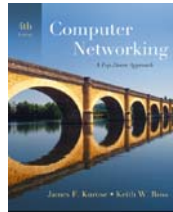


# Transport Layer



*Computer Networking: A Top Down Approach*  
4th edition.  
Jim Kurose, Keith Ross  
Addison-Wesley, July 2007.

These slides adapted from those made available by the text authors.

All material copyright 1996-2007  
J.F. Kurose and K.W. Ross, All Rights Reserved

## Transport Layer (Chapter 3 in KR)

### Our goals:

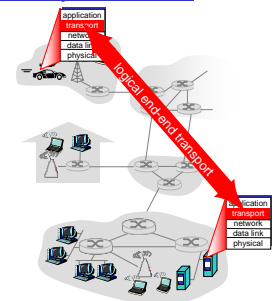
- understand principles behind transport layer services:
  - multiplexing/demultiplexing
  - reliable data transfer
  - flow control
  - congestion control
- learn about transport layer protocols in the Internet:
  - UDP: connectionless transport
  - TCP: connection-oriented transport
  - TCP congestion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

## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

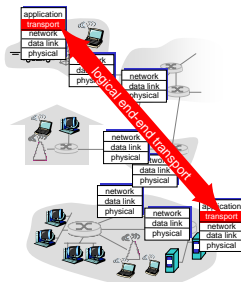
## Transport services and protocols

- provide *logical communication* between app **processes** running on different nodes
- 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



## 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



## Outline

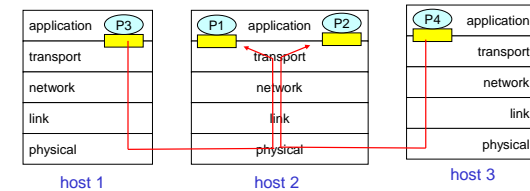
- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

## 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)

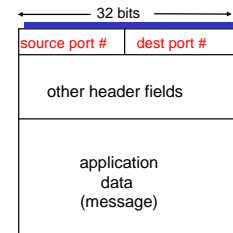
■ = socket ○ = process



Transport Layer 7

## How demultiplexing works

- receive packet from network layer
  - source, destination addresses
  - 1 segment per packet
  - source, destination port numbers
- addresses & port numbers direct segment to appropriate socket

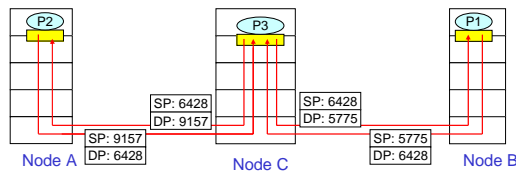


TCP/UDP segment format

Transport Layer 8

## Connectionless demultiplexing

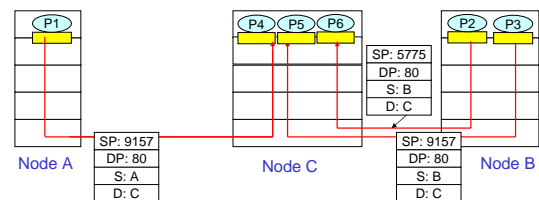
- UDP socket identified by (destination address, destination port number)
- Upon receipt of UDP segment, message is delivered to corresponding port
- Packets from different source address and/or port number are treated the same



Transport Layer 9

## Connection-oriented demux

- TCP socket identified by 4-tuple (source IP address, source port number, dest IP address, dest port number)
- All four values direct segment to appropriate socket
- Each node may support many TCP sockets/sessions



Transport Layer 10

## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

Transport Layer 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 app
- connectionless:
  - no handshaking between UDP sender, receiver
  - each UDP segment handled independently of others

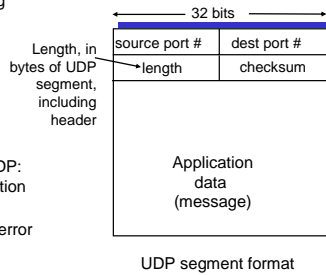
### Why is there a UDP?

- 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 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 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:

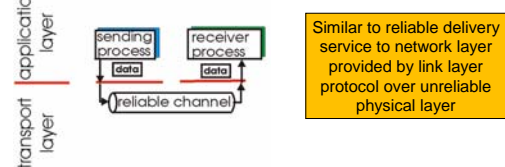
- compute checksum of received segment
- check if computed checksum equals checksum field value:
  - NO - error detected
  - YES - no error detected

## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

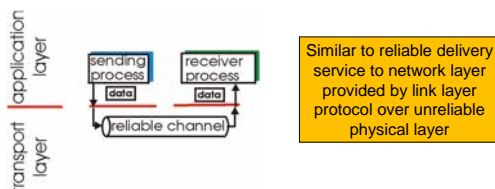
## Principles of Reliable data transfer

- Reliable message delivery between applications provided by transport layer
- Characteristics of unreliable channel below determines complexity of reliable data transfer protocol
- Sound familiar? Link layer provides reliable single-hop data transfer over unreliable physical layer



## Principles of Reliable data transfer

- Stop-and-wait, Go-Back-N, and Selective repeat protocols are used at transport layer.
- We've already covered these protocols, so we're not going to do it again.
- Review details of these protocols at link layer on your own.



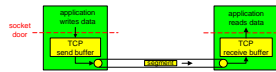
## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

## TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

- **point-to-point:**
  - one sender, one receiver
- **reliable, in-order byte stream:**
  - no "message boundaries"
- **pipelined:**
  - TCP congestion and flow control set window size
- **send & receive buffers**
- **full duplex data:**
  - bi-directional data flow in same connection
  - MSS: maximum segment size
- **connection-oriented:**
  - handshaking (exchange of control msgs) init's sender, receiver state before data exchange
- **flow controlled:**
  - sender will not overwhelm receiver



Transport Layer 19

## TCP seq. #'s and ACKs

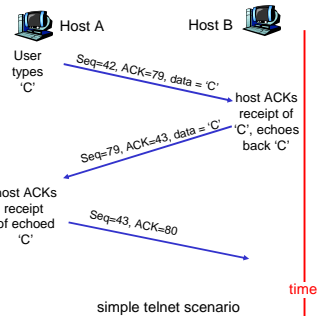
### Seq. #'s:

- 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



simple telnet scenario

Transport Layer 20

## 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?

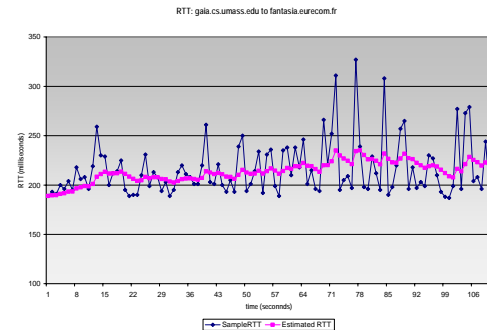
- **SampleRTT**: measured time from segment transmission until ACK receipt
  - ignore retransmissions
- **SampleRTT** will vary, want estimated RTT "smoother"
  - average several recent measurements, not just current SampleRTT

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

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

Transport Layer 21

## Example RTT estimation:



Transport Layer 22

## 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:

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

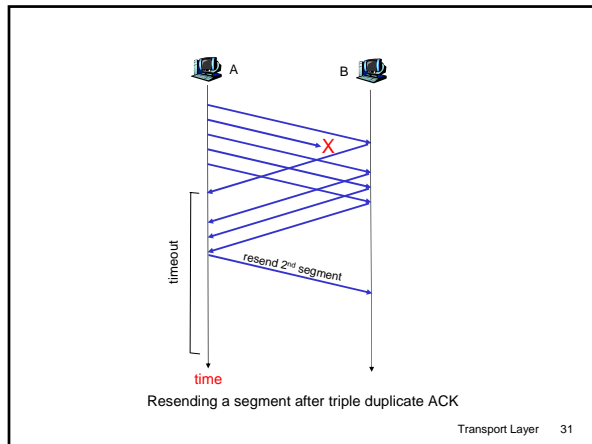
Transport Layer 23

## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

Transport Layer 24





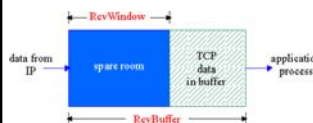
## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

Transport Layer 32

## TCP Flow Control

- receive side of TCP connection has a receive buffer:



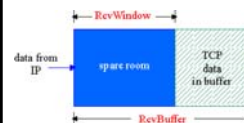
- app process may be slow at reading from buffer

**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 33

## 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 34

## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

Transport Layer 35

## 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
- server: contacted by client

### Three way handshake:

**Step 1:** client sends TCP SYN segment to server

- specifies initial seq #
- no data

**Step 2:** server 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 36

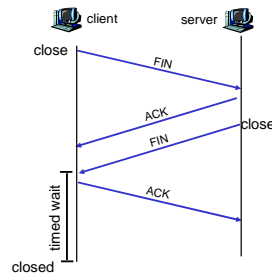
## TCP Connection Management (cont.)

### Closing a connection:

client closes socket

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

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



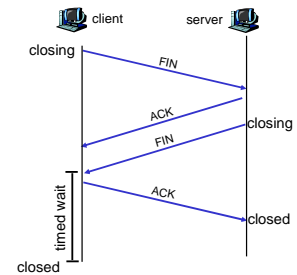
## TCP Connection Management (cont.)

**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.



## Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

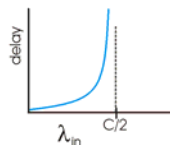
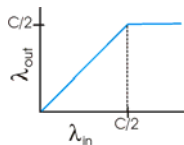
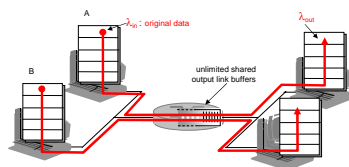
## Principles of Congestion Control

### 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)

## Causes/costs of congestion: scenario 1

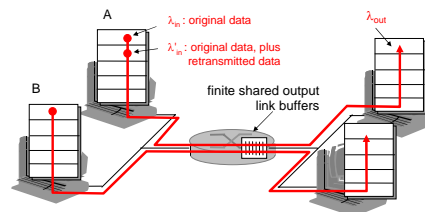
- two senders, two receivers
- one router, infinite buffers
- no retransmission



- large delays when congested
- maximum achievable throughput

## Causes/costs of congestion: scenario 2

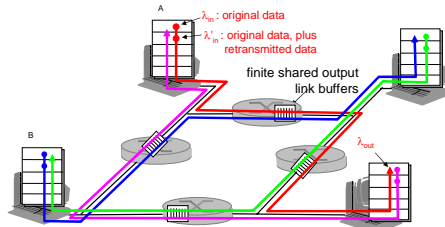
- one router, *finite* buffers
- sender retransmission of lost packet



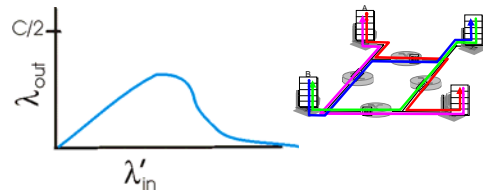
### Causes/costs of congestion: scenario 3

- four senders
- multihop paths
- timeout/retransmit

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



### Causes/costs of congestion: scenario 3



Another "cost" of congestion:

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

### 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:

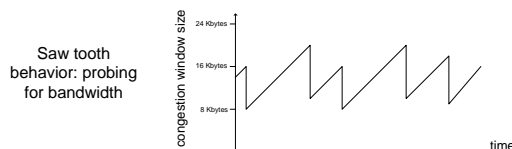
- routers provide feedback to end systems
  - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
  - explicit rate sender should send at

### Outline

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

### TCP congestion control: additive increase, multiplicative decrease

- Approach: increase transmission rate (window size), probing for usable bandwidth, until loss occurs
  - additive increase: increase CongWin by 1 MSS every RTT until loss detected
  - multiplicative decrease: cut CongWin in half after loss



Saw tooth behavior: probing for bandwidth

### TCP Congestion Control: details

- 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

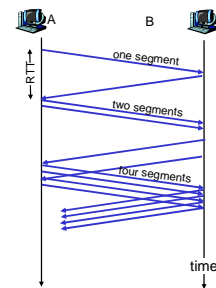


## TCP Slow Start

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

## 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



## 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:

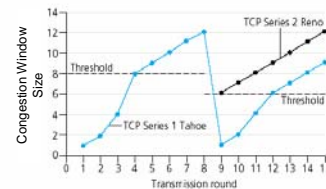
- 3 dup ACKs indicates network capable of delivering some segments
- timeout indicates a "more alarming" congestion scenario

## Refinement

- 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 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

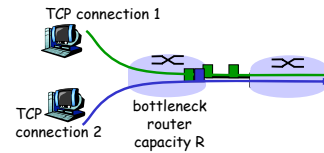
State	Event	TCP Sender Action	Commentary
Slow Start (SS)	ACK receipt for previously unacked data	$\text{CongWin} = \text{CongWin} + \text{MSS}$ , If ( $\text{CongWin} > \text{Threshold}$ ) set state to "Congestion Avoidance"	Resulting in a doubling of CongWin every RTT
Congestion Avoidance (CA)	ACK receipt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

## TCP throughput

- ❑ What's the average throughput of TCP as a function of window size and RTT?
  - Ignore slow start
- ❑ Let  $W$  be the window size when loss occurs.
- ❑ When window is  $W$ , throughput is  $W/RTT$
- ❑ Just after loss, window drops to  $W/2$ , throughput to  $W/2RTT$ .
- ❑ Average throughput:  $.75 W/RTT$

## TCP Fairness

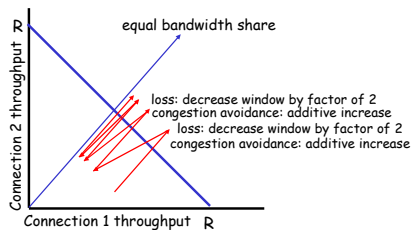
**Fairness goal:** if  $K$  TCP sessions share same bottleneck link of bandwidth  $R$ , each should have average rate of  $R/K$



## Why is TCP fair?

Two competing sessions:

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



## 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  $R/2$  !