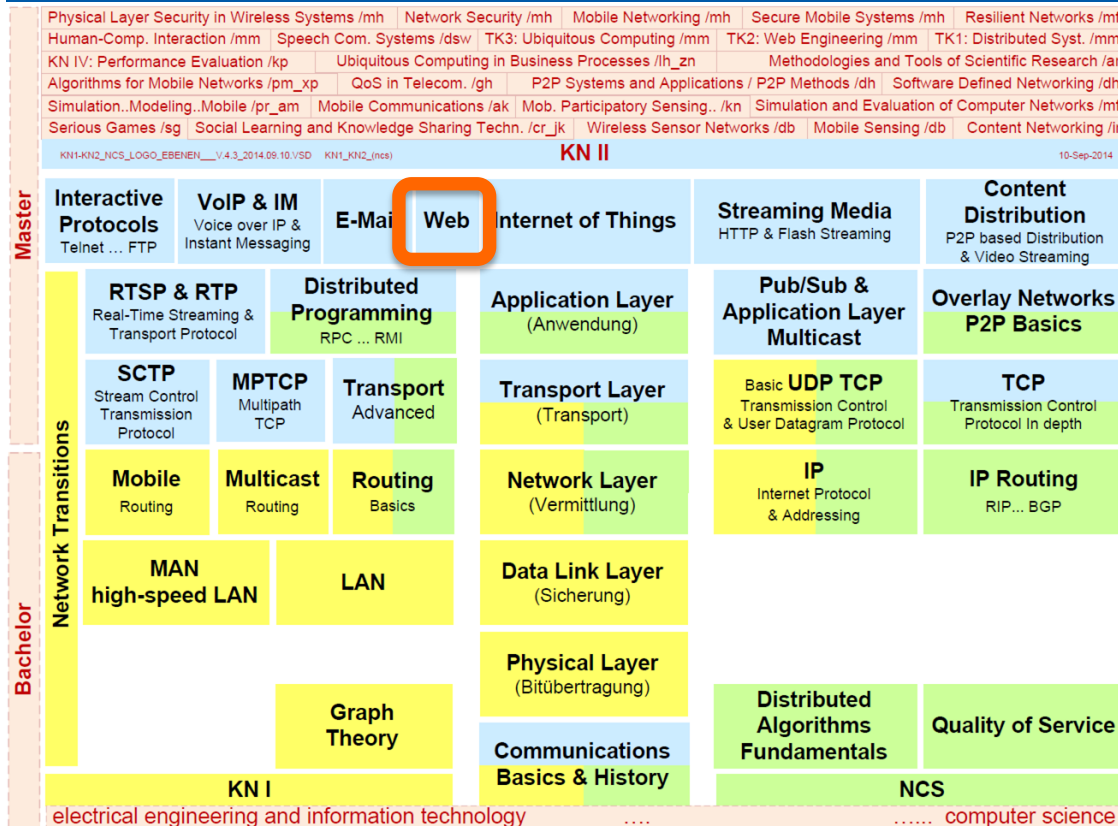


Communication Networks II

Web



TECHNISCHE
UNIVERSITÄT
DARMSTADT



10. September 2014

Overview

1 Background & History

1.1 Web – Predicted in the 1960's

1.2 Browser Statistics

1.3 Address - Uniform Resource Locators (URL)

2 Hypertext Markup Language (HTML)

3 Hypertext Transfer Protocol (HTTP)

3.1 Non-persistent HTTP

3.2 Persistent HTTP

3.3 HTTP Request Messages

3.4 HTTP Response Messages

3.5 HTTP Example

4 Some specific Web Issues

4.1 Cookies

4.2 Web Caches (Proxy Servers)

5 Domain Name System (DNS)

5.1 DNS Structure

5.2 DNS Name Resolution Example

5.3 DNS Records

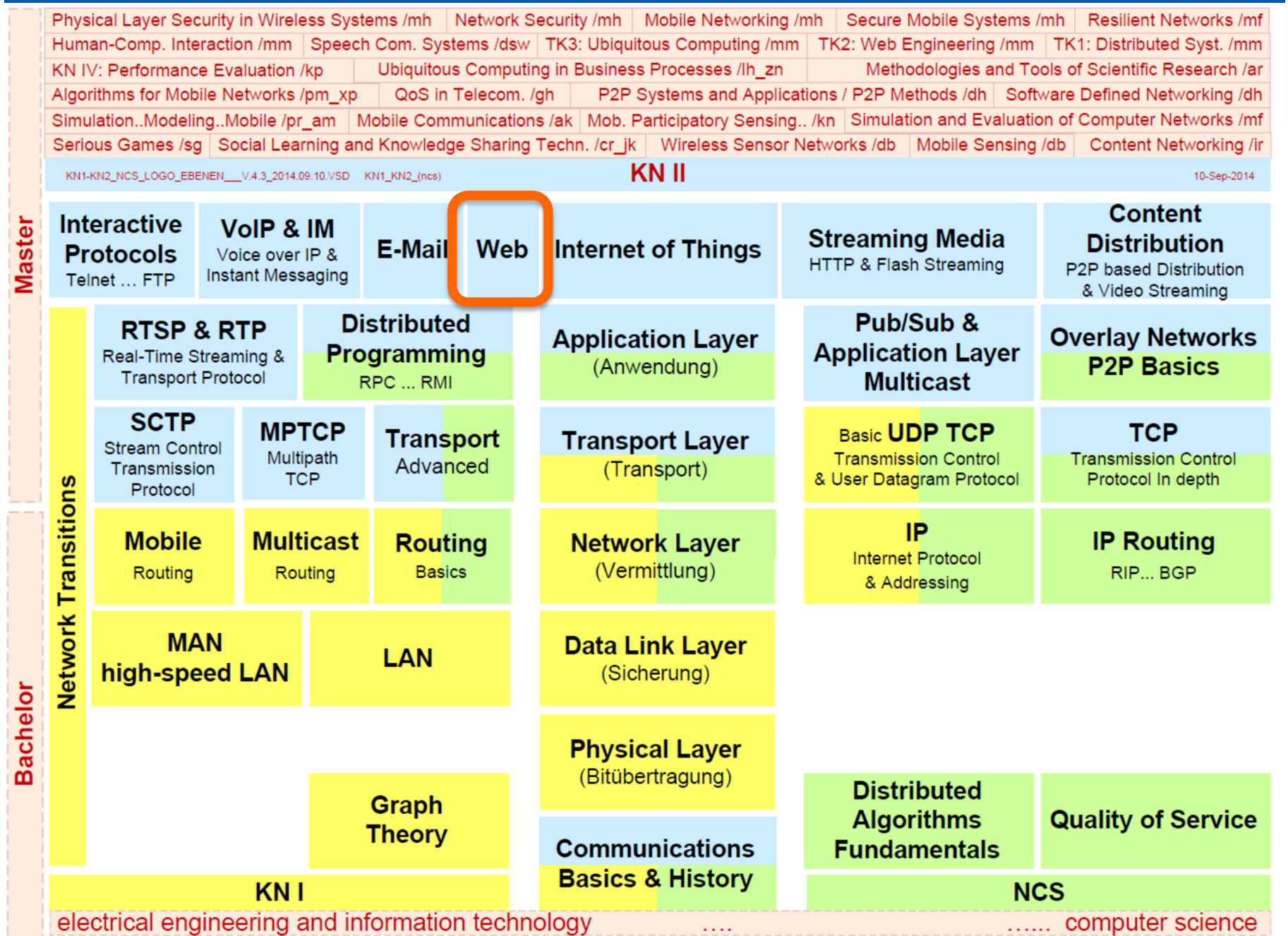
5.4 DNS Protocol

5.5 DNS Security Issues

6 Future Evolution: Semantic Web

6.1 Situation Today

6.2 Semantic Web



1 Background & History



The Computer as a Communication Device

In a few years, men will be able to communicate more effectively through a machine than face to face.

we believe that we are entering a technological age in which we will be able to interact with the richness of living information—not merely in the passive way that we have become accustomed to using books and libraries, but as active participants in an ongoing process, bringing something to it through our interaction with it, and not simply receiving something from it by our connection to it.

Such a medium is at hand—the programmed digital computer. Its presence can change the nature and value of communication even more profoundly than did the printing press and the picture tube, for, as we shall show, a well-programmed computer can provide direct access both to informational resources and to the *processes* for making use of the resources,

Web – Invented 1989

- 1989** Tim Berners-Lee (CERN, Geneva) publishes his first ideas
- 1993** approx. 50 web-servers
- 1993** NCSA distributes first version of Mosaic browser as shareware
- 1994** CERN and MIT found W3 Organization (see also <http://www.w3.org>)

...

- 1995** HTML defined as HTML 2.0. in RFC 1866
- 1996** HTTP 1.0 defined in RFC 1945
- 1996** HTML 3.2 consensus for 1996
- 1998** HTML 4.0 and very few variations
- 1999** HTTP 1.1 defined in RFC 2616

...

- ...** XML, CSS, ... RSS feeds, ...AJAX, REST, Mashups...

- Today** HTML 5.x ... web server everywhere

...

- Future** Semantic web to add “knowledge” to nodes (as metadata)

1.2 Browser Statistics



**Quiz on worldwide
usage**

Out of 100%

Firefox ?

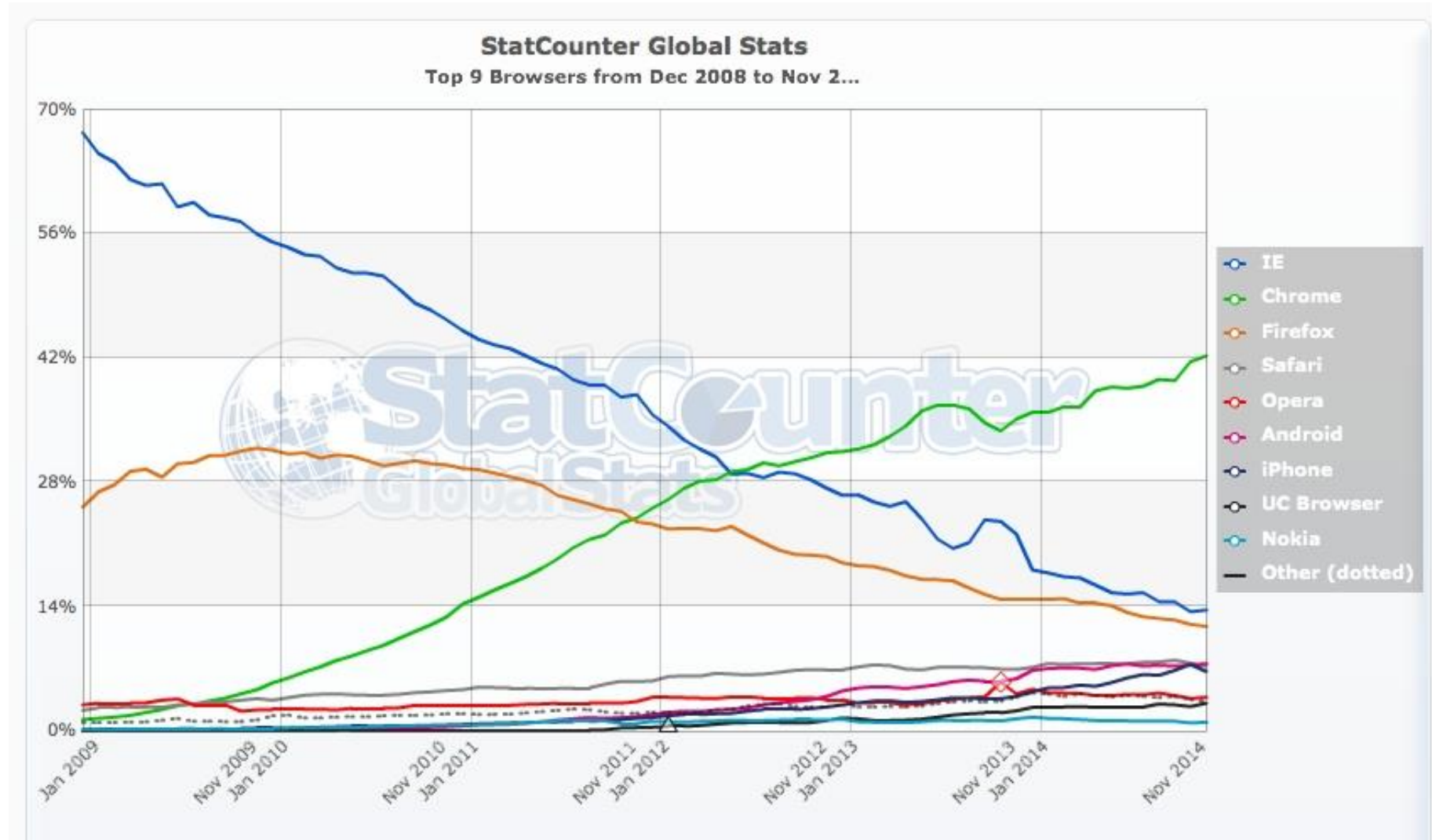
Chrome ?

Opera ?

Internet Explorer ?

Safari ?

Further ?



1.3 Address - Uniform Resource Locators (URL)

URL is the „address“ of a page

Format: <SCHEME>:<SCHEME-SPECIFIC-PART>

- `http://<host>:<port>/<path>?<query_string>`
- `ftp://<user>:<password>@<host>:<cwd1>/../<cwdN>/<name>;type=<typecode>`
- `mailto:<rfc822-addr-spec>`
- `telnet://<user>:<password>@<host>:<port>`
- `file://<host>/path`

	Used for	Example
Examples	Hypertext	http://www.kom.tu-darmstadt.de/
	FTP	ftp://ftp.ibr.cs.tu-bs.de/README
	Local file	file://home/rsteinmetz/.signature
	Sending email	mailto:kn2-lecture@kom.tu-darmstadt.de
	Remote login	telnet://www.w3.org:80

2 Hypertext Markup Language (HTML)



Web page consists of objects

- Object can be
 - HTML file
 - JPEG image
 - Java applet
 - Audio file,...



Web page consists of base HTML file

- Includes several referenced objects
- Each object is addressable by a URL (Uniform Resource Locator)
 - [http://www.kom.tu-darmstadt.de/.....](http://www.kom.tu-darmstadt.de/)



Web pages are accessed using hypertext transfer protocol (HTTP)

Hypertext Markup Language (HTML)

Hypertext Markup Language (HTML)

- Based on standard generalized markup language (SGML)
 - Standardized by ISO
- Defines markup tags for structuring text
- Browser interprets tags and generates page layout

Example HTML tags

- `<HEAD>...</HEAD>` page header
- `<BODY>...</BODY>` page body
- `<H1>...</H1>` headline
- `<P>...</P>` paragraph
- `` inserted image
- `...` link to another document

Example file

```
<html>
  <body>
    <h1>Hello World!</h1>
  </body>
</html>
```

Defining hyperlinks

- Format: ` hyperlinked text `
- Example code:
Click ` here ` to go to TU Darmstadt.
- Example layout:
Click [here](http://www.tu-darmstadt.de) to go to TU Darmstadt.

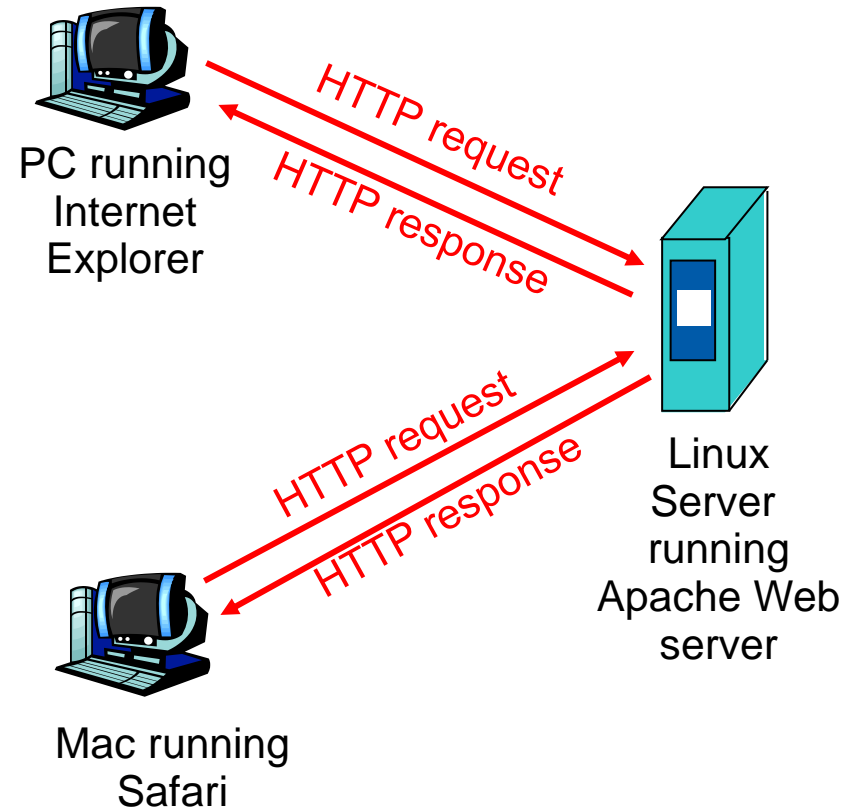
3 Hypertext Transfer Protocol (HTTP)

Hypertext transfer protocol (HTTP)

- Web's application layer protocol
- Client/server model
 - Client: browser that requests, receives, “displays” Web objects
 - Server: Web server sends objects
 - In response to requests

HTTP uses TCP

- Client initiates TCP connection
 - To serverport 80
- Server accepts TCP connection
- HTTP messages (application-layer protocol messages) exchanged
 - Between browser (HTTP client) and Web server (HTTP server)
- TCP connection closed



Hypertext Transfer Protocol (HTTP)

HTTP is a stateless protocol

- Server maintains no information about past client requests
- Side note:
protocols that maintain state are more complex
 - Past history (state) must be maintained
 - If server/client crashes
 - Their views of state may be inconsistent
 - State must be reconciled
 - Compare e.g. protocols in database world

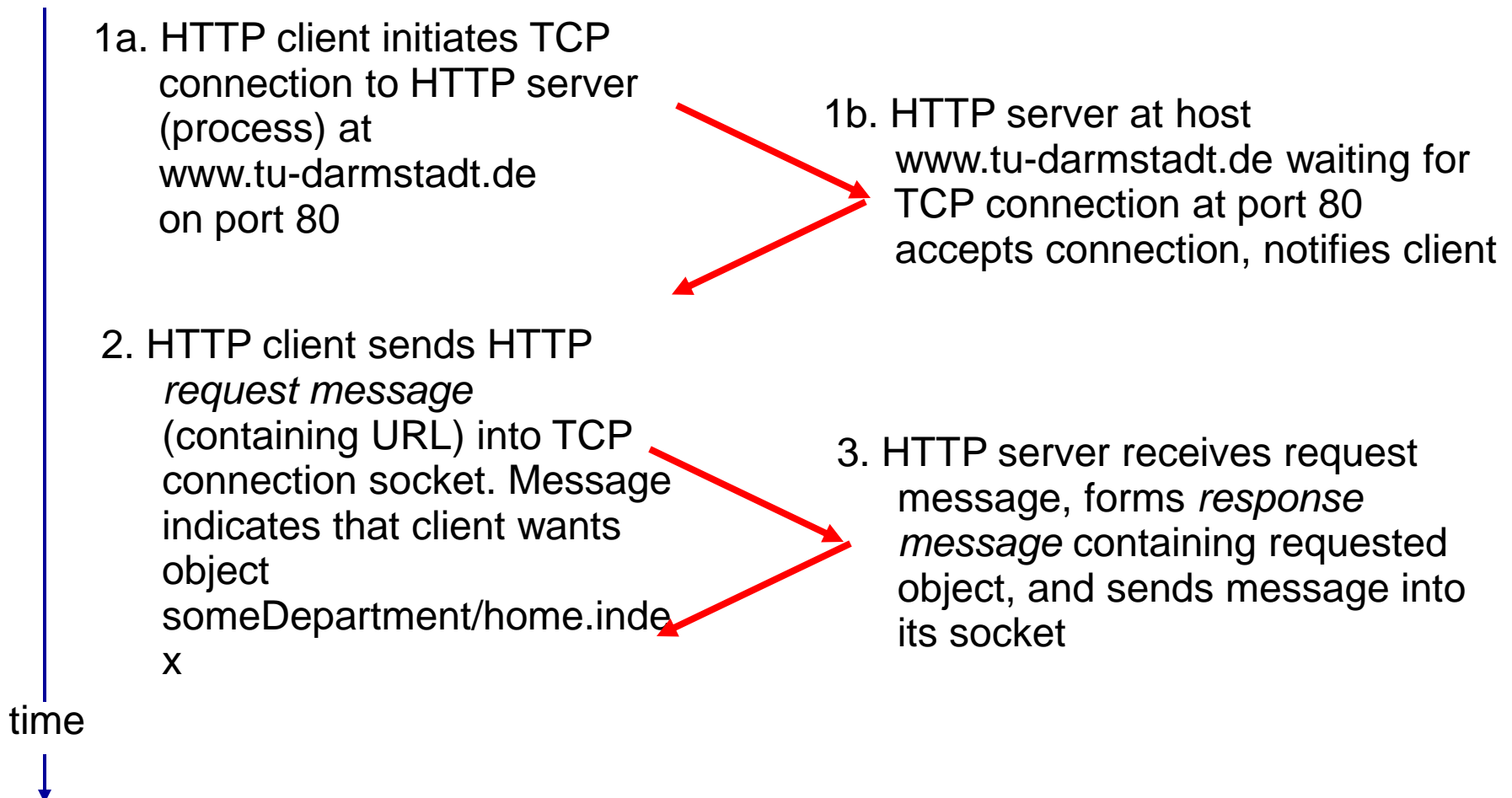
HTTP connections can be

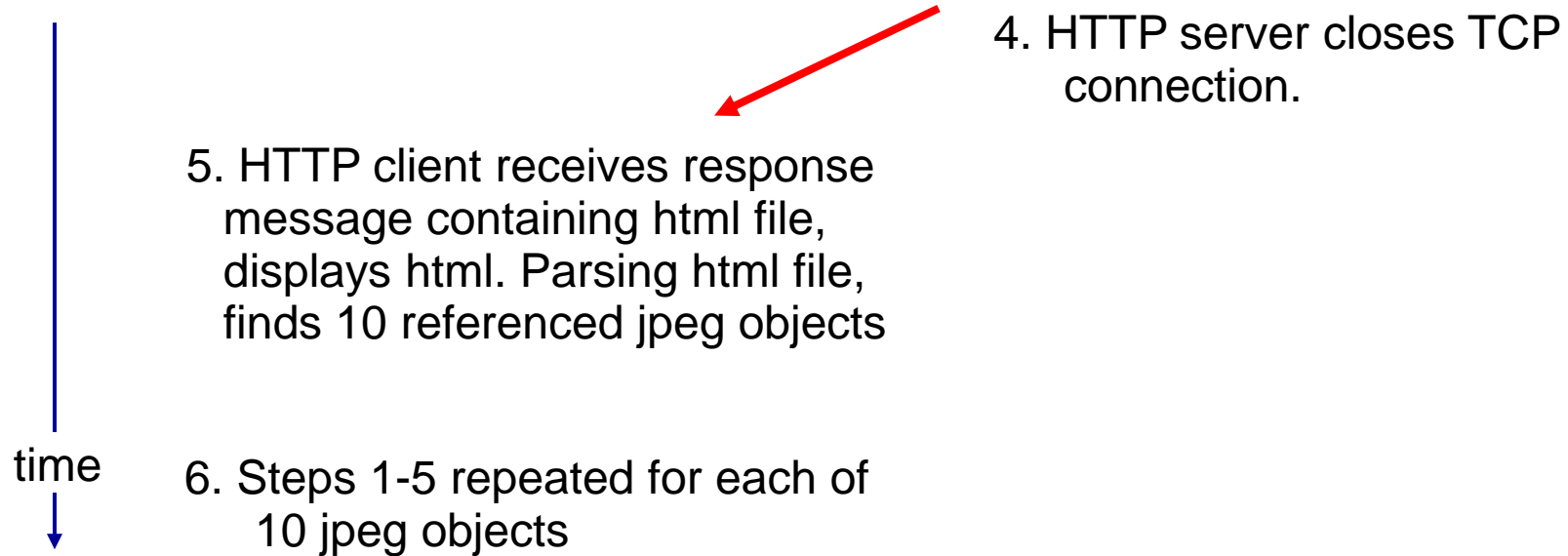
- Non-persistent (HTTP 1.0)
 - At most one object sent over one TCP connection
- Persistent (HTTP 1.1)
 - Multiple objects can be sent over one TCP connection

3.1 Non-persistent HTTP

Suppose user enters URL `http://www.tu-darmstadt.de`

- Suppose page contains text and references to 10 jpeg images





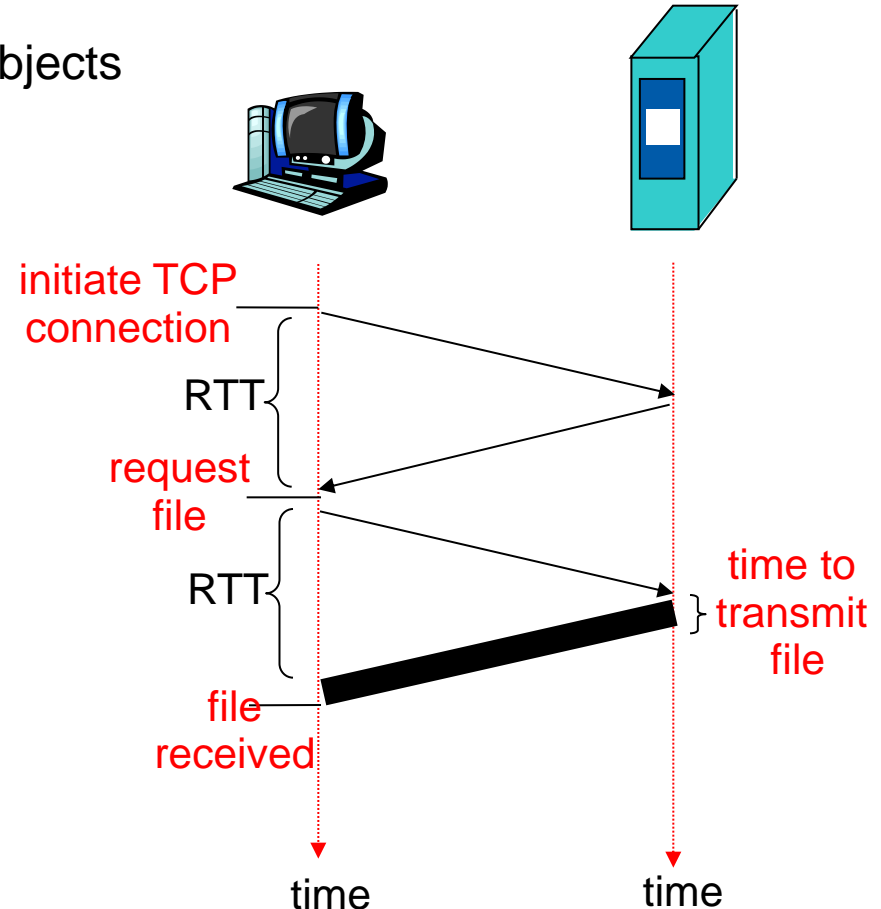
Non-persistent HTTP

A closer look at the response time

- For loading a web page with multiple objects
- Based on round trip time (RTT)
 - Between client and server

Response time consists of

- One RTT to initiate TCP connection
- One RTT for HTTP request
 - And first few bytes of HTTP response to return
- File transmission time
- Total = $2\text{RTT} + \text{transmission time}$



3.2 Persistent HTTP

Motivation: non-persistent HTTP issues

- Requires 2 additional RTTs per object
- Operating system overhead for each TCP connection
- Amplified by browsers opening parallel TCP connections
 - To fetch referenced objects

Persistent HTTP

- Server leaves connection open
 - After sending response
- Subsequent HTTP messages sent over open connection
- Client sends requests as soon as it encounters a referenced object
- Overhead reduced to ~one RTT for all referenced objects

3.3 HTTP Request Messages

Two types of HTTP messages: request, response

HTTP request message

- ASCII (human-readable format)

request line
(GET, POST,
HEAD commands)

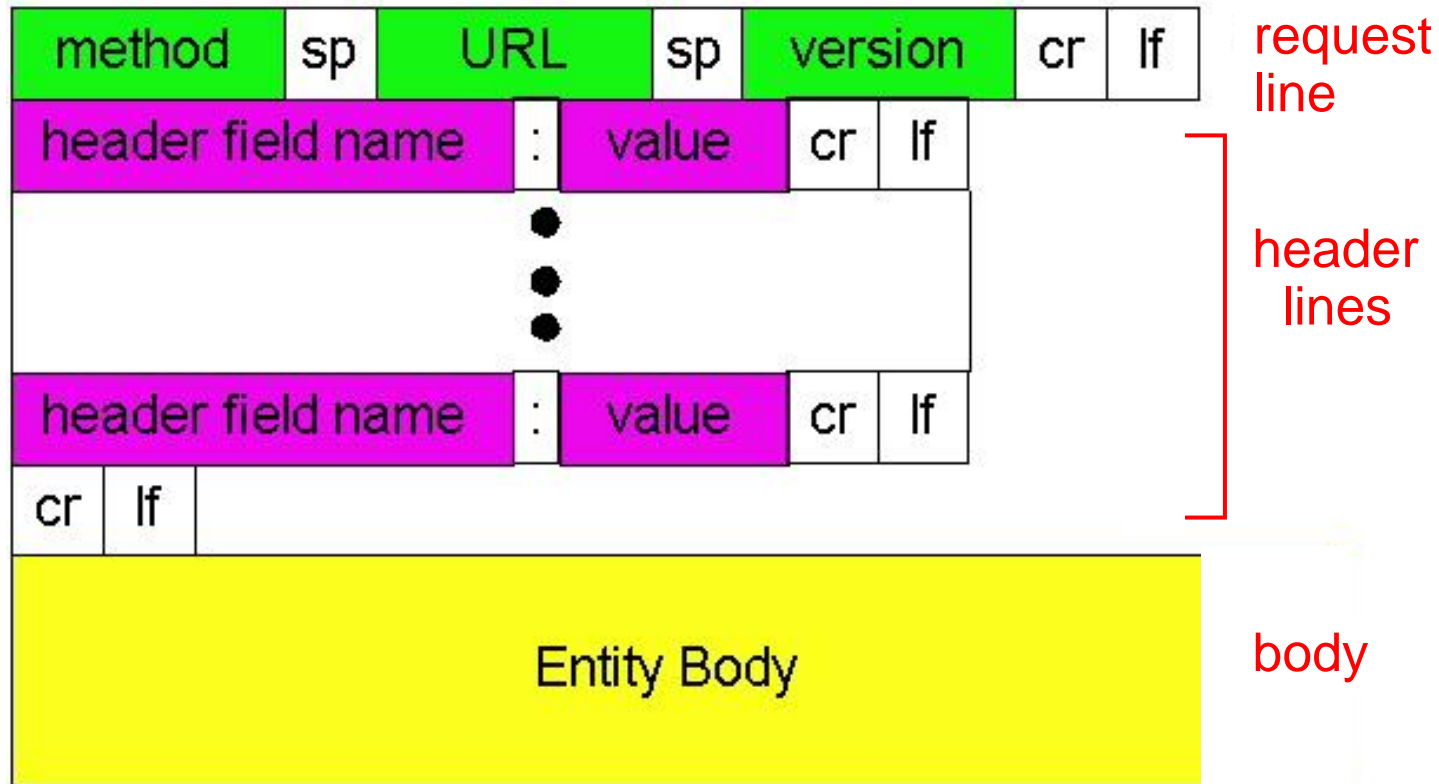
header
lines

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

General format of HTTP requests



HTTP method types

HTTP/1.0

- GET
 - Request file specified in URL field
- POST
 - Send request with additional info in entity body
- HEAD
 - Asks server to leave requested object out of response
 - Response only is metadata contained in HTTP message header

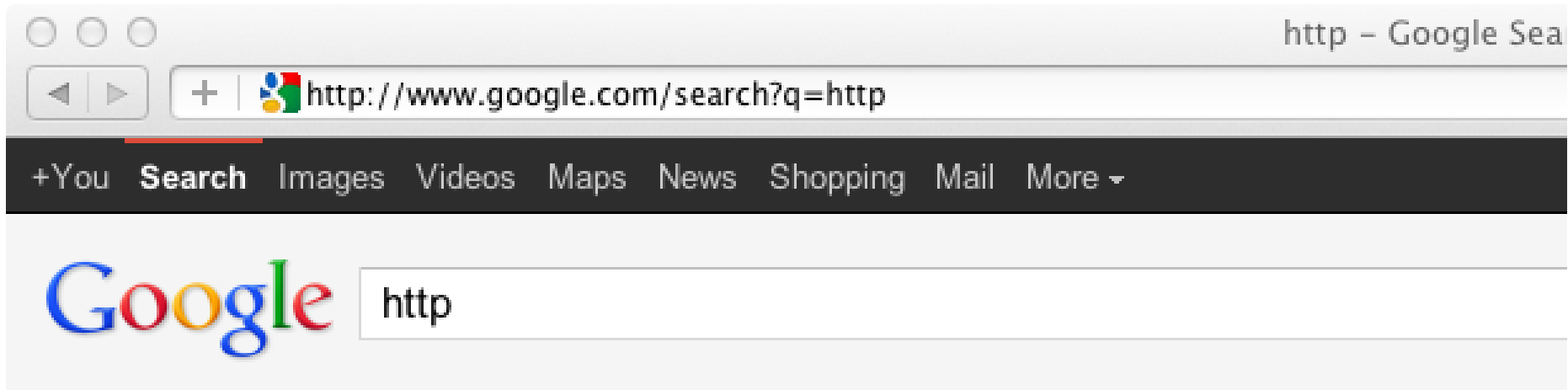
HTTP/1.1

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

HTTP knows two ways for uploading form input

1. form input based on GET method

- Input encoded in URL (query_string)
- Should be used for idempotent queries (i.e. all requests lead to same results)
 - I.e. reload is not harmful



Search

About 19,490,000,000 results (0.13 seconds)

Everything

[Hypertext Transfer Protocol - Wikipedia, the free encyclopedia](https://en.wikipedia.org/wiki/Hypertext_Transfer_Protocol)
en.wikipedia.org/wiki/Hypertext_Transfer_Protocol

HTTP knows two ways for uploading form input (cont'd)

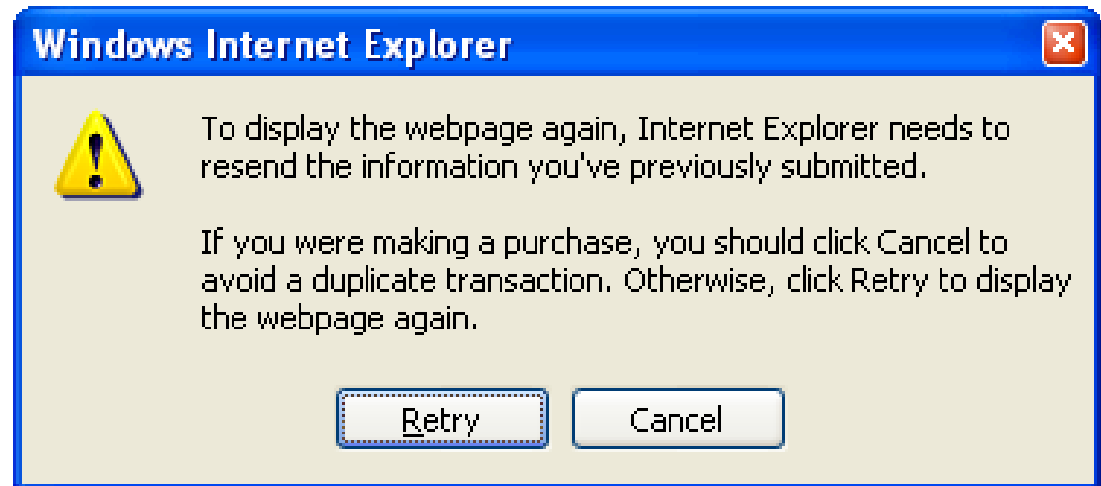
2. Form input based on POST method

- Input contained in entity body of HTTP request

- Should be used for non-idempotent queries

- I.e. reloading page would repeat actions (e.g. placing an order)
- Browser will display a warning message

- Should be also used
 - If input is non ASCII
 - If input exceeds certain length



3.4 HTTP Response Messages

HTTP response message

- ASCII (human-readable format)

status line
(protocol
status code
status phrase)

header
lines

```
HTTP/1.1 200 OK\r\n
Date: Sun, 26 Sep 2020 20:09:20 GMT\r\n
Server: Apache/2.0.52 (CentOS)\r\n
Last-Modified: Tue, 30 Oct 2019 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 ...
```

data, e.g.,
requested
HTML file

HTTP Response Messages

HTTP response status codes

- Appear in 1st line in server-client response message

Some sample codes

- 200 OK
 - Request succeeded, requested object contained in message
- 301 Moved Permanently
 - Requested object moved, new location specified in message
- 400 Bad Request
 - Request message not understood by server
- 404 Not Found
 - Requested document not found on this server
- 505 HTTP Version Not Supported

3.5 HTTP Example

Example: accessing simple web page

- File index.html at server glab013.g-lab.tu-darmstadt.de



Hello World!

- Source code received by client's web browser

Source of http://glab013.g-lab.tu-darmstadt.de/

```
<html><body><h1>Hello World!</h1>  
</body></html>
```

HTTP Example

Example: accessing simple web page (cont'd)

- Talking HTTP using Telnet

```
dhcp47:~ Andre$ telnet glab013.g-lab.tu-darmstadt.de 80
Trying 130.83.125.13...
Connected to glab013.g-lab.tu-darmstadt.de.
Escape character is '^]'.
GET /index.html HTTP/1.1
Host: glab013.g-lab.tu-darmstadt.de

HTTP/1.1 200 OK
Date: Fri, 06 Jan 2012 14:21:23 GMT
Server: Apache/2.2.20 (Ubuntu)
Last-Modified: Fri, 06 Jan 2012 13:41:40 GMT
ETag: "19e6-31-4b5dc356b8f8e"
Accept-Ranges: bytes
Content-Length: 49
Vary: Accept-Encoding
Content-Type: text/html

<html><body><h1>Hello World!</h1>
</body></html>
Connection closed by foreign host.
dhcp47:~ Andre$ █
```

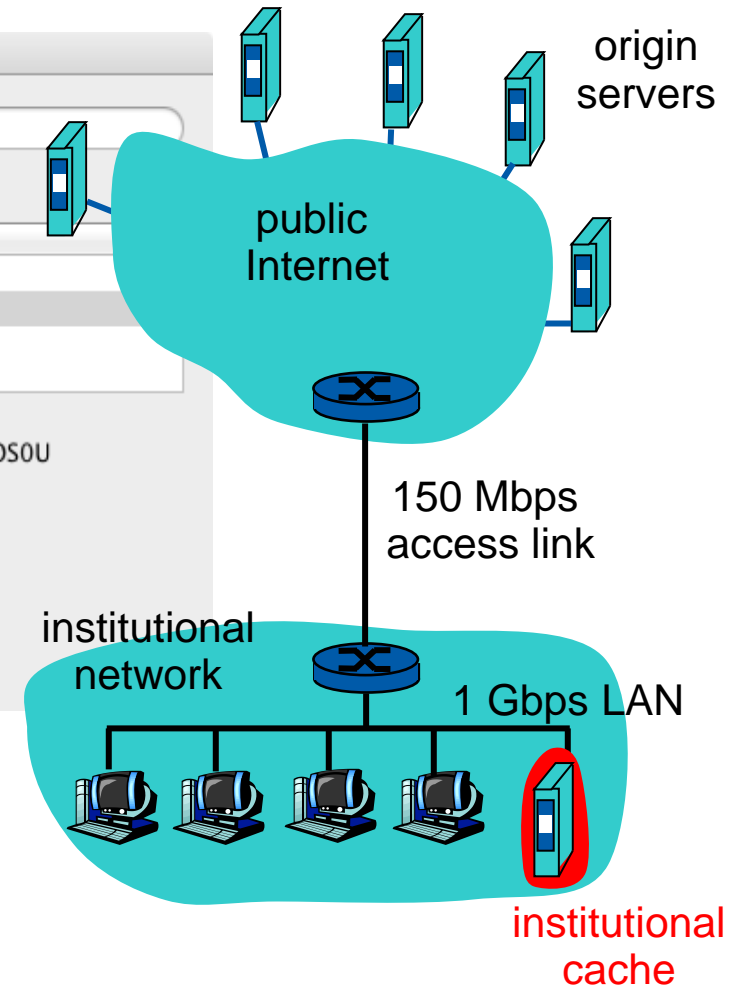
4 Some specific Web Issues



Cookies



Web Caches



4.1 Cookies

Cookies allow to maintain user state on server

- Recall: HTTP itself is stateless protocol
- Latest specification in RFC 6265 as of April 2011
 - First specification in RFC 2109 as of February 1997
- Used to remember user-specific server settings
 - Authorization, shopping carts, recommendations, session state
- Used by many web sites today
 - Take a look at the cookie list of your browser...

Four components required

- Cookie header line of HTTP response message
 - To transfer cookie from server to client
- Cookie header line of HTTP request message
 - To send back cookie from client to server
- Cookie file managed by client's browser
 - General format: name:value pairs
- Back-end database at web server
 - Storing user state

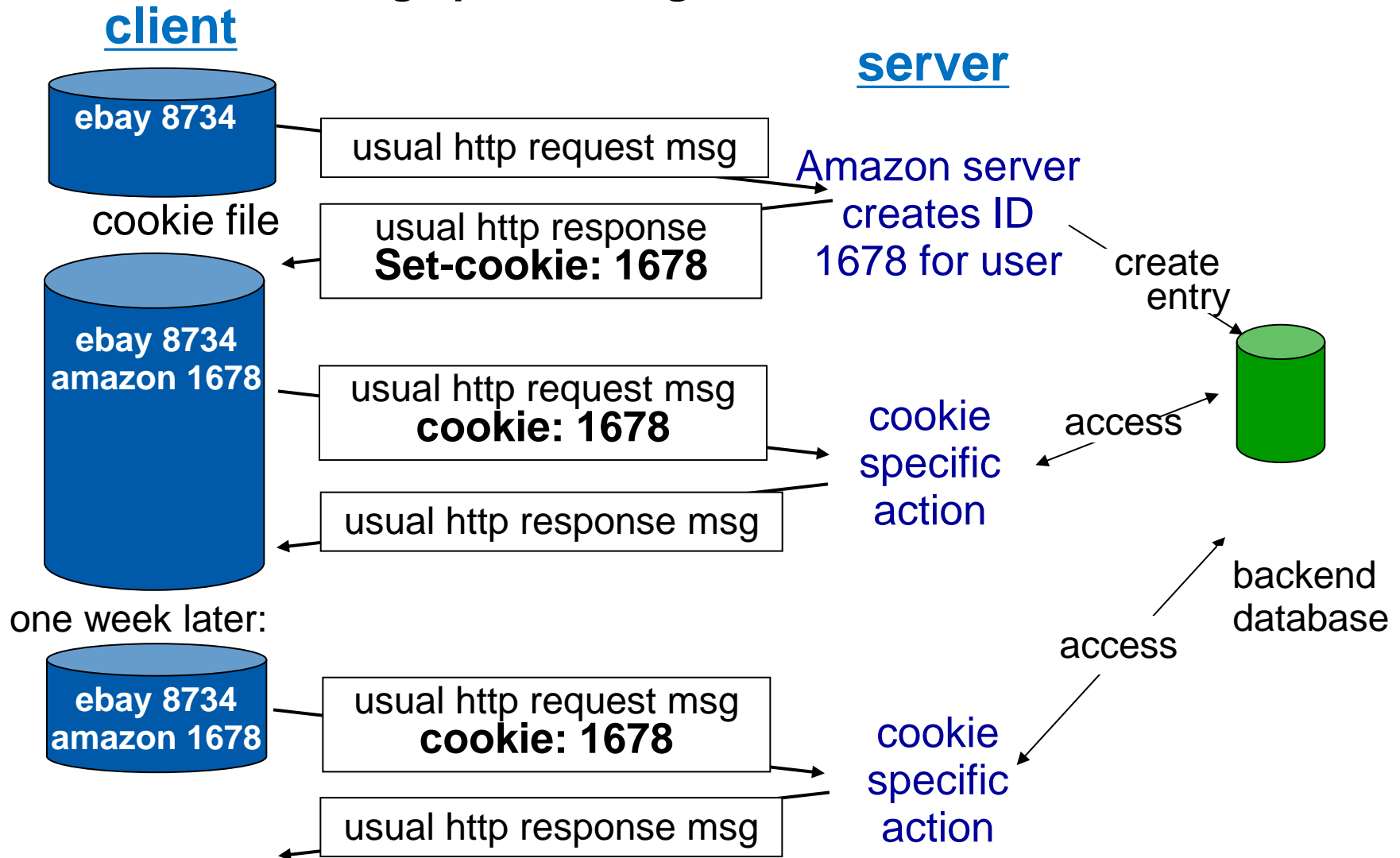
Cookies allow to maintain user state on server

- Recall:
HTTP itself is stateless protocol
- Latest specification in RFC 6265 as of April 2011
 - First specification in RFC 2109 as of February 1997
- Used to remember user-specific server settings
 - Authorization, shopping carts, recommendations, session state
- Used by many web sites today
 - Take a look at the cookie list of your browser...

Four components required

- Cookie header line of HTTP response message
 - To transfer cookie from server to client
- Cookie header line of HTTP request message
 - To send back cookie from client to server
- Cookie file managed by client's browser
 - General format: name:value pairs
- Back-end database at web server
 - Storing user state

Setting up and using cookies



Example: Firefox cookies

- Settings - privacy



Search:

The following cookies are stored on your computer:

Site	Cookie Name
▼ google.com	
google.com	PREF
▶ google.de	

Name: PREF
Content: ID=22dad3f0151dea87:FF=0:TM=1325862039:LM=1325862039:S=n-PRkShty2b3OS0U
Domain: .google.com
Path: /
Send For: Any type of connection
Expires: January 5, 2014 4:00:39 PM

Remove Cookie Remove All Cookies

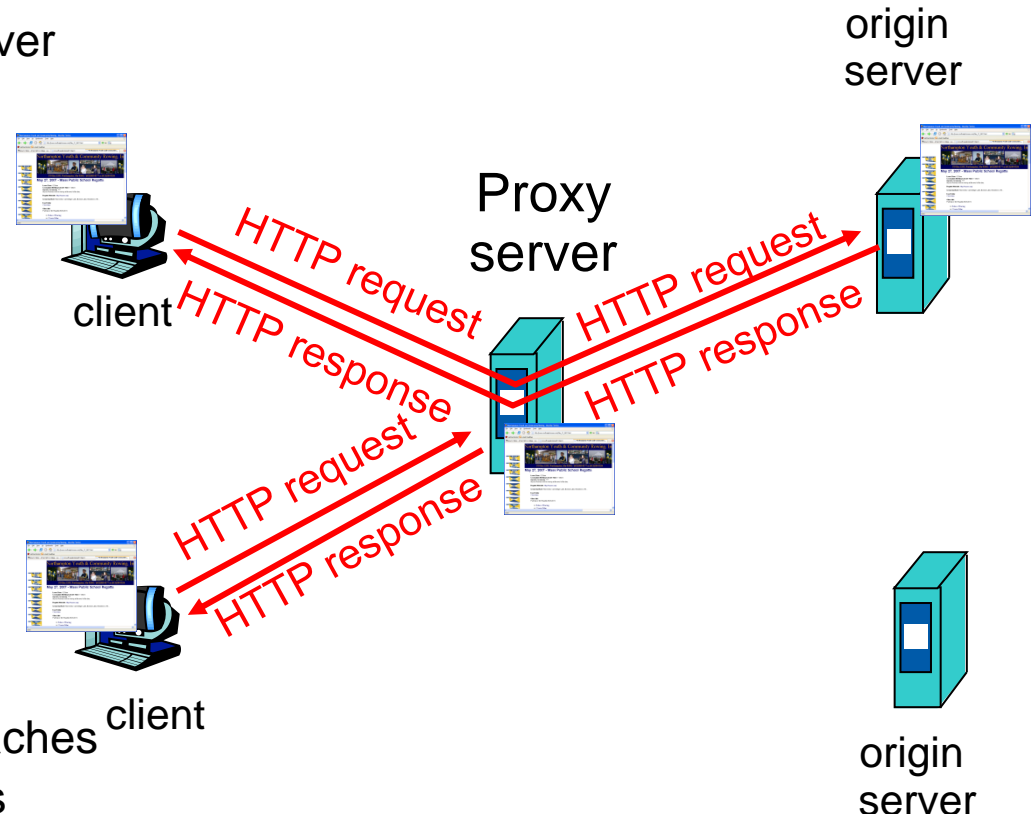
4.2 Web Caches (Proxy Servers)

Goal: to satisfy client request without involving origin server

- Typically cache is installed by ISP
 - University, company, residential ISP
- User sets in browser to access web via cache
- Browser sends all HTTP requests to cache
- Cache acts as both, client and server
- If Object in cache
 - Cache returns object
 - Else cache requests object
 - From origin server
 - Then returns object to client

Why web caching?

- To reduce response time for client request
- To reduce traffic on an institution's access link
- To increase Internet dense with caches
 - Enables “poor” content providers to effectively deliver content



Web Caches (Proxy Servers) - Example

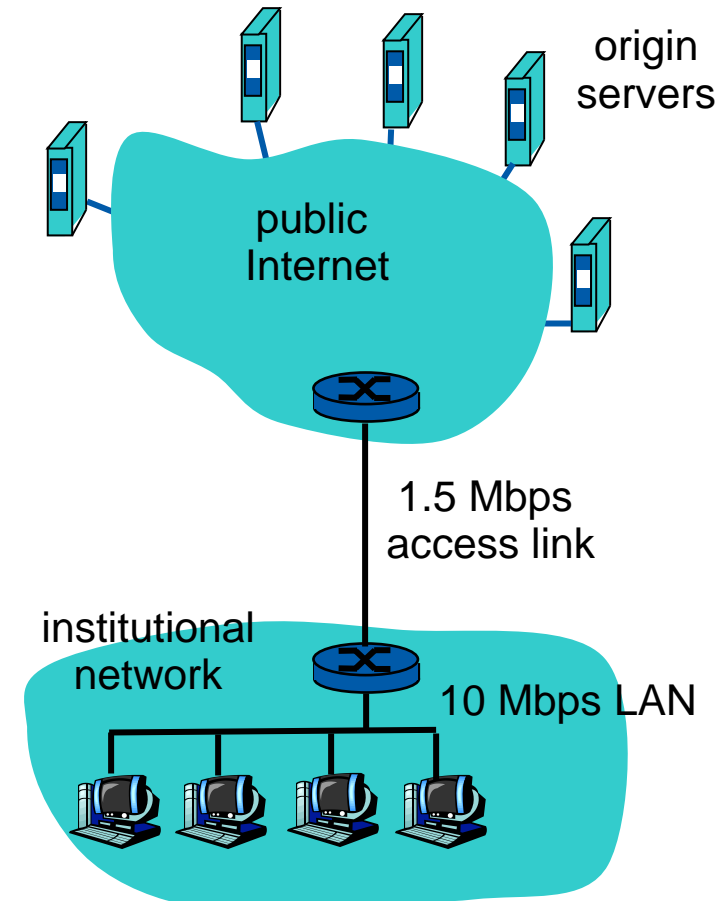
Problem Statement

Assumptions

- Average object size = 100.000 bits
- Average request rate from institution's browsers to origin servers = 15/sec
- Delay from institutional router to any origin server and back to router = 2 sec

Consequences

- Utilization on LAN = 15%
- Utilization on access link = 100%
- Total average delay
= Internet delay + access delay + LAN delay
= 2 sec + minutes + milliseconds



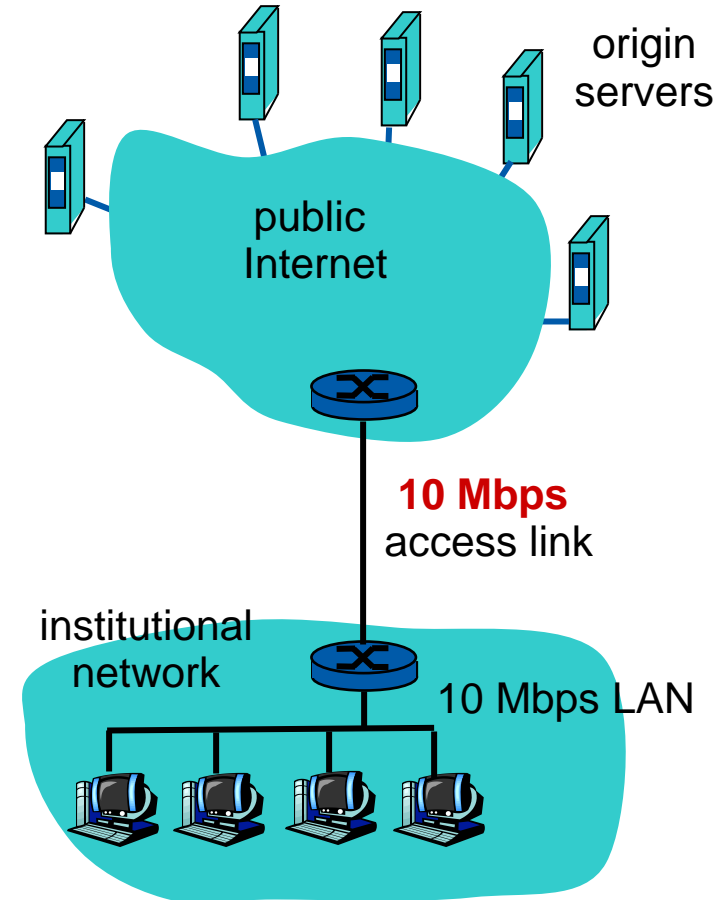
Web Caches (Proxy Servers) - Example

Possible 1.st solution

- Increase bandwidth of access link to e.g. 10 Mbps

Consequences

- Utilization on LAN = 15%
- Utilization on access link = 15%
- Total average delay
= Internet delay + access delay + LAN delay
= 2 sec + milliseconds + milliseconds
- But:
 - often a costly upgrade



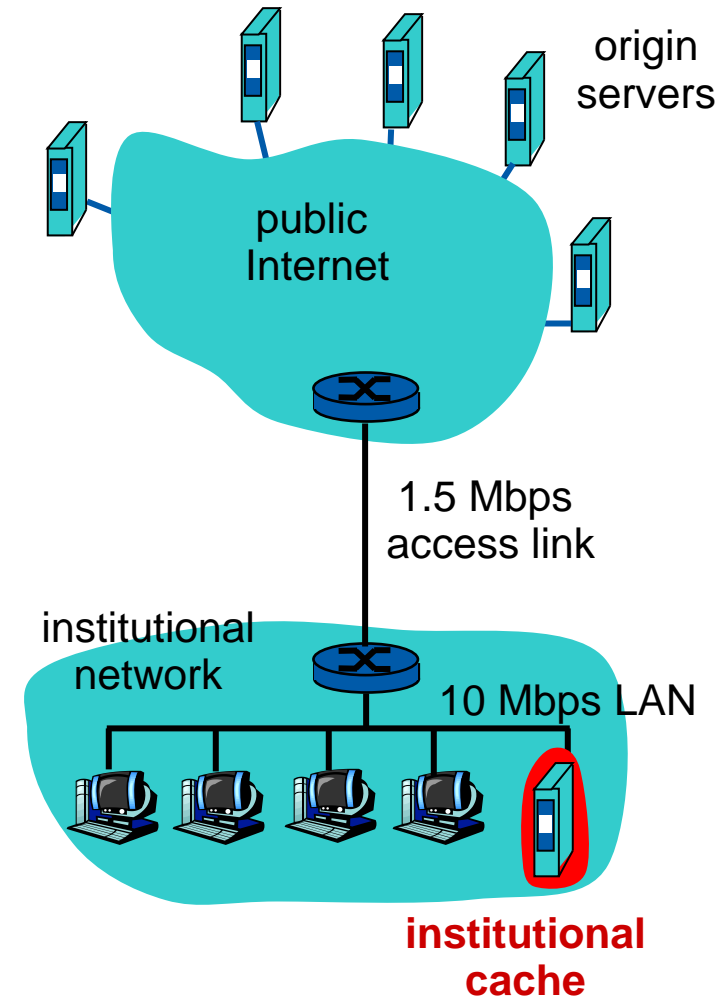
Web Caches (Proxy Servers) - Example

Possible cache solution

- Install cache

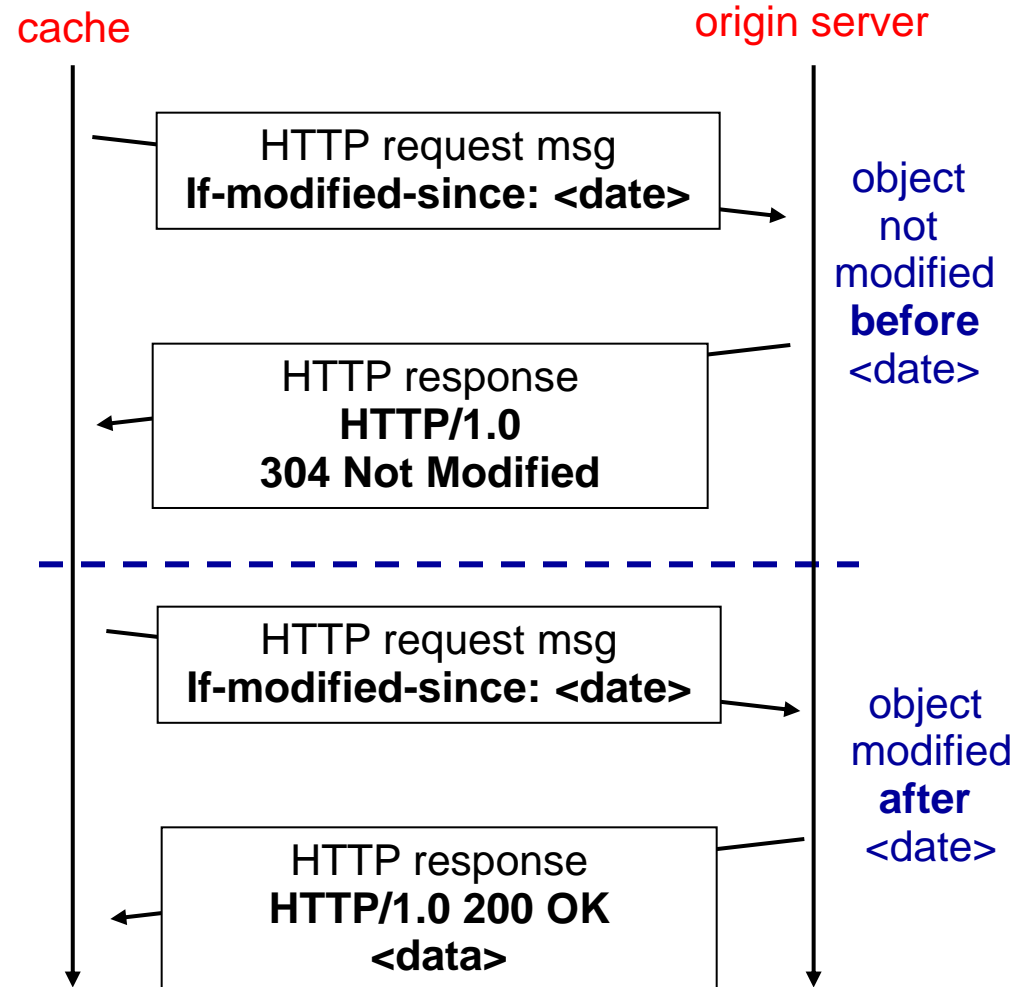
Consequence

- Suppose hit rate is 0.4
 - 40% requests satisfied almost immediately
 - 60% requests satisfied by origin server
- Utilization of access link reduced to 60%
 - Resulting in negligible delays (e.g. 10 msec)
- Total average delay
= Internet delay + access delay + LAN delay
= $0.6 \cdot (2.01) \text{ secs} + 0.4 \cdot \text{milliseconds}$
< 1.4 secs



Conditional GET

- Goal:
do not send object
if cache has up-to-date
cached version
- 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



5 Domain Name System (DNS)

Internet hosts, routers

- Addressed by machine readable IP address
 - Used for routing packets
- Addressed also by human readable “Name”
 - E.g. `www.google.com`
- How to map between IP address and name, and vice versa?
- At the early days
hosts.txt file was updated and transferred manually...

Domain Name System

- Distributed database
 - Implemented in hierarchy of many name servers
- Application-layer protocol
 - Used by host, routers, name servers
 - To communicate and resolve names (address/name translation)
- Note: DNS is one of Internet’s core functions
 - Implemented as application-layer protocol
 - Thus complexity at network’s “edge”

Domain Name System (DNS)

DNS services

- Hostname to IP address translation
- Host aliasing
 - Canonical, alias names
- Mail server aliasing
- Load distribution
 - E.g. replicated web servers
 - Set of IP addresses for one canonical name

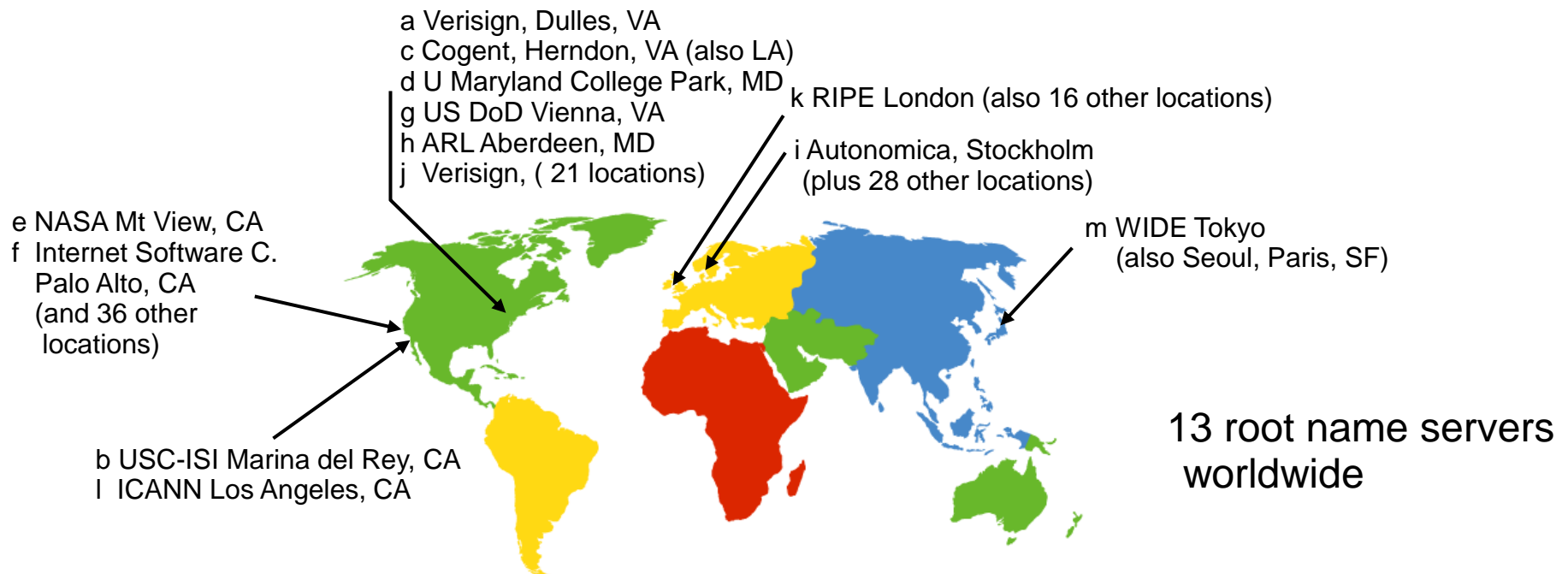
Why not centralize DNS?

- Single point of failure
- Traffic volume
- Distant centralized database
- Maintenance
 - Does not scale!

5.1 DNS Structure

DNS root name servers

- Contacted by local name server that can not resolve name
- Root name server
 - Contacts authoritative name server
 - If name mapping not known
 - Gets mapping
 - Returns mapping to local name server



DNS Structure

Top-level domain (TLD) servers

- Responsible for
 - .com, .org, .net, .edu, .aero, .jobs, .museums
 - And all top-level country domains, e.g.: de, uk, fr, ca, jp
- Network Solutions maintains servers for .com TLD
- Educause maintains servers for .edu TLD

Authoritative DNS servers

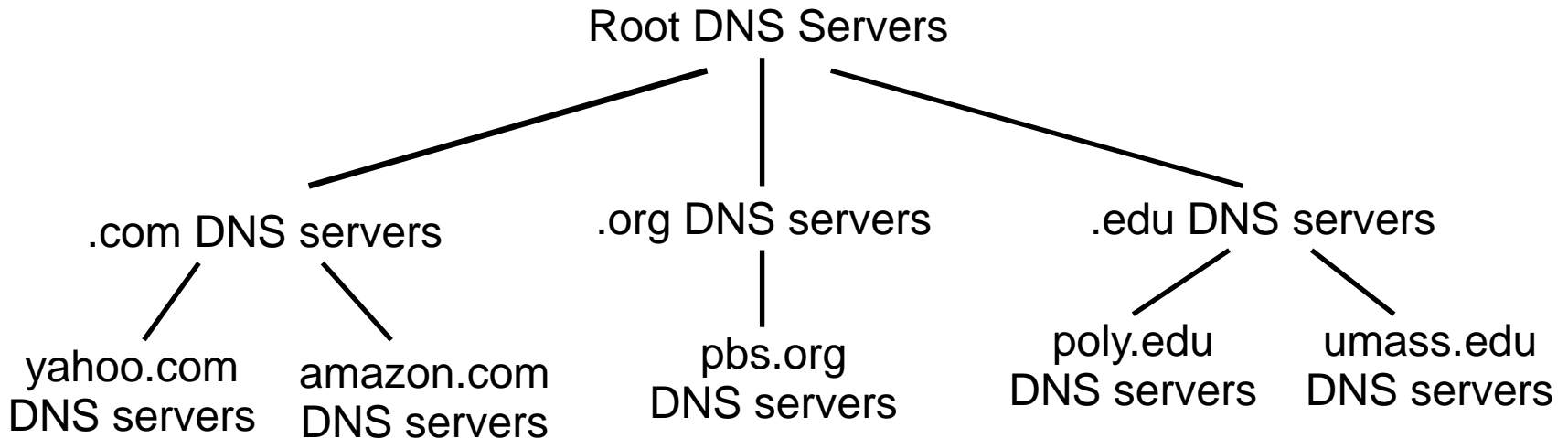
- Organization's DNS servers
 - Providing authoritative hostname to IP mappings for organization's servers (e.g., web, mail)
- Can be maintained by organization or service provider

Local name server

- Does not strictly belong to hierarchy
- Each ISP (residential, company, university, ...) has one
 - Also called default name server
- End system's DNS queries usually sent to local DNS server
 - Acts as proxy, forwards query into hierarchy

Client wants IP for **www.amazon.com**

- First approximation
 - Client queries a 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

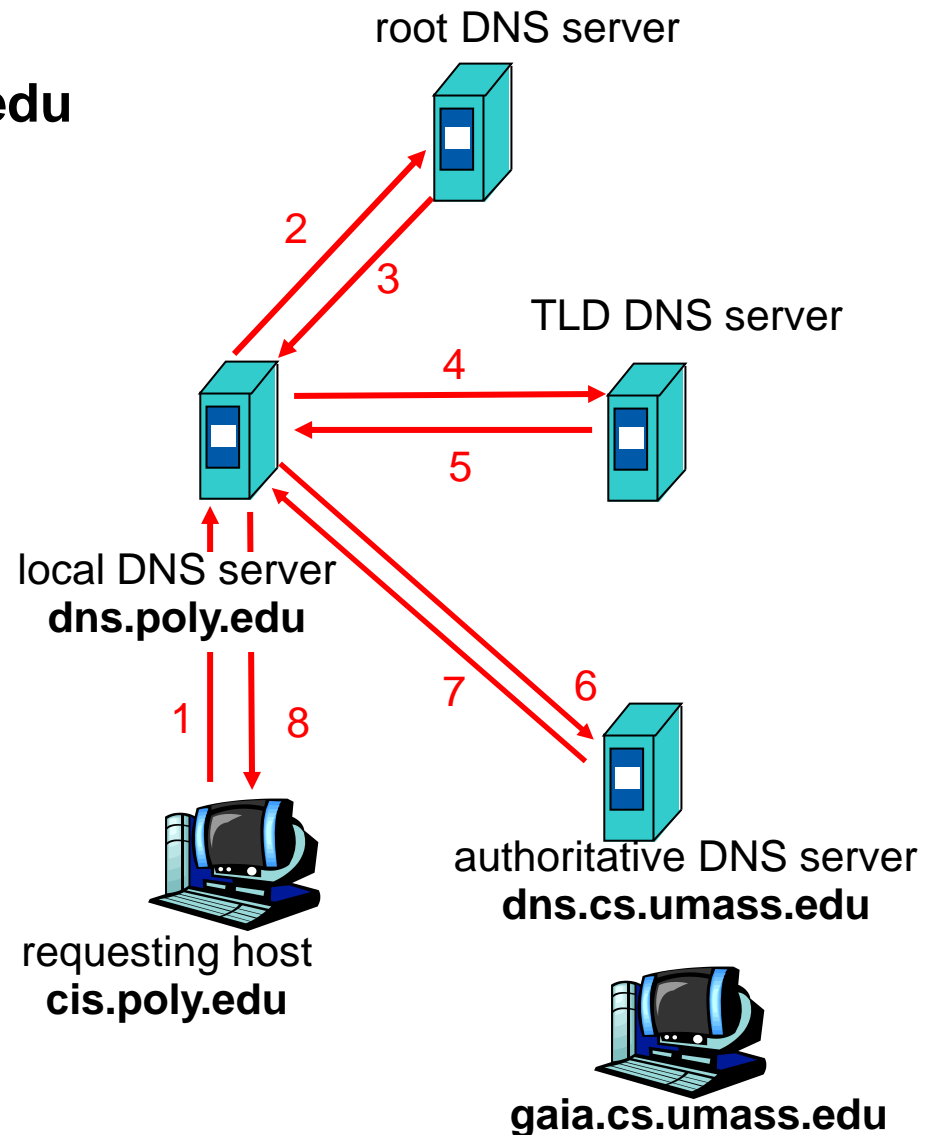


5.2 DNS Name Resolution Example

Host at **cis.poly.edu** wants
IP address for **gaia.cs.umass.edu**

Iterative query

- Local name server contacts server in hierarchy iteratively
- Contacted server replies
 - With name of server to contact
- “I don’t know this name, but ask this server”

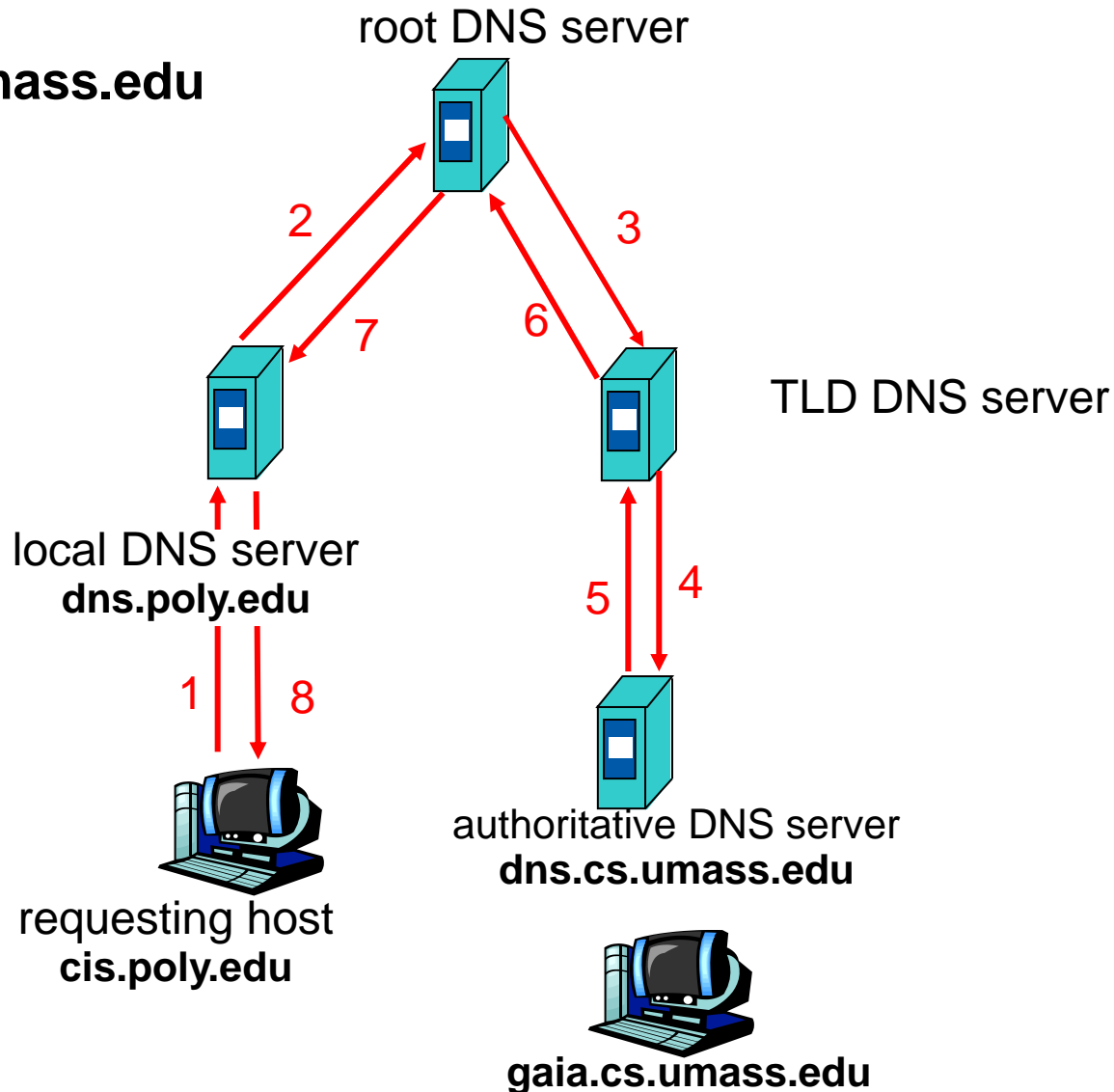


DNS Name Resolution Example

Host at **cis.poly.edu** wants
IP address for **gaia.cs.umass.edu**

Recursive query

- Each name server involved contacts next server in hierarchy
- Puts burden of name resolution on contacted name server
- Heavy load?



5.3 DNS Records

DNS's distributed database stores resource records (RR)

- Format: (*name*, *value*, *type*, *ttl*)
- Type A (authoritative)
 - *name* is hostname
 - *value* is IP address
- Type AAAA
 - Same as A for IPv6
- Type NS (name server)
 - *name* is domain (e.g. tu-darmstadt.de)
 - *value* is hostname of authoritative name server of domain
- Type CNAME (canonical name)
 - *name* is alias for real name
 - *value* is canonical name
 - E.g. www.tu-darmstadt.de is really vwcms-live01.hrz.tu-darmstadt.de
- Type MX (mail exchanger)
 - *value* is name of mail server associated with *name*

Records are cached and updated

- Once a name server learns a mapping it caches it
- Cache entries time out after some time
- TLD servers typically cached in local name servers
 - Thus root name servers not often visited
- Update/notify mechanisms defined April 1997 in RFC 2136

Inserting DNS records

- Example: new startup “Company”
- Registers name company.com at DNS registrar
 - Accredited by Internet Corporation for Assigned Names and Numbers ICANN
 - In Germany e.g. Deutsche Telekom, 1&1
- Provides names and IP addresses of authoritative name server
- Registrar admin inserts two resource records into .com TLD server
 - (company.com, dns1.company.com, NS)
 - (dns1.company.com, 123.234.345.1, A)
- company.com admin creates at authoritative name server
 - Type A record for www.company.com
 - Type MX record for company.com

5.4 DNS Protocol

DNS protocol specifies query and reply messages

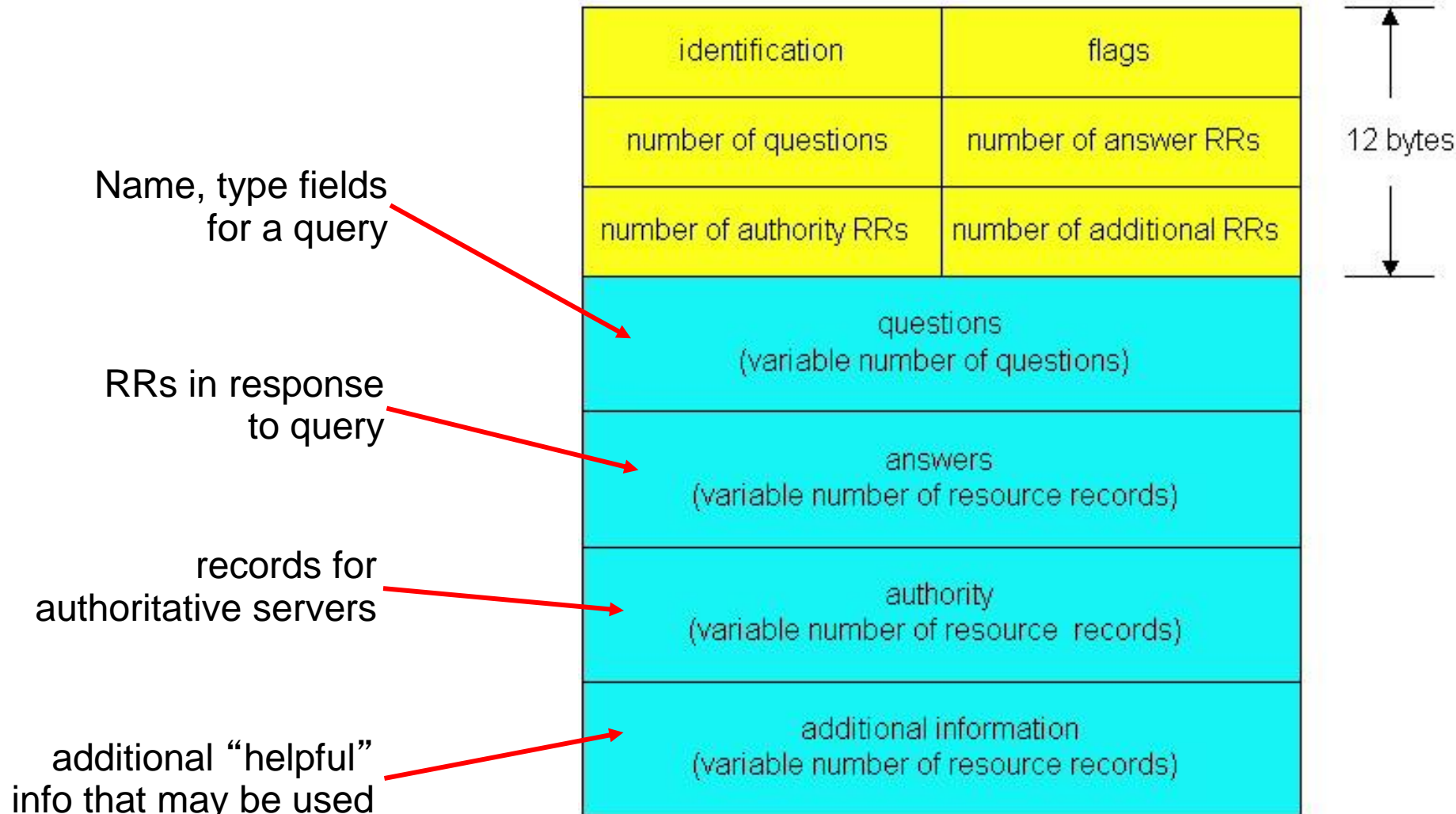
- Both with same message format

Message header

- Identification
 - 16 bit number generated for query
 - Reply to query uses same number
- Flags
 - Query or reply
 - Recursion desired
 - Recursion available
 - Reply is authoritative

identification	flags
number of questions	number of answer RRs
number of authority RRs	number of additional RRs
questions (variable number of questions)	
answers (variable number of resource records)	
authority (variable number of resource records)	
additional information (variable number of resource records)	

↑
12 bytes
↓



5.5 DNS Security Issues

DNS information is not authenticated

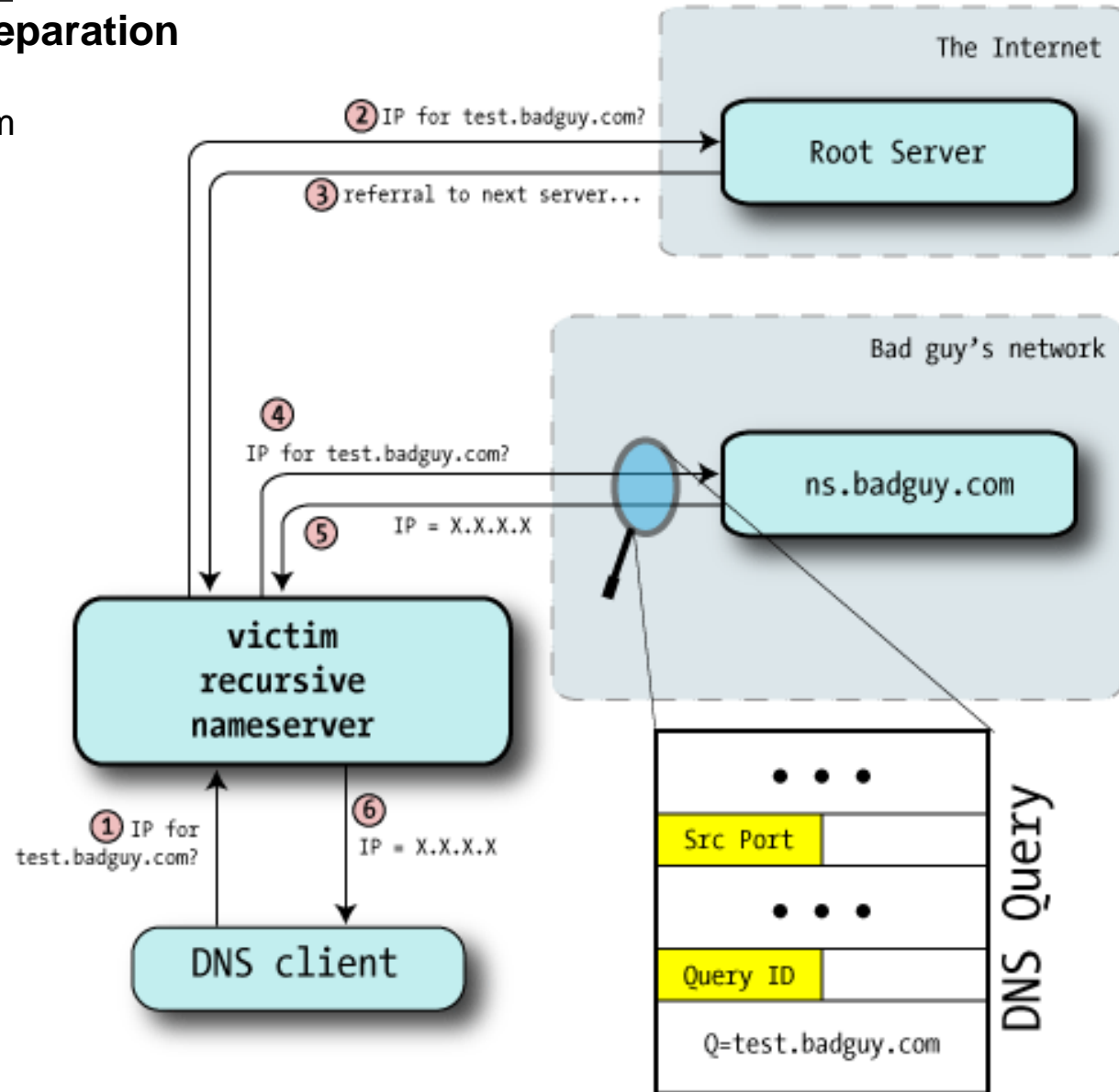
- In early days poisoning DNS caches was very easy
 - Attacker makes victim query for `www.attacker.com`
 - Query sent to attacker's name server (e.g. `ns.attacker.com`)
 - Attacker's name server replied with A record for `www.attacker.com`
 - Adds additional information (e.g. A record for `www.bank-of-victim.com`)
 - Victim trusts additional information
 - Victim's access to `www.bank-of-victim.com` redirected
- Today DNS implementations are more paranoid
 - E.g. additional information for foreign domains not accepted
- Injecting bogus ("malicious") information more complicated but still possible
- Authenticated DNS (DNSSEC) is yet to come
 - Established for root zone mid 2010
 - I.e. records for most top level domains can be validated

DNS Security Issues



DNS poisoning example - preparation

- Attacker wants to poison DNS entry for `www.bankofsteve.com`
- Challenge 1: DNS server only accepts authoritative replies
 - I.e. attacker has to blindly reply to legitimate request sent to authoritative name server of domain `bankofsteve.com`
- Challenge 2: DNS server only accepts replies with correct query ID
 - I.e. attacker has to guess query ID used
 - Assumption: query IDs are incremented linearly
 - Approach: learn query ID by forcing victim name server to send request to attackers name server

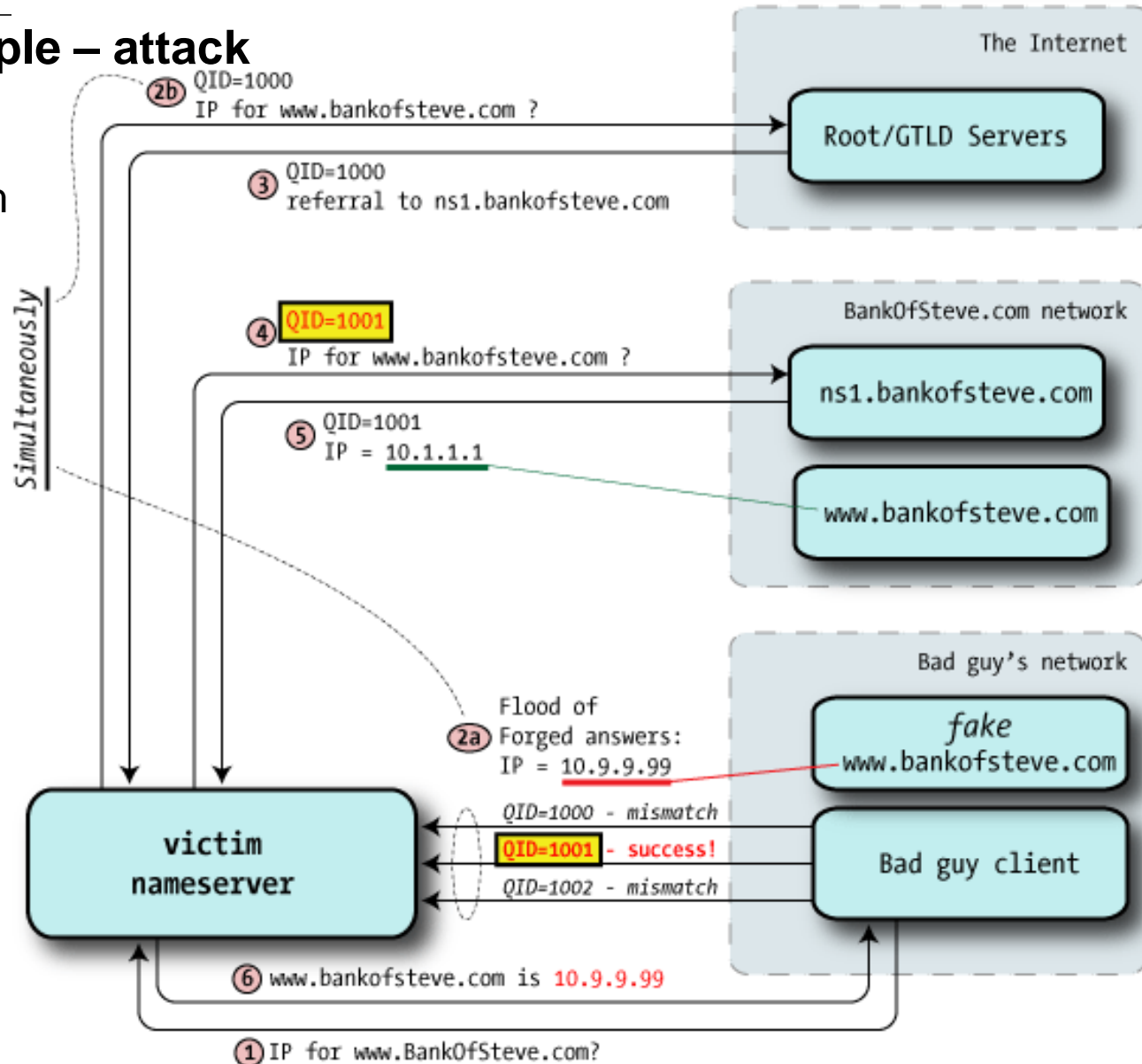


DNS Security Issues



DNS poisoning example – attack

- Attacker forces victim name server to query `www.bankofsteve.com`
- Attacker floods victim name server with faked replies
 - Using guessed query IDs
- Attacker succeeds if
 - Query ID guessed correctly
 - Reply arrives before legitimate reply



DNS poisoning – amplification

- So far: poisoned A record
 - One DNS reply will contain faked IP address
- One step further: poisoned NS record
 - All DNS requests will be handled by attacker's name server
 - Presented by Dan Kaminsky in 2008

DNS poisoning and https

- https uses digital certificates to uniquely bind IP address to URL
- But: www.bankofluxemburg.com vs. www.bankofiuxemburg.com ...

DNS poisoning – countermeasures

- Randomize query IDs to prevent guessing
 - But experiments showed chances are still good for $2^{16}=65536$ possible IDs
- Randomize query IDs and UDP port used to send query
 - But port translating NAT devices may nullify port randomization

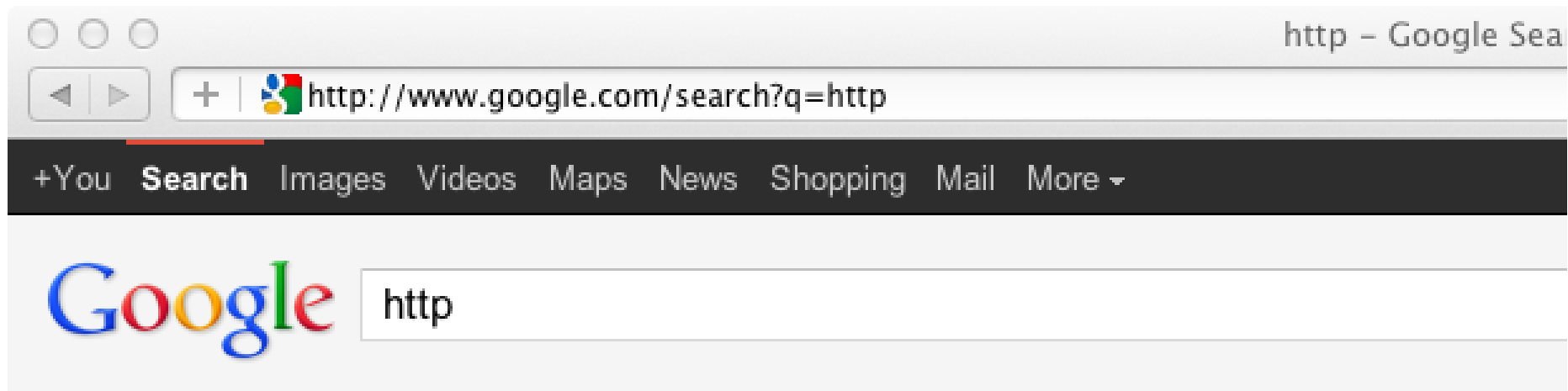
Recalling Licklider's visions

Obviously, collections of primary data can get too large to digest. There comes a time when the complexity of a communications process exceeds the available resources and the capability to cope with it; and at that point one has to simplify and draw conclusions.

6.1 Situation Today

Non-semantic web with a googol of pages

- Which are the relevant ones?
- Which are proved to be good?
- Which are still accessible?



Search

About 19,490,000,000 results (0.13 seconds)

Everything

[Hypertext Transfer Protocol - Wikipedia, the free encyclopedia](https://en.wikipedia.org/wiki/Hypertext_Transfer_Protocol)
en.wikipedia.org/wiki/Hypertext_Transfer_Protocol

“The Semantic Web

provides a common framework that allows DATA TO BE SHARED AND REUSED across application, enterprise, and community boundaries. It is a COLLABORATIVE EFFORT LED BY W3C with participation from a large number of researchers and industrial partners. It is based on the Resource Description Framework (RDF), which integrates a variety of applications using XML for syntax and URLs for naming.”

<http://www.w3.org/2001/sw/>

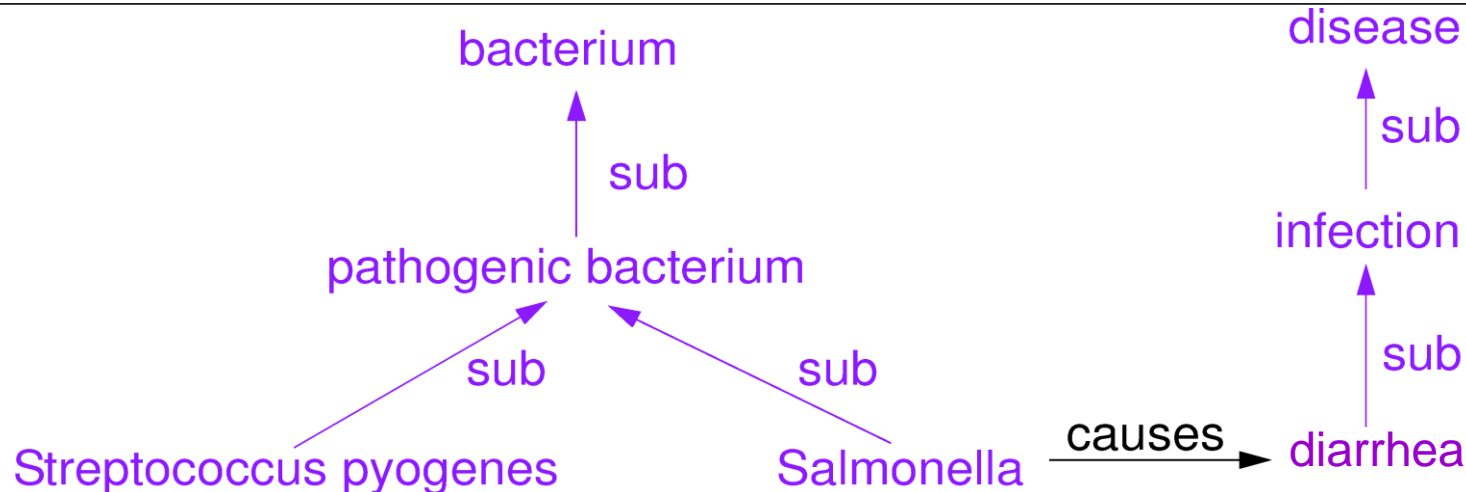
<http://www.semanticweb.org/>

”The Semantic Web

is an extension of the current web in which INFORMATION IS GIVEN WELL-DEFINED MEANING, better enabling computers and people to work in cooperation.”

**Tim Berners-Lee, James Hendler, Ora Lassila, The Semantic Web,
Scientific American, May 2001**

<http://www.scientificamerican.com/2001/0501issue/0501bern timers-lee.html>



Metadata = data about data

- Defined in a metadata scheme
- Clear vocabulary for entries

Ontology = a specification of a conceptualization, i.e. a shared conceptualization of a knowledge domain (Gruber)

- Concepts, instances as thematic entities of the knowledge domain
- Relations as semantic interconnections between concepts
 - Superconcepts - subconcepts, domain relations
 - Moreover: axioms, attributes, inference

E.g. semantic web as overlay on top of existing web resources

