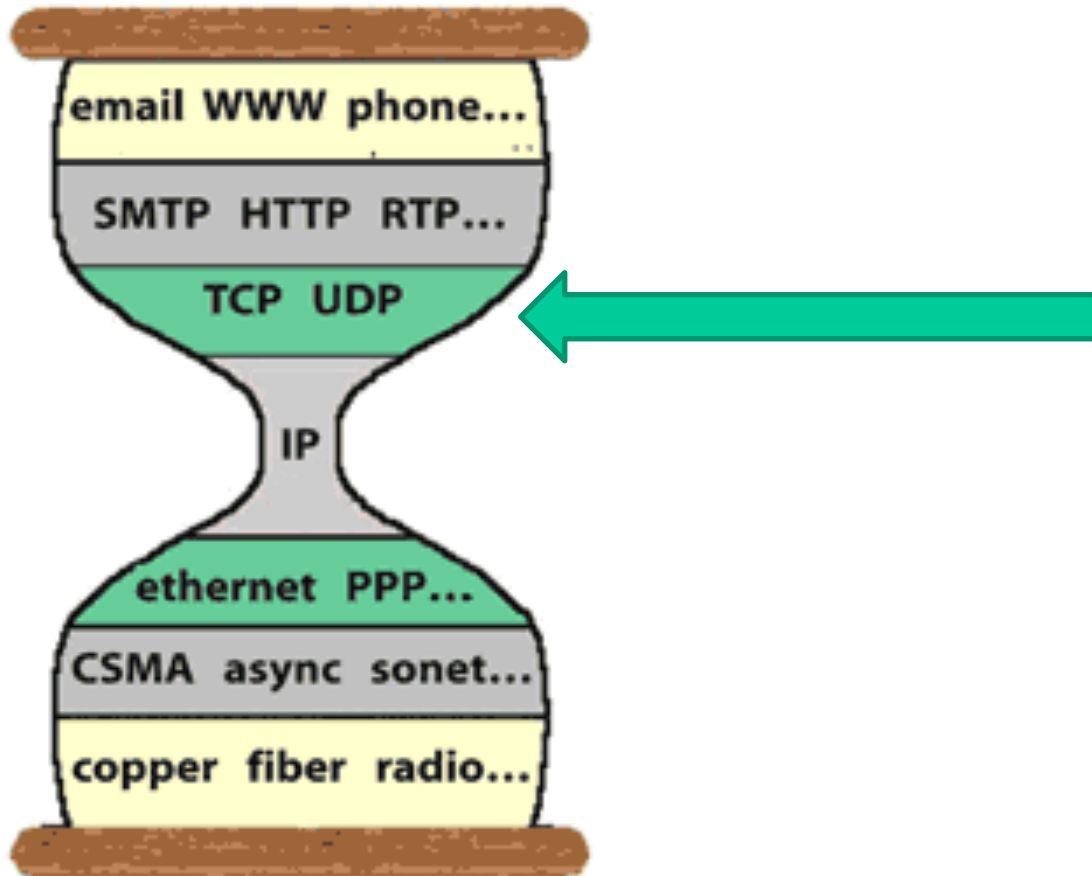


Wrapping Up Congestion Control

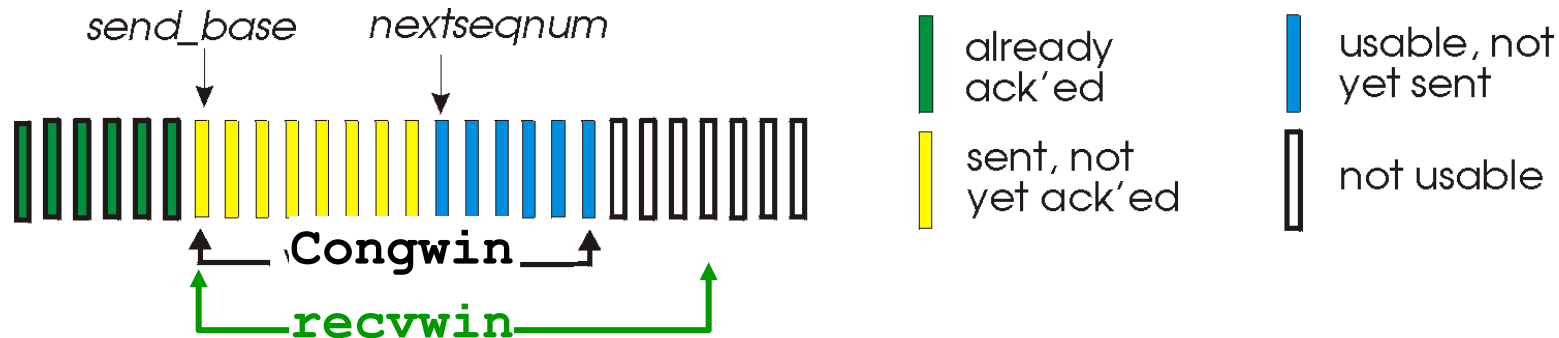
Chapter 3

3.5 TCP Congestion Control



TCP Congestion Control

- ◆ Add a “congestion control window” **cwnd** on top of flow-control window
- ◆ Sender limits: **LastByteSent - LastByteAcked** ≤ **cwnd**



- ◆ **cwnd** initialized to 1 packet, increases until congestion
 - How sender infers congestion: packet loss (timeout, or 3 dup. ACKs)
- ◆ Upon loss: *decrease* **cwnd**, then begin increasing again
 - Two phases: (1) **slow start**, (2) **congestion avoidance**
 - a threshold (**ssthresh**) defines the boundary between the two

“Slow” Start

♦ Objective: quickly gauge the pipeline size

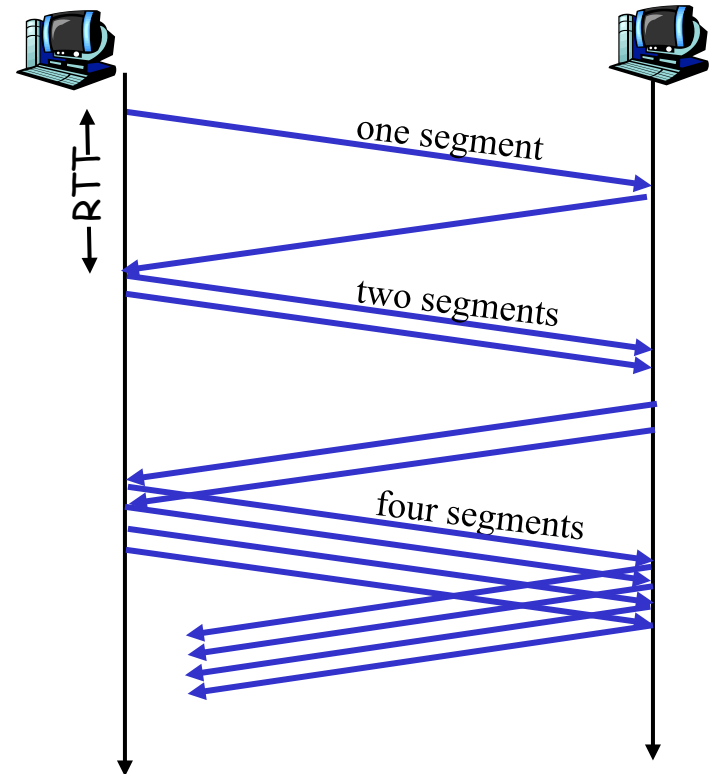
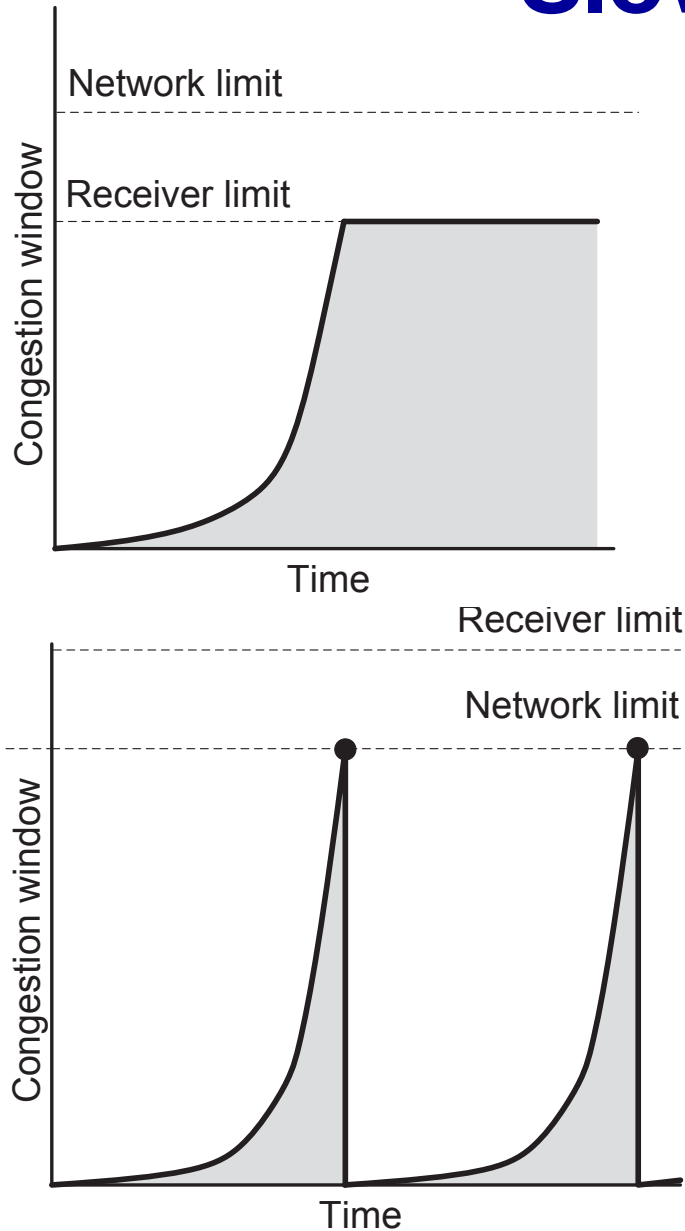
1. Start with $\text{cwnd} = 1 \text{ MSS}$
2. Send cwnd packets
3. If all packets got acked
 - $\text{cwnd} = 2 * \text{cwnd}$
 - goto 2
4. Else have gone too far
 - $\text{ss-thresh} = \text{cwnd} / 2$
 - $\text{cwnd} = 1 \text{ MSS}$
 - goto 2

Same, but using “self-clocking” of TCP

1. Start with $\text{cwnd} = 1 \text{ MSS}$
2. Send cwnd packets
3. If ack
 - $\text{cwnd} = \text{cwnd} + 1 \text{ MSS}$
 - send packets
4. If timeout
 - $\text{sshthresh} = \text{cwnd} / 2$
 - $\text{cwnd} = 1 \text{ MSS}$
 - goto 2

Multiplicative **I**ncrease, Reset to 1 **D**ecrease

Slow Start



Congestion Avoidance

- ◆ Objective: maintain steady state, probe for unused resources, avoid conflicts
- ◆ Send cwnd packets
- ◆ If all sent packets got ack
 - $\text{cwnd} = \text{cwnd} + 1 \text{ MSS}$
- ◆ Else
 - $\text{cwnd} = \text{cwnd} / 2$
- ◆ Additive Increase, Multiplicative Decrease (AIMD)

Send cwnd packets

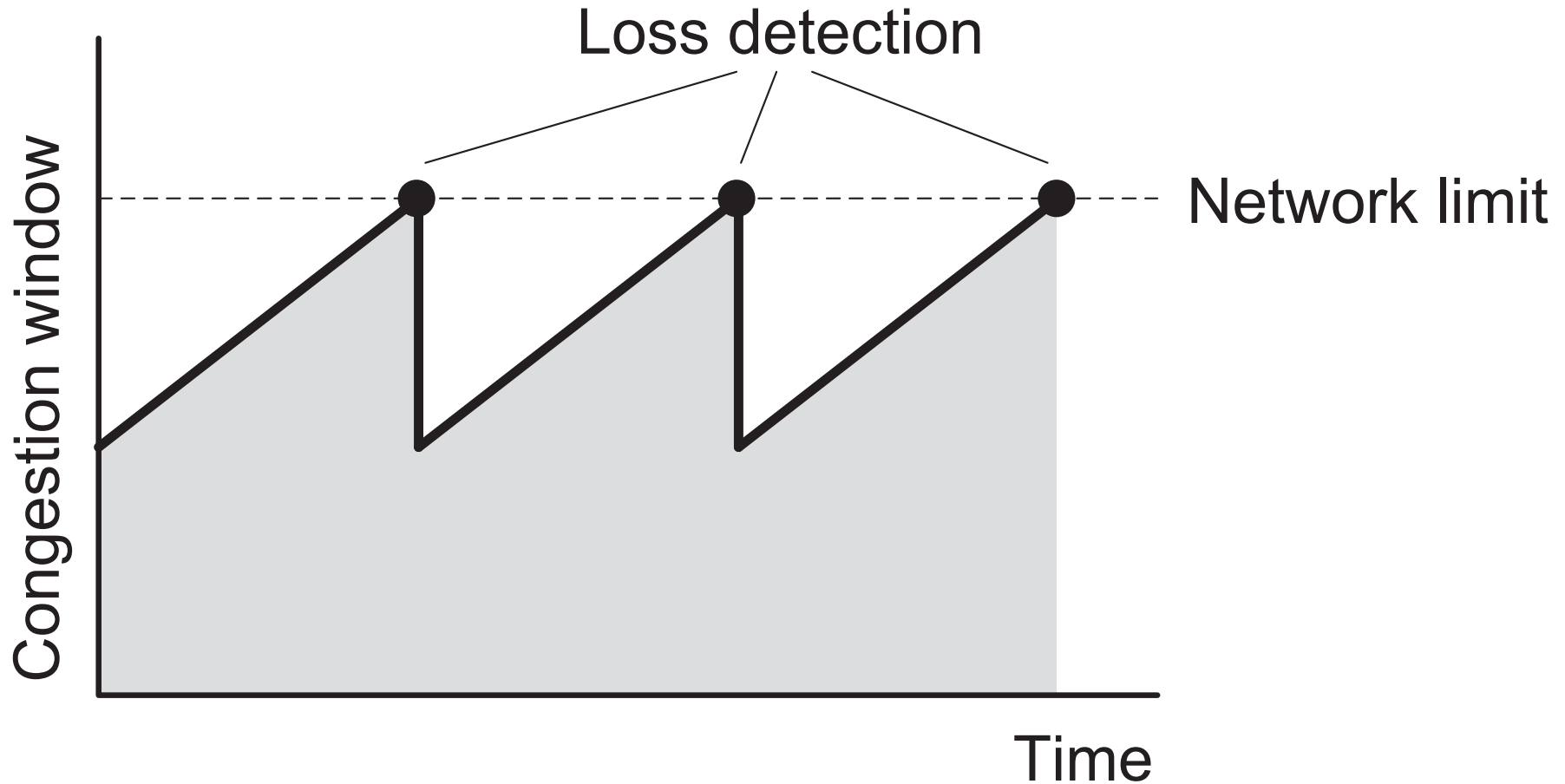
If ack

- $\text{cwnd} = \text{cwnd} + (1 \text{ MSS})^2 / \text{cwnd}$

If timeout

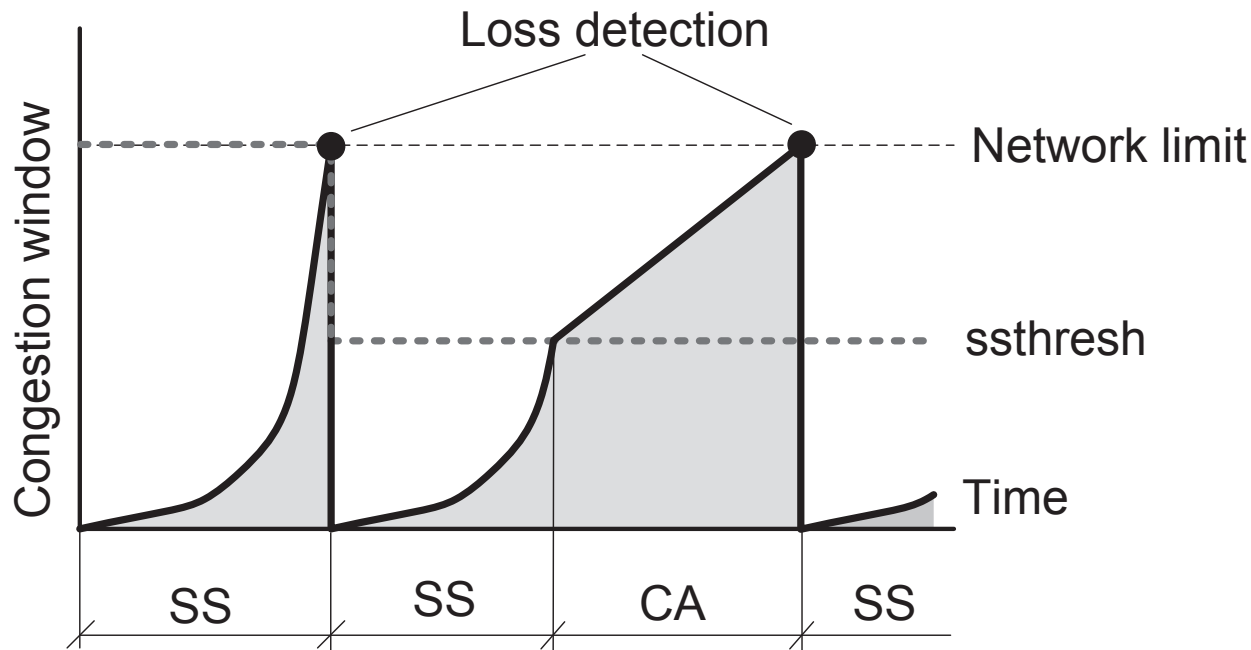
- $\text{cwnd} = \text{cwnd} / 2$

Congestion Avoidance



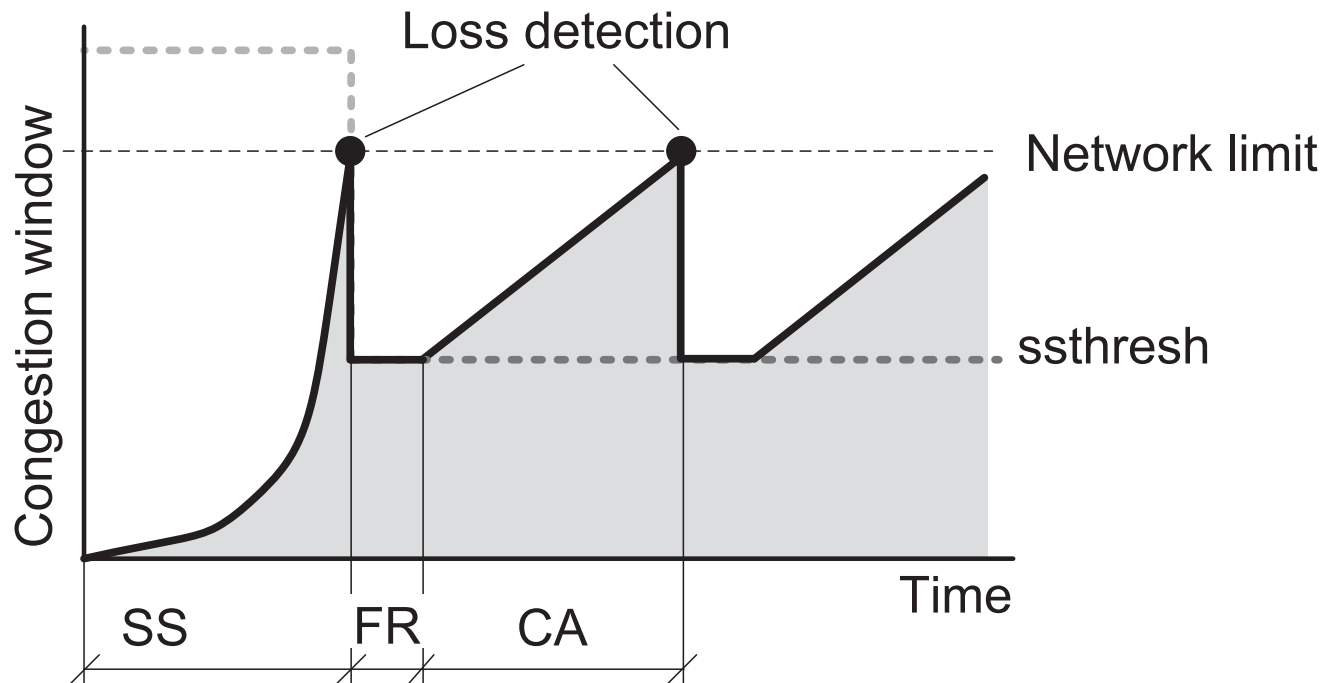
Combining Slow Start with Congestion Avoidance (Tahoe)

- ◆ Set initial **ss-thresh**
- ◆ When **cwnd** < **ss-thresh**, use Slow Start
 - ssthresh will get updated
- ◆ when **cwnd** ≥ **ss-thresh**, use Congestion Avoidance

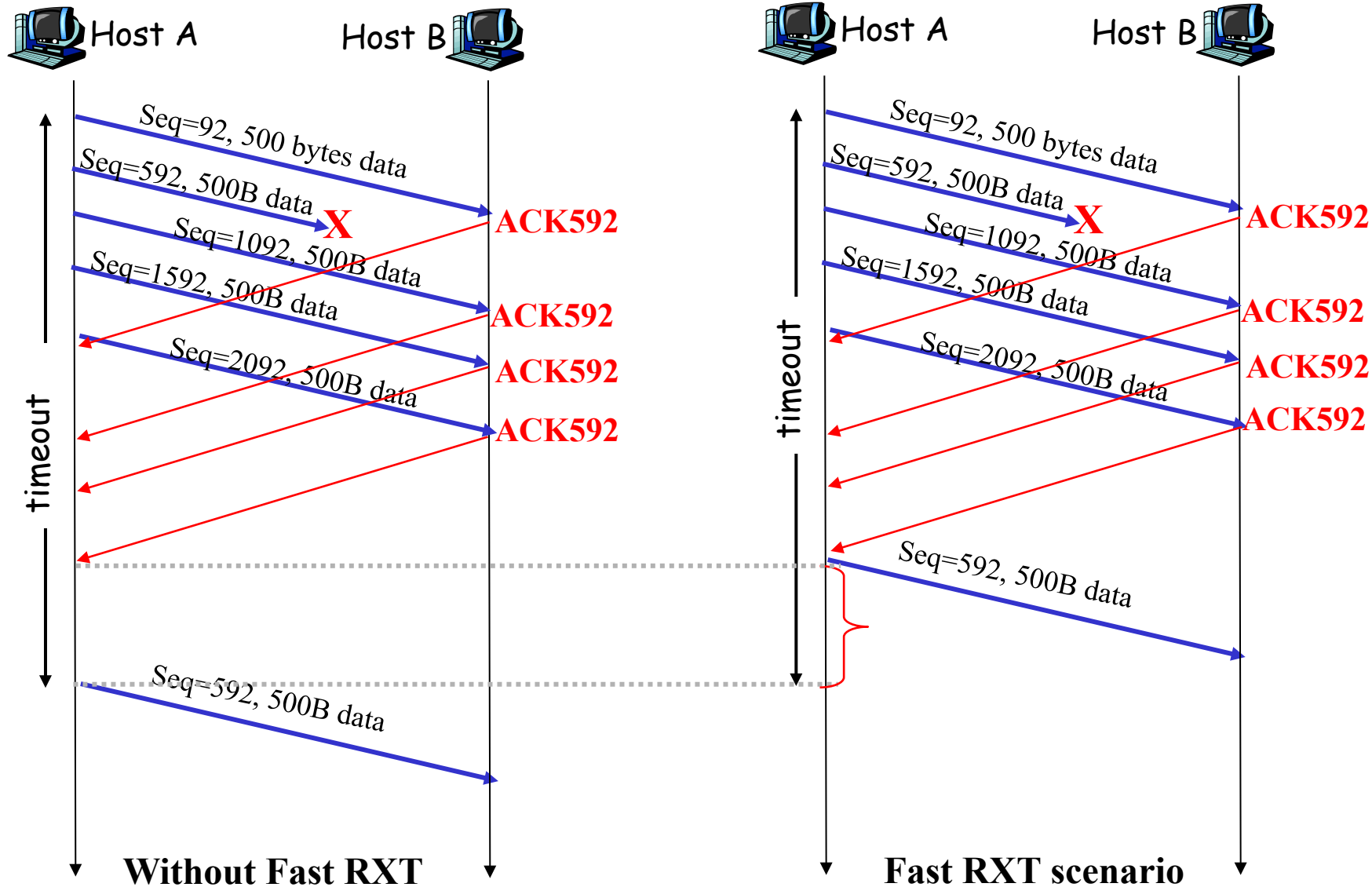


Slow Start and Congestion Avoidance (Reno)

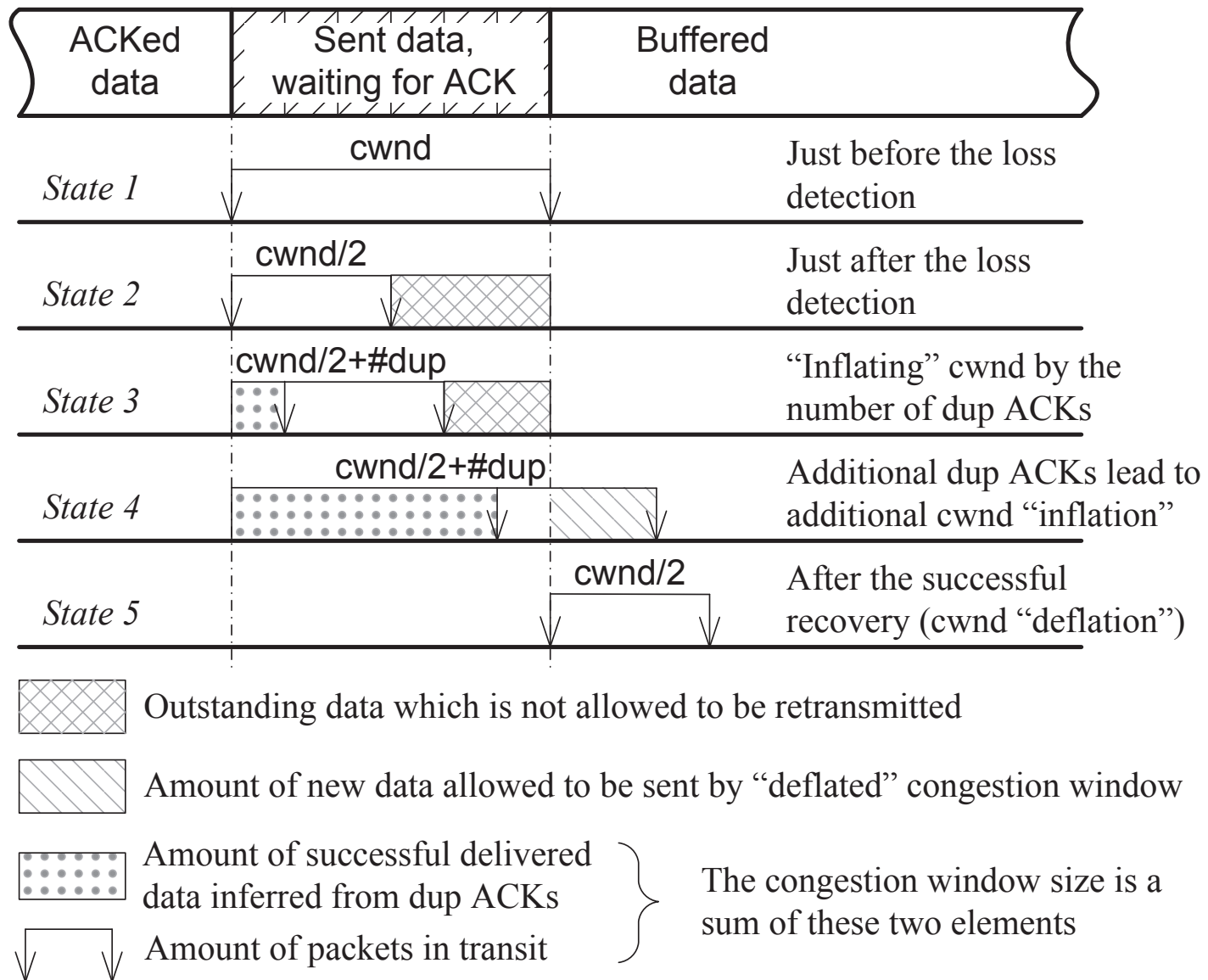
- ◆ If just one packet is lost, dropping cwnd to 1 is an overreaction
- ◆ What if loss detected through 3 dup acks, we just reduce cwnd by half, and quickly recover from the loss (Fast Recovery)



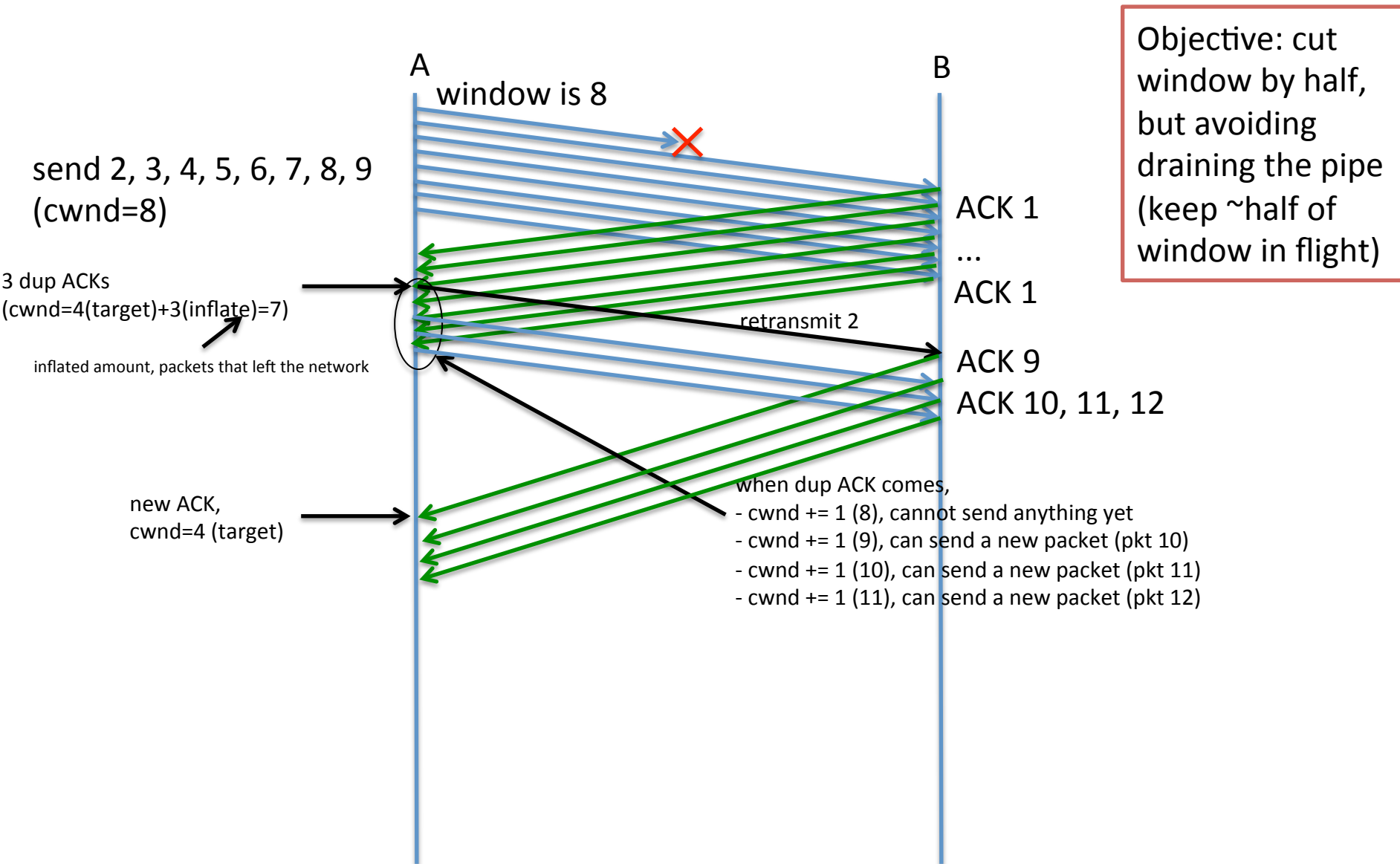
TCP fast retransmit example



Fast Recovery/Retransmit (Reno)



Fast Retransmit / Fast Recovery



Is TCP fair?

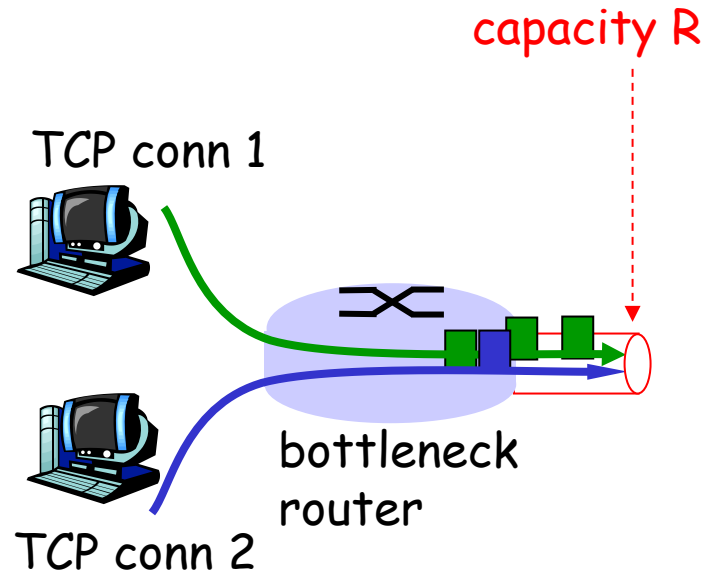
Fairness: if N TCP sessions share same bottleneck link, each should get 1/N of link capacity

Jain's fairness index

n is the number of users sharing the path

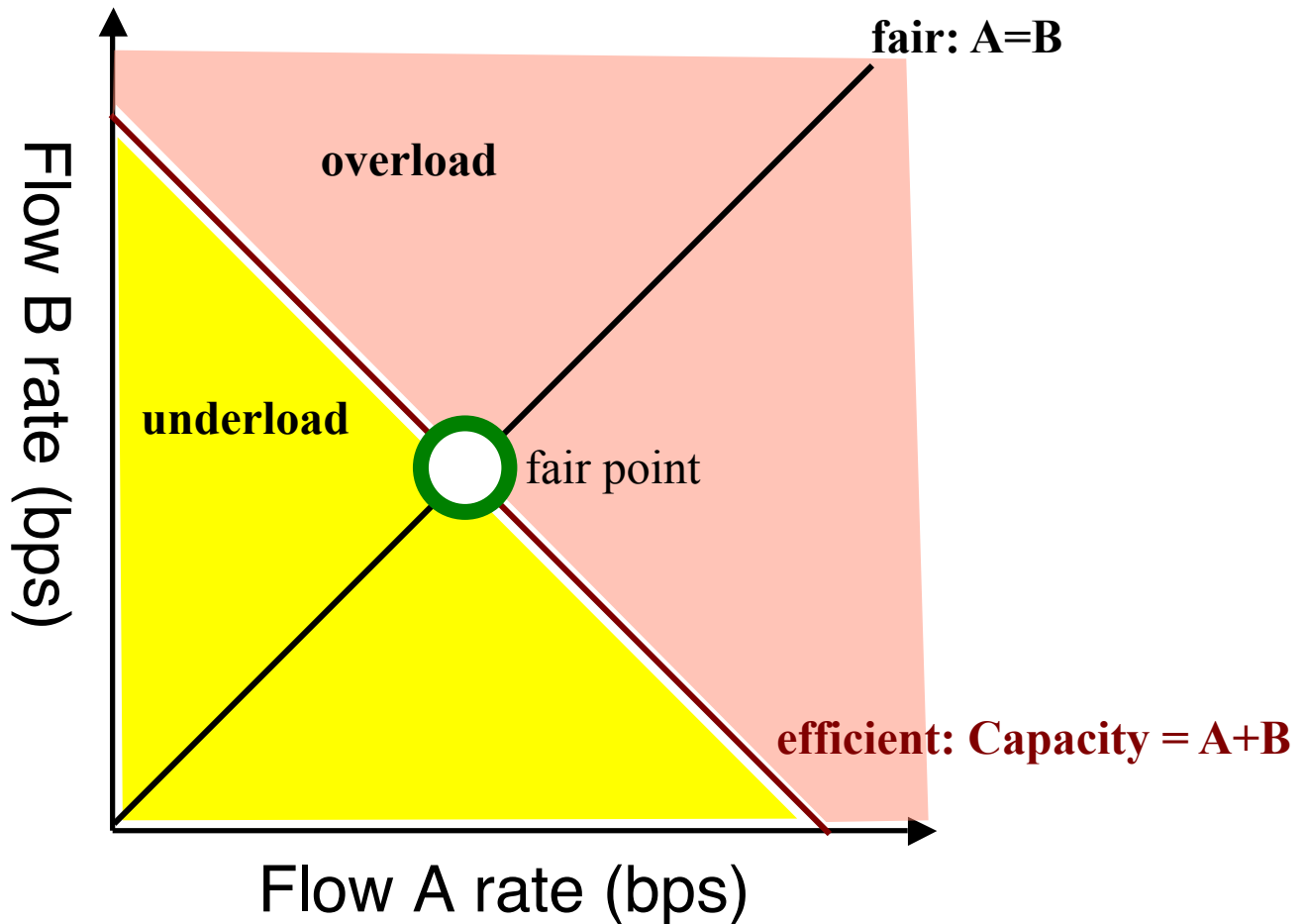
f_i is the network share of i^{th} user

$$F = \frac{(\sum_i^n f_i)^2}{n \cdot \sum_i f_i^2}$$

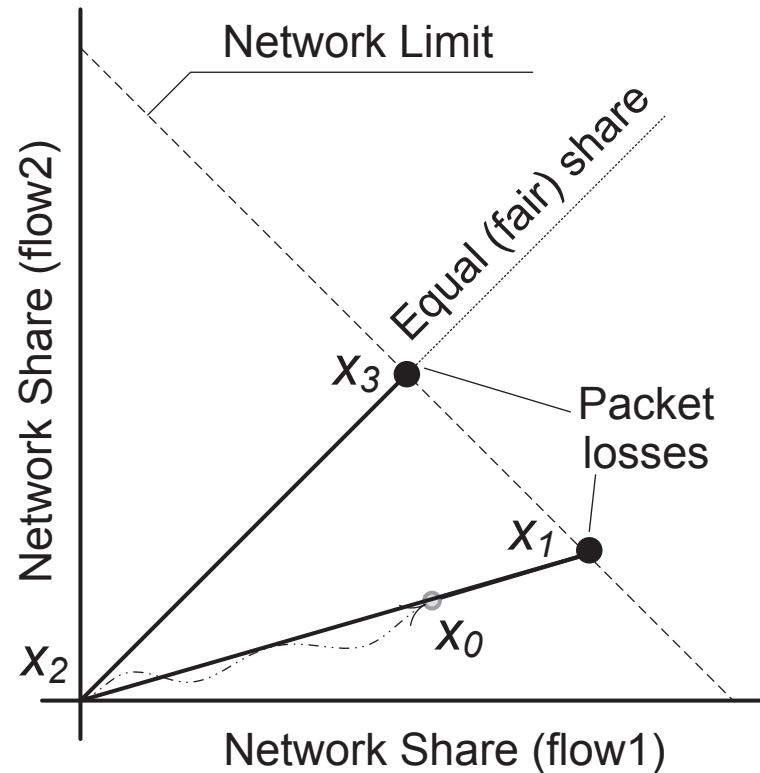


bigger cwnd, larger share TCP flow can “take”

Chiu Jain Phase Plots



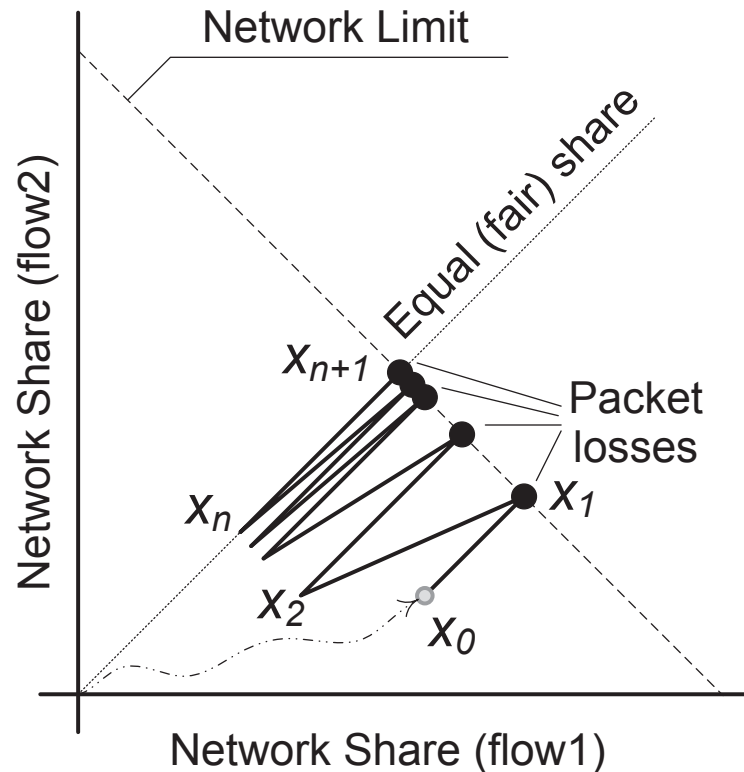
Is Slow Start of TCP Fair?



$x_0 \rightarrow x_1, \dots, x_n \rightarrow x_{n+1}$ multiplicative increase (both flows have the same increase rate of their congestion windows)

$x_1 \rightarrow x_2$ equalization of the congestion window sizes

Is Congestion Avoidance Is Fair?



$x_0 - x_1, \dots, x_n - x_{n+1}$ additive increase (both flows have the same increase rate of their congestion windows)

$x_1 - x_2, \dots, x_{n-1} - x_n$ multiplicative decrease (a flow with the larger congestion window decreases more than a flow with the smaller)

Summary: TCP sender congestion control

State	Event	TCP Sender Action	Commentary
Slow Start (SS)	Received ACK 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)	Received ACK 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 3 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 TCP throughput as a function of window size and RTT?
- ◆ Ignore slow start: let W = window-size when loss occurs
 - When window is W : throughput = W / RTT
 - Just after loss

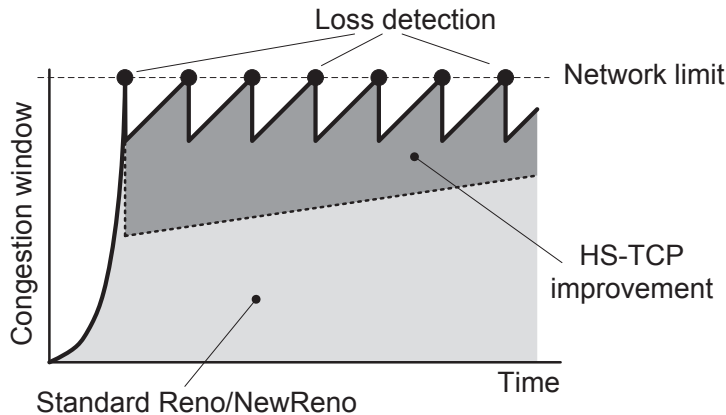
window $\rightarrow W/2$, throughput $\rightarrow W/2\text{RTT}$
 - Average throughput: $0.75 W/\text{RTT}$

Why Do We Need Other TCP Variants?

TCP Throughput (Mbps)	RTTs Between Losses	W	P
1	5.5	8.3	0.02
10	55.5	83.3	0.0002
100	555.5	833.3	0.000002
1000	5555.5	8333.3	0.00000002
10000	55555.5	83333.3	0.0000000002

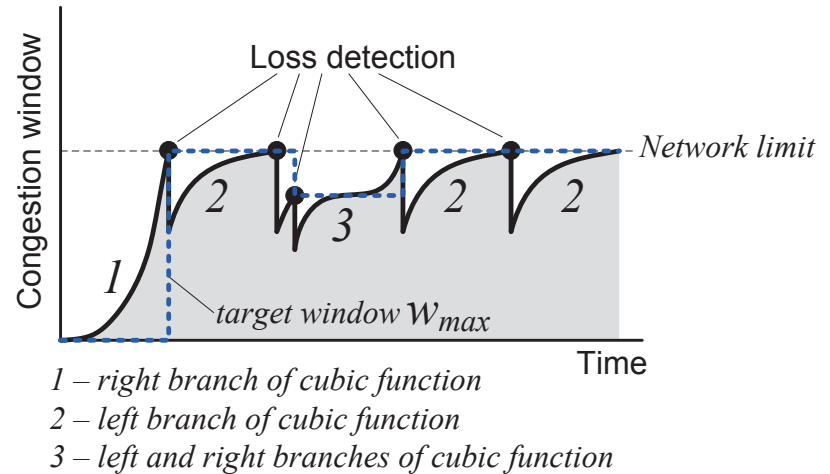
- ◆ HS TCP (High-Speed TCP)
 - OS X, available in Linux
- ◆ C-TCP (Compound TCP)
 - default on Windows, available in Linux
- ◆ CUBIC TCP
 - default in Linux

HS-TCP and CUBIC TCP



HS-TCP

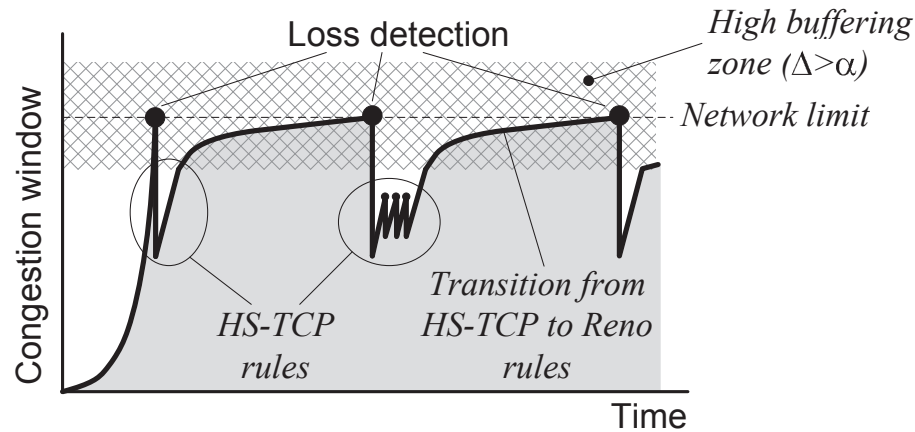
AIMD, but with increased additive increase and decreased multiplicative decrease



CUBIC-TCP

AIMD, but with increased additive increase as a cubic function (fast at first, slow later) and decreased multiplicative decrease

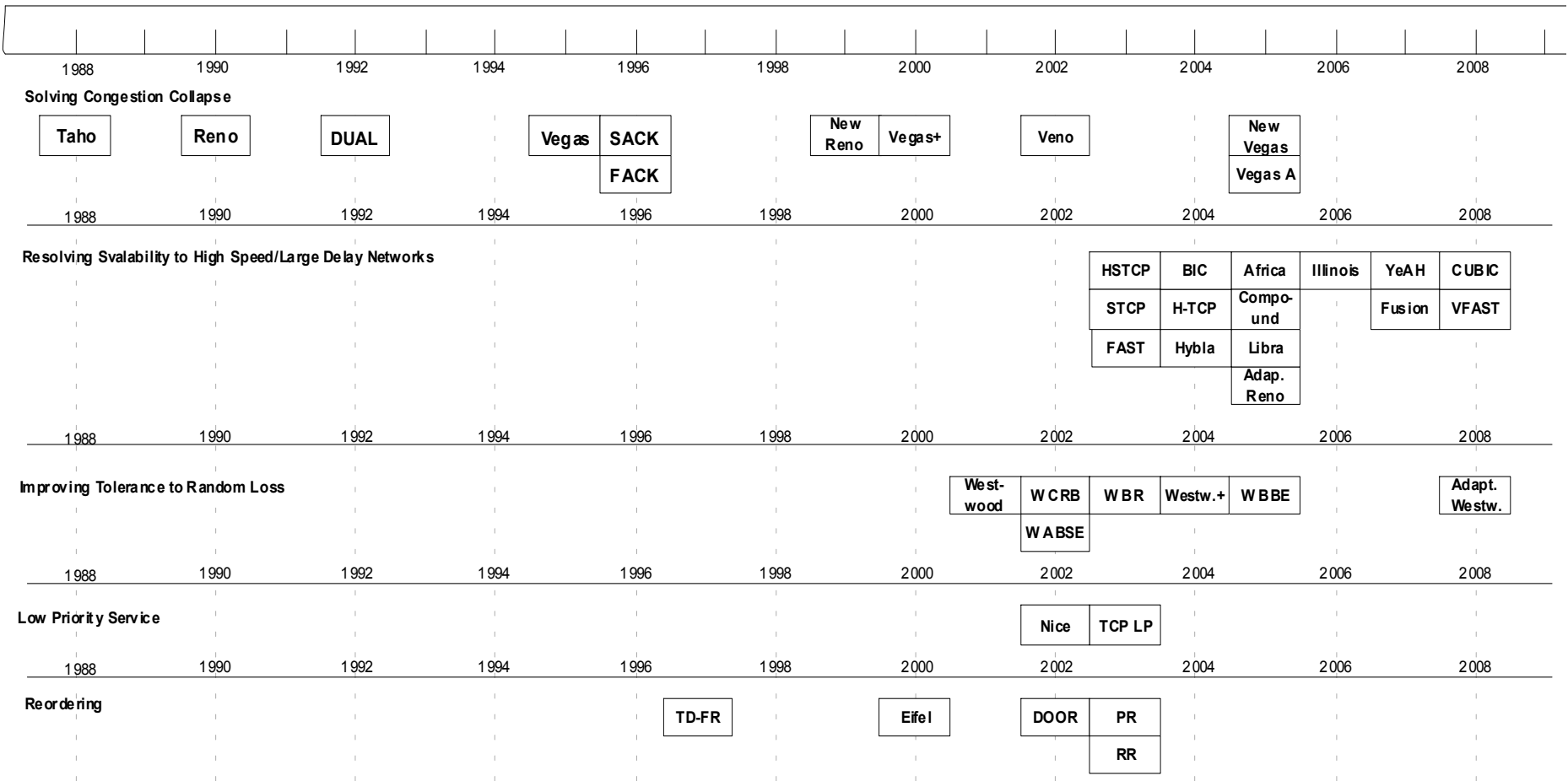
Compound TCP



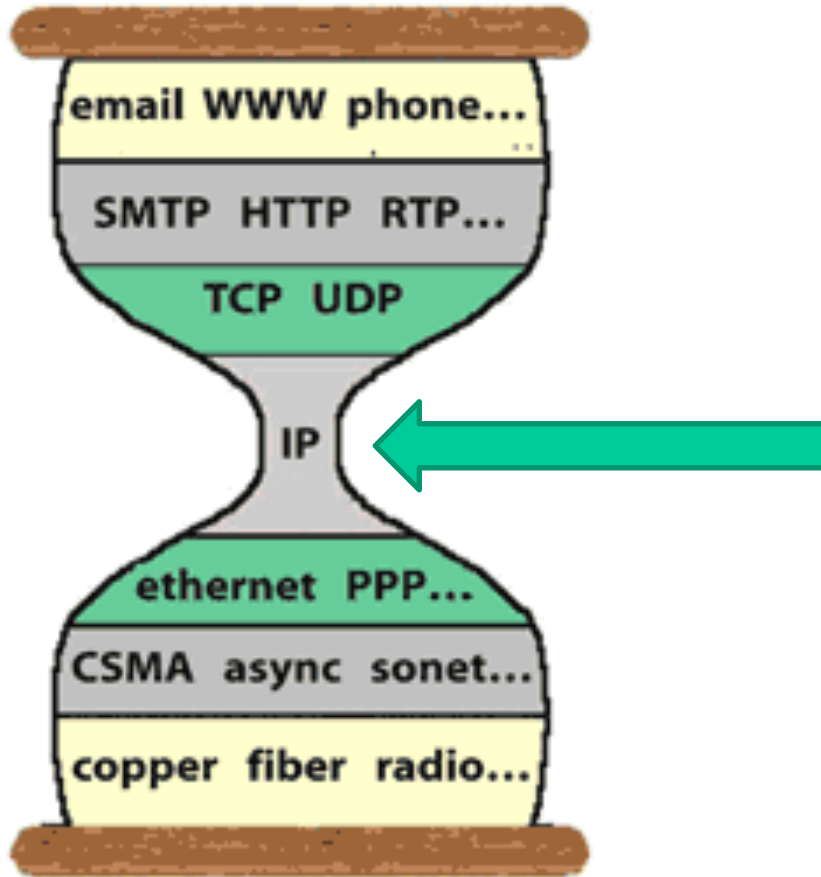
C-TCP

AIMD with multiple zones of different coefficients, depending on estimated buffering in the network

TCP Variants Timeline



Chapter 4: Network Layer



4.1 Introduction

4.2 Virtual circuit and datagram networks

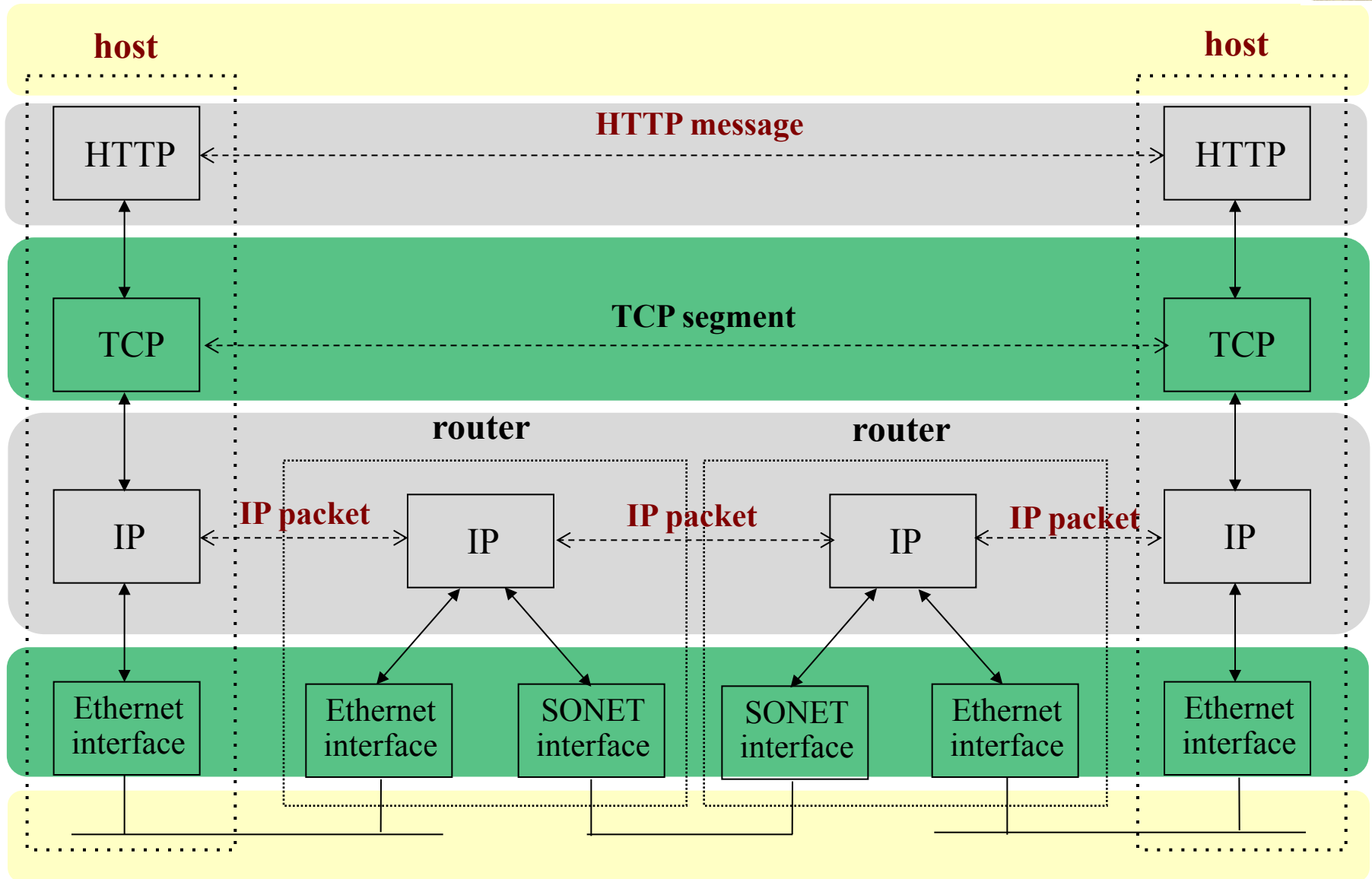
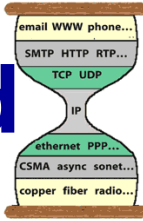
4.3 What's inside a router

4.4 IP: Internet Protocol

- IP datagram format
- IPv4 addressing
- ICMP
- IPv6

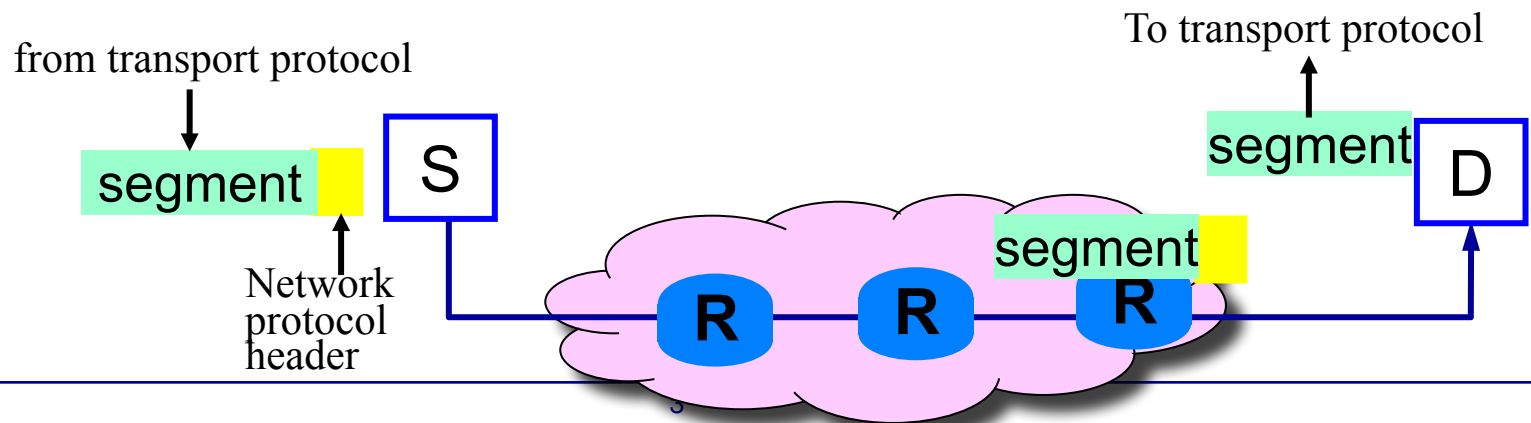
This and next week

Always keep the big picture in mind



Network layer

- ◆ Network layer protocols: in *every* host and router
- ◆ Moving data *segments* from sending to receiving host
 - **Routing**: calculate the best path to each destination
 - **Forwarding**: use FIB to move packets towards destinations
- ◆ Source host: encapsulates segments into *packets*
- ◆ Destination host: delivers segments to transport layer
- ◆ Each router along the path: examines the header field of a packet to determine where to forward it



Network layer:

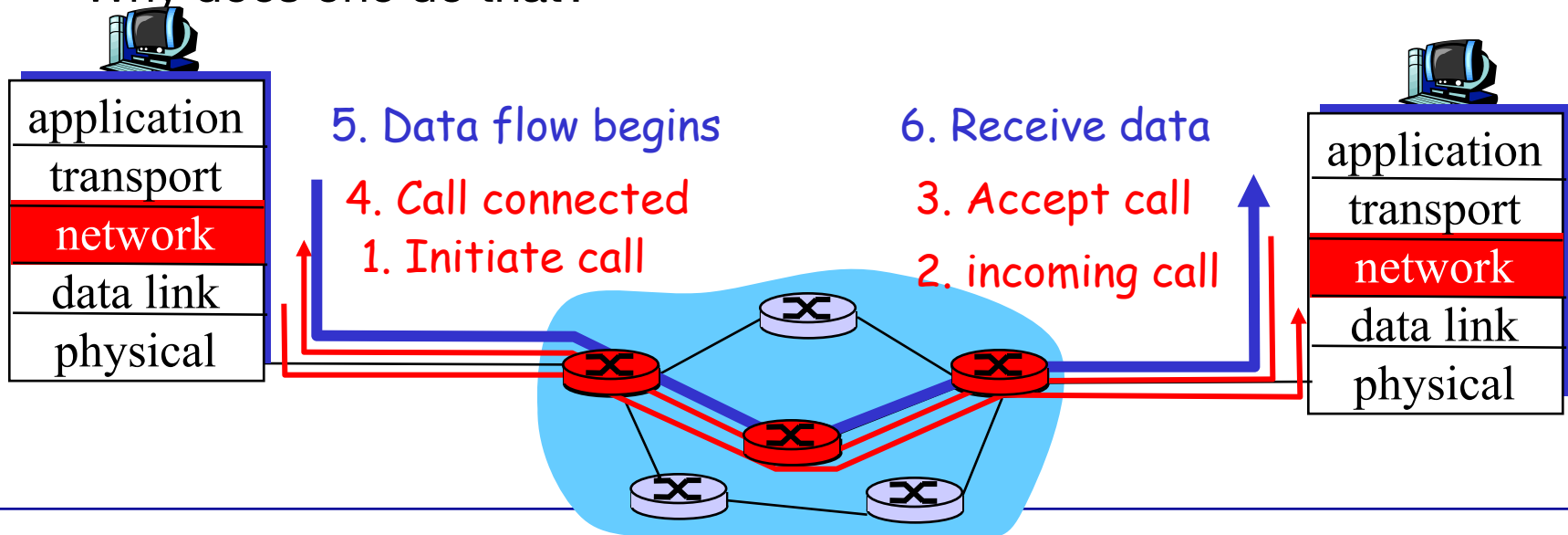
Connection vs. connection-less delivery

- ◆ *Virtual Circuit*: network sets up a connection, then delivers packets over the connection
 - works in a way very much like telephone circuit
- ◆ *Datagram*: Each packet finds its way to destination independently
- ◆ analogous to TCP vs. UDP at transport-layer, *except* that a given network provides one or the other, but not both (as in transport layer)

Virtual circuit (VC) Network

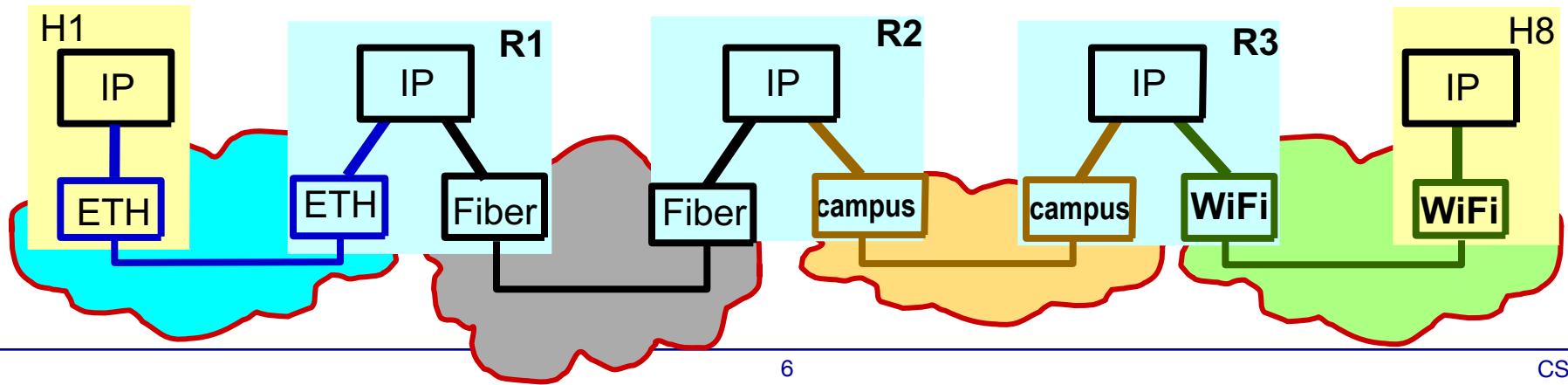
FYI

- ◆ Use a signaling protocol to setup connection first
- ◆ every router on the path maintains “state” for each passing virtual connection
 - link, router resources (bandwidth, buffers) allocated to the VC
- ◆ each packet carries VC identifier (*not destination host address*)
- ◆ VC number changes on each link
 - Why does one do that?

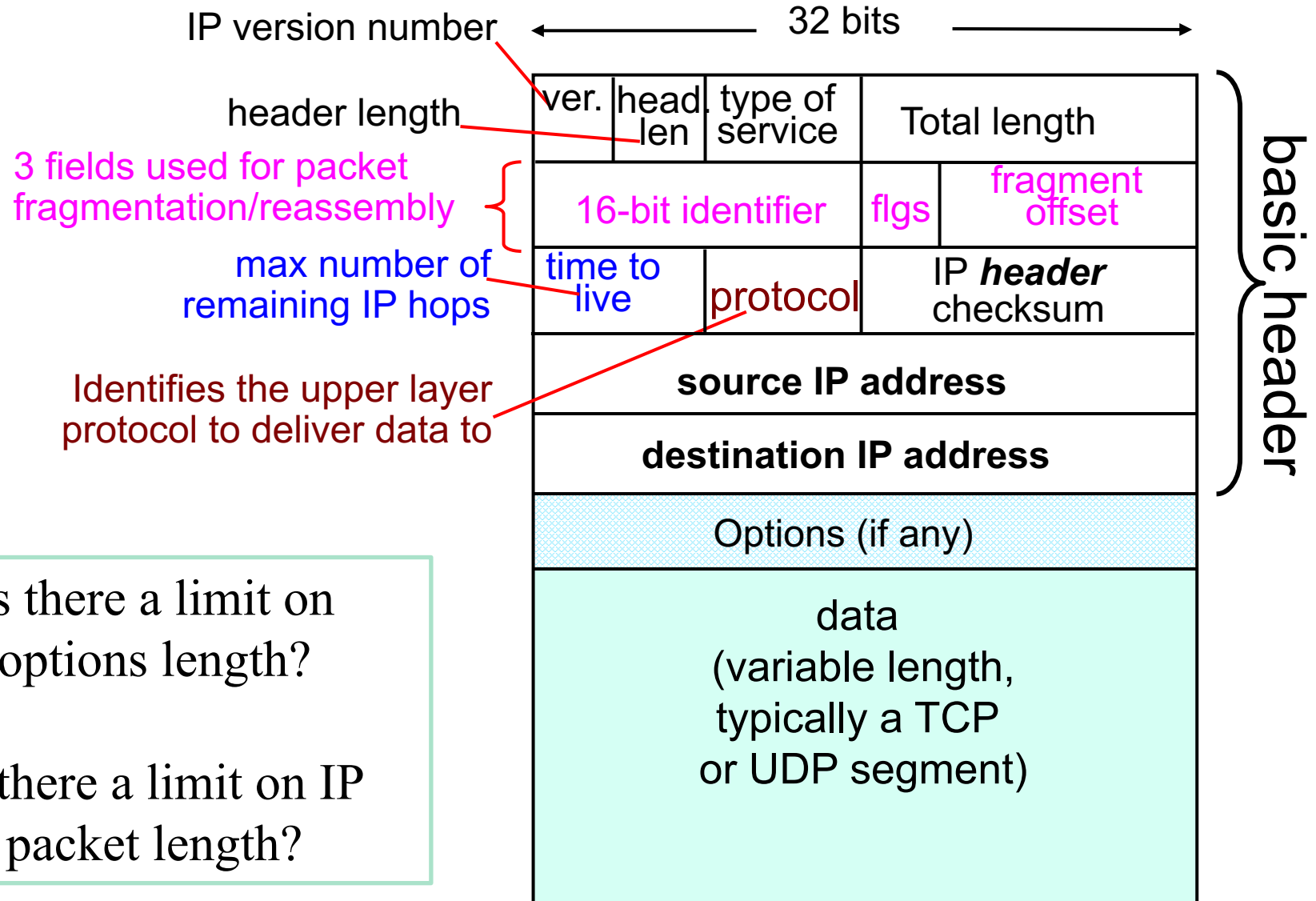


Internet: A Datagram Network

- ♦ **hosts** are connected to **subnets**
- ♦ subnets are interconnected by **routers**
- ♦ All hosts and routers speak **IP**
- ♦ **IP** provides two basic functions
 - Assigning globally unique addresses to all connected points
 - Datagram delivery from source to destination hosts

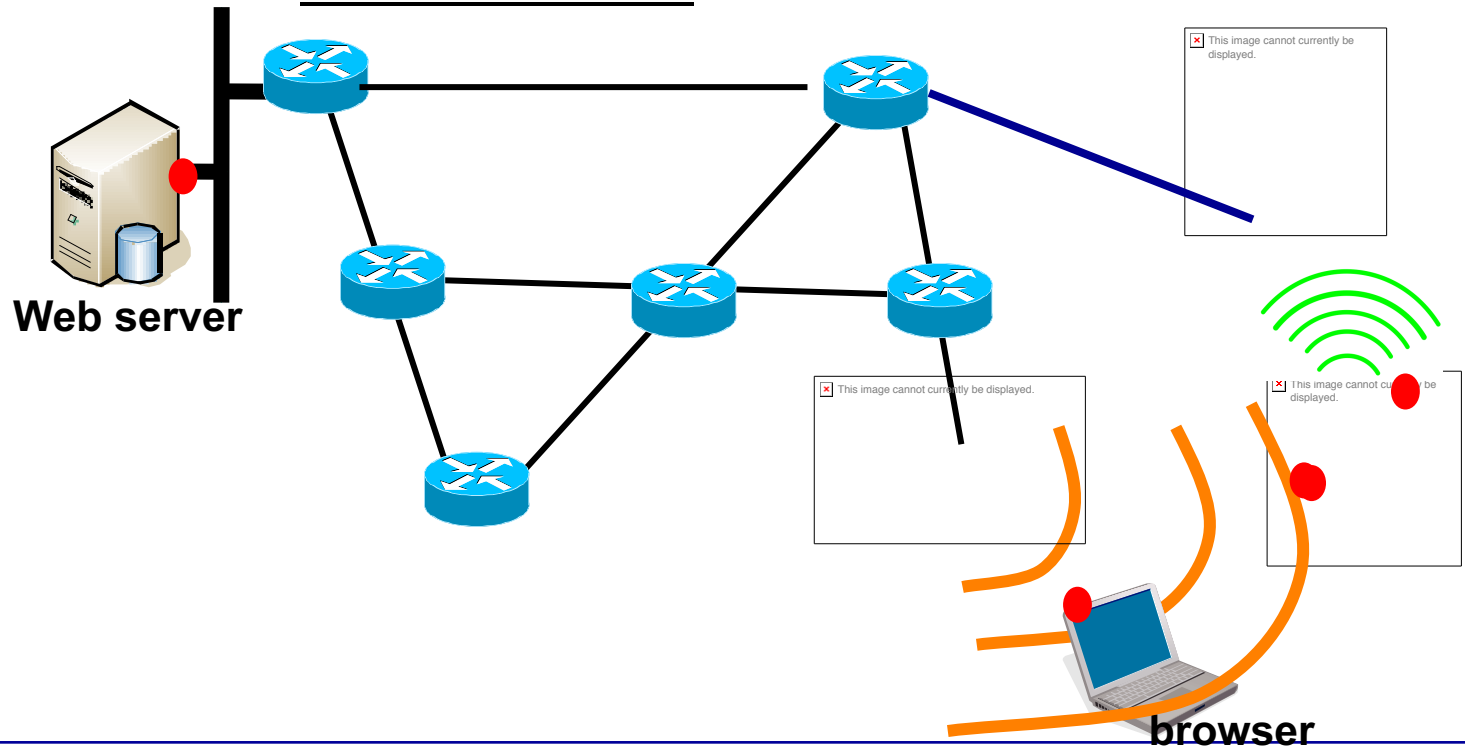


IP datagram format



Two further clarifications

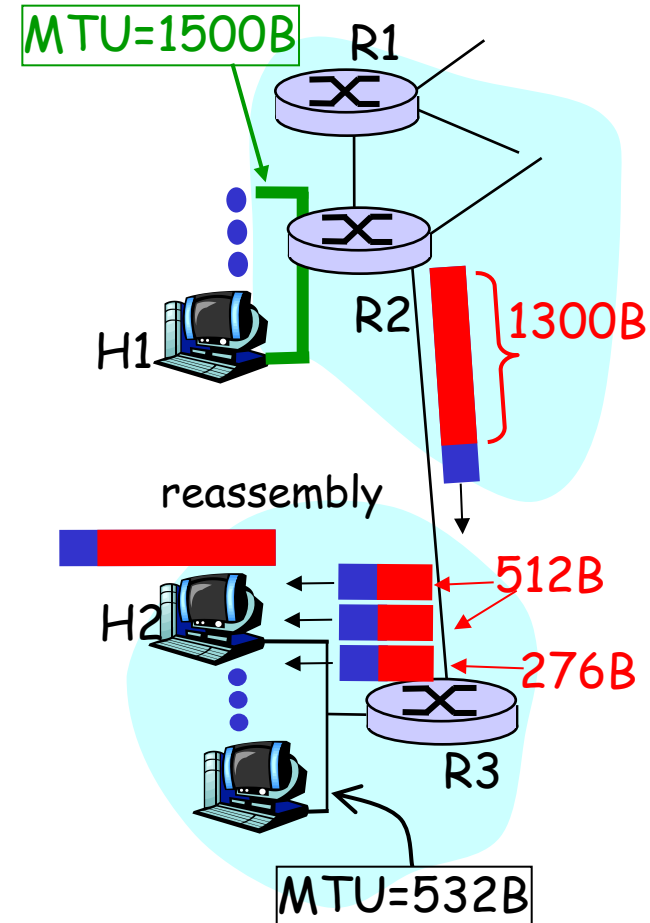
- ◆ IP assigns a globally unique IP address to each **attachment point**
- ◆ When delivering a packet that is too big for next hop: IP routers fragment the packet and then reassemble it at destination host



IP Fragmentation & Reassembly

- ◆ Different subnets may have different MTUs (**Maximum Transmission Unit**)
- ◆ Sending host uses its local MTU size
- ◆ Routers *fragment* IP packets if the next link has a smaller MTU
 - chop packets to the MTU size of next link
 - further fragmentation down the path possible
- ◆ packet fragments are reassembled at destination host

H1 sending an IP packet of 1300 byte data to H2:



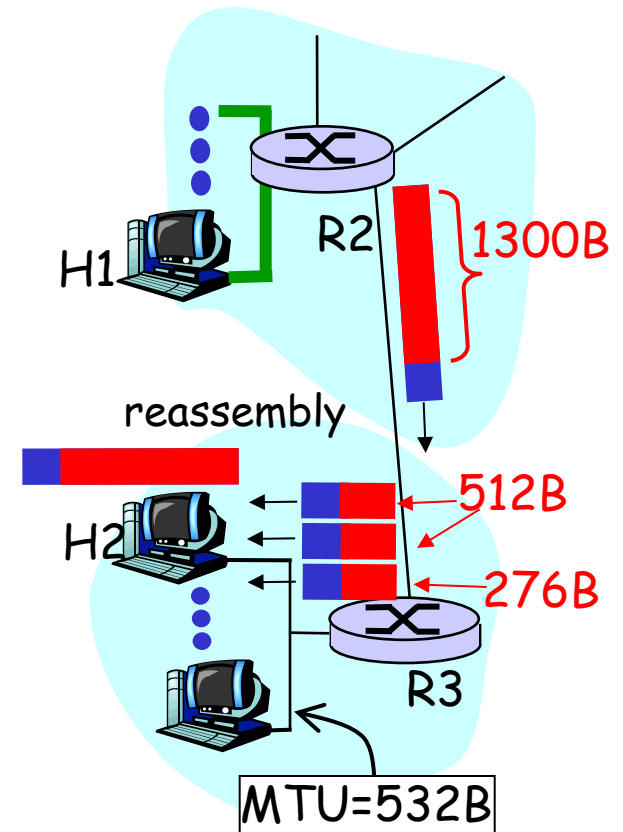
IP Fragmentation: An example

4	5	TOS	1320
7394	000	0	
rest of the IP header			
data (1300 bytes)			

4	5	TOS	532
7394	001	0	
rest of the IP header			
data (512 bytes)			

4	5	TOS	532
7394	001	64	
rest of the IP header			
data (512 bytes)			

4	5	TOS	296
7394	000	128	
rest of the IP header			
data (276 bytes)			

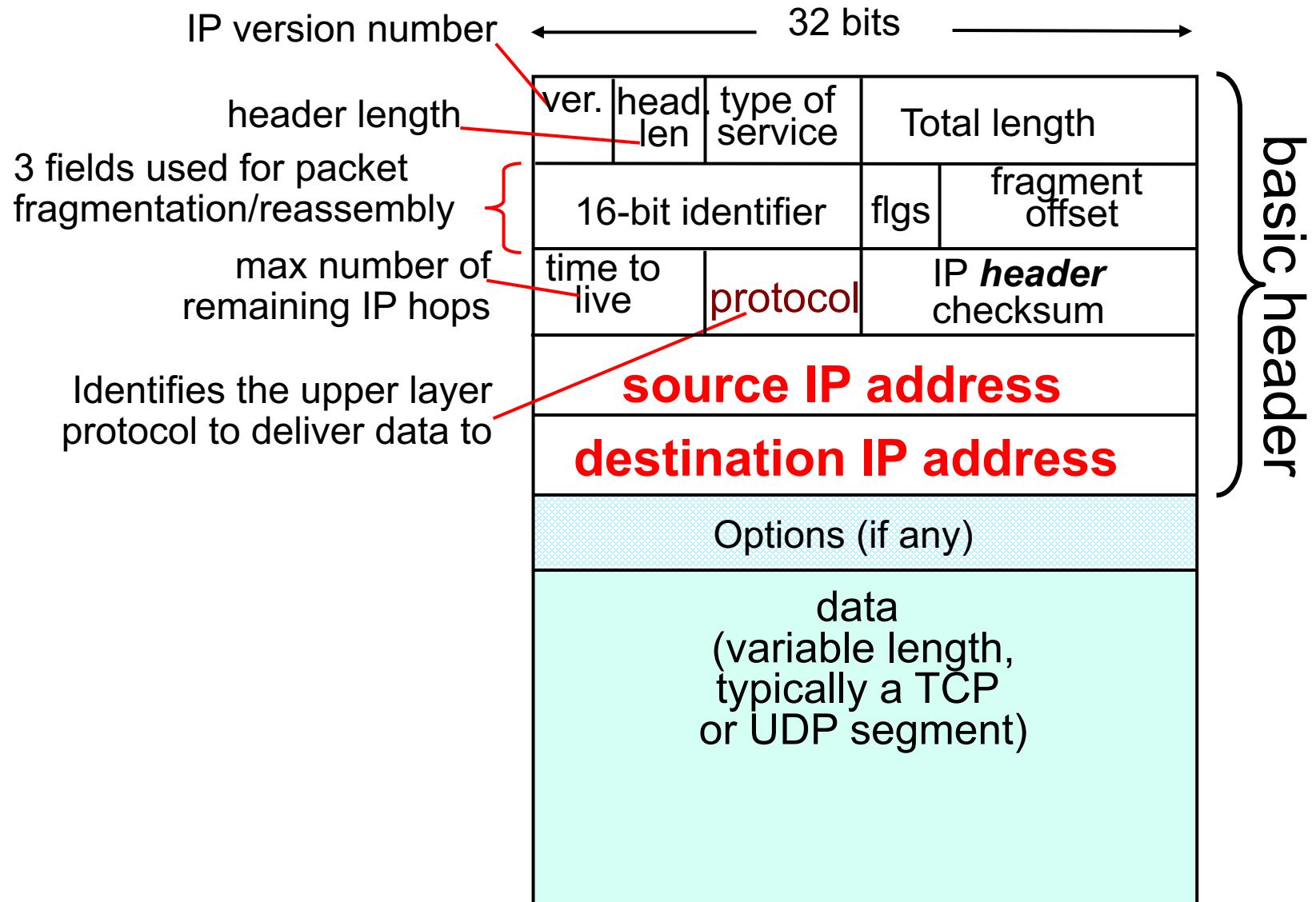


Destination host: examining IP header fields:

identifier: tell all pieces in the same packet

flag bit: if MF=0, the last fragment
offset: tell whether there are holes missing in the middle

IP datagram format



IP Packets in Wireshark

ver.	head len	type of service	Total length	
16-bit identifier			flgs	fragment offset
time to live	protocol		IP header checksum	
source IP address				
destination IP address				
Options (if any)				
data (variable length, typically a TCP or UDP segment)				

Thunderbolt Ethernet: en3

ip.addr==131.179.128.22

No.	Time	Source	Destination	Protocol	Length	Info
332	1...	131.179.128.22	131.179.196.220	TCP	78	80 → 50143 [SYN, ACK] S
333	1...	131.179.196.220	131.179.128.22	TCP	66	50143 → 80 [ACK] Seq=1
334	1...	131.179.196.220	131.179.128.22	HTTP	821	GET /classes/spring16/c

- Frame 334: 821 bytes on wire (6568 bits), 821 bytes captured (6568 bits) on interface 0
- Ethernet II, Src: Apple_c0:61:b6 (68:5b:35:c0:61:b6), Dst: All-MSRP-routers_c4 (00:00:0c:0

▼ Internet Protocol Version 4, Src: 131.179.196.220, Dst: 131.179.128.22

- 0100 = Version: 4
- 0101 = Header Length: 20 bytes
- Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)
- Total Length: 807
- Identification: 0x7bc0 (31680)
- Flags: 0x00
- Fragment offset: 0
- Time to live: 64
- Protocol: TCP (6)
- Header checksum: 0x0000 [validation disabled]
- Source: 131.179.196.220
- Destination: 131.179.128.22
- [Source GeoIP: Unknown]
- [Destination GeoIP: Unknown]
- Transmission Control Protocol, Src Port: 50143 (50143), Dst Port: 80 (80), Seq: 1, Ack: 1, Len: 755
- Hypertext Transfer Protocol

```

0000  00 00 0c 07 ac c4 68 5b 35 c0 61 b6 08 00 45 00  .....h[ 5.a...E.
0010  03 27 7b c0 00 00 40 06 00 00 83 b3 c4 dc 83 b3  ..{...@. ....
0020  80 16 c3 df 00 50 00 d1 e6 37 cf 60 4d 07 80 18  ....P.. .7.`M...
0030  10 15 4f 73 00 00 01 01 08 0a 31 e6 53 89 8e c8  ..0s.... ..1.S...
0040  bf f7 47 45 54 20 2f 63 6c 61 73 73 65 73 2f 73  ..GET /c lasses/s
0050  70 72 69 6e 67 31 36 2f 63 73 31 31 38 2f 20 48  pring16/ cs118/ H
0060  54 54 50 2f 31 2e 31 0d 0a 48 6f 73 74 3a 20 77  TTP/1.1. .Host: w
  
```

Internet Protocol Version 4 (ip), 20 bytes

Packets: 480 · Displayed: 25 (5.2%) · Dropped: 0 (0.0%) Profile: Default

IPv4 Address structure

32-bits, uniquely identifies a host or router *interface*

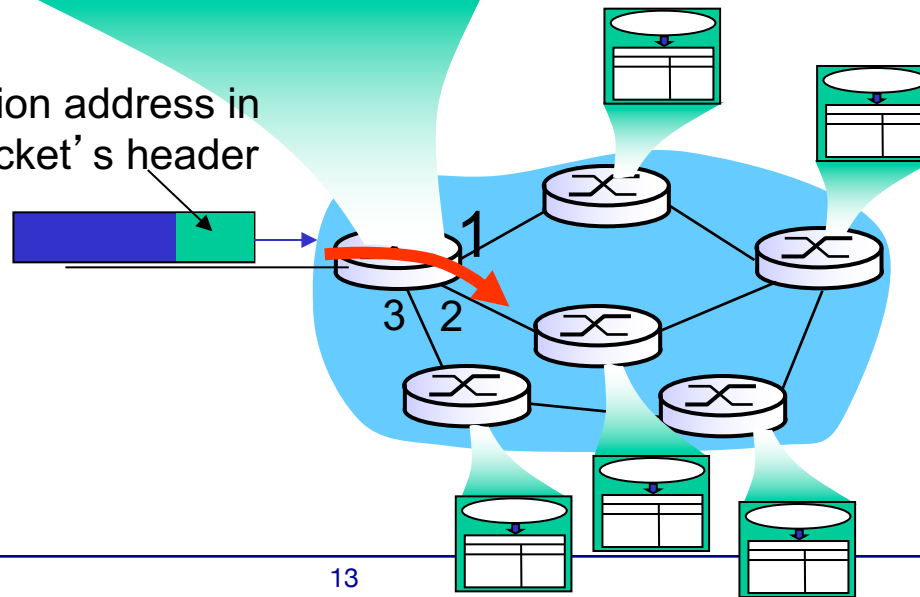
– *interface*: connecting point between host/router and physical link

173.1.1.10 = $\underbrace{10101101}_{173} \underbrace{00000001}_1 \underbrace{00000001}_1 \underbrace{00001010}_{10}$

4 billion IP addresses, so rather than list individual destination address, list **range** of addresses (aggregate table entries)

local forwarding table	
dest address	output link
address-range 1	3
address-range 2	2
address-range 3	2
address-range 4	1

IP destination address in arriving packet's header



How to define IP ranges in routing table?

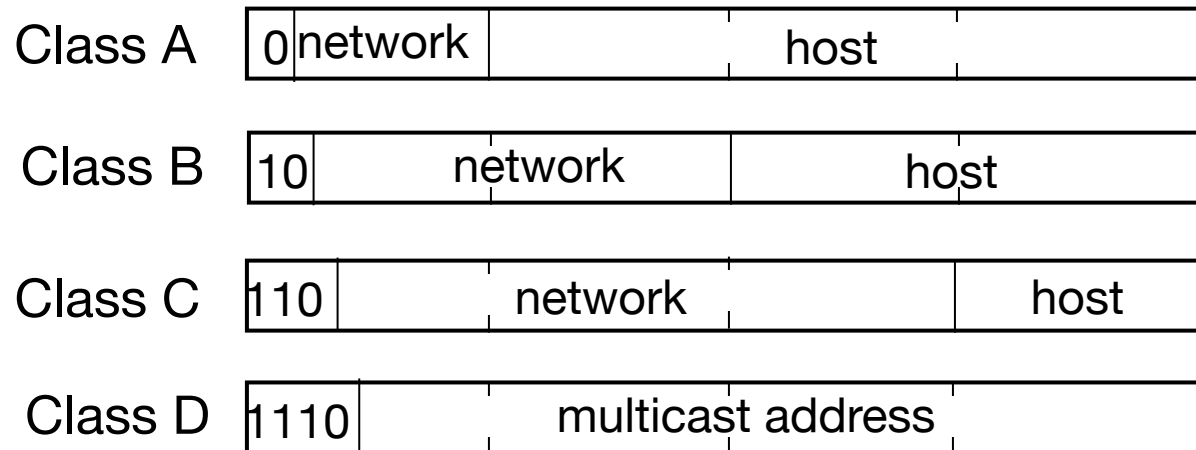
- ◆ Should be compact
- ◆ Should be fast and easy to process

Designate a number of bits for the “network”.

Routing table will contain only “networks”

How many bits for network ID

- ◆ Original IP design (RFC791): **class-based address**



Calculate the address ranges

Today's IPv4 Address Structure

- ◆ Two changes added later:
 - **CIDR: C**lassless **I**nter**D**omain **R**outing (today)
 - network portion can take any arbitrary number of bits



- **Subnetting:** An organization gets one address block, then split the host part into two parts: **subnet** and **host** part

Classless InterDomain Routing



- ◆ Internet Service Providers (ISPs), and some large user sites, get blocks of IP addresses from the Regional Internet Registries (RIRs)
- ◆ Internet customers get a sub-block from their ISP's address block

ISP's block	<u>11001000 00010111 00010000</u> 00000000
Organization 0	<u>11001000 00010111 00010000</u> 00000000
Organization 1	<u>11001000 00010111 00010010</u> 00000000
Organization 2	<u>11001000 00010111 00010100</u> 00000000
...
Organization 7	<u>11001000 00010111 00011110</u> 00000000

CIDR Address Format



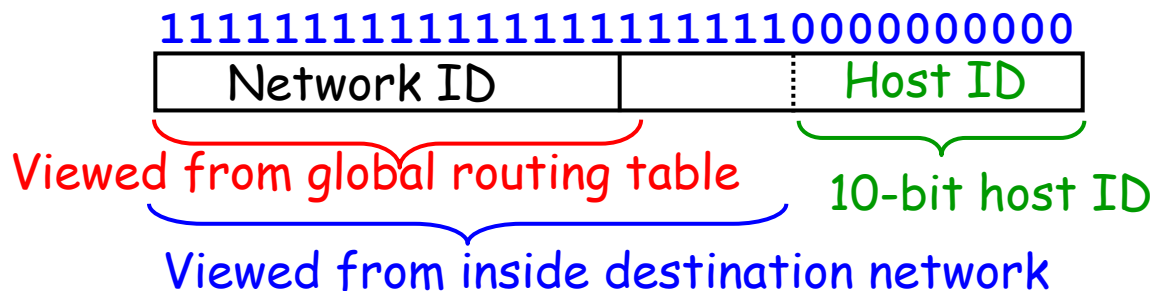
- ◆ $a.b.c.d/x$, x = #bits in network ID portion of the address
- ◆ address: $a.b.c.d$, network mask: $2^{32} - 2^{(32-x)}$

200.23.16.0/23

address 200.23.16.0, netmask 255.255.254.0

IP Subnet

- ♦ **subnet mask**: indicates the portion of the address that is considered as “network ID” *by the local site*
 - Does not need to align with byte boundary



- ♦ subnets are **invisible** outside the local site
 - backbone routers only know how to forward packets to the networkID
 - Within the organization:
 - routers store: [subnet, **mask**, next hop]
 - Each host is configured with an IP address and a subnet mask

Special Addresses

- ◆ 0.0.0.0/8
 - “this network”
- ◆ 255.255.255.255/32
 - broadcast address of “this network”
- ◆ first address of the network (e.g., 192.168.1.0 for 192.168.1.0/24)
 - network address (cannot be used for end-hosts)
- ◆ last address of the network (e.g., 192.168.1.255 for 192.168.1.0/24)
 - broadcast address for the network