# Application layer: overview

- Principles of network applications
- Web and HTTP
- **E-mail, SMTP, IMAP**
- The Domain Name System DNS

- P2P applications
- video streaming and content distribution networks
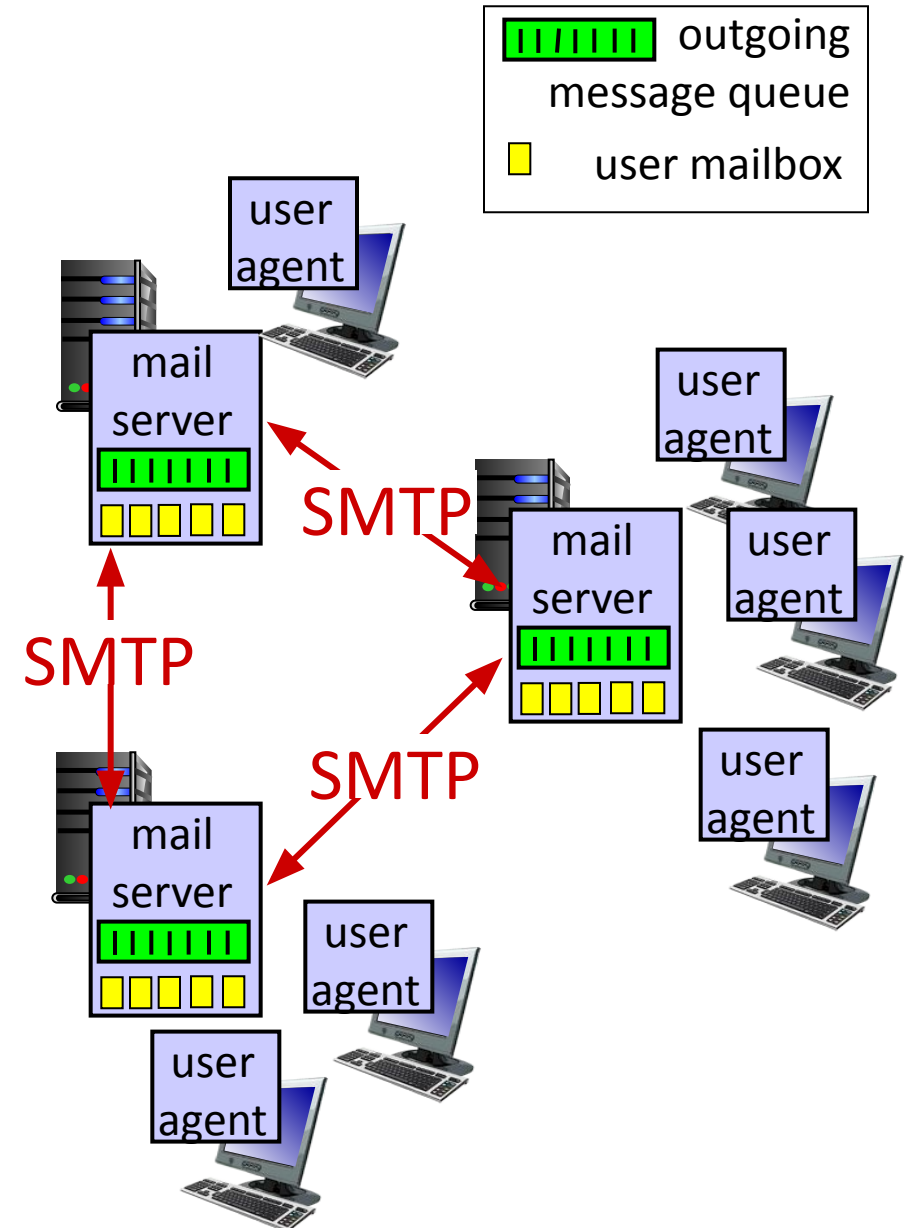- socket programming with UDP and TCP

# E-mail

## Three major components:

- user agents
- mail servers
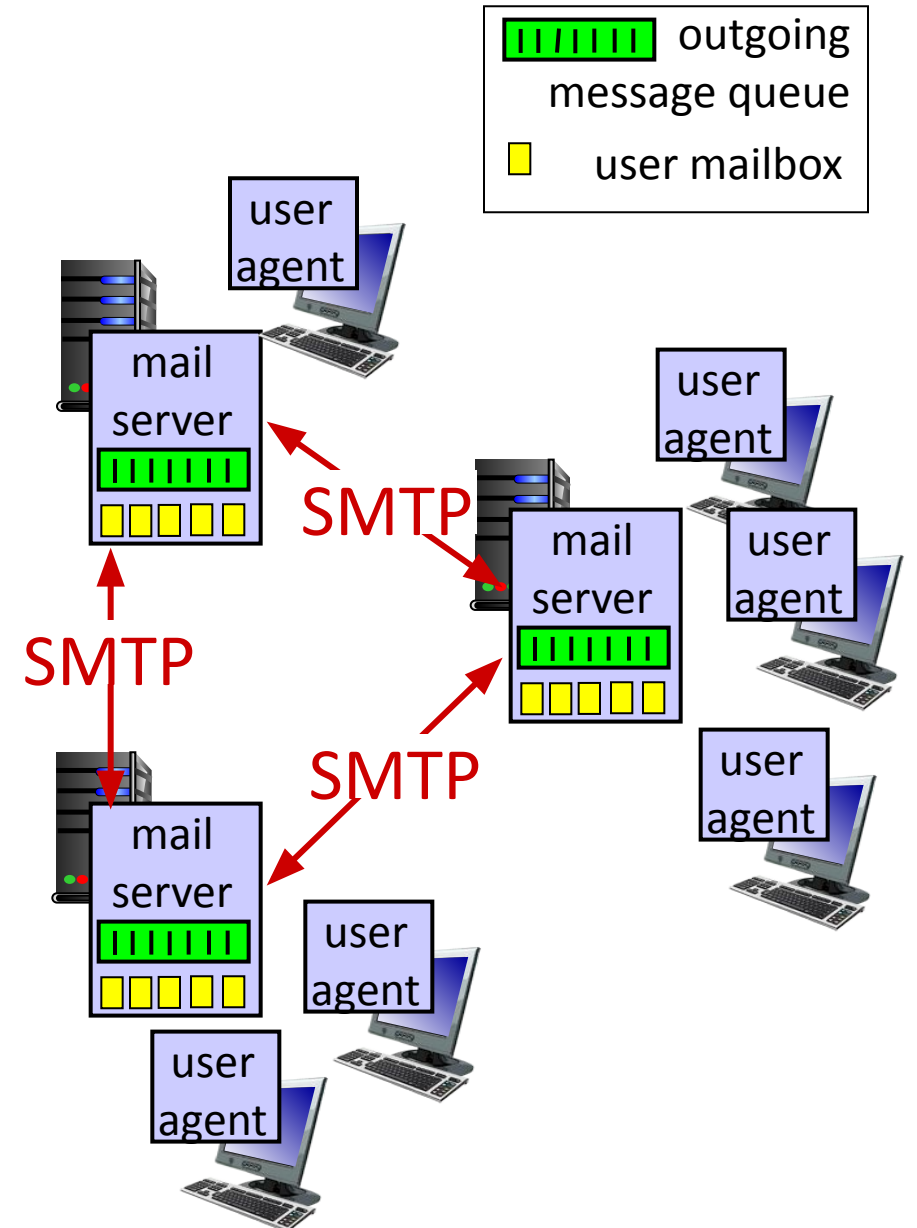- simple mail transfer protocol: SMTP

## User Agent

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

# E-mail: mail servers

mail servers:

- *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
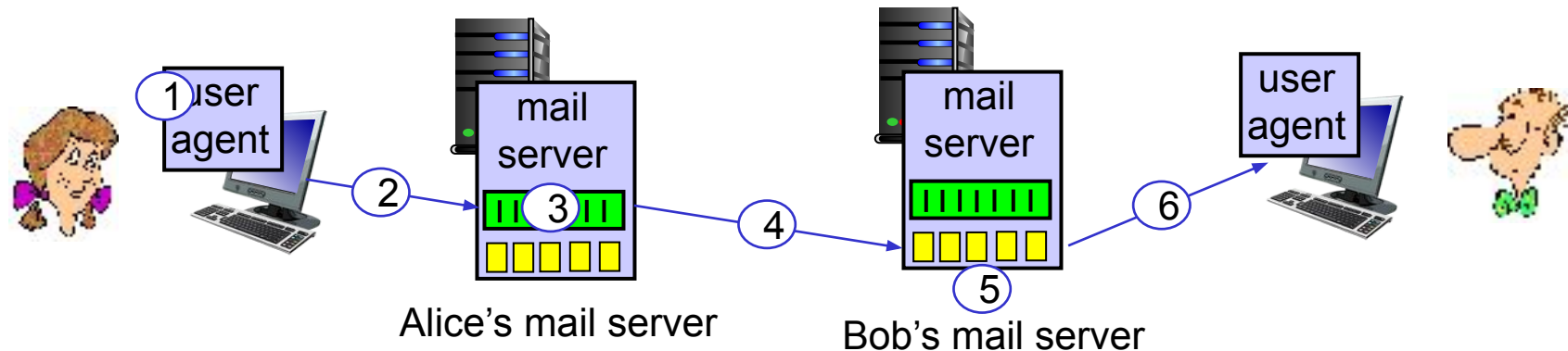  - "server": receiving mail server

# E-mail: the RFC (5321)

- uses TCP to reliably transfer email message from client (mail server initiating connection) to server, port 25
- direct transfer: sending server (acting like client) 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 ASCI

# Scenario: Alice sends e-mail to Bob

1) Alice uses UA to compose e-mail 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



Alice's mail server

Bob's mail server

# Sample SMTP interaction

```
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 e-mail client (reader)

*Note: this will only work if <servername> allows telnet connections to port 25 (this is becoming increasingly rare because of security concerns)*

# SMTP: closing observations
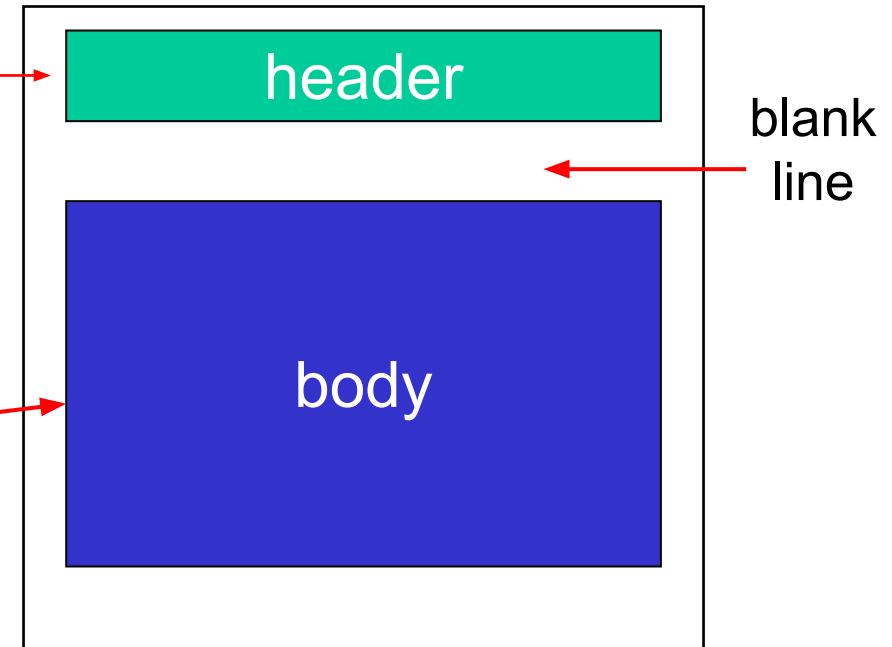
*comparison with HTTP:*

- HTTP: pull

- SMTP: push

- both have ASCII command/response interaction, status codes

- HTTP: each object encapsulated in its own response message

- SMTP: multiple objects sent in multipart message

- 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
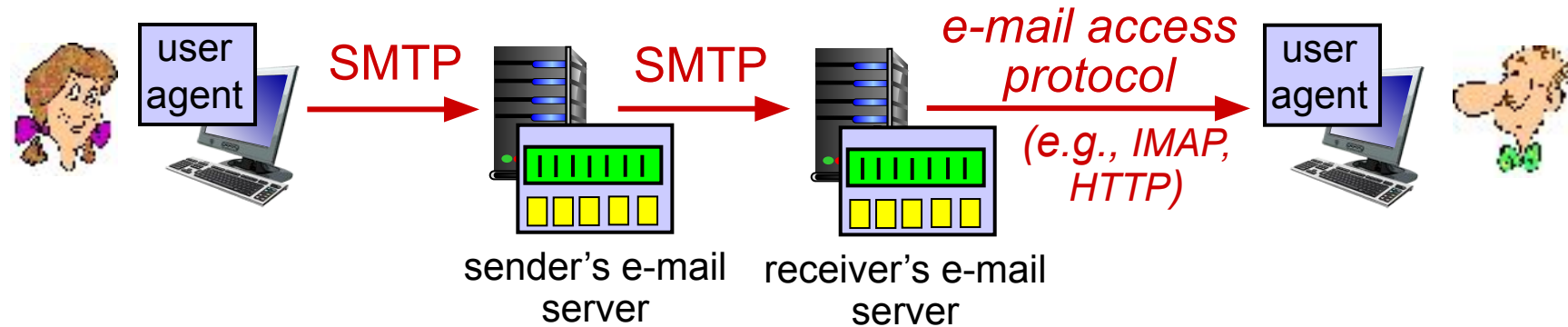
# Mail message format

SMTP: protocol for exchanging e-mail messages, defined in RFC 531 (like HTTP)

RFC 822 defines *syntax* for e-mail message itself (like HTML)

- header lines, e.g.,
  - To:
  - From:
  - Subject:

  these lines, within the body of the email message area different from SMTP MAIL FROM:, RCPT TO: commands!

- Body: the "message" , ASCII characters only



header

blank line

body

# Mail access protocols



- **SMTP:** delivery/storage of e-mail messages to receiver's server

- **mail access protocol:** retrieval from server
  - **IMAP:** Internet Mail Access Protocol [RFC 3501]: messages stored on server, IMAP provides retrieval, deletion, folders of stored messages on server

- **HTTP:** gmail, Hotmail, Yahoo!Mail, etc. provides web-based interface on top of STMP (to send), IMAP (or POP) to retrieve e-mail messages

# Application Layer: Overview

- Principles of network applications
- Web and HTTP
- E-mail, SMTP, IMAP
- **The Domain Name System DNS**

- P2P applications
- video streaming and content distribution networks
- 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
- "name", e.g., cs.umass.edu - used by humans

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

*Domain Name System:*

- *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"

# DNS: services, structure

## DNS services

- hostname to IP address translation

- host aliasing
  - canonical, alias names

- mail server aliasing

- load distribution
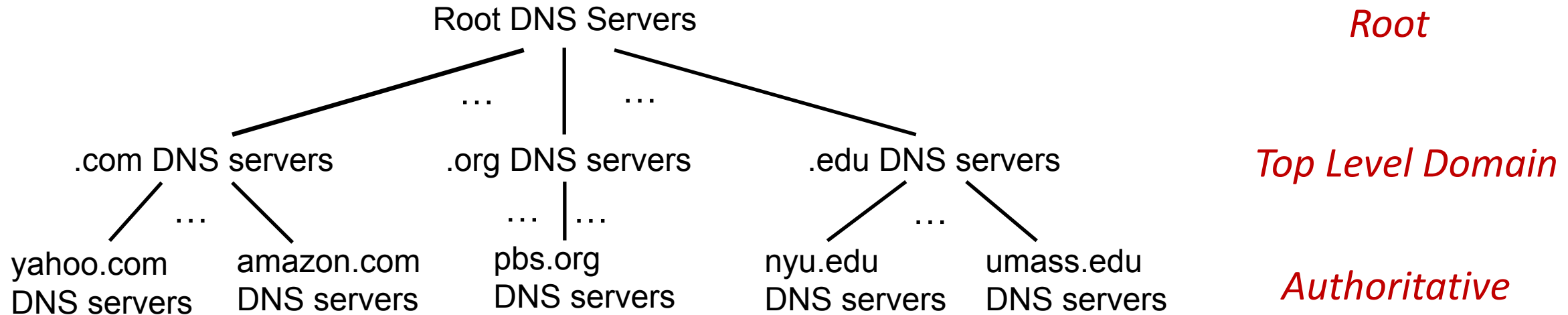  - replicated Web servers: many IP addresses correspond to one name

*Q: Why not centralize DNS?*

- single point of failure
- traffic volume
- distant centralized database
- maintenance

*A: doesn't scale!*

- Comcast DNS servers alone: 600B DNS queries per day
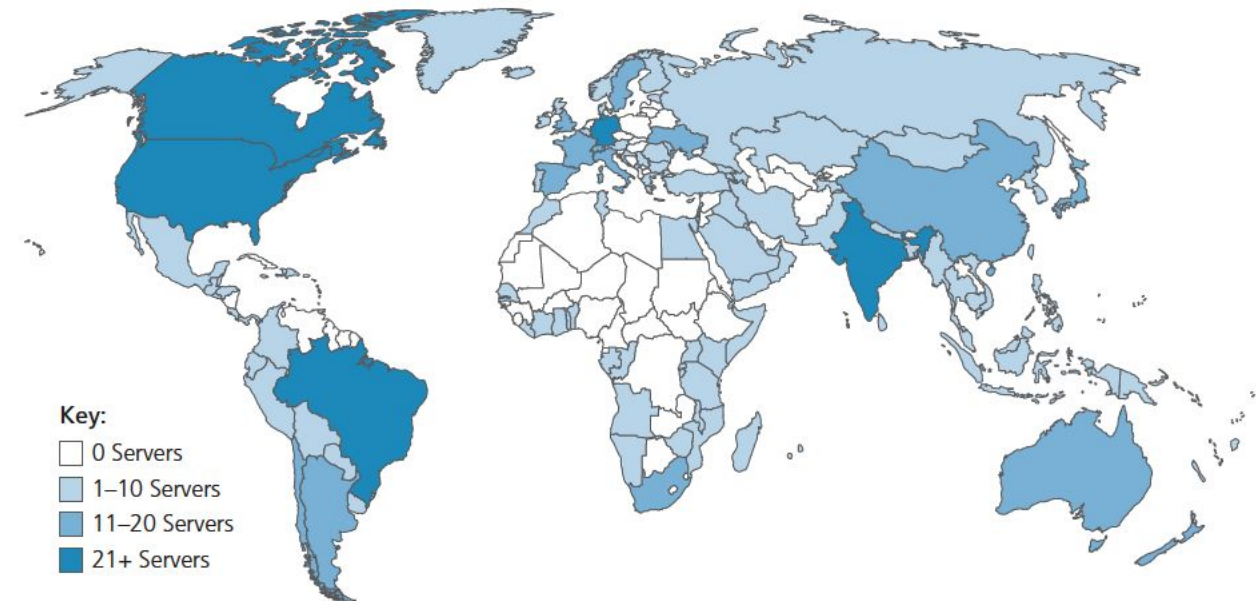
# DNS: a distributed, hierarchical database



Root DNS Servers — *Root*

.com DNS servers   .org DNS servers   .edu DNS servers — *Top Level Domain*

yahoo.com DNS servers   amazon.com DNS servers   pbs.org DNS servers   nyu.edu DNS servers   umass.edu DNS servers — *Authoritative*

Client wants IP address for www.amazon.com; 1$^{st}$ 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

# DNS: root name servers

- official, contact-of-last-resort by name servers that can not resolve name

- *incredibly important* Internet function
  - Internet couldn't function without it!
  - DNSSEC – provides security (authentication and message integrity)

- ICANN (Internet Corporation for Assigned Names and Numbers) manages root DNS domain

13 logical root name "servers" worldwide each "server" replicated many times (~200 servers in US)



Key:
- ☐ 0 Servers
- ☐ 1–10 Servers
- ☐ 11–20 Servers
- ☐ 21+ Servers

# TLD: authoritative servers

Top-Level Domain (TLD) servers:

- responsible for .com, .org, .net, .edu, .aero, .jobs, .museums, and all top-level country domains, e.g.: .cn, .uk, .fr, .ca, .jp
- Network Solutions: authoritative registry for .com, .net TLD
- Educause: .edu TLD

Authoritative DNS servers:

- 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
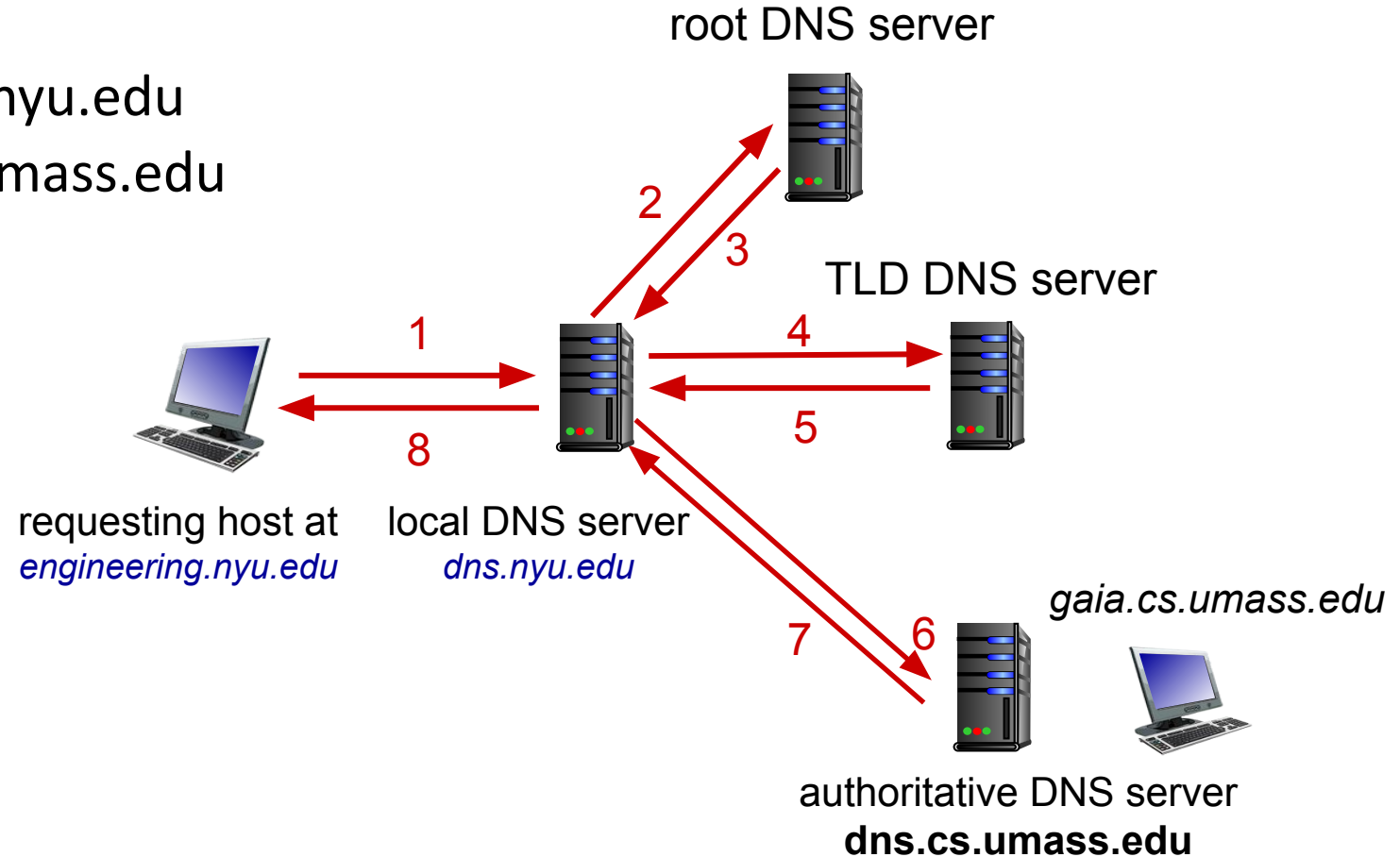
# Local DNS name servers

- does not strictly belong to hierarchy

- each ISP (residential ISP, company, university) has one
  - also called "default name server"

- when host makes DNS query, query is sent to its local DNS server
  - 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: iterated query

Example: host at engineering.nyu.edu
wants IP address for gaia.cs.umass.edu

Iterated query:
- contacted server replies with name of server to contact
- "I don't know this name, but ask this server"



root DNS server

TLD DNS server

requesting host at
*engineering.nyu.edu*

local DNS server
*dns.nyu.edu*

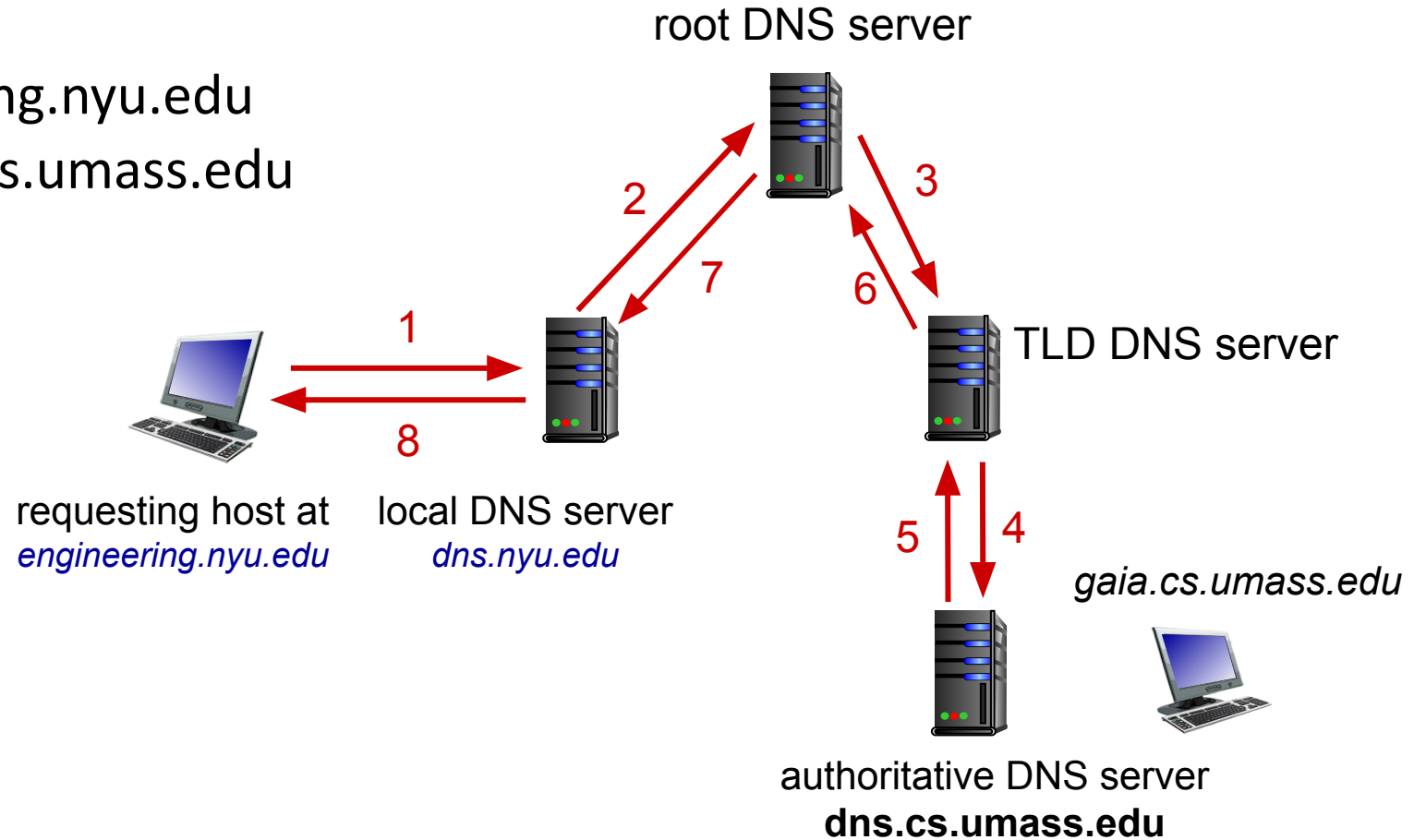*gaia.cs.umass.edu*

authoritative DNS server
**dns.cs.umass.edu**

# DNS name resolution: recursive query

Example: host at engineering.nyu.edu wants IP address for gaia.cs.umass.edu

Recursive query:
- puts burden of name resolution on contacted name server
- heavy load at upper levels of hierarchy?

root DNS server

TLD DNS server

local DNS server
*dns.nyu.edu*

requesting host at
*engineering.nyu.edu*

*gaia.cs.umass.edu*

authoritative DNS server
**dns.cs.umass.edu**

1
2
3
4
5
6
7
8

# Caching, Updating DNS Records

- once (any) name server learns mapping, it *caches* mapping
  - cache entries timeout (disappear) after some time (TTL)
  - TLD servers 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 records

DNS: distributed database storing resource records (RR)

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
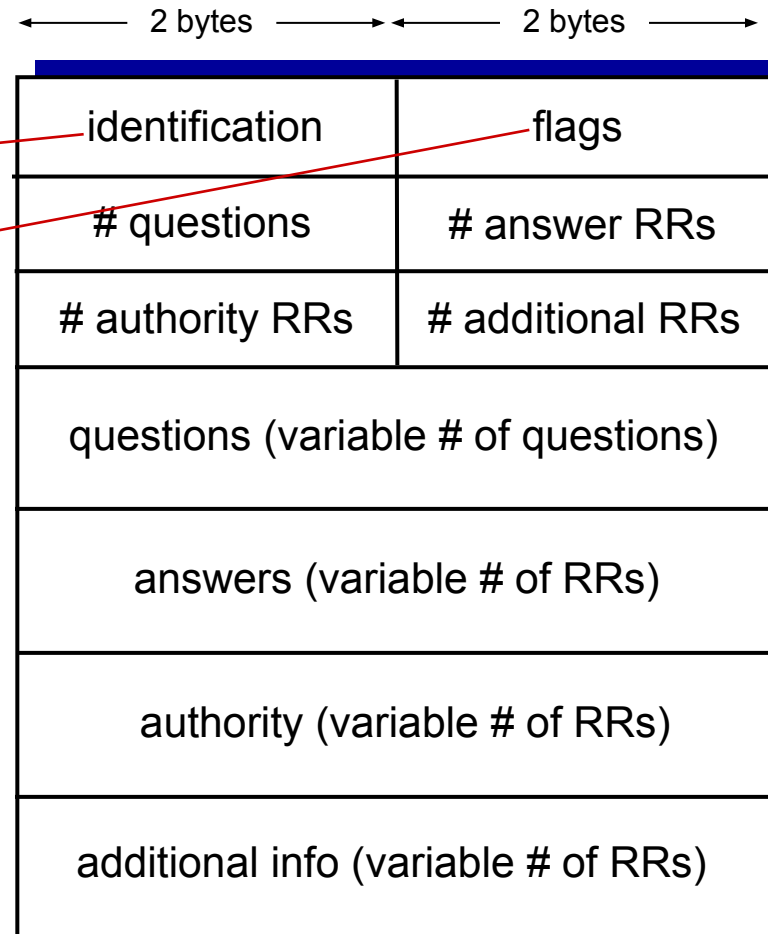- `value` is canonical name

type=MX
- `value` is name of mailserver associated with `name`

# DNS protocol messages

DNS *query* and *reply* messages, both have same *format:*

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

| ← 2 bytes → | ← 2 bytes → |
|:---:|:---:|
| identification | flags |
| # questions | # answer RRs |
| # authority RRs | # additional RRs |
| questions (variable # of questions) ||
| answers (variable # of RRs) ||
| authority (variable # of RRs) ||
| additional info (variable # of RRs) ||

# DNS protocol messages

DNS *query* and *reply* messages, both have same *format:*



name, type fields for a query —— questions (variable # of questions)

RRs in response to query —— answers (variable # of RRs)

records for authoritative servers —— authority (variable # of RRs)

additional " helpful" info that may be used —— additional info (variable # of RRs)

Header fields: identification | flags | # questions | # answer RRs | # authority RRs | # additional RRs

2 bytes | 2 bytes

# Inserting records into DNS

Example: new startup "Network Utopia"

- register name networkuptopia.com at *DNS registrar* (e.g., Network Solutions)
  - provide names, IP addresses of authoritative name server (primary and secondary)
  - registrar inserts NS, A RRs into .com TLD server:

    ```
    (networkutopia.com, dns1.networkutopia.com, NS)
    (dns1.networkutopia.com, 212.212.212.1, A)
    ```

- create authoritative server locally with IP address `212.212.212.1`
  - type A record for www.networkuptopia.com
  - type MX record for networkutopia.com

# DNS security

## 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 DNS queries
- DNS poisoning
  - send bogus relies to DNS server, which caches

## Exploit DNS for DDoS

- send queries with spoofed source address: target IP
- requires amplification
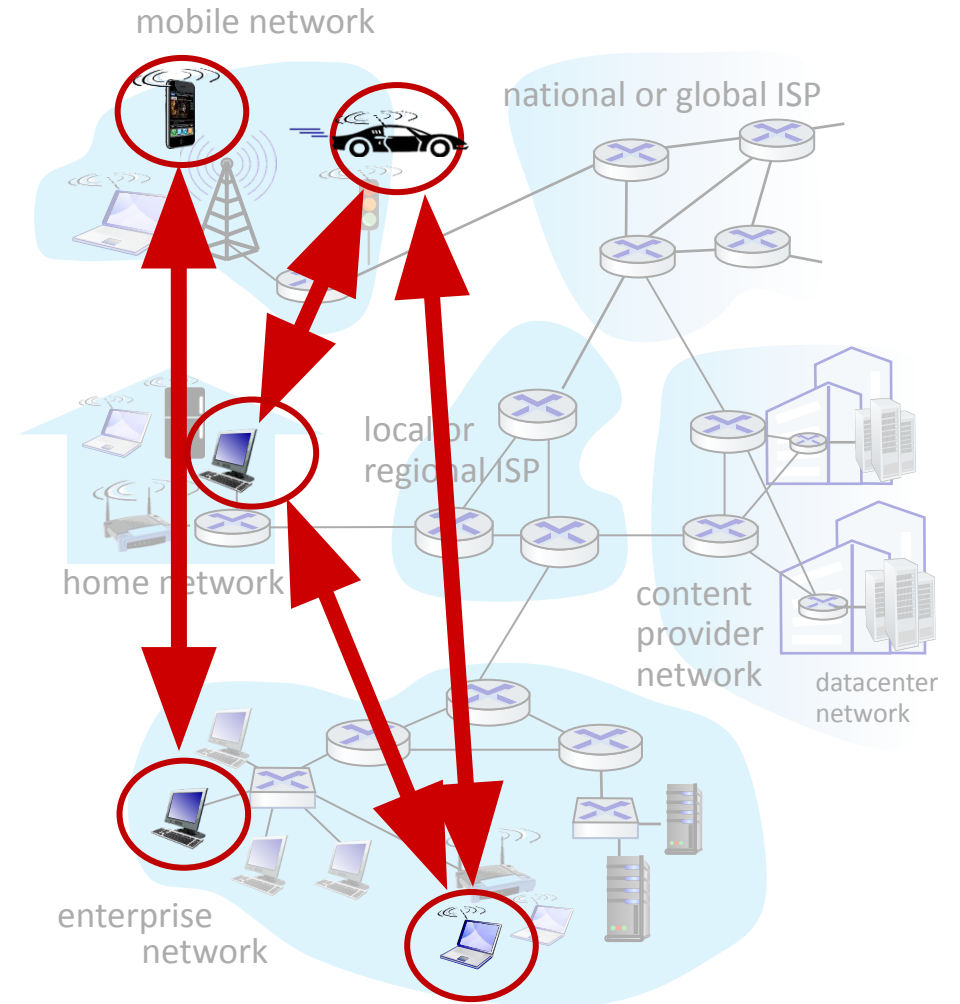
DNSSEC
[RFC 4033]

# Application Layer: Overview

- Principles of network applications

- Web and HTTP

- E-mail, SMTP, IMAP

- The Domain Name System DNS

- **P2P applications**

- video streaming and content distribution networks

- socket programming with UDP and TCP
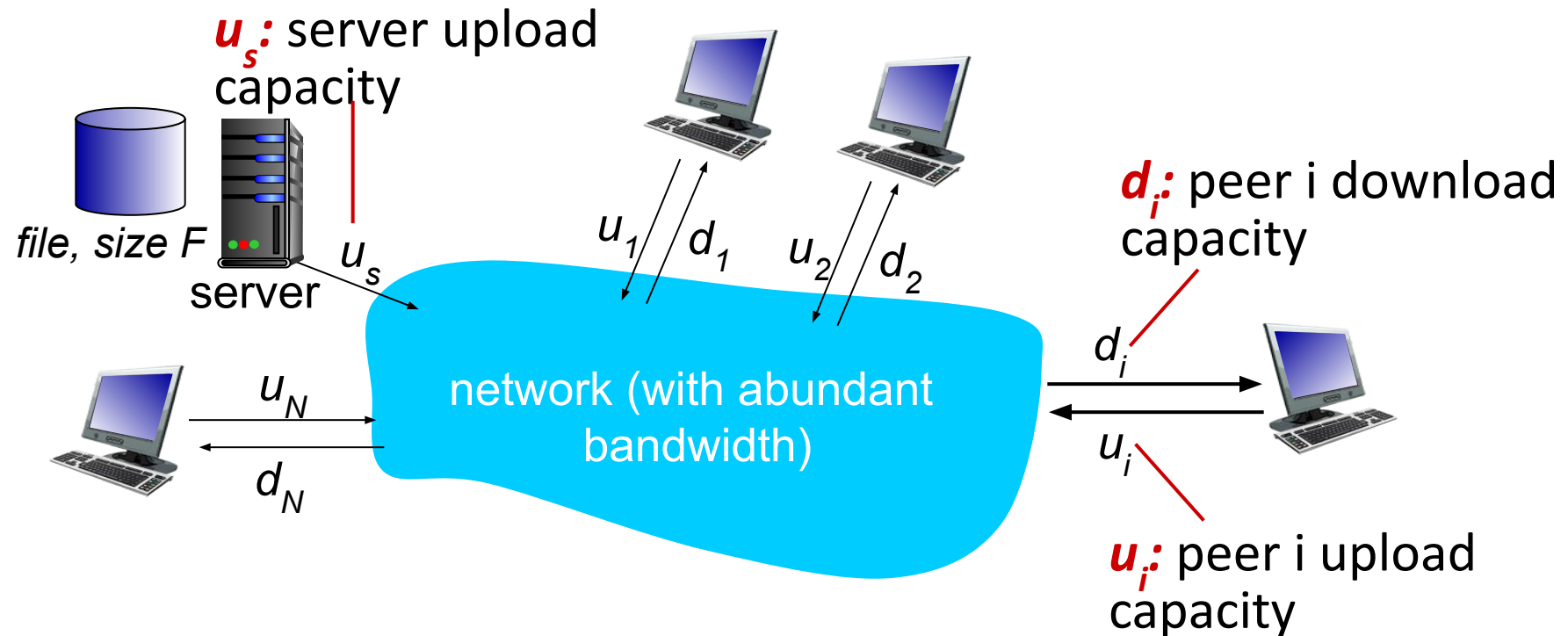
# Peer-to-peer (P2P) architecture

- *no* always-on server

- arbitrary end systems directly communicate

- peers request service from other peers, provide service in return to other peers
  - *self scalability* – new peers bring new service capacity, and new service demands

- peers are intermittently connected and change IP addresses
  - complex management

- examples: P2P file sharing (BitTorrent), streaming (KanKan), VoIP (Skype)

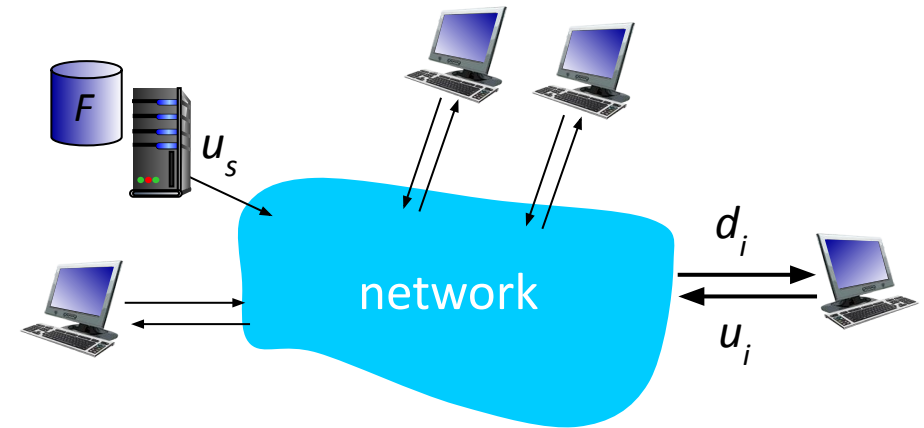# File distribution: client-server vs P2P

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

- peer upload/download capacity is limited resource



$u_s$: server upload capacity

*file, size F*

server

$u_s$

$u_1$ $d_1$ $u_2$ $d_2$

network (with abundant bandwidth)

$u_N$

$d_N$

$d_i$: peer i download capacity

$d_i$

$u_i$

$u_i$: peer i upload capacity

# 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
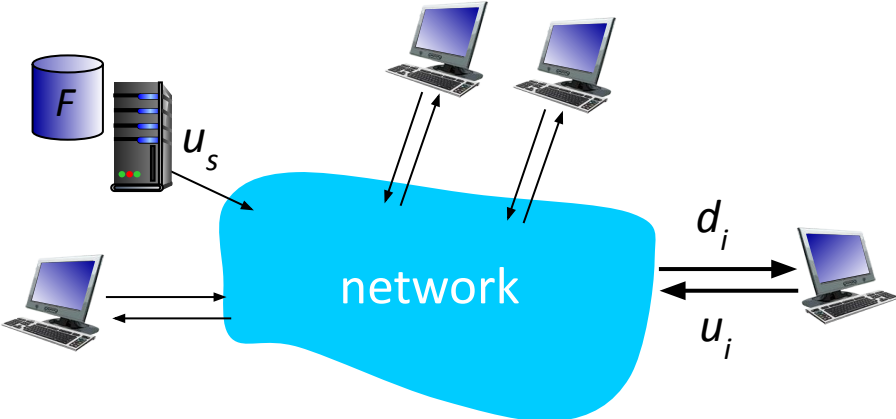  - min client download time: $F/d_{min}$



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
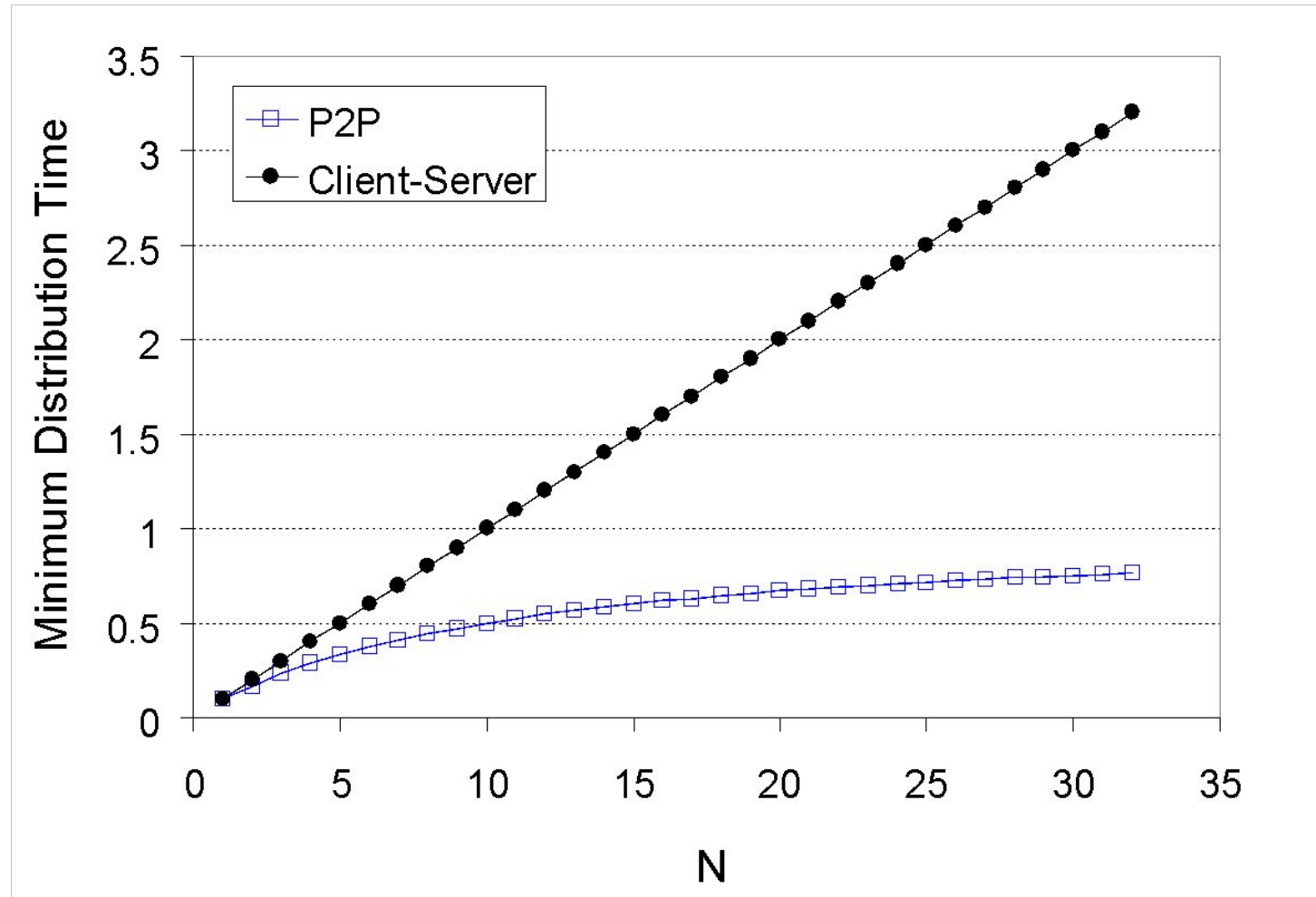  - max upload rate (limiting max download rate) is $u_s + \Sigma u_i$



time to distribute F to N clients using P2P approach

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

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,  $u_s = 10u$,  $d_{min} \geq u_s$
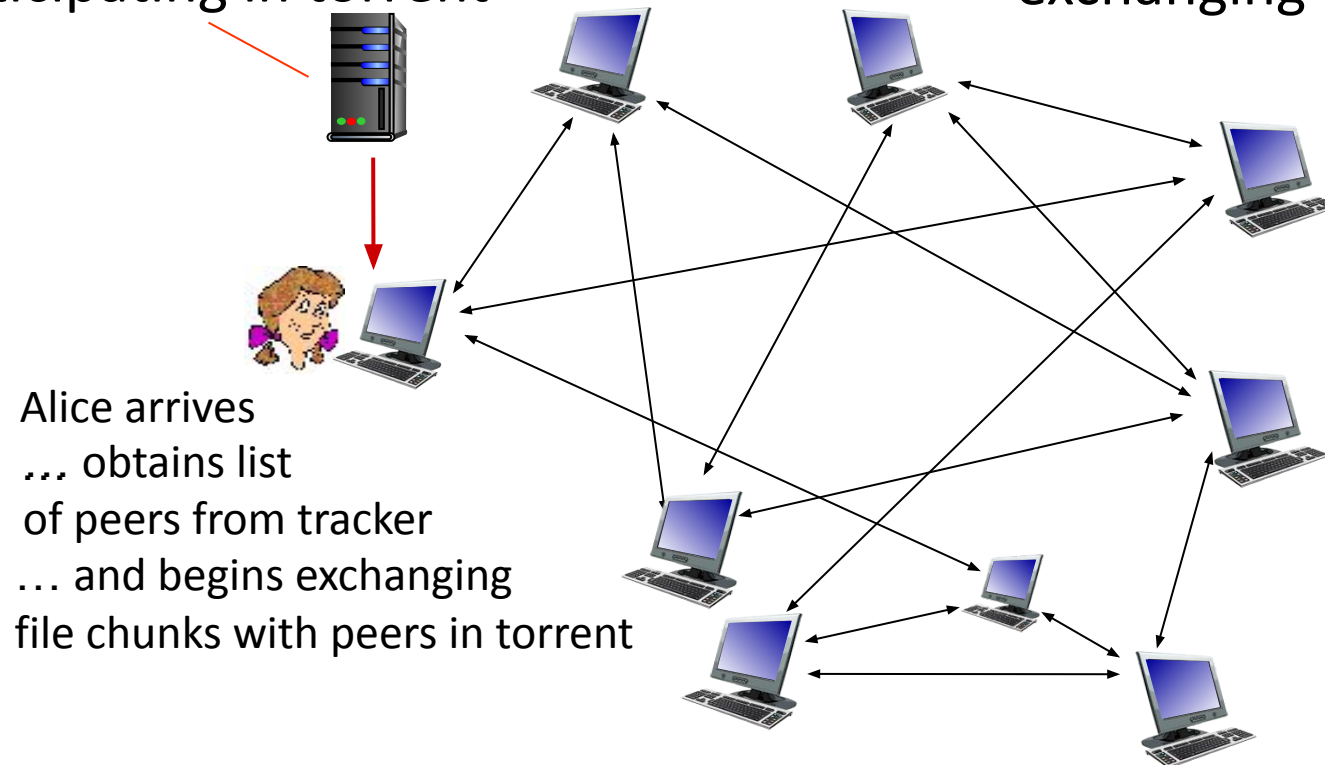
# P2P file distribution: BitTorrent

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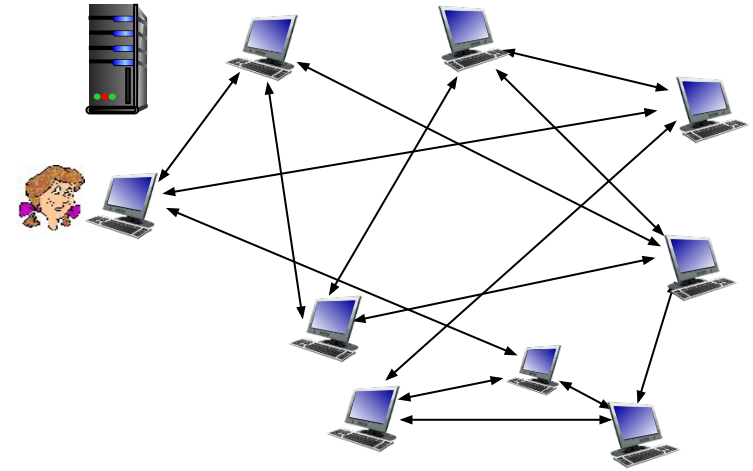- peers in torrent send/receive file chunks

*tracker:* tracks peers participating in torrent

*torrent:* group of peers exchanging  chunks of a file

Alice arrives
*…* obtains list
of peers from tracker
… and begins exchanging
file chunks with peers in torrent

# 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

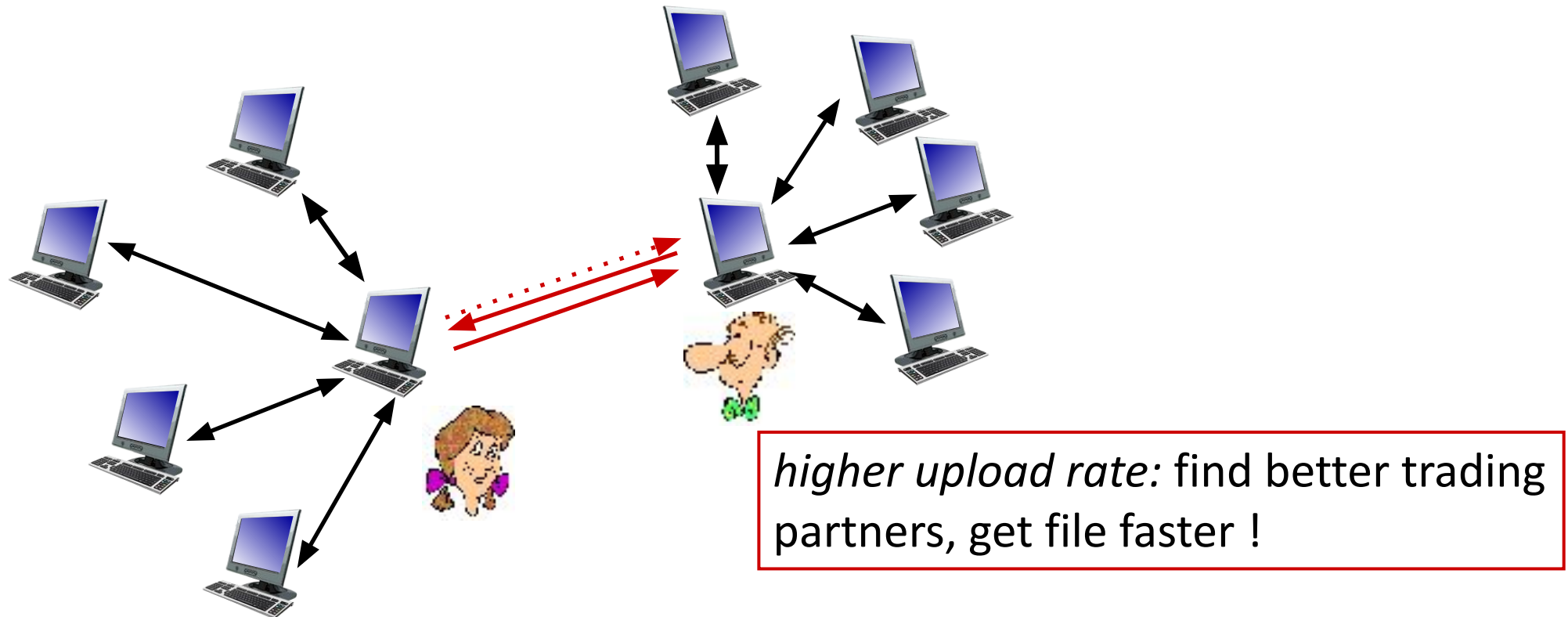# BitTorrent: requesting, sending file chunks

## 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 every10 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

*higher upload rate:* find better trading partners, get file faster !

# Application layer: overview

- Principles of network applications

- Web and HTTP

- E-mail, SMTP, IMAP

- The Domain Name System DNS

- P2P applications

- **video streaming and content distribution networks**

- socket programming with UDP and TCP
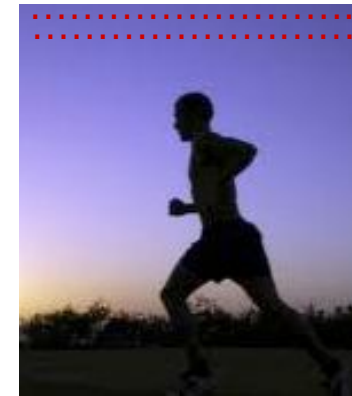
# Video Streaming and CDNs: context

- stream video traffic: major consumer of Internet bandwidth
  - Netflix, YouTube, Amazon Prime: 80% of residential ISP traffic (2020)
- 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*

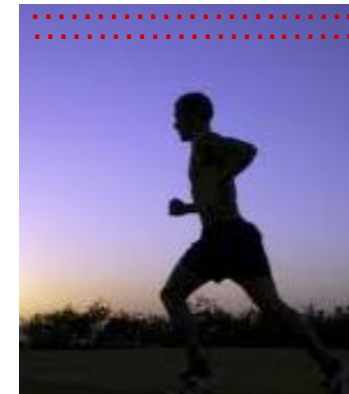*temporal coding example:* instead of sending complete frame at i+1, send only differences from frame i



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 1 (CD-ROM) 1.5 Mbps
  - MPEG2 (DVD) 3-6 Mbps
  - MPEG4 (often used in Internet,  64Kbps – 12 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)
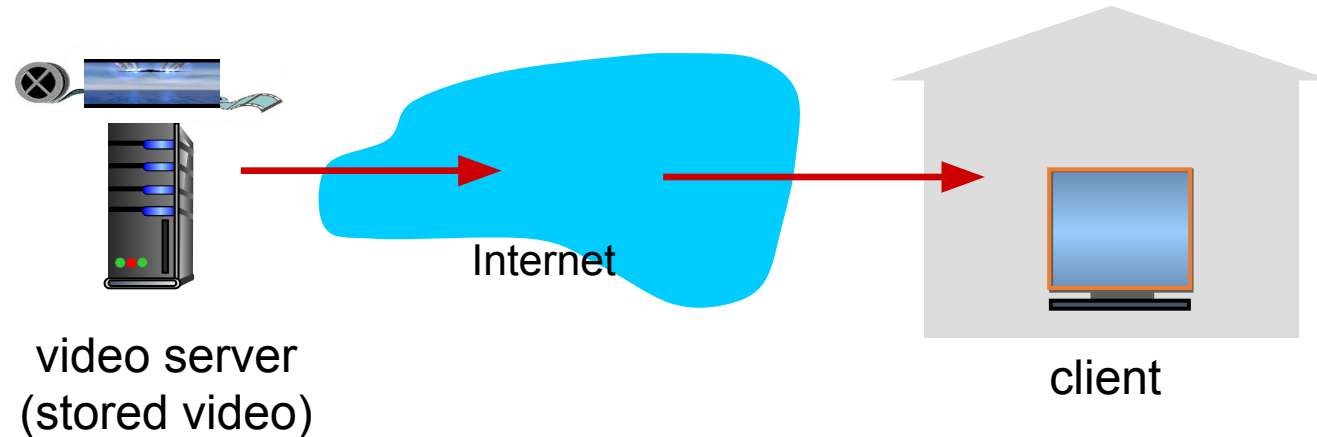
frame *i*

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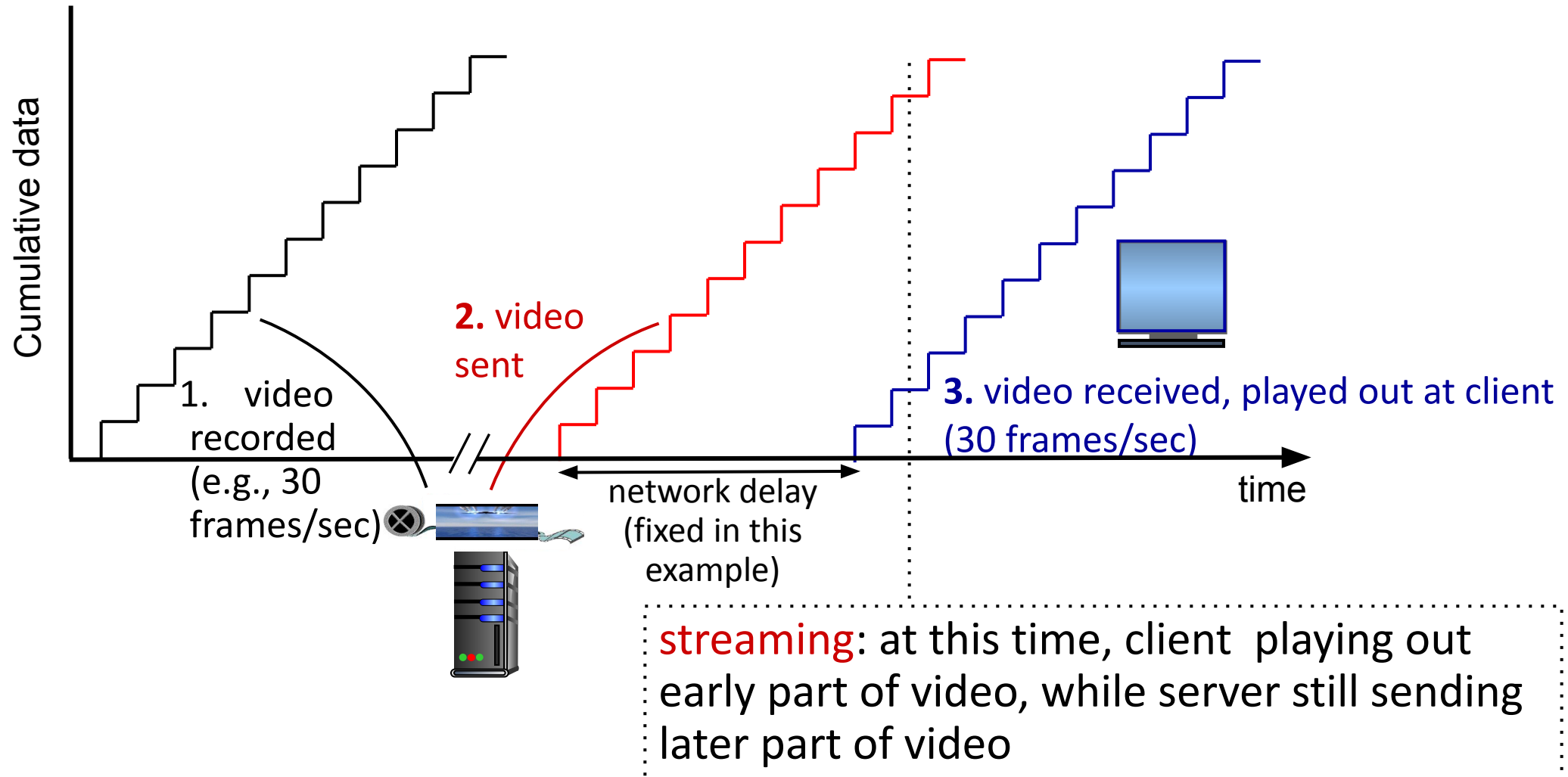
frame *i+1*

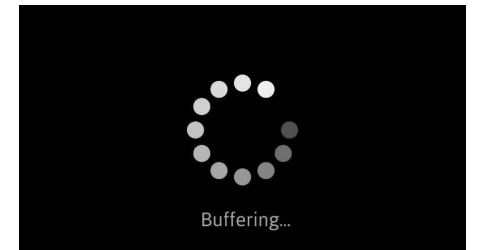# Streaming stored video

simple scenario:



Main challenges:

- server-to-client bandwidth will *vary* over time, with changing network congestion levels (in house, in access network, in network core, at video server)

- packet loss and delay due to congestion will delay playout, or result in poor video quality

# Streaming stored video



Cumulative data (y-axis) vs time (x-axis)

1. video recorded (e.g., 30 frames/sec)

**2.** video sent

network delay (fixed in this example)

**3.** video received, played out at client (30 frames/sec)

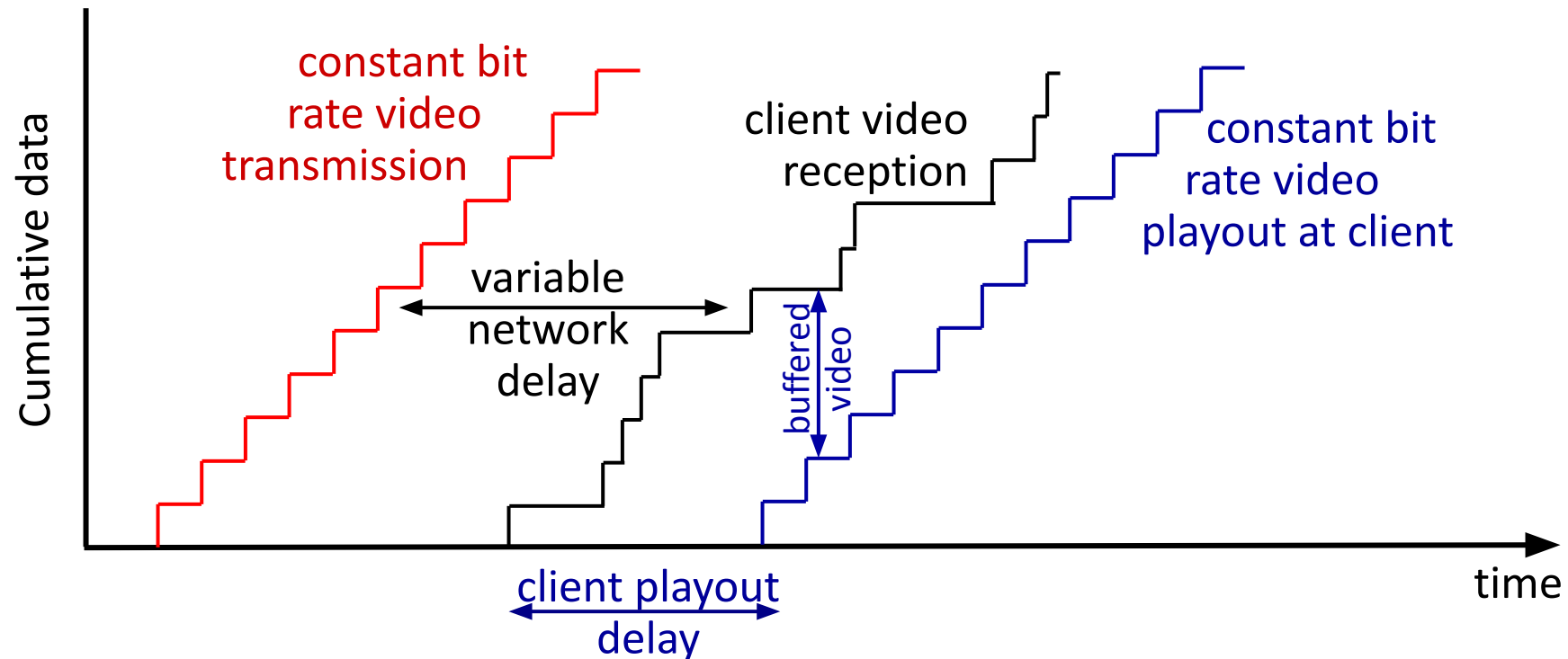**streaming**: at this time, client playing out early part of video, while server still sending later part of video

# Streaming stored video: challenges

- continuous playout constraint: once client playout begins, playback must match original timing
  - … but network delays are variable (jitter), so will need client-side buffer to match playout requirements

- other challenges:
  - client interactivity: pause, fast-forward, rewind, jump through video
  - video packets may be lost, retransmitted

Buffering...

# Streaming stored video: playout buffering



- *client-side buffering and playout delay:* compensate for network-added delay, delay jitter
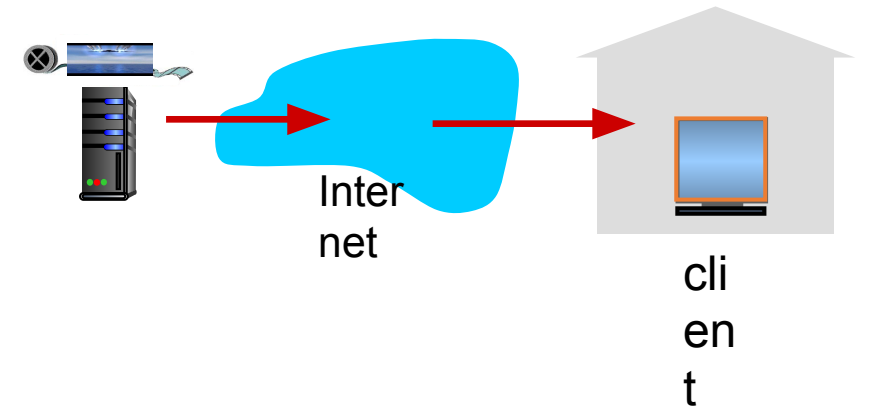
# Streaming multimedia: DASH

- *DASH: Dynamic, Adaptive Streaming over HTTP*

- *server:*
  - divides video file into multiple chunks
  - each chunk stored, encoded at different rates
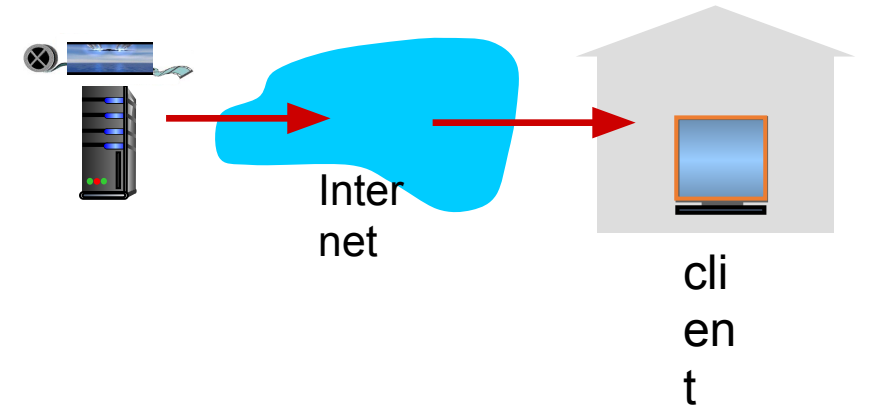  - *manifest file:* provides URLs for different chunks

- *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)

Inter net

cli en t

# Streaming multimedia: DASH

- *"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)

Streaming video = encoding + DASH + playout buffering

# Content distribution networks (CDNs)

- *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 (CDNs)

- *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
    - Akamai: 240,000 servers deployed in more than 120 countries (2015)
  - *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: "over the top"*

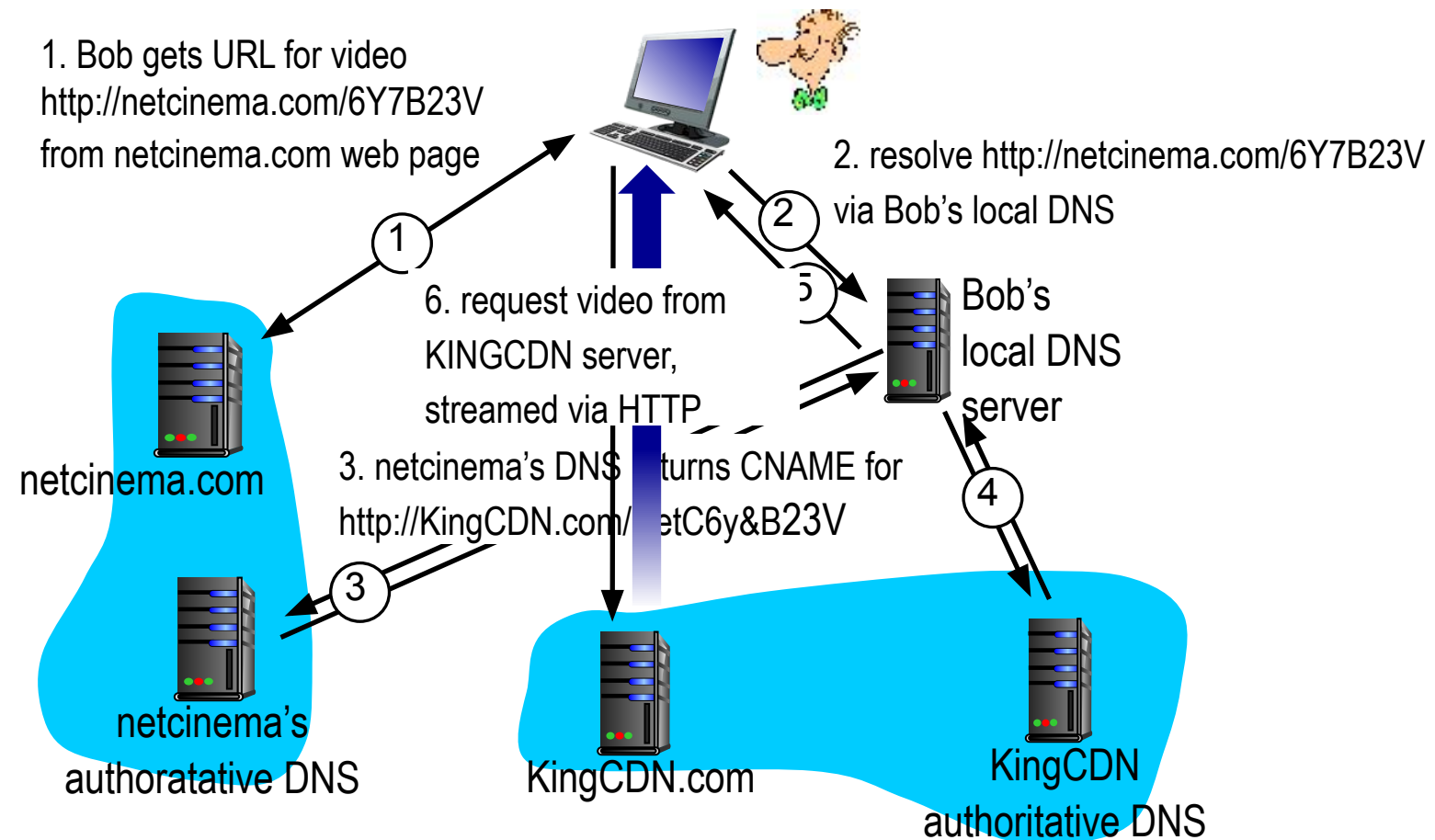Internet host-host communication as a service

*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?

# CDN content access: a closer look

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

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

1. Bob gets URL for video
http://netcinema.com/6Y7B23V
from netcinema.com web page

②

2. resolve http://netcinema.com/6Y7B23V
via Bob's local DNS

⑤

Bob's local DNS server

6. request video from
KINGCDN server,
streamed via HTTP

netcinema.com

3. netcinema's DNS returns CNAME for
http://KingCDN.com/NetC6y&B23V

③

④

netcinema's
authoratative DNS

KingCDN.com

KingCDN
authoritative DNS

# Case study: Netflix



Netflix registration, accounting servers

Amazon cloud

upload copies of multiple versions of video to CDN servers

CDN server

CDN server

CDN server

Bob browses Netflix video

② 

Manifest file, requested returned for specific video

③

①

Bob manages Netflix account

④

DASH server selected, contacted, streaming begins
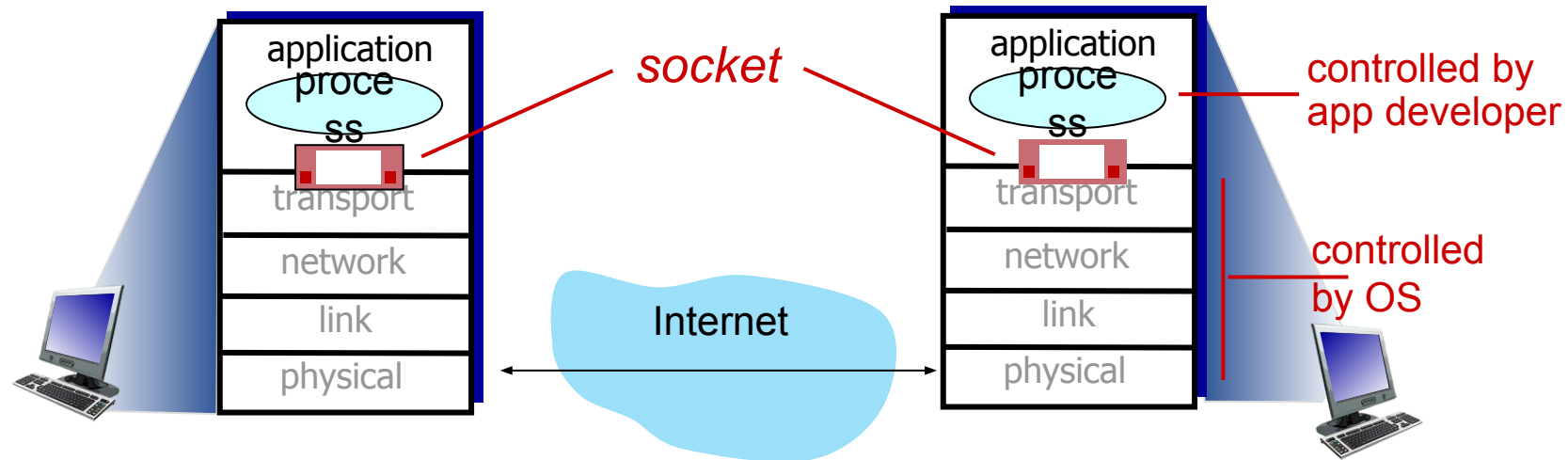
# Application Layer: Overview

- Principles of network applications

- Web and HTTP

- E-mail, SMTP, IMAP

- The Domain Name System DNS

- P2P applications

- video streaming and content distribution networks

- **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

# TCP vs UDP

**UDP: User Datagram Protocol**

- no acknowledgements
- no retransmissions
- out of order, duplicates possible
- connectionless, i.e., app indicates destination for each packet

**TCP: Transmission Control Protocol**

- reliable byte-stream channel (in order, all arrive, no duplicates)
- similar to file I/O
- flow control
- connection-oriented
- bidirectional

# TCP vs UDP

TCP is used for services with a large data capacity, and a persistent connection

UDP is more commonly used for quick lookups, and single use query-reply actions.
Some common examples of TCP and UDP with their default ports:

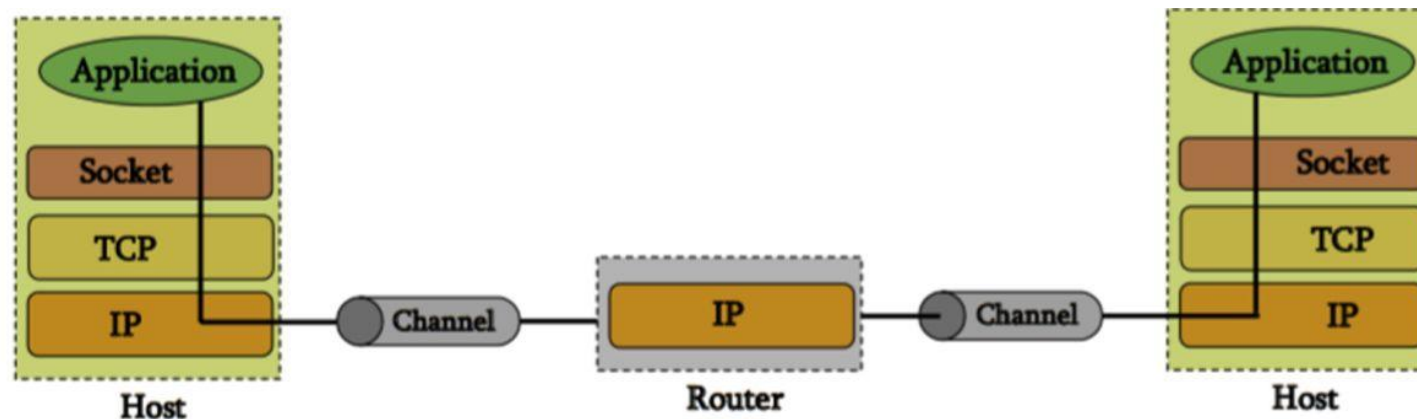| DNS lookup | UDP | 53 |
| FTP | TCP | 21 |
| HTTP | TCP | 80 |
| POP3 | TCP | 110 |
| Telnet | TCP | 23 |

# Berkley Sockets

Universally known as Sockets

It is an abstraction through which an application may send and receive data

Provide generic access to interprocess communication services (e.g. IPX/SPX, Appletalk, TCP/IP)

Standard API for networking

# Sockets

Uniquely identified by: an internet address, an end-to-end protocol (e.g. TCP or UDP), a port number
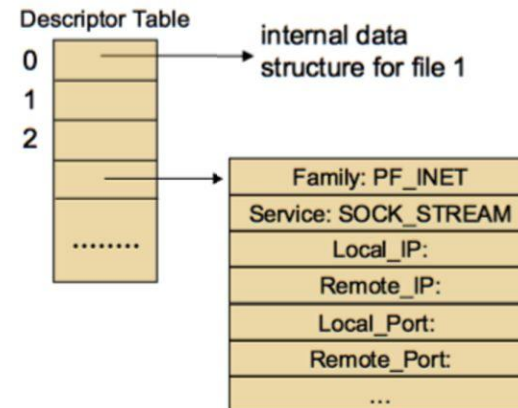
Two types of (TCP/IP) sockets:

Stream sockets (e.g. uses TCP) - provide reliable byte-stream service

Datagram sockets (e.g. uses UDP): provide best-effort datagram service, messages up to 65.500 bytes

Socket extend the convectional UNIX I/O facilities:

file descriptors for network communication, extended the read and write system calls

Descriptor Table

| 0 | |
|---|---|
| 1 | |
| 2 | |

→ internal data structure for file 1

........

| Family: PF_INET |
|---|
| Service: SOCK_STREAM |
| Local_IP: |
| Remote_IP: |
| Local_Port: |
| Remote_Port: |
| ... |

# Sockets

# Client-Server Communication

**Server**

• passively waits for and responds to clients

• passive socket


**Client**

• initiates the communication

• must know the address and the port of the server

• active socket

# Sockets - Procedures

| Procedure | Meaning |
|-----------|---------|
| Socket | Create a new communication endpoint |
| Bind | Attach a local address to a socket |
| Listen | Announce willingness to accept connections |
| Accept | Block caller until a connection request arrives |
| Connect | Actively attempt to establish a connection |
| Send | Send some data over the connection |
| Receive | Receive some data over the connection |
| Close | Release the connection |

# Client-Server Communication

# Socket creation in C: socket ()

**fint sockid = socket(family, type, protocol);**

**sockid**: socket descriptor, an integer (like a file-handle)

**family**: integer, communication domain, e.g.,

PF_INET, IPv4 protocols, Internet addresses (typically used)

PF_UNEX, Local communication, File addresses

**type**: communication type

SOCK_STREAM - reliable, 2-way, connection-based service

SOCK_DGRAM - unreliable, connectionless, messages of maximum length
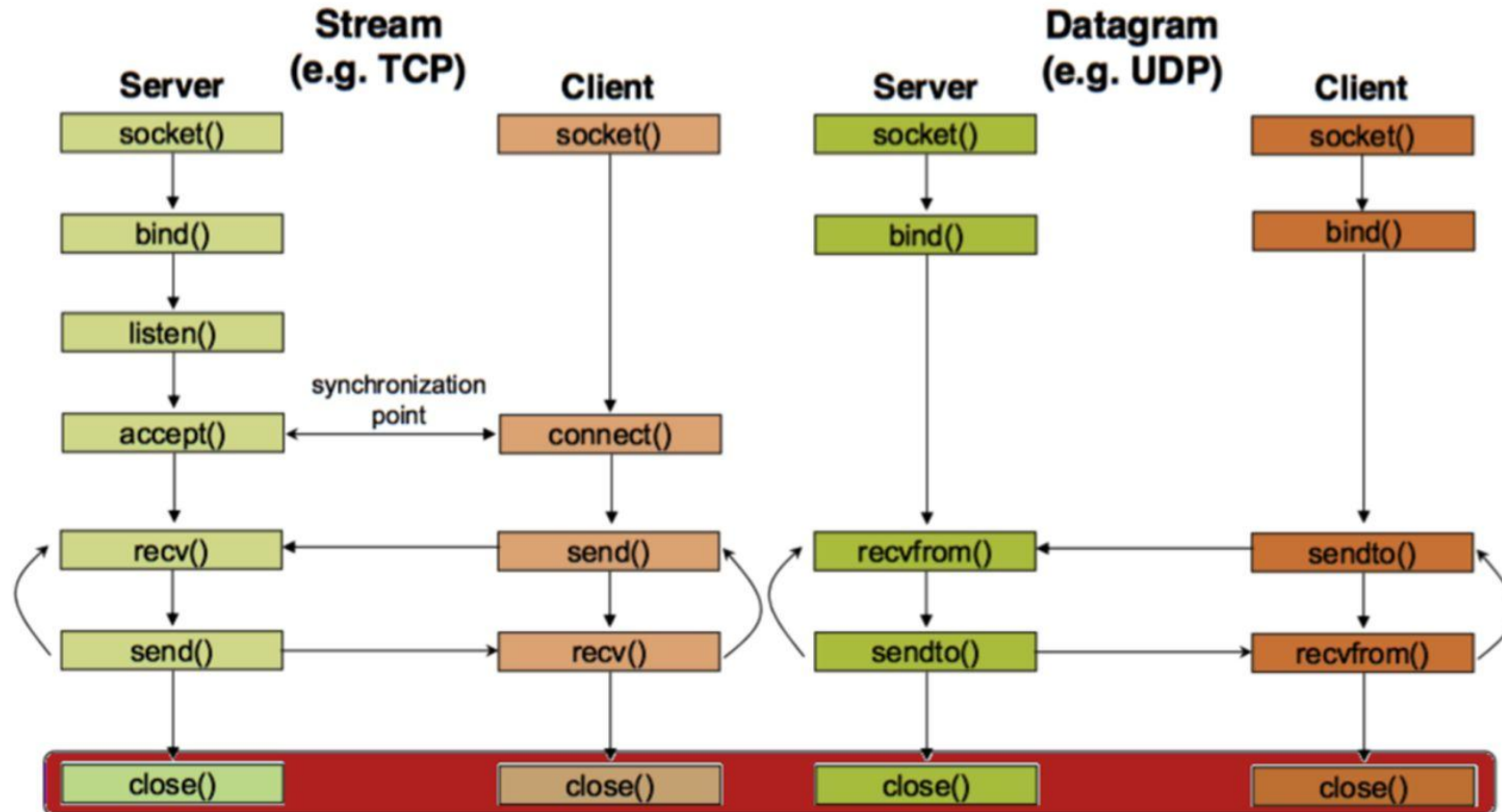
**protocol**: specifies protocol

IPPROTO_TCP IPPROT0_UDP

usually set to 0 (i.e., use default protocol)

**upon failure returns -1**

**NOTE**: socket call does not specify where data will be coming from, nor where it will be going to - it just creates the interface!

# Client-Server Communication

# Socket close in C: close ()

When finished using a socket, the socket should be closed

**status = close(sockid);**

    sockid: the file descriptor (socket being closed)

    status: 0 if successful, -1 if error

Closing a socket

    closes a connection (for stream socket)

    frees up the port used by the socket

# Specifying Addresses
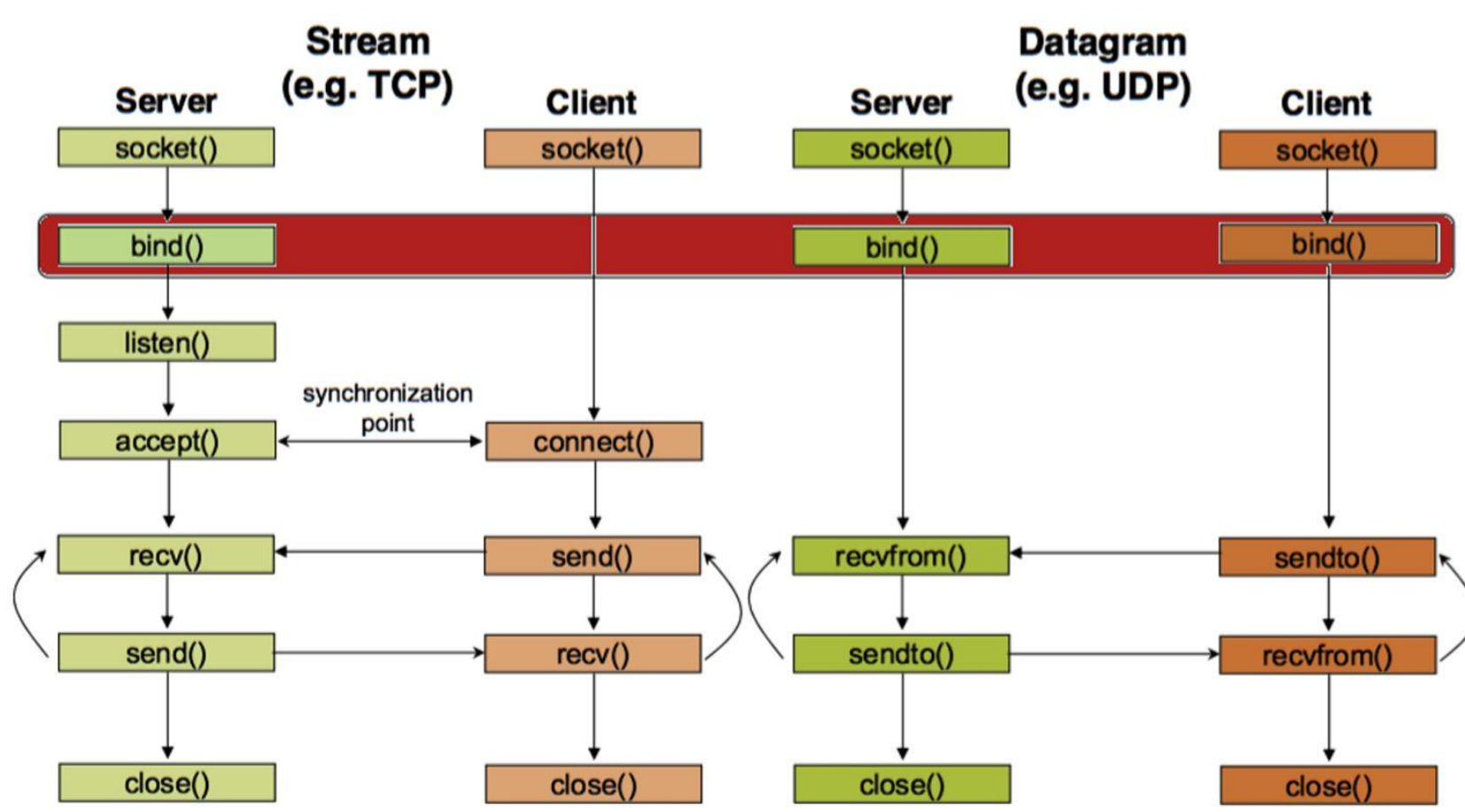
Socket API defines a generic data type for addresses:

**struct sockaddr {**

**unsigned short sa__family;** /* Address family (e.g. AF_INET) 7 **char sa_data [14] ;** /* Family-specific address information 7

**}**

Particular form of the sockaddr used for TCP/IP addresses:

**struct in_addr {**

**unsigned long s_addr;** /* Internet address (32 bits) 7

**}**

**struct sockaddr_in {**

**unsigned short sin_family;** /* Internet protocol (AF_INET) 7 **unsigned short sin_port;** /* Address port (16 bits) 7 **struct in_addr sin_addr;** /* Internet address (32 bits) 7 **char sin_zero [ 8 ] ;** /* Not used 7

**}**

Important: sockaddr_in can be casted to a sockaddr

# Client-Server Communication

# Assign address to socket: bind ()

associates and reserves a port for use by the socket

int status = bind(sockid, fiaddrport, size);

sockid: integer, socket descriptor

addrport: struct sockaddr, the (IP) address and port of the machine

for TCP/IP server, internet address is usually set to INADDR_ANY, i.e., chooses any incoming interface

size: the size (in bytes) of the addrport structure

status: upon failure -1 is returned

# bind () - Example with TCP

```
int soclcid;

struct sockaddr_in addrport;

soclcid = socket (PF_INET , SOCK_STREAM, 0) ;


addrport. si n__f ami ly = AF_INET;
addrport.sin_port = htons(5100);
addrport.sin_addr.s_addr = htonl(INADDR_ANY);
if(bind(sockid,   (struct   sockaddr   *)   &addrport,
sizeof(addrport))!= -1)    {
…}
```

# Skipping    the    bind    ()
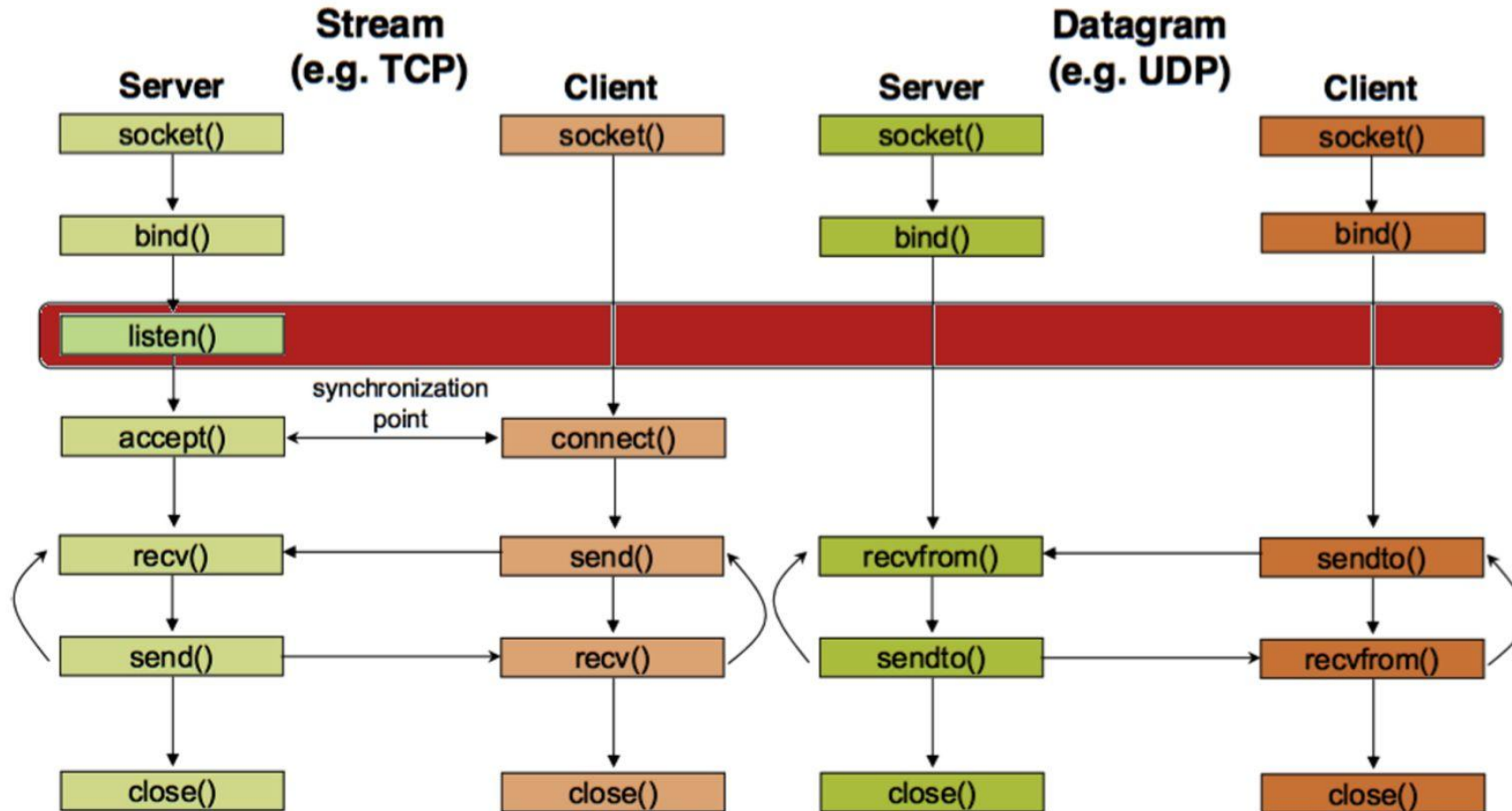
bind() can be skipped for both types of sockets

**Datagram socket:**

• if only sending, no need to bind. The OS finds a port each time the socket sends a packet

• if receiving, need to bind

**Stream socket:**

• destination determined during connection setup

• don't need to know port sending from (during connection setup, receiving end is informed of port)

# Client-Server Communication

# listen ()

Instructs TCP protocol implementation to listen for connections

**int status = listen(sockid, queueLimit);**

sockid: integer, socket descriptor

queuelen: integer, # of active participants that can "wait" for a connection
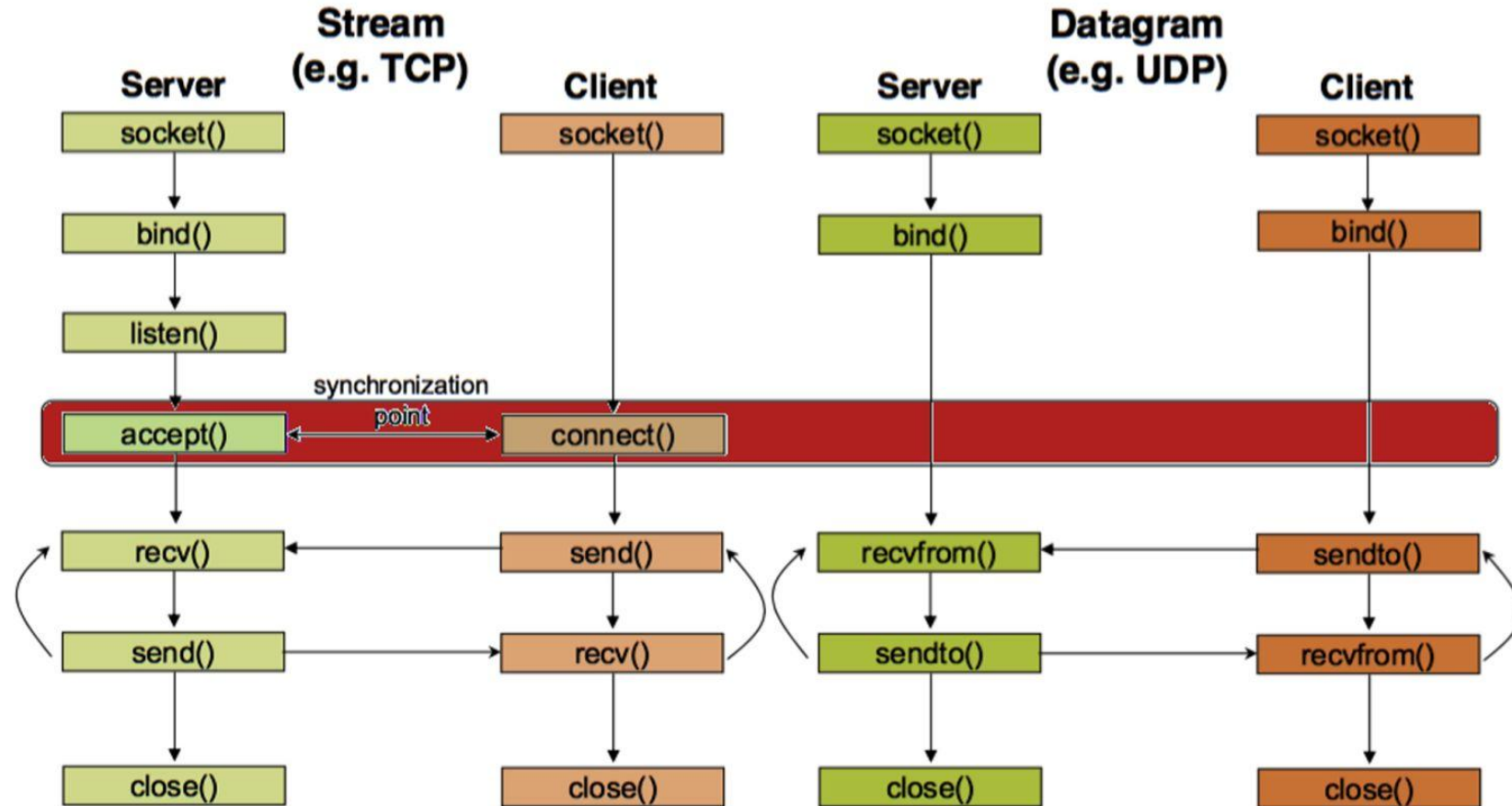
status: 0 if listening, -1 if error

listen () is non-blocking: returns immediately

The listening socket (sockid)
is never used for sending and receiving
is used by the server only as a way to get new sockets

# Client-Server Communication

# Establish Connection: connect ()

The client establishes a connection with the server by calling connect()

**int status = connect(sockid, &foreignAddr, addrlen);**

sockid: integer, socket to be used in connection

foreignAddr: struct sockaddr: address of the passive participant

addrlen: integer, sizeof(name)

status: 0 if successful connect, -1 otherwise

connect () is blocking

# Incoming Connection: accept ()

The server gets a socket for an incoming client connection by calling accept()

**int s = accept(sockid, ficlientAddr, SaddrLen);**

s: integer, the new socket (used for data-transfer)

sockid: integer, the orig. socket (being listened on)

clientAddr: struct sockaddr, address of the active participant

filled in upon return

addrLen: sizeof(clientAddr): value/result parameter
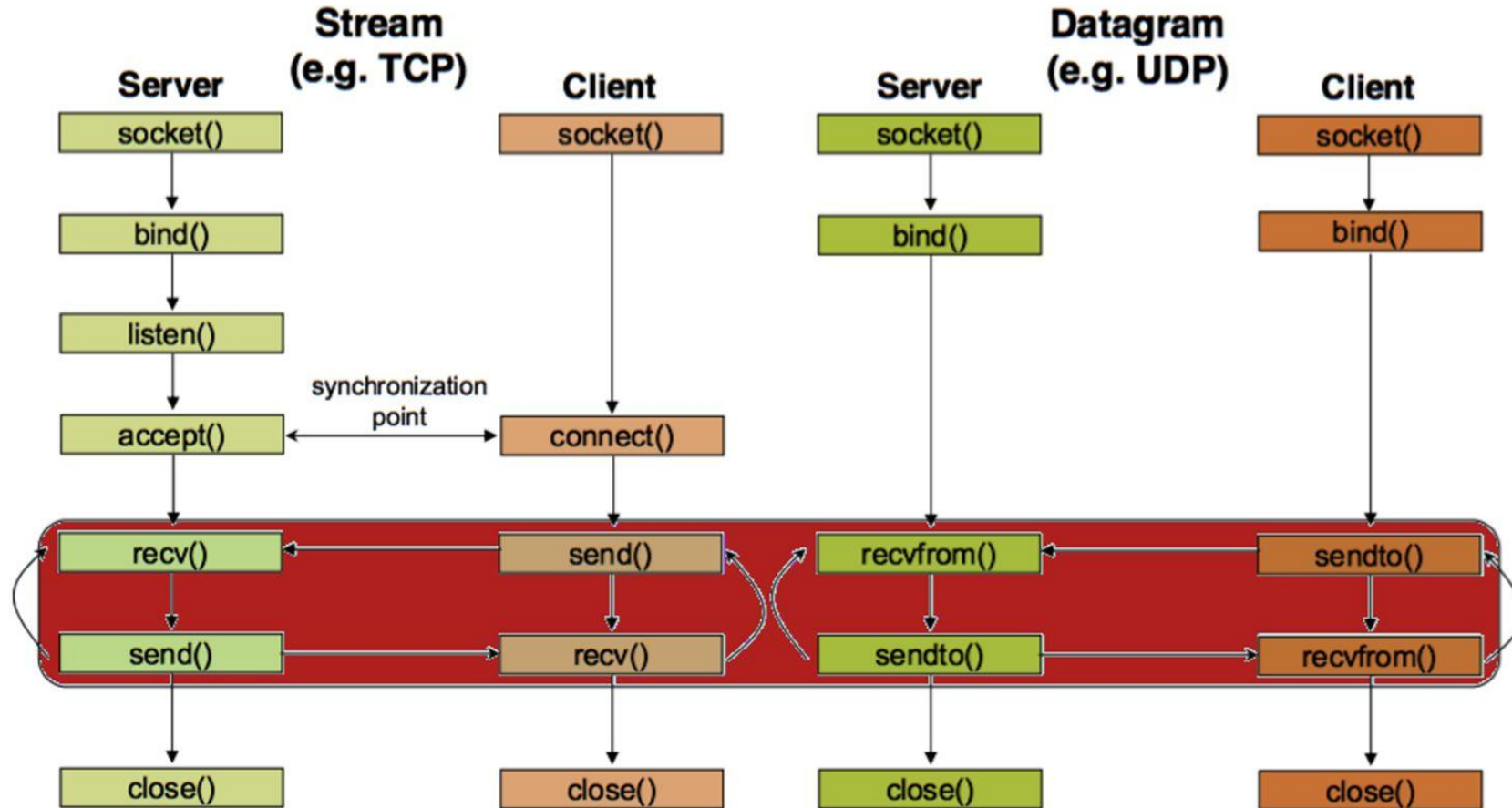
must be set appropriately before call

adjusted upon return

accept()

is blocking: waits for connection before returning

dequeues the next connection on the queue for socket (sockid)

# Client-Server Communication

# Exchanging data with stream socket

**int count = send(sockid, msg, msgLen, flags);**

msg: const void[], message to be transmitted

msgLen: integer, length of message (in bytes) to transmit

flags: integer, special options, usually just 0

count: # bytes transmitted (-1 if error)

**int count = recv(sockid, recvBuf, bufLen, flags);**

recvBuf: void[], stores received bytes

bufLen: # bytes received

flags: integer, special options, usually just 0

count: # bytes received (-1 if error)

Calls are blocking

returns only after data is sent / received

# Exchanging data with datagram socket

**int count = sendto(sockid, msg, msgLen, flags, &foreignAddr, addrlen);**

    msg, msgLen, flags, count: same with send ()

    foreignAddr: struct sockaddr, address of the destination

    addrLen: sizeof(foreignAddr)

**int count = recvfrom(sockid, recvBuf, bufLen, flags, &clientAddr, addrlen) ;**

    recvBuf, bufLen, flags, count: same with recv ()

    clientAddr: struct sockaddr, address of the client

    addrLen: sizeof(clientAddr)

Calls are blocking

    returns only after data is sent / received

# 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

# 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 ⟶

attach server name, port to message; send into socket ⟶ message = raw_input('Input lowercase sentence:')

clientSocket.sendto(message.encode(),

read reply characters from socket into string ⟶                    (serverName, serverPort))

modifiedMessage, serverAddress =

print out received string and close socket ⟶                    clientSocket.recvfrom(2048)

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))

loop forever ⟶ print ("*The server is ready to receive*")

Read from UDP socket into message, getting ⟶ while True:
client's address (client IP and port)

    message, clientAddress = serverSocket.recvfrom(2048)

send upper case string back to this client ⟶ modifiedMessage = message.decode().upper()

    serverSocket.sendto(modifiedMessage.encode(),

        clientAddress)

# 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
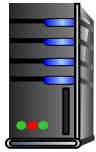
**Client contacts server by:**

- 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
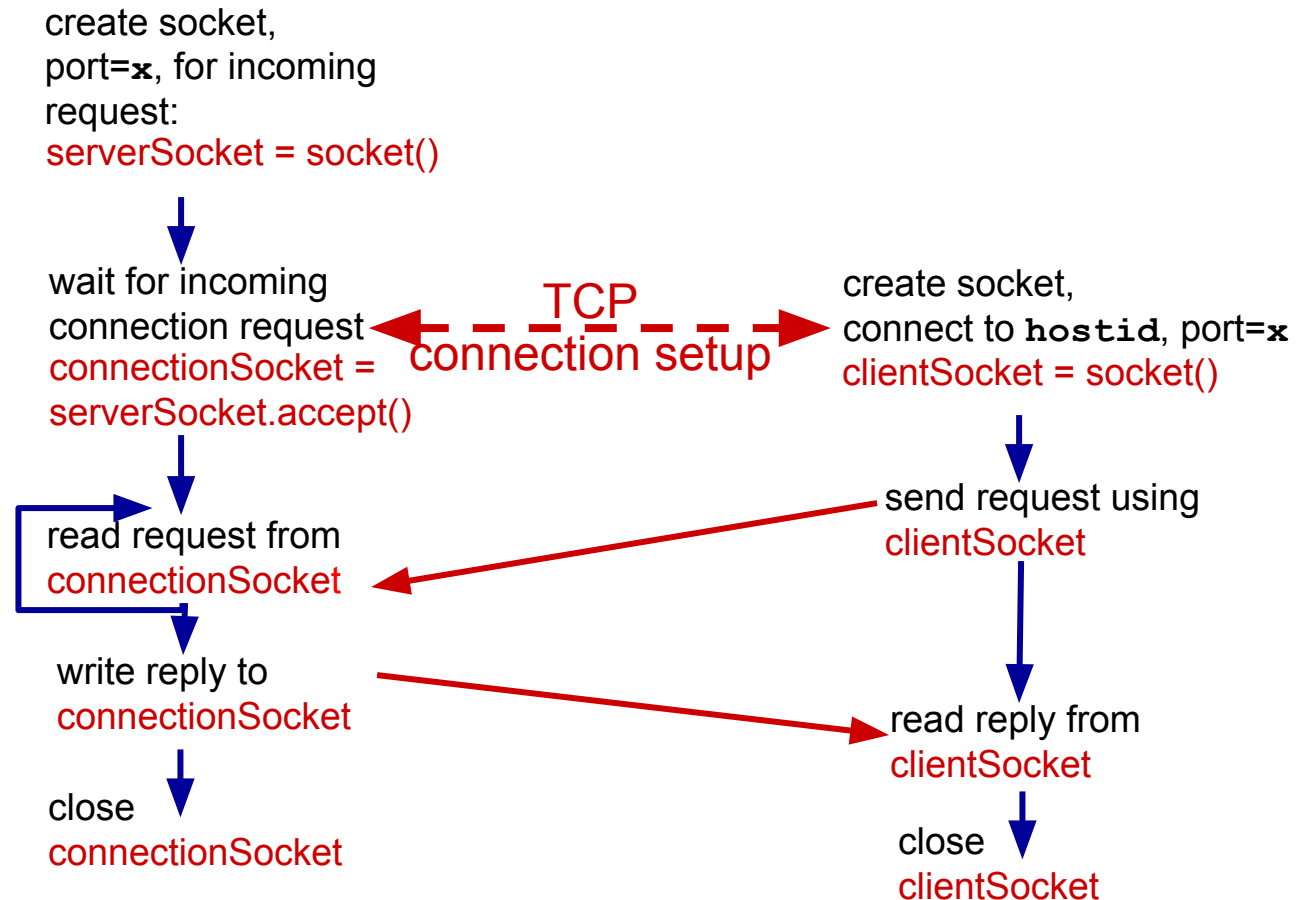
**Application viewpoint**

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

# Client/server socket interaction: TCP

server (running on hostid)　　　client

**create socket,**
**port=x, for incoming**
**request:**
serverSocket = socket()

↓

wait for incoming
connection request ← − − TCP − − → create socket,
connectionSocket =  connection setup  connect to **hostid**, port=**x**
serverSocket.accept()  clientSocket = socket()

↓　　　　　　　　　　　　　　　　　　↓

read request from ← send request using
connectionSocket  clientSocket

↓

write reply to
connectionSocket ──→ read reply from
 clientSocket

↓　　　　　　　　　　　　　　　　　　↓

close  close
connectionSocket  clientSocket

# Example app: TCP client

*Python TCPClient*

from socket import *

serverName = 'servername'

serverPort = 12000

create TCP socket for server, remote port 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 →

server begins listening for
incoming TCP requests →

loop forever →

server waits on accept() for incoming
requests, new socket created on return →

read bytes from socket (but
not address as in UDP) →

close connection to this client (but *not*
welcoming socket) →

```
from socket import *
serverPort = 12000
serverSocket = socket(AF_INET,SOCK_STREAM)
serverSocket.bind(('',serverPort))
serverSocket.listen(1)
print 'The server is ready to receive'
while True:
    connectionSocket, addr = serverSocket.accept()

    sentence = connectionSocket.recv(1024).decode()
    capitalizedSentence = sentence.upper()
    connectionSocket.send(capitalizedSentence.
                                        encode())
    connectionSocket.close()
```

# Topic 2: Summary

our study of network application layer is 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, IMAP
  - DNS
  - P2P: BitTorrent
- video streaming, CDNs
- socket programming: TCP, UDP sockets

# Topic 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:
  - centralized vs. decentralized
  - stateless vs. stateful
  - scalability
  - reliable vs. unreliable message transfer
  - "complexity at network edge"