

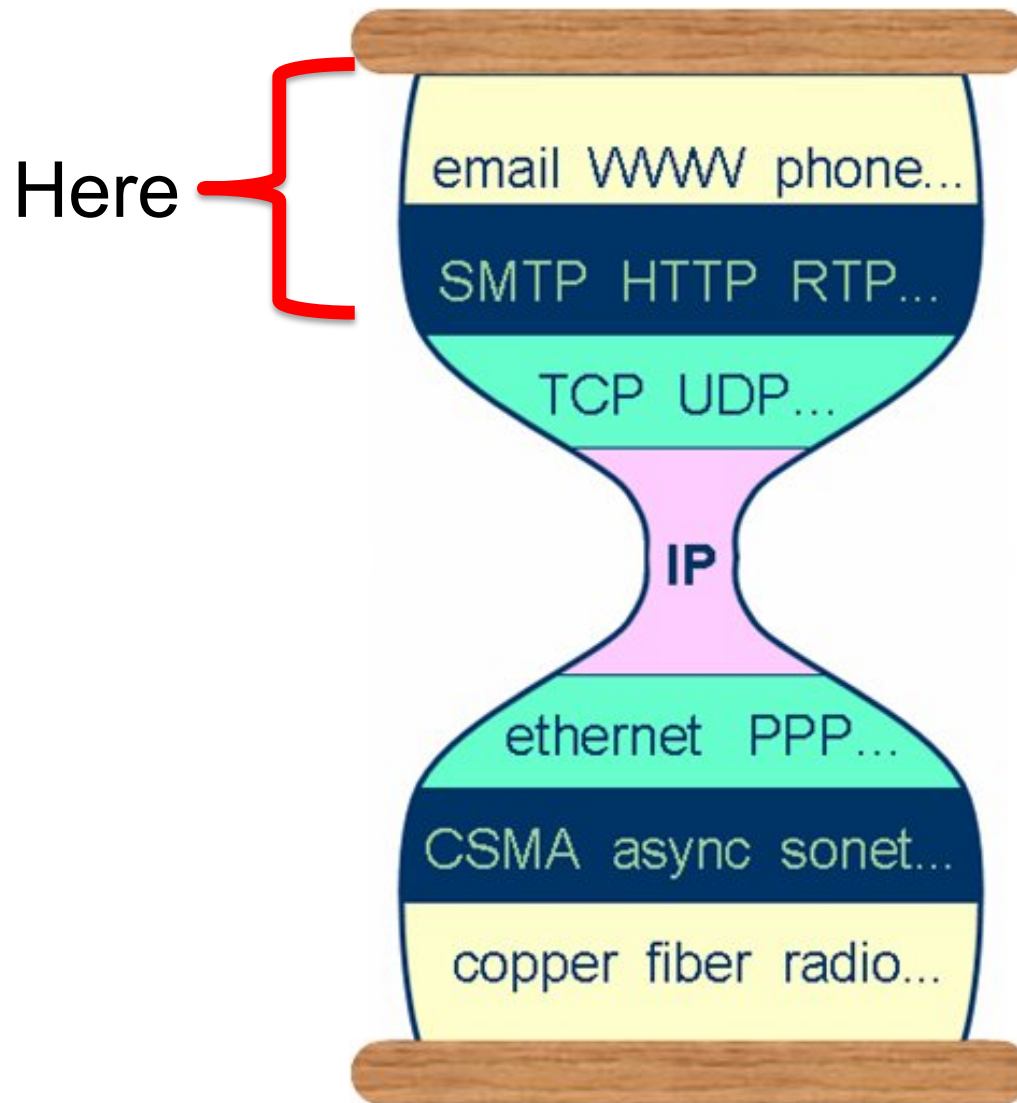
Chapter II: Application Layer

UG3 Computer Communications & Networks
(COMN)

MAHESH MARINA
mahesh@ed.ac.uk

Slides copyright of Kurose and Ross

Internet hourglass



Some network apps

- e-mail
- web
- text messaging
- remote login
- 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
- ...
- ...

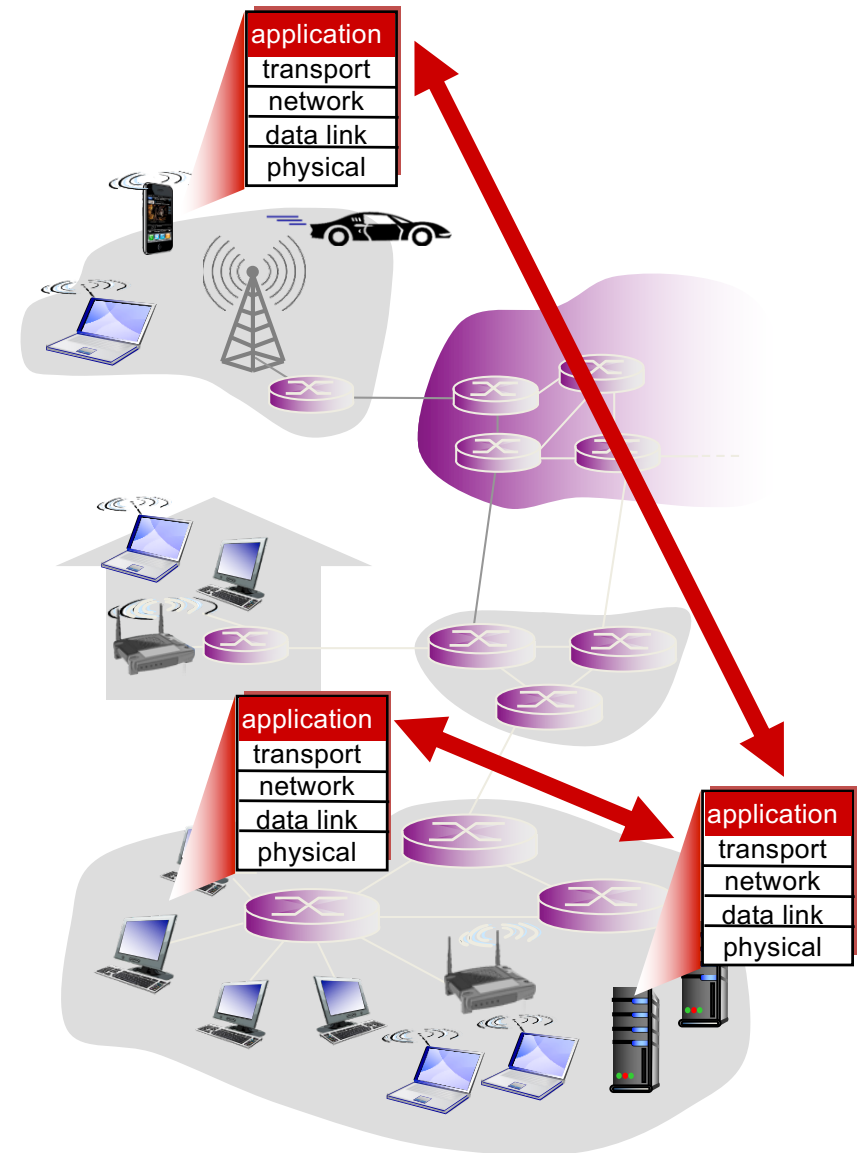
Creating a network app

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
- applications on end systems allows for rapid app development, propagation

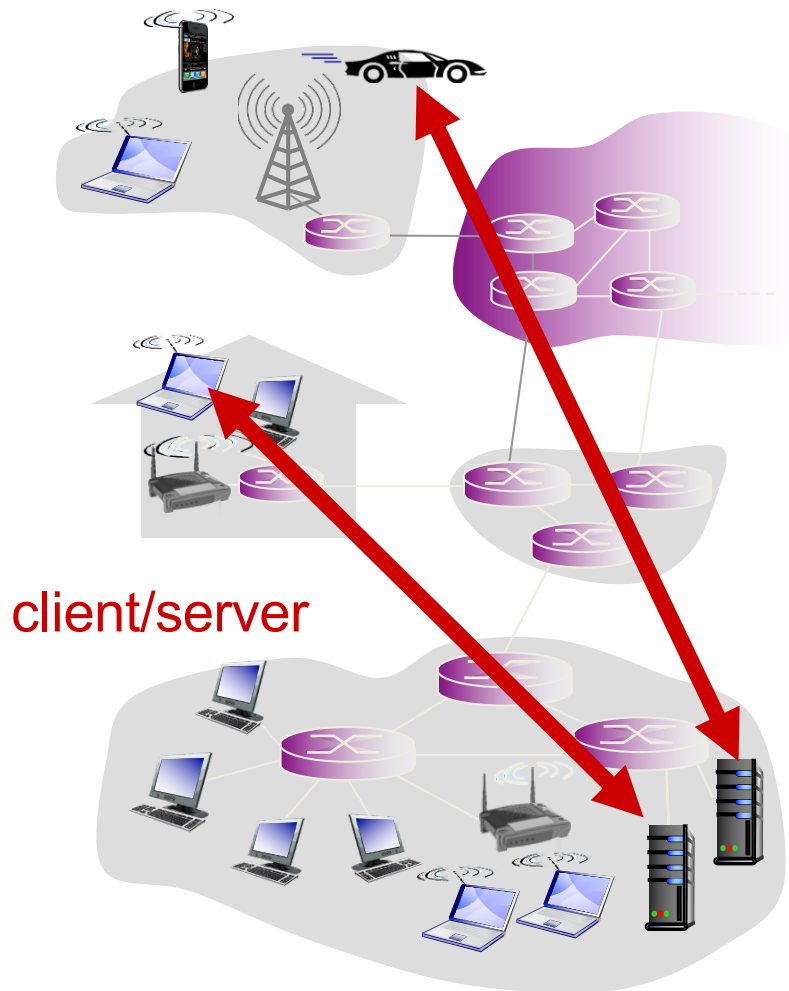


Application architectures

possible structure of applications:

- client-server
- peer-to-peer (P2P)

Client-server architecture



server:

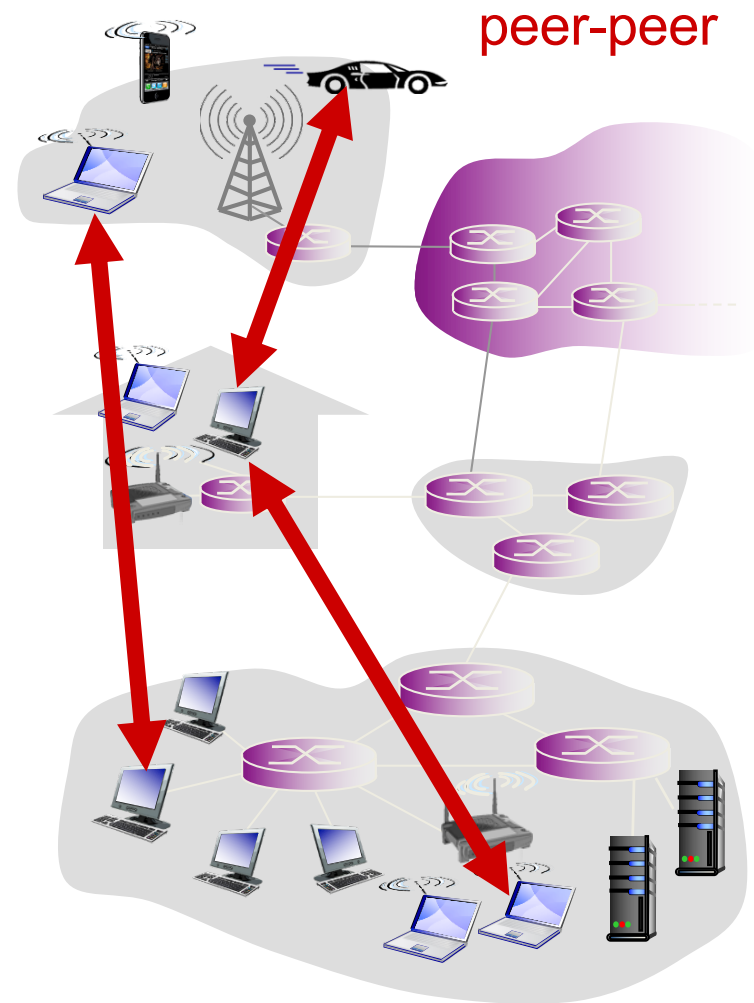
- always-on host
- permanent IP address
- data centers for scaling

clients:

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

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, as well as new service demands
- peers are intermittently connected and change IP addresses
 - complex management



Processes communicating

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

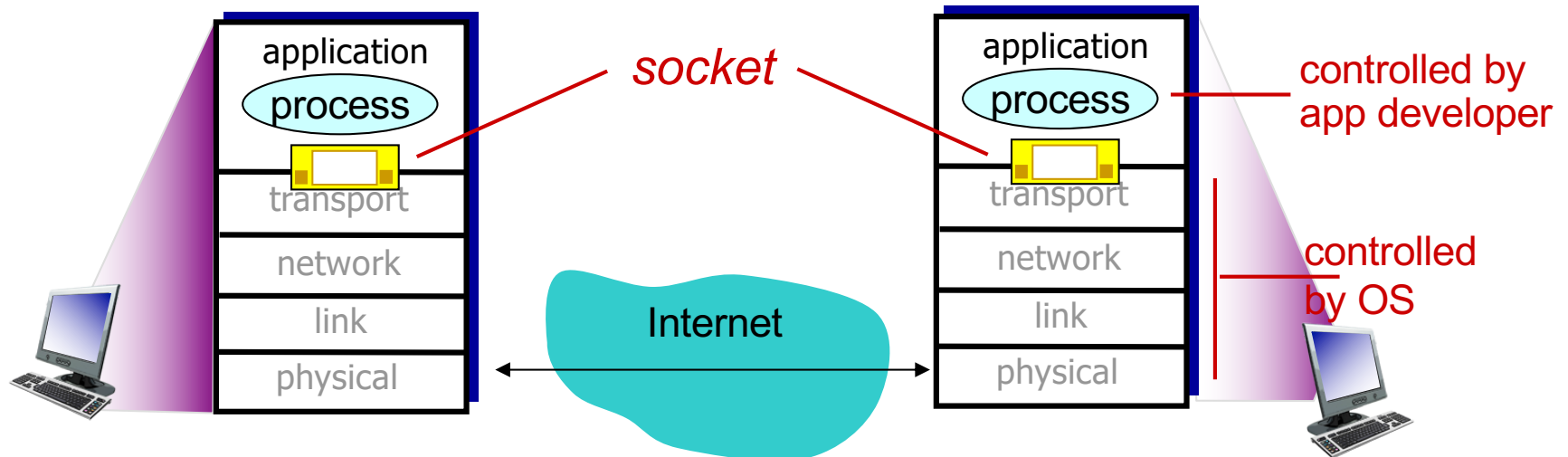
client process: process that initiates communication

server process: process that waits to be contacted

- ❖ aside: applications with P2P architectures have client processes & server processes

Sockets

- 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



Addressing processes

- 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 www.inf.ed.ac.uk web server:
 - **IP address:** 129.215.33.176
 - **port number:** 80
- more shortly...

App-layer protocol defines

- **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**:
 - defined in RFCs
 - allows for interoperability
 - e.g., HTTP, SMTP
 - **proprietary protocols**:
 - e.g., Skype

What transport service does an app need?

data integrity

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

timing

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

throughput

- ❖ 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

- ❖ encryption, data integrity,
...

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

Internet transport protocols services

TCP service:

- *reliable transport* between sending and receiving process
- *flow control*: sender won't overwhelm receiver
- *congestion control*: throttle sender when network overloaded
- *does not provide*: timing, minimum throughput guarantee, security
- *connection-oriented*: setup required between client and server processes

UDP service:

- *unreliable data transfer* between sending and receiving process
- *does not provide*: reliability, flow control, congestion control, timing, throughput guarantee, security, or connection setup,

Q: why bother? Why is there a UDP?

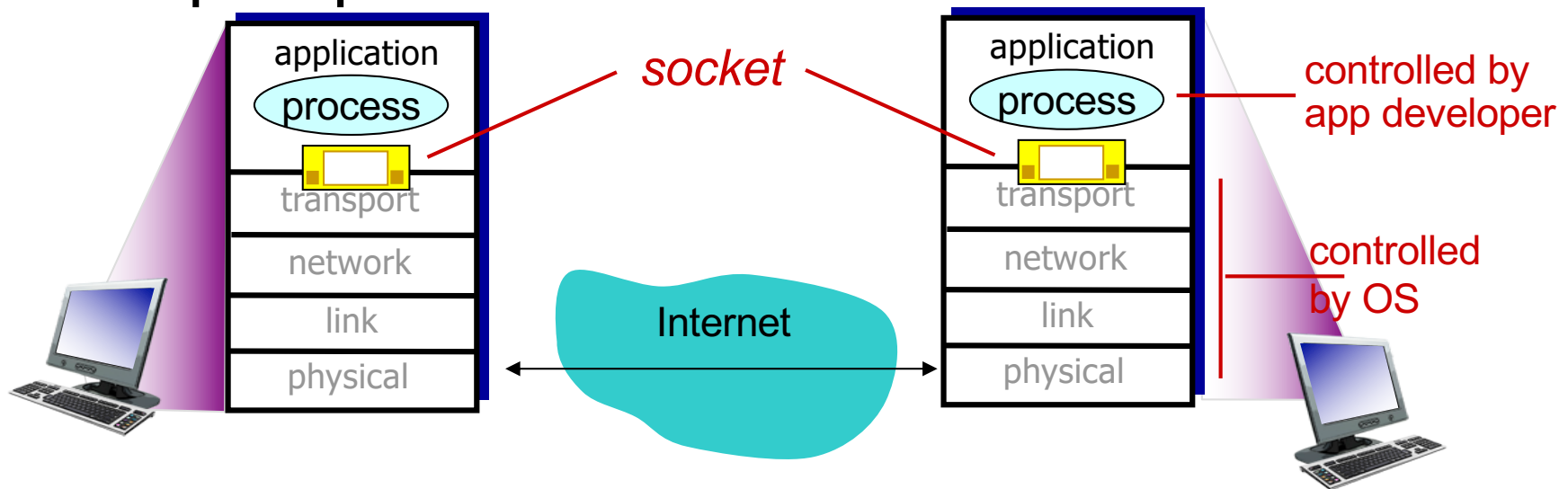
Internet apps: application, transport protocols

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

Socket programming

goal: learn how to build network applications that communicate using sockets

socket: door between application process and end-to-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 the data to the server.
2. The server receives the data and converts characters to uppercase.
3. The server sends the modified data to the client.
4. The client receives the modified data and displays the 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
- rcvr 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 `server IP`)

Client

create socket at port = x:

`serverSocket =
DatagramSocket(x)`

read datagram from
`serverSocket`

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

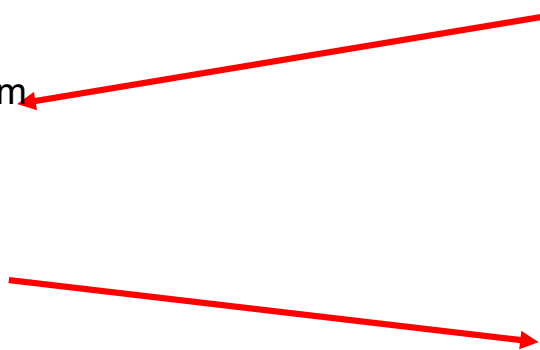
create socket:

`clientSocket =
DatagramSocket()`

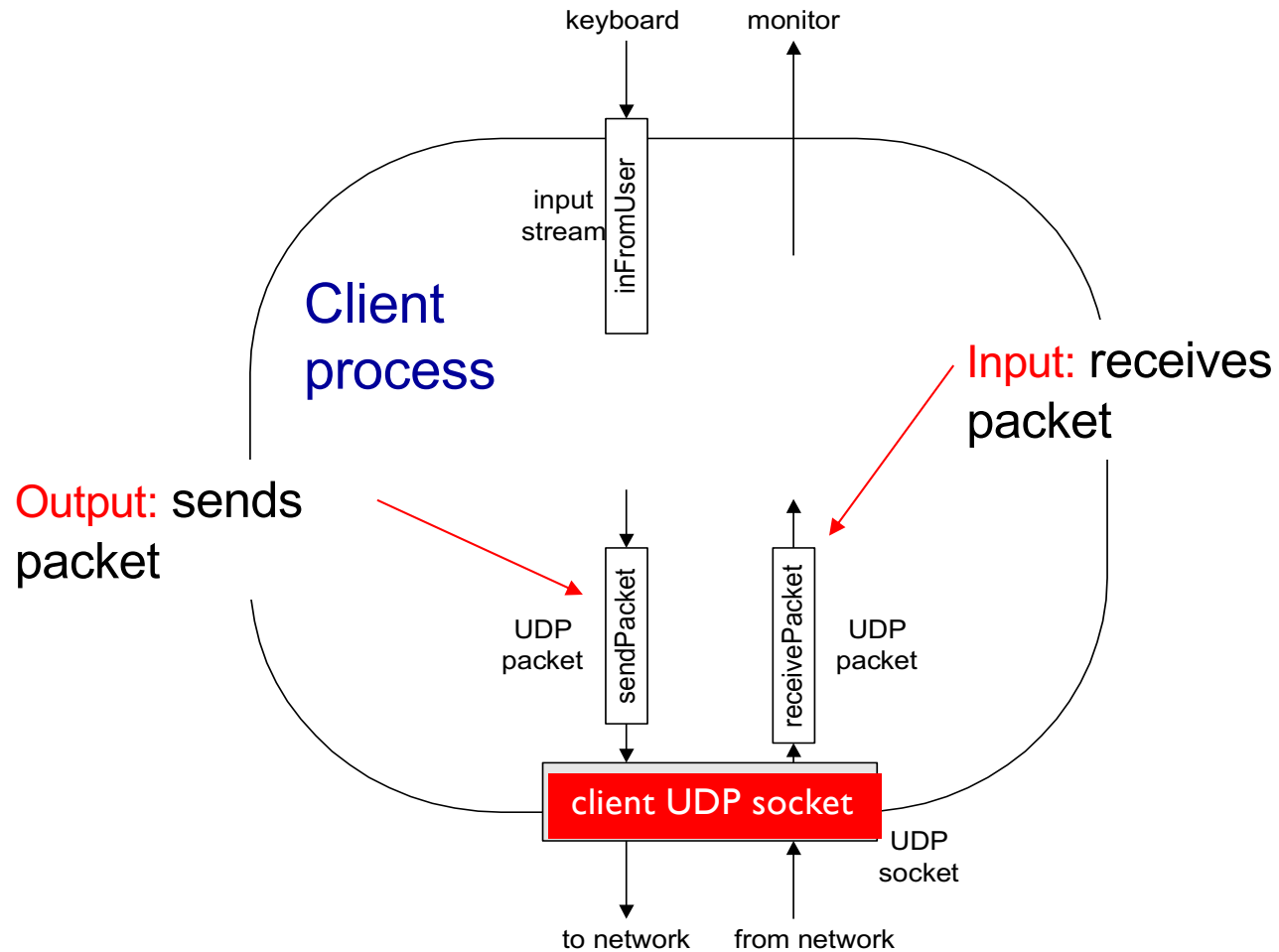
Create datagram with server IP and
port=x; send datagram via
`clientSocket`

read datagram from
`clientSocket`

close
`clientSocket`



Example: Java client (UDP)



Example: Java client (UDP)

```
import java.io.*;
import java.net.*;
```

```
class UDPClient {
    public static void main(String args[]) throws Exception
    {
```

create
input stream

```
        BufferedReader inFromUser =
            new BufferedReader(new InputStreamReader(System.in));
```

create
client socket

```
        DatagramSocket clientSocket = new DatagramSocket();
```

translate
hostname to IP
address using DNS

```
        InetAddress IPAddress = InetAddress.getByName("hostname");

        byte[] sendData = new byte[1024];
        byte[] receiveData = new byte[1024];
```

```
        String sentence = inFromUser.readLine();
```

```
        sendData = sentence.getBytes();
```

Example: Java client (UDP), cont.

```
create datagram with  
  data-to-send,  
  length, IP addr, port } DatagramPacket sendPacket =  
                          → new DatagramPacket(sendData, sendData.length, IPAddress, 9876);  
  
send datagram  
  to server } clientSocket.send(sendPacket);  
  
read datagram  
  from server } DatagramPacket receivePacket =  
                → new DatagramPacket(receiveData, receiveData.length);  
  
                clientSocket.receive(receivePacket);  
  
                String modifiedSentence =  
                  new String(receivePacket.getData());  
  
                System.out.println("FROM SERVER:" + modifiedSentence);  
                clientSocket.close();  
                }  
            }
```

Example: Java server (UDP)

```
import java.io.*;
import java.net.*;

class UDPServer {
    public static void main(String args[]) throws Exception
    {
        create datagram socket at port 9876 → DatagramSocket serverSocket = new DatagramSocket(9876);

        byte[] receiveData = new byte[1024];
        byte[] sendData = new byte[1024];

        while(true)
        {
            create space for received datagram → DatagramPacket receivePacket =
                new DatagramPacket(receiveData, receiveData.length);
            receive datagram → serverSocket.receive(receivePacket);
        }
    }
}
```

Example: Java server (UDP), cont

```
String sentence = new String(receivePacket.getData());
```

get IP addr
port #, of
sender

```
InetAddress IPAddress = receivePacket.getAddress();
```

```
int port = receivePacket.getPort();
```

```
String capitalizedSentence = sentence.toUpperCase();
```

```
sendData = capitalizedSentence.getBytes();
```

create datagram
to send to client

```
DatagramPacket sendPacket =  
    new DatagramPacket(sendData, sendData.length, IPAddress,  
                        port);
```

write out
datagram
to socket

```
serverSocket.send(sendPacket);
```

```
}
```

```
}
```

```
}
```

end of while loop,
loop back and wait for
another datagram

Connectionless demultiplexing

❖ *recall*: created socket has host-local port #:

```
DatagramSocket mySocketI =  
new DatagramSocket(12534);
```

❖ *recall*: when creating datagram to send into UDP socket, must specify

- destination IP address
- destination port #

❖ when host receives UDP segment:

- checks destination port # in segment
- directs UDP segment to socket with that port #



IP datagrams with *same dest. port #*, but different source IP addresses and/or source port numbers will be directed to *same socket* at dest

Connectionless demux: example

DatagramSocket

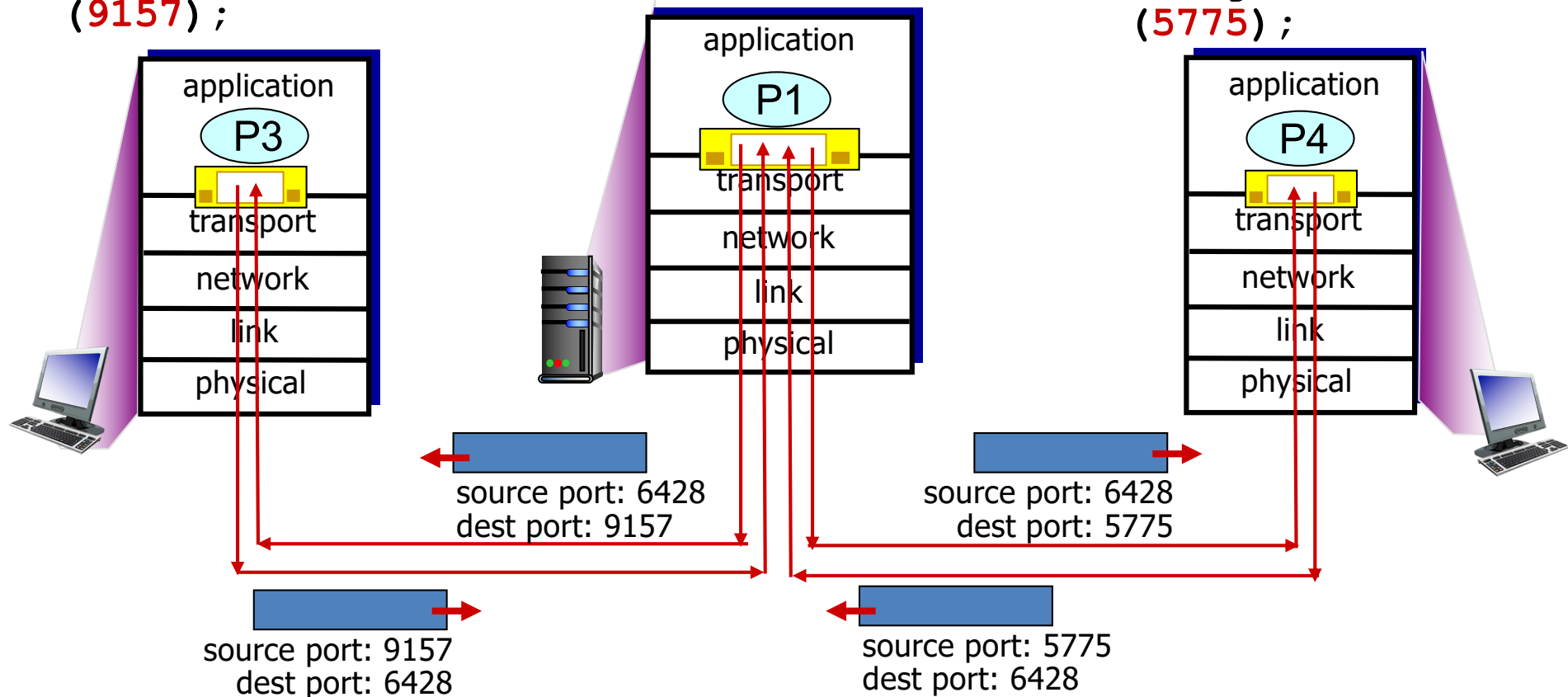
```
serverSocket = new
```

```
DatagramSocket
```

```
(6428);
```

```
DatagramSocket  
mySocket2 = new  
DatagramSocket  
(9157);
```

```
DatagramSocket  
mySocket1 = new  
DatagramSocket  
(5775);
```



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:

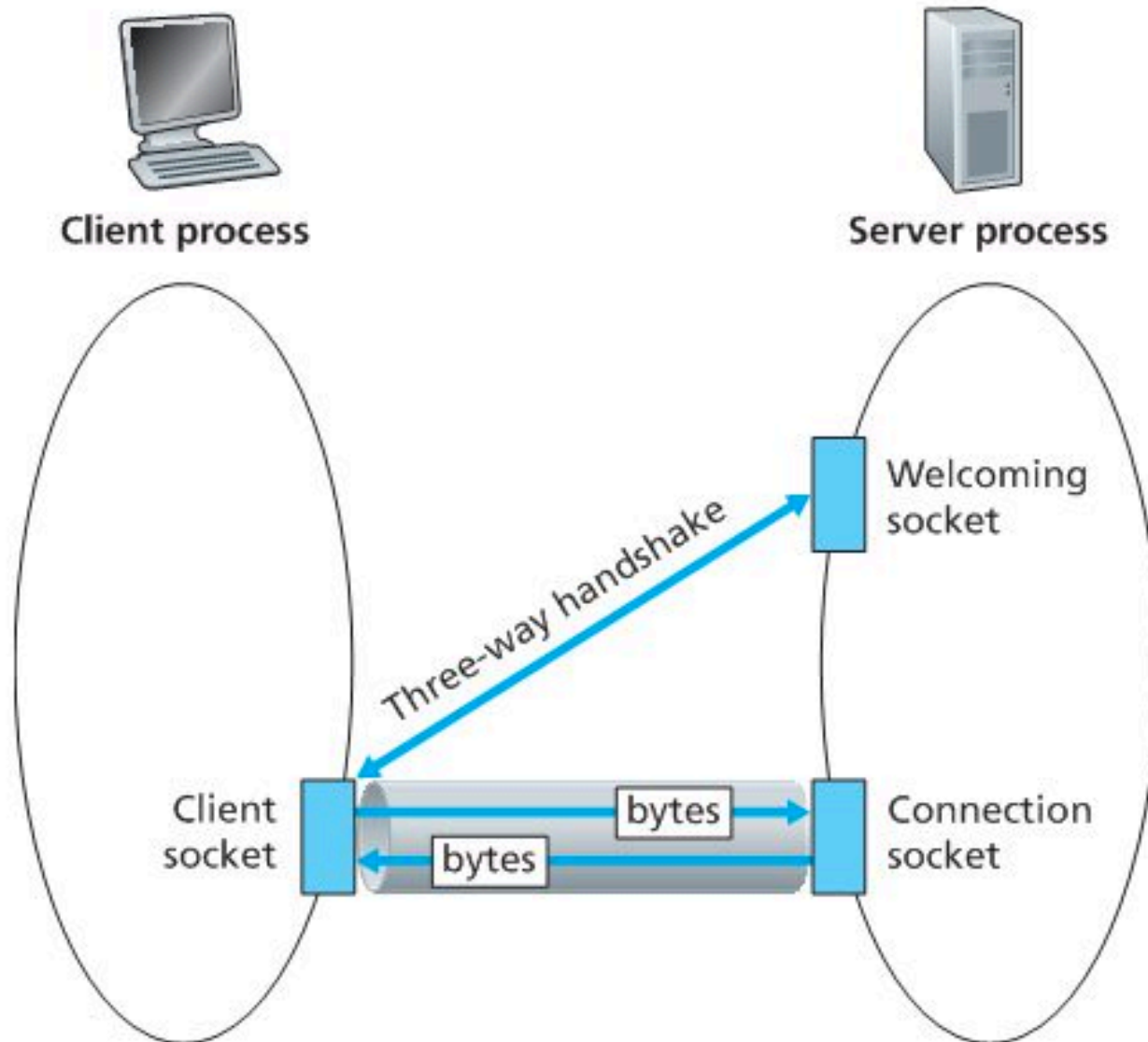
- Creating TCP socket, specifying IP address, port number of server process
- *when client creates socket:* client TCP establishes connection to server TCP

- 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
 - 4-tuple (clarified shortly) used to distinguish clients

application viewpoint:

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

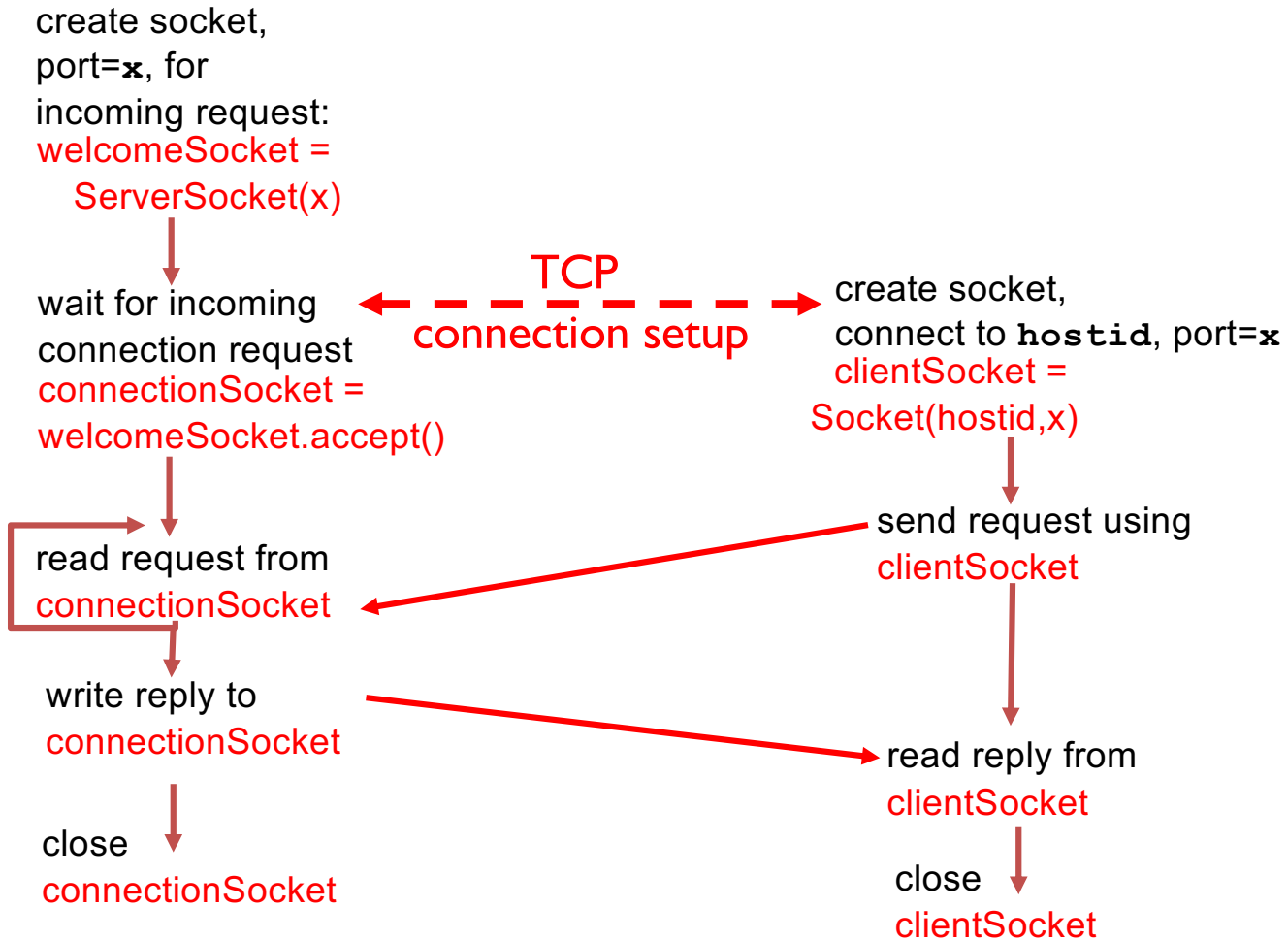
Illustration of TCP socket in client/server



Client/server socket interaction:TCP

Server (running on `hostid`)

Client



Example: Java client (TCP)

```
import java.io.*;
import java.net.*; ← This package defines Socket()
                    and ServerSocket() classes
class TCPClient {
```

```
    public static void main(String argv[]) throws Exception
    {
```

```
        String sentence;
        String modifiedSentence;
```

create
input stream →

```
        BufferedReader inFromUser =
            new BufferedReader(new InputStreamReader(System.in));
```

create
clientSocket object
of type Socket,
connect to server →

```
        Socket clientSocket = new Socket("hostname", 6789);
```

create
output stream
attached to socket →

```
        DataOutputStream outToServer =
            new DataOutputStream(clientSocket.getOutputStream());
```

server name,
e.g., www.ed.ac.uk

server port #

hostname, 6789

Example: Java client (TCP), cont.

```
        create  
        input stream attached to socket → BufferedReader inFromServer =  
                                         new BufferedReader(new  
                                         InputStreamReader(clientSocket.getInputStream()));  
  
                                         sentence = inFromUser.readLine();  
  
        send line to server → outToServer.writeBytes(sentence + '\n');  
  
        read line from server → modifiedSentence = inFromServer.readLine();  
  
                                         System.out.println("FROM SERVER: " + modifiedSentence);  
  
        close socket → clientSocket.close();  
  
        }  
    }
```

Example: Java server (TCP)

```
import java.io.*;
import java.net.*;
```

```
class TCPServer {
```

```
    public static void main(String argv[]) throws Exception
    {
```

```
        String clientSentence;
        String capitalizedSentence;
```

create
welcoming socket
at port 6789

```
        ServerSocket welcomeSocket = new ServerSocket(6789);
```

wait, on welcoming
socket accept() method
for client contact create,
new socket on return

```
        while(true) {
```

```
            Socket connectionSocket = welcomeSocket.accept();
```

create input
stream, attached
to socket

```
            BufferedReader inFromClient =
                new BufferedReader(new
                    InputStreamReader(connectionSocket.getInputStream()));
```


Example: Java server (TCP), cont

create output
stream, attached
to socket

→ `DataOutputStream outToClient =
new DataOutputStream(connectionSocket.getOutputStream());`

read in line
from socket

→ `clientSentence = inFromClient.readLine();`

`capitalizedSentence = clientSentence.toUpperCase() + '\n';`

write out line
to socket

→ `outToClient.writeBytes(capitalizedSentence);`

`connectionSocket.close();`

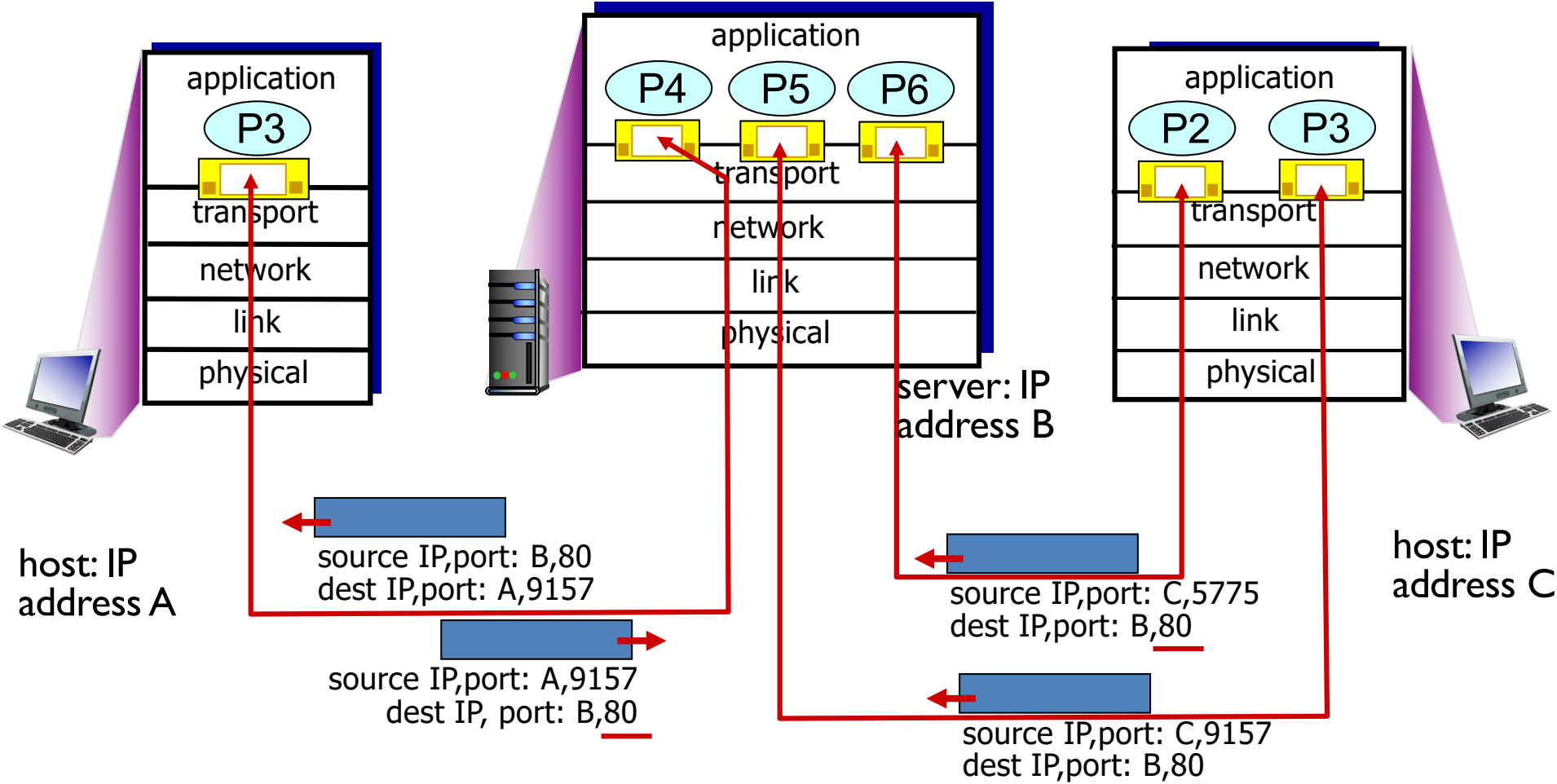
}
}
}

end of while loop,
loop back and wait for
another client connection

Connection-oriented demux

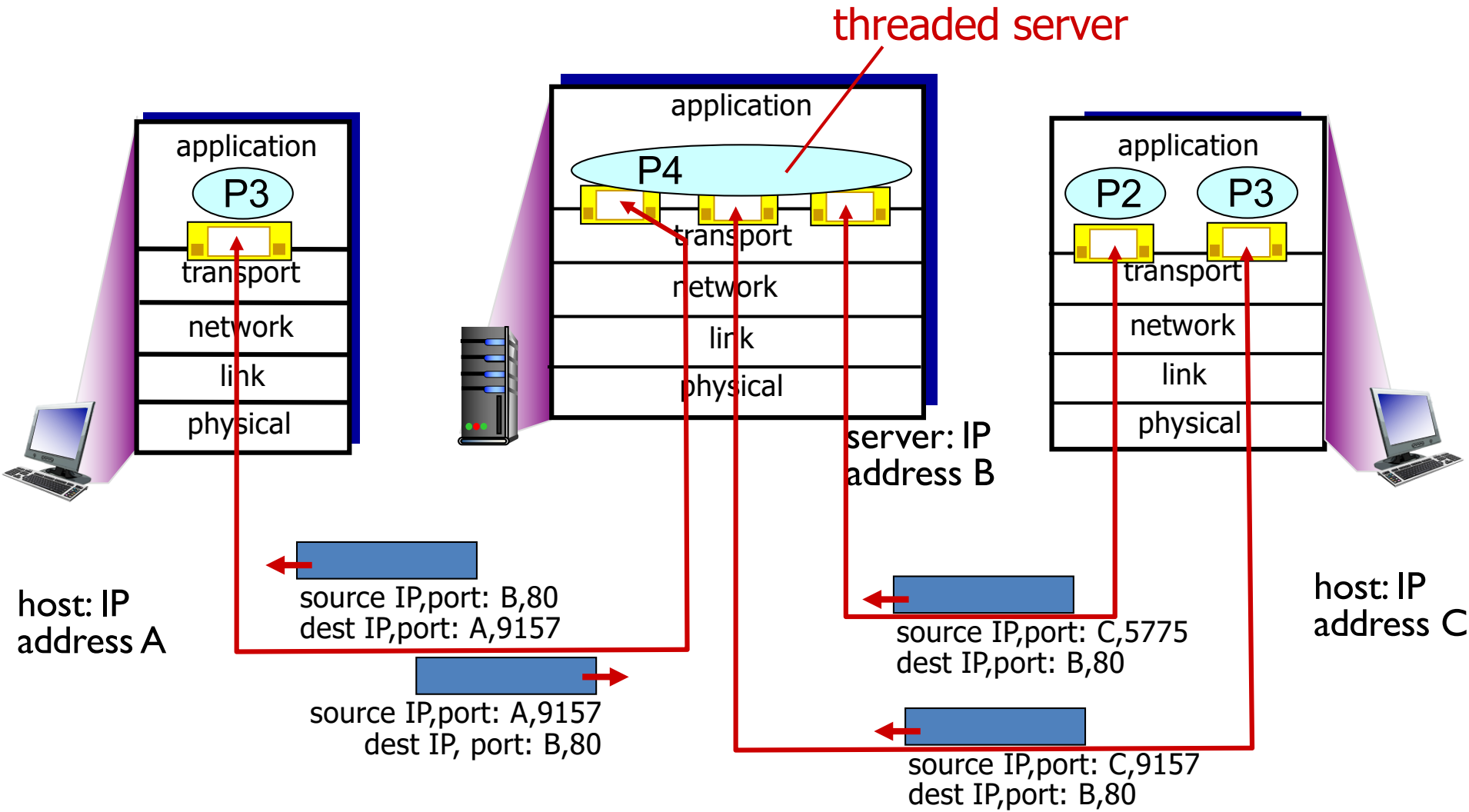
- ❖ TCP socket identified by 4-tuple:
 - source IP address
 - source port number
 - dest IP address
 - dest port number
- ❖ demux: receiver uses all four values to direct segment to appropriate socket
- ❖ server host may support many simultaneous TCP sockets:
 - each socket identified by its own 4-tuple
- ❖ web servers have different sockets for each connecting client
 - non-persistent HTTP (coming up shortly) will have different socket for each request

Connection-oriented demux: example



three segments, all destined to IP address: B,
dest port: 80 are demultiplexed to *different* sockets

Connection-oriented demux: example



Web and HTTP

First, a review...

- *web page* consists of *objects*
- object can be HTML file, JPEG image, Java applet, audio file,...
- web page consists of *base HTML-file* which includes *several referenced objects*
- each object is addressable by a *URL*, e.g.,

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

host name

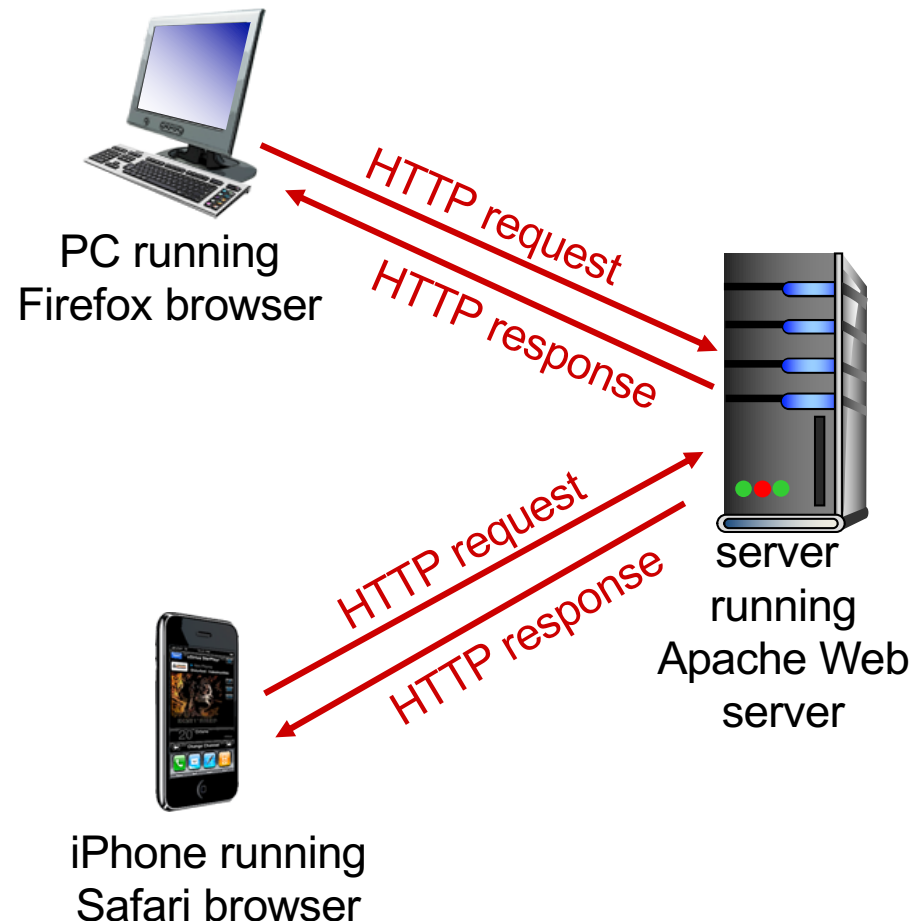
path name

```
view-source:https://www.bbc.co.uk
Secure | view-source:https://www.bbc.co.uk
1 <!DOCTYPE html>
2 <!--[if lte IE 9]>
3   <html lang="en-GB" class="no-js no-flexbox no-flexboxlegacy">
4 <![endif]-->
5 <!--[if gt IE 9]><!-->
6   <html lang="en-GB" class="no-js">
7 <!--![endif]-->
8 <head>
9   <meta http-equiv="X-UA-Compatible" content="IE=edge,chrome=1"/><script type="text/javascript">
10     var _sf_startpt = (new Date()).getTime();
11   </script><meta content="text/html; charset=UTF-8" http-equiv="Content-Type"><meta content="The best of the BBC, with the latest news and
headlines, weather, TV & radio highlights and much more from across the whole of BBC Online" name="description"><meta content="BBC, Briti
Broadcasting Corporation, BBCi, News, Sport, iPlayer, TV, Radio, Food, Music, Business, Arts, Bitesize, Lifestyle, Entertainment, Headlines"
name="keywords"><meta property="og:title" content="BBC - Home"><meta property="og:type" content="website"><meta property="og:description"
content="The best of the BBC, with the latest news and sport headlines, weather, TV & radio highlights and much more from across the whol
BBC Online"><meta property="og:site_name" content="BBC Homepage"><meta property="og:locale" content="en_GB"><meta property="article:author"
content="https://www.facebook.com/bbc"><meta property="og:article:section" content="Home"><meta property="og:url" content="http://www.bbc.co.
"><meta property="og:image" content="//homepage.files.bbc.co.uk/s/homepage-v5/2563/images/bbc_homepage.png"><meta name="twitter:card"
content="summary_large_image"><meta name="twitter:site" content="@bbccouk"><meta name="twitter:title" content="BBC - Home"><meta
name="twitter:description" content="The best of the BBC, with the latest news and sport headlines, weather, TV & radio highlights and muc
more from across the whole of BBC Online"><meta name="twitter:creator" content="@bbccouk"><meta name="twitter:image:src"
content="//homepage.files.bbc.co.uk/s/homepage-v5/2563/images/bbc_homepage.png"><meta name="twitter:image:alt" content="BBC Homepage"><meta
name="twitter:domain" content="www.bbc.co.uk"><link rel="canonical" href="https://www.bbc.co.uk" /><script type="text/javascript">(function (
{window.bbcirection = { geo: true }})());</script><!-- Nav Env: live -->
12 <!-- Analytics Web Module: 83 -->
13 <!-- NavID Web Module: 0.2.0-143 -->
14 <!-- Searchbox Web Module: 133 -->
15 <!-- Promo Web Module: 0.0.0-239.4080e99 -->
16 <meta name="viewport" content="width=device-width, initial-scale=1.0"><meta property="fb:admins" content="100004154058350"><link rel="stylesh
href="https://nav.files.bbc.co.uk/orbit/1.0.0-519.0b4da2b/css/orb-ltr.min.css"><!--[if (lt IE 9) & (!IEMobile)]>
17 <link rel="stylesheet" href="https://nav.files.bbc.co.uk/orbit/1.0.0-519.0b4da2b/css/orb-ie-ltr.min.css">
18 <![endif]--><script type="text/javascript">/*!CDATA[*]
19   window.orb = {
20     lang: 'en',
21     bbcBaseUrl: 'http://www.bbc.co.uk',
22     staticHost: 'https://nav.files.bbc.co.uk/orbit/1.0.0-519.0b4da2b',
23     figUrl: 'https://fig.bbc.co.uk/frameworks/fig/2/fig.js',
24     partialCookieOvenUrl: 'https://cookie-oven.api.bbc'
25   };
26
27   document.documentElement.className += (document.documentElement.className? ' ' : '') + 'orb-js';
28   window.orb.worldwideNavlinks = '<li class="orb-nav-home"><a href="http://www.bbc.com/">Home</a></li><li class="orb-nav-newsdotcom"><a
href="http://www.bbc.com/news">News</a></li><li class="orb-nav-sport"><a href="http://www.bbc.com/sport/">Sport</a></li><li class="orb-nav-
weather"><a href="http://www.bbc.com/weather/">Weather</a></li><li class="orb-nav-shop"><a href="http://shop.bbc.com/">Shop</a></li><li
class="orb-nav-earthdotcom"><a href="http://www.bbc.com/earth/">Earth</a></li><li class="orb-nav-travel-dotcom"><a
BUSINESS FOOTBALL US & CANADA
```

HTTP overview

HTTP: hypertext transfer protocol

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



HTTP overview (continued)

uses TCP:

- client initiates TCP connection (creates socket) to server, port 80
- server accepts TCP connection from client
- HTTP messages (application-layer protocol messages) exchanged between browser (HTTP client) and Web server (HTTP server)
- TCP connection closed

HTTP is “stateless”

- server maintains no information about past client requests

aside

protocols that maintain “state” are complex!

- ❖ past history (state) must be maintained
- ❖ if server/client crashes, their views of “state” may be inconsistent, must be reconciled

HTTP connections

non-persistent HTTP

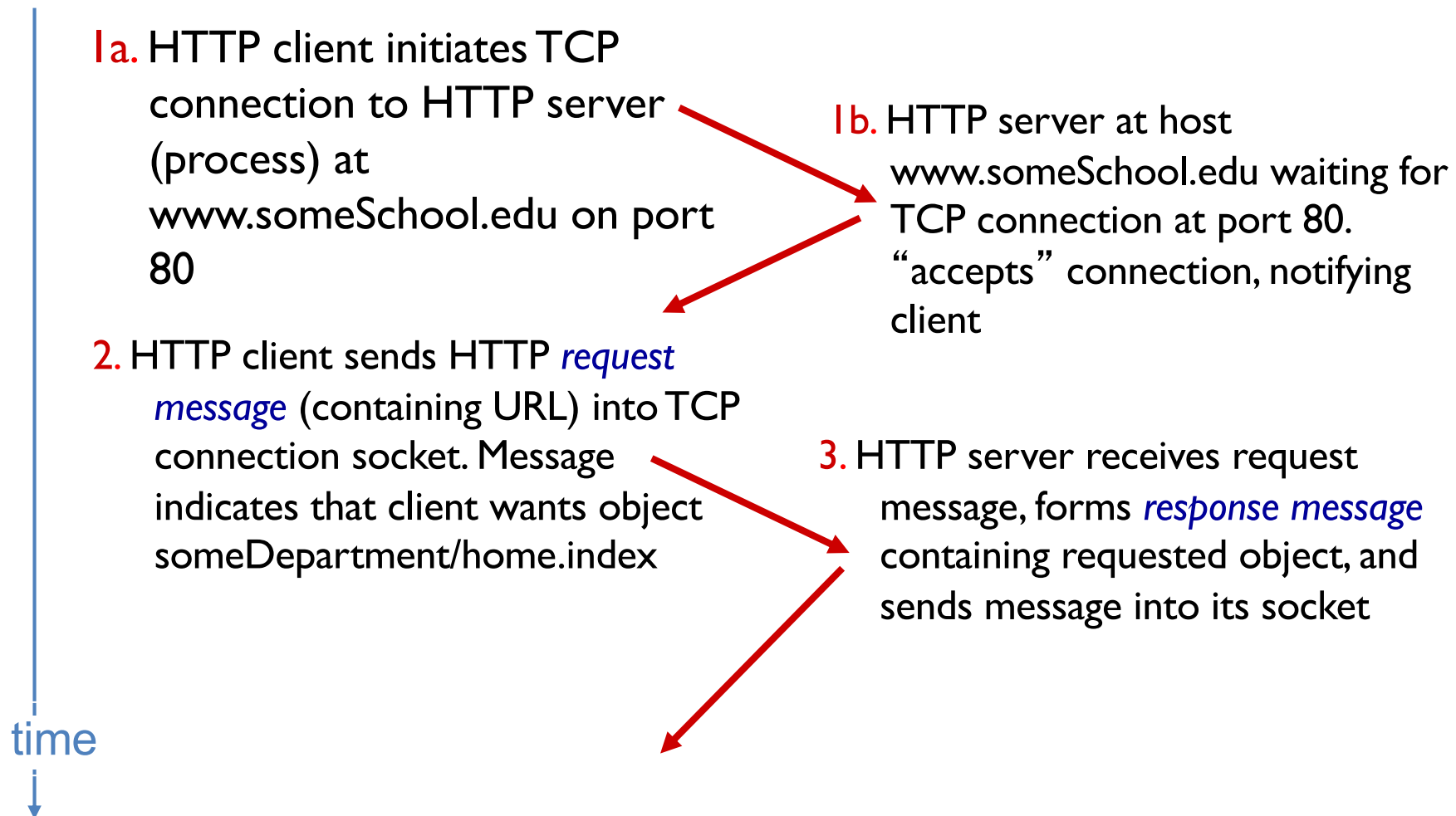
- at most one object sent over TCP connection
 - connection then closed
- downloading multiple objects requires multiple connections

persistent HTTP

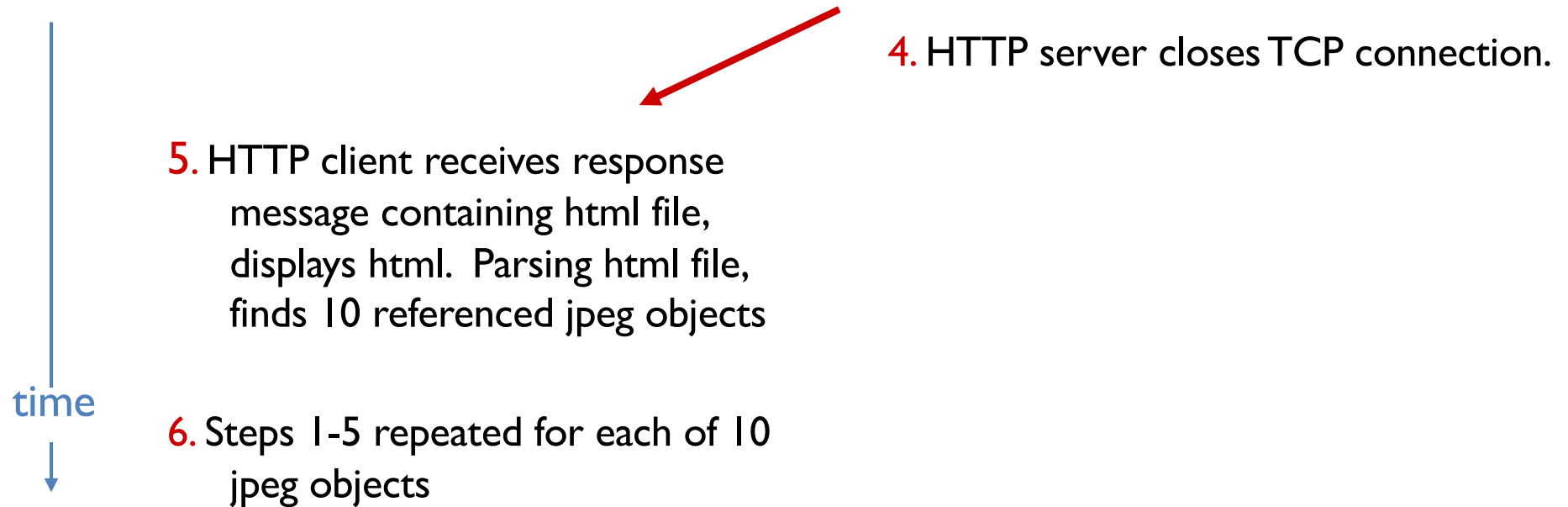
- multiple objects can be sent over single TCP connection between client and server

Non-persistent HTTP

suppose user enters URL: `www.someSchool.edu/someDepartment/home.index` (contains text, references to 10 jpeg images)



Non-persistent HTTP (cont.)

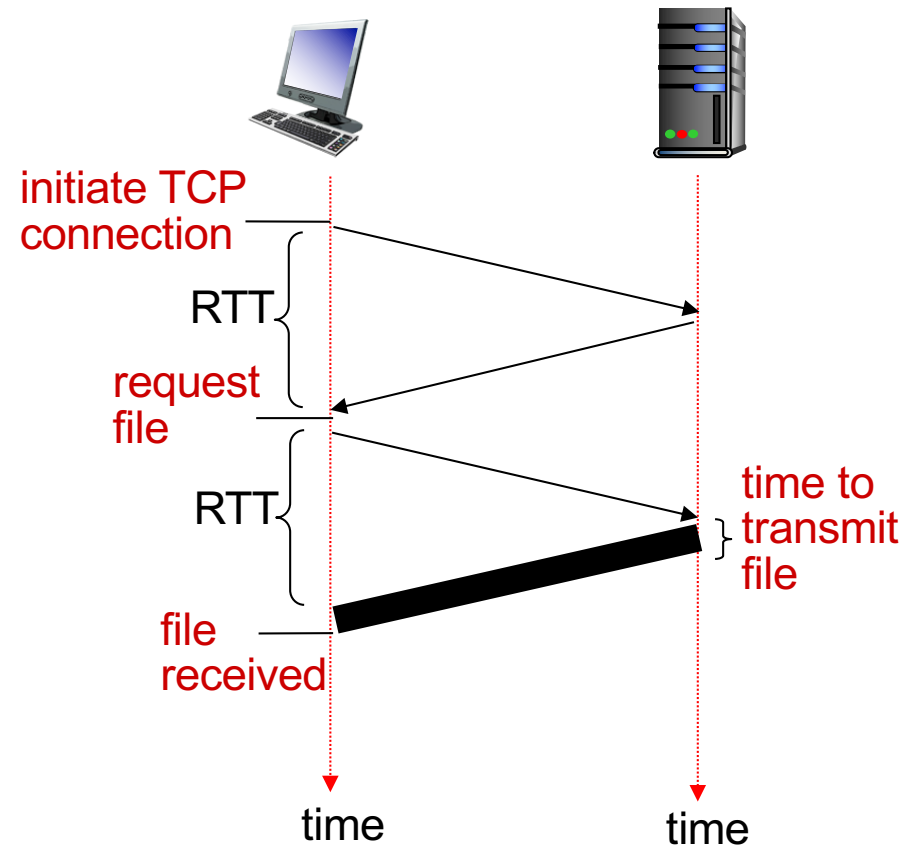


Non-persistent HTTP: response time

RTT (definition): time for a small packet to travel from client to server and back

HTTP response time:

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



Persistent HTTP

non-persistent HTTP issues:

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

persistent HTTP:

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

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

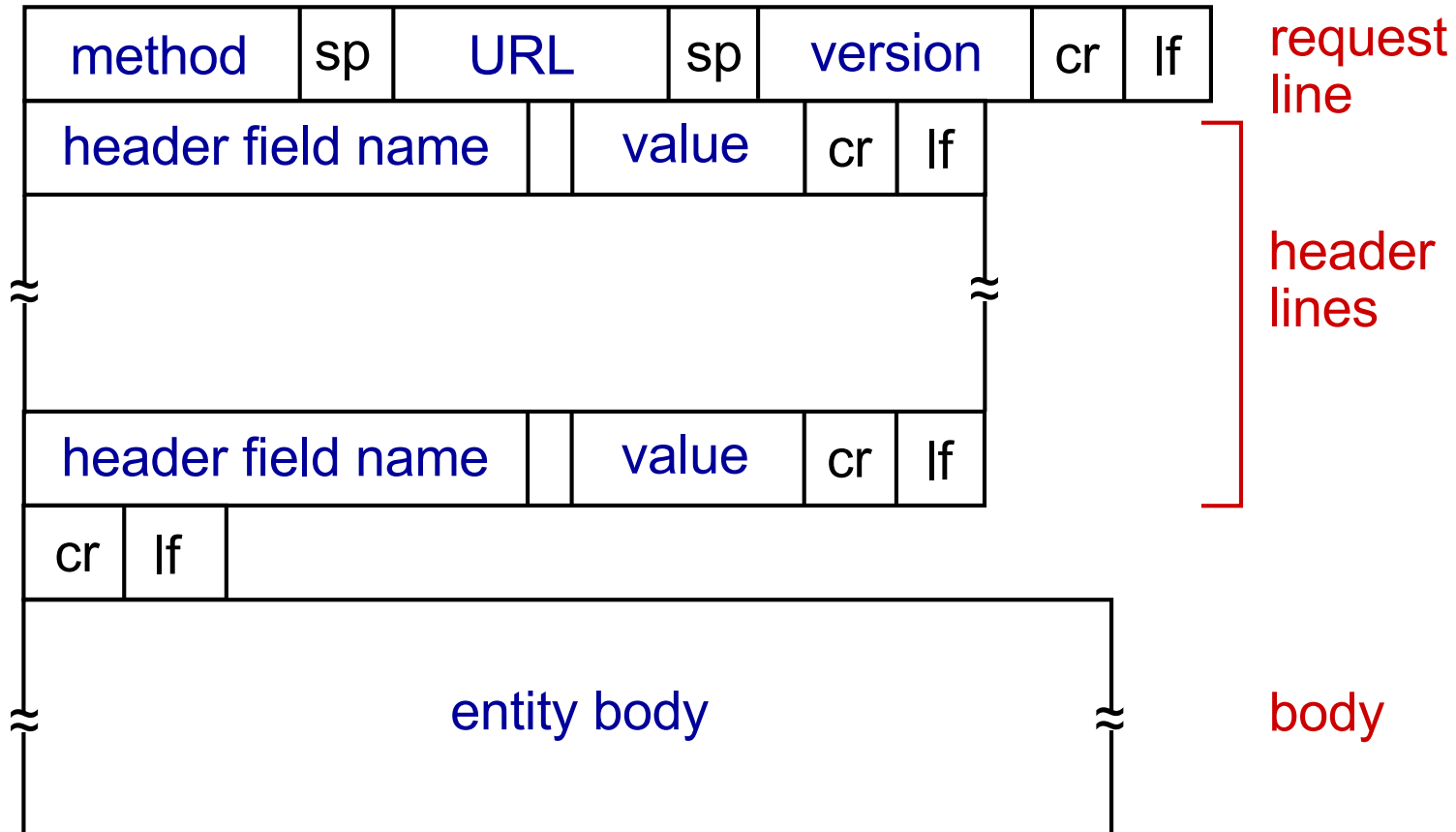
carriage return,
line feed at start
of line indicates
end of header lines

```
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

HTTP request message: general format



Uploading form input

POST method:

- web page often includes form input
- input is uploaded to server in entity body

URL method:

- uses GET method
- input is uploaded in URL field of request line:

`www.somesite.com/animalsearch?monkeys&banana`

Method types

HTTP/1.0:

- GET
- POST
- HEAD
 - asks server to leave requested object out of response

HTTP/1.1:

- GET, POST, HEAD
- PUT
 - uploads file in entity body to path specified in URL field
- DELETE
 - deletes file specified in the URL field

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 ...
```

HTTP response status codes

- ❖ status code appears in 1st line in server-to-client response message.
- ❖ some sample codes:

200 OK

- request succeeded, requested object later in this msg

301 Moved Permanently

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

400 Bad Request

- request msg not understood by server

404 Not Found

- requested document not found on this server

505 HTTP Version Not Supported

User-server state: cookies

many Web sites use cookies

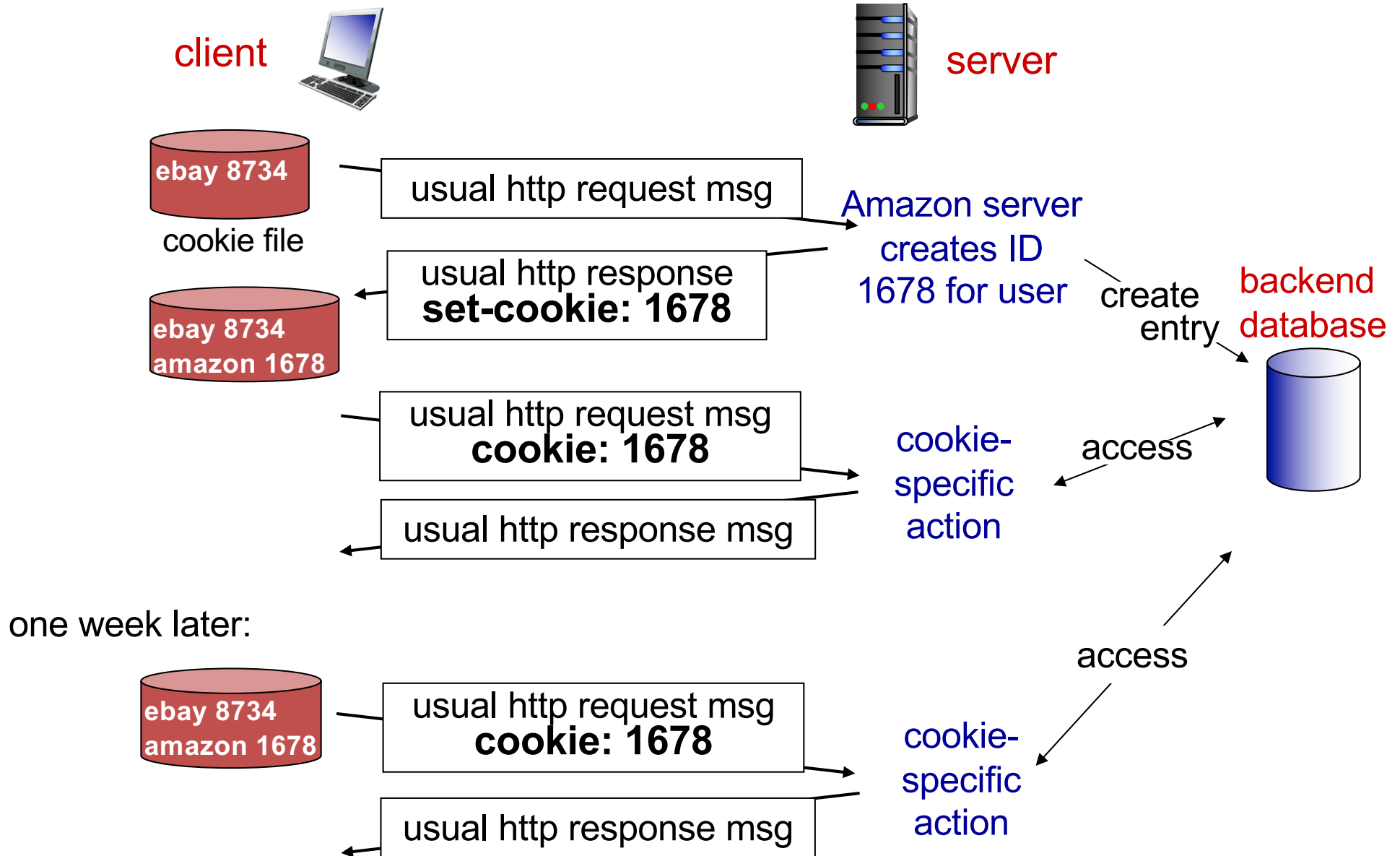
four components:

- 1) cookie header line of HTTP *response* message
- 2) cookie header line in next HTTP *request* message
- 3) cookie file kept on user's host, managed by user's browser
- 4) back-end database at Web site

example:

- Susan always access Internet from PC
- visits specific e-commerce site for first time
- when initial HTTP requests arrives at site, site creates:
 - unique ID
 - entry in backend database for ID

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

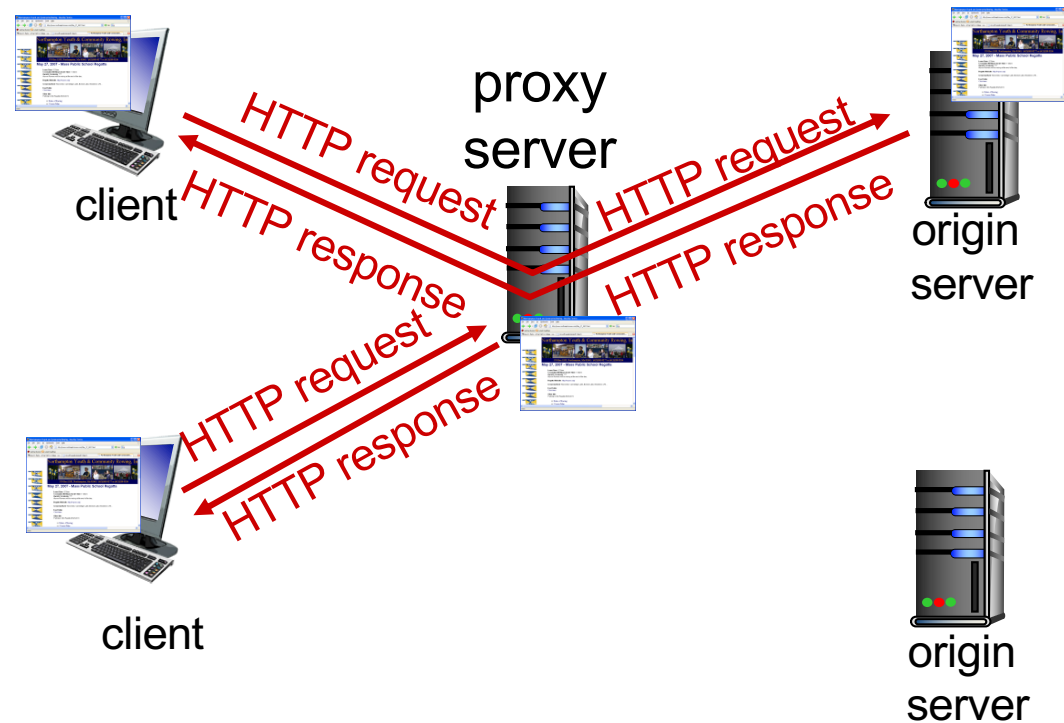
how to keep “state”:

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

Web caches (proxy server)

goal: satisfy client request without involving origin server

- 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
- Internet dense with caches: enables “poor” content providers to effectively deliver content (so too does P2P file sharing)

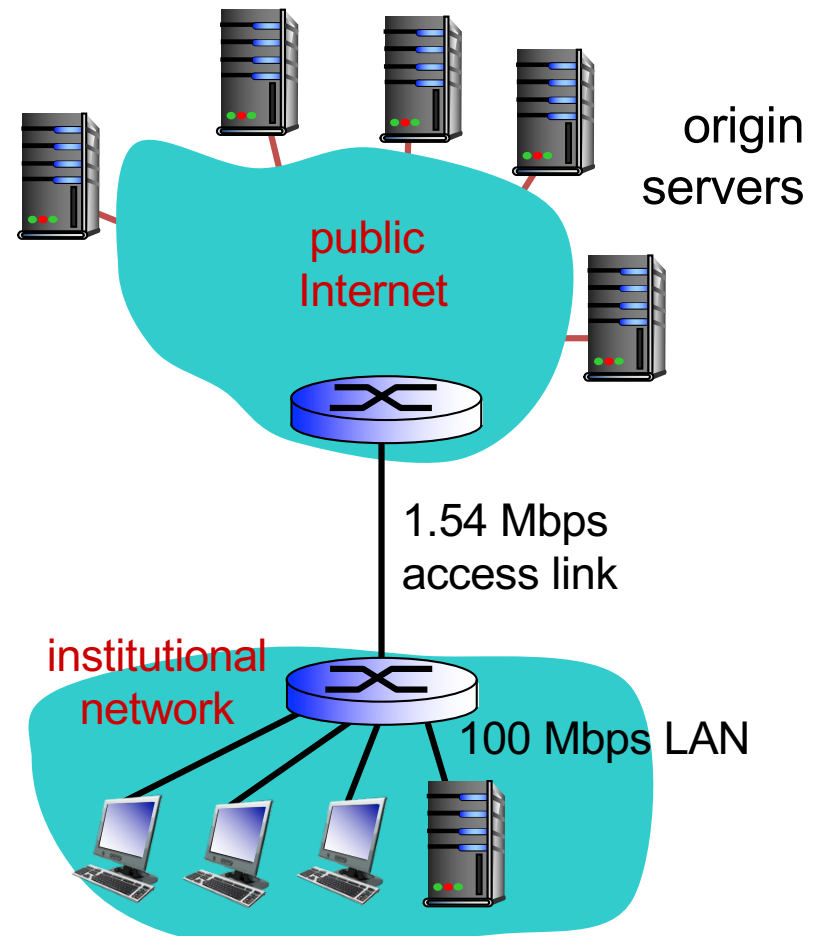
Caching example:

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: 1.5%
- ❖ access link utilization = **97%** *problem!*
- ❖ total delay = Internet delay + access delay + LAN delay
= 2 sec + minutes + usecs



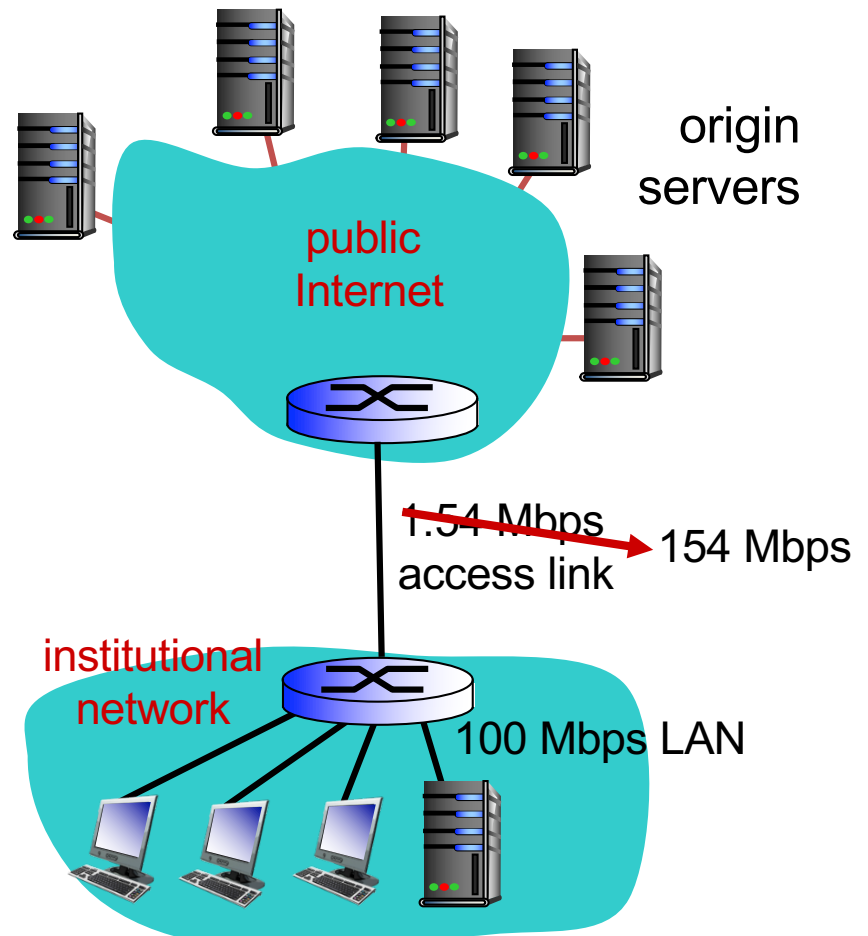
Caching example: fatter access link

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~~ → 154 Mbps

consequences:

- ❖ LAN utilization: 1.5%
- ❖ access link utilization = ~~97%~~ → 0.97%
- ❖ total delay = Internet delay + access delay + LAN delay
= 2 sec + ~~minutes~~ → usecs msec



Cost: increased access link speed (not cheap!)

Caching example: install local cache

assumptions:

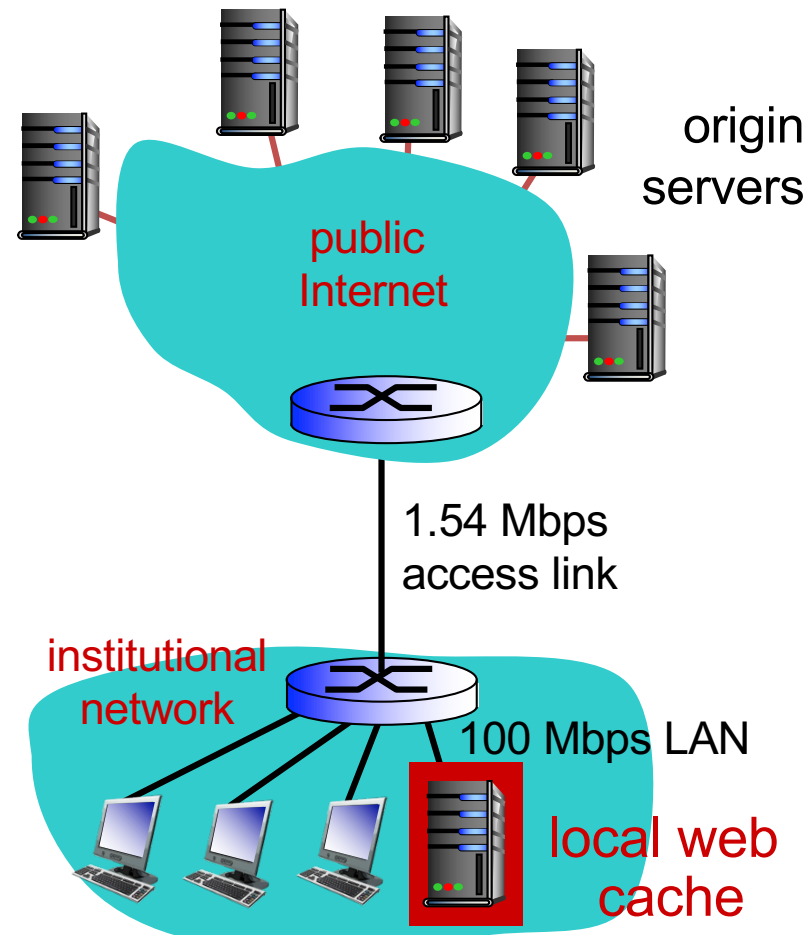
- ❖ avg object size: 100K bits
- ❖ avg request rate from browsers to origin servers: 15/sec
- ❖ avg data rate to browsers: 1.50 Mbps
- ❖ RTT from institutional router to any origin server: 2 sec
- ❖ access link rate: 1.54 Mbps

consequences:

- ❖ LAN utilization: 1.5%
- ❖ access link utilization = ?
- ❖ total delay = ?

How to compute link utilization, delay?

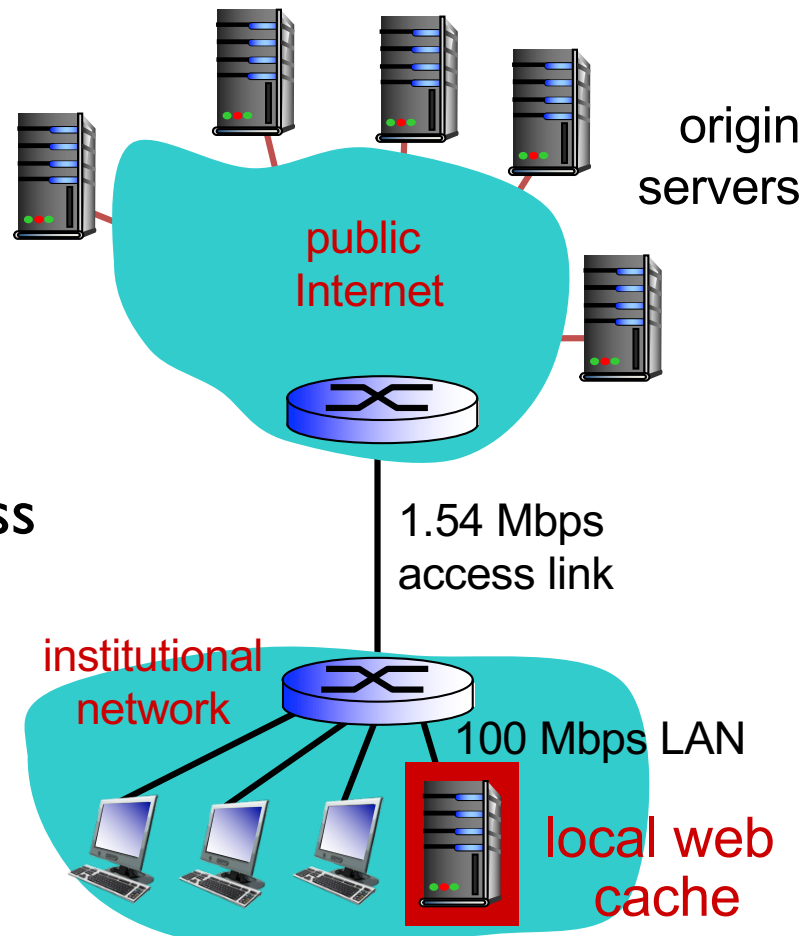
Cost: web cache (cheap!)



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 = $0.6 * 1.50 \text{ Mbps} = .9 \text{ Mbps}$
 - utilization = $0.9 / 1.54 = .58$
- ❖ 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!)

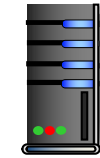


Conditional GET

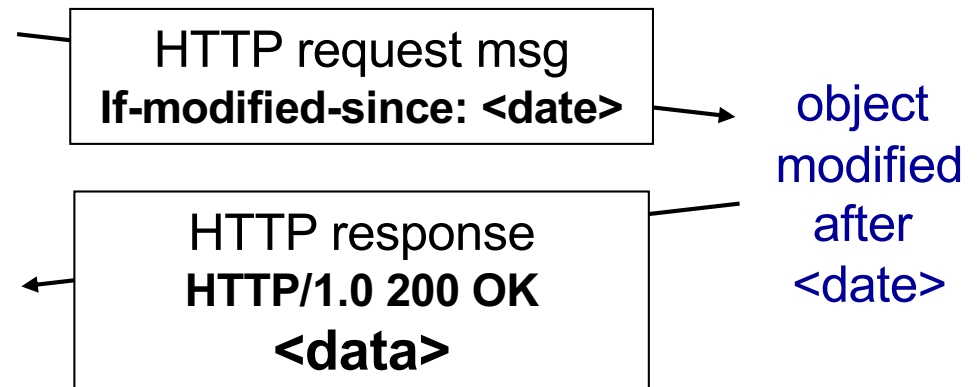
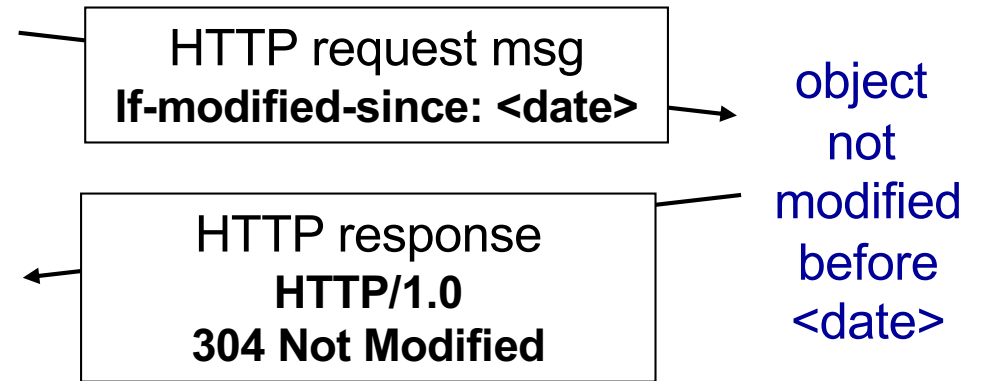
client



server



- **Goal:** don't send object if cache has up-to-date cached version
 - no object transmission delay
 - lower link utilization
- **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



DNS: domain name system

people: many identifiers:

- SSN, name, passport #

Internet hosts, routers:

- IP address (32 bit) - used for addressing datagrams
- “name”, e.g.,
www.yahoo.com - 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

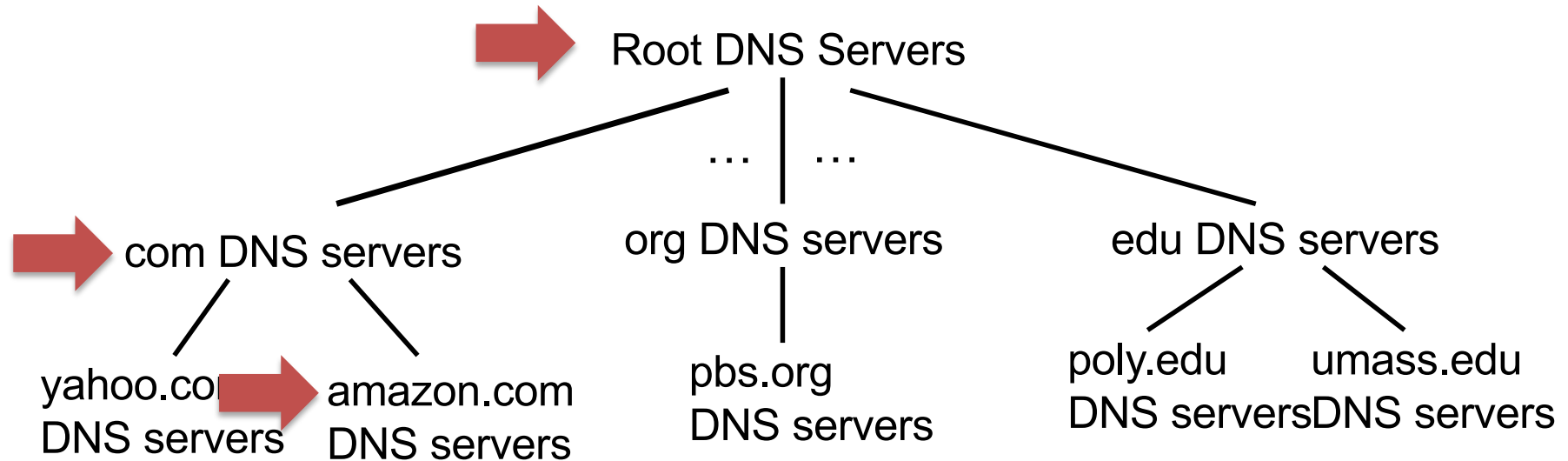
- hostname to IP address translation
- host aliasing
 - canonical, alias names
- mail server aliasing
- load distribution
 - replicated Web servers:
many IP addresses
correspond to one name

why not centralize DNS?

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

A: doesn't scale!

DNS: a distributed, hierarchical database

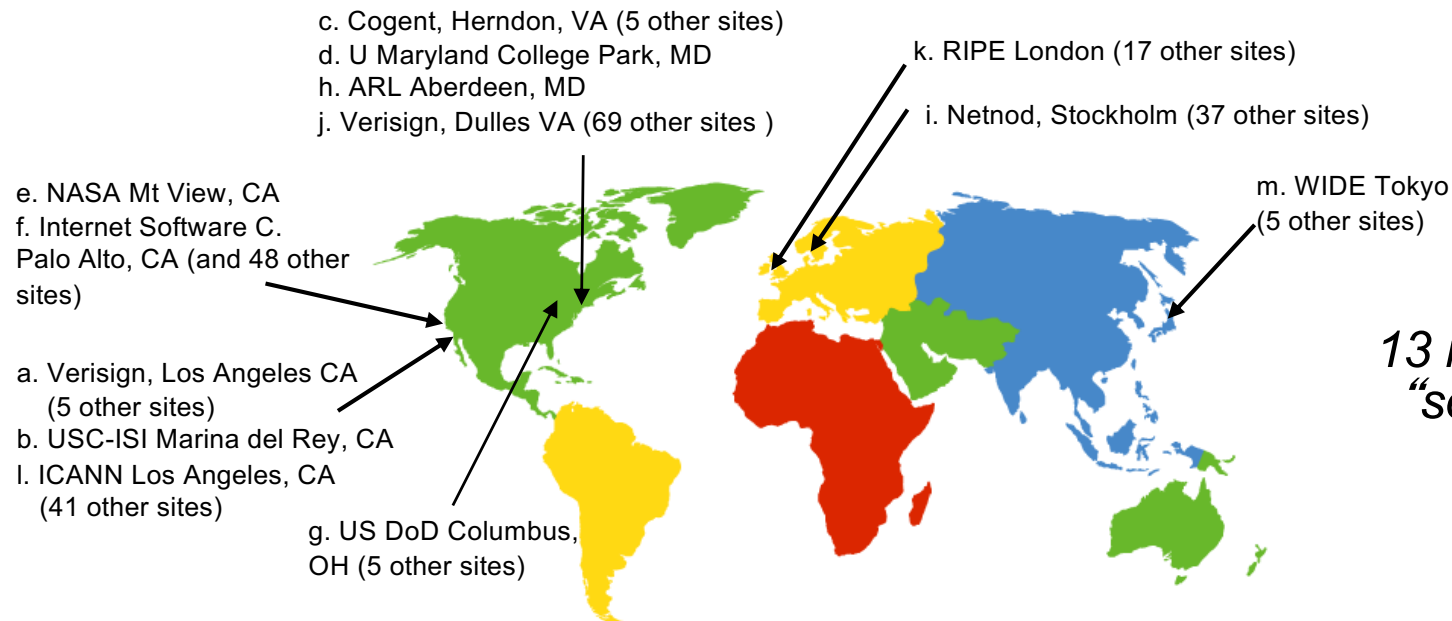


client wants IP for www.amazon.com; 1st approx:

- 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

- contacted by local name server that cannot resolve name
- root name server:
 - contacts TLD name server if name mapping not known
 - gets mapping
 - returns mapping to local name server



*13 root name
"servers" worldwide*

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.: uk, fr, ca, jp
- Network Solutions maintains servers for .com TLD
- Educause for .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 server

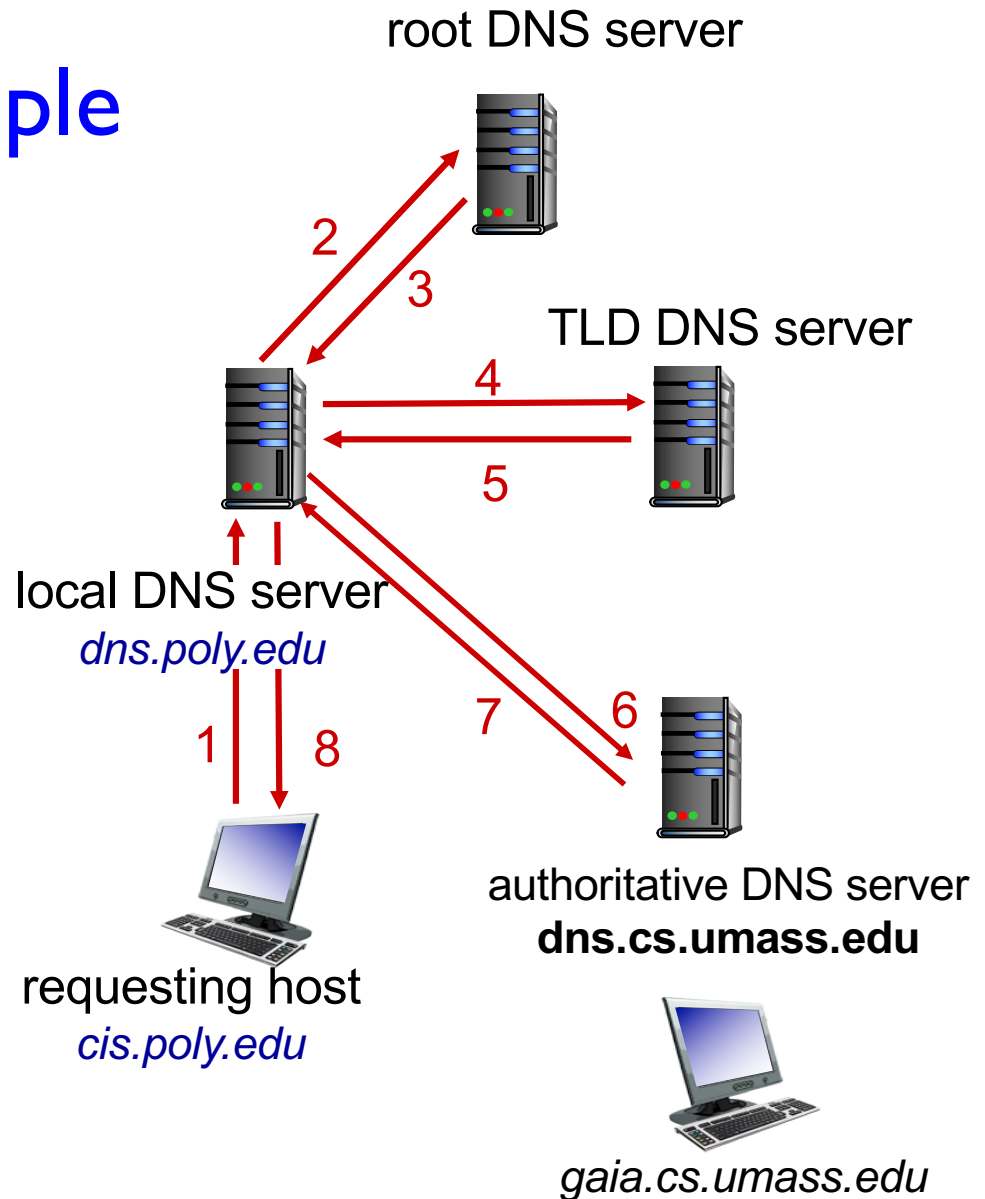
- does not strictly belong to hierarchy
- each ISP (residential ISP, company, university) has one
 - also called “default name server”
- when host makes DNS query, query is sent to its local DNS server
 - has local cache of recent name-to-address translation pairs (but may be out of date!)
 - acts as proxy, forwards query into hierarchy
- Try “nslookup <domain-name>” on a DICE machine

DNS name resolution example

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

iterated query:

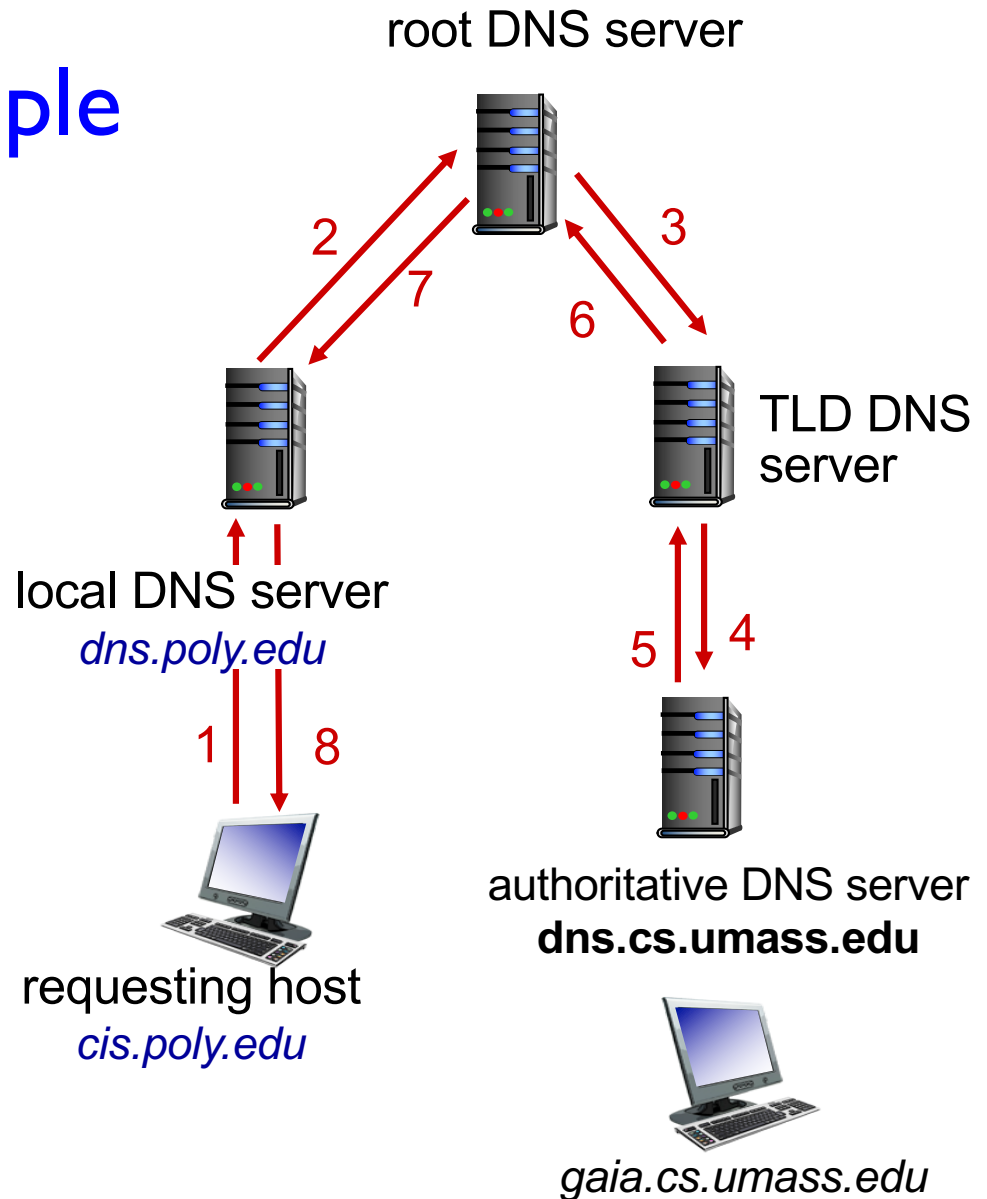
- ❖ contacted server replies with name of server to contact
- ❖ “I don’t know this name, but ask this server”



DNS name resolution example

recursive query:

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



DNS: caching, updating 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 db 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

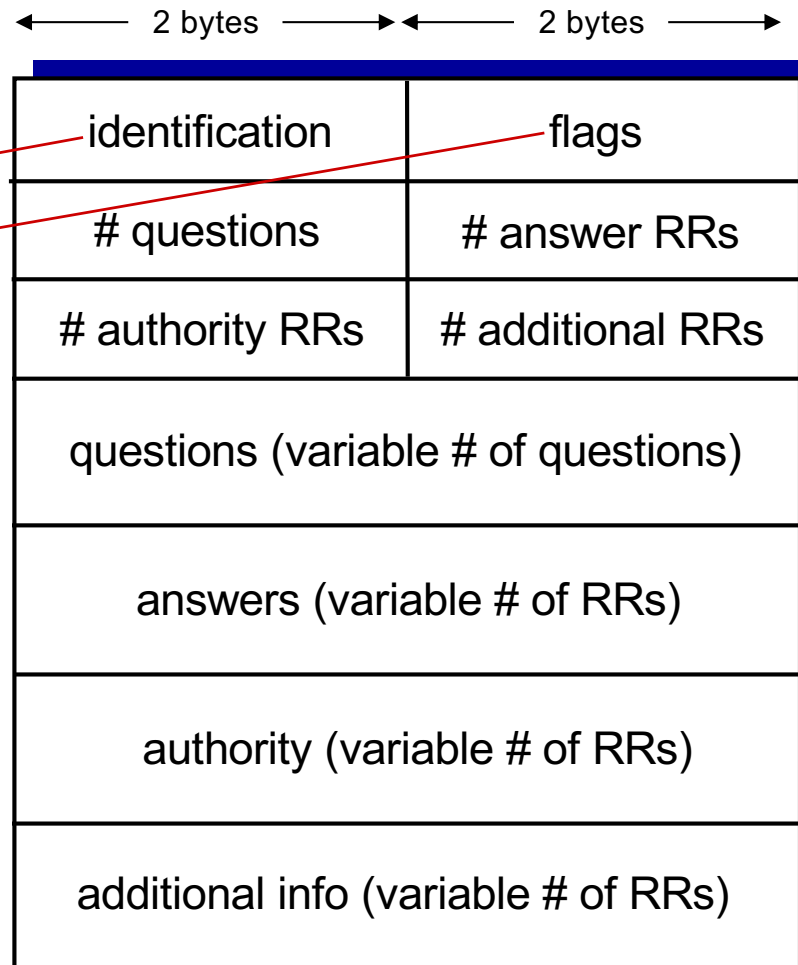
- **value** is name of mailserver associated with **name**

DNS protocol, messages

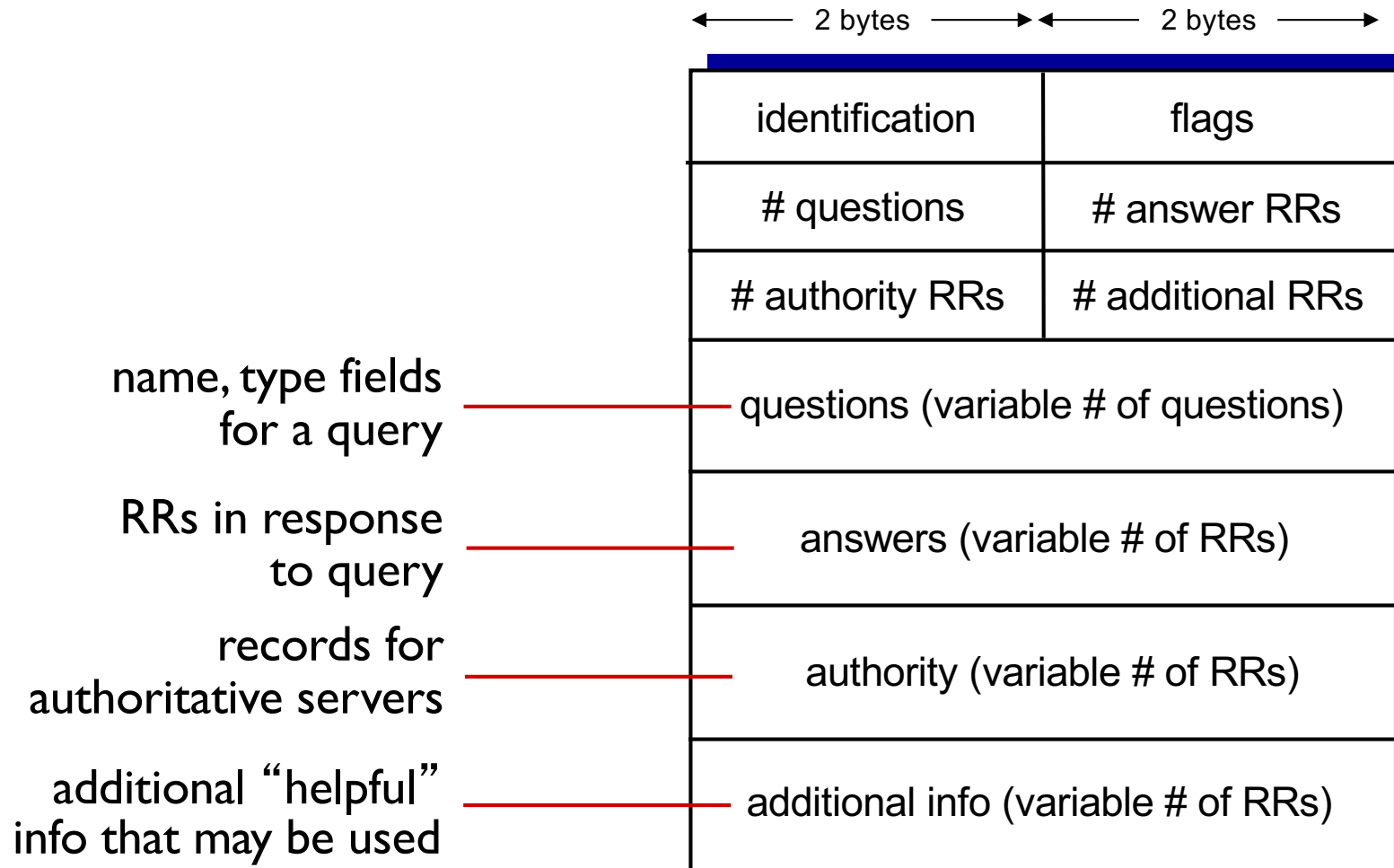
- *query* and *reply* messages, both with same *message format*

msg header

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



DNS protocol, messages



Inserting records into DNS

- example: new startup “Network Utopia”
- register name networkutopia.com at *DNS registrar* (e.g., Network Solutions)
 - provide 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 www.networkutopia.com