# Computer Networks

**Lecture 13:      Transport layer + DNS**

# Round Trip Time Estimation

Sample

Data

ACK

- Original TCP round-trip estimator
    - RTT estimated as a moving average
    - new_rtt = α (old_rtt) + (1 − α)(new_sample)
    - Recommended α: 0.8-0.9 (0.875 for most TCPs)
- RTO = 2 * new_rtt (i.e. TCP is conservative)

# RTT Sample Ambiguity

Initial Send

RTO

Sample

Retry

ACK

Initial Send

RTO

ACK

Sample?

Retry

☐ Karn's algorithm: ignore samples for retransmitted segments

# TCP Congestion Control

- **The network is congested if the load in the network is higher than its capacity.**
- Each TCP connection has a window
  - Controls the number of unACKed packets
- Sending rate is ~ window/RTT
- Idea: vary the window size to control the send rate
- Introduce a congestion window at the sender
  - Congestion control is sender-side problem

# Two Basic Components

1. ## Detect congestion
   - Packet dropping is most reliably signal
     - Delay-based methods are hard and risky
   - How do you detect packet drops? ACKs
     - Timeout after not receiving an ACK
     - Several duplicate ACKs in a row (ignore for now)

2. ## Rate adjustment algorithm
   - Modify *cwnd*
   - Probe for bandwidth
   - Responding to congestion

# Rate Adjustment

- Recall: TCP is ACK clocked
  - Congestion = delay = long wait between ACKs
  - No congestion = low delay = ACKs arrive quickly
- Basic algorithm
  - Upon receipt of ACK: increase cwnd
    - Data was delivered, perhaps we can send faster
    - *cwnd* growth is proportional to RTT
  - On loss: decrease cwnd
    - Data is being lost, there must be congestion
- Question: increase/decrease functions to use? !!!!
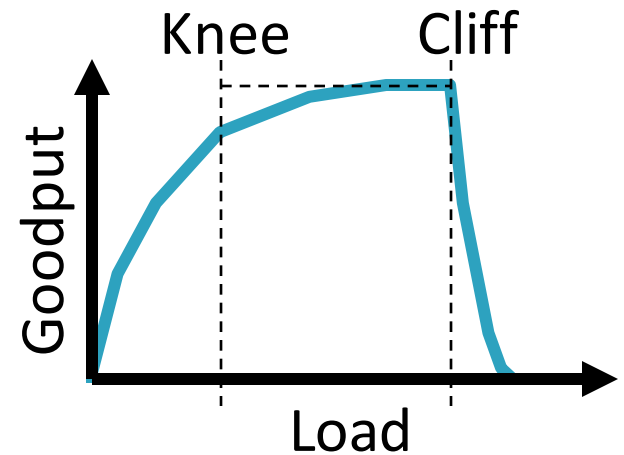
# Implementing Congestion Control

- Maintains three variables:
  - *cwnd*:  congestion window
  - *adv_wnd*: receiver advertised window
  - *ssthresh*:  threshold size (used to update *cwnd*)
- For sending, use: *wnd = min(cwnd, adv_wnd)*
- Two phases of congestion control
  1. Slow start (*cwnd < ssthresh*)
     - Probe for bottleneck bandwidth
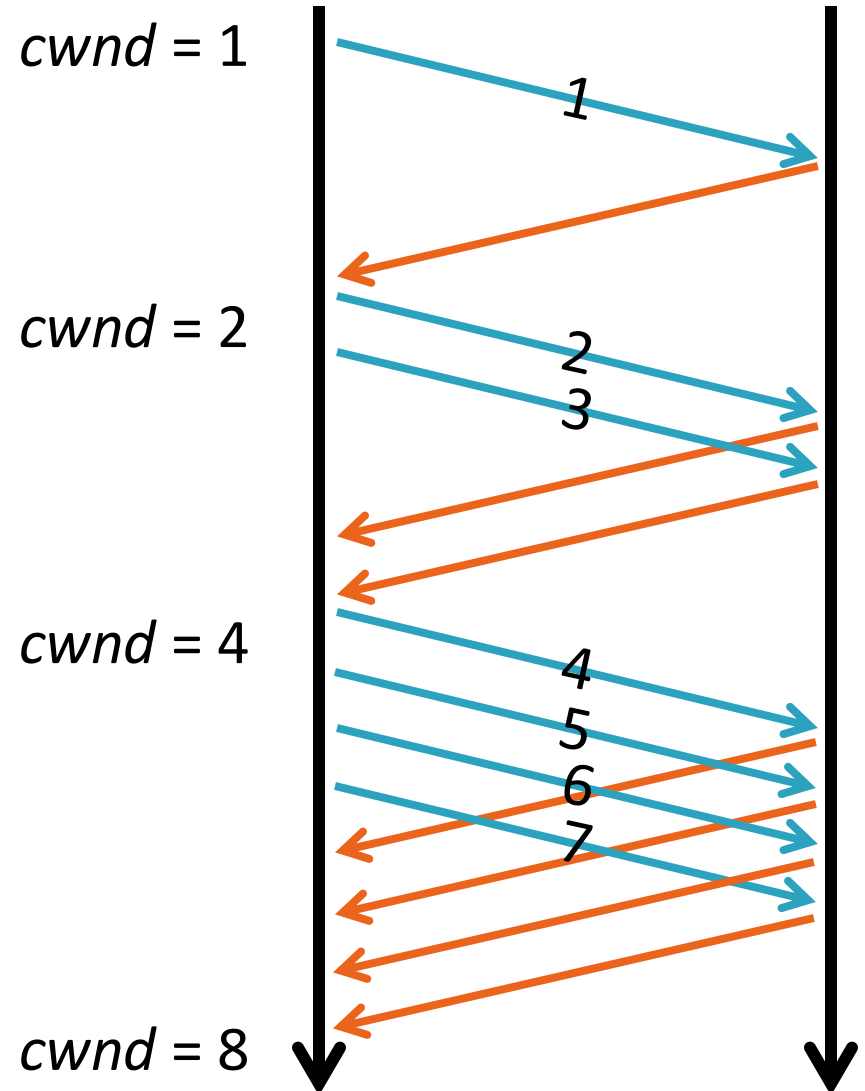  2. Congestion avoidance (*cwnd >= ssthresh*)
     - AIMD

# Slow Start

- Goal: reach knee quickly
- Upon starting (or restarting) a connection
  - *cwnd* =1
  - *ssthresh = adv_wnd*
  - Each time a segment is ACKed, *cwnd*++
- Continues until…
  - *ssthresh* is reached
  - Or a packet is lost
- Slow Start is not actually slow
  - *cwnd* increases exponentially

# Slow Start Example

- *cwnd* grows rapidly
- Slows down when…
  - *cwnd >= ssthresh*
  - Or a packet drops



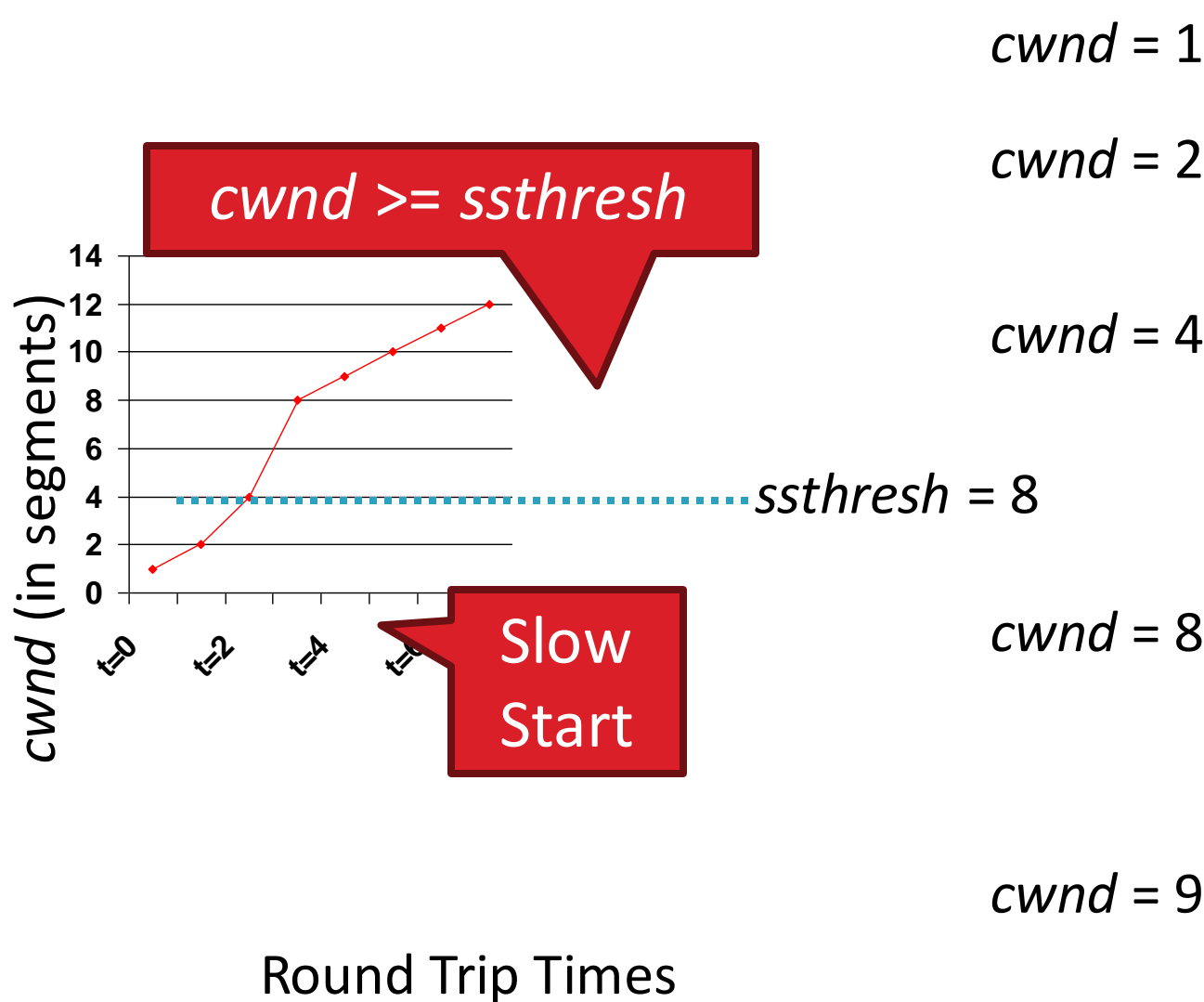*cwnd* = 1

1

*cwnd* = 2

2
3

*cwnd* = 4

4
5
6
7

*cwnd* = 8

# Congestion Avoidance

- Additive Increase Multiplicative Decrease (AIMD) mode

- *ssthresh* is lower-bound guess about location of the knee

- **If** *cwnd >= ssthresh* **then**
  each time a segment is ACKed
  increment *cwnd by 1/cwnd  (cwnd += 1/cwnd).*

- So *cwnd* is increased by one only if all segments have been acknowledged
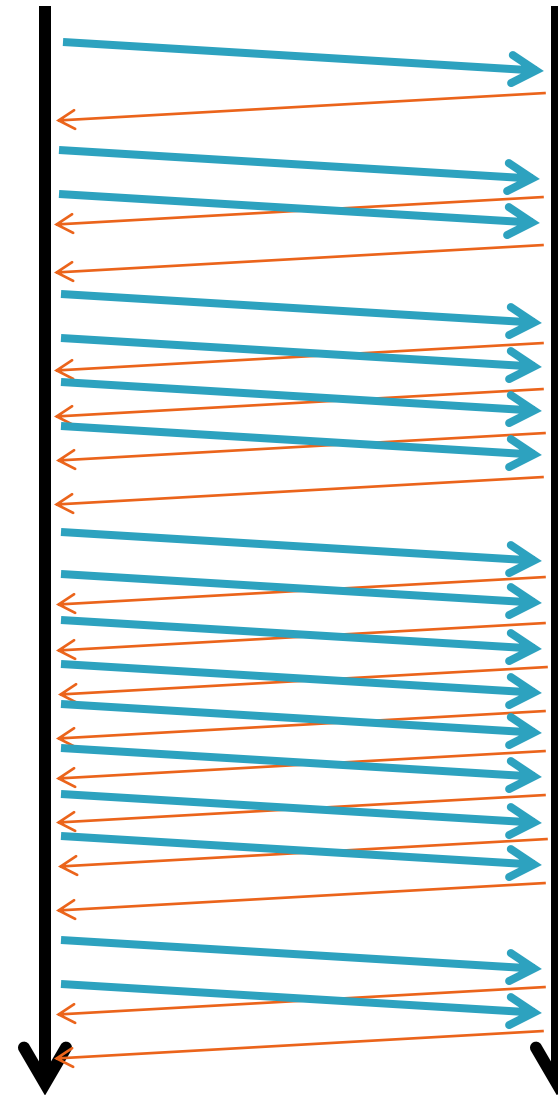
# Congestion Avoidance Example

cwnd >= ssthresh

Slow Start

*cwnd* = 1

*cwnd* = 2

*cwnd* = 4

ssthresh = 8

*cwnd* = 8

*cwnd* = 9

Round Trip Times

*cwnd* (in segments)
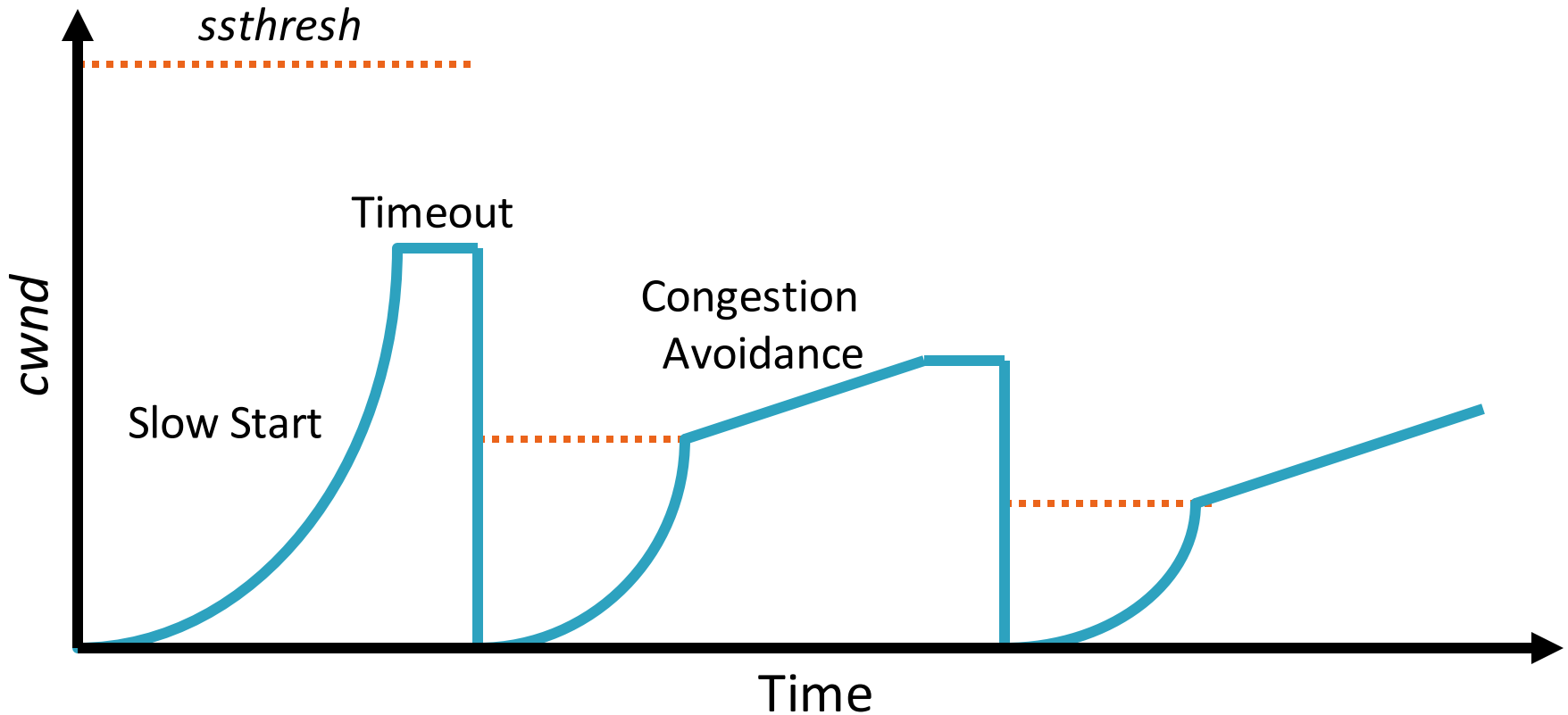
14
12
10
8
6
4
2
0

t=0   t=2   t=4   t=6

# The Big Picture – TCP Tahoe
## (the original TCP)

# Outline

- UDP
- TCP
- Congestion Control
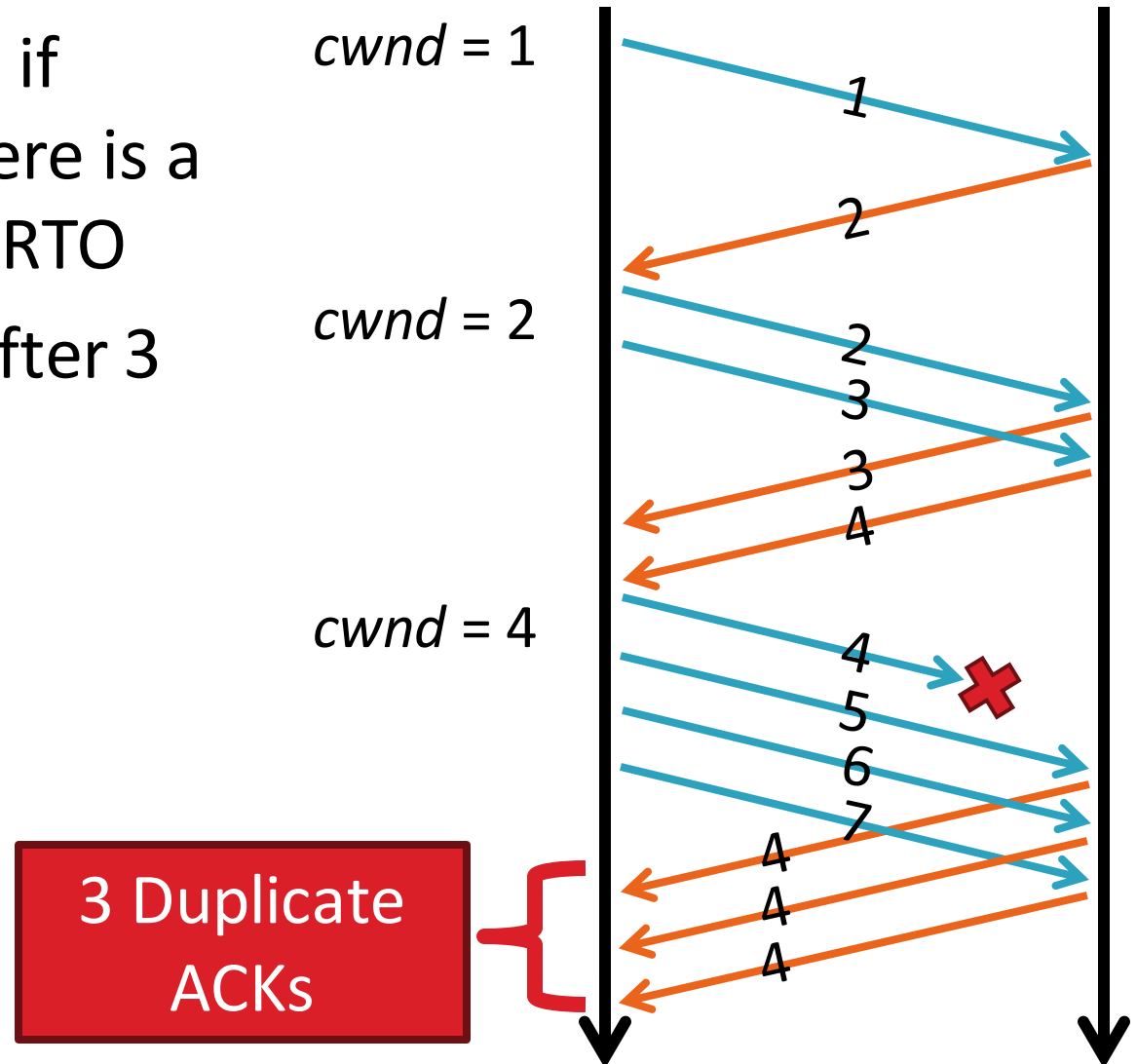- **Evolution of TCP**
- Problems with TCP

# The Evolution of TCP

- Thus far, we have discussed TCP Tahoe
  - Original version of TCP
- However, TCP was invented in 1974!
  - Today, there are many variants of TCP
- Early, popular variant: TCP Reno
  - Tahoe features, plus…
  - Fast retransmit
    - 3 duplicate ACKs? -> retransmit (don't wait for RTO)
  - Fast recovery
    - On loss: set cwnd = cwnd/2 (ssthresh = new cwnd value)

# TCP Reno: Fast Retransmit

- Problem: in Tahoe, if segment is lost, there is a long wait until the RTO
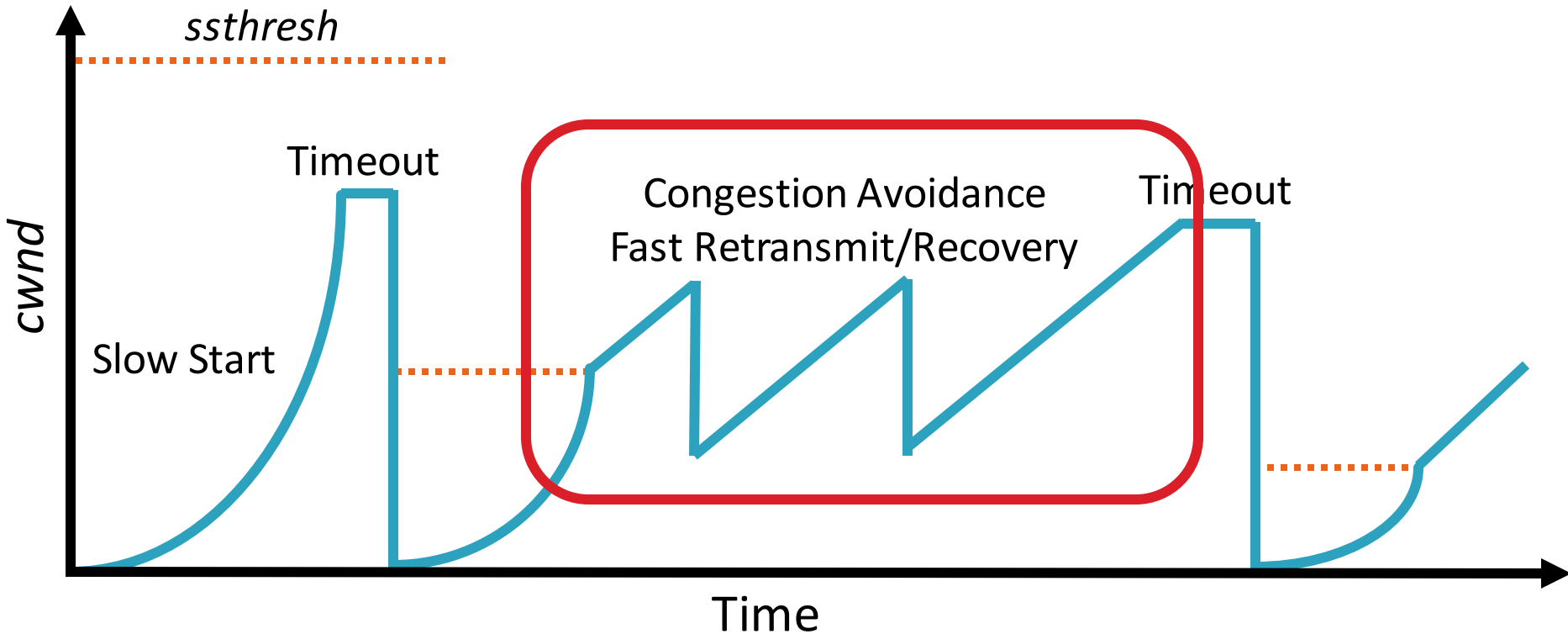- Reno: retransmit after 3 duplicate ACKs

*cwnd* = 1

1

2

*cwnd* = 2

2
3
3
4

*cwnd* = 4

4
5
6
7
4
4
4

**3 Duplicate ACKs**

# TCP Reno: Fast Recovery

- After a fast-retransmit set *cwnd* to *cwnd/2*
  - Also reset ssthresh to the new halved cwnd value
  - i.e. don't reset *cwnd* to 1
  - Avoid unnecessary return to slow start
  - Prevents expensive timeouts
- But when RTO expires still do *cwnd* = 1
  - Return to slow start, same as Tahoe
  - Indicates packets aren't being delivered at all
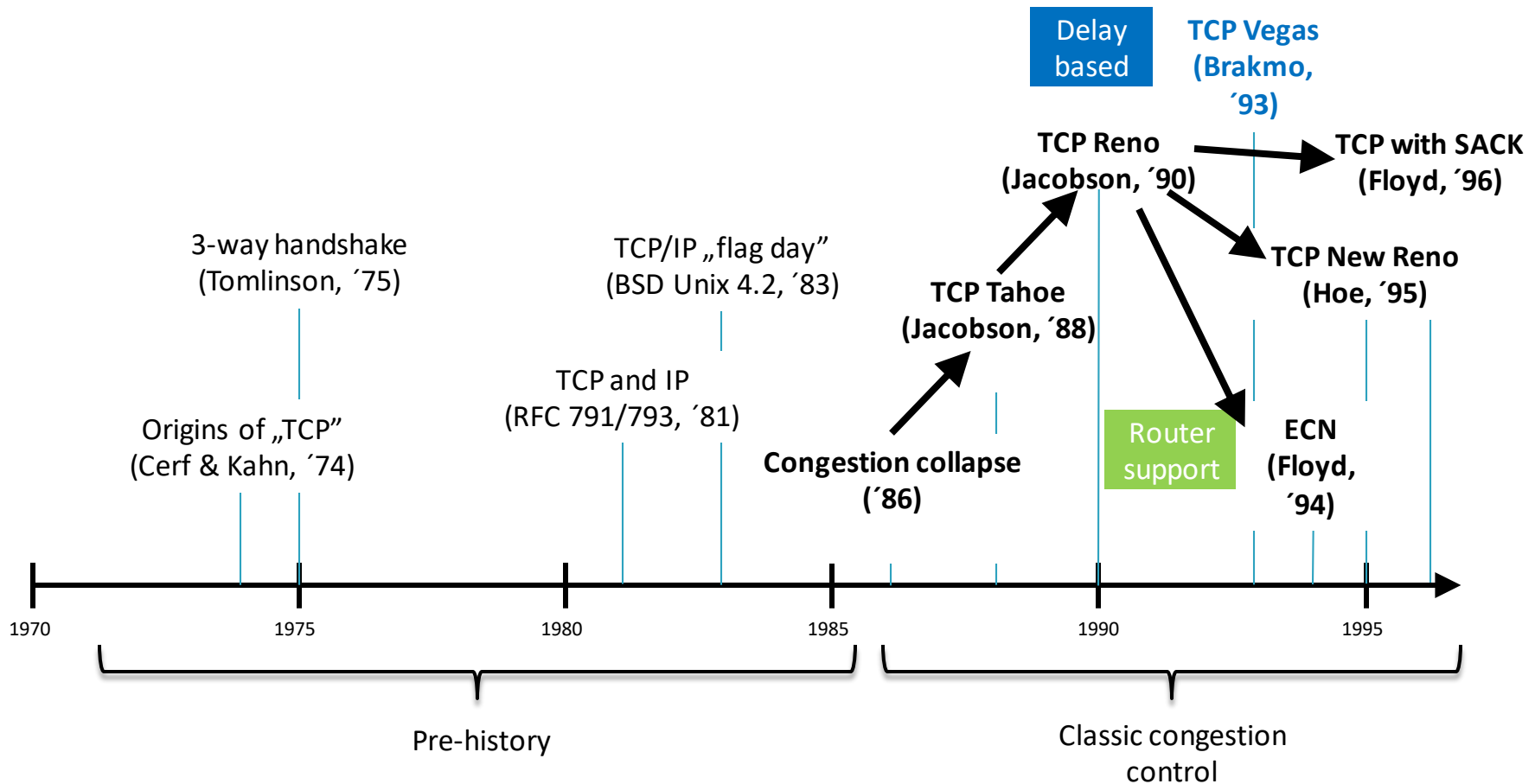  - i.e. congestion must be really bad

# Fast Retransmit and Fast Recovery

- At steady state, *cwnd* oscillates around the optimal window size

- TCP always forces packet drops

# Many TCP Variants…
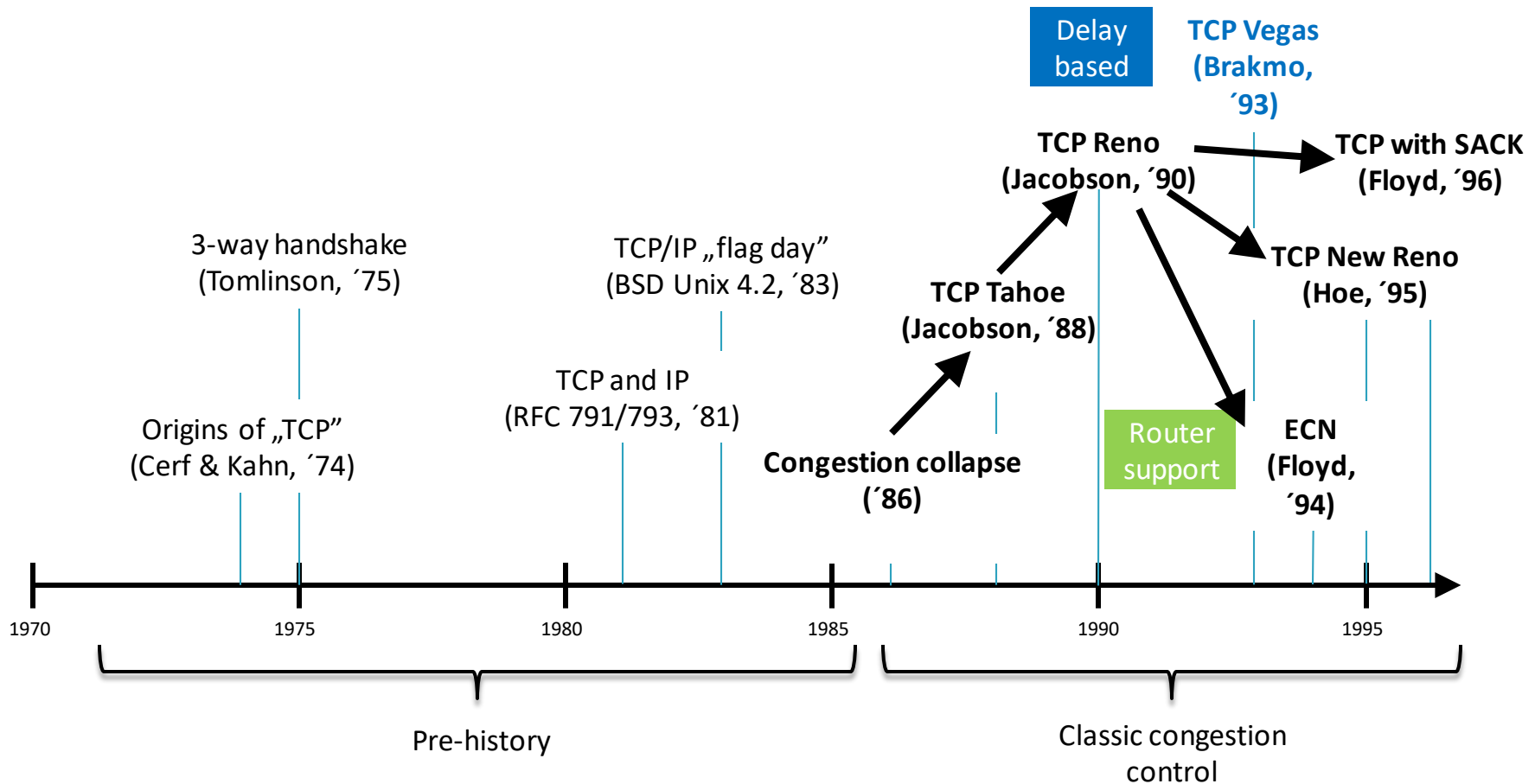
- Tahoe: the original
  - Slow start with AIMD
  - Dynamic RTO based on RTT estimate
- Reno:
  - fast retransmit (3 dupACKs)
  - fast recovery (cwnd = cwnd/2 on loss)
- NewReno: improved fast retransmit
  - Each duplicate ACK triggers a retransmission
  - Problem: >3 out-of-order packets causes pathological retransmissions
- Vegas: delay-based congestion avoidance
- And many, many, many more…

# Transport layer evolution



Delay based

TCP Vegas (Brakmo, ´93)

TCP Reno (Jacobson, ´90)

TCP with SACK (Floyd, ´96)

3-way handshake (Tomlinson, ´75)

TCP/IP „flag day" (BSD Unix 4.2, ´83)

TCP New Reno (Hoe, ´95)

TCP Tahoe (Jacobson, ´88)

TCP and IP (RFC 791/793, ´81)

Origins of „TCP" (Cerf & Kahn, ´74)

Congestion collapse (´86)

Router support

ECN (Floyd, ´94)

1970    1975    1980    1985    1990    1995

Pre-history

Classic congestion control

# Transport layer evolution

# Transport layer (r)evolution

QUIC
(Roskind et al., ´12)

Delay based

TCP Vegas
(Brakmo, ´93)

Background traffic

TCP LEDBAT
(IETF, ´08)

BBRv2
(ICCRG, ´19)

TCP Reno
(Jacobson, ´90)

TCP with SACK
(Floyd, ´96)

Compound TCP
(Windows, ´07)

COPA
(Arun et al., ´18)

TCP New Reno
(Hoe, ´95)

FAST TCP
(Low et al., ´04)

BBRv1
(ICCRG, ´17)

P Tahoe
bson, ´88)

TCP CUBIC
(Linux, ´06)

TCP PCC
(Dong et al., ´15)

pse

Router support

ECN
(Floyd, ´94)

TCP BIC
(Linux, ´04)

DCTCP
(Alizadeh et al., ´10)

TCP Prague
(Linux patch, ´19)

1990        1995        2000        2005        2010        2015        2020

Classic congestion control

Diversification

Congestion control revolution
CC in application layers

# TCP in the Real World

- What are the most popular variants today?
  - Key problem: TCP performs poorly on high bandwidth-delay product networks (like the modern Internet)
  - Compound TCP (Windows)
    - Based on Reno
    - Uses two congestion windows: delay based and loss based
    - Thus, it uses a *compound* congestion controller
  - TCP CUBIC (Linux)
    - Enhancement of BIC (Binary Increase Congestion Control)
    - Window size controlled by cubic function
    - Parameterized by the time $T$ since the last dropped packet

# High Bandwidth-Delay Product

- Key Problem: TCP performs poorly when
  - The capacity of the network (bandwidth) is large
  - The delay (RTT) of the network is large
  - Or, when bandwidth * delay is large
    - b * d = maximum amount of in-flight data in the network
    - a.k.a. the bandwidth-delay product
- Why does TCP perform poorly?
  - Slow start and additive increase are slow to converge
  - TCP is ACK clocked
    - i.e. TCP can only react as quickly as ACKs are received
    - Large RTT → ACKs are delayed → TCP is slow to react

# Goals

- Fast window growth
  - Slow start and additive increase are too slow when bandwidth is large
  - Want to converge more quickly
- Maintain fairness with other TCP varients
  - Window growth cannot be too aggressive
- Improve RTT fairness
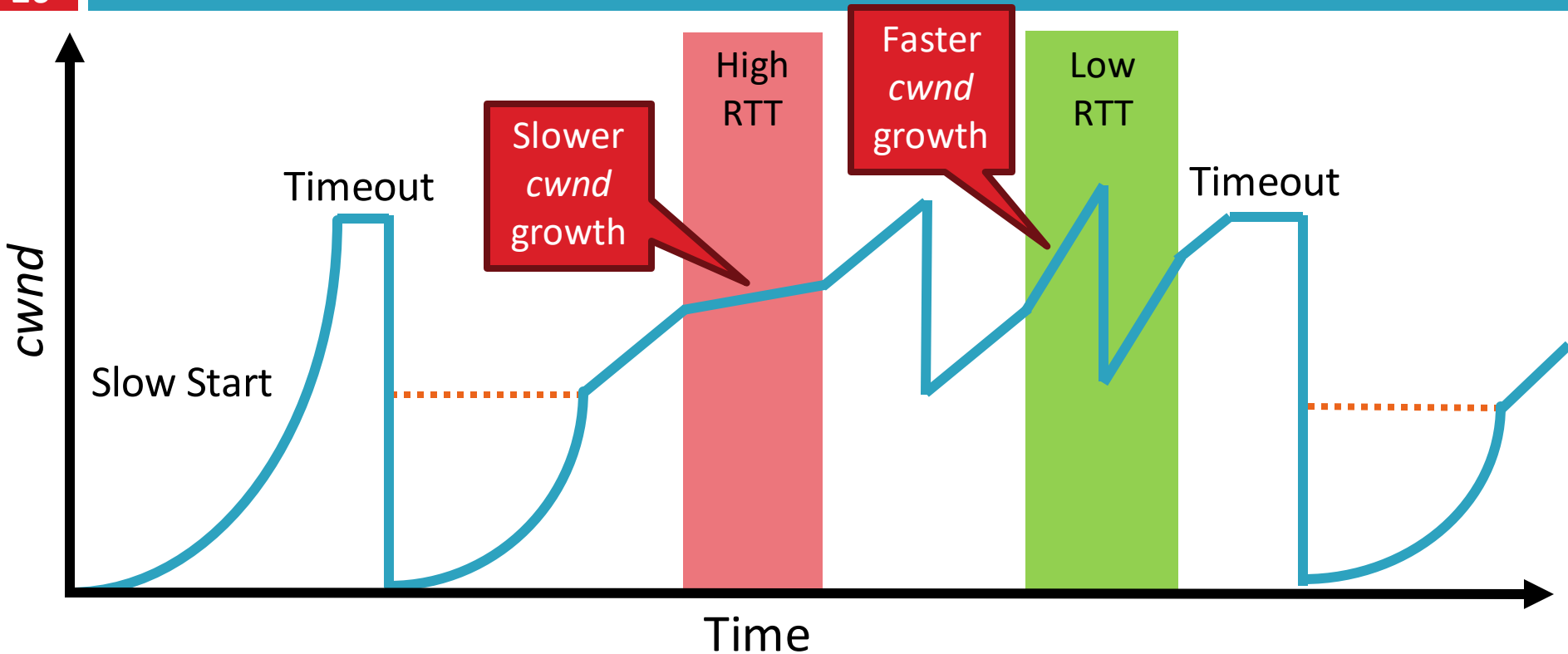  - TCP Tahoe/Reno flows are not fair when RTTs vary widely
- Simple implementation

# Compound TCP Implementation

- Default TCP implementation in Windows (before Win 10)
- Key idea: split *cwnd* into two separate windows
    - Traditional, loss-based window
    - New, delay-based window
- *wnd* = min(*cwnd* + *dwnd*, *adv_wnd*)
    - *cwnd* is controlled by AIMD
    - *dwnd* is the delay window
- Rules for adjusting *dwnd:*
    - If RTT is increasing, decrease *dwnd* (*dwnd* >= 0)
    - If RTT is decreasing, increase *dwnd*
    - Increase/decrease are proportional to the rate of change

# Compound TCP Example

- Aggressiveness corresponds to changes in RTT
- Advantages: fast ramp up, more fair to flows with different RTTs
- Disadvantage: must estimate RTT, which is very challenging

# TCP CUBIC Implementation

- Default TCP implementation in Linux
- Replace AIMD with cubic function

$$W_{cubic} = C(T - K)^3 + W_{max} \tag{1}$$

$$C \text{ is a scaling constant, and } K = \sqrt[3]{\frac{W_{max}\beta}{C}}$$

- B $\rightarrow$ a constant fraction for multiplicative increase
- T $\rightarrow$ time since last packet drop
- W_max $\rightarrow$ cwnd when last packet dropped

# TCP CUBIC Implementation

☐

☐

# TCP CUBIC Example

- Less wasted bandwidth due to fast ramp up
- Stable region and slow acceleration help maintain fairness
    - Fast ramp up is more aggressive than additive increase
    - To be fair to Tahoe/Reno, CUBIC needs to be less aggressive

# Outline

- UDP
- TCP
- Congestion Control
- Evolution of TCP
- Problems with TCP

# Issues with TCP

- The vast majority of Internet traffic is TCP

- However, many issues with the protocol

    - Poor performance with small flows

    - Really poor performance on wireless networks

    - Susceptibility to denial of service

# Small Flows

- Problem: TCP is biased against short flows
    - 1 RTT wasted for connection setup (SYN, SYN/ACK)
    - *cwnd* always starts at 1

- Vast majority of Internet traffic is short flows
    - Mostly HTTP transfers, <100KB
    - Most TCP flows never leave slow start!

- Proposed solutions (driven by Google):
    - Increase initial *cwnd* to 10
    - TCP Fast Open: use cryptographic hashes to identify receivers, eliminate the need for three-way handshake

# Wireless Networks

- Problem: Tahoe and Reno assume loss = congestion
  - True on the WAN, bit errors are very rare
  - False on wireless, interference is very common
- TCP throughput ~ 1/sqrt(drop rate)
  - Even a few interference drops can kill performance
- Possible solutions:
  - Break layering, push data link info up to TCP
  - Use delay-based congestion detection (TCP Vegas)
  - Explicit congestion notification (ECN)

# Denial of Service

- Problem: TCP connections require state
  - Initial SYN allocates resources on the server
  - State must persist for several minutes (RTO)
- SYN flood: send enough SYNs to a server to allocate all memory/meltdown the kernel
- Solution: SYN cookies
  - Idea: don't store initial state on the server
  - Securely insert state into the SYN/ACK packet (sequence number field)
  - Client will reflect the state back to the server

Further topics

# Typical Internet Queuing

- FIFO + drop-tail
  - Simplest choice
  - Used widely in the Internet
- FIFO (first-in-first-out)
  - Implies single class of traffic
- Drop-tail
  - Arriving packets get dropped when queue is full regardless of flow or importance
- Important distinction:
  - FIFO: scheduling discipline
  - Drop-tail: drop policy

# Buffer sizing

- Network is a shared resource
  - Many flows using the same bottleneck

- Temporal overloads should be
  - Buffers are needed

- Buffers are needed for good u

- Drawbacks of large buffers
  - Increased end-to-end delay

Slide from Nick McKeown's presentation at Stanford "pre-workshop" meeting on Buffer Sizing.

1988 — Congestion Avoidance and Control — VJ & MK

1994 — High Performance TCP in ANSNET — CV & CS

2004 — Sizing Router Buffers — GA, IK, NM

2020

$$B \geq \frac{2T \times C}{\sqrt{N}}$$ where $N$ is the number of long-lived flows

$B$

$C$

$Min\ RTT = 2T$

Slide from Nick McKeown's presentation at Stanford "pre-workshop" meeting on Buffer Sizing.

Congestion Avoidance and Control
VJ & MK

High Performance TCP in ANSNET
CV & CS

Sizing Router Buffers
GA, IK, NM

Routers with Very Small Buffers
ME, YG, AG, NM, TR

Experimental Study of Router Buffers
NB, YG, MG, NM, GS

$B$

$A$        $C$        $B$

$$Min\ RTT = 2T$$

Slide from Nick McKeown's presentation at Stanford "pre-workshop" meeting on Buffer Sizing.

# 2010s – reducing queuing delay

- ☐ Active Queue Management (AQM)
  - ▪ Goal is to reduce the average queuing delay, but allow temporal overshoots
  - ▪ Proactively starts dropping or marking packets to reduce queuing delay



Traditional AQM

PI²
(De Schepper et al., ´16)

PED
(Francini, ´11)

SRED
(Ott et al., ´99)

LRED
(Wang et al., ´04)

PIE
…et al., ´13)

AVQ
(Kunniyur et al., ´01)

GSP
(Lautenschlaeger et al., ´15)

RED
(Floyd et al., ´93)

FRED
(Lin et al., ´97)

AFD
(Pan et al., ´03)

CoDel
(TSVWG, ´12)

1990    1995    2000    2005    2010    2015

AQM
with stateless resource sharing

RFQ
(Cao et al., ´00)

PPV
(Nadas et al., ´16)

ABC
(Menth et al., ´18)

# RED Algorithm

- Maintain running average of queue length
- If $avgq < min_{th}$ do nothing
  - Low queuing, send packets through
- If $avgq > max_{th}$, drop packet
  - Protection from misbehaving sources
- Else mark packet in a manner proportional to queue length
  - Notify sources of incipient congestion
  - E.g. by ECN IP field or dropping packets with a given probability

# RED Operation

Max thresh

Min thresh

Average Queue Length

P(drop)

1.0

$max_P$

$min_{th}$

$max_{th}$

Avg queue length

# RED Algorithm

- Maintain running average of queue length

- For each packet arrival

  - Calculate average queue size (avg)

  - If $\text{min}_{th} \le \text{avgq} < \text{max}_{th}$

    - Calculate probability $P_a$

    - With probability $P_a$

      - Mark the arriving packet: drop or set-up ECN

    - Else if $\text{max}_{th} \le \text{avg}$

      - Mark the arriving packet: drop, ECN

# Data Center TCP: DCTCP

# Generality of Partition/Aggregate

- The foundation for many large-scale web applications.
  - Web search, Social network composition, Ad selection, etc.

- Example: **Facebook**

**Partition/Aggregate ~ Multiget**
- Aggregators: **Web Servers**
- Workers: **Memcached Servers**



Internet

**Web Servers**

**Memcached Protocol**

**Memcached Servers**

# Workloads

□ Partition/Aggregate
   **(Query)**                                    → **Delay-sensitive**

□ Short messages [50KB-1MB]
   **(Coordination, Control state)**              → **Delay-sensitive**

□ Large flows [1MB-50MB]
   **(Data update)**                              → **Throughput-sensitive**

# Impairments

- Incast

- Queue Buildup

- Buffer Pressure

# Incast

**Worker 1**

**Worker 2**

**Worker 3**

**Worker 4**

- Synchronized mice collide.
  - ➢ **Caused by Partition/Aggregate.**

**Aggregator**

$RTO_{min} = 300$ ms

**TCP timeout**

# Queue Buildup

- Big flows buildup queues.
  - **Increased latency for short flows.**

**Receiver**

**Sender 2**

- Measurements in Bing cluster
  - **For 90% packets: RTT < 1ms**
  - **For 10% packets: 1ms < RTT < 15ms**

# Data Center Transport Requirements

1.  **High Burst Tolerance**

    – Incast due to Partition/Aggregate is common.

2.  **Low Latency**

    – Short flows, queries

3. **High Throughput**

    – Continuous data updates, large file transfers

**The challenge is to achieve these three together.**

# DCTCP: The TCP/ECN Control Loop

**Sender 1**

**ECN = Explicit Congestion Notification**

**ECN Mark (1 bit)**

**Receiver**

**Sender 2**

# DCTCP: Two Key Ideas

1. React in proportion to the **extent** of congestion, not its **presence**.
   - ✓ Reduces **variance** in sending rates, lowering queuing requirements.

| ECN Marks | TCP | DCTCP |
|---|---|---|
| 1 0 1 1 1 1 0 1 1 1 | Cut window by **50%** | Cut window by **40%** |
| 0 0 0 0 0 0 0 0 0 1 | Cut window by **50%** | Cut window by **5%** |

2. Mark based on **instantaneous** queue length.
   - ✓ Fast feedback to better deal with bursts.

# Data Center TCP Algorithm

## Switch side:

☐ Mark packets when **Queue Length > K.**

**B**  **Mark**  **K**  **Don't Mark**

## Sender side:

– Maintain running average of *fraction* of packets marked **($\alpha$)**.

**In each RTT:**

$$F = \frac{\# \; of \; marked \; ACKs}{Total \; \# \; of \; ACKs} \qquad \alpha \leftarrow (1 - g)\alpha + gF$$

➢ **Adaptive window decreases:** $Cwnd \leftarrow \left(1 - \frac{\alpha}{2}\right)Cwnd$

– Note: decrease factor between 1 and 2.

# DNS

# Layer 8 (The Carbon-based nodes)

- If you want to…
  - Call someone, you need to ask for their phone number
    - You can't just dial "P R O F  G I L L"
  - Mail someone, you need to get their address first
- What about the Internet?
  - If you need to reach Google, you need their IP
  - Does anyone know Google's IP?
- Problem:
  - People can't remember IP addresses
  - Need human readable names that map to IPs

# Internet Names and Addresses

- Addresses, e.g. 129.10.117.100
  - Computer usable labels for machines
  - Conform to structure of the network
- Names, e.g. www.northeastern.edu
  - Human usable labels for machines
  - Conform to organizational structure
- How do you map from one to the other?
  - Domain Name System (DNS)

# History

- Before DNS, all mappings were in *hosts.txt*
  - */etc/hosts* on Linux
  - *C:\Windows\System32\drivers\etc\hosts* on Windows
- Centralized, manual system
  - Changes were submitted to SRI via email
  - Machines periodically FTP new copies of *hosts.txt*
  - Administrators could pick names at their discretion
  - Any name was allowed
    - alans_server_at_sbu_pwns_joo_lol_kthxbye

# Towards DNS

- Eventually, the *hosts.txt* system fell apart
  - Not scalable, SRI couldn't handle the load
  - Hard to enforce uniqueness of names
    - e.g MIT
      - Massachusetts Institute of Technology?
      - Melbourne Institute of Technology?
  - Many machines had inaccurate copies of *hosts.txt*
- Thus, DNS was born

# Outline

- ❑ DNS Basics

- ❑ DNS Security

- ❑ DNS and Censorship

# DNS at a High-Level

- Domain Name System
- Distributed database
  - No centralization
- Simple client/server architecture
  - UDP port 53, some implementations also use TCP
  - Why?
- Hierarchical namespace
  - As opposed to original, flat namespace
  - e.g. .com → google.com → mail.google.com

# Naming Hierarchy

Root

net   edu   com   gov   mil   org   uk   fr   etc.

neu   mit

ccs   ece   husky

www   login   mail

- Top Level Domains (TLDs) are at the top
- Maximum tree depth: 128
- Each Domain Name is a subtree
  - .edu → neu.edu → ccs.neu.edu → www.ccs.neu.edu
- Name collisions are avoided
  - neu.com vs. neu.edu

# Hierarchical Administration

Verisign

Root

ICANN

net    edu    com    gov    mil    org    uk    fr    etc.

neu    mit

ccs

www    login    mail

- Tree is divided into zones
  - Each zone has an administrator
  - Responsible for the part of the hierarchy

- Example:
  - CCIS controls *.ccs.neu.edu
  - NEU controls *.neu.edu

# Server Hierarchy

- Functions of each DNS server:
    - Authority over a portion of the hierarchy
        - No need to store all DNS names
    - Store all the records for hosts/domains in its zone
        - May be replicated for robustness
    - Know the addresses of the root servers
        - Resolve queries for unknown names
- Root servers know about all TLDs
    - The buck stops at the root servers

# Root Name Servers

- Responsible for the Root Zone File
  - Lists the TLDs and who controls them
  - ~272KB in size

```
com.                    172800    IN        NS        a.gtld-servers.net.
com.                    172800    IN        NS        b.gtld-servers.net.
com.                    172800    IN        NS        c.gtld-servers.net.
```

- Administered by ICANN
  - 13 root servers, labeled A→M
  - 6 are anycasted, i.e. they are globally replicated
- Contacted when names cannot be resolved
  - In practice, most systems cache this information

# Map of the Roots

E
A
B
D
L
G
H

Anycast instances

C F I J K M

based on root-servers.org
2006-12-29

# Local Name Servers

Where is google.com?

Northeastern

- Each ISP/company has a local, default name server
- Often configured via DHCP
- Hosts begin DNS queries by contacting the local name server
- Frequently cache query results

# Authoritative Name Servers

Where is www.neu.edu?

www.neu.edu = 155.33.17.68

www.neu.edu

Northeastern

Root

edu

neu

Authority for 'edu'

Authority for 'neu.edu'

☐ Stores the name→IP mapping for a given host

# Basic Domain Name Resolution

- Every host knows a local DNS server
  - Sends all queries to the local DNS server
- If the local DNS can answer the query, then you're done
  1. Local server is also the authoritative server for that name
  2. Local server has cached the record for that name
- Otherwise, go down the hierarchy and search for the authoritative name server
  - Every local DNS server knows the root servers
  - Use cache to skip steps if possible
    - e.g. skip the root and go directly to .edu if the root file is cached

# Recursive DNS Query

www.google.com

Where is www.google.com?

- Puts the burden of resolution on the contacted name server

- How does asgard know who to forward responses too?
  - Random IDs embedded in DNS queries

asgard.ccs.neu.edu

ns1.google.com

Root

com

# Iterated DNS query

www.google.com

Where is www.google.com?

- Contact server replies with the name of the next authority in the hierarchy
- "I don't know this name, but this other server might"
- This is how DNS works today

asgard.ccs.neu.edu

ns1.google.com

com

Root

# DNS Propagation

- How many of you have purchased a domain name?
  - Did you notice that it took ~72 hours for your name to become accessible?
  - This delay is called DNS Propagation

www.my-new-site.com

asgard.ccs.neu.edu → Root → com → ns.godaddy.com

- Why would this process fail for a new DNS name?

# Caching vs. Freshness

- DNS Propagation delay is caused by caching

Where is www.my-new-site.

That name does not exist.

- Cached Root Zone File
- Cached .com Zone File
- Cached .net Zone File
- Etc.

asgard.ccs.neu.edu

Root

com

- Zone files may be cached for 1-72 hours

www.my-new-site.com

ns.godaddy.com

# DNS Resource Records

- DNS queries have two fields: name and type

- Resource record is the response to a query

  - Four fields: (name, value, type, TTL)

  - There may be multiple records returned for one query

- What do the name and value mean?

  - Depends on the type of query and response

# DNS Types

- Type = A / AAAA
  - Name = domain name
  - Value = IP address
  - A is IPv4, AAAA is IPv6

**Query**
Name: www.ccs.neu.edu
Type: A

**Resp.**
Name: www.ccs.neu.edu
Value: 129.10.116.81

- Type = NS
  - Name = partial domain
  - Value = name of DNS server for this domain
  - "Go send your query to this other server"

**Query**
Name: ccs.neu.edu
Type: NS

**Resp.**
Name: ccs.neu.edu
Value: 129.10.116.51

# DNS Types, Continued

- Type = CNAME
  - Name = hostname
  - Value = canonical hostname
  - Useful for aliasing
  - CDNs use this

| Query | Name: foo.mysite.com<br>Type: CNAME |
|---|---|
| Resp. | Name: foo.mysite.com<br>Value: bar.mysite.com |

- Type = MX
  - Name = domain in email address
  - Value = canonical name of mail server

| Query | Name: ccs.neu.edu<br>Type: MX |
|---|---|
| Resp. | Name: ccs.neu.edu<br>Value: amber.ccs.neu.edu |

# Reverse Lookups

- What about the IP→name mapping?
- Separate server hierarchy stores reverse mappings
  - Rooted at in-addr.arpa and ip6.arpa
- Additional DNS record type: PTR
  - Name = IP address
  - Value = domain name
- Not guaranteed to exist for all IPs

| Query | Name: 129.10.116.51 Type: PTR |
|-------|-------------------------------|
| Resp. | Name: 129.10.116.51 Value: ccs.neu.edu |

# DNS as Indirection Service

- DNS gives us very powerful capabilities
  - Not only easier for humans to reference machines!

- Changing the IPs of machines becomes trivial
  - e.g. you want to move your web server to a new host
  - Just change the DNS record!

# Aliasing and Load Balancing

□ One machine can have many aliases

www.reddit.com

www.foursquare.com

www.huffingtonpost.com

david.choffnes.com

alan.mislo.ve

*.blogspot.com

□ One domain can map to multiple machines

www.google.com

# Content Delivery Networks

DNS responses may vary based on geography, ISP, etc

# Outline

- HTTP Connection Basics

- HTTP Protocol

- Cookies, keeping state + tracking

# Web and HTTP

*First, a review…*

- *web page* consists of *objects*

- object can be HTML file, JPEG image, Java applet, audio file,…

- web page consists of *base HTML-file* which includes *several referenced objects*

- each object is addressable by a *URL,* e.g.,

```
www.someschool.edu/someDept/pic.gif
```

host name          path name

# HTTP overview

## HTTP: hypertext transfer protocol

- Web's application layer protocol
- client/server model
  - *client:* browser that requests, receives, (using HTTP protocol) and "displays" Web objects
  - *server:* Web server sends (using HTTP protocol) objects in response to requests



PC running
Firefox browser

HTTP request

HTTP response

server running Apache Web server

HTTP request

HTTP response

iphone running
Safari browser

Application Layer

# HTTP overview (continued)

*uses TCP:*

- □ client initiates TCP connection (creates socket) to server, port 80

- □ server accepts TCP connection from client

- □ HTTP messages (application-layer protocol messages) exchanged between browser (HTTP client) and Web server (HTTP server)

- □ TCP connection closed

*HTTP is "stateless" (in theory...)*

- □ server maintains no information about past client requests

*aside*

protocols that maintain "state" are complex!

- ❖ past history (state) must be maintained

- ❖ if server/client crashes, their views of "state" may be inconsistent, must be reconciled

# HTTP connections

## non-persistent HTTP

- at most one object sent over TCP connection
  - ⬚ connection then closed
- downloading multiple objects required multiple connections

## persistent HTTP

- multiple objects can be sent over single TCP connection between client, server

Application Layer

# Example Web Page

**Harry Potter Movies**

As you all know, the new HP book will be out in June and then there will be a new movie shortly after that…

"Harry Potter and the Bathtub Ring"

hpface.jpg

castle.gif

page.html

Client          Server

**Non-Persistent HTTP**

TCP SYN

G

page.html

TCP FIN

TCP SYN

G

hpface.jpg

TCP FIN

TCP SYN

G

castle.gif

88 TCP FIN

The "classic" approach in HTTP/1.0 is to use one HTTP request per TCP connection, serially.

Client          Server          Concurrent (parallel) TCP
                                connections can be used
TCP SYN                         to make things faster.

G                    C              S    C              S

page.html        S                      S

TCP FIN          G                      G

                    hpface.jpg           castle.gif

                 F                      F

# Persistent HTTP

## non-persistent HTTP issues:

- requires 2 RTTs per object
- OS overhead for *each* TCP connection
- browsers often open parallel TCP connections to fetch referenced objects

## persistent HTTP:

- server leaves connection open after sending response
- subsequent HTTP messages between same client/server sent over open connection
- client sends requests as soon as it encounters a referenced object
- as little as one RTT for all the referenced objects

Application Layer

# Non-persistent HTTP: response time

RTT: time for a packet to travel from client to server and back

HTTP response time:

☐ one RTT to initiate TCP connection

☐ one RTT for HTTP request and first few bytes of HTTP response to return

  ▫ This assumes HTTP GET piggy backed on the ACK

☐ file transmission time

☐ non-persistent HTTP response time =

  2RTT+ file transmission time



initiate TCP connection

RTT

request file

RTT

time to transmit file

file received

time          time

Client          Server

**Persistent HTTP**

TCP SYN

G

page.html

G

hpface.jpg

G

castle.gif

Timeout

TCP FIN

The "persistent HTTP" approach can re-use the same TCP connection for Multiple HTTP transfers, one after another, serially. Amortizes TCP overhead, but maintains TCP state longer at server.

Client    Server

TCP SYN

G

page.html

GG

hpface.jpg

castle.gif

Timeout

The "pipelining" feature
in HTTP/1.1 allows
requests to be issued
asynchronously on a
persistent connection.
Requests must be
processed in proper order.
Can do clever packaging.

93 TCP FIN

# Outline

- HTTP Connection Basics

- HTTP Protocol

- Cookies, keeping state + tracking

# HTTP request message

□ two types of HTTP messages: *request, response*

□ HTTP request message:

  ▫ ASCII (human-readable format)

carriage return character

line-feed character

request line
(GET, POST,
HEAD commands)

```
GET /index.html HTTP/1.1\r\n
Host: www-net.cs.umass.edu\r\n
User-Agent: Firefox/3.6.10\r\n
Accept: text/html,application/xhtml+xml\r\n
Accept-Language: en-us,en;q=0.5\r\n
Accept-Encoding: gzip,deflate\r\n
Accept-Charset: ISO-8859-1,utf-8;q=0.7\r\n
Keep-Alive: 115\r\n
Connection: keep-alive\r\n
\r\n
```

header
lines

carriage return,
line feed at start
of line indicates
end of header lines

Application Layer

# HTTP request message: general format

| method | sp | URL | sp | version | cr | lf |
|--------|-----|-----|-----|---------|-----|-----|

request line

| header field name | | value | cr | lf |
|-------------------|---|-------|-----|-----|

~ ~

| header field name | | value | cr | lf |
|-------------------|---|-------|-----|-----|

header lines

| cr | lf |
|-----|-----|

| entity body |
|-------------|

~ entity body ~    body

Application Layer

# Uploading form input

## POST method:

- web page often includes form input

- input is uploaded to server in entity body

## URL method:

- uses GET method

- input is uploaded in URL field of request line:
  `www.somesite.com/animalsearch?monkeys&banana`

Application Layer

# Method types

## HTTP/1.0:

□ GET

□ POST

□ HEAD

  ▫ asks server to leave requested object out of response

## HTTP/1.1:

□ GET, POST, HEAD

□ PUT

  ▫ uploads file in entity body to path specified in URL field

□ DELETE

  ▫ deletes file specified in the URL field

Application Layer

# HTTP response message

status line
(protocol
status code
status phrase)

header
lines

data, e.g.,
requested
HTML file

```
HTTP/1.1 200 OK\r\n
Date: Sun, 26 Sep 2010 20:09:20 GMT\r\n
Server: Apache/2.0.52 (CentOS)\r\n
Last-Modified: Tue, 30 Oct 2007 17:00:02
   GMT\r\n
ETag: "17dc6-a5c-bf716880"\r\n
Accept-Ranges: bytes\r\n
Content-Length: 2652\r\n
Keep-Alive: timeout=10, max=100\r\n
Connection: Keep-Alive\r\n
Content-Type: text/html; charset=ISO-8859-
   1\r\n
\r\n
data data data data data ...
```

Application Layer

# HTTP response status codes

❖ status code appears in 1st line in server-to-client response message.

❖ some sample codes:

## 200 OK
- request succeeded, requested object later in this msg

## 301 Moved Permanently
- requested object moved, new location specified later in this msg (Location:)

## 400 Bad Request
- request msg not understood by server

## 404 Not Found
- requested document not found on this server

## 505 HTTP Version Not Supported

# Trying out HTTP (client side) for yourself

1. Telnet to your favorite Web server:

**telnet cis.poly.edu 80**  opens TCP connection to port 80 (default HTTP server port) at cis.poly.edu. anything typed in sent to port 80 at cis.poly.edu

2. type in a GET HTTP request:

**GET /~ross/ HTTP/1.1**
**Host: cis.poly.edu**

by typing this in (hit carriage return twice), you send this minimal (but complete) GET request to HTTP server

3. look at response message sent by HTTP server!

(or use Wireshark to look at captured HTTP request/response)

Outline

- HTTP Connection Basics
- HTTP Protocol
- Cookies, keeping state + tracking

# User-server state: cookies

many Web sites use cookies

*four components:*

1) cookie header line of HTTP *response* message

2) cookie header line in next HTTP *request* message

3) cookie file kept on user's host, managed by user's browser

4) back-end database at Web site

example:

☐ Susan always access Internet from PC

☐ visits specific e-commerce site for first time

☐ when initial HTTP requests arrives at site, site creates:

⬚ unique ID

⬚ entry in backend database for ID

Application Layer

# Cookies: keeping "state" (cont.)

client      server

**ebay 8734**

cookie file

usual http request msg →

Amazon server creates ID 1678 for user

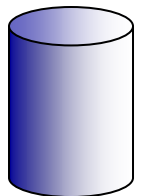← usual http response **set-cookie: 1678**

create entry

backend database

**ebay 8734
amazon 1678**

usual http request msg **cookie: 1678** →

cookie-specific action

access

← usual http response msg

one week later:

**ebay 8734
amazon 1678**

usual http request msg **cookie: 1678** →

access

cookie-specific action

Application Layer ← usual http response msg

# Cookies (continued)

## *what cookies can be used for:*

- authorization
- shopping carts
- recommendations
- user session state (Web e-mail)

## *how to keep "state":*

- ❖ protocol endpoints: maintain state at sender/receiver over multiple transactions
- ❖ cookies: http messages carry state

---

*aside*

*cookies and privacy:*

- ❖ cookies permit sites to learn a lot about you
- ❖ you may supply name and e-mail to sites

---

Application Layer

# Cookies + Third Parties

□ Example page (from Wired.com)



Elijah Wood in *Open Windows*. courtesy Cinedigm
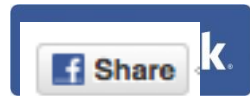
# How it works

And it's not just Facebook!



Wired.com

GET article.html

GET sharebutton.gif
Cookie: FBCOOKIE

Facebook now knows you visited this Wired article.
Works for all pages where 'like'/'share' button is embedded!