

---

# TCP

## Flow Control and Congestion Control

EECS 489 Computer Networks

<http://www.eecs.umich.edu/courses/eecs489/w07>

Z. Morley Mao

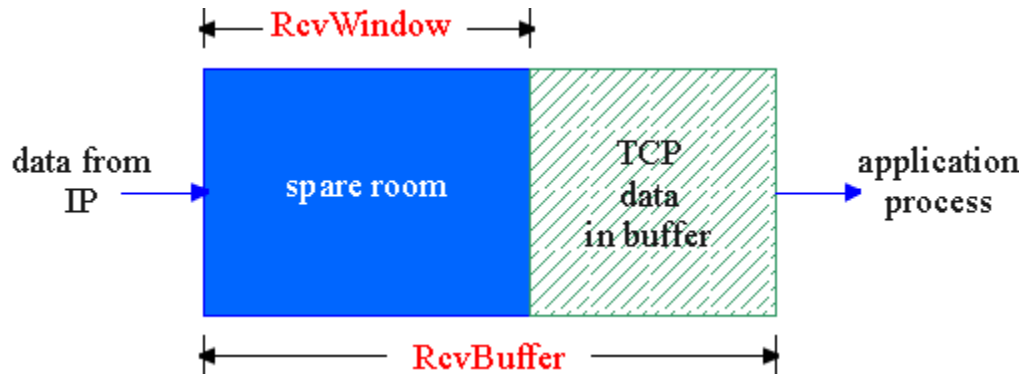
Monday Feb 5, 2007

# TCP Flow Control

- receive side of TCP connection has a receive buffer:

## flow control

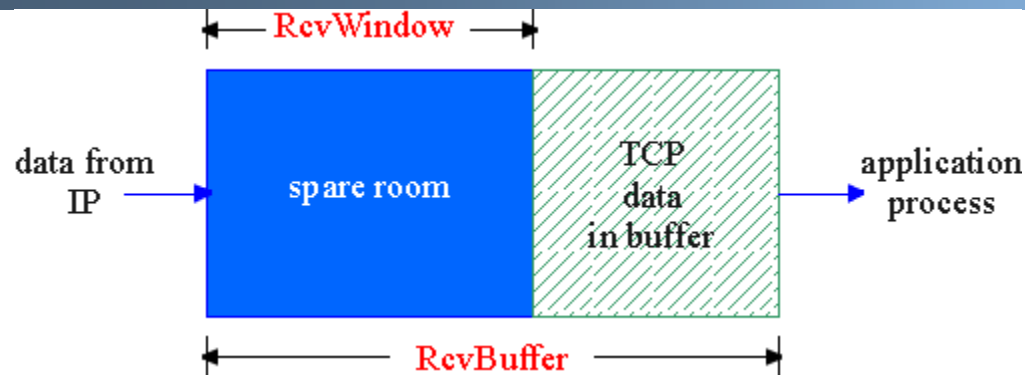
sender won't overflow receiver's buffer by transmitting too much, too fast



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

app process may be slow at reading from buffer

# 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

# 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

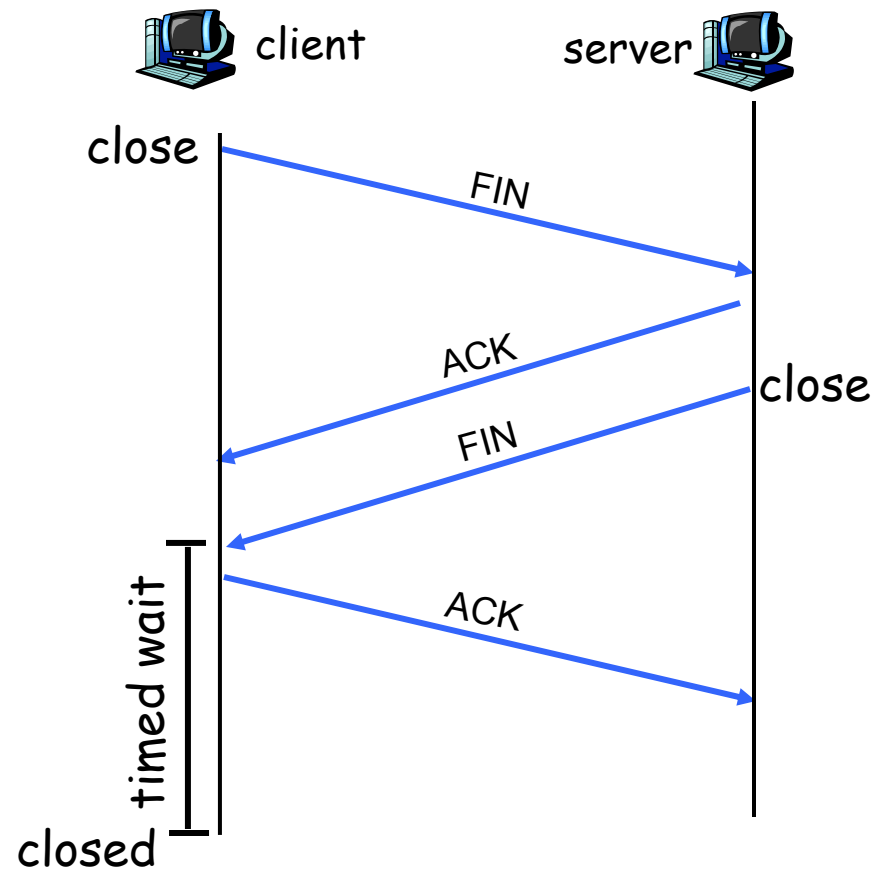
# TCP Connection Management (cont.)

## 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.



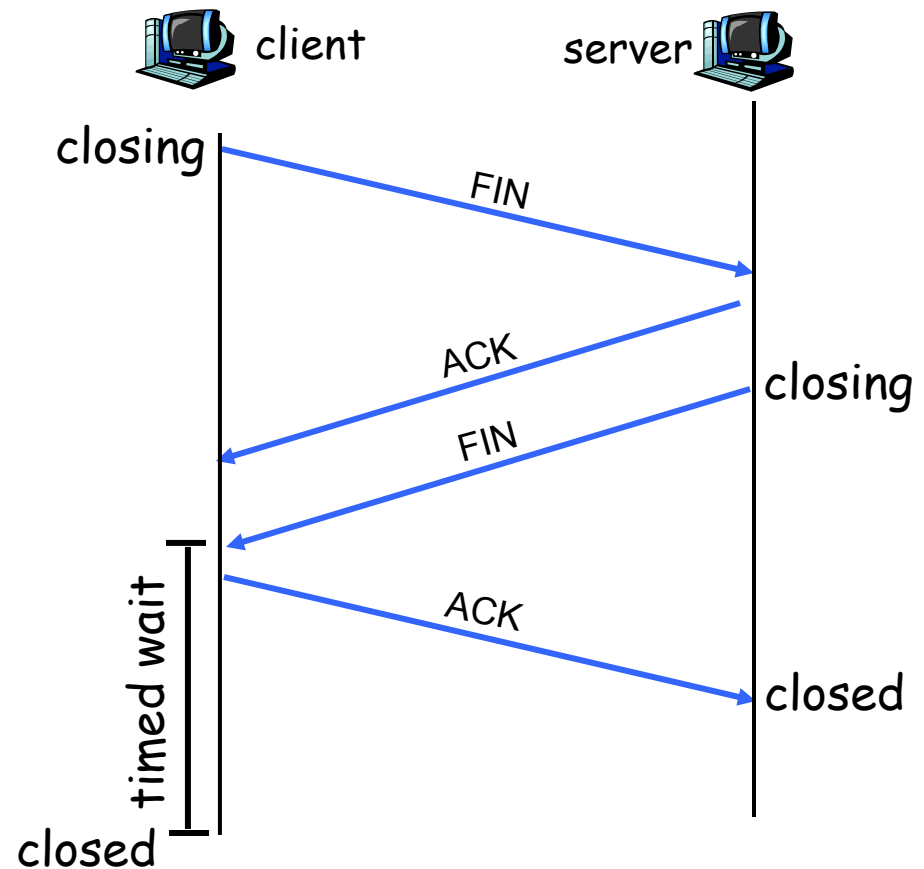
# TCP Connection Management (cont.)

Step 3: client receives FIN,  
replies with ACK.

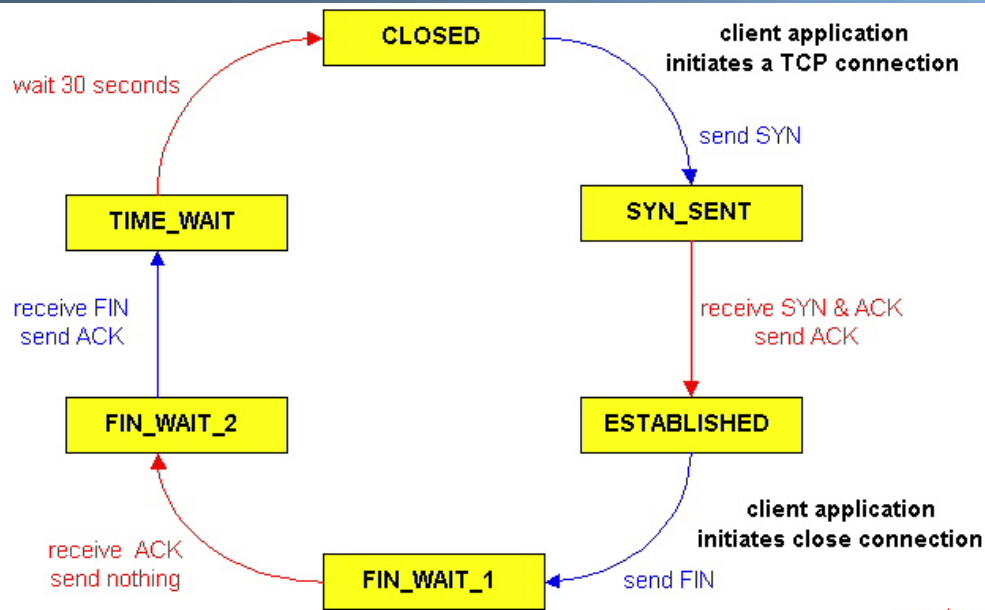
- Enters “timed wait” - will respond with ACK to received FINs

Step 4: server, receives ACK.  
Connection closed.

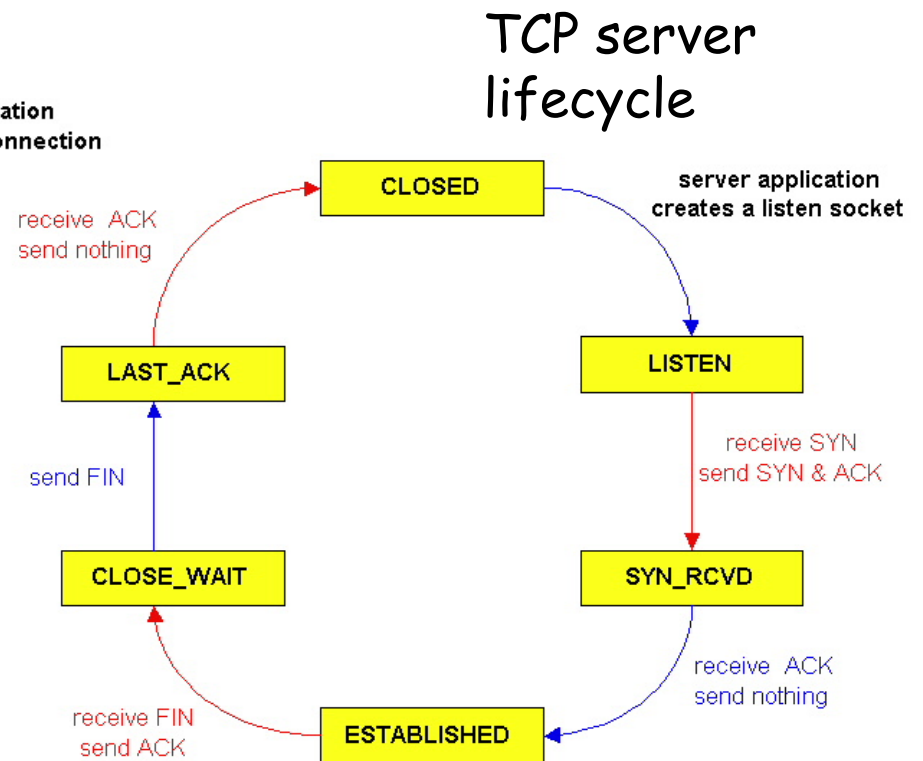
Note: with small modification,  
can handle simultaneous  
FINs.



# TCP Connection Management (cont)



TCP client lifecycle



TCP server lifecycle

# Principles of Congestion Control

---

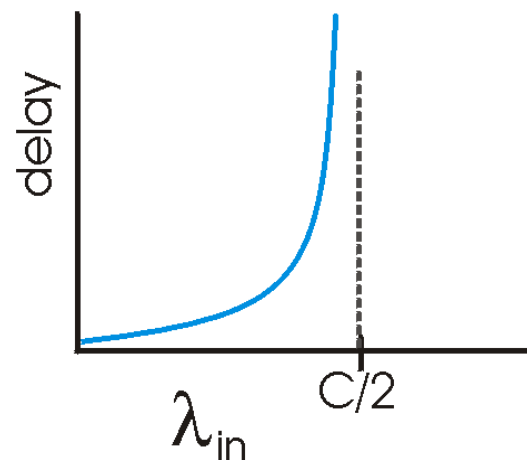
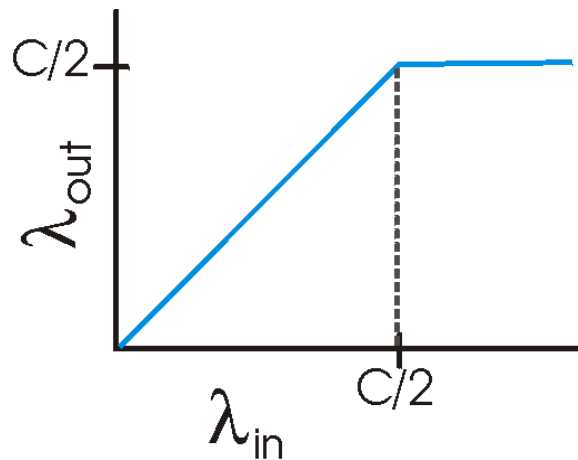
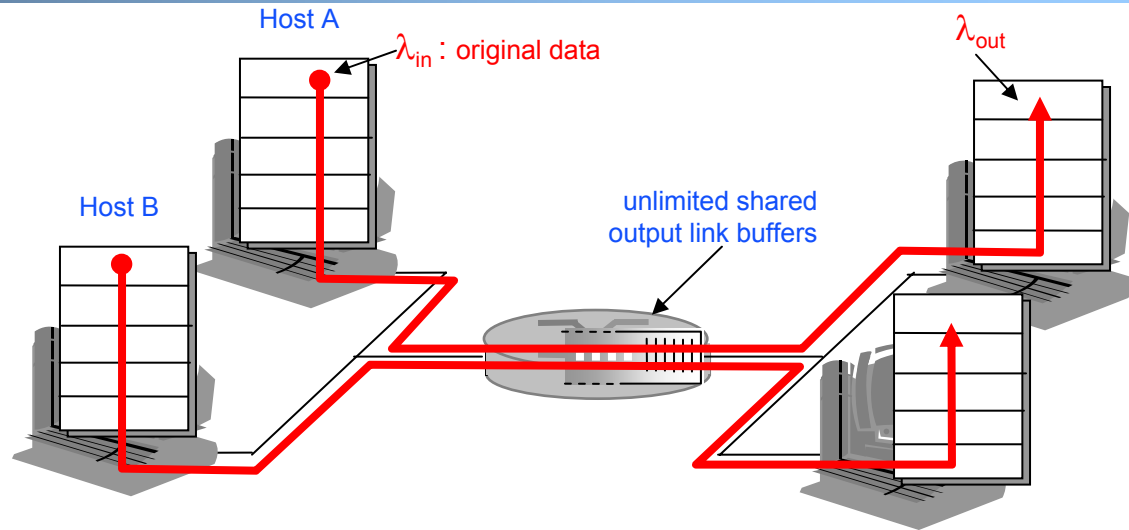
## Congestion:

- 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!



# Causes/costs of congestion: scenario 1

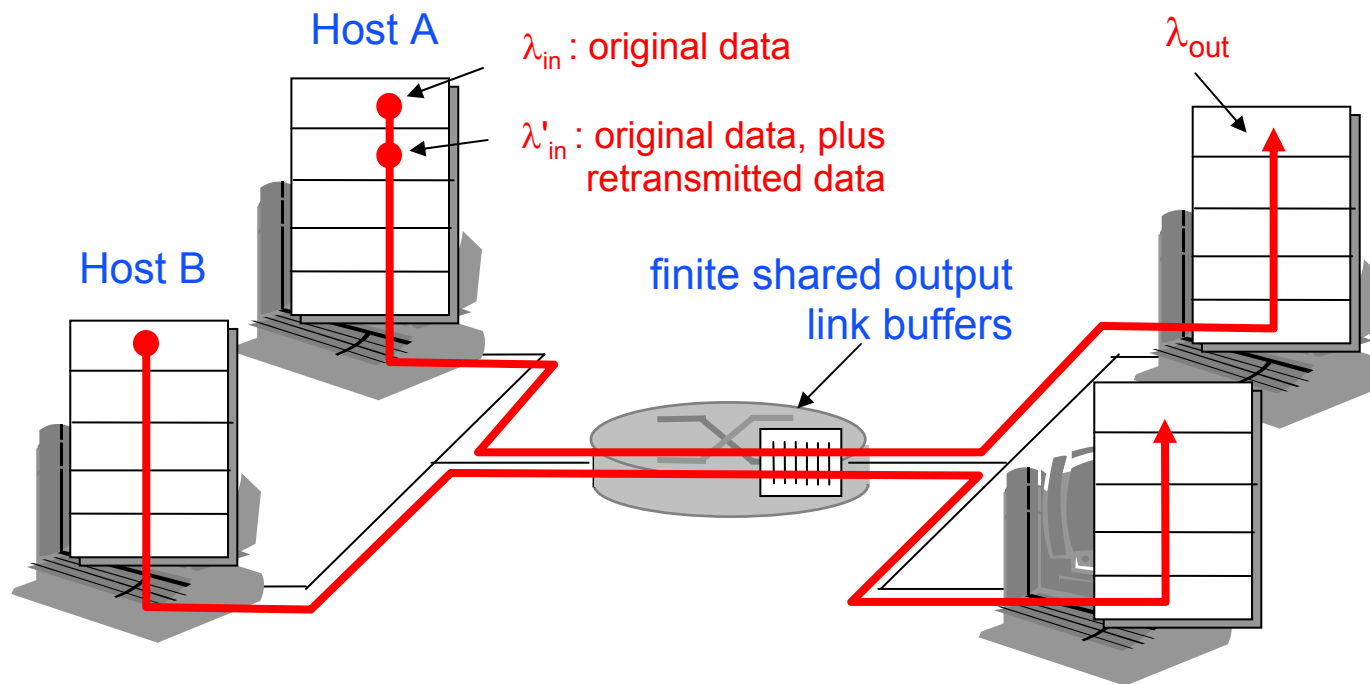
- two senders, two receivers
- one router, infinite buffers
- no retransmission



- large delays when congested
- maximum achievable throughput

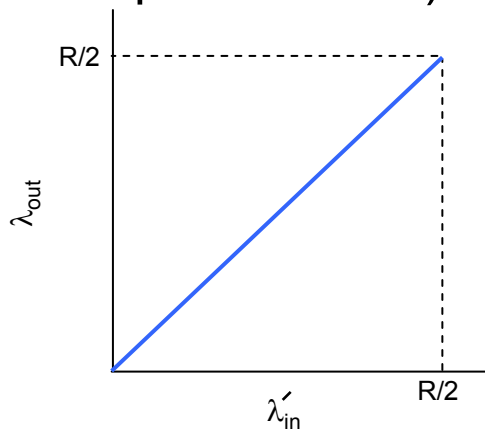
# Causes/costs of congestion: scenario 2

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

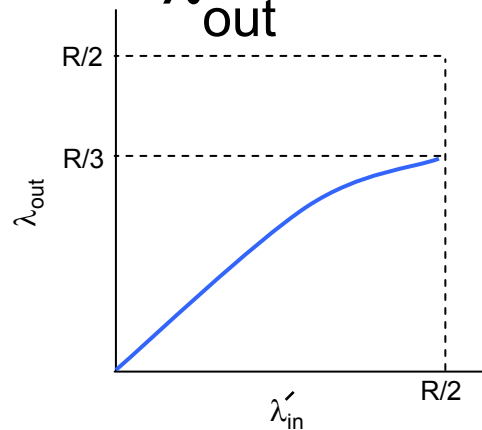


# Causes/costs of congestion: scenario 2

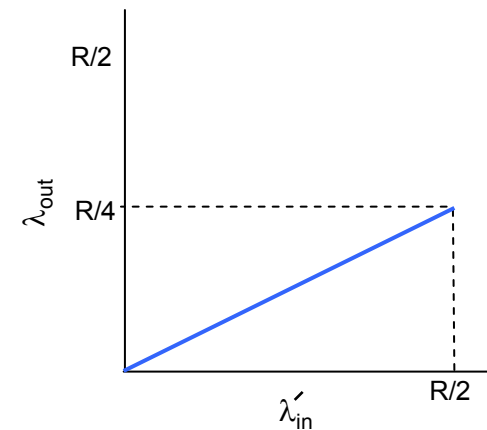
- always:  $\lambda_{in} = \lambda_{out}(\text{goodput})$
- “perfect” retransmission only when loss:  $\lambda'_{in} > \lambda_{out}$
- retransmission of delayed (not lost) packet makes  $\lambda'_{in}$  larger (than perfect case) for same  $\lambda_{out}$



a.



b.



c.

“costs” of congestion:

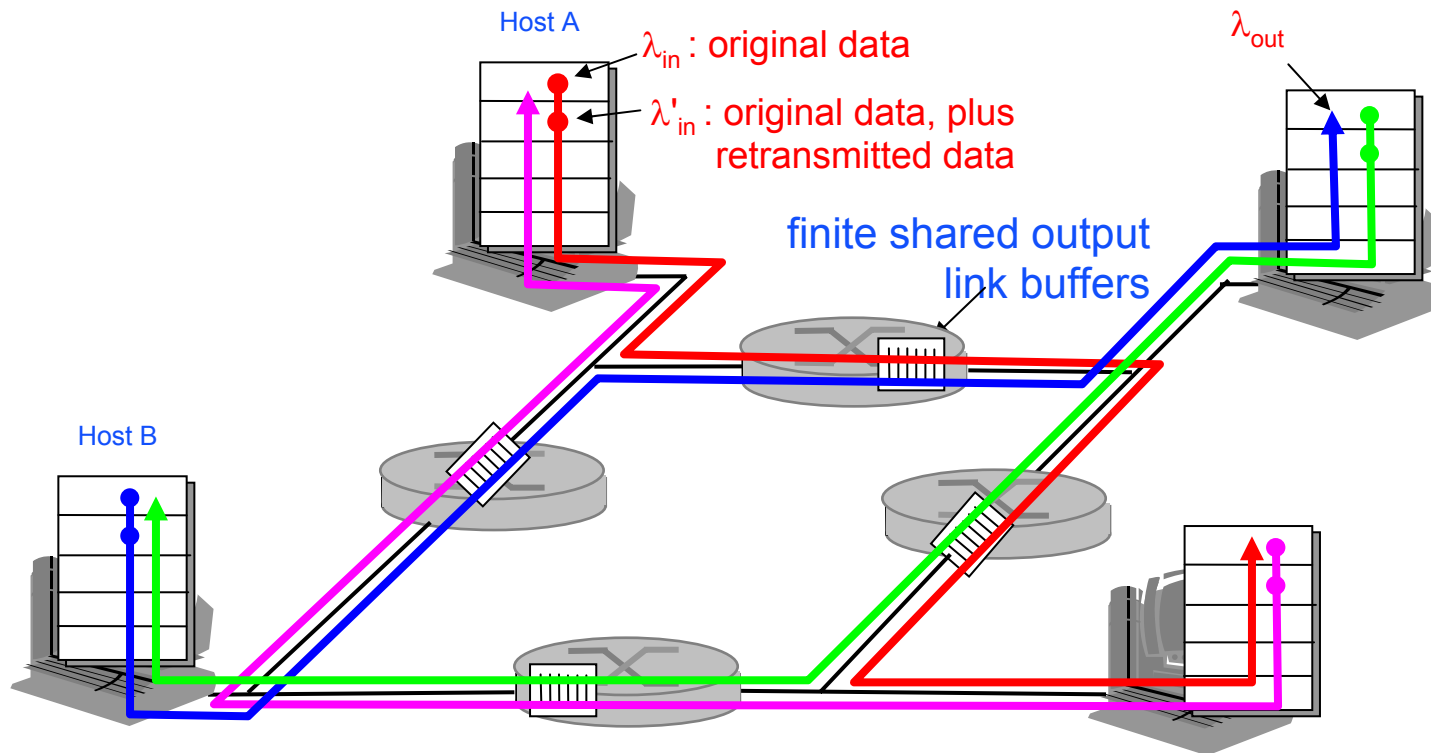
more work (retrans) for given “goodput”

unnneeded retransmissions: link carries multiple copies of pkt

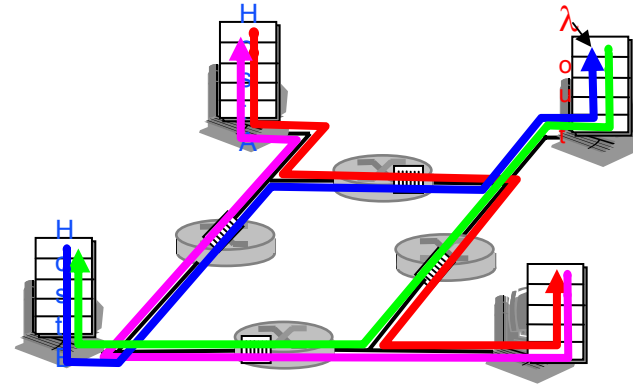
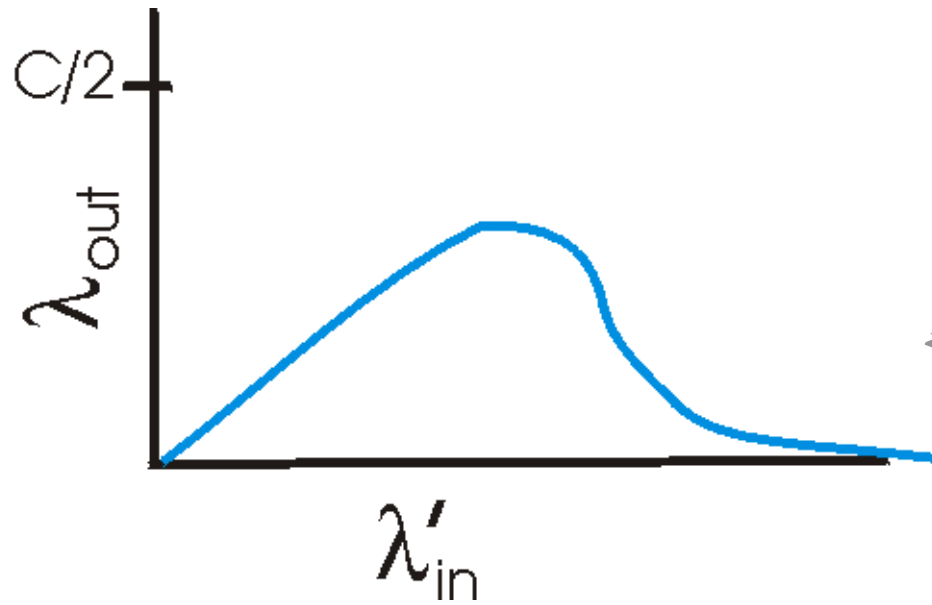
# Causes/costs of congestion: scenario 3

- 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

# Case study: ATM ABR congestion control

---

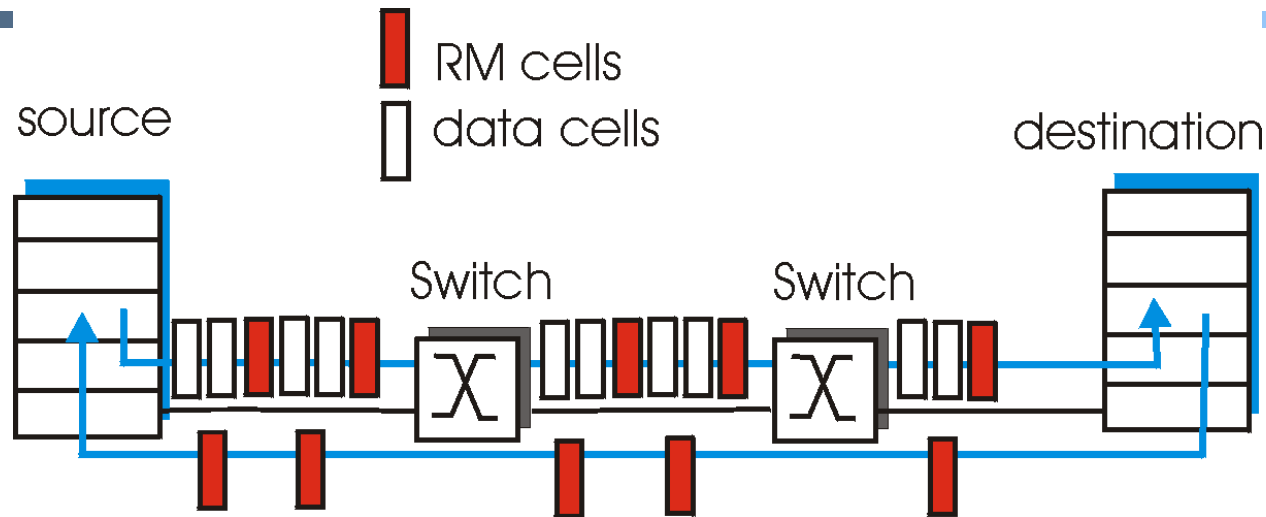
## ABR: available bit rate:

- “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

# 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 minimum 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



# TCP Congestion Control

- end-end control (no network assistance)

- sender limits transmission:

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

- Roughly,

- CongWin is dynamic, function of perceived network congestion

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

## How does sender perceive congestion?

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

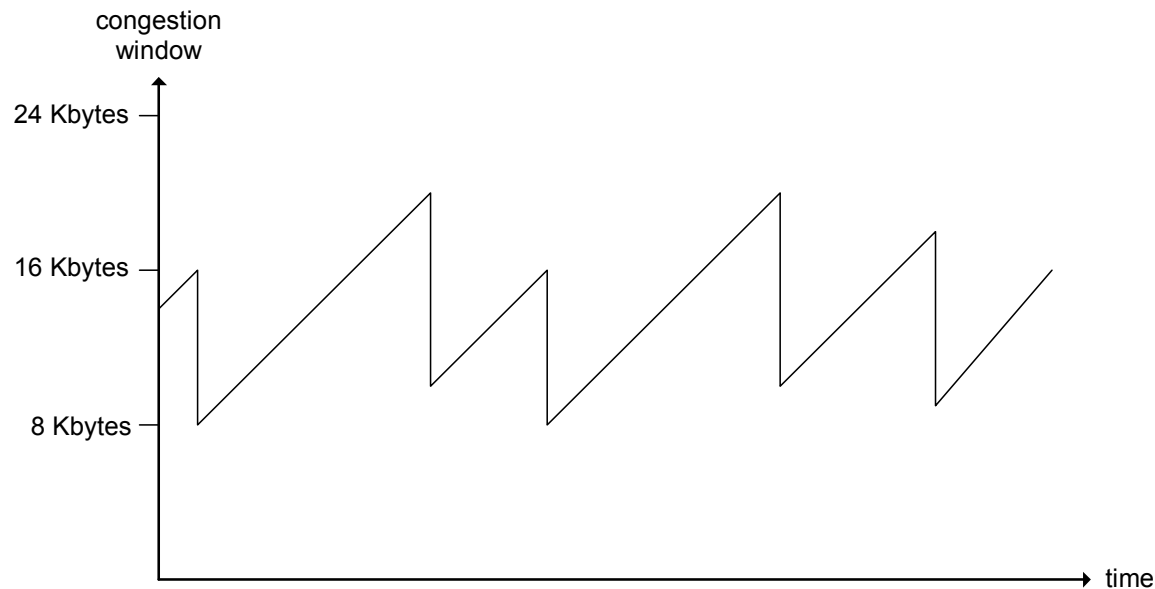
## three mechanisms:

- AIMD
- slow start
- conservative after timeout events

# TCP 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*



Long-lived TCP connection

# TCP Slow Start

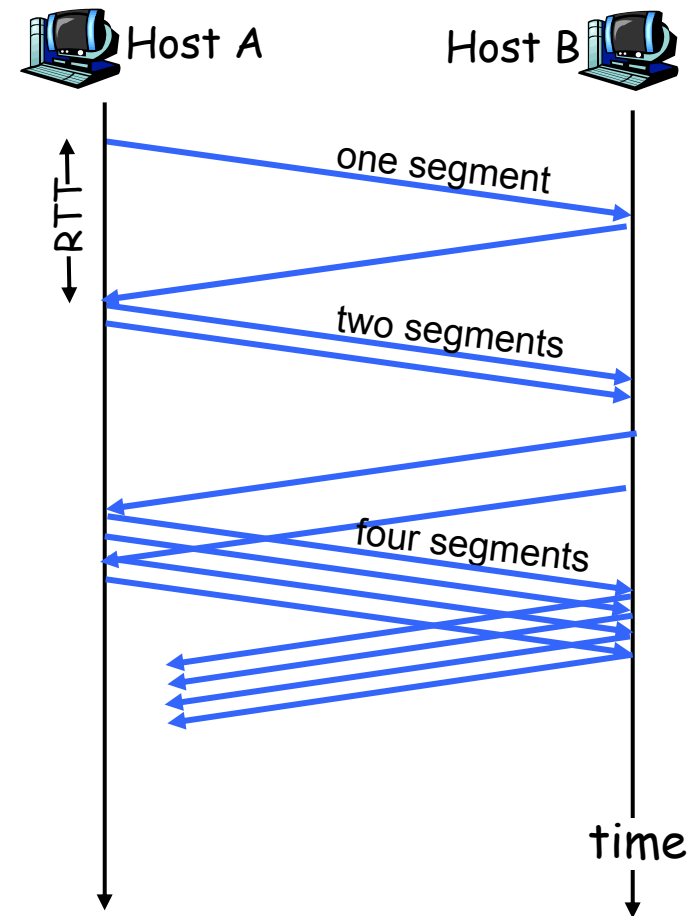
---

- 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

## Philosophy:

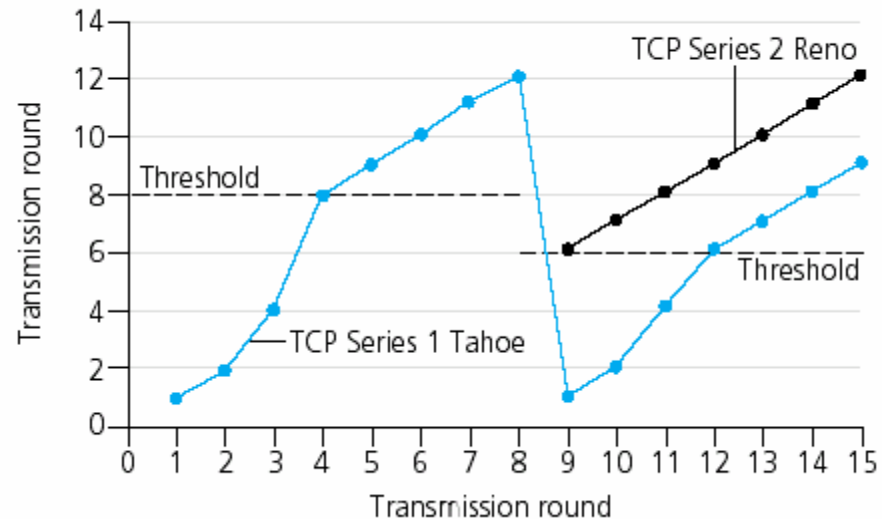
- After 3 dup ACKs:
  - `CongWin` is cut in half
  - window then grows linearly
- But after timeout event:
  - `CongWin` instead set to 1 MSS;
  - window then grows exponentially
  - to a threshold, then grows linearly

- 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

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

---

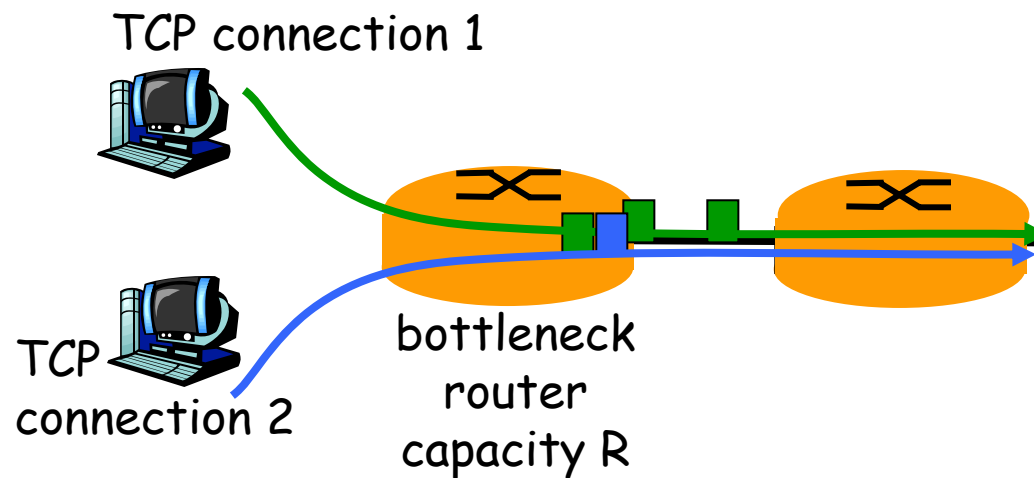
- Example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- Requires window size  $W = 83,333$  in-flight segments
- Throughput in terms of loss rate:

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

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

# TCP Fairness

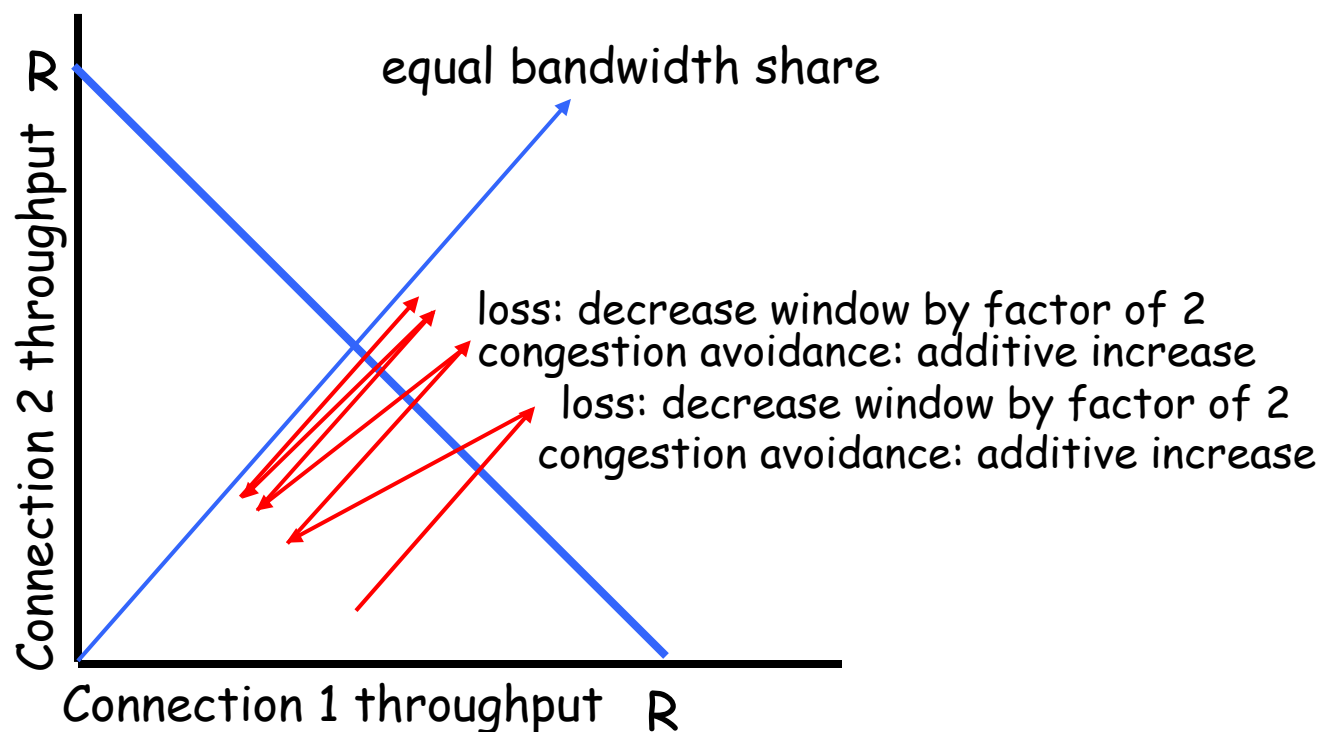
**Fairness goal:** if  $K$  TCP sessions share same bottleneck link of bandwidth  $R$ , each should have average rate of  $R/K$



# 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 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  $R/2$  !

# Delay modeling

---

Q: How long does it take to receive an object from a Web server after sending a request?

Ignoring congestion, delay is influenced by:

- TCP connection establishment
- data transmission delay
- slow start

Notation, assumptions:

- Assume one link between client and server of rate  $R$
- $S$ : MSS (bits)
- $O$ : object size (bits)
- no retransmissions (no loss, no corruption)

Window size:

- First assume: fixed congestion window,  $W$  segments
- Then dynamic window, modeling slow start

# TCP Delay Modeling: Slow Start (1)

Now suppose window grows according to slow start

Will show that the delay for one object is:

$$Latency = 2RTT + \frac{O}{R} + P \left[ RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R}$$

where  $P$  is the number of times TCP idles at server:

$$P = \min\{Q, K - 1\}$$

- where  $Q$  is the number of times the server idles if the object were of infinite size.
- and  $K$  is the number of windows that cover the object.

# TCP Delay Modeling: Slow Start (2)

## Delay components:

- 2 RTT for connection estab and request
- $O/R$  to transmit object
- time server idles due to slow start

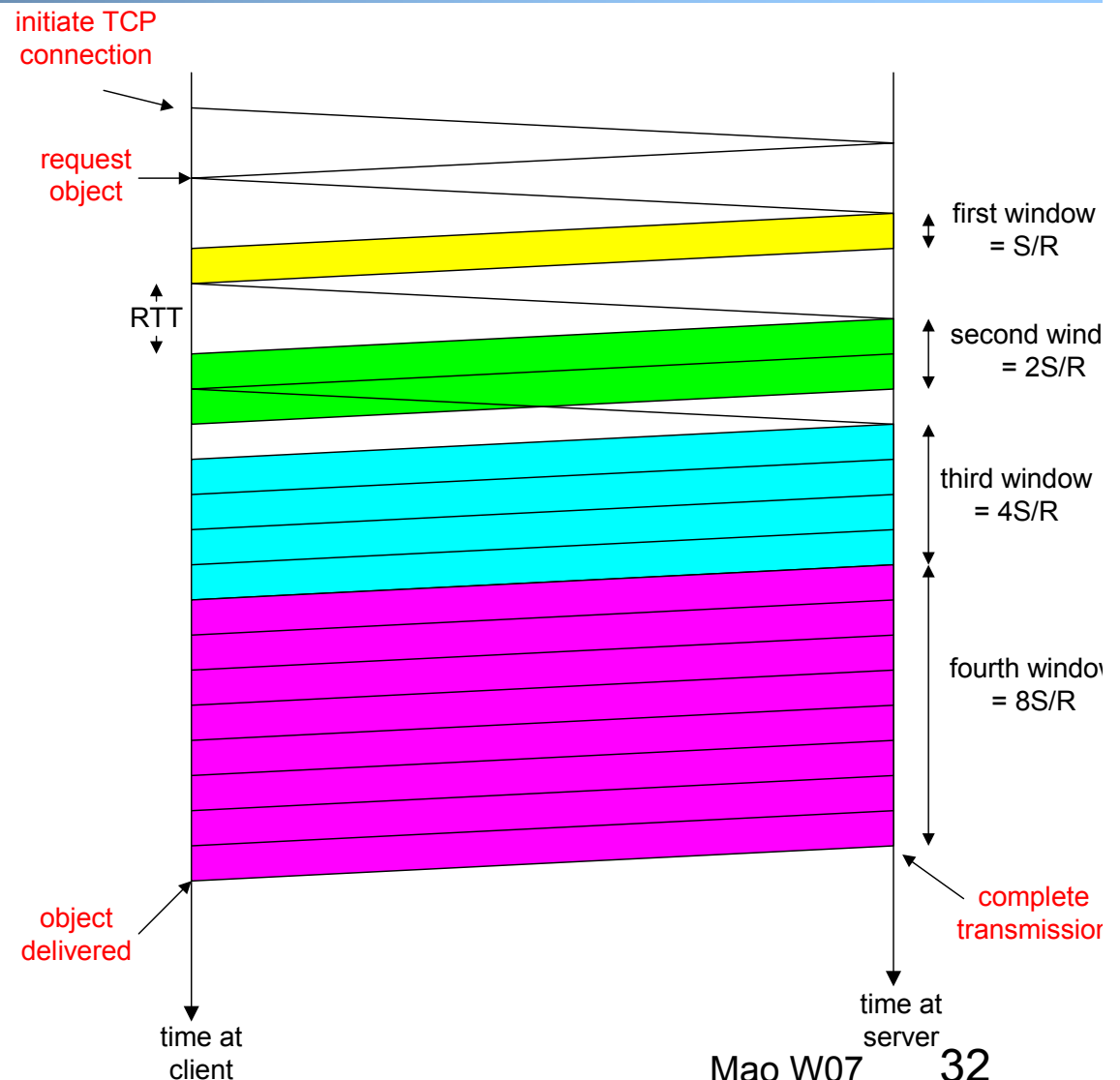
## Server idles:

$$P = \min\{K-1, Q\} \text{ times}$$

## Example:

- $O/S = 15$  segments
- $K = 4$  windows
- $Q = 2$
- $P = \min\{K-1, Q\} = 2$

Server idles  $P=2$  times





# TCP Delay Modeling (3)

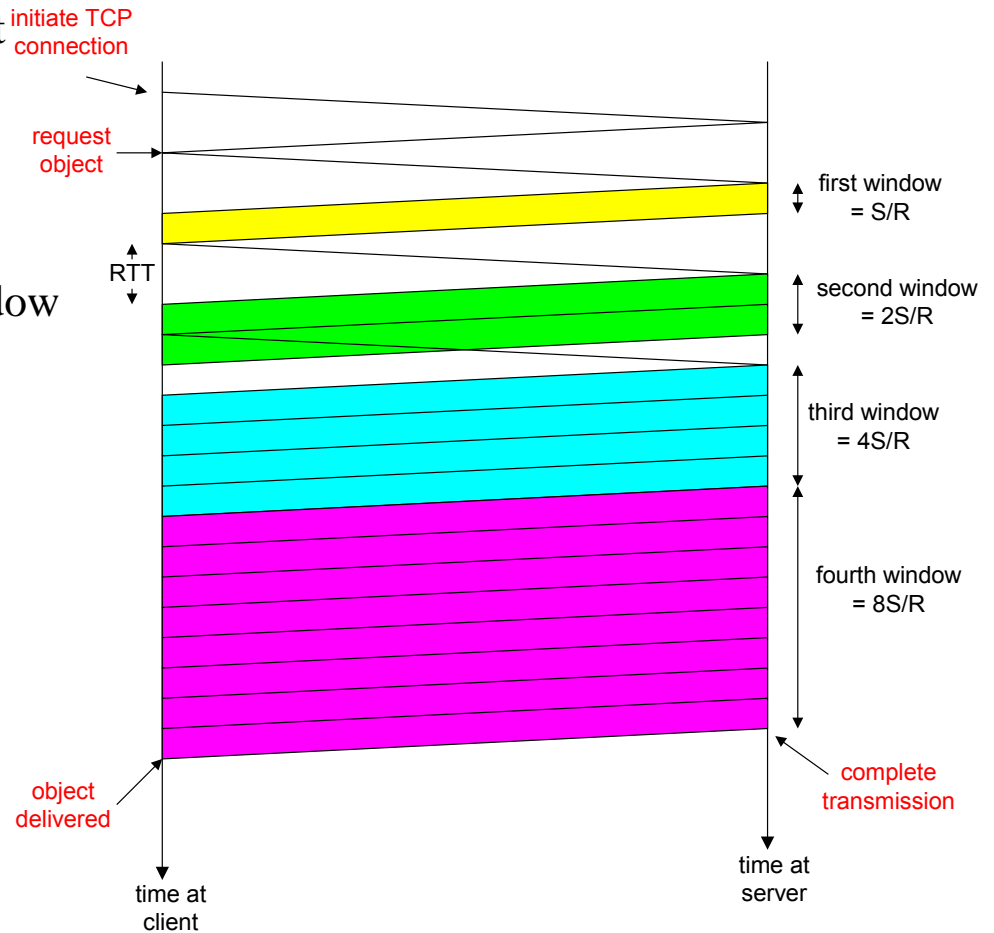
$\frac{S}{R} + RTT$  = time from when server starts to send segment

until server receives acknowledgement

$2^{k-1} \frac{S}{R}$  = time to transmit the  $k$ th window

$\left[ \frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right]^+$  = idle time after the  $k$ th window

$$\begin{aligned} \text{delay} &= \frac{O}{R} + 2RTT + \sum_{p=1}^P \text{idleTime}_p \\ &= \frac{O}{R} + 2RTT + \sum_{k=1}^P \left[ \frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right] \\ &= \frac{O}{R} + 2RTT + P \left[ RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R} \end{aligned}$$



# TCP Delay Modeling (4)

Recall  $K$  = number of windows that cover object

How do we calculate  $K$  ?

$$\begin{aligned} K &= \min\{k : 2^0 S + 2^1 S + \dots + 2^{k-1} S \geq O\} \\ &= \min\{k : 2^0 + 2^1 + \dots + 2^{k-1} \geq O/S\} \\ &= \min\{k : 2^k - 1 \geq \frac{O}{S}\} \\ &= \min\{k : k \geq \log_2(\frac{O}{S} + 1)\} \\ &= \left\lceil \log_2(\frac{O}{S} + 1) \right\rceil \end{aligned}$$

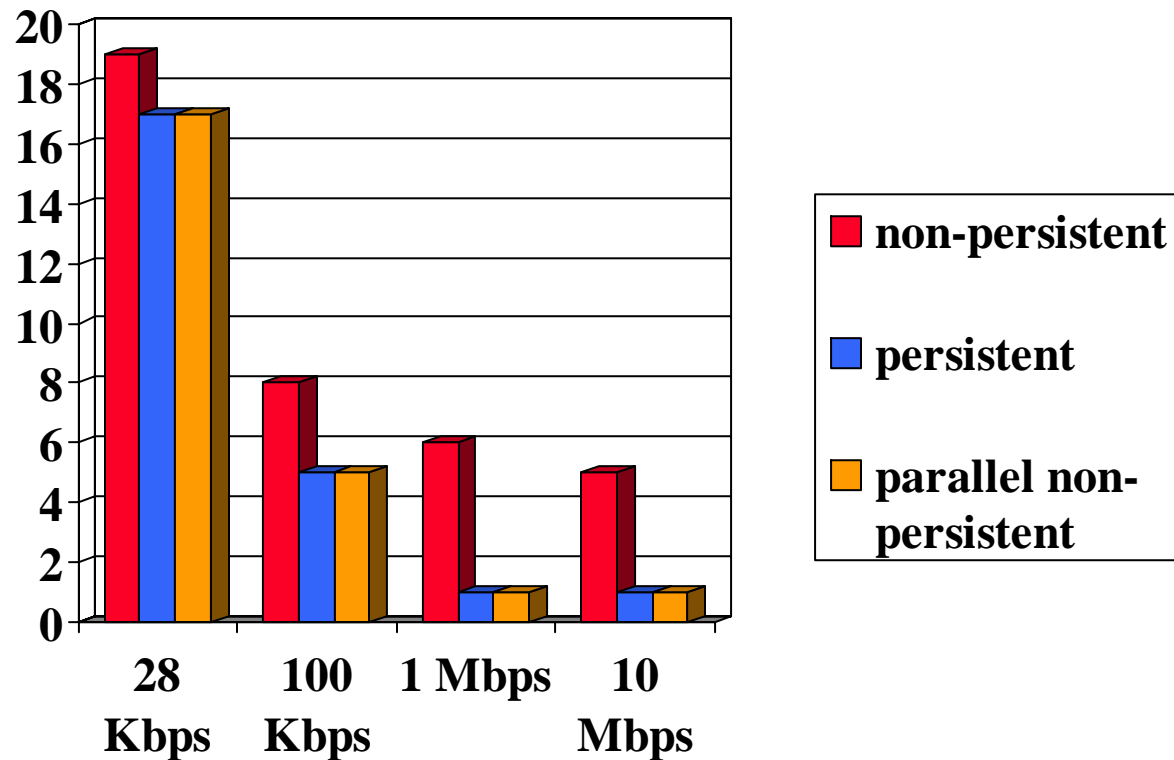
Calculation of  $Q$ , number of idles for infinite-size object, is similar (see HW).

# HTTP Modeling

- Assume Web page consists of:
  - 1 base HTML page (of size  $O$  bits)
  - $M$  images (each of size  $O$  bits)
- Non-persistent HTTP:
  - $M+1$  TCP connections in series
  - *Response time =  $(M+1)O/R + (M+1)2RTT + \text{sum of idle times}$*
- Persistent HTTP:
  - 2  $RTT$  to request and receive base HTML file
  - 1  $RTT$  to request and receive  $M$  images
  - *Response time =  $(M+1)O/R + 3RTT + \text{sum of idle times}$*
- Non-persistent HTTP with  $X$  parallel connections
  - Suppose  $M/X$  integer.
  - 1 TCP connection for base file
  - $M/X$  sets of parallel connections for images.
  - *Response time =  $(M+1)O/R + (M/X + 1)2RTT + \text{sum of idle times}$*

# HTTP Response time (in seconds)

RTT = 100 msec, O = 5 Kbytes, M=10 and X=5

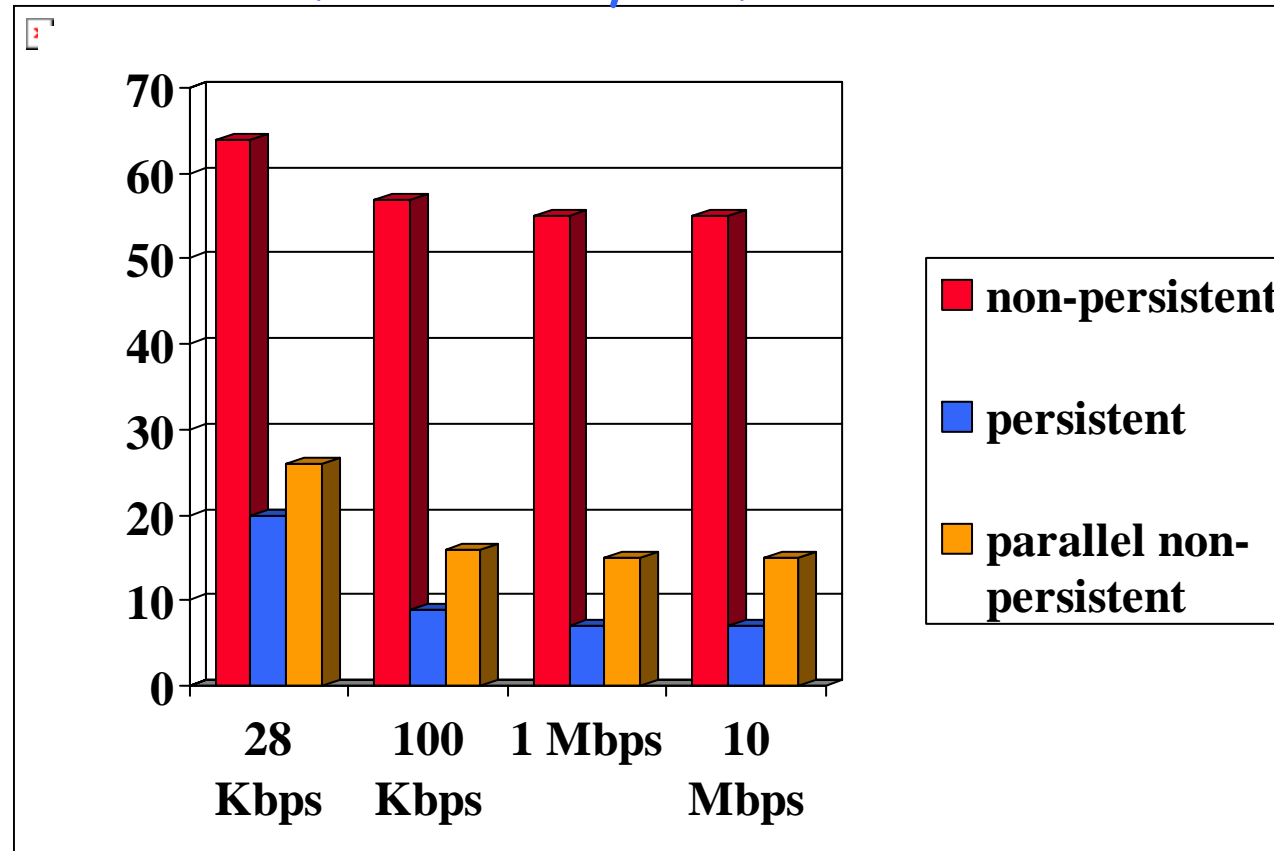


For low bandwidth, connection & response time dominated by transmission time.

Persistent connections only give minor improvement over parallel connections.

# HTTP Response time (in seconds)

RTT = 1 sec, O = 5 Kbytes, M=10 and X=5



For larger RTT, response time dominated by TCP establishment & slow start delays. Persistent connections now give important improvement: particularly in high delay•bandwidth networks.

# Issues to Think About

---

- What about short flows? (setting initial cwnd)
  - most flows are short
  - most bytes are in long flows
- How does this work over wireless links?
  - packet reordering fools fast retransmit
  - loss not always congestion related
- High speeds?
  - to reach 10gbps, packet losses occur every 90 minutes!
- Fairness: how do flows with different RTTs share link?

# Security issues with TCP

---

- Example attacks:
  - Sequence number spoofing
  - Routing attacks
  - Source address spoofing
  - Authentication attacks

# Network Layer

---

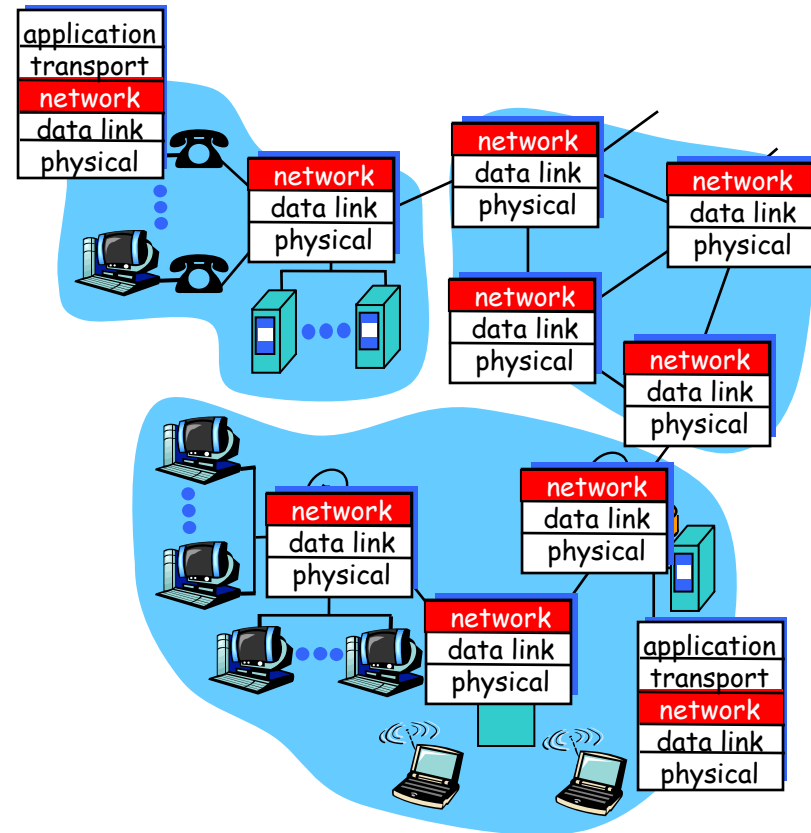
## goals:

- understand principles behind network layer services:
  - routing (path selection)
  - dealing with scale
  - how a router works
  - advanced topics: IPv6, mobility
- instantiation and implementation in the Internet



# Network layer

- transport segment from sending to receiving host
- on sending side encapsulates segments into datagrams
- on rcving side, delivers segments to transport layer
- network layer protocols in *every* host, router
- Router examines header fields in all IP datagrams passing through it



# Key Network-Layer Functions

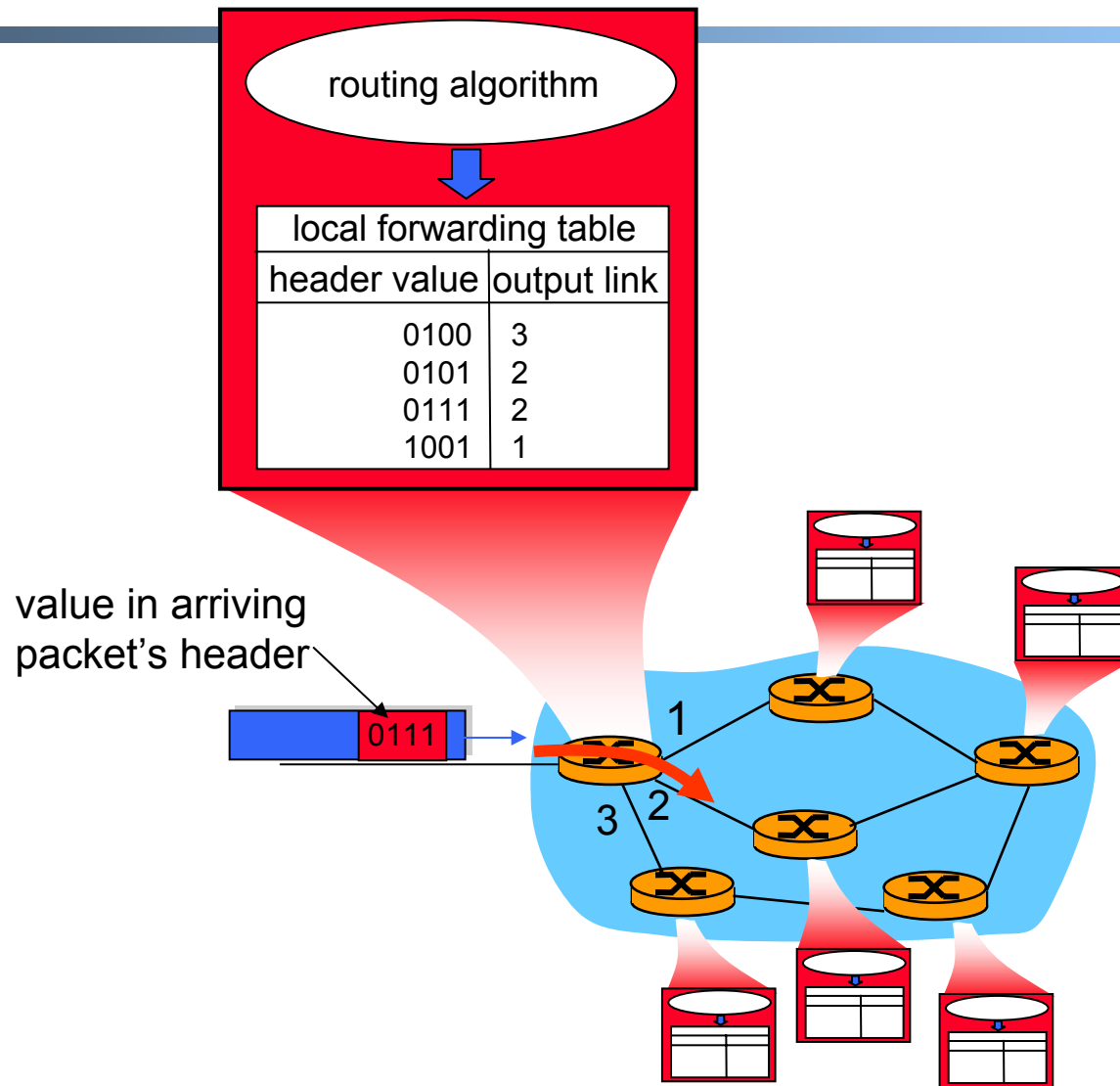
---

- *forwarding*: move packets from router's input to appropriate router output
- *routing*: determine route taken by packets from source to dest.
  - *Routing algorithms*

## analogy:

- *routing*: process of planning trip from source to dest
- *forwarding*: process of getting through single interchange

# Interplay between routing and forwarding



# Connection setup

---

- 3<sup>rd</sup> important function in *some* network architectures:
  - ATM, frame relay, X.25
- Before datagrams flow, two hosts and intervening routers establish virtual connection
  - Routers get involved
- Network and transport layer cnctn service:
  - **Network:** between two hosts
  - **Transport:** between two processes

# Network service model

---

**Q:** What *service model* for “channel” transporting datagrams from sender to rcvr?

## Example services for individual datagrams:

- guaranteed delivery
- Guaranteed delivery with less than 40 msec delay

## Example services for a flow of datagrams:

- In-order datagram delivery
- Guaranteed minimum bandwidth to flow
- Restrictions on changes in inter-packet spacing

## Network layer service models:

Network Architecture	Service Model	Guarantees ?				Congestion feedback
		Bandwidth	Loss	Order	Timing	
Internet	best effort	none	no	no	no	no (inferred via loss)
ATM	CBR	constant rate	yes	yes	yes	no congestion
ATM	VBR	guaranteed rate	yes	yes	yes	no congestion
ATM	ABR	guaranteed minimum	no	yes	no	yes
ATM	UBR	none	no	yes	no	no

# Network layer connection and connection-less service

---

- Datagram network provides network-layer connectionless service
- VC network provides network-layer connection service
- Analogous to the transport-layer services, but:
  - **Service:** host-to-host
  - **No choice:** network provides one or the other
  - **Implementation:** in the core

# Virtual circuits

“source-to-dest path behaves much like telephone circuit”

- performance-wise
- network actions along source-to-dest path

- call setup, teardown for each call *before* data can flow
- each packet carries VC identifier (not destination host address)
- *every* router on source-dest path maintains “state” for each passing connection
- link, router resources (bandwidth, buffers) may be *allocated* to VC



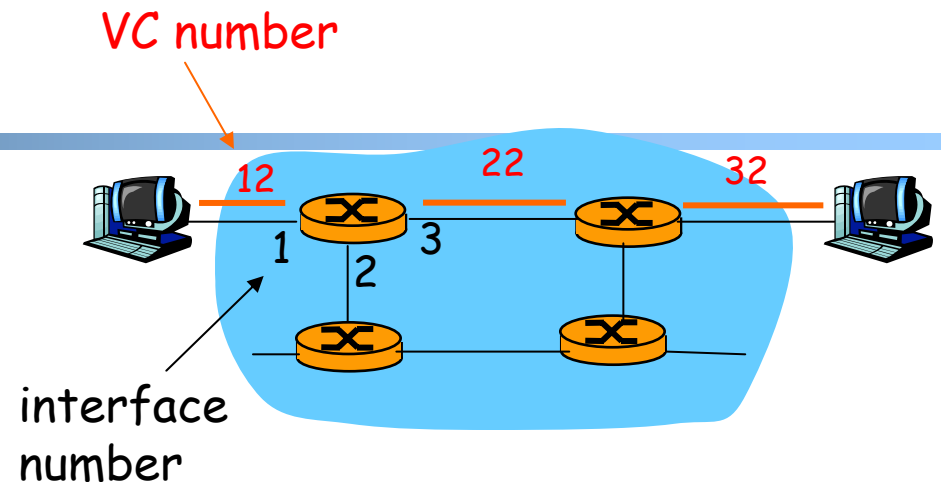
# VC implementation

---

A VC consists of:

1. Path from source to destination
  2. VC numbers, one number for each link along path
  3. Entries in forwarding tables in routers along path
- Packet belonging to VC carries a VC number.
  - VC number must be changed on each link.
    - New VC number comes from forwarding table

# Forwarding table



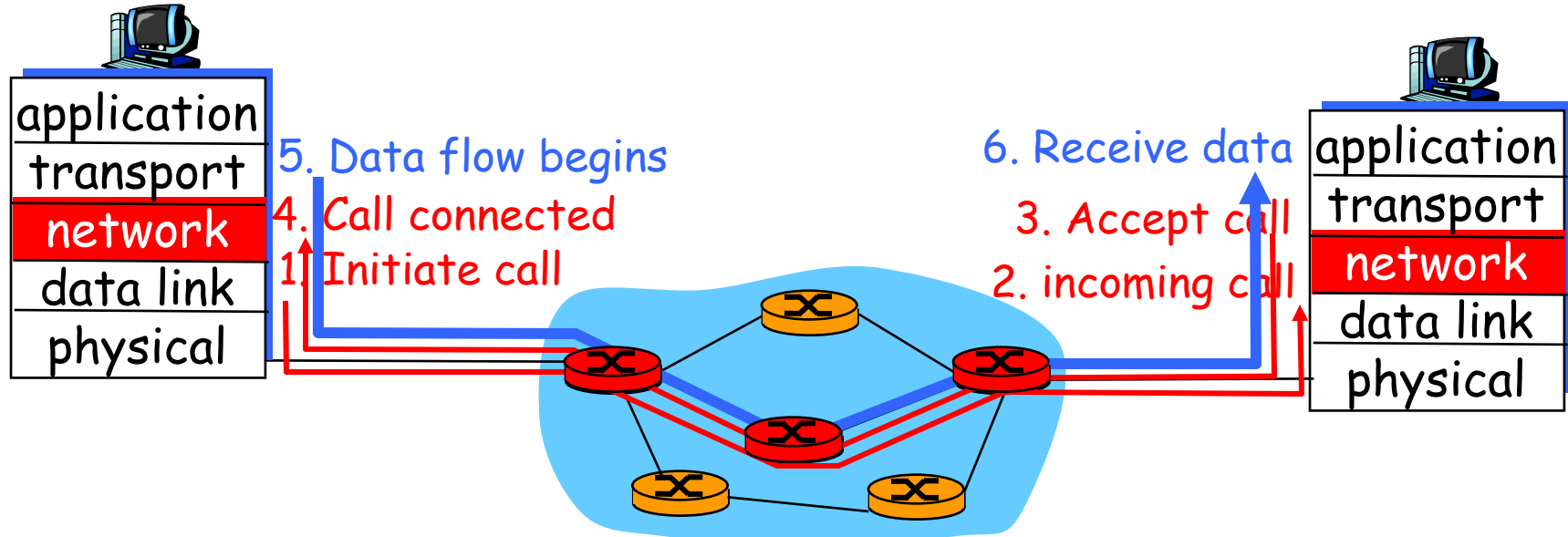
Forwarding table in northwest router:

Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	2	22
2	63	1	18
3	7	2	17
1	97	3	87
...	...	...	...

**Routers maintain connection state information!**

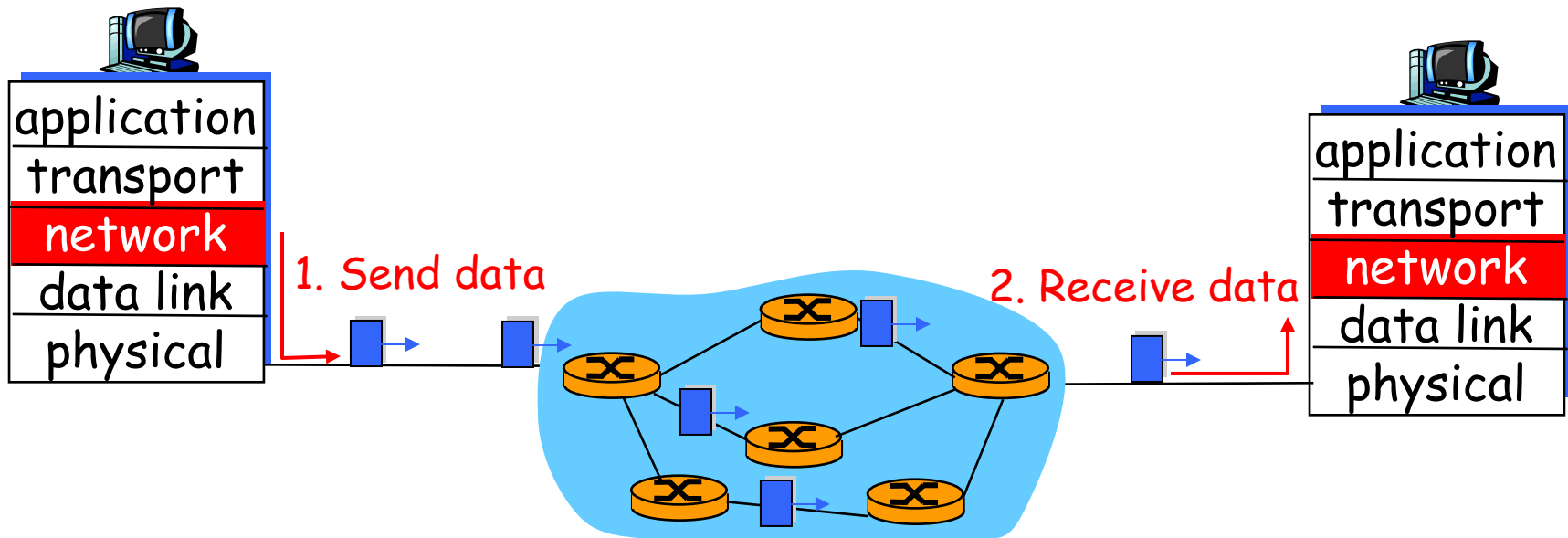
# Virtual circuits: signaling protocols

- used to setup, maintain teardown VC
- used in ATM, frame-relay, X.25
- not used in today's Internet



# Datagram networks

- no call setup at network layer
- routers: no state about end-to-end connections
  - no network-level concept of “connection”
- packets forwarded using destination host address
  - packets between same source-dest pair may take different paths



# Forwarding table

4 billion  
possible entries

<u>Destination Address Range</u>	<u>Link Interface</u>
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111	0
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 11111111	1
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 11111111	2
otherwise	3

# Longest prefix matching

---

<u>Prefix Match</u>	<u>Link Interface</u>
11001000 00010111 00010	0
11001000 00010111 00011000	1
11001000 00010111 00011	2
otherwise	3

## Examples

DA: 11001000 00010111 00010110 10100001

Which interface?

DA: 11001000 00010111 00011000 10101010

Which interface?

# Datagram or VC network: why?

---

## Internet

- data exchange among computers
  - “elastic” service, no strict timing req.
- “smart” end systems (computers)
  - can adapt, perform control, error recovery
  - simple inside network, complexity at “edge”
- many link types
  - different characteristics
  - uniform service difficult

## ATM

- evolved from telephony
- human conversation:
  - strict timing, reliability requirements
  - need for guaranteed service
- “dumb” end systems
  - telephones
  - complexity inside network