

Chapter 2

Application Layer

应用层

A note on the use of these Powerpoint slides:

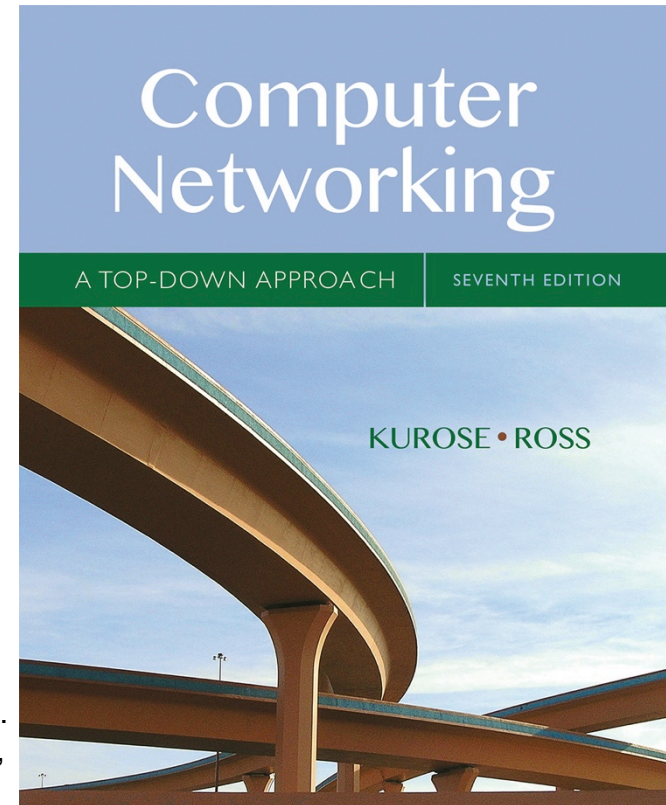
We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- If you use these slides (e.g., in a class) that you mention their source (after all, we'd like people to use our book!)
- If you post any slides on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

© All material copyright 1996-2016

J.F Kurose and K.W. Ross, All Rights Reserved



Computer Networking: A Top Down Approach

7th edition

Jim Kurose, Keith Ross

Pearson/Addison Wesley

April 2016

Some network apps

- e-mail (smtp, pop, IMAP)
- web (http)
- text messaging
- remote login (telnet)
- P2P file sharing
- multi-user network games
- streaming stored video (YouTube, Hulu, Netflix)
- voice over IP (e.g., Skype)
- real-time video conferencing
- social networking
- search (google, baidu)
- ...
- ...

学习有关网络应用的原理和实现方面的知识。

Chapter 2: application layer

our goals:

- conceptual, implementation aspects of network application protocols
 - transport-layer service models
 - client-server paradigm
 - process
 - peer-to-peer paradigm
 - content distribution networks
- learn about protocols by examining popular application-level protocols
 - HTTP
 - FTP
 - SMTP / POP3 / IMAP
 - DNS
- creating network applications
 - socket API

Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

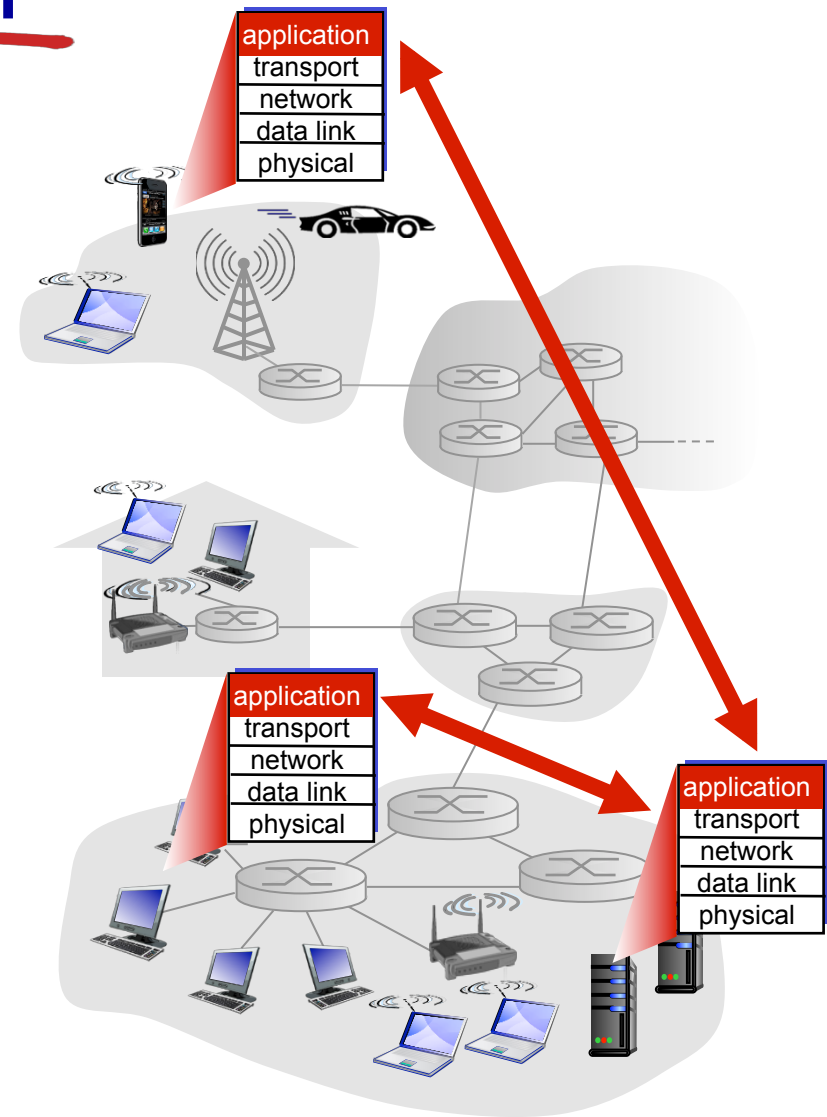
Creating a network app

write programs that:

- run on (different) *end systems*
- communicate over network
- e.g., web server software communicates with browser software

no need to write software
for network-core devices

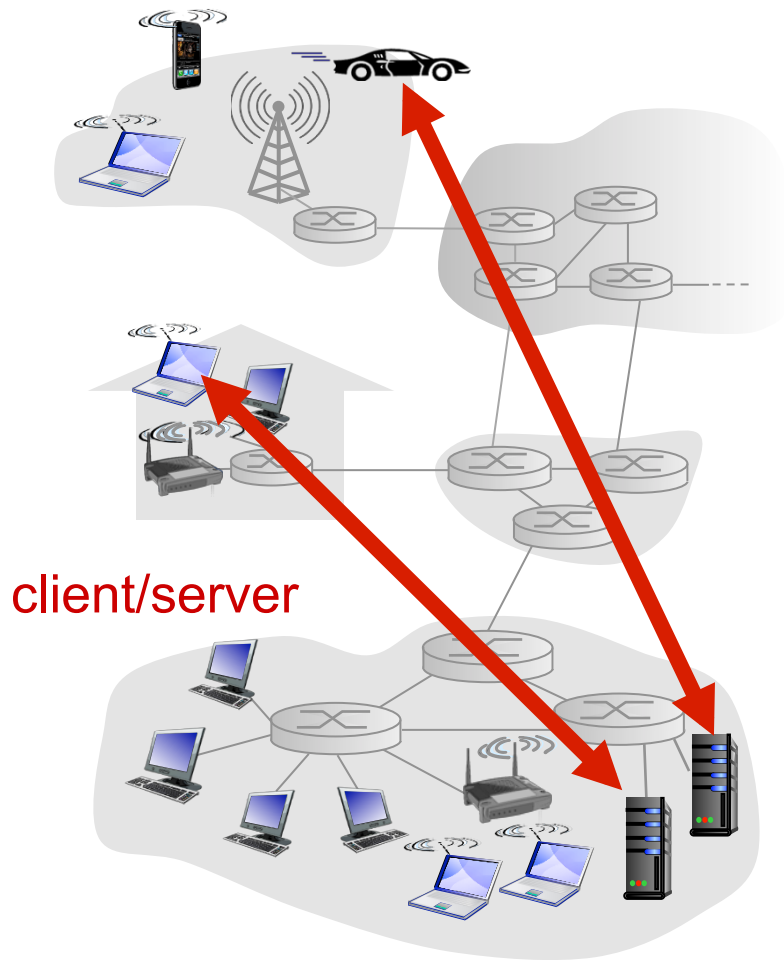
- network-core devices do not run user applications (1-3层)
- applications on end systems (1-5层) allows for rapid app development, propagation



2.1.1 Application architectures

- 从应用程序开发者的角度看，网络体系结构是固定的，并为应用程序提供了特定的服务集合。
- 应用程序体系结构由应用程序研发者设计，规定了如何在端系统上组织该应用程序。
- possible structure of applications:
 - client-server
 - peer-to-peer (P2P)

I. Client-server architecture



server:

- always-on host serving client
- permanent well-known IP address
- data centers for scaling
(搜索引擎、电子商务、电子邮件)
数十万台服务器, 开销大, 供电、维护、带宽...

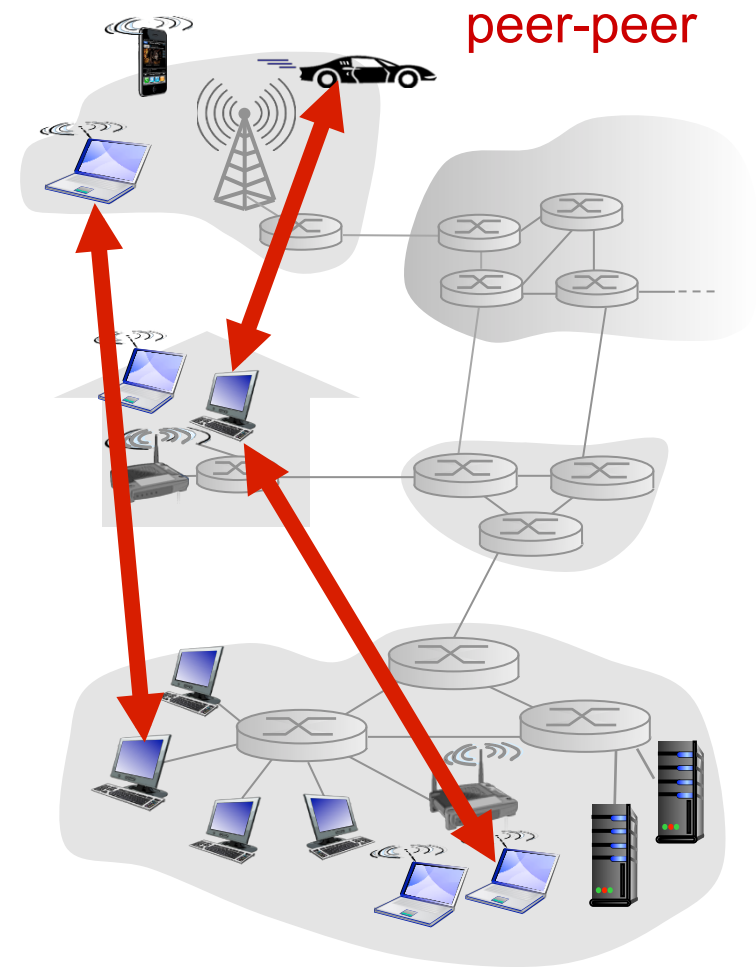
clients:

- communicate with server
- may be intermittently connected
- may have dynamic IP addresses
- do not communicate directly with each other

2. P2P architecture

对等模式

- **no always-on server**（或仅用于追踪客户）
- **arbitrary end systems directly communicate**
- 流量密集型应用，如文件共享、下载、视频会议...
- **peers request service from other peers, provide service in return to other peers**
 - **self scalability** – new peers bring new service capacity, as well as new service demands
- **peers are intermittently connected and change IP addresses**
 - complex management
 - 面临安全性、性能、可靠性挑战



2.1.2 Processes communicating

I. Client and Server Process

- **process**: program running within a host
- within same host, two processes communicate using **inter-process communication** (defined by OS)
- processes in different hosts communicate by exchanging **messages**

clients, servers

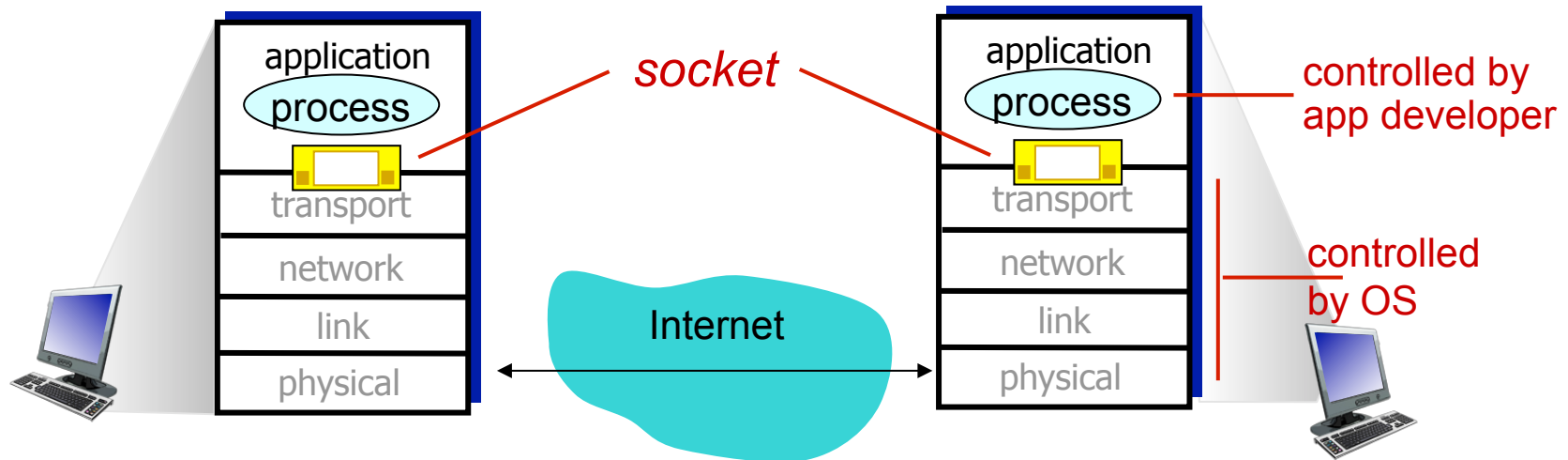
client process: process that initiates communication

server process: process that waits to be contacted

- aside: applications with P2P architectures have client processes & server processes
- 根据角色互为客户、服务器，一个进程可能既是客户机又是服务器

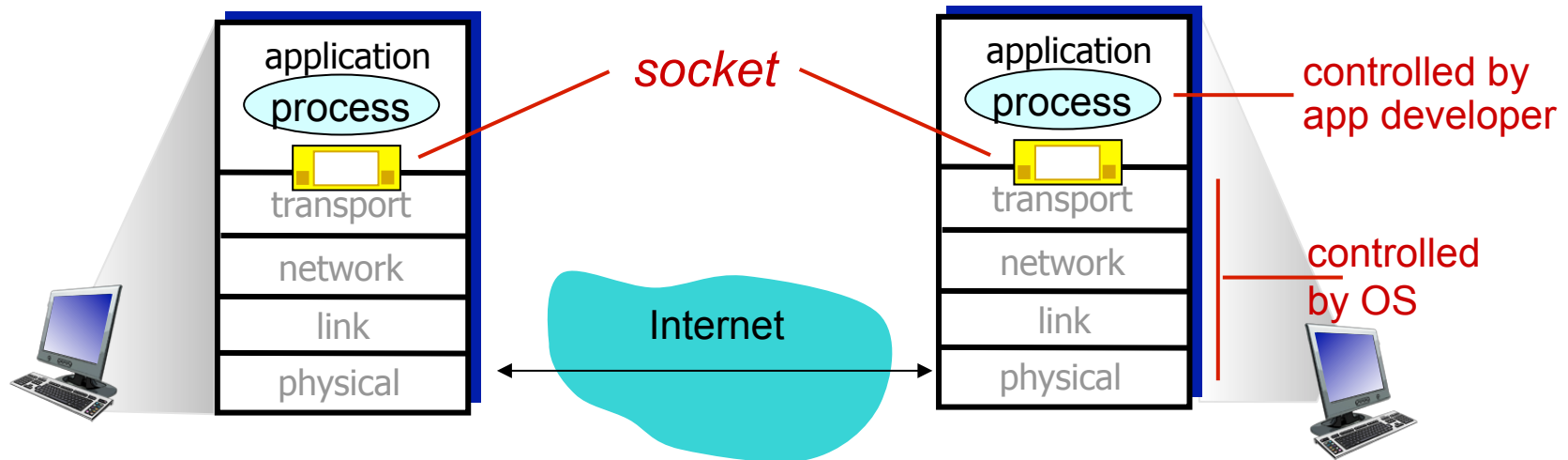
2. Sockets 套接字

- process sends/receives messages to/from its **socket**
- socket analogous to door
 - sending process shoves message out door
 - sending process relies on transport infrastructure on other side of door to deliver message to socket at receiving process



2. Sockets 套接字

- **套接字**，是同一台主机内应用层与传输层之间的接口，也称为应用程序和网络之间的应用程序编程接口API。
- 应用程序开发者可以控制套接字在应用层端的一切，对于套接字在传输层端的控制仅限于：
 - 选择传输层协议；
 - 设定几个传输层参数，如最大缓存和最大报文长度等。



3. Addressing processes 进程寻址

- to receive messages, process must have *identifier*
- host device has unique 32-bit IP address
- Q: does IP address of host on which process runs suffice for identifying the process?
- A: no, *many* processes can be running on same host
- *identifier* includes both **IP address** and **port numbers** associated with process on host.
- example port numbers:
 - HTTP server: 80
 - mail server: 25
- to send HTTP message to gaia.cs.umass.edu web server:
 - **IP address**: 128.119.245.12
 - **port number**: 80
- more shortly...

2.1.3 What transport service does an app need?

- 根据以下四方面应用程序服务要求，来选择传输层协议。

data integrity（正确、不丢）

- some apps (e.g., file transfer, web transactions) require 100% reliable data transfer
- other apps (e.g., audio) can tolerate some loss

timing（定时）

- some apps (e.g., Internet telephony, interactive games) require low delay(<100ms) to be “effective”

throughput（吞吐量）

- some apps (e.g., multimedia) require minimum amount of throughput to be “effective”
- other apps (“elastic apps”) make use of whatever throughput they get

security（安全性）

- encryption加密,解密, data integrity完整,站

Transport service requirements: common apps

application	data loss	throughput	time sensitive
file transfer	no loss	elastic	no
e-mail	no loss	elastic	no
Web documents	no loss	elastic	no
real-time audio/video	loss-tolerant	audio: 5kbps-1Mbps video: 10kbps-5Mbps	yes, 100' s msec
stored audio/video	loss-tolerant	same as above	yes, few secs
interactive games	loss-tolerant	few kbps up	yes, 100' s msec
text messaging	no loss	elastic	yes and no

2.1.4 Internet transport protocols services

TCP service:

- **reliable transport** between sending and receiving process (无差错、按序、无丢失和冗余)
- **flow control**: sender won't overwhelm receiver
- **congestion control**: throttle sender when network overloaded (限制发送进程)
- **does not provide**: timing, minimum throughput guarantee, security (无加密)
- **connection-oriented**: setup required between client and

UDP service:

- **unreliable data transfer** between sending and receiving process (差错、丢失、乱序、冗余)
- **does not provide**: reliability, flow control, congestion control, timing, throughput guarantee, security, or connection setup,

Q: why bother? Why is there a UDP?

Securing TCP (TCP 安全性)

TCP & UDP

- no encryption
- cleartext passwds sent into socket traverse Internet in cleartext

SSL (Secure Sockets Layer 安全套接字层)

- provides encrypted TCP connection 加密
- data integrity 数据完整性
- end-point authentication 端点鉴别

SSL is at app layer

- SSL不是传输层第三种协议
- apps use SSL libraries, that “talk” to TCP
- SSL sockets API works like TCP sockets API

SSL socket API

- cleartext passwords sent into socket traverse Internet encrypted, then transfer to TCP socket
- see Chapter 8

Internet apps: application, transport protocols

application	application layer protocol	underlying transport protocol
e-mail	SMTP [RFC 2821]	TCP
remote terminal access	Telnet [RFC 854]	TCP
Web	HTTP [RFC 2616]	TCP
file transfer	FTP [RFC 959]	TCP
streaming multimedia	HTTP (e.g., YouTube), RTP [RFC 1889]	TCP or UDP
Internet telephony	SIP, RTP, proprietary (e.g., Skype)	TCP or UDP

- 因特网能够为时延敏感的应用提供满意的服务，但是并不提供任何定时或带宽的保证。
- 因特网电话等容丢失应用喜欢**UDP**，但是防火墙经常配置成阻挡**UDP**流量，所以用**TCP**做备份。

2.1.5 App-layer protocol defines

- 应用层协议定义了以下方面:
- **types of messages exchanged**,
 - e.g., request, response
- **message syntax**:
 - what fields in messages & how fields are delineated
- **message semantics**
 - meaning of information in fields
- **rules** for when and how processes send & respond to messages

open protocols:

公共域协议

- defined in RFCs
- allows for interoperability
- e.g., HTTP, SMTP

proprietary protocols:

专用协议

- e.g., Skype

应用层协议是网络应用的一部分。

- 如: web应用包含HTML文档格式、浏览器、服务器、应用层协议HTTP

Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

Web and HTTP

First, a review... Web页面的组成

- *web page* consists of *objects*
- *object* can be HTML file, JPEG image, Java applet, audio file,...
基本文件 和 引用对象
- 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`

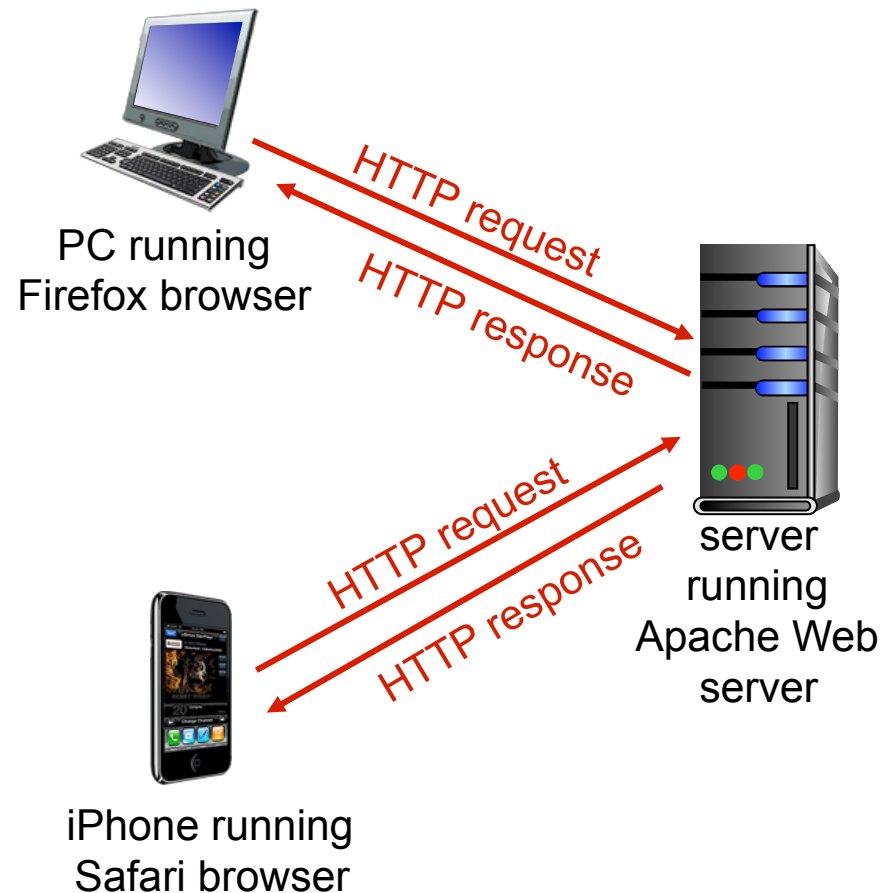
host name

path name

2.2.1 HTTP overview

超文本传输协议 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



HTTP overview (continued)

uses TCP 作为支撑传输协议:

- client initiates TCP connection (creates socket) to server, port 80 (服务器端)
- 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
- 报文通过套接字发送出去后, 即脱离客户控制并进入TCP的控制。

HTTP is “stateless”

- 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

2.2.2 HTTP connections

non-persistent HTTP

非持续连接

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

persistent HTTP

持续连接

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

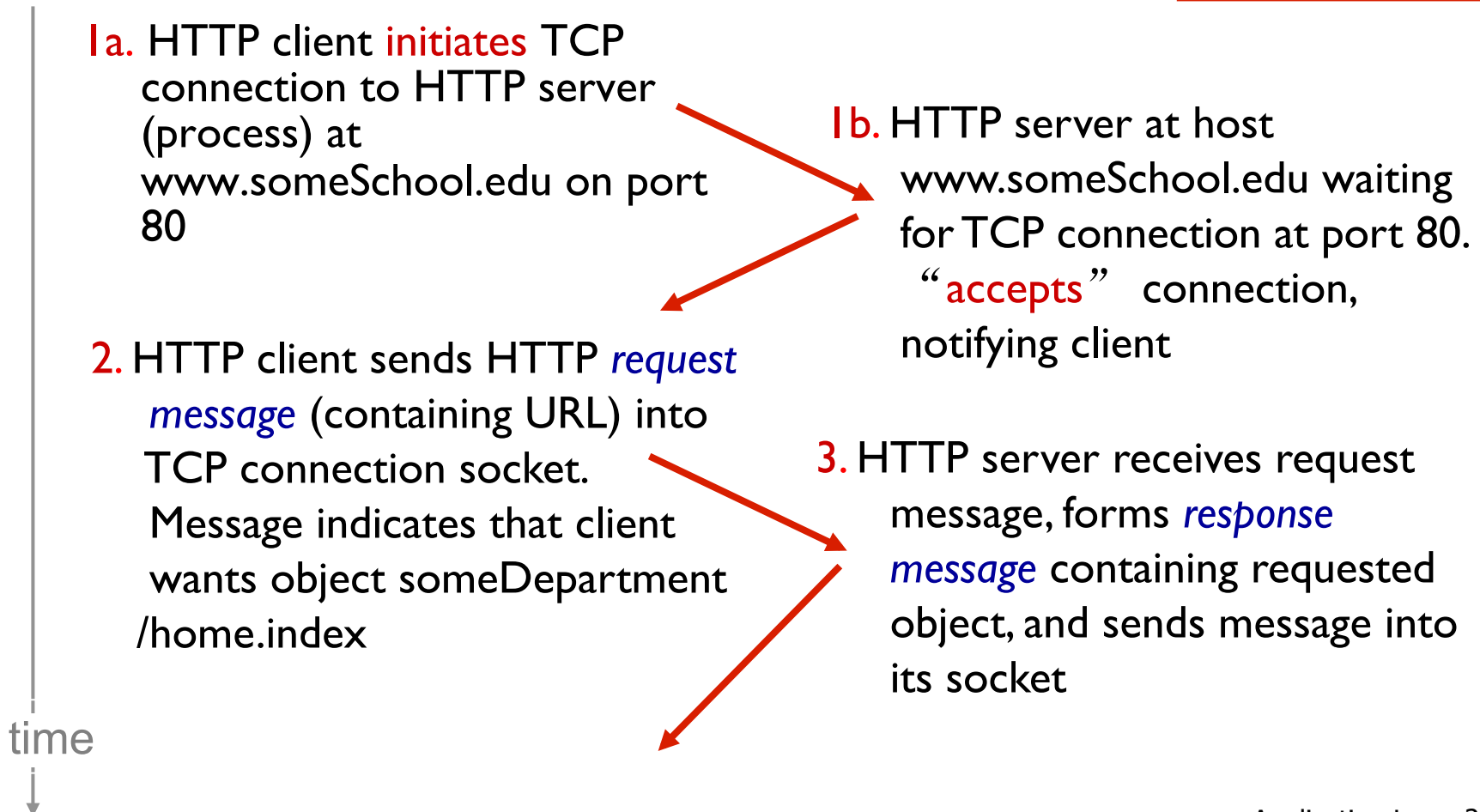
HTTP在默认方式下使用持续连接，如果需要也可以配置为非持续连接。

I. Non-persistent HTTP

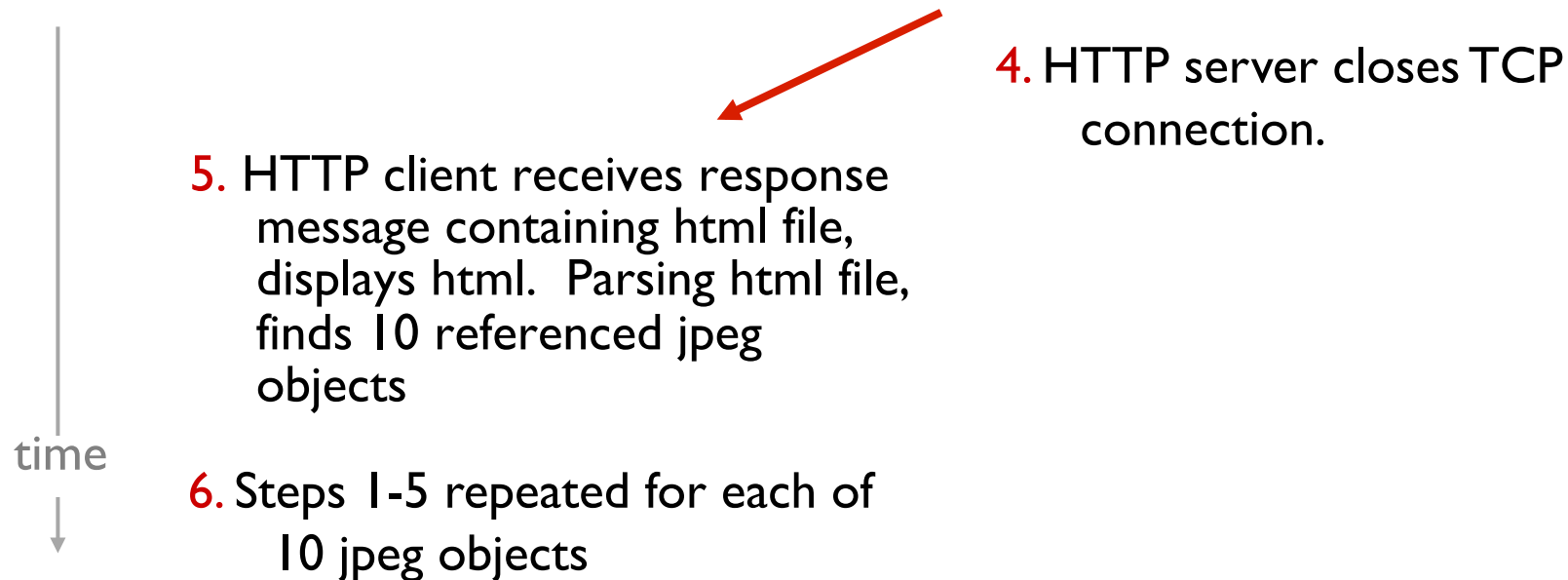
suppose user enters URL:

`www.someSchool.edu/someDepartment/home.index`

(contains text,
references to 10
jpeg images)



Non-persistent HTTP (cont.)



在非持续连接情况下，每个TCP服务器发送一个对象后就关闭。在此例中要产生11个TCP连接。

现代浏览器能够用并行的方式工作，默认情况下打卡5~10个并行的TCP连接。

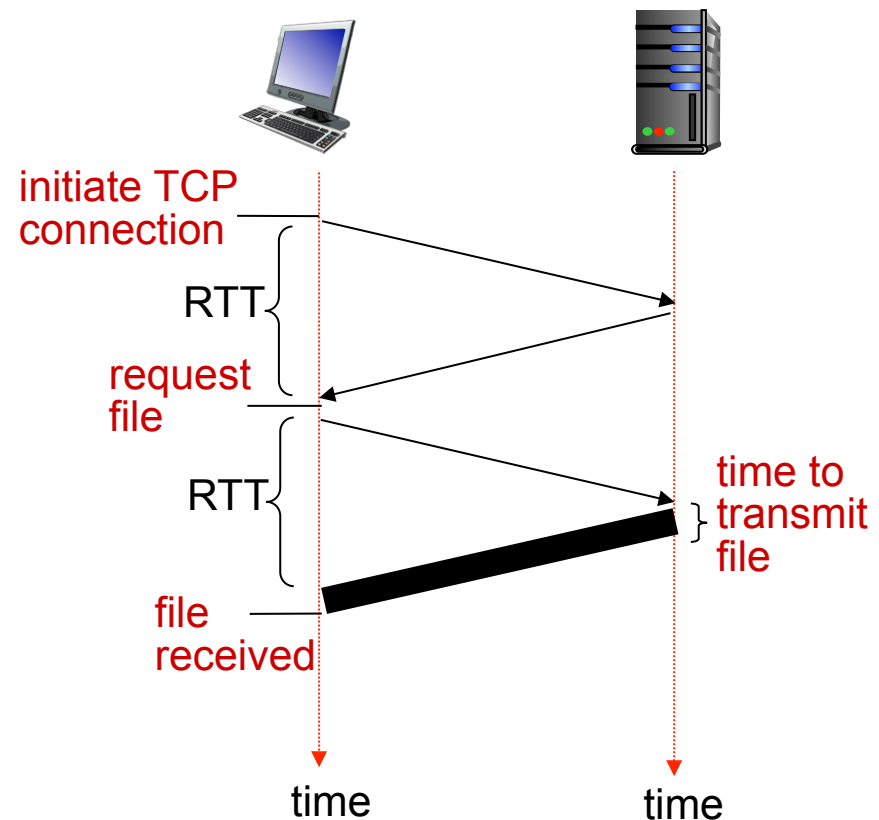
Non-persistent HTTP: response time

RTT (往返时间): time for a small packet to travel from client to server and back.

(分组传播时延、分组在中间节点排队时延、分组处理时延)

HTTP response time:

- one RTT to initiate TCP connection
- one RTT for HTTP request and first few bytes of HTTP response to return
- file transmission time
- **non-persistent HTTP response time = $2RTT + \text{file transmission time}$**



2. Persistent HTTP

non-persistent HTTP issues:

- requires 2 RTTs **per object**
一个建连接，一个传对象；
- 含10个对象的网页一共需要22 RTTs。
- 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**
- 含10个对象的网页一共需要2 RTTs。
- 通常配置一个超时间隔，一条连接超过这个间隔未使用，服务器就关闭这个连接。

2.2.3 HTTP request message

- two types of HTTP messages: *request, response*
- **HTTP request message:**
 - ASCII (human-readable format)

request line
(GET, POST,
HEAD commands)

header
lines

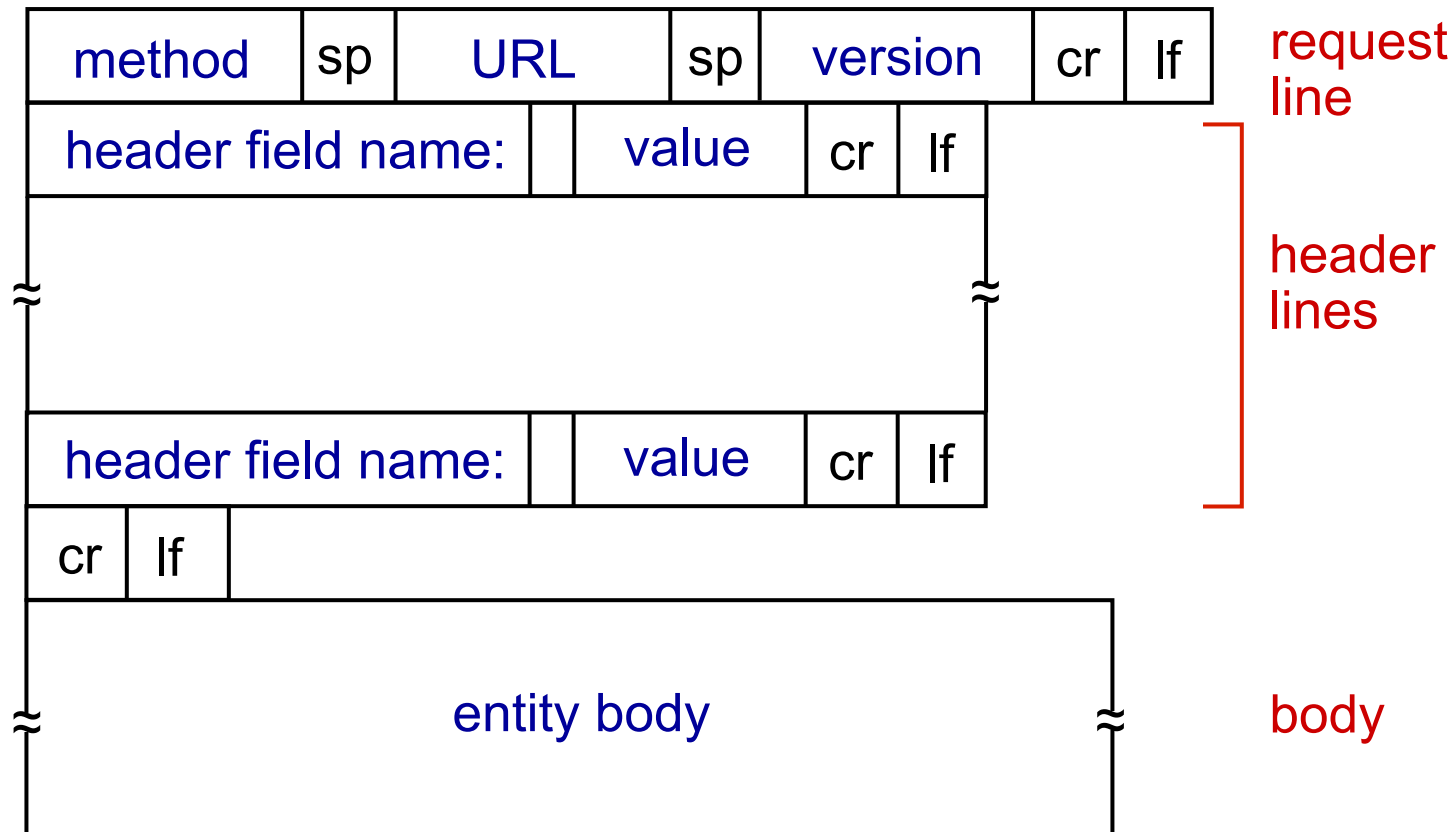
carriage return,
line feed at start
of line indicates
end of header lines

```
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
```

carriage return character
line-feed character

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

HTTP request message: general format



Uploading form input 上传表单数据

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`

Method types

HTTP/1.0:

- GET
 - 表单数据添加在URL末尾
- POST
 - 表单数据放在实体字段中
- HEAD
 - 只发送响应报文，而不返回请求对象，用于应用程序的调试追踪。

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

2.2.3 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 ...
```

* Check out the online interactive exercises for more
examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

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:) 客户软件将自动获取新的URL

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 gaia.cs.umass.edu 80` { opens TCP connection to port 80
(default HTTP server port)
at gaia.cs.umass.edu.
anything typed in will be sent
to port 80 at gaia.cs.umass.edu

2. type in a GET HTTP request:

`GET /kurose_ross/interactive/index.php HTTP/1.1`
`Host: gaia.cs.umass.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)

2.2.4 User-server state: cookies

- HTTP是无状态的，但是很多应用中服务器希望识别用户身份进一步提供服务。

many Web sites use cookies

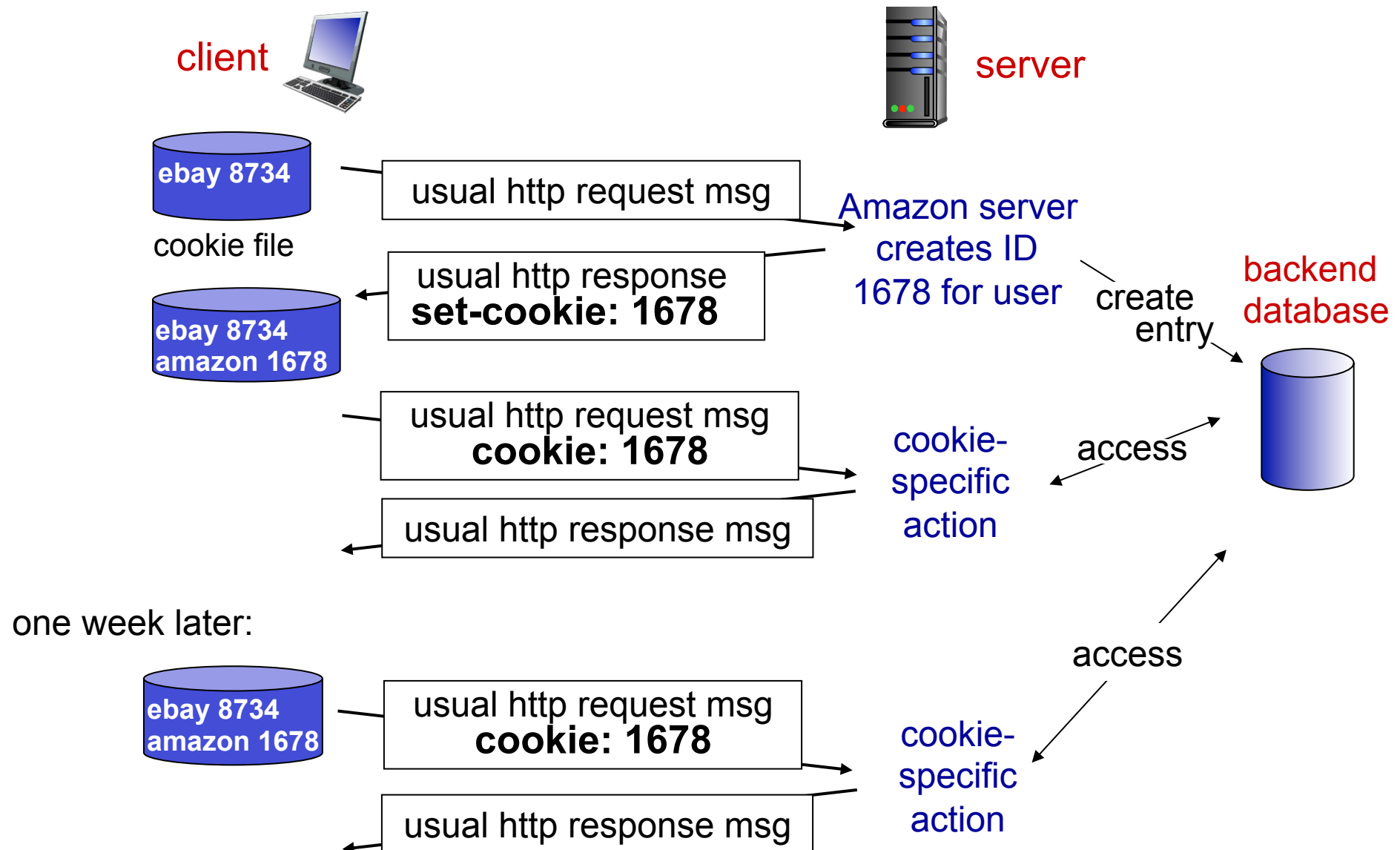
four components:

- 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 (session ID)
 - entry in backend database for ID
 - send response message to client browser including this ID in Set-cookie field

Cookies: keeping “state” (cont.)



Cookies (continued)

what cookies can be used for:

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

cookies and privacy: aside

- cookies permit sites to learn a lot about you
- you may supply name and e-mail to sites
- don't supply credit information to cookies

how to keep “state” :

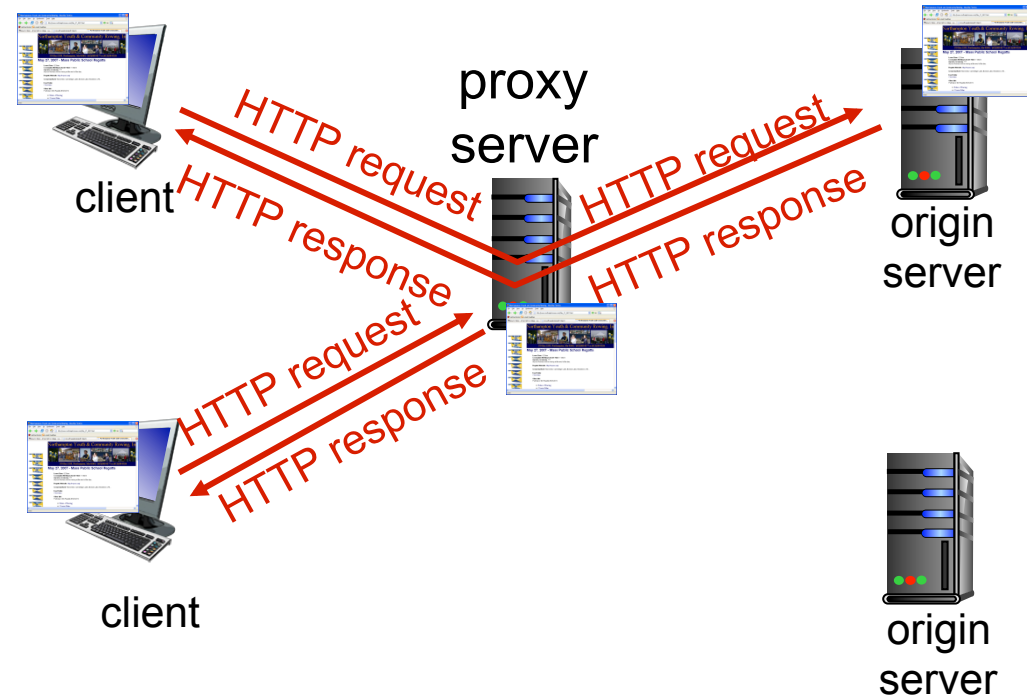
- protocol endpoints: maintain state at sender/receiver over multiple transactions
- cookies: http messages carry state

2.2.5 Web caches (proxy server)

goal: satisfy client request without involving origin server

Web缓存的存储空间中保存最近请求过的对象的副本。

- user sets browser: Web accesses via cache
- browser sends all HTTP requests to cache
 - **object in cache:** cache returns object
 - **else** cache requests object from origin server, then returns object to client



More about Web caching

- cache acts as both client and server
 - server for original requesting client
 - client to origin server
- typically cache is installed by ISP (university, company, residential ISP)

why Web caching?

- **reduce response time** for client request
- **reduce traffic** on an institution's access link
- **reduce Internet dense** with caches: enables “poor” content providers to effectively deliver content (so too does P2P file sharing)

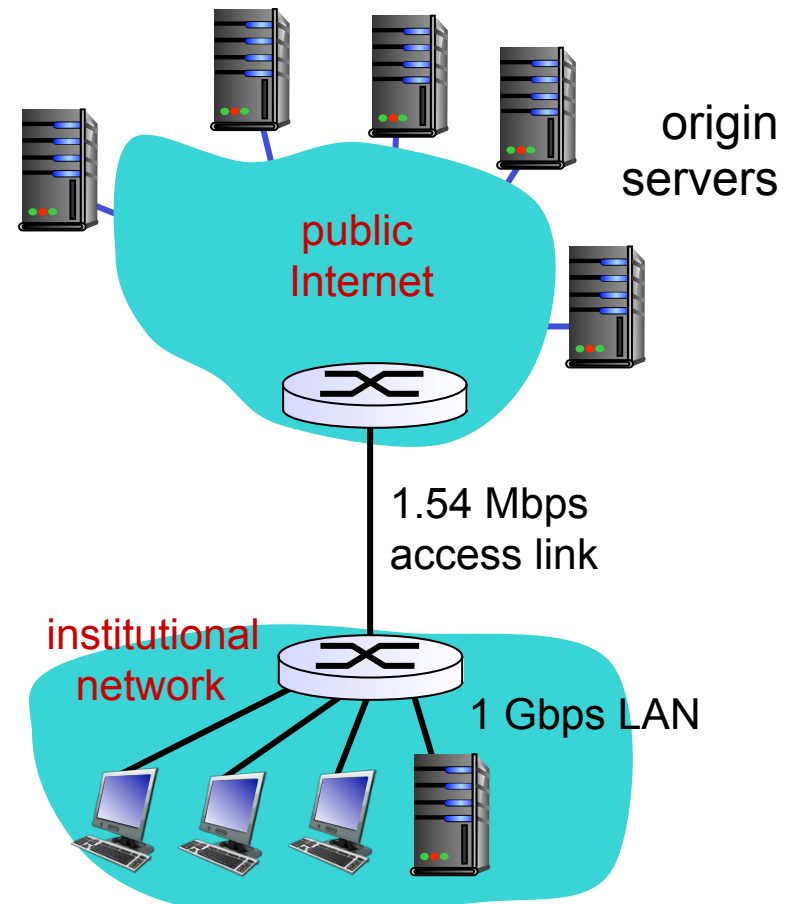
Caching example:

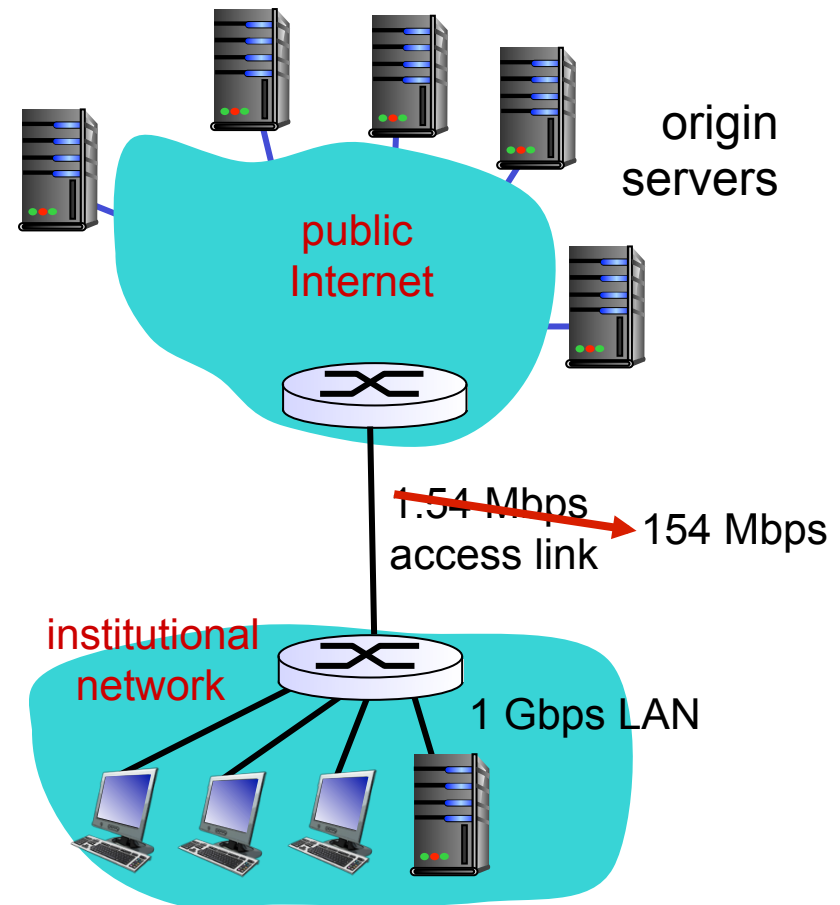
assumptions:

- avg object size: 100K bits
- avg request rate from browsers to origin servers: 15/sec
- avg data rate to browsers: 1.50 Mbps
- RTT from institutional router to any origin server: 2 sec
- access link rate: 1.54 Mbps

consequences:

- LAN utilization: 0.15%
- access link utilization = **99%** *problem!*
- total delay = Internet delay + access delay + LAN delay
= 2 sec + minutes + usecs





Caching example: install local cache

assumptions:

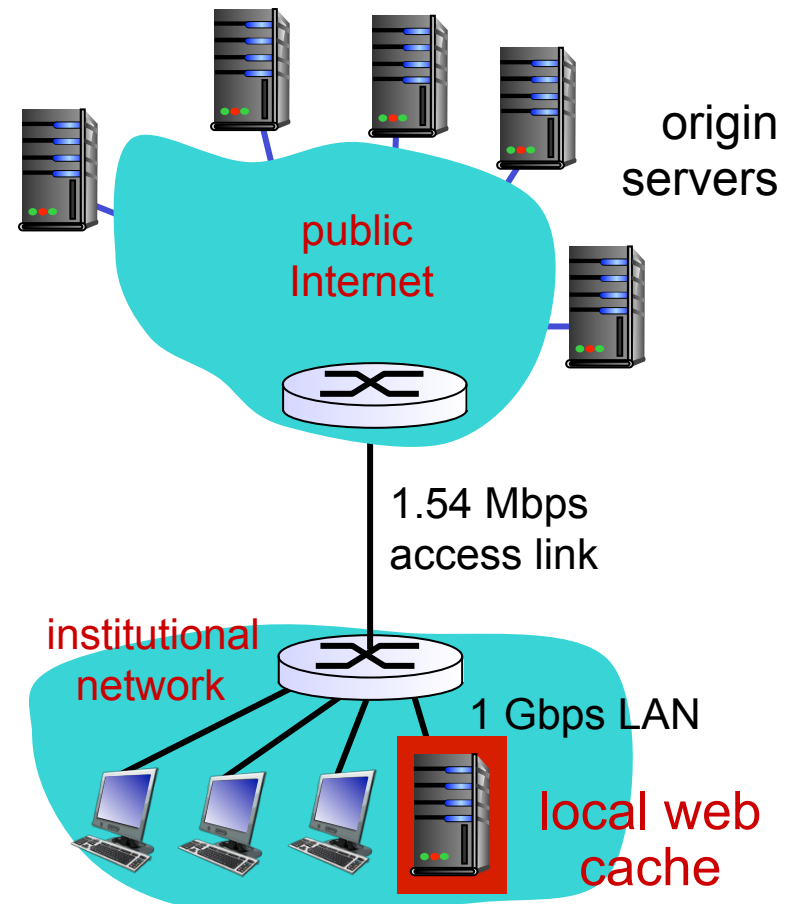
- avg object size: 100K bits
- avg request rate from browsers to origin servers: 15/sec
- avg data rate to browsers: 1.50 Mbps
- RTT from institutional router to any origin server: 2 sec
- access link rate: 1.54 Mbps

consequences:

- LAN utilization: 0.15%
- access link utilization = ?
- total delay = ?

How to compute link utilization, delay?

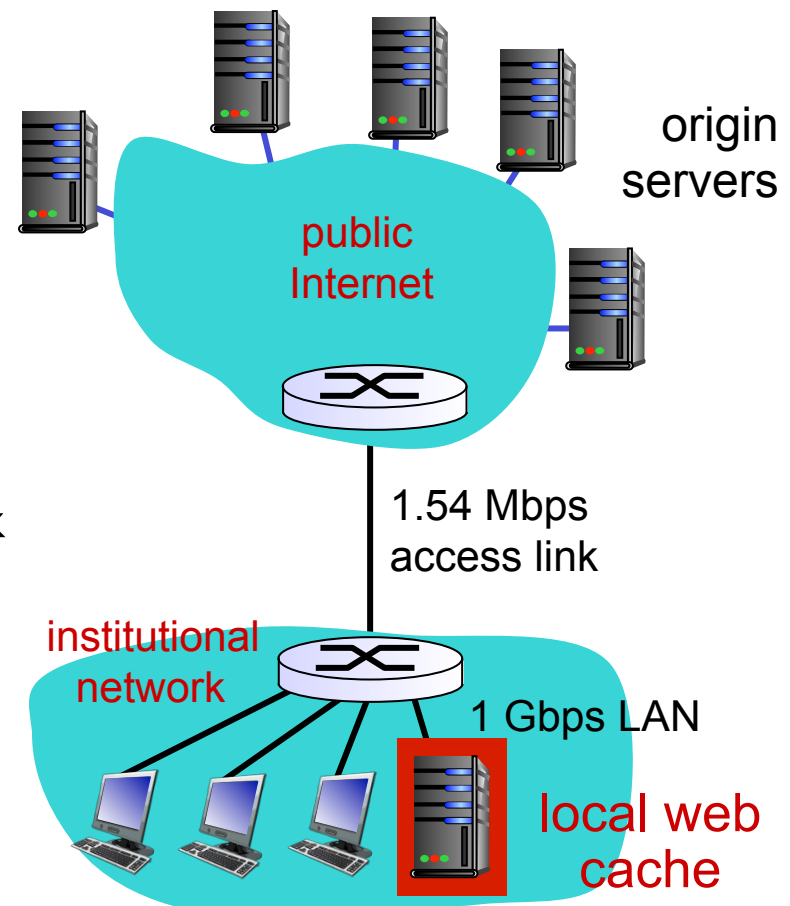
Cost: web cache (cheap!)



Caching example: install local cache

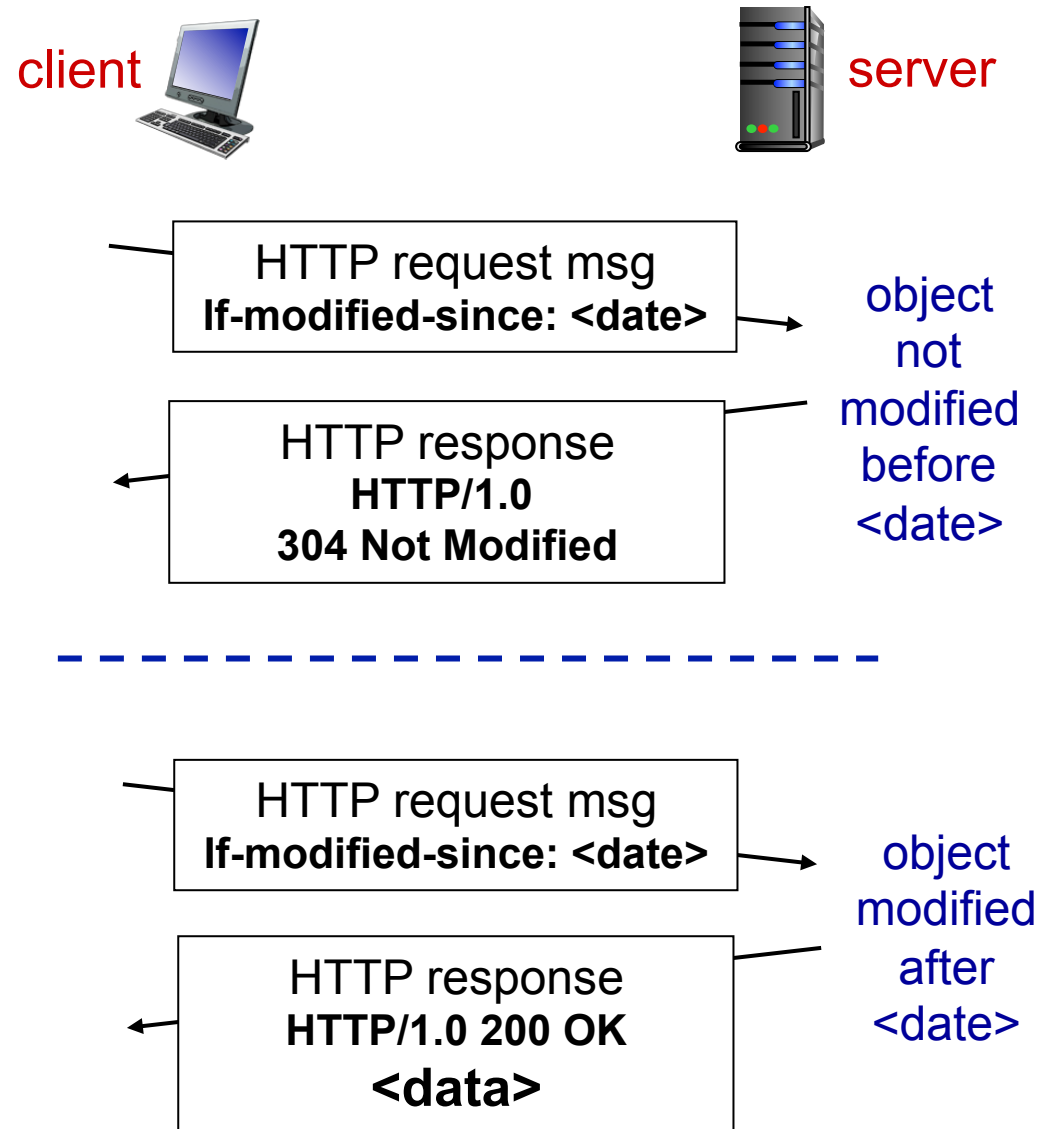
Calculating access link utilization, delay with cache:

- suppose cache hit rate is 0.4
 - 40% requests satisfied at cache,
 - 60% requests satisfied at origin
- access link utilization:
 - 60% of requests use access link
- data rate to browsers over access link
 $= 0.6 * 1.50 \text{ Mbps} = 0.9 \text{ Mbps}$
 - utilization $= 0.9 / 1.54 = 0.58$
- total delay
 - $= 0.6 * (\text{delay from origin servers}) + 0.4 * (\text{delay when satisfied at cache})$
 - $= 0.6 (2.01) + 0.4 (\sim \text{msecs}) = \sim 1.2 \text{ secs}$
 - less than with 154 Mbps link (and cheaper too!)



2.2.6 Conditional GET

- **Goal:** 检验web缓存中保存的对象副本是否陈旧过期
- **条件GET方法:**
 - 请求报文使用GET方法
 - 请求报文中含有 If-modified-since 首部行
- **cache:** specify date of cached copy in HTTP request
If-modified-since: <date>
- **server:** response contains no object if cached copy is up-to-date:
HTTP/1.0 304 Not Modified



Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

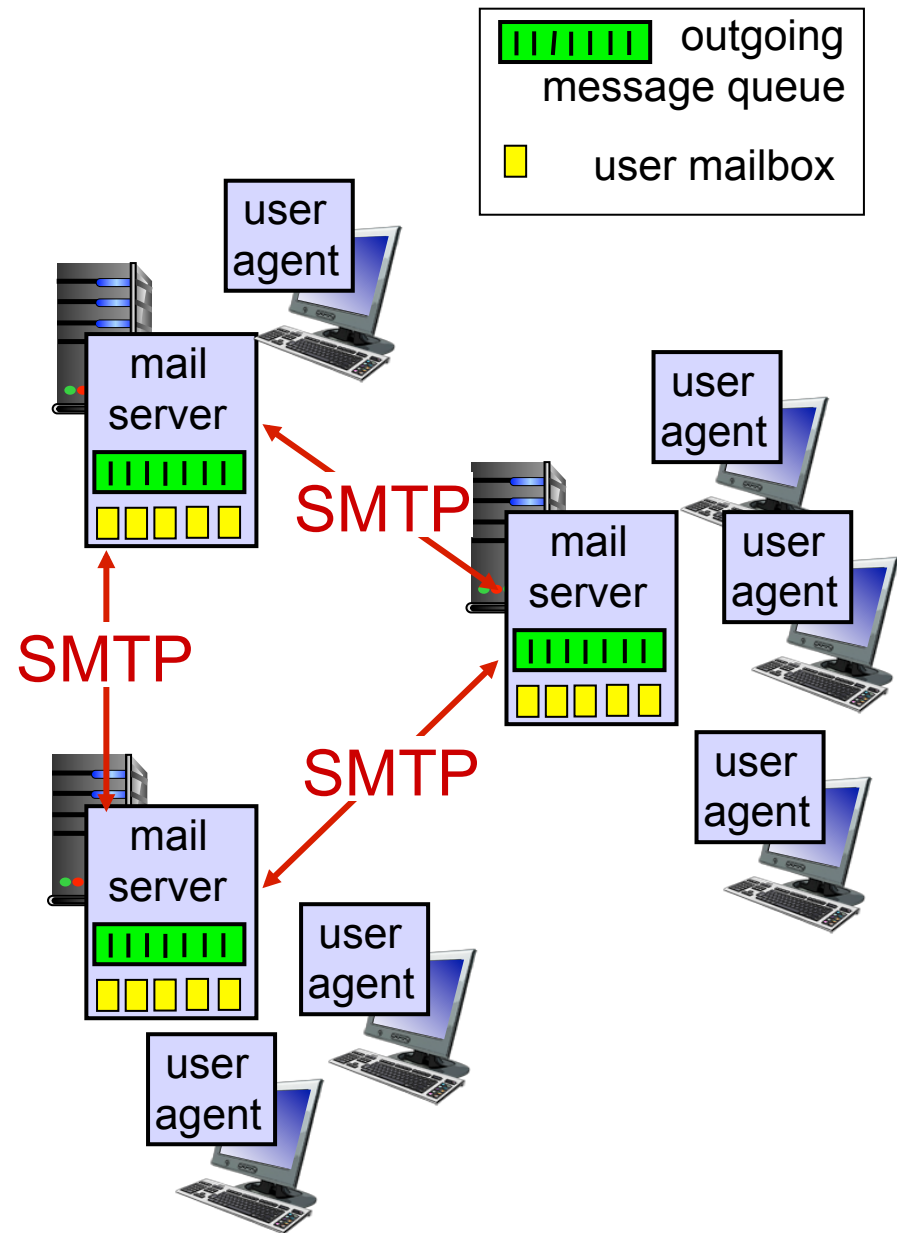
Electronic mail

Three major components:

- user agents 用户代理
- mail servers 邮件服务器
- simple mail transfer protocol: SMTP 简单邮件传输协议

User Agent

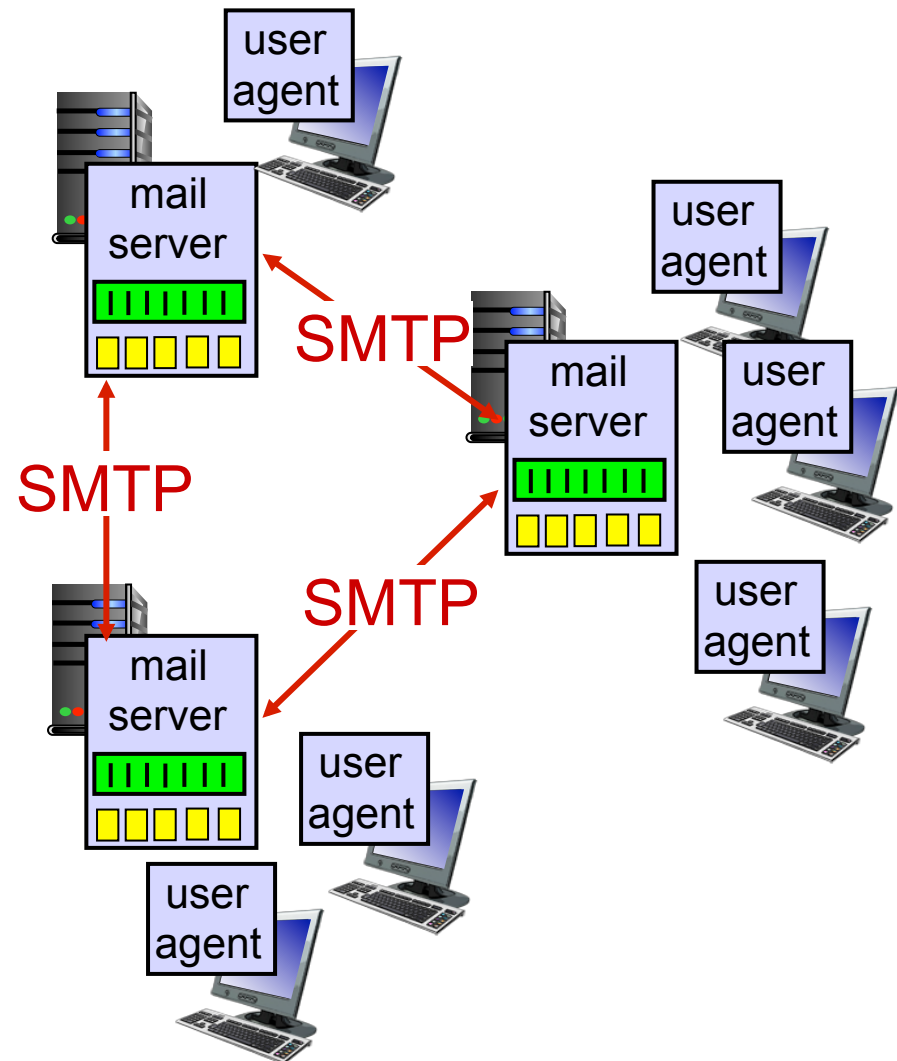
- a.k.a. “mail reader”
- composing, editing, reading mail messages
- e.g., Outlook, Thunderbird, iPhone mail client
- outgoing, incoming messages stored on server



Electronic mail: mail servers

mail servers: 核心

- **mailbox** contains incoming messages for user
- **message queue** of outgoing (to be sent) mail messages
- **SMTP protocol** between mail servers to send email messages
 - client: sending mail server
 - server: receiving mail server
- 邮件发送过程:
 - 邮件从发送方代理开始;
 - 传输到发送方邮件服务器 (进入外出报文队列, 30min);
 - 再传输到接收方邮件服务器;

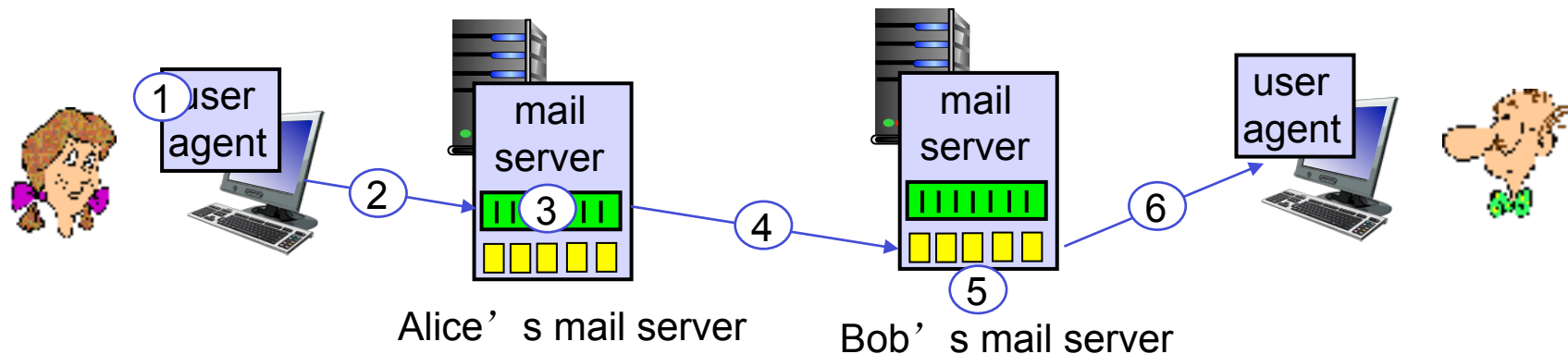


2.3.1 Electronic Mail: SMTP [RFC 2821]

- uses TCP to reliably transfer email message from client（发送方邮件服务器） to server（接收方邮件服务器），port 25
- **direct transfer:** sending server to receiving server
- **three phases of transfer**
 - handshaking (greeting)
 - transfer of messages
 - closure
- command/response interaction (like HTTP)
 - **commands:** ASCII text
 - **response:** status code and phrase
- messages must be in 7-bit ASCII，现代应用中需要将发送数据编码为ASCII码，到达对方在解码还原为多媒体数据

Scenario: Alice sends message to Bob

- 1) Alice uses UA to compose message “to”
`bob@someschool.edu`
- 2) Alice’s UA sends message to her mail server; message placed in message queue
- 3) client side of SMTP opens TCP connection with Bob’s mail server
- 4) SMTP client sends Alice’s message over the TCP connection
- 5) Bob’s mail server places the message in Bob’s mailbox
- 6) Bob invokes his user agent to read message



Sample SMTP interaction

Tcp连接建立后。。。持续连接在一个TCP上可发多条报文

```
S: 220 hamburger.edu
C: HELO crepes.fr
S: 250 Hello crepes.fr, pleased to meet you
C: MAIL FROM: <alice@crepes.fr>
S: 250 alice@crepes.fr... Sender ok
C: RCPT TO: <bob@hamburger.edu>
S: 250 bob@hamburger.edu ... Recipient ok
C: DATA
S: 354 Enter mail, end with "." on a line by itself
C: Do you like ketchup?
C: How about pickles?
C: .
S: 250 Message accepted for delivery
C: QUIT
S: 221 hamburger.edu closing connection
```

Try SMTP interaction for yourself:

- `telnet servername 25`
- see 220 reply from server
- enter HELO, MAIL FROM, RCPT TO, DATA, QUIT commands

above lets you send email without using email client (reader)

2.3.2 comparison with HTTP

- SMTP uses persistent connections
- SMTP requires message (header & body) to be in 7-bit ASCII
- SMTP server uses CRLF.CRLF to determine end of message

comparison with HTTP:

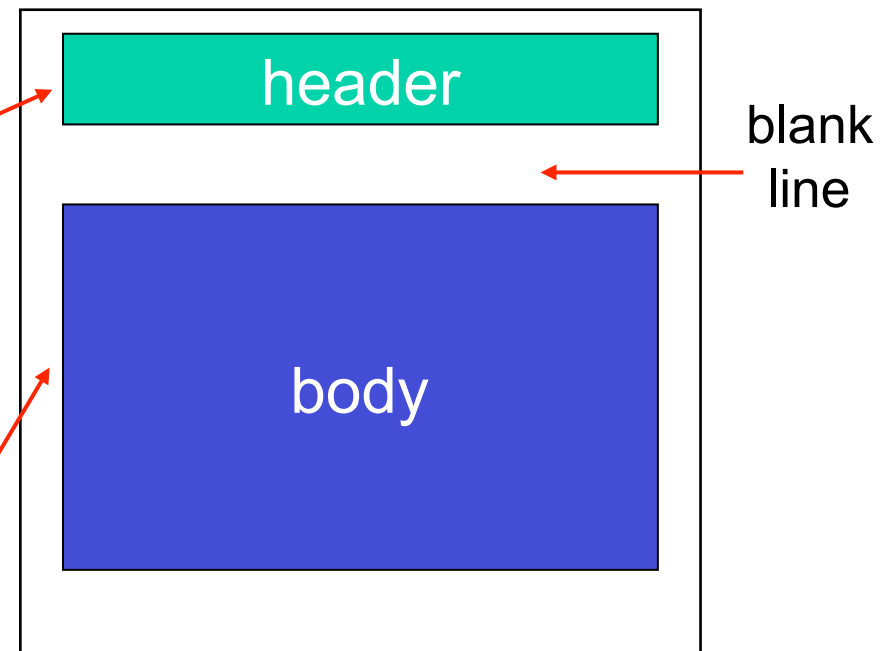
- HTTP: pull
- SMTP: push
- both have ASCII command /response interaction, status codes
- SMTP must be 7-bit ASCII
- HTTP: each object encapsulated in its own response message
- SMTP: multiple objects sent in one multipart message

2.3.3 Mail message format

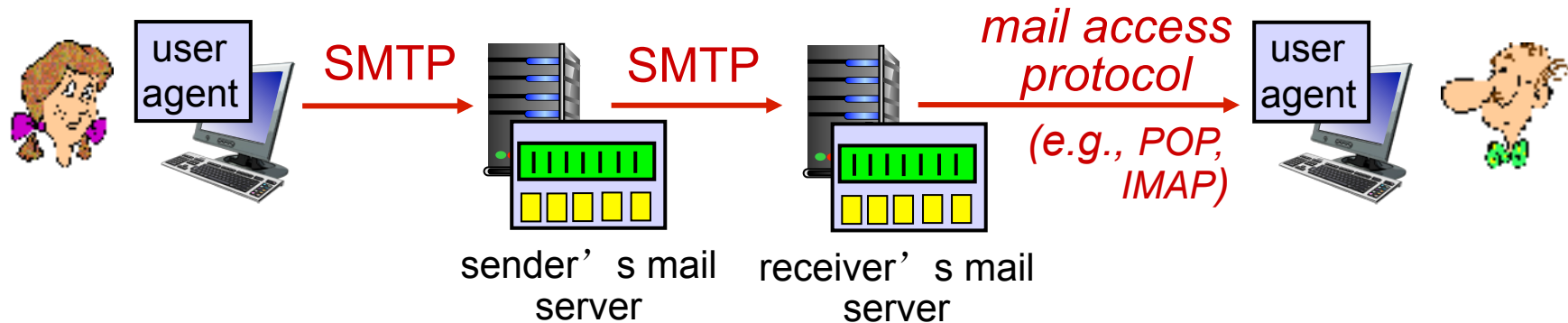
SMTP: protocol for exchanging email messages

RFC 822: standard for text message format:

- header lines, e.g.,
 - To:
 - From:
 - Subject: 可选
- different* from SMTP MAIL FROM, RCPT TO: commands!
- blank line
- Body: the “message”
 - **ASCII characters only**



2.3.4 Mail access protocols



- **SMTP**: delivery/storage to receiver's server
- mail access protocol: retrieval from server
 - **POP**: Post Office Protocol [RFC 1939]: authorization, download
 - **IMAP**: Internet Mail Access Protocol [RFC 1730]: more features, including manipulation of stored messages on server
 - **HTTP**: gmail, Hotmail, Yahoo! Mail, etc.

I. POP3 protocol

authorization phase

- client commands:
 - **user**: declare username
 - **pass**: password
- server responses
 - **+OK**
 - **-ERR**

transaction phase, client:

- **list**: list message numbers
- **retr**: retrieve message by number
- **delete**: delete 标记删除
- **quit** 退出后真正删除

```
S: +OK POP3 server ready
C: user bob
S: +OK
C: pass hungry
S: +OK user successfully logged on
```

```
C: list
S: 1 498
S: 2 912
S: .
C: retr 1
S: <message 1 contents>
S: .
C: dele 1
C: retr 2
S: <message 1 contents>
S: .
C: dele 2
C: quit
S: +OK POP3 server signing off
```

2. POP3 (more) and IMAP

more about POP3

- previous example uses POP3 “download and delete” mode
 - Bob cannot re-read e-mail if he changes client
- POP3 “download-and-keep” : copies of messages on different clients
- POP3 is stateless across sessions

IMAP

- keeps all messages in one place: at server
- allows user to organize messages in folders
- keeps user state across sessions:
 - names of folders and mappings between message IDs and folder name

Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

DNS: domain name system

people: many identifiers:

- SSN, name, passport #

Internet hosts, routers:

- **IP address** (32 bit)
 - used for addressing datagrams,
0~255的点分十进制数
- “**name**”, e.g.,
www.yahoo.com
 - used by humans

Q: how to map between IP address and name, and vice versa ?

Domain Name System

域名系统:

- **distributed database**
implemented in hierarchy of many *name servers*
- **application-layer protocol:** hosts, name servers communicate to **resolve** names (address/name translation)
 - note: core Internet function, implemented as application-layer protocol
 - complexity at network's “edge”
- 运行在UDP上，使用53号端口
Application Layer 2-58

2.4.1 DNS services

- DNS通常是由其它应用层协议使用的，如HTTP、SMTP、FTP，将用户提供的主机名解析为IP地址。

(1)hostname to IP address translation

(2)host aliasing 主机别名

- canonical hostname 规范主机名(比较复杂难记)
- 调用DNS获得alias names(更好记)对应的规范主机名或IP

(3)mail server aliasing 邮件服务器别名

- 电子邮件应用程序调用DNS获得邮件服务器别名对应的规范主机名或IP地址

(4)load distribution 负载均衡

- replicated Web servers 冗余服务器: many IP addresses correspond to one name
- 当用户发出对某繁忙站点域名的DNS请求时，DNS服务器用整个IP集合响应并每次都循环这个集合，用户向集合里IP地址排在前面的服务器发送请求。

2.4.2 DNS 工作机理

DNS 工作过程

- ①用户主机上的应用程序需要将一个主机名转换为IP地址;
- ②用户主机调用**DNS**客户端, 并指明需要转换的主机名;
- ③用户主机**DNS**客户端向网络发送一个**DNS**查询报文;
- ④经过若干秒到若干毫秒的时延, 用户主机**DNS**客户端收到一个回答报文, 回答报文中保护映射结果;
- ⑤**DNS**客户端将映射结果传递到调用**DNS**的应用程序。

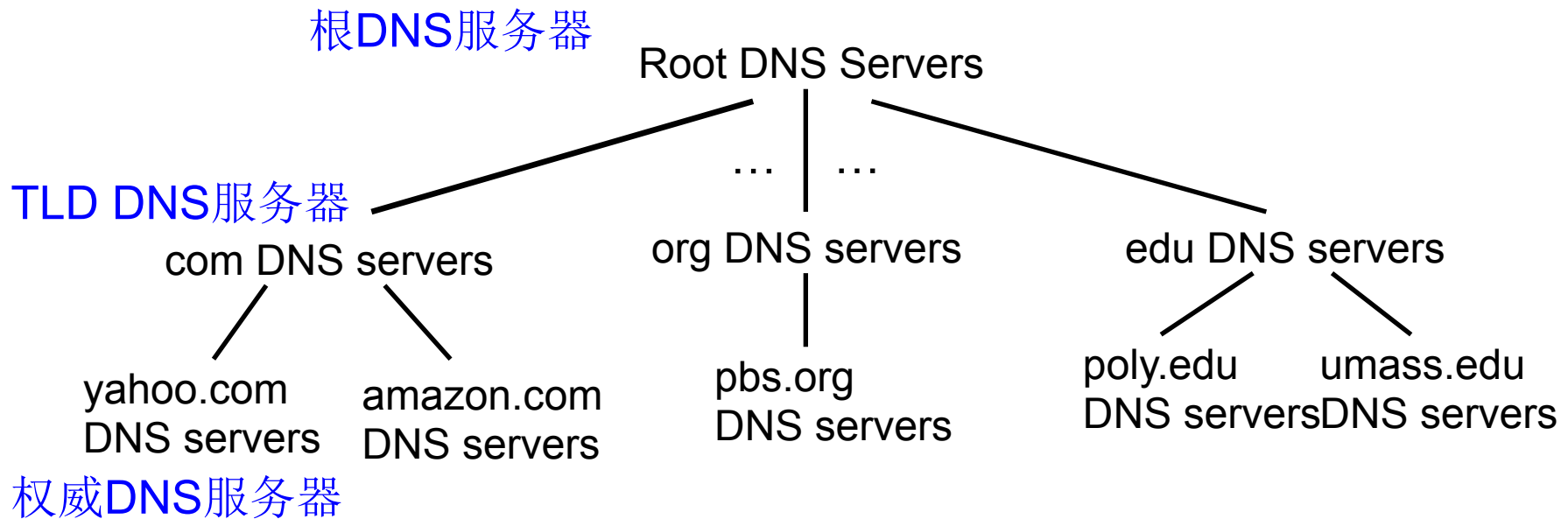
why not centralize DNS?

- single point of failure 单点故障, 集中服务器坏因特网瘫痪
- traffic volume 通信容量, 单个**DNS**服务器不能处理所有查询
- distant centralized database 远距离通信, 严重的时延, 拥塞、低速
- maintenance 维护, 中央数据库将过于庞大, 频繁更新

A: *doesn't scale!*

I. DNS: a distributed, hierarchical database

- 为了解决扩展性问题，DNS使用了大量DNS服务器，以层次方式组织并分布在全球范围内。

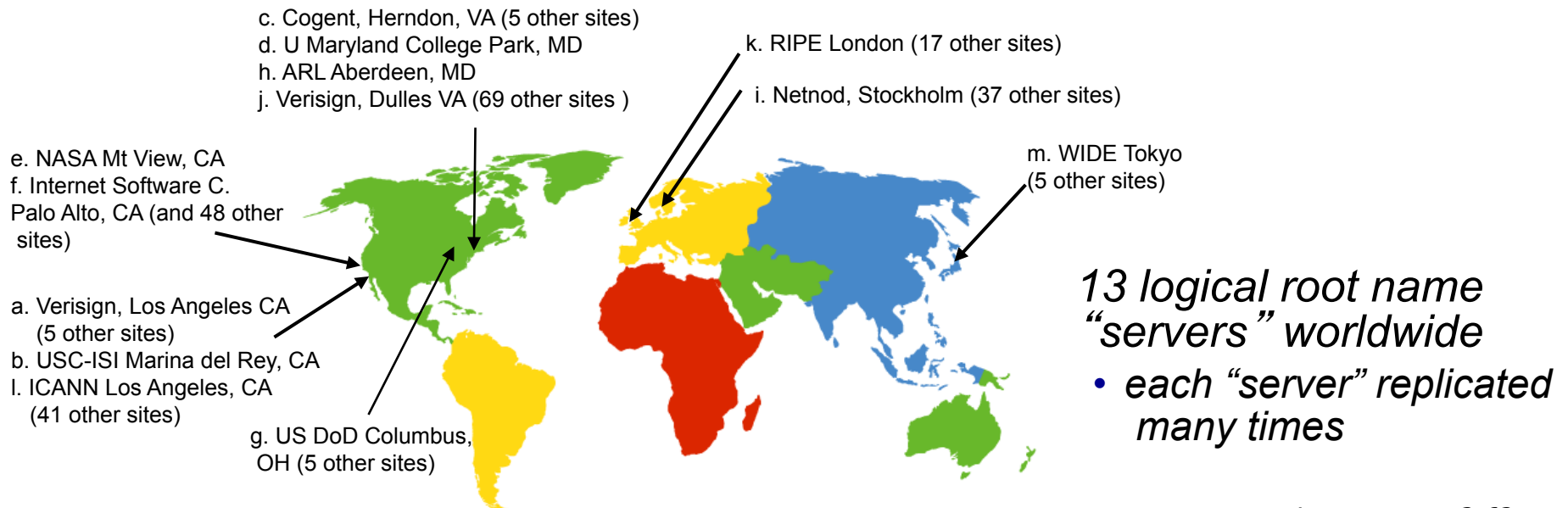


client wants IP for www.amazon.com; 1st approximation:

- client queries root server to find com DNS server
- client queries .com DNS server to get amazon.com DNS server
- client queries amazon.com DNS server to get IP address for
www.amazon.com

(I)根DNS服务器: root name servers

- 400多个根DNS服务器遍及全球，由13个不同组织管理；
- contacted by local name server that can not resolve name;
- root name server:
 - contacts **authoritative name server** if name mapping not known ;
 - gets mapping 提供TLD服务器的IP地址映射；
 - returns mapping to local name server。



(2) 顶级域(TLD)DNS和权威DNS Server

top-level domain (TLD) servers:

- responsible for com, org, net, edu, aero, jobs, museums, and all top-level country domains, e.g.: uk, fr, ca, jp
- Network Solutions maintains servers for .com TLD
- Educause for .edu TLD
- TLD服务器提供了权威DNS服务器的IP地址。

authoritative DNS servers:

- 因特网上具有公共可访问主机的组织机构必须提供公共可访问的DNS记录, 将公共可访问主机名映射为IP;
- organization's own DNS server(s), providing authoritative hostname to IP mappings for organization's named hosts
- can be maintained by organization or service provider (付费, 将记录存储在ISP的一个权威DNS服务器中)

(3) Local DNS server 本地DNS服务器

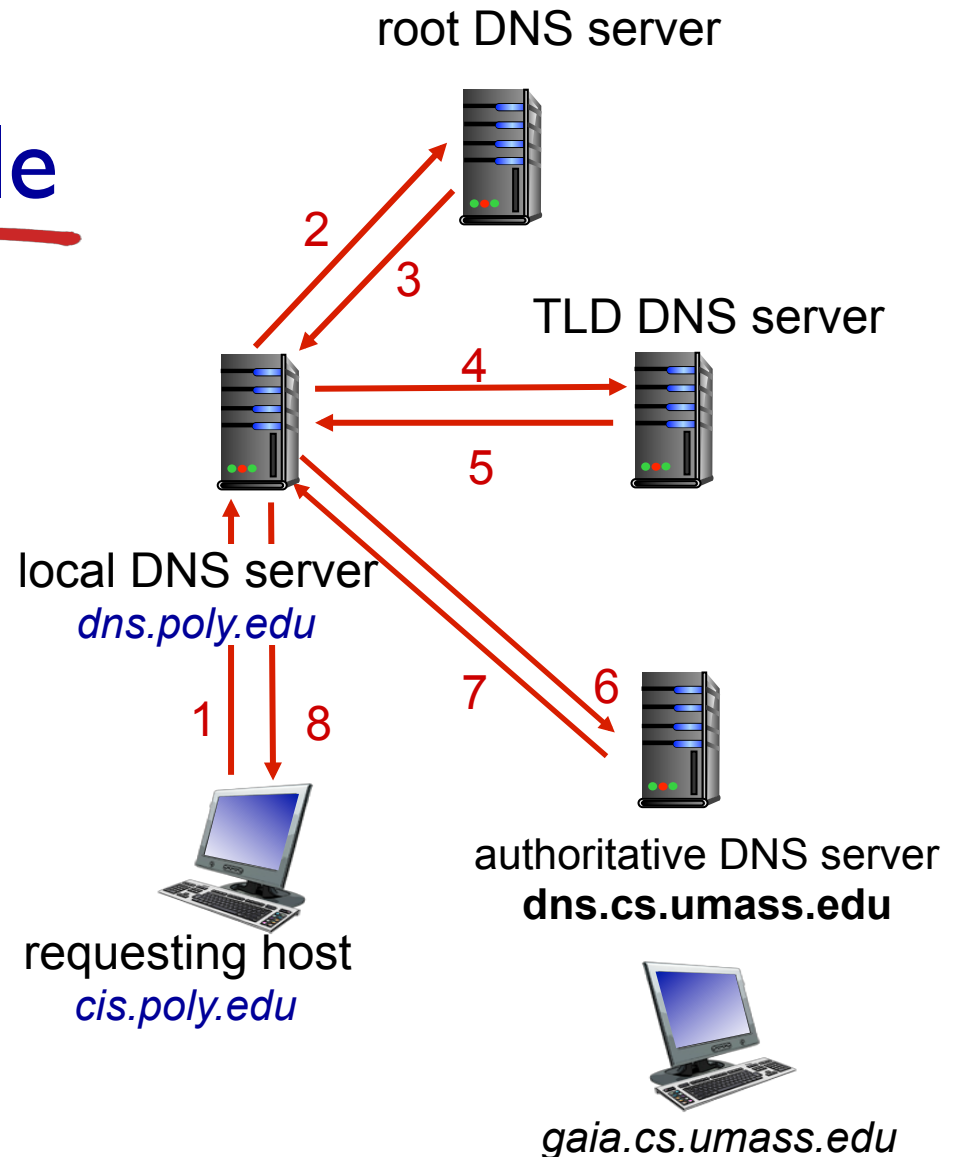
- does not strictly belong to hierarchy
- each ISP (residential ISP, company, university) has one
 - also called “default name server” 默认DNS服务器
- when host makes DNS query, query is sent to its local DNS server
 - 主机入网时, ISP会分配一个IP地址并同时分配一个“默认DNS服务器”, 一般“邻近”主机(在一个LAN或间隔不超过几个router);
 - has local cache of recent name-to-address translation pairs (but may be out of date!)
 - acts as proxy, forwards query into hierarchy

DNS name resolution example

- host at cis.poly.edu wants IP address for gaia.cs.umass.edu

iterated query 迭代查询:

- contacted server replies with name of server to contact
- “I don’t know this name, but ask this server”

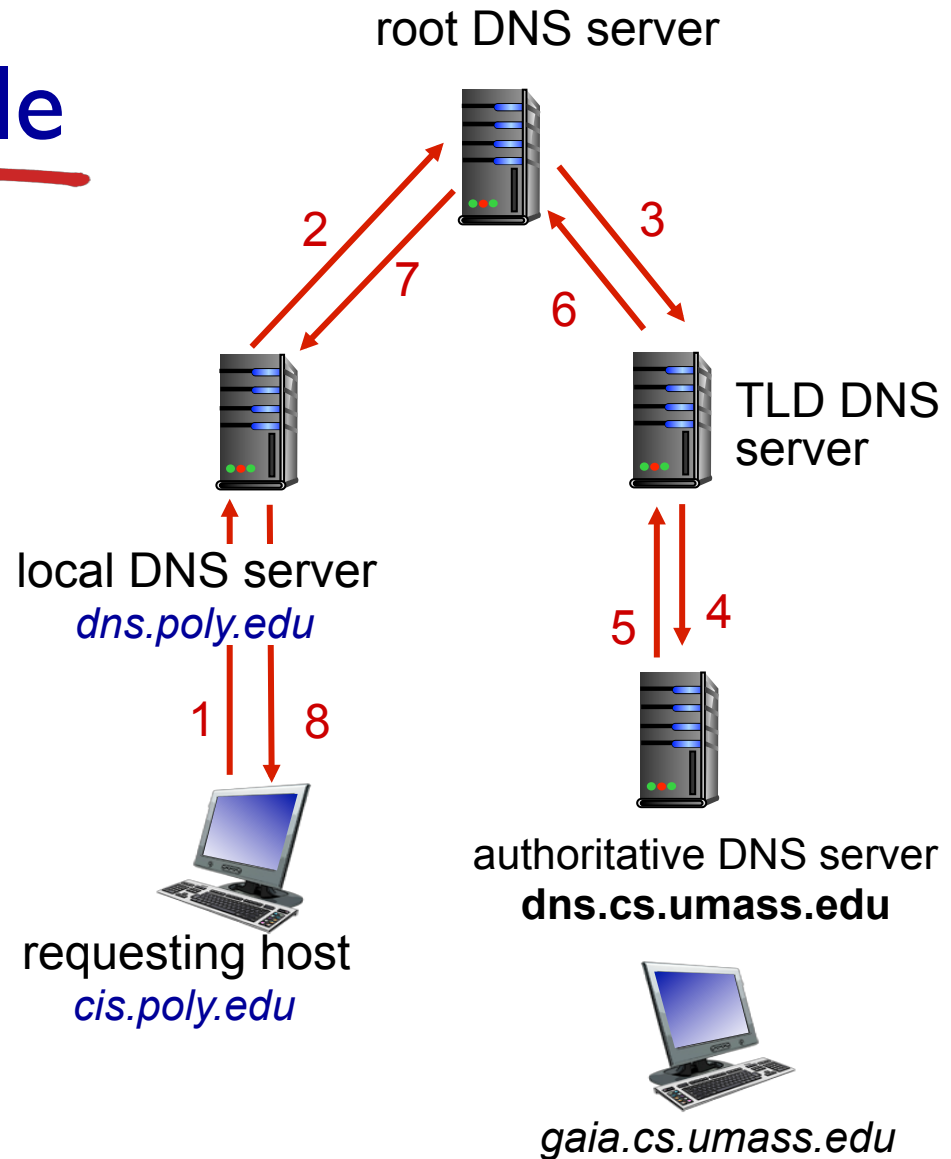


迭代查询又称重指引，当服务器使用迭代查询时能使其他服务器返回一个最佳的查询点提示或主机地址，若此最佳的查询点中包含需要查询的主机地址，则返回主机地址信息，若此时服务器不能够直接查询到主机地址，则是按照提示的指引依次查询，直到服务器给出的提示中包含所需要查询的主机地址。

DNS name resolution example

recursive query 递归查询:

- puts burden of name resolution on contacted name server
- heavy load at upper levels of hierarchy?
- 通常根服务器或者流量较大的域名服务器都不使用递归查询，因为大量的递归查询会导致服务器过载。
- 通常根服务器设置迭代模式，低层DNS服务器设置递归和迭代的混合模式。



2. DNS: caching, updating records

- 改善时延性能及因特网内产生的DNS流量。
- once (any) name server learns mapping, it *caches* mapping
 - cache entries timeout (disappear) after some time (TTL) 2days
 - TLD servers IP address typically cached in local name servers
 - thus root name servers not often visited
- cached entries may be *out-of-date* (best effort name-to-address translation!)
 - if name host changes IP address, may not be known Internet-wide until all TTLs expire
- update/notify mechanisms proposed IETF standard
 - RFC 2136 DNS的动态更新选项

2.4.3 DNS records

- **DNS**: distributed database storing resource records (RR)资源记录。
- RR提供了主机名到IP地址的映射，每个DNS响应报文中可包含一条或多条RR。 RR是个四元组。

RR format: (name, value, type, ttl)

type=A

- **name** is hostname
- **value** is IP address

type=NS

- **name** is domain (e.g., foo.com)
- **value** is hostname of authoritative name server for this domain

type=CNAME

- **name** is alias name for some “canonical” (the real) name
- **www.ibm.com** is really **servereast.backup2.ibm.com**
- **value** is canonical name 规范主机名

type=MX

- **value** is canonical name of mailserver associated with **name**
- 邮件服务器和web服务器可用相同别名

2.4.3 DNS records

RR format: (name, value, type, ttl)

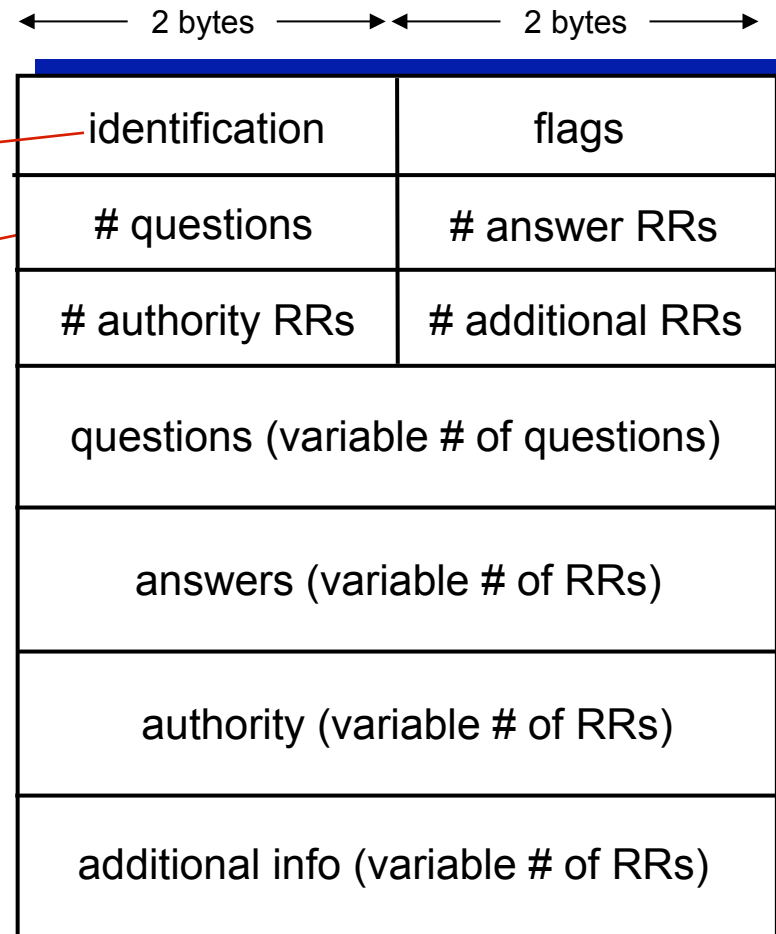
- 如果当前DNS是收到的请求中所请求主机名的权威DNS服务器，或者缓存中有所请求主机名的A记录，则返回一条A记录；
- 如果当前DNS服务器不是所请求主机名的权威DNS服务器，则返回一条NS记录提供所请求主机的域及其权威DNS服务器主机名，同时再返回一条A记录包含该权威DNS服务器的IP地址。

1. DNS protocol, messages

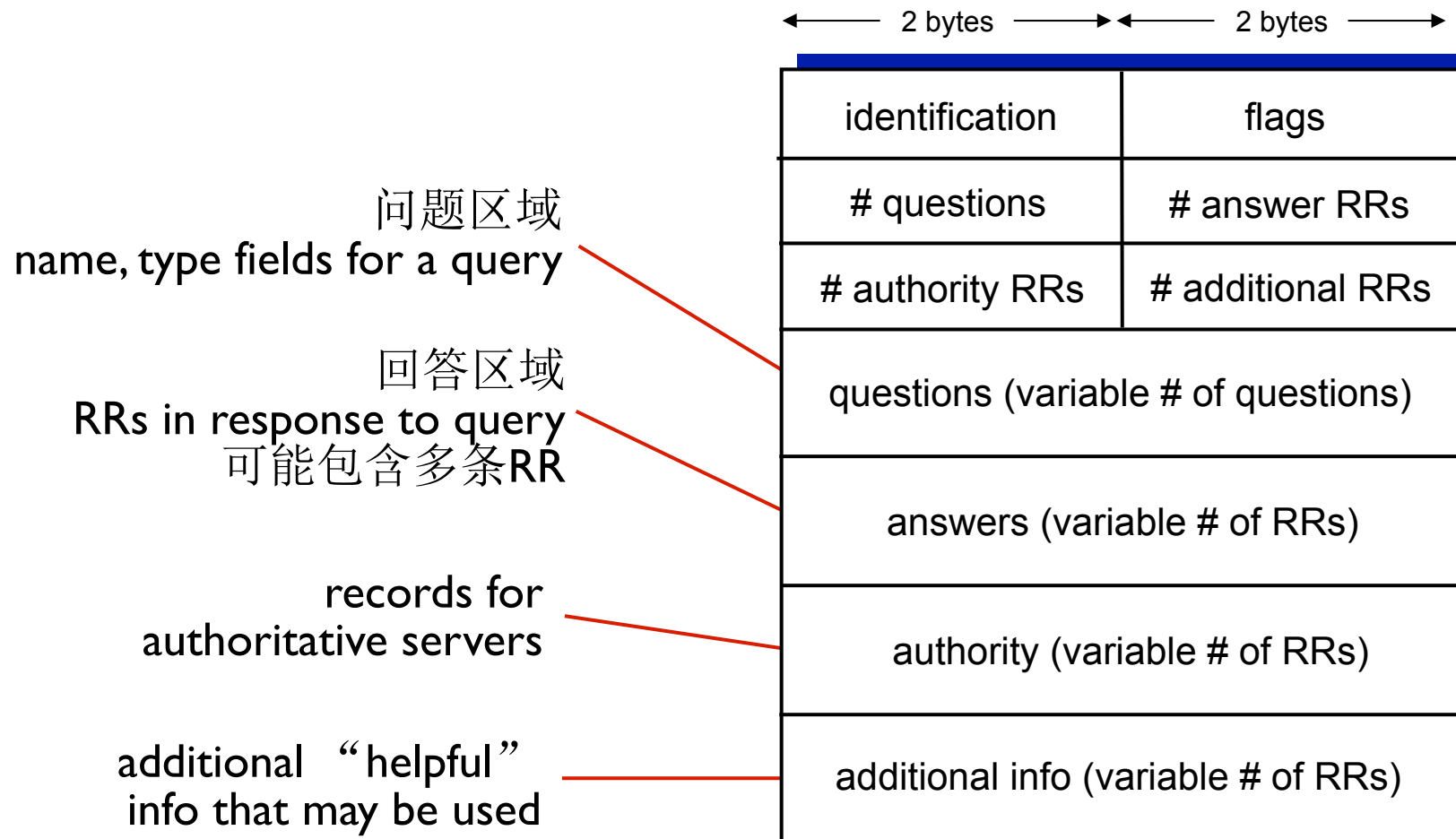
- *query* and *reply* messages, both with same *message format*

message header 前12字节

- **identification:** 16 bit # for query, reply to query uses same #
- **flags:**
 - query 0 or reply 1
 - recursion desired
 - recursion available
 - reply is authoritative



I. DNS protocol, messages



2. Inserting records into DNS

- example: new startup “Network Utopia”
- register name networkutopia.com at **DNS registrar** (e.g., Network Solutions, 1999年前独家)
- **登记注册机构 DNS regisitar**是一个商业实体，验证域名的唯一性，并将域名输入DNS数据库。
 - provide your host names, IP addresses of authoritative name server (primary and secondary)
 - registrar inserts two RRs into .com TLD server:
(networkutopia.com, dns1.networkutopia.com, NS)
(dns1.networkutopia.com, 212.212.212.1, A)
- create authoritative server type A record for www.networkutopia.com; type MX record for mail.networkutopia.com

Attacking DNS

DDoS attacks

- bombard root servers with traffic
 - not successful to date
 - traffic filtering
 - local DNS servers cache IPs of TLD servers, allowing root server bypass
- bombard TLD servers
 - potentially more dangerous

redirect attacks

- man-in-middle
 - Intercept queries
- DNS poisoning
 - Send bogus replies to DNS server, which caches

exploit DNS for DDoS

- send queries with spoofed source address: target IP
- requires amplification

Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks

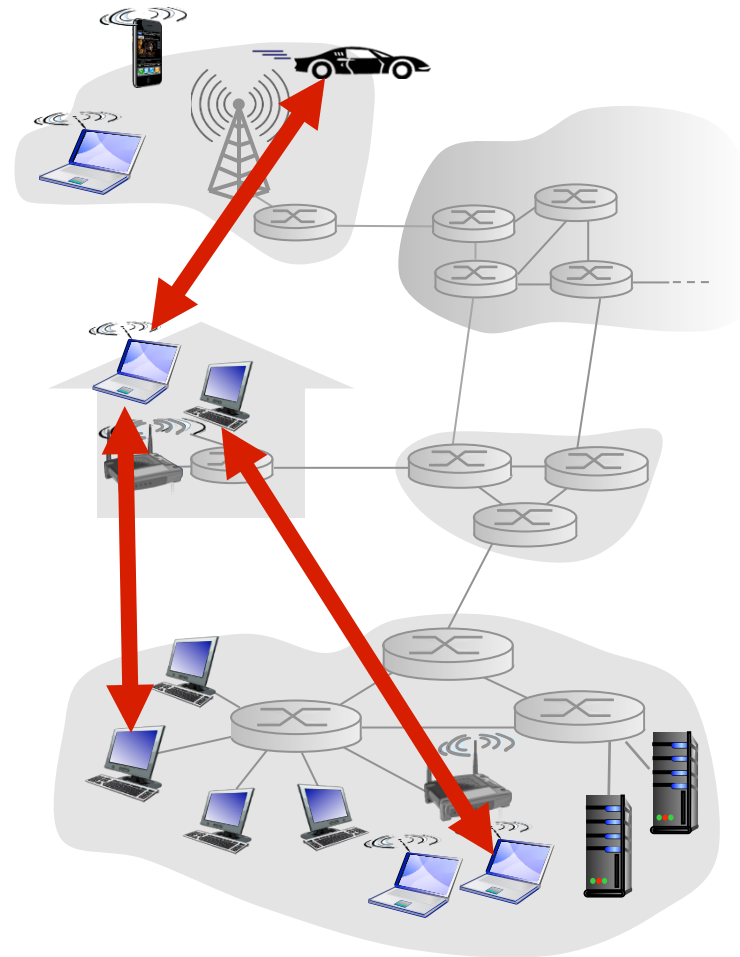
2.7 socket programming with UDP and TCP

Pure P2P architecture

- *no* always-on server
- arbitrary end systems directly communicate
- peers are intermittently connected and change IP addresses

examples:

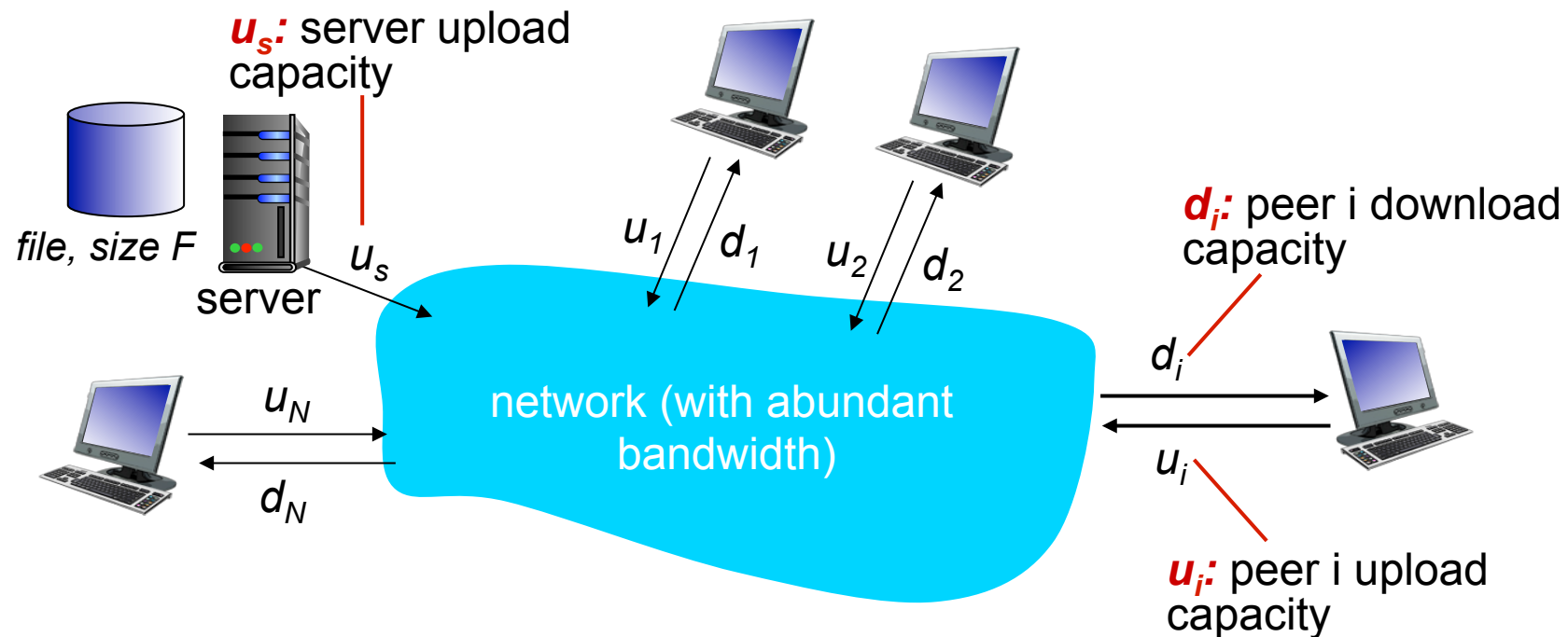
- file distribution (BitTorrent)
- Streaming (KanKan)
- VoIP (Skype)



File distribution: client-server vs P2P

Question: how much time to distribute file (size F) from one server to N peers?

- peer upload/download capacity is limited resource



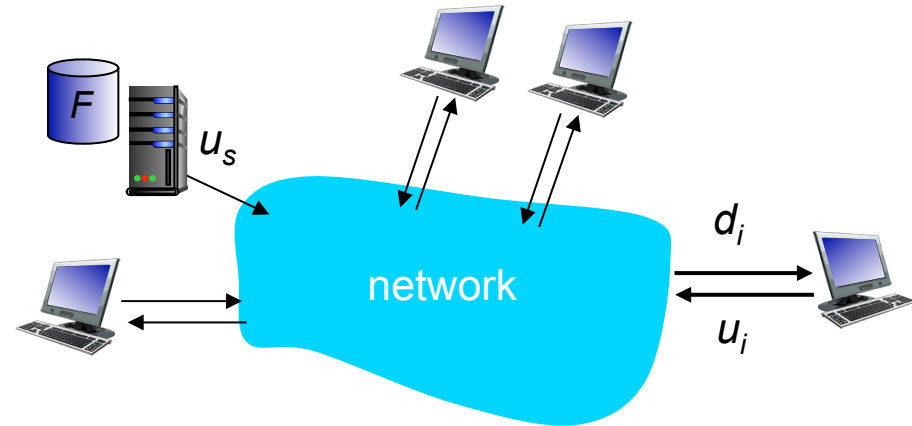
File distribution time: client-server

- **server transmission:** must sequentially send (upload) N file copies:

- time to send one copy: F/u_s
- time to send N copies: NF/u_s

- **client:** each client must download file copy

- d_{min} = min client download rate
- min client download time: F/d_{min}



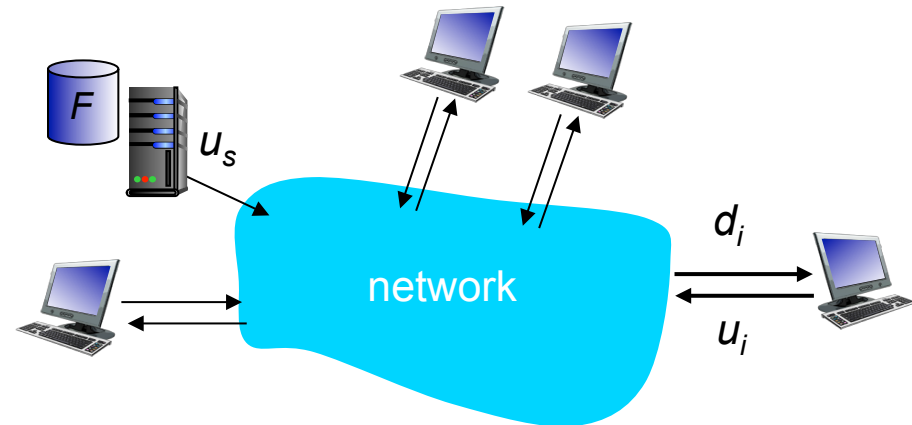
*time to distribute F
to N clients using
client-server approach*

$$D_{c-s} \geq \max\{NF/u_s, F/d_{min}\}$$

increases linearly in N

File distribution time: P2P

- **server transmission:** must upload at least one copy
 - time to send one copy: F/u_s
- **client:** each client must download file copy
 - min client download time: F/d_{\min}
- **clients:** as aggregate must download NF bits
 - max upload rate (limiting max download rate) is $u_s + \sum u_i$



*time to distribute F
to N clients using
P2P approach*

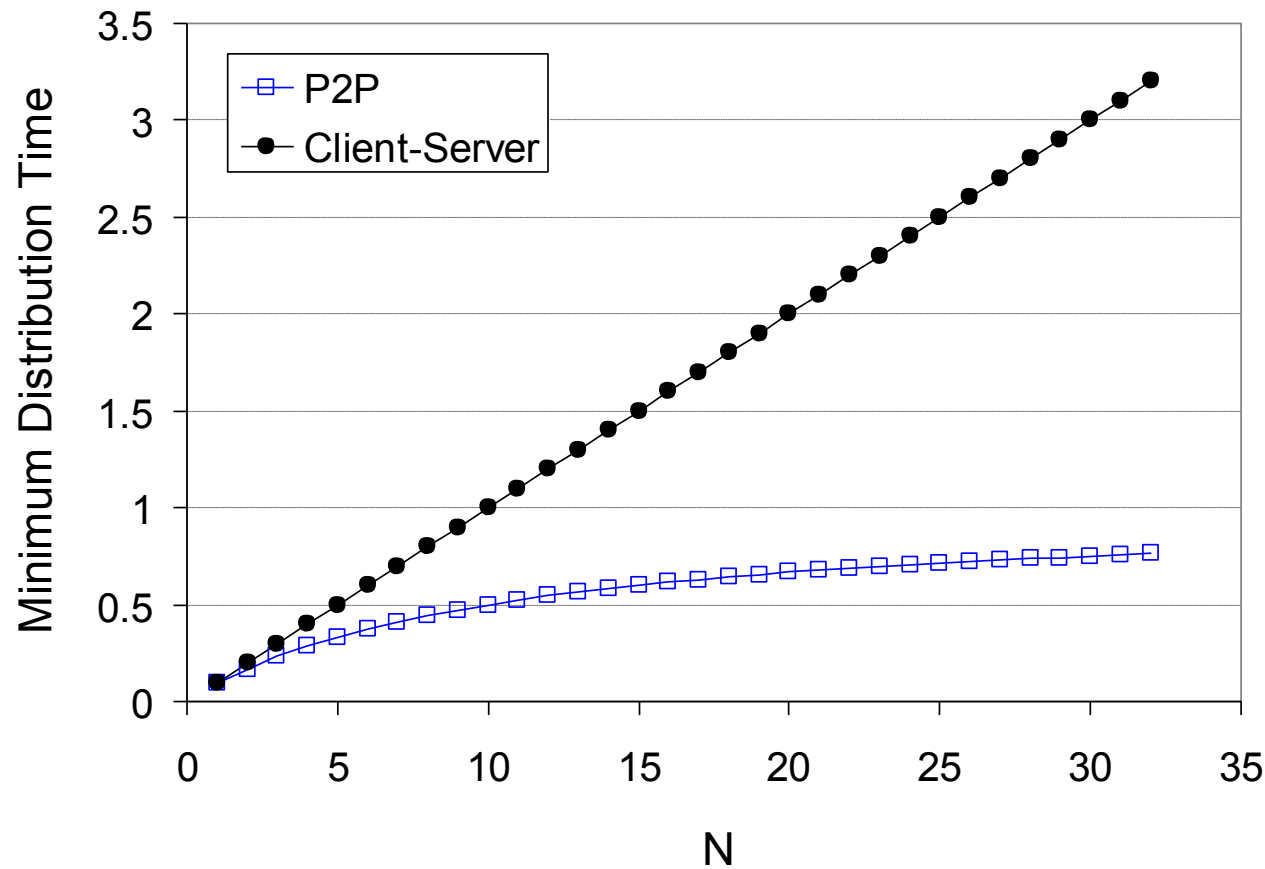
$$D_{P2P} \geq \max\{F/u_s, F/d_{\min}, NF/(u_s + \sum u_i)\}$$

increases linearly in N ...

... but so does this, as each peer brings service capacity

Client-server vs. P2P: example

client upload rate = u , $F/u = 1$ hour, $u_s = 10u$, $d_{min} \geq u_s$

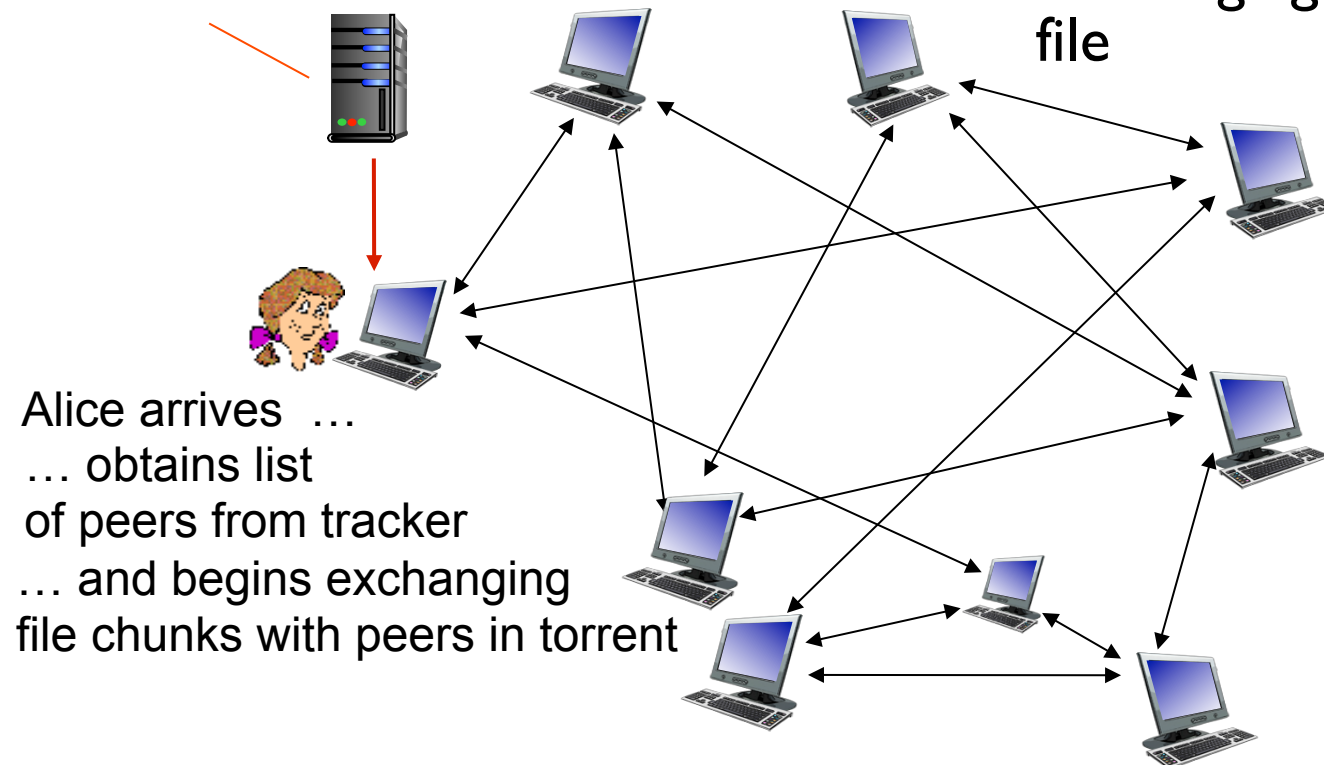


P2P file distribution: BitTorrent

- file divided into 256Kb chunks
- peers in torrent send/receive file chunks

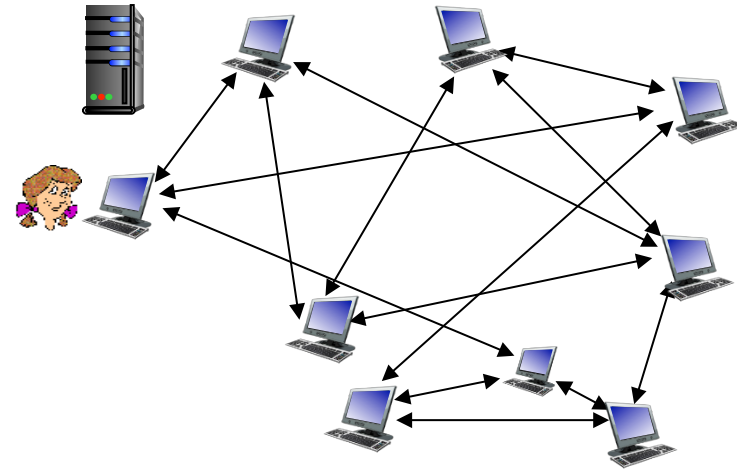
tracker: tracks peers participating in torrent

torrent: group of peers exchanging chunks of a file



P2P file distribution: BitTorrent

- peer joining torrent:
 - has no chunks, but will accumulate them over time from other peers
 - registers with tracker to get list of peers, connects to subset of peers (“neighbors”)
- while downloading, peer uploads chunks to other peers
- peer may change peers with whom it exchanges chunks
- **churn**: peers may come and go
- once peer has entire file, it may (selfishly) leave or (altruistically) remain in torrent



BitTorrent: requesting, sending file chunks

requesting chunks:

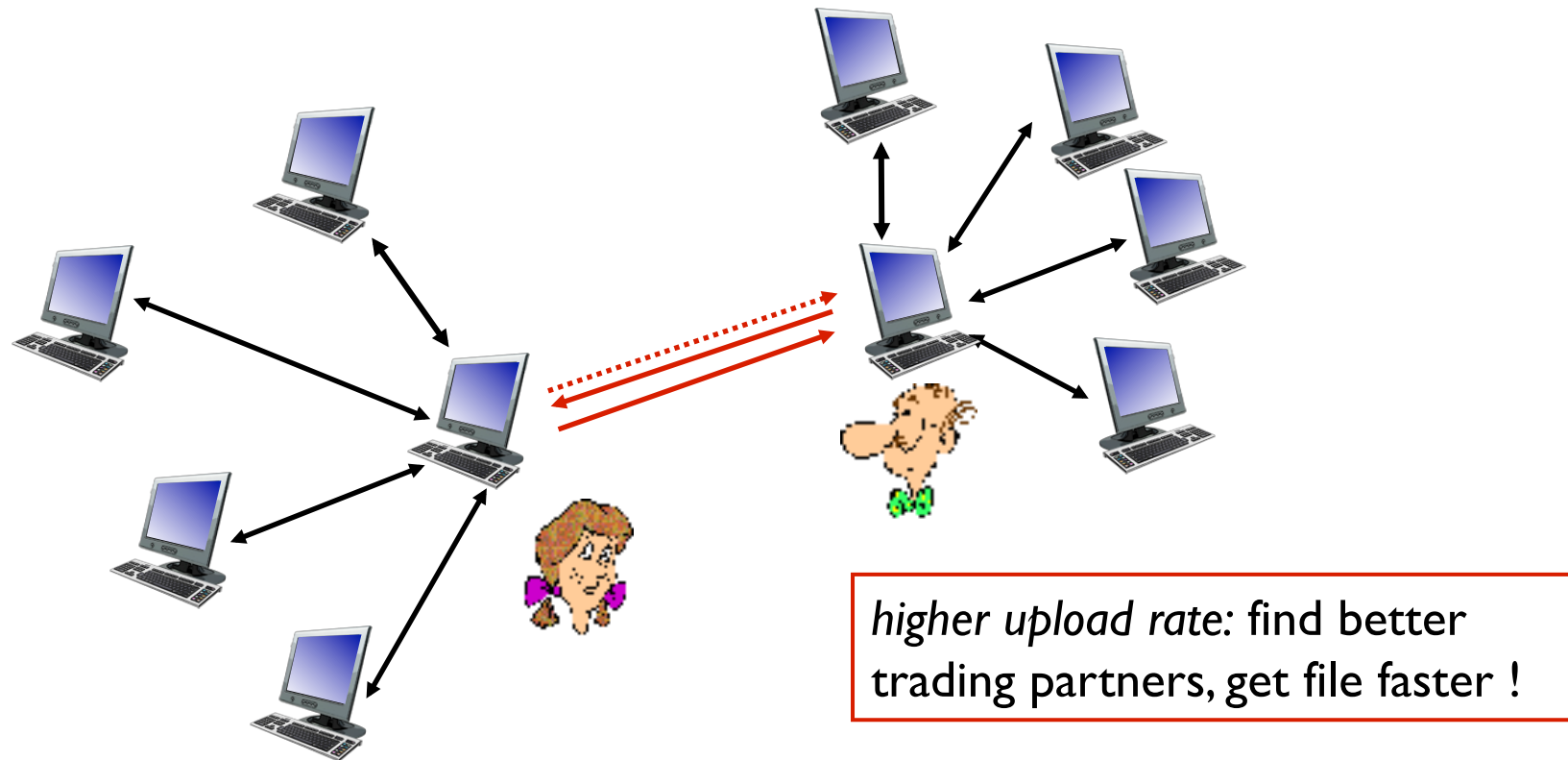
- at any given time, different peers have different subsets of file chunks
- periodically, Alice asks each peer for list of chunks that they have
- Alice requests missing chunks from peers, rarest first

sending chunks: tit-for-tat

- Alice sends chunks to those four peers currently sending her chunks *at highest rate*
 - other peers are choked by Alice (do not receive chunks from her)
 - re-evaluate top 4 every 10 secs
- every 30 secs: randomly select another peer, starts sending chunks
 - “optimistically unchoke” this peer
 - newly chosen peer may join top 4

BitTorrent: tit-for-tat

- (1) Alice “optimistically unchokes” Bob
- (2) Alice becomes one of Bob’s top-four providers; Bob reciprocates
- (3) Bob becomes one of Alice’s top-four providers



Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

2.6 video streaming and content distribution networks (CDNs)

2.7 socket programming with UDP and TCP

Video Streaming and CDNs: context

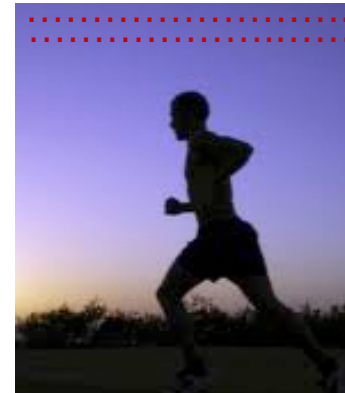
- video traffic: major consumer of Internet bandwidth
 - Netflix, YouTube: 37%, 16% of downstream residential ISP traffic
 - ~1B YouTube users, ~75M Netflix users
- challenge: scale - how to reach ~1B users?
 - single mega-video server won't work (why?)
- challenge: heterogeneity
 - different users have different capabilities (e.g., wired versus mobile; bandwidth rich versus bandwidth poor)
- *solution*: distributed, application-level infrastructure



Multimedia: video

- video: sequence of images displayed at constant rate
 - e.g., 24 images/sec
- digital image: array of pixels
 - each pixel represented by bits
- coding: use redundancy *within* and *between* images to decrease # bits used to encode image
 - spatial (within image)
 - temporal (from one image to next)

spatial coding example: instead of sending N values of same color (all purple), send only two values: color value (*purple*) and number of repeated values (N)



frame i

temporal coding example: instead of sending complete frame at $i+1$, send only differences from frame i

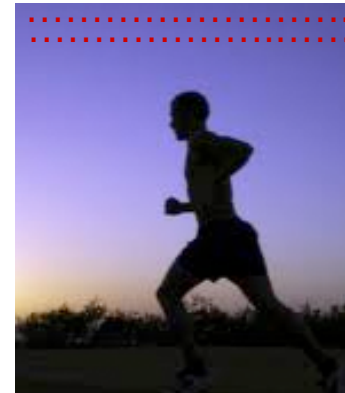


frame $i+1$

Multimedia: video

- **CBR: (constant bit rate):**
video encoding rate fixed
- **VBR: (variable bit rate):**
video encoding rate changes
as amount of spatial,
temporal coding changes
- **examples:**
 - MPEG I (CD-ROM) 1.5 Mbps
 - MPEG2 (DVD) 3-6 Mbps
 - MPEG4 (often used in Internet, < 1 Mbps)

spatial coding example: instead of sending N values of same color (all purple), send only two values: color value (*purple*) and number of repeated values (N)



frame i

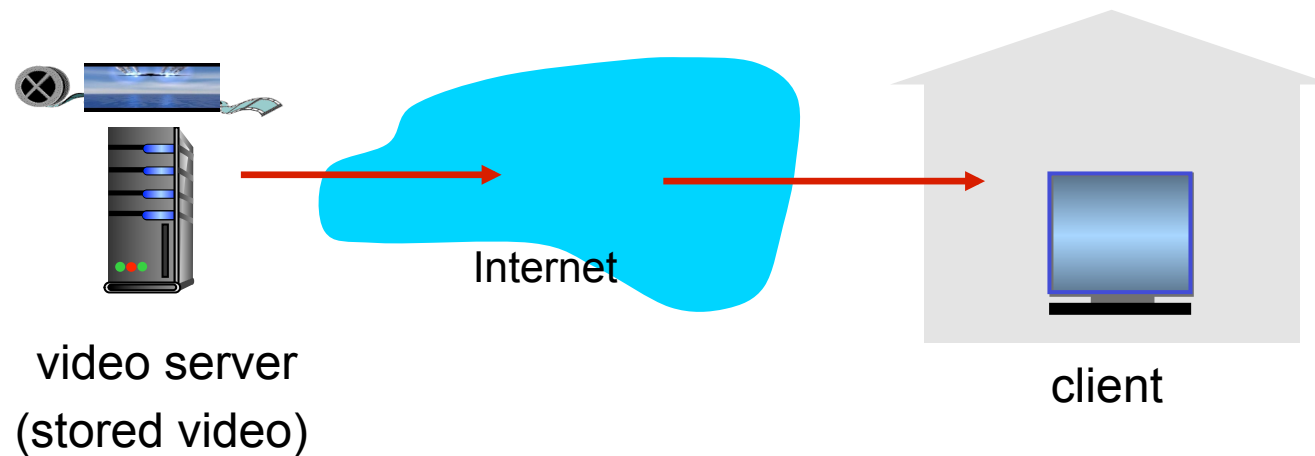
temporal coding example:
instead of sending complete frame at $i+1$, send only differences from frame i



frame $i+1$

Streaming stored video:

simple scenario:



Streaming multimedia: DASH

- **DASH:** *D*ynamic, *A*daptive *S*treaming over *H*TTP
- **server:**
 - divides video file into multiple chunks
 - each chunk stored, encoded at different rates
 - *manifest file:* provides URLs for different chunks
- **client:**
 - periodically measures server-to-client bandwidth
 - consulting manifest, requests one chunk at a time
 - chooses maximum coding rate sustainable given current bandwidth
 - can choose different coding rates at different points in time (depending on available bandwidth at time)

Streaming multimedia: DASH

- *DASH: Dynamic, Adaptive Streaming over HTTP*
- “intelligence” at client: client determines
 - *when* to request chunk (so that buffer starvation, or overflow does not occur)
 - *what encoding rate* to request (higher quality when more bandwidth available)
 - *where* to request chunk (can request from URL server that is “close” to client or has high available bandwidth)

Content distribution networks

- *challenge*: how to stream content (selected from millions of videos) to hundreds of thousands of *simultaneous* users?
- *option 1*: single, large “mega-server”
 - single point of failure
 - point of network congestion
 - long path to distant clients
 - multiple copies of video sent over outgoing link

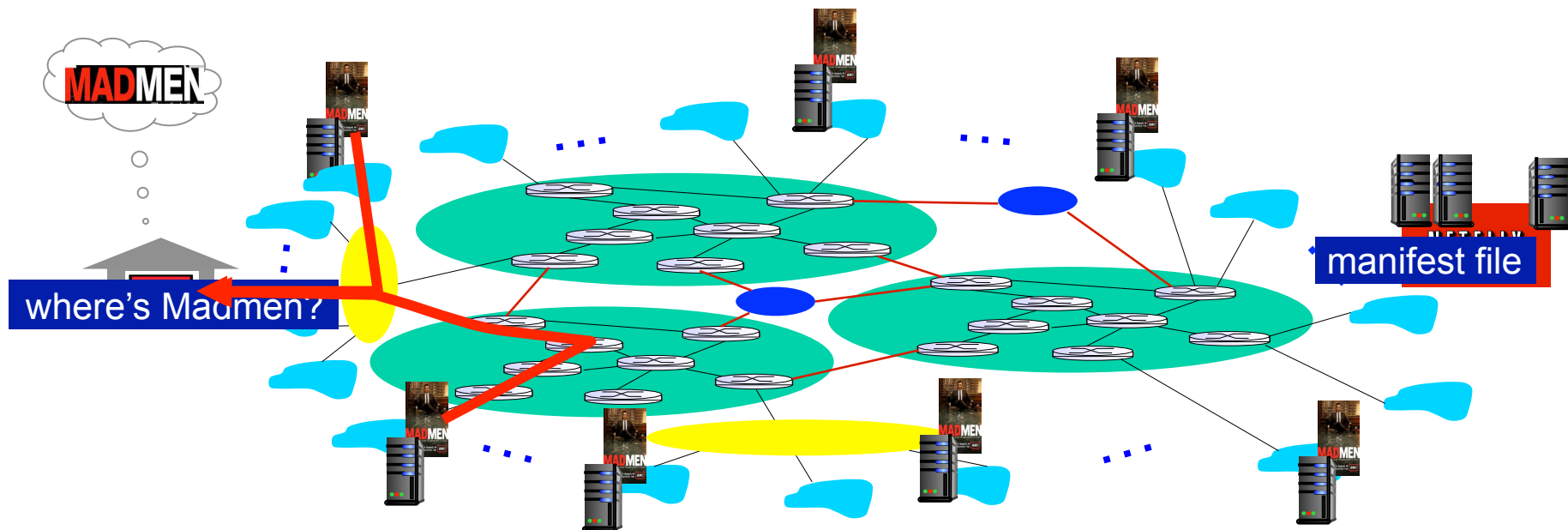
....quite simply: this solution *doesn't scale*

Content distribution networks

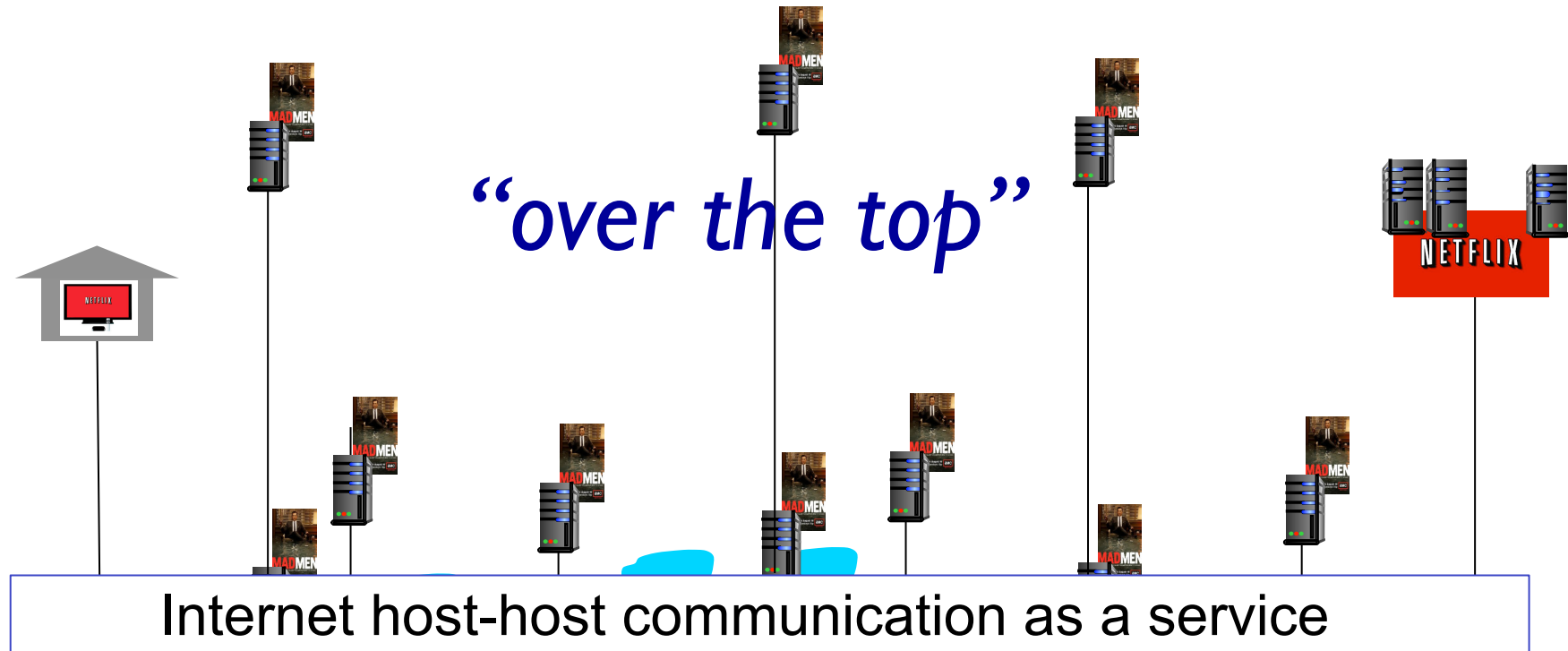
- *challenge*: how to stream content (selected from millions of videos) to hundreds of thousands of simultaneous users?
- *option 2*: store/serve multiple copies of videos at multiple geographically distributed sites (*CDN*)
 - *enter deep*: push CDN servers deep into many access networks
 - close to users
 - used by Akamai, 1700 locations
 - *bring home*: smaller number (10's) of larger clusters in POPs near (but not within) access networks
 - used by Limelight

Content Distribution Networks (CDNs)

- CDN: stores copies of content at CDN nodes
 - e.g. Netflix stores copies of MadMen
- subscriber requests content from CDN
 - directed to nearby copy, retrieves content
 - may choose different copy if network path congested



Content Distribution Networks (CDNs)



OTT challenges: coping with a congested Internet

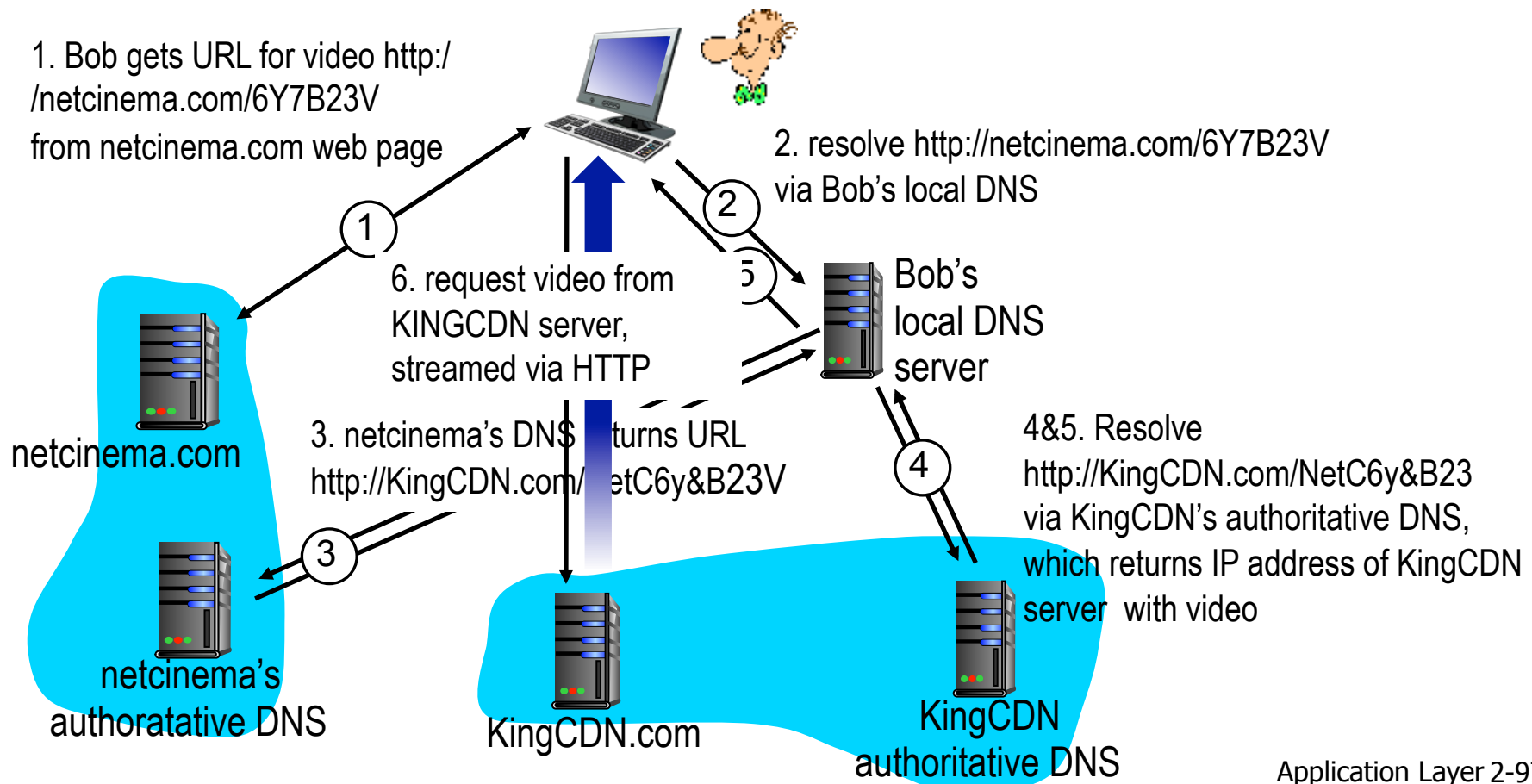
- from which CDN node to retrieve content?
- viewer behavior in presence of congestion?
- what content to place in which CDN node?

more .. in chapter 7

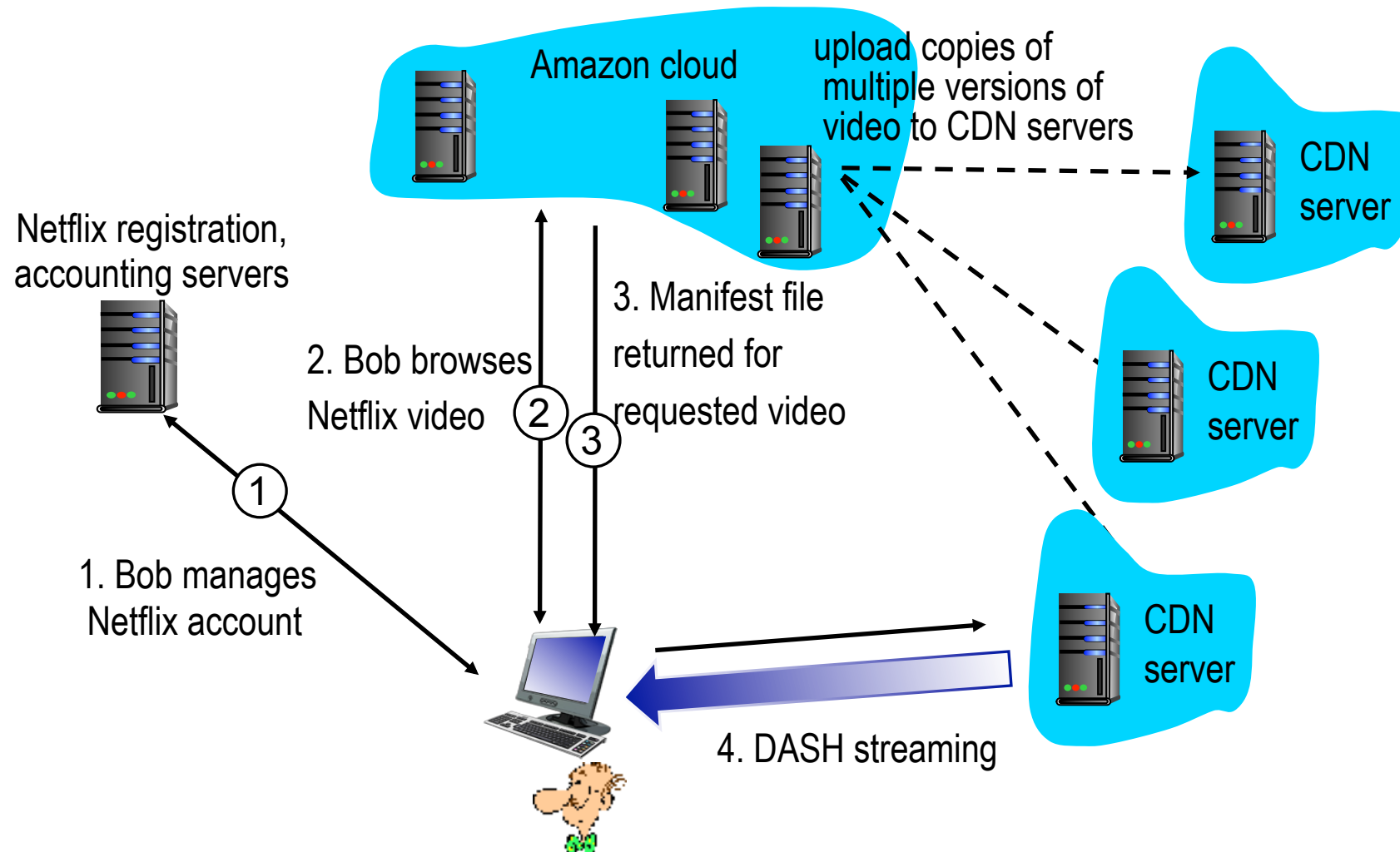
CDN content access: a closer look

Bob (client) requests video `http://netcinema.com/6Y7B23V`

- video stored in CDN at `http://KingCDN.com/NetC6y&B23V`



Case study: Netflix



Chapter 2: outline

2.1 principles of network applications

2.2 Web and HTTP

2.3 electronic mail

- SMTP, POP3, IMAP

2.4 DNS

2.5 P2P applications

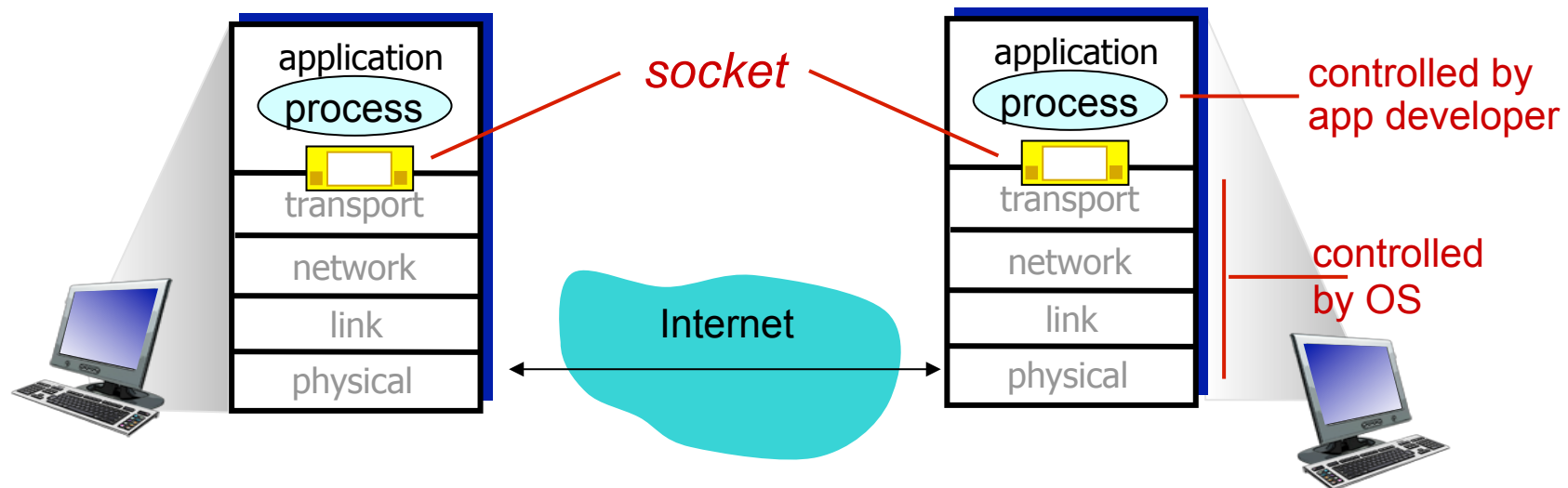
2.6 video streaming and content distribution networks

2.7 socket programming with UDP and TCP

Socket programming

goal: learn how to build client/server applications that communicate using sockets

socket: door between application process and end-end-transport protocol



Socket programming

Two socket types for two transport services:

- **UDP:** unreliable datagram
- **TCP:** reliable, byte stream-oriented

Application Example:

1. client reads a line of characters (data) from its keyboard and sends data to server
2. server receives the data and converts characters to uppercase
3. server sends modified data to client
4. client receives modified data and displays line on its screen

2.7.1 Socket programming *with UDP*

UDP: no “connection” between client & server

- no handshaking before sending data
- sender explicitly attaches IP destination address and port # to each packet
- receiver extracts sender IP address and port# from received packet

UDP: transmitted data may be lost or received out-of-order

Application viewpoint:

- UDP provides *unreliable* transfer of groups of bytes (“datagrams”) between client and server

Client/server socket interaction: UDP

server (running on serverIP)

create socket, port= x:
`serverSocket =
socket(AF_INET,SOCK_DGRAM)`

↓
read datagram from
`serverSocket`

↓
write reply to
`serverSocket`
specifying
client address,
port number

client

create socket:
`clientSocket =
socket(AF_INET,SOCK_DGRAM)`

↓
Create datagram with server IP and
port=x; send datagram via
`clientSocket`

↓
read datagram from
`clientSocket`

↓
close
`clientSocket`

Example app: UDP client

Python UDPClient

include Python's socket
library

→ from socket import *

serverName = 'hostname'

serverPort = 12000

create UDP socket for
server

→ clientSocket = socket(AF_INET,
SOCK_DGRAM)

get user keyboard
input

→ message = raw_input('Input lowercase sentence:')

Attach server name, port to
message; send into socket

→ clientSocket.sendto(message.encode(),
(serverName, serverPort))

read reply characters from
socket into string

→ modifiedMessage, serverAddress =
clientSocket.recvfrom(2048)

print out received string
and close socket

→ print modifiedMessage.decode()
clientSocket.close()

Example app: UDP server

Python UDPServer

```
from socket import *
```

```
serverPort = 12000
```

create UDP socket → `serverSocket = socket(AF_INET, SOCK_DGRAM)`

bind socket to local port
number 12000 → `serverSocket.bind(("", serverPort))`

```
print ( " The server is ready to receive" )
```

loop forever → `while True:`

Read from UDP socket into
message, getting client's
address (client IP and
port) → `message, clientAddress = serverSocket.recvfrom(2048)`

```
    modifiedMessage = message.decode().upper()
```

send upper case string
back to this client → `serverSocket.sendto(modifiedMessage.encode(),
clientAddress)`

2.7.2 Socket programming *with TCP*

client must contact server

- server process must first be running
- server must have created socket (door) that welcomes client's contact

client contacts server by:

- Creating TCP socket, specifying IP address, port number of server process
- **when client creates socket:** client TCP establishes connection to server TCP

- when contacted by client, *server TCP creates new socket* for server process to communicate with that particular client
 - allows server to talk with multiple clients
 - source port numbers used to distinguish clients (more in Chap 3)

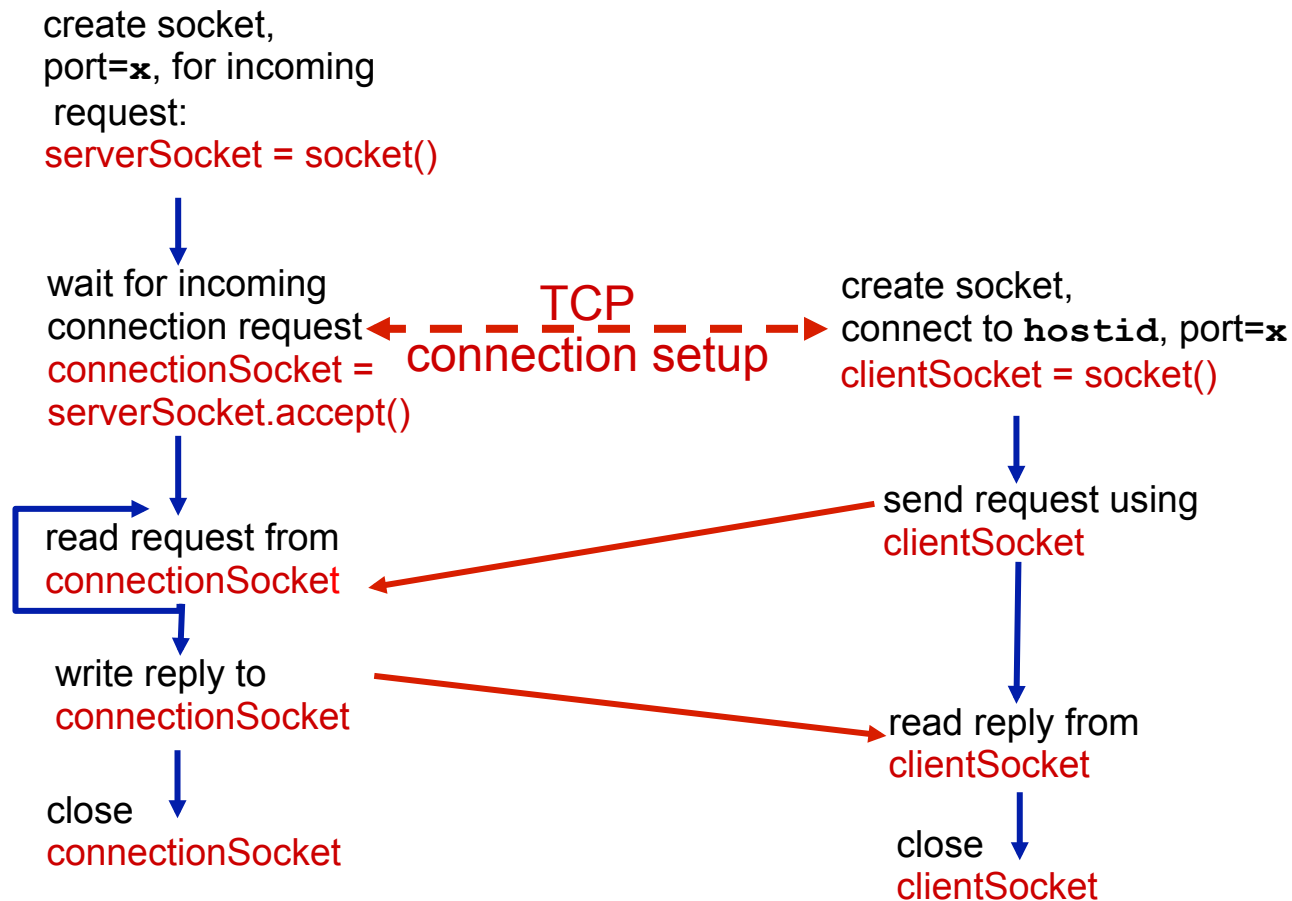
application viewpoint:

TCP provides reliable, in-order byte-stream transfer (“pipe”) between client and server

Client/server socket interaction: TCP

server (running on `hostid`)

client



Example app: TCP client

Python TCPClient

create TCP socket for
server, remote port 12000

```
from socket import *
```

```
serverName = 'servername'
```

```
serverPort = 12000
```

```
clientSocket = socket(AF_INET, SOCK_STREAM)
```

```
clientSocket.connect((serverName, serverPort))
```

```
sentence = raw_input('Input lowercase sentence:')
```

No need to attach server
name, port

```
clientSocket.send(sentence.encode())
```

```
modifiedSentence = clientSocket.recv(1024)
```

```
print ('From Server:', modifiedSentence.decode())
```

```
clientSocket.close()
```

Example app: TCP server

Python TCPServer

create TCP welcoming socket	→	from socket import * serverPort = 12000 serverSocket = socket(AF_INET, SOCK_STREAM) serverSocket.bind(('', serverPort))
server begins listening for incoming TCP requests	→	serverSocket.listen(1) print 'The server is ready to receive'
loop forever	→	while True:
server waits on accept() for incoming requests, new socket created on return	→	connectionSocket, addr = serverSocket.accept()
read bytes from socket (but not address as in UDP)	→	sentence = connectionSocket.recv(1024).decode() capitalizedSentence = sentence.upper()
close connection to this client (but <i>not</i> welcoming socket)	→	connectionSocket.send(capitalizedSentence. encode()) connectionSocket.close()

Chapter 2: summary

our study of network apps now complete!

- application architectures
 - client-server
 - P2P
- application service requirements:
 - reliability, bandwidth, delay
- Internet transport service model
 - connection-oriented, reliable: TCP
 - unreliable, datagrams: UDP
- specific protocols:
 - HTTP
 - SMTP, POP, IMAP
 - DNS
 - P2P: BitTorrent
- video streaming, CDNs
- socket programming:
TCP, UDP sockets

Chapter 2: summary

most importantly: learned about protocols!

- typical request/reply message exchange:
 - client requests info or service
 - server responds with data, status code
- message formats:
 - *headers*: fields giving info about data
 - *data*: info(payload) being communicated

important themes:

- control vs. messages
 - in-band, out-of-band
- centralized vs. decentralized
- stateless vs. stateful
- reliable vs. unreliable message transfer
- “complexity at network edge”