of CongWin by 1 MSS every RTT
SS or CA	Loss event detected by triple duplicate ACK	Threshold = CongWin/2, CongWin = Threshold, Set state to "Congestion Avoidance"	Fast recovery, implementing multiplicative decrease. CongWin will not drop below 1 MSS.
SS or CA	Timeout	Threshold = CongWin/2, CongWin = 1 MSS, Set state to "Slow Start"	Enter slow start
SS or CA	Duplicate ACK	Increment duplicate ACK count for segment being acked	CongWin and Threshold not changed

Transport Layer 3-97

## Summary: TCP Congestion Control

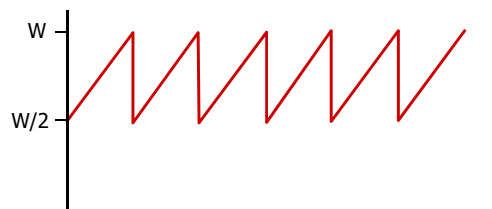


Transport Layer 3-98

## TCP throughput

- avg. TCP thruput as function of window size, RTT?
  - ignore slow start, assume always data to send
- **W: window size** (measured in bytes) where loss occurs
  - avg. window size (# in-flight bytes) is  $\frac{3}{4} W$
  - avg. thruput is  $\frac{3}{4} W$  per RTT

$$\text{avg TCP thruput} = \frac{3}{4} \frac{W}{RTT} \text{ bytes/sec}$$



Transport Layer 3-99

## TCP Futures: TCP over "long, fat pipes"

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

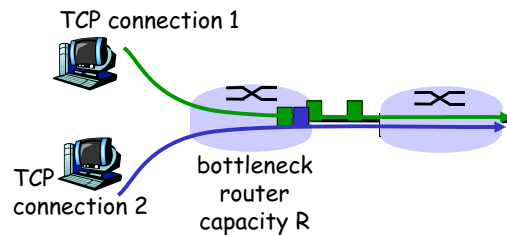
$$\frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$

- $\rightarrow L = 2 \cdot 10^{-10}$  **Wow**
- New versions of TCP for high-speed

Transport Layer 3-100

## TCP Fairness

**Fairness goal:** if  $K$  TCP sessions share same bottleneck link of bandwidth  $R$ , each should have average rate of  $R/K$

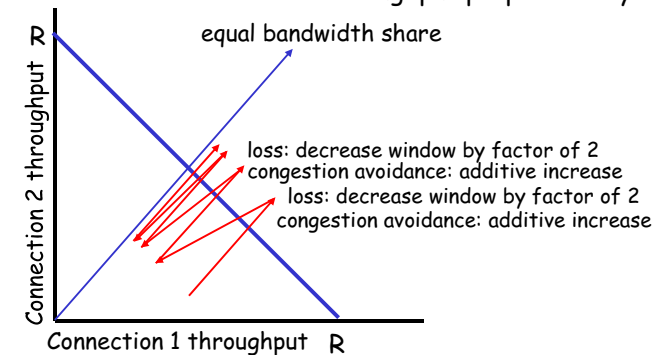


Transport Layer 3-101

## Why is TCP fair?

Two competing sessions:

- Additive increase gives slope of 1, as throughput increases
- multiplicative decrease decreases throughput proportionally



Transport Layer 3-102

## 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 cncctions between 2 hosts.
- Web browsers do this
- Example: link of rate  $R$  supporting 9 cncctions;
  - new app asks for 1 TCP, gets rate  $R/10$
  - new app asks for 11 TCPs, gets  $R/2$  !

Transport Layer 3-103