

## Week 3 Transport Layer



These slides are modified from the slides made available by Kurose and Ross.

*Computer Networking: A Top Down Approach Featuring the Internet, 2nd edition.*  
Jim Kurose, Keith Ross  
Addison-Wesley, July 2002.

Transport Layer 3-1

## Week 3: Transport Layer

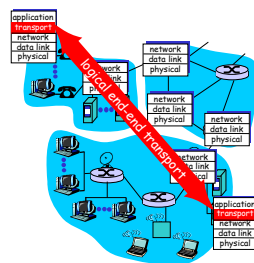
### Our goals:

- understand principles behind transport layer services:
  - multiplexing/demultiplexing
  - reliable data transfer
  - flow control
  - congestion control
- learn about transport layer protocols in the Internet:
  - UDP: connectionless transport
  - TCP: connection-oriented transport
  - TCP congestion control

Transport Layer 3-2

## Transport services and protocols

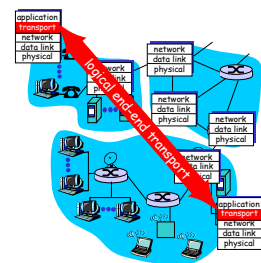
- provide *logical communication* between app processes running on different hosts
- transport protocols run in end systems
  - send side: breaks app messages into **segments**, passes to network layer
  - rcv side: reassembles segments into messages, passes to app layer
- more than one transport protocol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



Transport Layer 3-3

## 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



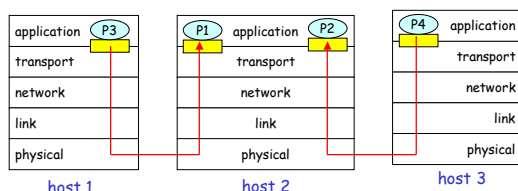
Transport Layer 3-4

## 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)

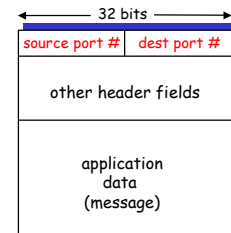
□ = socket    ○ = process



Transport Layer 3-5

## How demultiplexing works

- host receives IP datagrams
  - each datagram has source IP address, destination IP address
  - each datagram carries 1 transport-layer segment
  - each segment has source, destination port number (recall: well-known port numbers for specific applications)
- host uses IP addresses & port numbers to direct segment to appropriate socket



TCP/UDP segment format

Transport Layer 3-6

## Connectionless demultiplexing

- Create sockets with port numbers:

```
DatagramSocket mySocket1 = new
    DatagramSocket(99111);
DatagramSocket mySocket2 = new
    DatagramSocket(99222);
```

- UDP socket identified by two-tuple:

(dest IP address, dest 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 numbers directed to same socket

Transport Layer 3-7

## Connection-oriented demux

- TCP socket identified by 4-tuple:

- source IP address
- source port number
- dest IP address
- dest port number

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

- 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

Transport Layer 3-8

## 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 app

- 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 3-9

## UDP: more

- often used for streaming multimedia apps

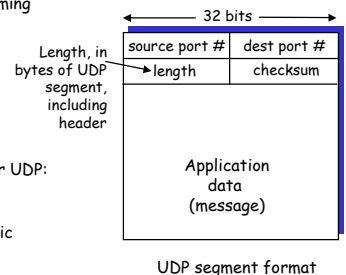
- loss tolerant
- rate sensitive

- other UDP uses

- DNS
- SNMP

- reliable transfer over UDP: add reliability at application layer

- application-specific error recovery!



Transport Layer 3-10

## 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

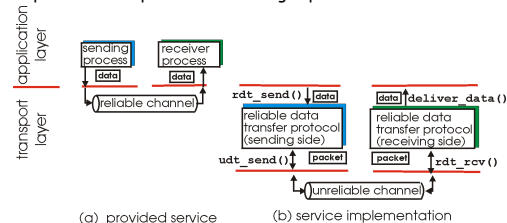
### Receiver:

- 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*

Transport Layer 3-11

## Principles of Reliable data transfer

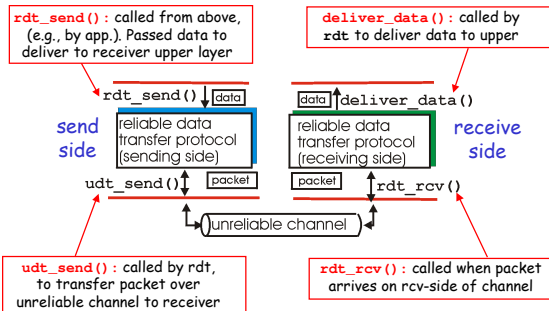
- important in app., transport, link layers
- top-10 list of important networking topics!



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

Transport Layer 3-12

## Reliable data transfer: getting started

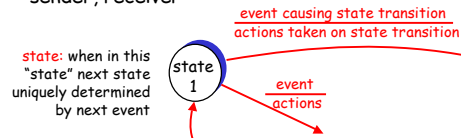


Transport Layer 3-13

## Reliable data transfer: getting started

We'll:

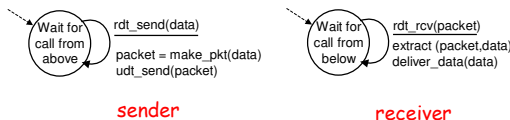
- incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)
- consider only unidirectional data transfer
  - but control info will flow on both directions!
- use finite state machines (FSM) to specify sender, receiver



Transport Layer 3-14

## 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



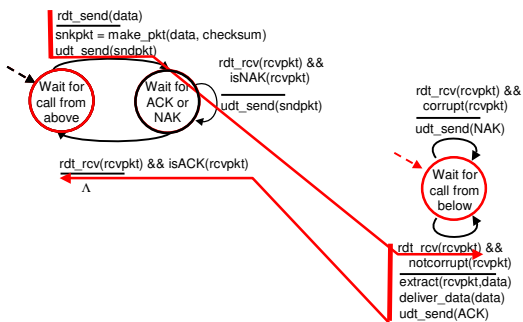
Transport Layer 3-15

## 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
  - human scenarios using ACKs, NAKs?
- new mechanisms in rdt2.0 (beyond rdt1.0):
  - error detection
  - receiver feedback: control msgs (ACK,NAK) rcvr->sender

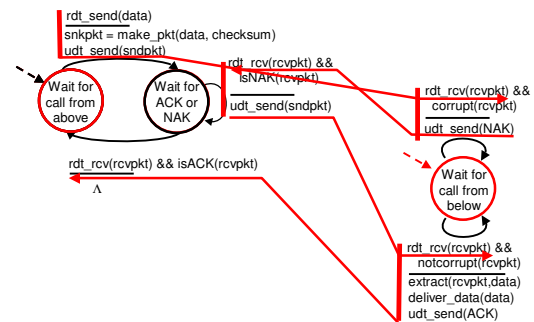
Transport Layer 3-16

## rdt2.0: operation with no errors



Transport Layer 3-17

## rdt2.0: error scenario



Transport Layer 3-18

## 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

### 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

Transport Layer 3-19

## rdt2.2: a NAK-free protocol

- same functionality as rdt2.1, using NAKs 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 3-20

## 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?

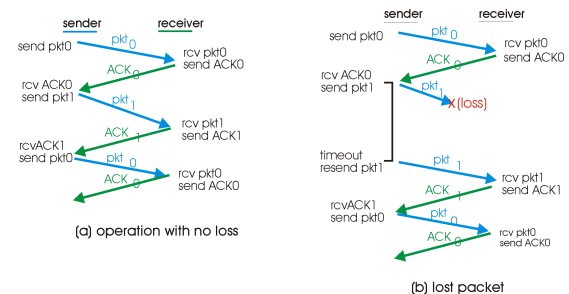
- sender waits until certain data or ACK lost, then retransmits
- yuck: drawbacks?

### 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

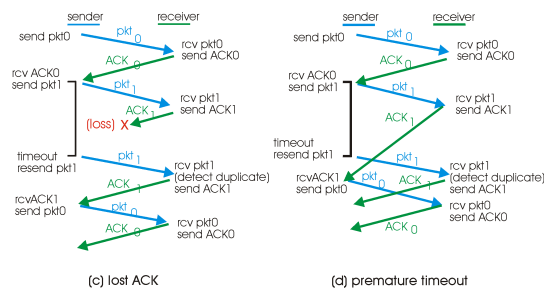
Transport Layer 3-21

## rdt3.0 in action



Transport Layer 3-22

## rdt3.0 in action



Transport Layer 3-23

## Performance of rdt3.0

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

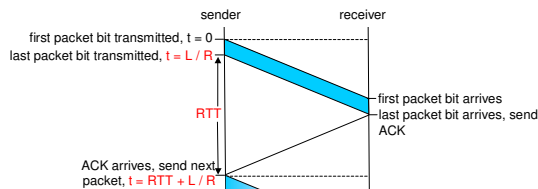
$$T_{\text{transmit}} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8\text{kb/pkt}}{10^9 \text{ b/sec}} = 8 \text{ microsec}$$

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

- $U_{\text{sender}}$ : utilization - fraction of time sender busy sending
- 1KB pkt every 30 msec → 33kB/sec throughput over 1 Gbps link
- network protocol limits use of physical resources!

Transport Layer 3-24

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



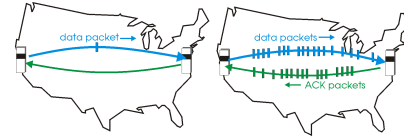
$$U_{\text{sender}} = \frac{L/R}{RTT + L/R} = \frac{.008}{30,008} = 0.00027$$

Transport Layer 3-25

## 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

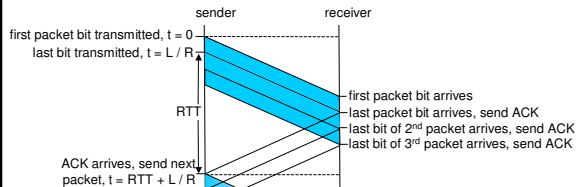


(a) a stop-and-wait protocol in operation (b) a pipelined protocol in operation

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

Transport Layer 3-26

## Pipelining: increased utilization



$$U_{\text{sender}} = \frac{3 * L/R}{RTT + L/R} = \frac{.024}{30,008} = 0.0008$$

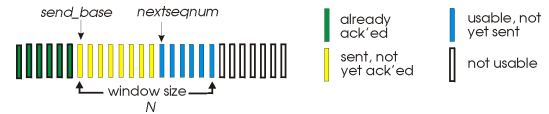
Increase utilization by a factor of 3!

Transport Layer 3-27

## Go-Back-N

**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 deceive duplicate ACKs (see receiver)
- timer for each in-flight pkt**
- timeout(n):** retransmit pkt n and all higher seq # pkts in window

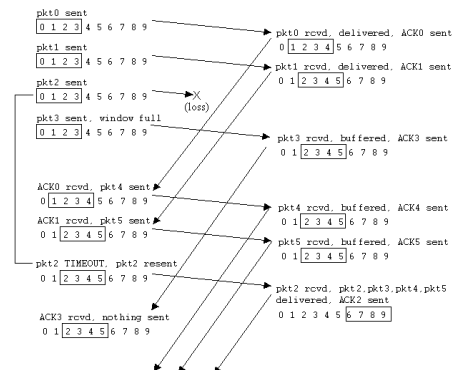
Transport Layer 3-28

## 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 3-29

## Selective repeat in action

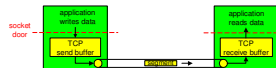


† Layer 3-30

## TCP: Overview

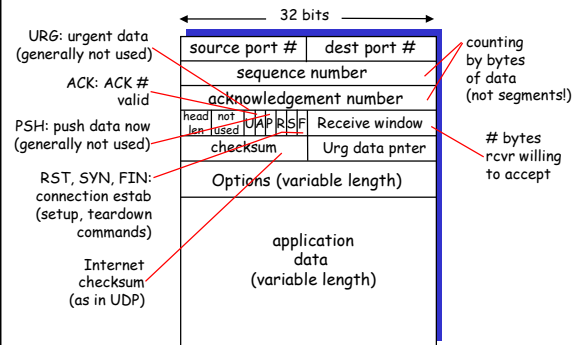
RFCs: 793, 1122, 1323, 2018, 2581

- **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 3-31

## TCP segment structure



Transport Layer 3-32

## TCP seq. #'s and ACKs

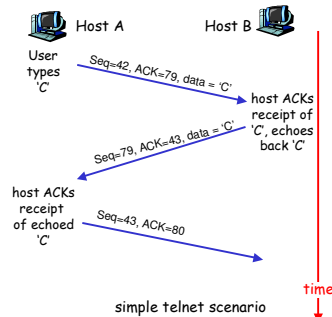
### Seq. #'s:

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

### ACKs:

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

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



Transport Layer 3-33

## TCP Round Trip Time and Timeout

### Q: how to set TCP timeout value?

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

### Q: how to estimate RTT?

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

Transport Layer 3-34

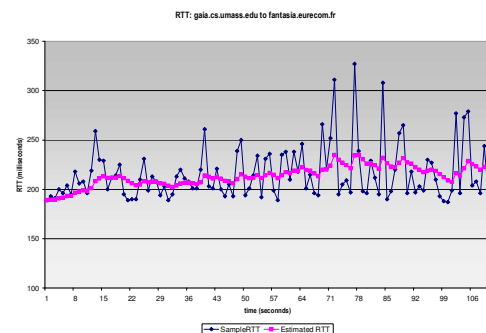
## TCP Round Trip Time and Timeout

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

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

Transport Layer 3-35

## Example RTT estimation:



Transport Layer 3-36

## 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:

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

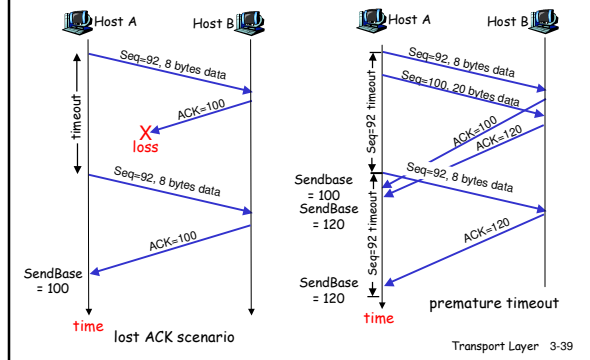
Transport Layer 3-37

## TCP reliable data transfer

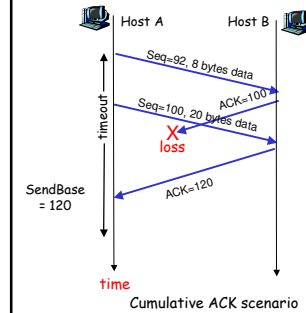
- TCP creates rdt service on top of IP's unreliable service
- Pipelined segments
- Cumulative acks
- TCP uses single retransmission timer
- Retransmissions are triggered by:
  - timeout events
  - duplicate acks
- Initially consider simplified TCP sender:
  - ignore duplicate acks
  - ignore flow control,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

Transport Layer 3-38

## TCP: retransmission scenarios



## TCP retransmission scenarios (more)



## 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 3-41

## Fast retransmit algorithm:

```

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

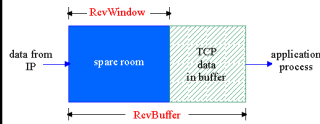
a duplicate ACK for already ACKed segment

fast retransmit

Transport Layer 3-42

## TCP Flow Control

- receive side of TCP connection has a receive buffer:



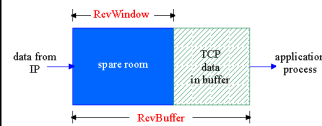
- app process may be slow at reading from buffer

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

- speed-matching service: matching the send rate to the receiving app's drain rate

Transport Layer 3-43

## TCP Flow control: how it works



(Suppose TCP receiver discards out-of-order segments)

- spare room in buffer  
=  $RcvWindow$   
=  $RcvBuffer - [LastByteRcvd - LastByteRead]$

- Rcvr advertises spare room by including value of  $RcvWindow$  in segments
- Sender limits unACKed data to  $RcvWindow$ 
  - guarantees receive buffer doesn't overflow

Transport Layer 3-44

## 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();`

### Three way handshake:

- Step 1:** client host sends TCP SYN segment to server
  - specifies initial seq #
  - no data
- Step 2:** server host receives SYN, replies with SYNACK segment
  - server allocates buffers
  - specifies server initial seq. #
- Step 3:** client receives SYNACK, replies with ACK segment, which may contain data

Transport Layer 3-45

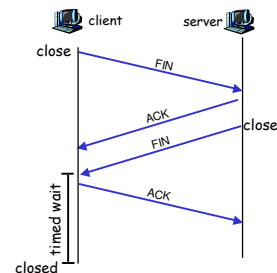
## TCP Connection Management (cont.)

### Closing a connection:

client closes socket:  
`clientSocket.close();`

**Step 1:** client end system sends TCP FIN control segment to server

**Step 2:** server receives FIN, replies with ACK. Closes connection, sends FIN.



Transport Layer 3-46

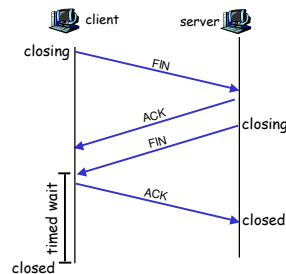
## TCP Connection Management (cont.)

**Step 3:** client receives FIN, replies with ACK.

- Enters "timed wait" - will respond with ACK to received FINs

**Step 4:** server, receives ACK. Connection closed.

**Note:** with small modification, can handle simultaneous FINs.



Transport Layer 3-47

## 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 3-48



## TCP Congestion Control

- end-end control (no network assistance)

- sender limits transmission:

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{CongWin}$$

- 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:

- AIMD
- slow start
- conservative after timeout events

Transport Layer 3-49