

HTTP: the language of Web communication

Claudia Hauff

TI1506: Web and Database Technology

ti1506-ewi@tudelft.nl

Course overview [Web]

1. http: the language of Web communication

2. HTML & Web app design

3. JavaScript: interactions in the browser

4. node.js: JavaScript on the server

5. CSS: Lets make the app pretty

6. Ajax: asynchronous JavaScript

7. Personalization: Cookies & sessions

8. Securing your application

Source: <https://vimeo.com/110256895>

At the end of this lecture, you should be able to ...

- **Describe** how Web servers and clients interact with each other (via TCP/IP and HTTP)
- **Write** HTTP messages that request Web resources from Web servers and understand the responses
- **Describe** the different components of URLs and their purpose
- **Understand** and **employ** basic HTTP authentication
- **Explain** the difference between HTTP and HTTPS

World Wide Web
vs.
Internet

#<<<>>>
#copyright

Your [continued donations](#) keep Wikipedia running!

Lynx (web browser)

From Wikipedia, the free encyclopedia

Jump to: [navigation](#), [search](#)

CAPTION: Lynx

Wikipedia Main Page displayed in Lynx

Wikipedia Main Page displayed in Lynx

Maintainer: [Thomas Dickey](#)

Stable release: [2.8.5](#) (February 4, 2004) [\[\[+/-\]\]](#)

Preview release: [2.8.6](#) (?) [\[\[+/-\]\]](#)

OS: [Cross-platform](#)

Use: [web browser](#)

License: [GPL](#)

Website: [lynx.isc.org](#)

Lynx is a text-only [Web browser](#) and [Internet Gopher](#) client for use on cursor-addressable, character cell [terminals](#).

Browsing in Lynx consists of highlighting the chosen link using cursor keys, or having all links on a page numbered and entering the chosen link's number. Current versions support [SSL](#) and many [HTML](#) features. Tables are linearized (scrunched together one cell after another without tabular structure), while frames are identified by name and can be explored as if they were separate pages.

Lynx is a product of the Distributed Computing Group within Academic Computing Services of the [University of Kansas](#), and was initially developed in 1992 by a team of students at the university ([Lou Montulli](#), Michael Grobe and Charles Rezac) as a hypertext browser used solely to distribute campus information as part of a [Campus-Wide Information Server](#). In 1993 Montulli added an Internet interface and released a new version (2.0) of the browser [\[1\]](#) [\[2\]](#) [\[3\]](#).

The Web: a brief history

World Wide Web: a global system of interconnected hypertext documents available via the Internet
(envisioned already in 1945)



- **1960s:** Precursor to the Internet (ARPANET) devised by the US department of Defense
 - Initial services: electronic mail, file transfer
- Late **1980s:** Internet opened to commercial interests
- **1989:** WWW created by Tim Berners-Lee (CERN)
- **1994:** Netscape released its first Web browser
- **1995:** Microsoft released Internet Explorer v1
- **1998:** Google was founded
- **2002:** Mozilla released Firefox v1



Location: [Guided Tour](#) [What's New](#) [Questions](#) [Net Search](#) [Net Directory](#) [Newsgroups](#)

Mosaic Netscape version 0.9 beta

Copyright © 1994 Mosaic Communications Corporation,
All rights reserved.

This is *BETA* software subject to the license agreement set forth in the README file.
Please read and agree to all terms before using this software.

Report any problems to win_cbug@mcom.com.



Mosaic Communications, Mosaic Netscape, and the Mosaic Communications logo are trademarks of Mosaic Communications Corporation.

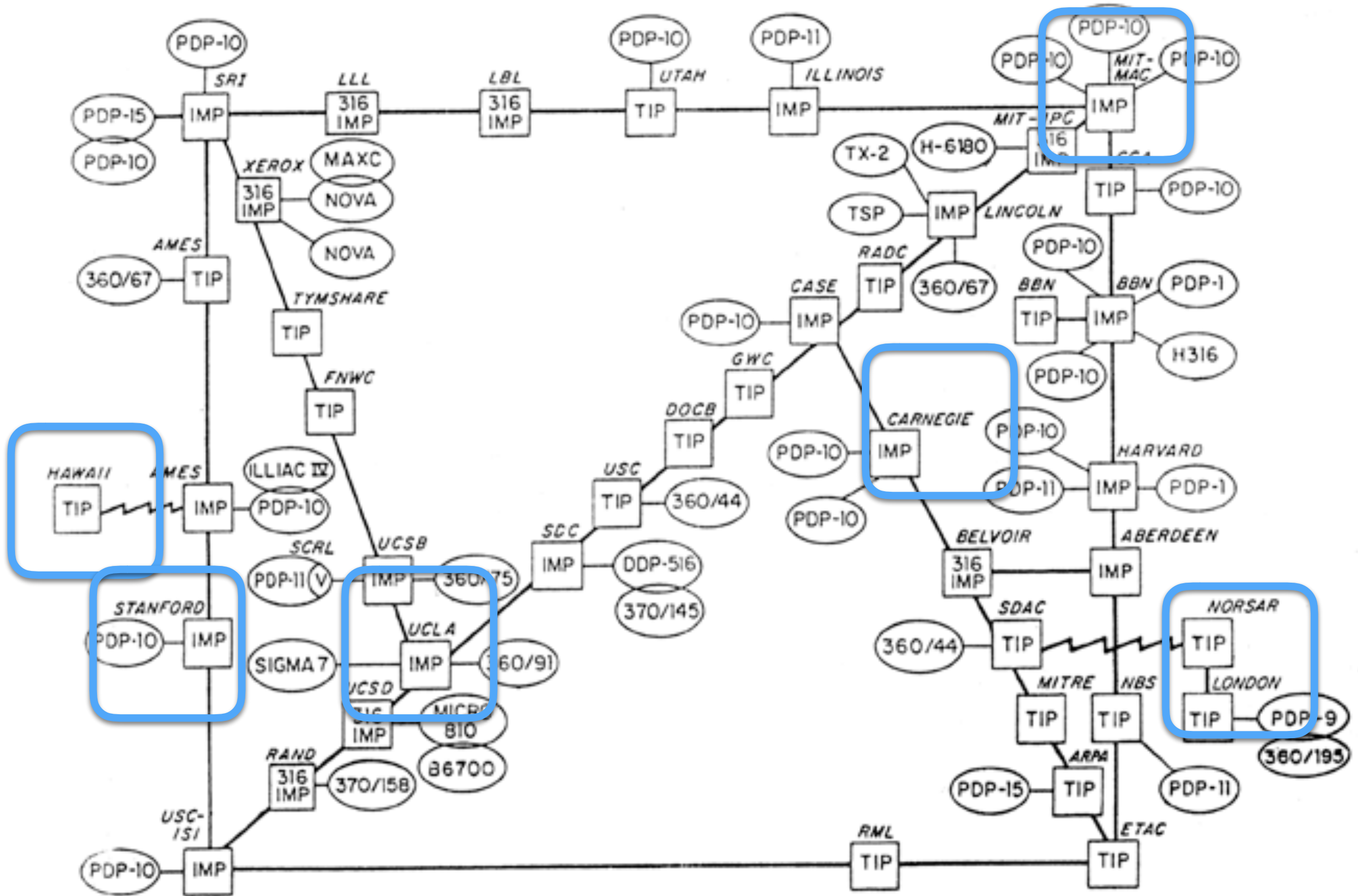
Any provision of Mosaic Software to the U.S. Government is with "Restricted rights" as follows: Use, duplication or disclosure by the Government is subject to restrictions set forth in subparagraphs (a) through (d) of the Commercial Computer Restricted Rights clause at FAR 52.227-19 when applicable, or in subparagraph (c) (1) (ii) of the Rights in Technical Data and Computer Software clause at DFARS 252.227-7013, and in similar clauses in the NASA FAR Supplement. Contractor/manufacturer is Mosaic Communications Corporation, 650 Castro Street, Suite 500, Mountain View, California, 94041.

Key aspects of the Internet

Internet: interconnected computer networks (sub-networks) that span the globe; communicating through a common standard (TCP/IP)

- Sub-networks function **autonomously**
- **No centralised** (global) control instance
- Devices **dynamically** join and leave the network
- Devices interact through **agreed-upon open standards**; anyone can create a new device
- **Easy** to use: server/client software is widely available

ARPA NETWORK, LOGICAL MAP, SEPTEMBER 1973



State of the Internet in 1973

Image source: <http://qz.com/860873/a-1973-map-of-the-internet-charted-by-darpa/>

Two important organisations

- Internet Engineering Task Force (IETF)

“The mission of the IETF is to make the Internet work better by producing high quality, relevant technical documents that influence the way people design, use, and manage the Internet. “

Request for Comments (RFC)

- World Wide Web Consortium (W3C)

“The W3C mission is to lead the World Wide Web to its full potential by developing protocols and guidelines that ensure the long-term growth of the Web.”

An introduction to HTTP messages



HTTP 1.1

RFC 2068

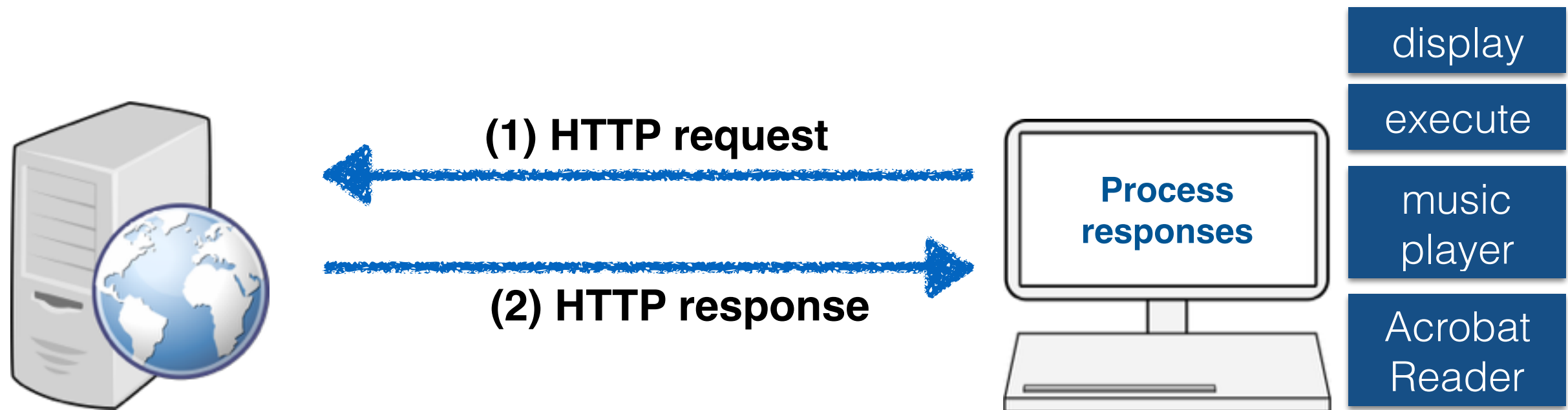
1997

HTTP/2

RFC 7540

2015

Web servers and clients



- Servers wait for data requests
- Answer thousands of clients simultaneously
- Host **web resources**
- Clients are most often Web browsers
- Telnet

Web resource: any kind of content with an identity, including static files (e.g. text, images, video), software programs, Web cam gateway, etc.

Network communication

- Conceptual model **Open Systems Interconnection** (OSI) from 1995
- Network protocols are matched to different “layers”
- Model still in use, but not the associated protocols
- Many network protocols exist, three are of interest to us:

IP: Internet Protocol

TCP: Transmission Control Protocol

HTTP: Hypertext Transfer Protocol

others include: SSH, IMAP, STMP, FTP

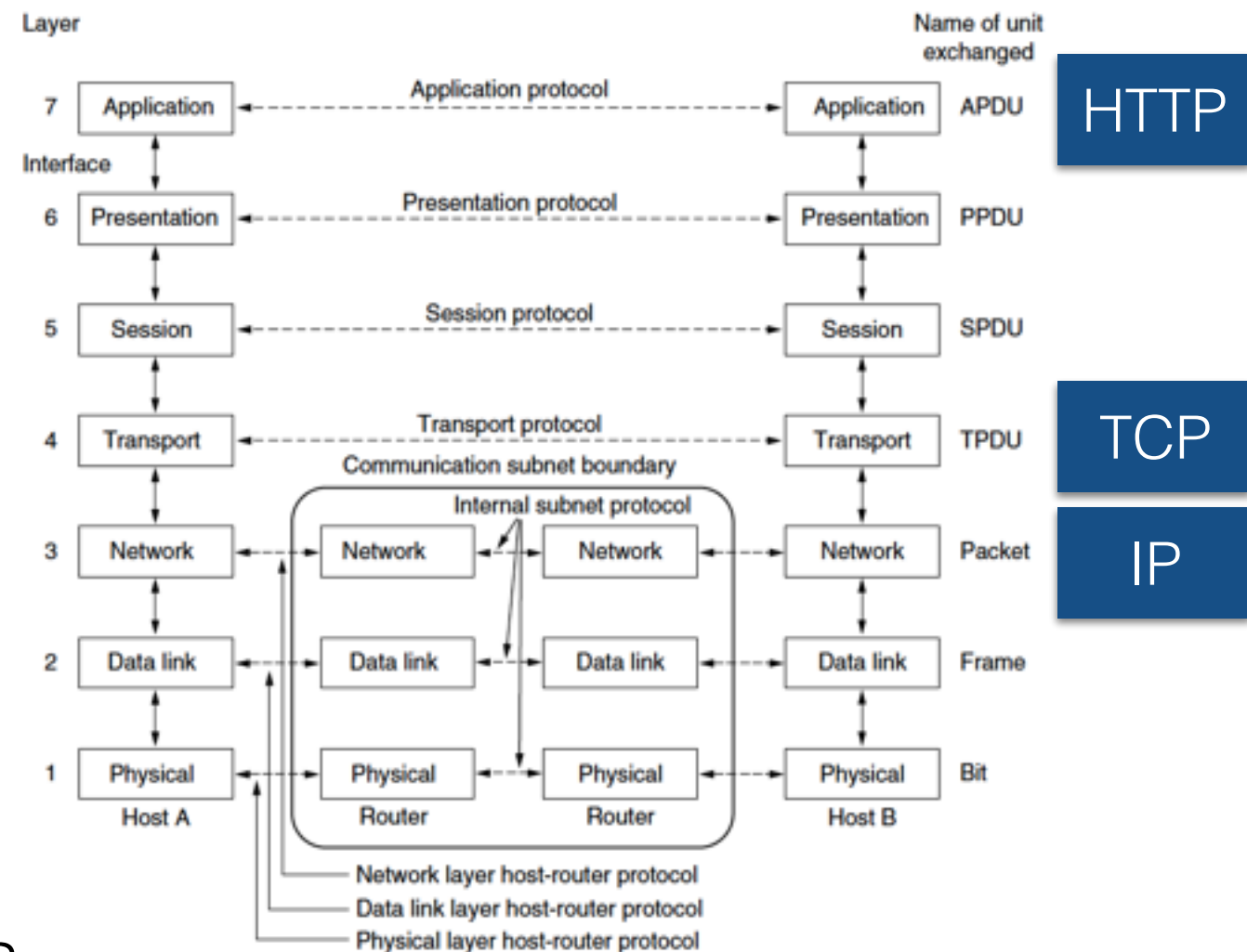


Image src: Computer Networks (5th edition), Tanenbaum & Whetherall, p. 42

HTTP uses **reliable** data-transmission protocols.

Web Developer Tool Demo

Insights

Modern web sites consist of many resources.

A cascade of HTTP transactions is required.

Different MIME types are typically found.

Development tools of all major browsers provide a lot of information!

HTTP request message

plain text, line-oriented character sequences

```
GET / HTTP/1.1
Host: www.tudelft.nl
User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X
10.9; rv:31.0) Gecko/20100101 Firefox/31.0
Accept: text/html,application/xhtml+xml,application/
xml;q=0.9,*/*;q=0.8
Accept-Language: en-gb,en;q=0.5
Accept-Encoding: gzip, deflate
DNT: 1
Cookie:
__utma=1.20923577936111.16111.19805.2;utmcmd=(none);
```


HTTP response message

```
HTTP/1.1 200 OK
```

```
Date: Fri, 01 Aug 2014 13:35:55 GMT
```

```
Content-Type: text/html; charset=utf-8
```

```
Content-Length: 5994
```

```
Connection: keep-alive
```

```
Set-Cookie: fe_typo_user=d5e20a55a4a92e0;  
path=/; domain=tudelft.nl
```

```
[...]
```

```
Server: TU Delft Web Server
```

```
....
```

```
....
```

start line

header fields

name:value

body
(optional)

HTTP headers dissected

Primary entity header fields

Content-Type	Entity type
Content-Length	Length/size of the message
Content-Language	Language of the entity sent (e.g. English)
Content-Encoding	Data transformations applied to the entity
Content-Location	Alternative location of the entity
Content-Range	For partial entities, range defines the pieces sent
Content-MD5	Checksum of the content
Last-Modified	Date on which this entity was created/modified
Expires	Date at which the entity will become stale
Allow	Lists the legal request methods for the entity

Important: Entity bodies only contain **raw** data, **header** information required to **interpret** the data.

Content-Type

- **MIME** types are attached to all HTTP object data

Multipurpose Internet **Mail** Extensions

historic reasons

- MIME type determines clients reaction to the received data

- Pattern: [primary object type]/[subtype]

e.g. `text/plain`, `text/html`, `image/jpeg`,
`video/quicktime`, `application/vnd.ms-`
`powerpoint`

Content types are diverse

Most popular
<code>text/html</code>
<code>image/jpeg</code>
<code>text/xml</code>
<code>application/rss+xml</code>
<code>text/plain</code>
<code>application/xml</code>
<code>text/calendar</code>
<code>application/pdf</code>
<code>application/atom+xml</code>
<code>unknown/unknown</code>

Least popular
<code>application/pgp-keys</code>
<code>application/x-httpd-php4</code>
<code>chemical/x-pdb</code>
<code>model/mesh</code>
<code>application/x-perl</code>
<code>audio/x_mpegurl</code>
<code>application/bib</code>
<code>application/postscript</code>
<code>application/x-msdos-program</code>

Content-Length

- Indicates the **size** of the entity body in the message
- Necessary to detect premature message truncation (e.g. due to a server crash, faulty proxy)
- Essential for **persistent connections** to discover where one HTTP message ends and the next one begins

Persistent connections reuse the same TCP connection for multiple HTTP request/response messages.

Content-Encoding

- Commonly either gzip, compress (Unix compression), deflate (zlib compression) or identity (no encoding)
 - Servers aim to use **encodings** that clients understand
 - Clients send a list of acceptable encodings in the `Accept-Encoding` request header
- ```
Accept-Encoding: gzip, deflate
```
- **Compression** saves network bandwidth (fewer bits to transmit) but increases processing costs to decompress

# Content-MD5

Message Digest

- HTTP messages are sent via TCP/IP (ensuring reliable transport)
- **But**: the Internet is huge, many servers interact to transport a message with different implementations of established protocols (which may be **buggy**)
- **Sanity check**: Sender generates a **MD5 checksum** of the content (hashed into a 128 bit value) to detect unintended modifications of the content

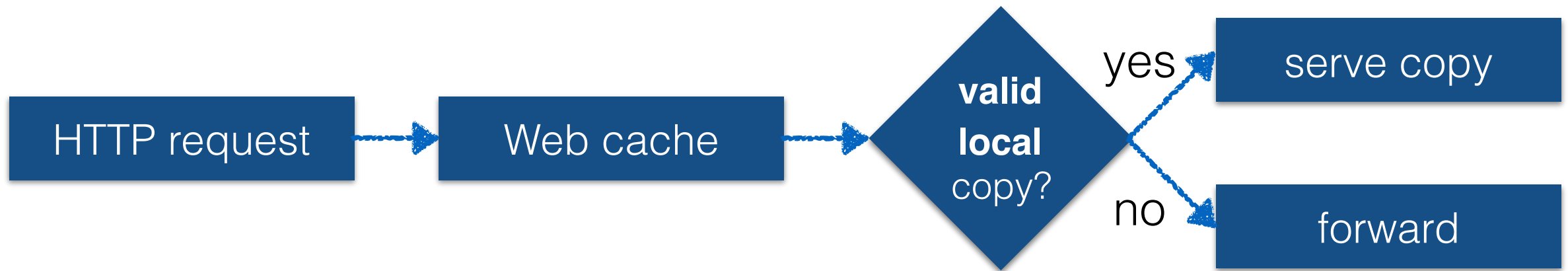


**MD5 is cryptographically broken**



# Expires

- **Web caches** keep copies of *popular* documents



- Advantages of Web caches
  - A. **Reduction** of redundant data transfer
  - B. **Reduction** of network bottlenecks
  - C. **Reduction** demand on origin servers
  - D. **Reduced** distance delay
- **Expires** indicates when fetched resource is no longer valid and needs to be retrieved from the origin server

# Expires & Cache-Control

- Content on the origin server can change
- Caches need to ensure that their copies are **in sync** with the origin server
- Caches can revalidate their copies at any time (**inefficient**)
- `Expires` in HTTP response header indicates a **document's expiration date** in **absolute** terms — date determines when the cache revalidates
- `Cache-Control`: indicates a document's expiration date in **relative** terms (number of seconds since being sent)

# Last-Modified

- Contains the date when the document was last **altered** (in HTTP response)
- No indication about the amount of changes in the document
- Often used in combination with **If-Modified-Since** for cache revalidation requests — origin server only returns the document if it changed since the given date (otherwise *304 - Not Modified*)
- Last-Modified dates are **not reliable**

# Remember: HTTP response message

```
HTTP/1.1 200 OK
```

start line

```
Date: Fri, 01 Aug 2014 13:35:55 GMT
```

```
Content-Type: text/html; charset=utf-8
```

```
Content-Length: 5994
```

```
Connection: keep-alive
```

```
Set-Cookie: fe_typo_user=d5e20a55a4a92e0;
path=/; domain=tudelft.nl
```

```
....
```

```
Server: TU Delft Web Server
```

header fields

name: value

```
....
```

body  
(optional)

# Common status codes

|            |                              |
|------------|------------------------------|
| <b>1xx</b> | Informational                |
| <b>2xx</b> | Success (200 OK)             |
| <b>3xx</b> | Redirected                   |
| <b>4xx</b> | Client error (404 Not Found) |
| <b>5xx</b> | Server error                 |

In practice only a few codes per category are supported

## **10.4.3 402 Payment Required**

This code is reserved for future use.

[A more detailed overview.](#)



# HTTP methods

GET / HTTP/1.1

Host: www.tudelft.nl

User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10.9; rv:31.0) Gecko/20100101 Firefox/31.0

Accept: text/html,application/xhtml+xml;q=0.9,\*/\*;q=0.8

Accept-Language: en-gb,en;q=0.5

Accept-Encoding: gzip, deflate

# Common HTTP methods

|                |                                                         |
|----------------|---------------------------------------------------------|
| <b>GET</b>     | Get a document from the Web server.                     |
| <b>HEAD</b>    | Get the header of a document from the Web server.       |
| <b>POST</b>    | Send data from the client to the server for processing. |
| <b>PUT</b>     | Save the body of the request on the server.             |
| <b>TRACE</b>   | Trace the message through proxy servers to the server.  |
| <b>OPTIONS</b> | Determine what methods can operate on a server.         |
| <b>DELETE</b>  | Remove a document from a Web server.                    |

Servers may implement more or fewer methods than shown.

# Demo: telnet

Telnet opens a **TCP connection** to a Web server; chars are typed directly into the port. The server treats telnet as Web client, the returned data is displayed onscreen.

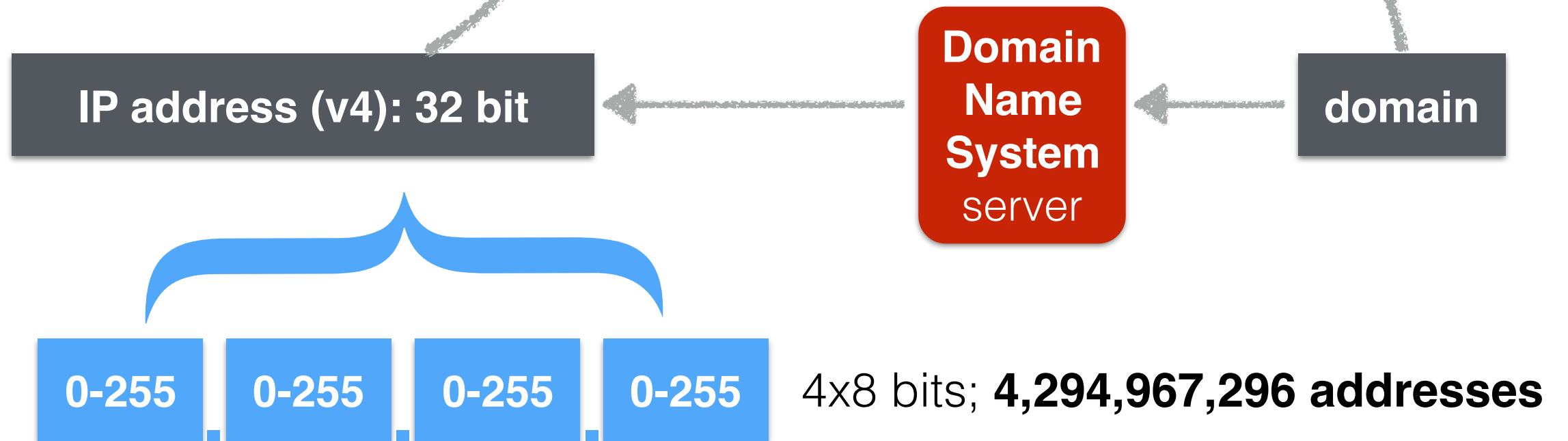
A number of HTTP requests may be required to end up at the final (wanted) page.

Often Web servers treat Web-browser requests differently from machine-generated requests.

# From domain to IP address

```
claudiahauff@Cludias-MacBook-Air:~ $ telnet microsoft.com 80
Trying 134.170.185.46...
Connected to microsoft.com.
Escape character is '^]'.

```



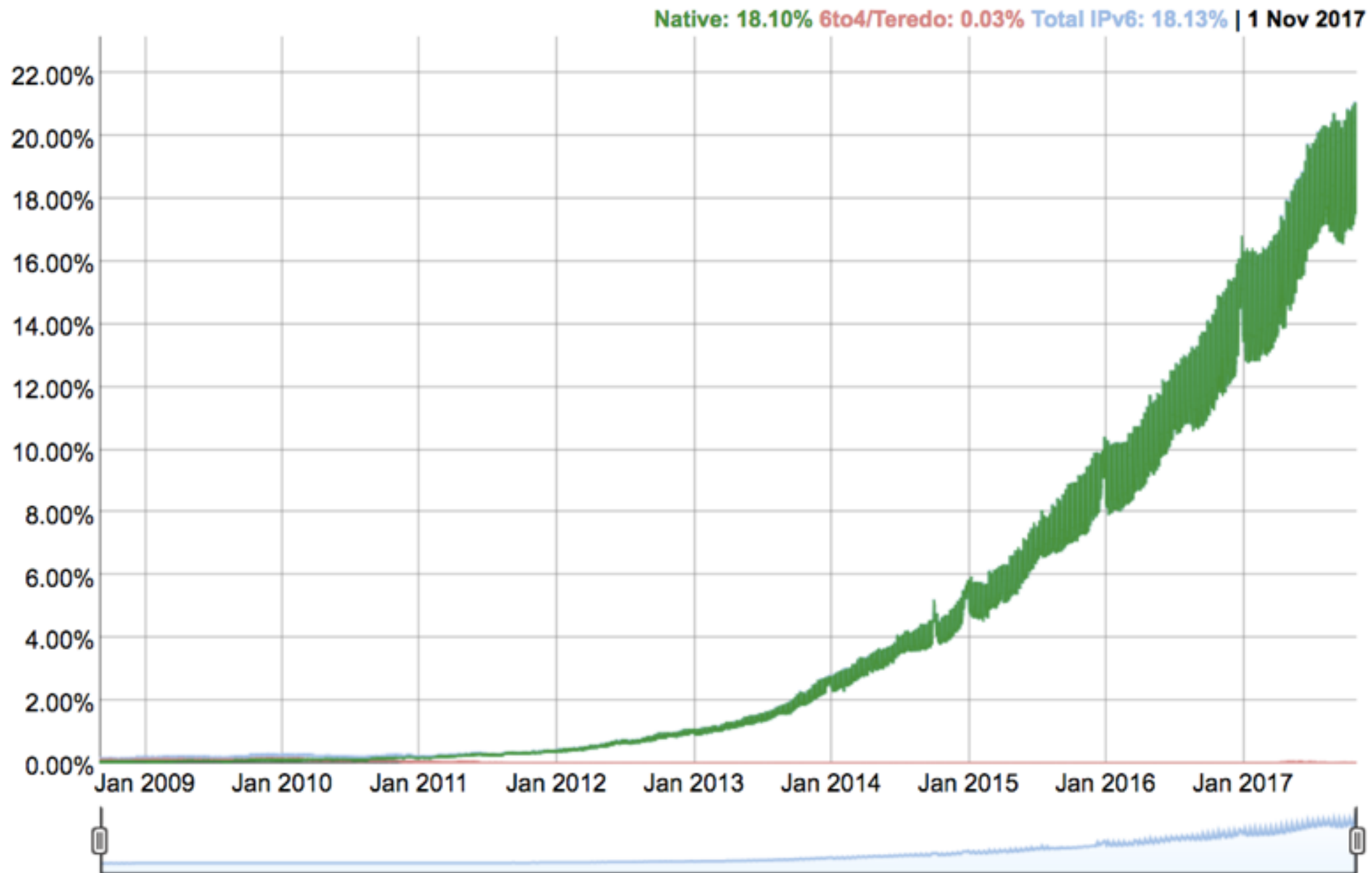
IP address (v6): 128 bit

2A41 0000 3341 3AAA 0000 FFFF 86AA B92E

# From domain to IP address

## IPv6 Adoption

We are continuously measuring the availability of IPv6 connectivity among Google users. The graph shows the percentage of users that access Google over IPv6.





# HTTP 1.1

RFC 2068

1997

# HTTP/2

RFC 7540

2015

# HTTP/2 in one slide

based on **SPDY**, a protocol developed at Google

“The primary goals for HTTP/2 are to reduce latency by enabling **full request and response multiplexing**, minimize protocol overhead via **efficient compression of HTTP header fields**, and add support for **request prioritization** and **server push**.”

Quote: <https://hpbn.co/http2/>

**HTTP/1.1: concurrency achieved through multiple connections**

**HTTP/1.1: no header compression**

**HTTP/1.1: no prioritization of requests**

**HTTP/1.1: everything works on a request basis**

- Find a **lab partner** as soon as possible!  
Know your team's assessment **cluster**!
- Work through **Chapter 2** (intro to HTML) of the Web development book **before** this Thursday's lecture!
- Start working on this lab **assignment.**