



WWW (World Wide Web) Basics

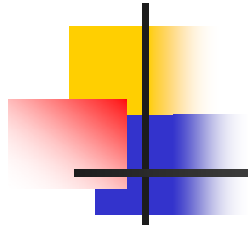
BUPT/QMUL
2010-12-07



Agenda

- Brief introduction to WWW
- WWW Components
- WWW Standards
- Web Applications
- Summary

Refer to Chapter 27, textbook



Brief Introduction To WWW

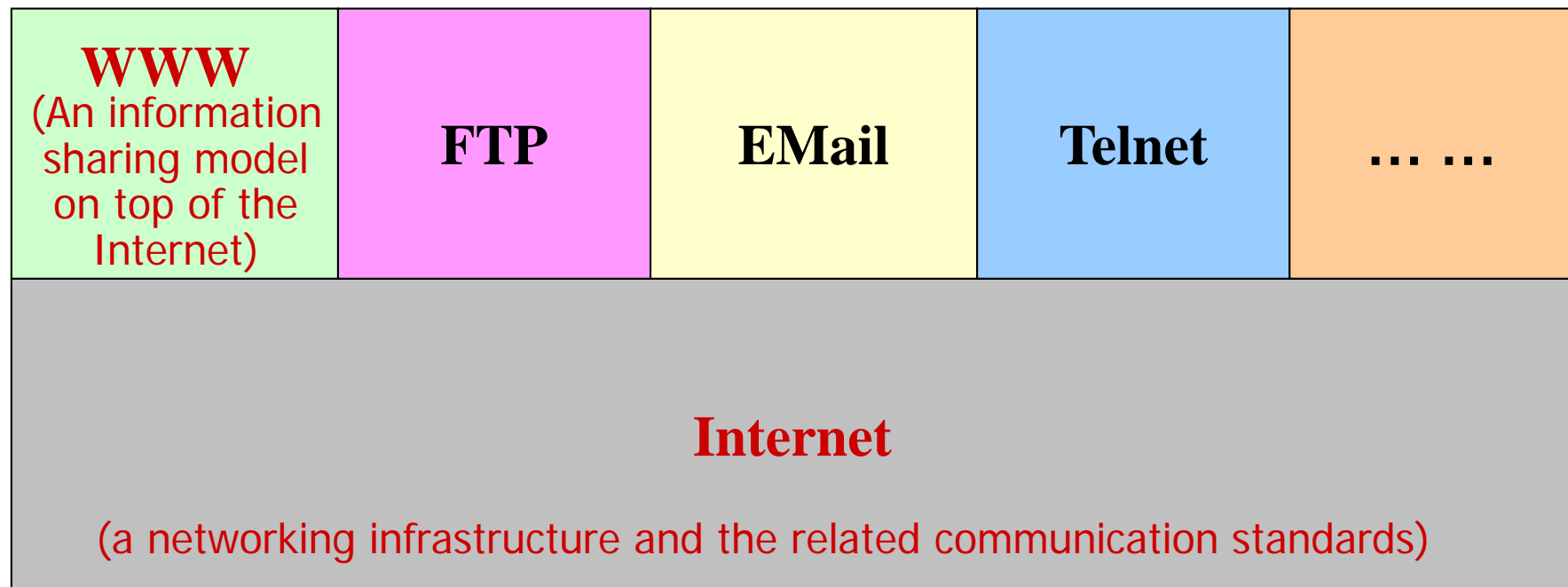


What Is WWW?

- World Wide Web
 - WWW, the Web, W3
- A technical definition
 - All the [resources and users](#) on the Internet that are using the Hypertext Transfer Protocol ([HTTP](#)).
 - A system of [interlinked hypertext](#) documents accessed via the Internet. – *Wikipedia*
- A broader definition from W3C (World Wide Web Consortium)
 - The World Wide Web is the [universe of](#) network-accessible [information](#), an embodiment of human knowledge.



WWW vs. Internet





History Of WWW

CERN

The world's largest
particle physics laboratory

... where the web was born!

1989-03, Tim Berners Lee

proposed the idea of sharing information through
hypertext in CERN

1989-12 , Tim Berners Lee

named his invention WWW (World Wide Web)

1990-11

The first (text-based) prototype was operational

1991-12

The first public demonstration was given at
Hypertext ' 91 in San Antonio - Texas

1993-02, Marc Andreessen

The first GUI browser – Mosaic, at NCSA, Illinois

1994-95, Netscape, Microsoft

Netscape Navigator, Internet Explorer

Other browsers

Mozilla, Opera, Lynx, ELinks, Safari, ...

Other technologies

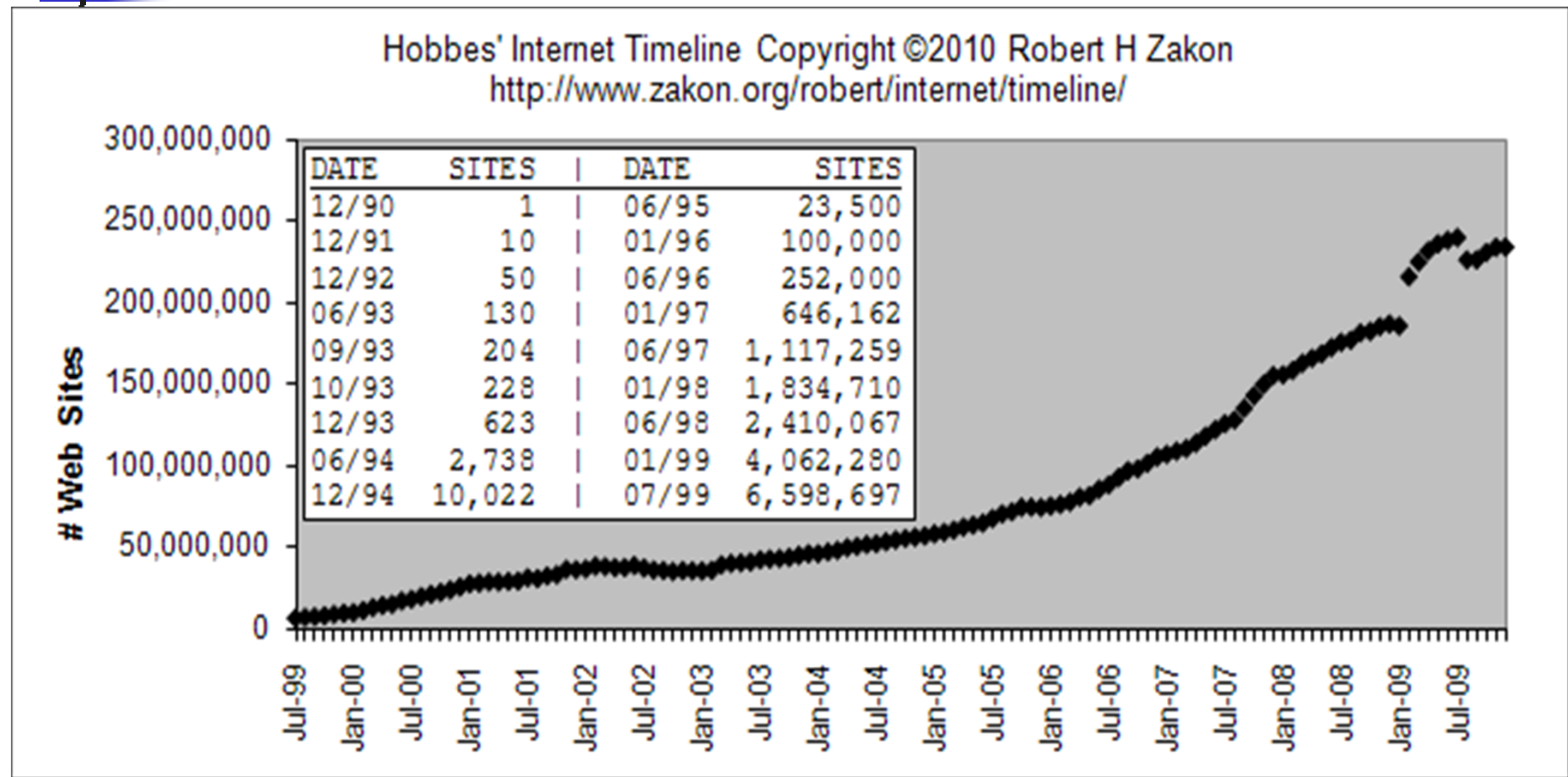
HTML, JAVA, VRML, Web 2.0, ...



Features of WWW

- Global
- Open
- Interactive
- Dynamic
- Platform-independent
- Multimedia
- ...

WWW Growth





WWW Terminologies

- **The Web**

- Is a true information superhighway

- **URL** (Uniform Resource Locator)

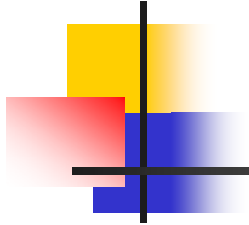
- Designates a specific webpage on a specific webserver

- **HTTP** (HyperText Transfer Protocol)

- An application-level transfer protocol standard

- **HTML** (HyperText Markup Language)

- A document format standard



WWW Components



WWW Components

- Structural Components

- Clients/browsers – various implementations
- Servers – run on sophisticated hardware
- Caches – used to improve response time
- Internet – the global infrastructure which facilitates data transfer

- Semantic Components

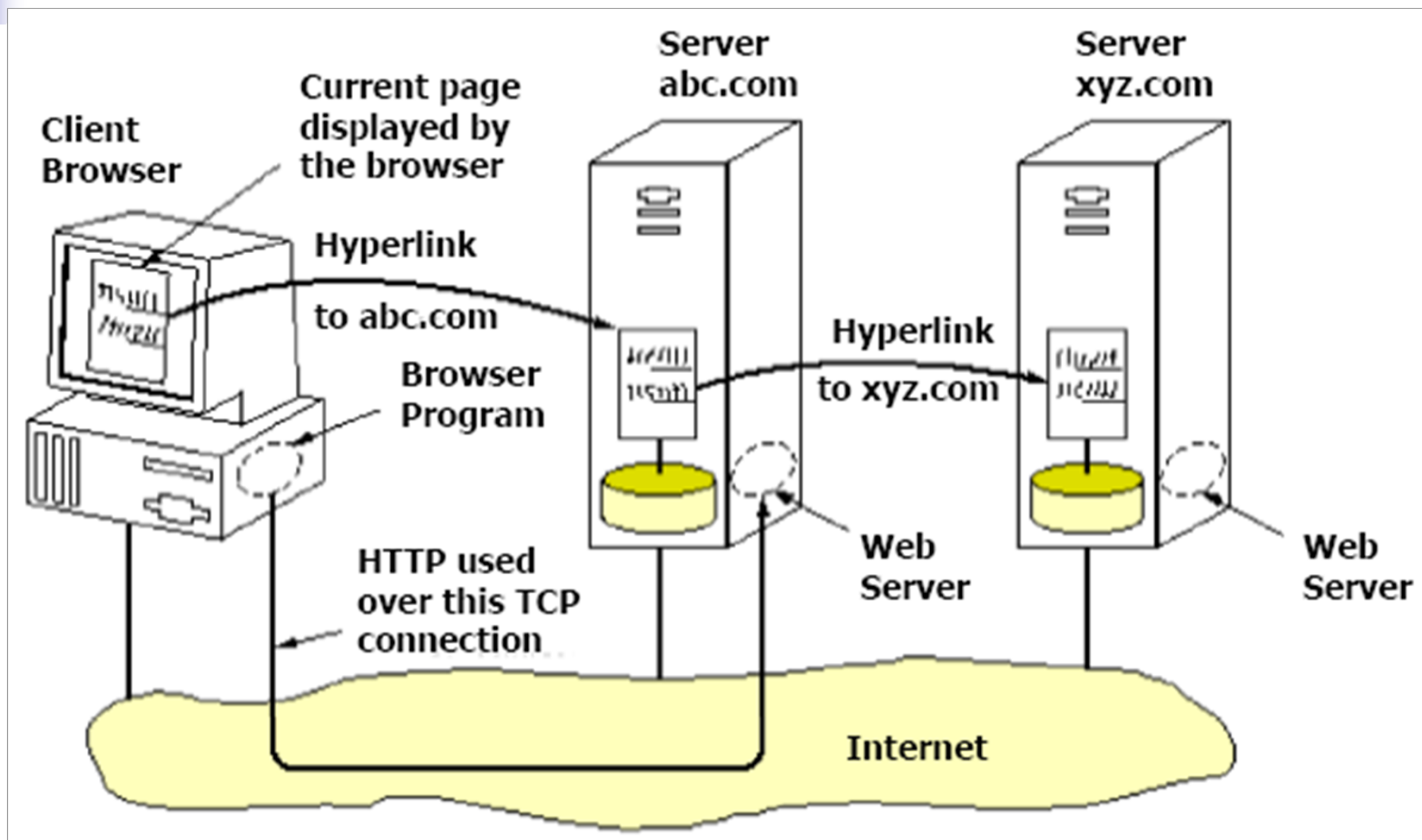
- Hyper Text Transfer Protocol (HTTP)
- Hyper Text Markup Language (HTML)
 - eXtensible Markup Language (XML)
- Uniform Resource Locators (URLs)
 - Uniform Resource Identifiers (URIs)

The Web



- The Web is actually an information superhighway
- The Web is a collection of electronic documents that are linked together like a spider web
- The Web is basically an information system that links data from many different Internet services under one set of protocols
- **Web clients**, also called **browsers**, **interpret HTML** delivered from **Web servers**
- These documents use **hypertext links** to connect different **documents** and information **resources** together; click on a link and the client software retrieves the linked document or jumps to a specific position in the current document
- **HTTP** is easily modified to incorporate **new data formats** and uses
- The Web model successfully unites the diverse Internet resources under a **single** system, relying on servers and Web-browsers to “**negotiate**” or handle data compatibility

The Web Access Model



WWW Clients

- The Web is designed like all the client/server applications
 - The **client** is called a “**browser**”
 - The **server** is where the **data is stored** and it is software that runs on **well known port (80)** ... usually
- The browser and server talk using a protocol – **HTTP**
- We already know from past experience that this architecture gives us client options
 - Netscape, Internet Explorer, Maxthon, Mozilla, Firefox, Lynx, ...



Firefox

Web Browsers Statistics(1)

Browser Statistics Month by Month

<http://www.w3schools.com>

2010	Internet Explorer	Firefox	Chrome	Safari	Opera
October	29.7 %	44.1%	19.2%	3.9%	2.2%
September	31.1 %	45.1%	17.3%	3.7%	2.2%
August	30.7 %	45.8%	17.0%	3.5%	2.3%
July	30.4 %	46.4%	16.7%	3.4%	2.3%
June	31.0 %	46.6%	15.9%	3.6%	2.1%

Site Usage

Goal Set 1

Google Analytics - Browsers re

Visits

88,323

% of Site Total: 100.00%

Pages/Visit

1.62

Site Avg: 1.62 (0.00%)

Avg. Time on Site

00:01:36

Site Avg: 00:01:36 (0.00%)

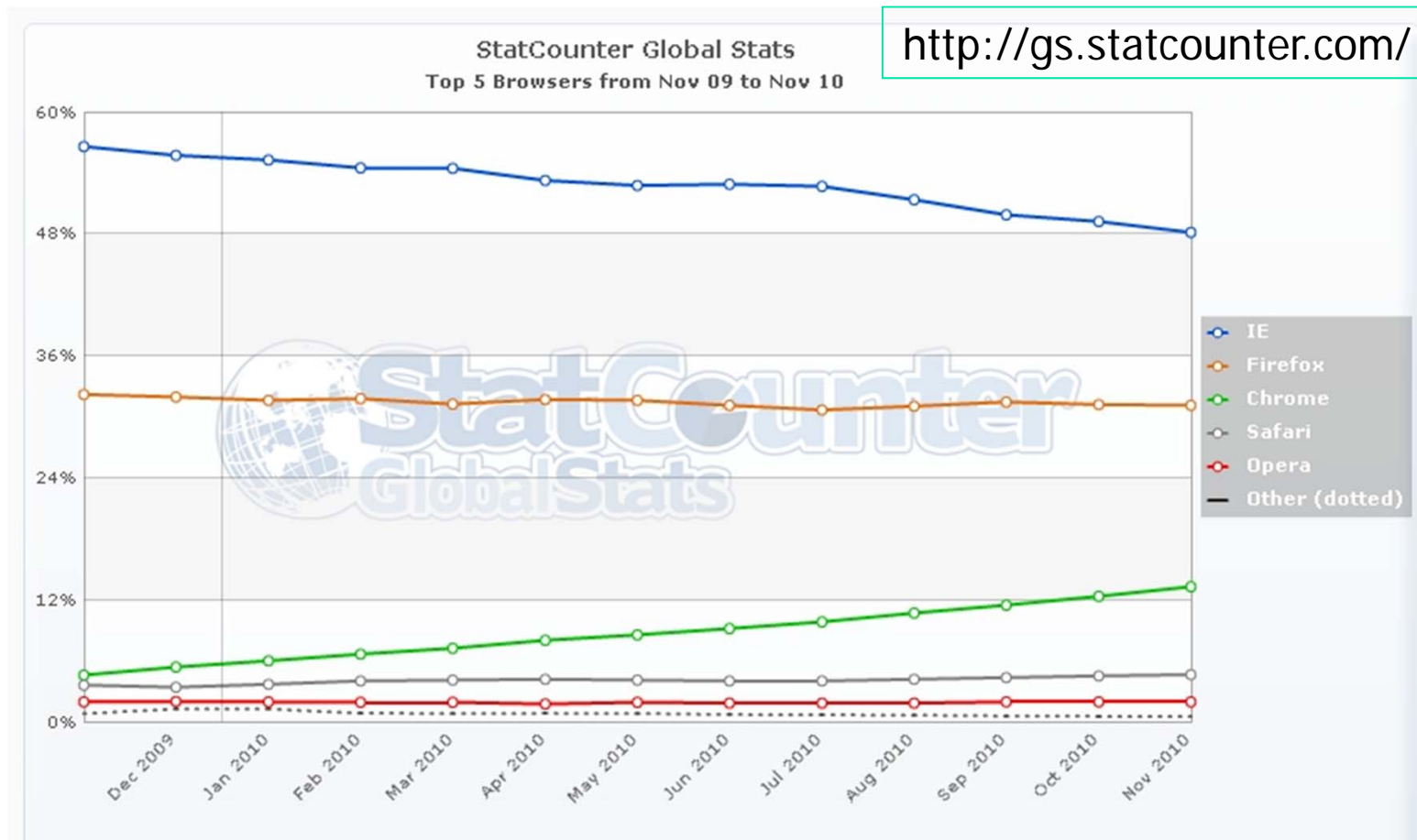
	Browser	Visits	Visits
1.	Internet Explorer	58,686	66.44%
2.	Firefox	20,655	23.39%
3.	Safari	6,370	7.21%
4.	Chrome	1,511	1.71%
5.	Opera	599	0.68%
6.	Mozilla	255	0.29%
7.	Opera Mini	38	0.04%
8.	Camino	32	0.04%
9.	Playstation Portable	32	0.04%
10.	Netscape	31	0.04%

http://devon.freepgs.com/imgv10

/GA_Browsers_2009.png

http://devon.freepgs.com/imgv10/GA_Browsers_2009.png

Web Browsers Statistics(2)



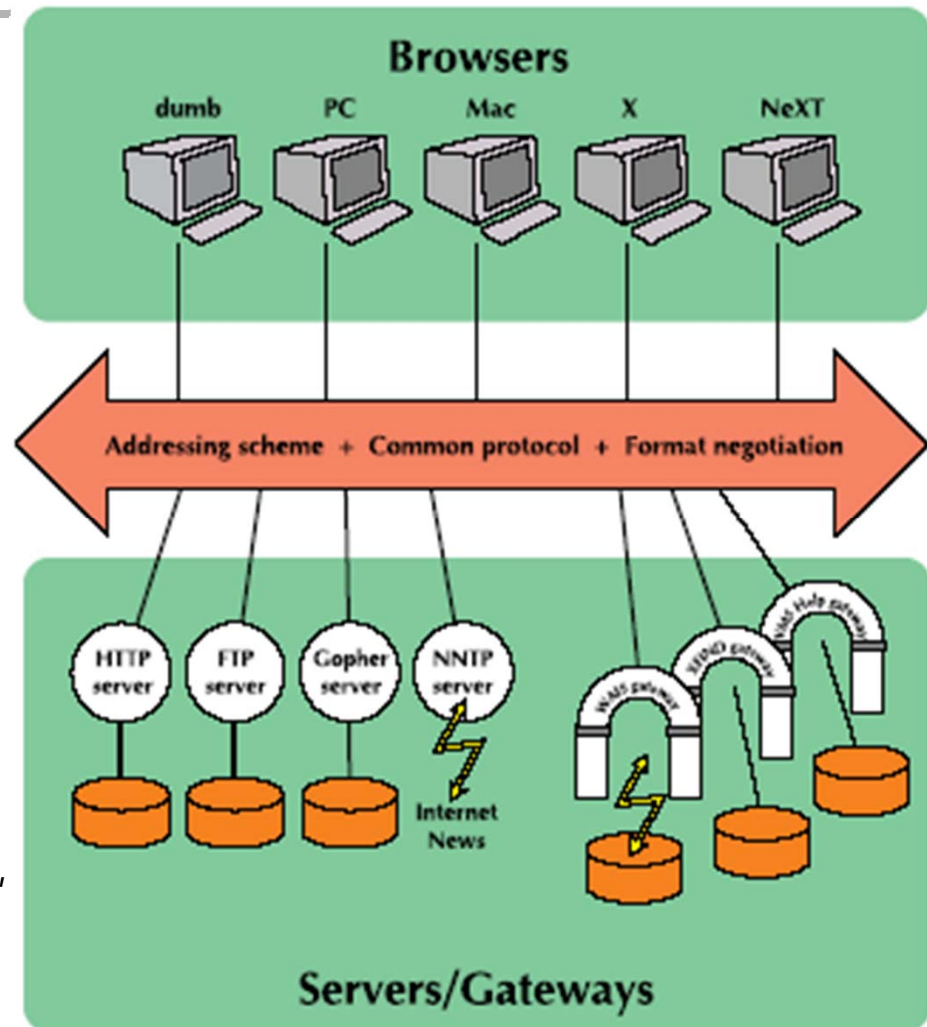


Basic Client Properties (1)

- All the different browsers show us the **same information** but they **display it differently** (depending on their **capabilities**)
- In front of each Web address there is an **http://** to indicate to the browser that it is talking **HTTP**, the protocol of the Web
- A user on a client machine uses a browser to download a Web page by either entering either a **URL** or clicking on a **HyperLink**

Basic Client Properties (2)

- Web browsers are often called **Universal Clients** because most can talk other protocols besides HTTP
 - ftp://home.domain: to use our Web browser as an FTP client
 - telnet://home.domain
 - gopher://host.domain
 - ...
- The Web is capable of accessing data on many different Internet services:
 - Web pages, FTP, Email service, Gopher menus, file directories, Wide Area Information Service (WAIS) databases, Finger Services, UseNet, Telnet services, HTML, plain ASCII, etc.





WWW Servers (1)

- The server is **software** that is running on a remote location. Its job is to **make “pages” available** to the client - so when a client requests a page the server responds appropriately
- Web servers are typically on Unix or Windows NT boxes rather than on individual PCs
- Popular Web Servers:
 - On Unix - Apache, On Windows NT - IIS (Internet Information Server), Both - Netscape's Web Server



WWW Servers (2)

- Every Web site has a server process listening to **TCP** port **80** for incoming connections from clients – normally browsers
- After a connection has been established, the client sends one **request** and the server sends one **response**
- Then the **connection is released**
- The protocol that defines the legal request and response is **HTTP**
- The operation is **Stateless**



URLs (Uniform Resource Locators)

- The global address of a Web page is described by its URL
- URLs identify the **protocol** you want to talk, the **site** (domain name or IP Address) you want to go to, and possible the **item** you want to see
- They have the form:
 - **protocol://hostname [:port]/directory/item-you-want**

**Resources can be dynamically-generated query results
(see POST command later)**



Structure Of URLs

- A URL consists of three parts:
 - The protocol – for example **http** or **ftp**
 - The DNS name of the host
 - The directory and file name



- Protocol: `http` by default
- Port: `80` by default
- `Index.html`, `index.htm`, `default.htm`, `default.asp` etc. are assumed if no file-name given



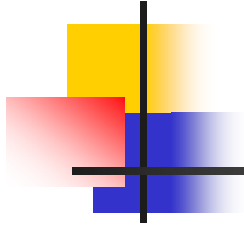
Some URLs Examples

Protocol	Use	Example
http	Web pages	http://www.elec.qmul.ac.uk
ftp	File transfer	ftp://elec.qmul.ac.uk/pub/info.doc
file	Local files	file:///D:/src/multim/filter.txt
news	News	news://comp.sys.os.linux
gopher	Gopher	gopher://gopher.tc.umn.edu/11/lib
mailto	E-mail	mailto:cip@elec.qmul.ac.uk
telnet	Remote login	telnet://www.elec.qmul.ac.uk



Other Related Terminologies

- **URI** (Uniform Resource Identifier)
- **URN** (Uniform Resource Name)
- What's the relationship between URI, URL and URN?
 - See RFC 3305 for more description



WWW Standards



WWW Standards

- URL

- RFC 1630, RFC 1738
- Many RFCs define the URL used for telnet, gopher, mailto, POP, IMAP, etc.

- HTML

- RFC 2854

- HTTP

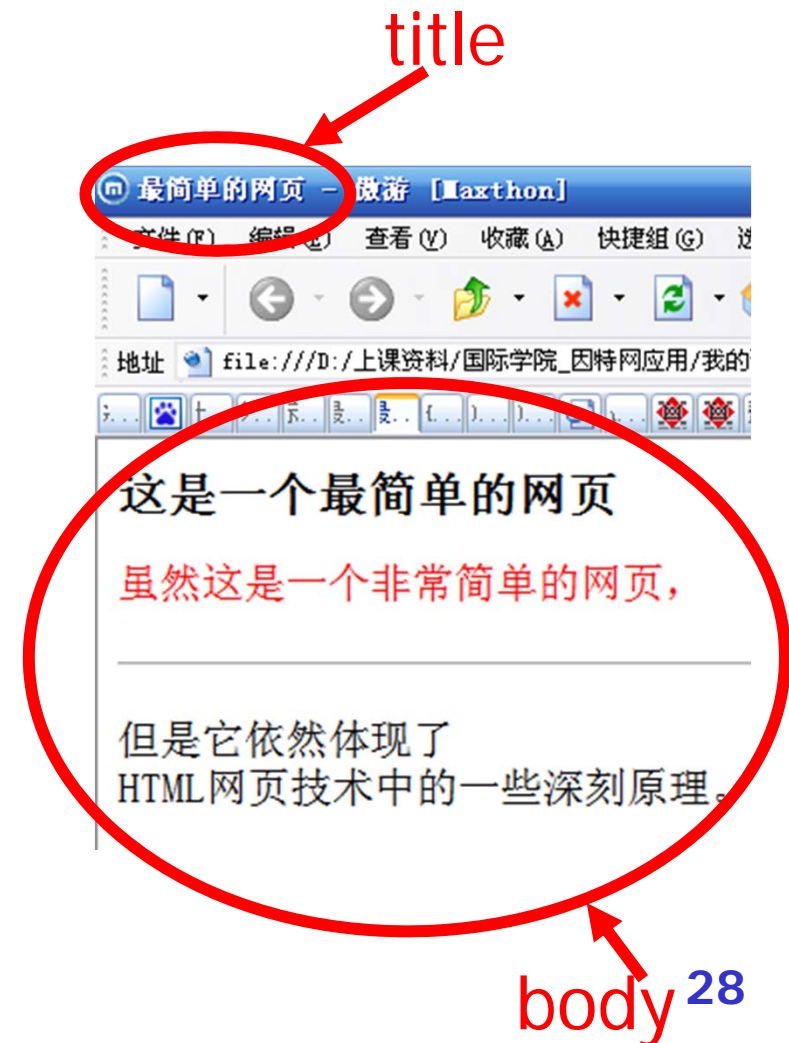
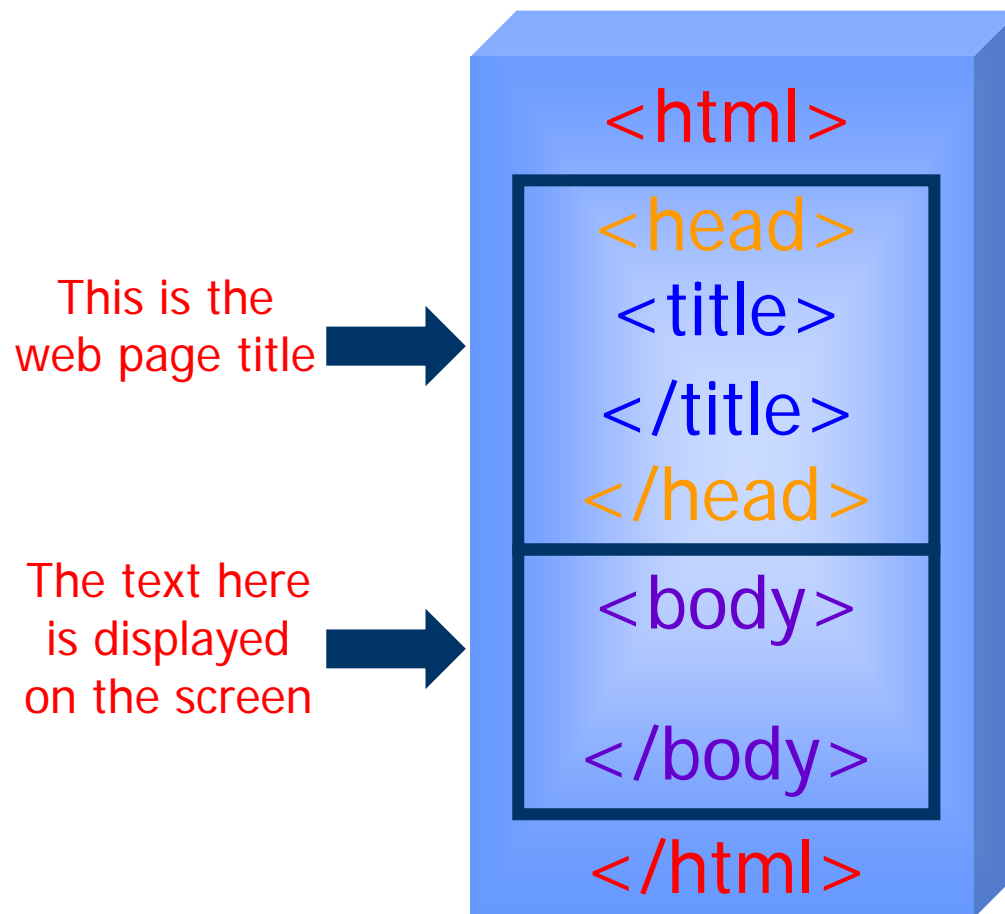
- RFC 2616: defines HTTP/1.1
- RFC 2617: defines HTTP Authentication (Basic and Digest Access Authentication)



HTML – HTML standards

- HTML is the agreed upon markup language for the Web
- Currently several versions are available
 - HTML 1.0 - most basic tags
 - HTML 2.0 - forms support
 - HTML 3.0 - vendor specific tags crept in
 - HTML 3.2 - current standard, scaled-back 3.0
 - HTML 4.0 - current recommended
 - XHTML (eXtensible HyperText Markup Language) 1.0/1.1/2.0 – XML based, more extensible, more flexible
- Depending on the browser you use and what version you use, pages can look different because different browsers support different HTML versions
- Differences between HTML and XHTML
 - <http://www.w3.org/MarkUp/2004/xhtml-faq>

HTML – A Basic Web Page





HTML – Tags (1)

- An HTML **page** is basically an **ASCII text** page with **various tags** inserted to format the page
- The tags can be in UPPER or lower case.
 - **** == **** (the former stands out better though)
- Used to mark text up for display by the browser
 - to divide the document into **logical units** or indicate the semantics of a piece of text
 - to format the display of information like **** to start bold **** to end bold
 - to link to other items like ****



HTML – Tags (2)

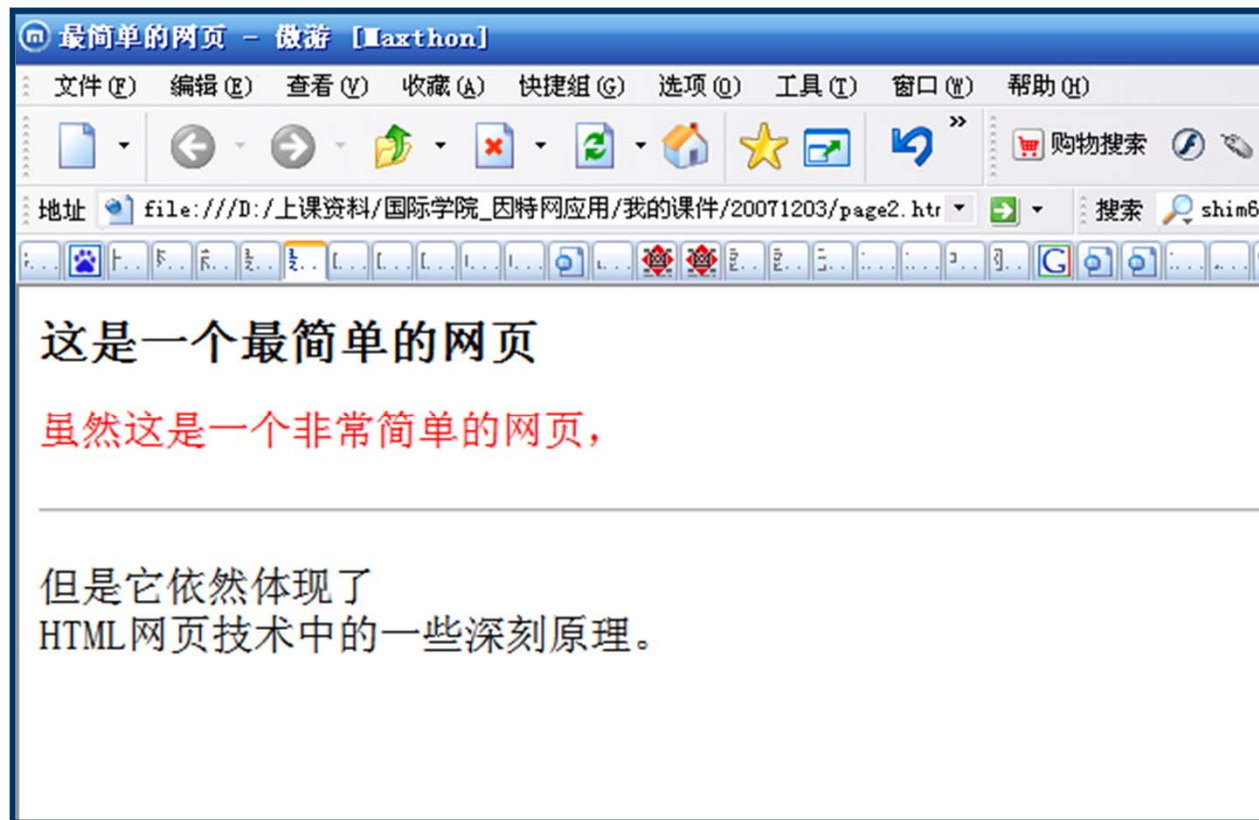
- Tags are **not case sensitive**
- **Blank lines and spaces are ignored** when interpreting HTML document
- Typical tag is: `<h1>This is a heading</h1>`
 - Most tags enclose the marked up text, but there are some that do not need an end tag
- **Anchor tag** is used to “**link**” documents
 - ` ILS Home Page `



HTML – Basic Tags

- **<P>** - paragraph
- **** - bold
- **<I>** - italics
- **<H1>**, **<H2>**, ..., **<H6>** - headers
- **<A>** - anchor, to create a link to another place
- **** - image
- Many tags have ending tags but not always mandatory or necessary
 - For example: This is a **bold section** and this is not
 - When viewed in a browser would be:
This is a **bold section** and this is not
- Some tags have named parameters
 - For example: ****
 - The tag is ****, **SRC** and **ALT** are parameters

HTML – A Simple Example (1)



HTML – A Simple Example (2)

Standards followed by this file

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.0  
Transitional//EN">
```

```
<html>
```

```
<head>
```

```
<meta http-equiv="Content-Type"  
content="text/html; charset=gb2312">
```

```
<title>最简单的网页</title>
```

Name displayed in
the title bar

```
</head>
```

```
<body>
```

```
<h3>这是一个最简单的网页</h3>
```

Heading

```
<p><font color=red>虽然这是一个非常简单的网页，  
</font></p><hr><p>但是它依然体现了<br>HTML网页技术  
中的一些深刻原理。</p>
```

horizontal line

new line

This is the
web page title

This is a
HTML file

This is the
content to be
displayed



HTML – Tags For Images

- To include a graphic image into a Web document, the `` markup is used. It can take several optional parameters, and can refer to a variety of image formats. Basic usage is:
 - ``
- It could also look something like:
 - ``



HTML – Tags For Hyperlinks

- The **hyperlink** is the basis of the entire "linked documents" idea of the Web. The HTML looks like:
 - ` Queen Mary, University of London `
- In a browser, this would normally be displayed something like:
 - Queen Mary, University of London
- When the user moves the cursor (mouse pointer) over this text, and "clicks", the browser fetches the URL named in the HREF parameter and displays it instead of the current page

HTML – Extension To The Example (1)

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.0 Transitional//EN">
<html>
<head>
<meta http-equiv="Content-Type"
content="text/html; charset=gb2312">
<title>最简单的网页</title>
</head>
<body>
<h3>这是一个最简单的网页</h3>
<IMG SRC="typingman.gif" align=middle width=200 height=150
alt="Photo">
<br><A HREF="http://www.mayan.cn/IA"> Couseware download for
Internet Application </A>
<p><font color=red>虽然这是一个非常简单的网页,</font></p><hr><p>但是
它依然体现了<br>HTML网页技术中的一些深刻原理。</p>
</body>
</html>
```

HTML – Extension To The Example (2)



这是一个最简单的网页



[Couseware download for Internet Application](#)

虽然这是一个非常简单的网页，

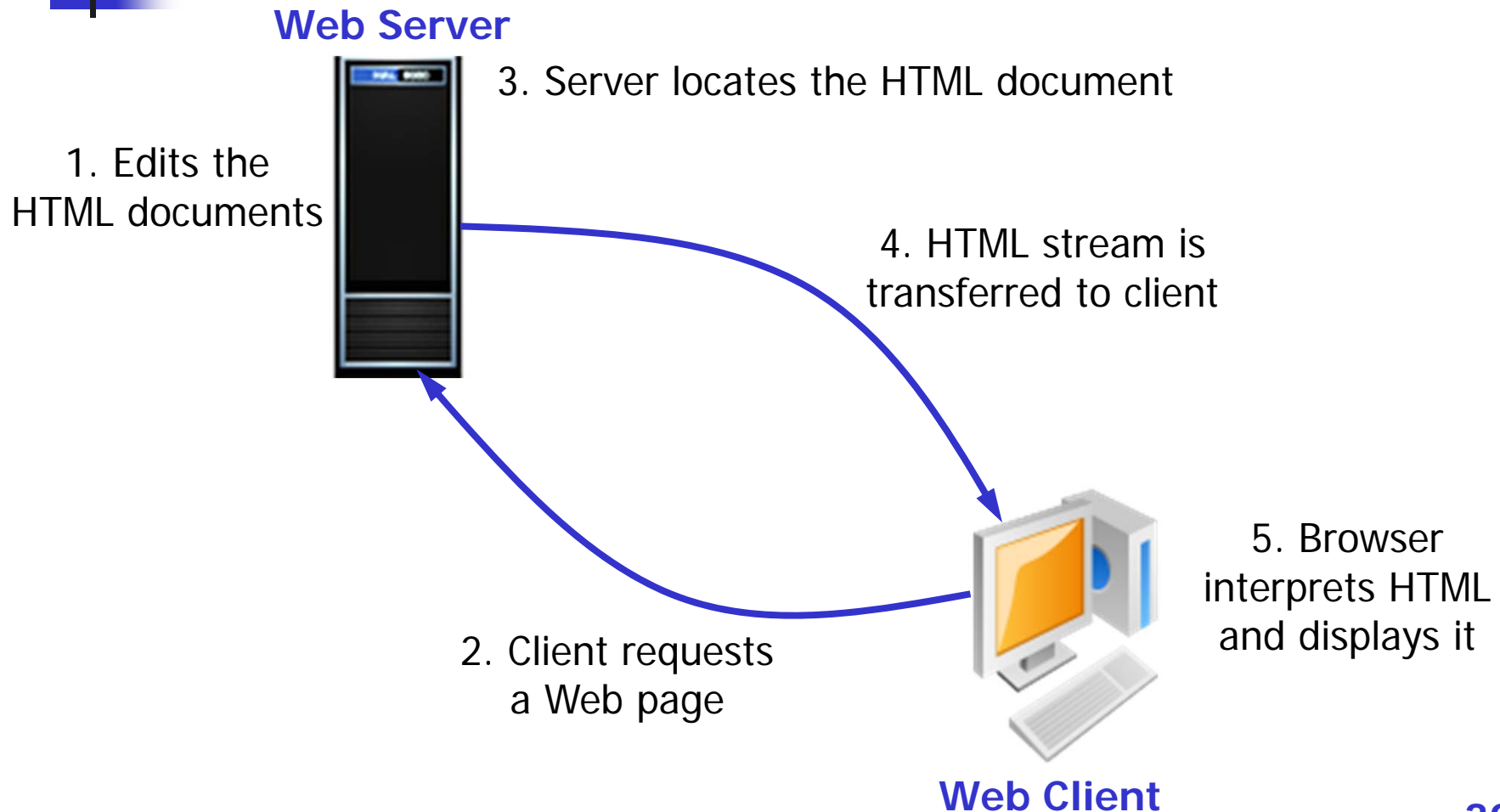
但是它依然体现了
HTML网页技术中的一些深刻原理。



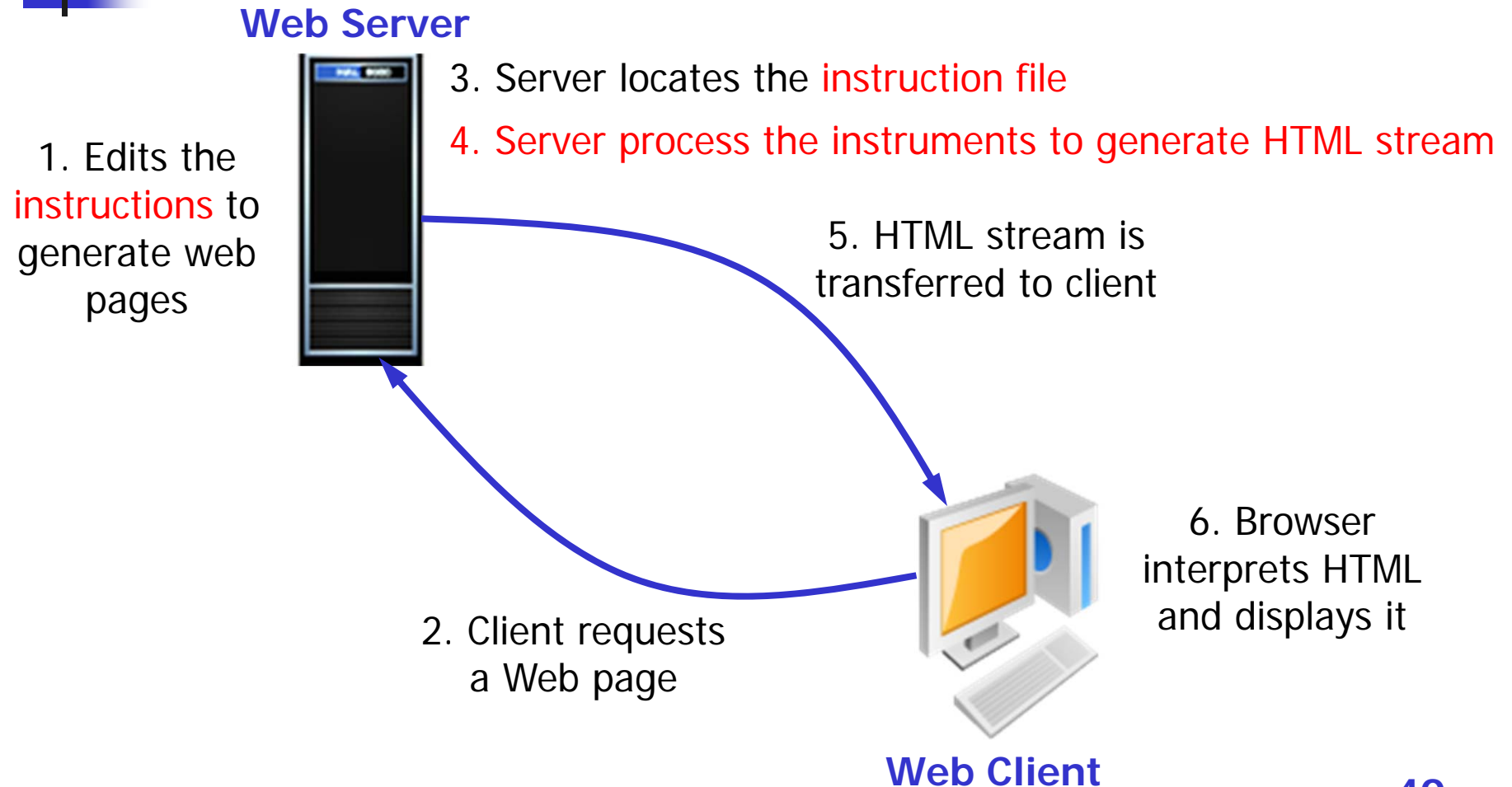
Static vs. Dynamic

- At the beginning, WWW was made up of **static documents**
 - Each URL corresponded to a single file stored on some hard disk
 - Edit in HTML format
 - .html, .htm
- Today - many of WWW documents are **built at request time**
 - The URL doesn't correspond to a single file
 - Examples: website access counter, WWW based date-time server, BBS, ...
 - Generated dynamically by ASP, JSP, VB Script, PHP, CGI or other programs
 - .asp, .shtm, .php, .cgi etc.
- Why dynamic documents?
 - automation of web site maintenance
 - customized advertising
 - database access
 - shopping carts
 - date and time service
 - jobs for ElecEng students

Procedure Of Static Pages

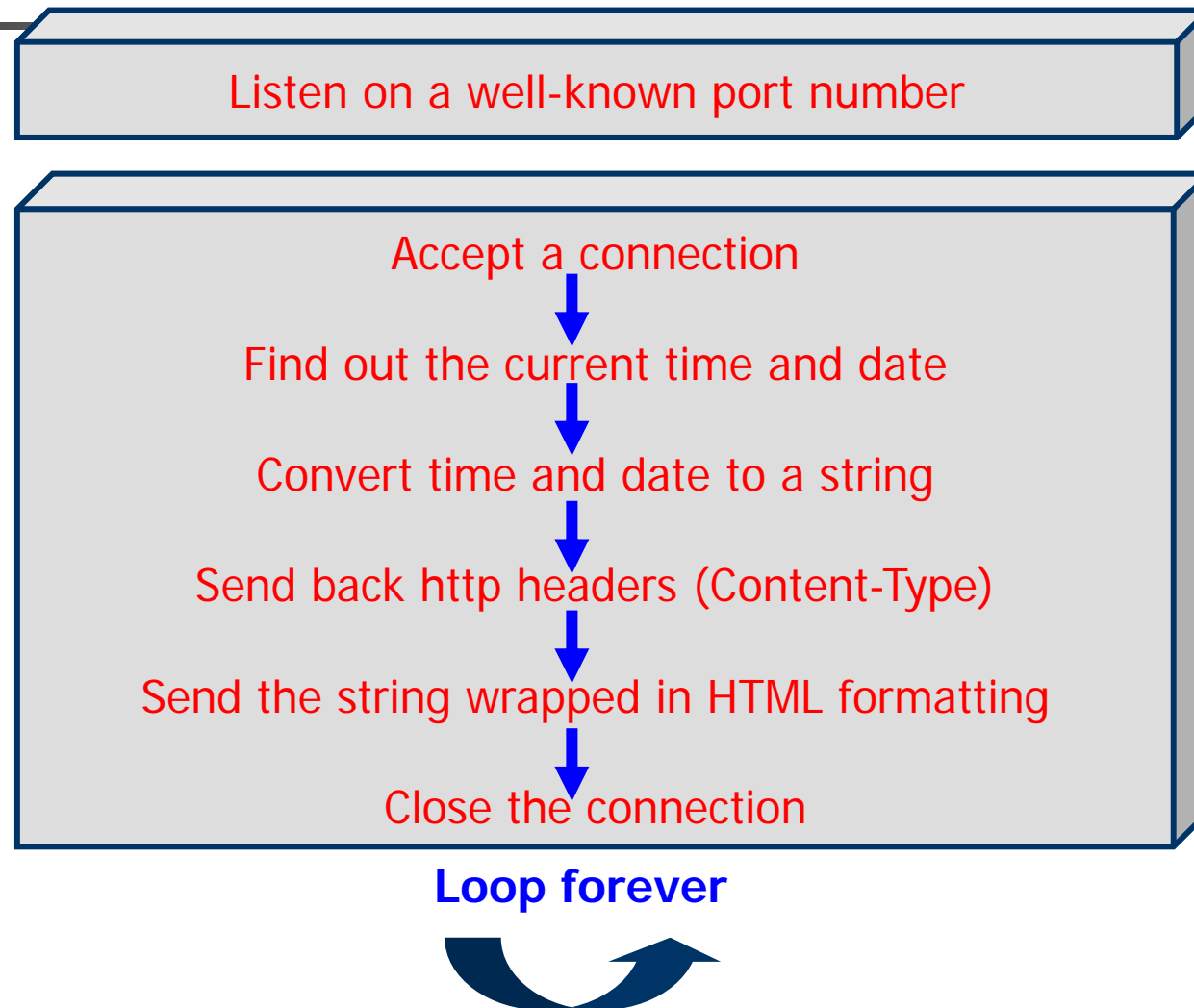


Procedure Of Server-based Dynamic Pages





Example: WWW based time and date server

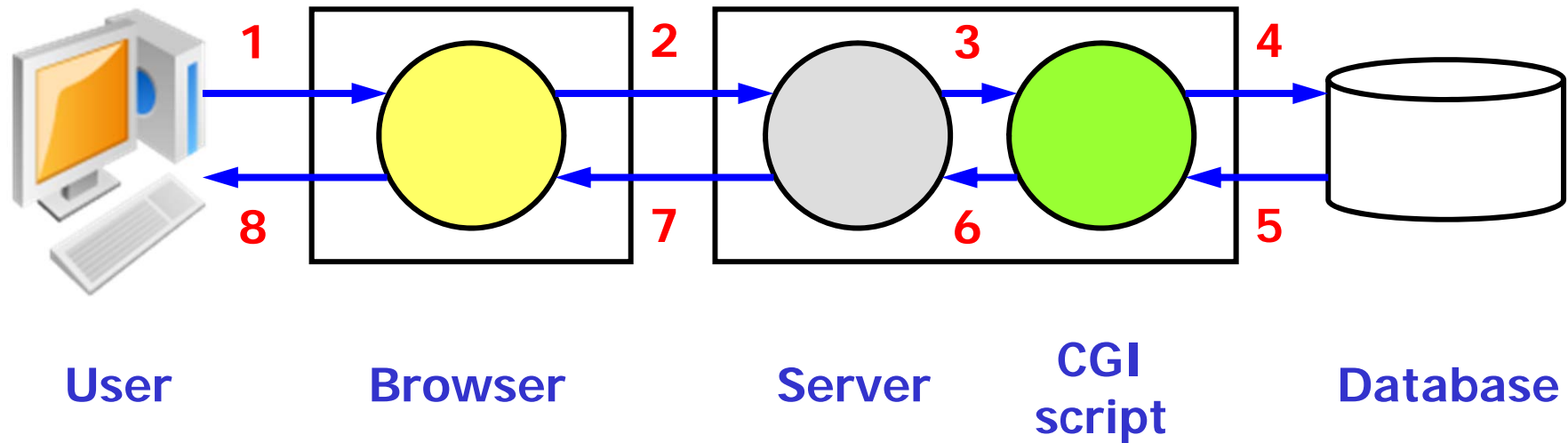




CGI (Common Gateway Interface)

- The Common Gateway Interface (CGI) is a standard for **interfacing external applications with information servers**, such as HTTP or Web servers
- A plain HTML document that the Web client retrieves is static, which means it exists in a constant state: a text file that doesn't change. A CGI program, on the other hand, is **executed in realtime**, so that it can **output dynamic information** for the server
- The Web server executes a CGI program to transmit information to the database engine, receive the results and display them to the client. This is an example of a gateway. Currently version is 1.1
- A CGI program is basically the equivalent of letting the world run a program on your system. For safety, security precautions are taken

CGI – Procedure



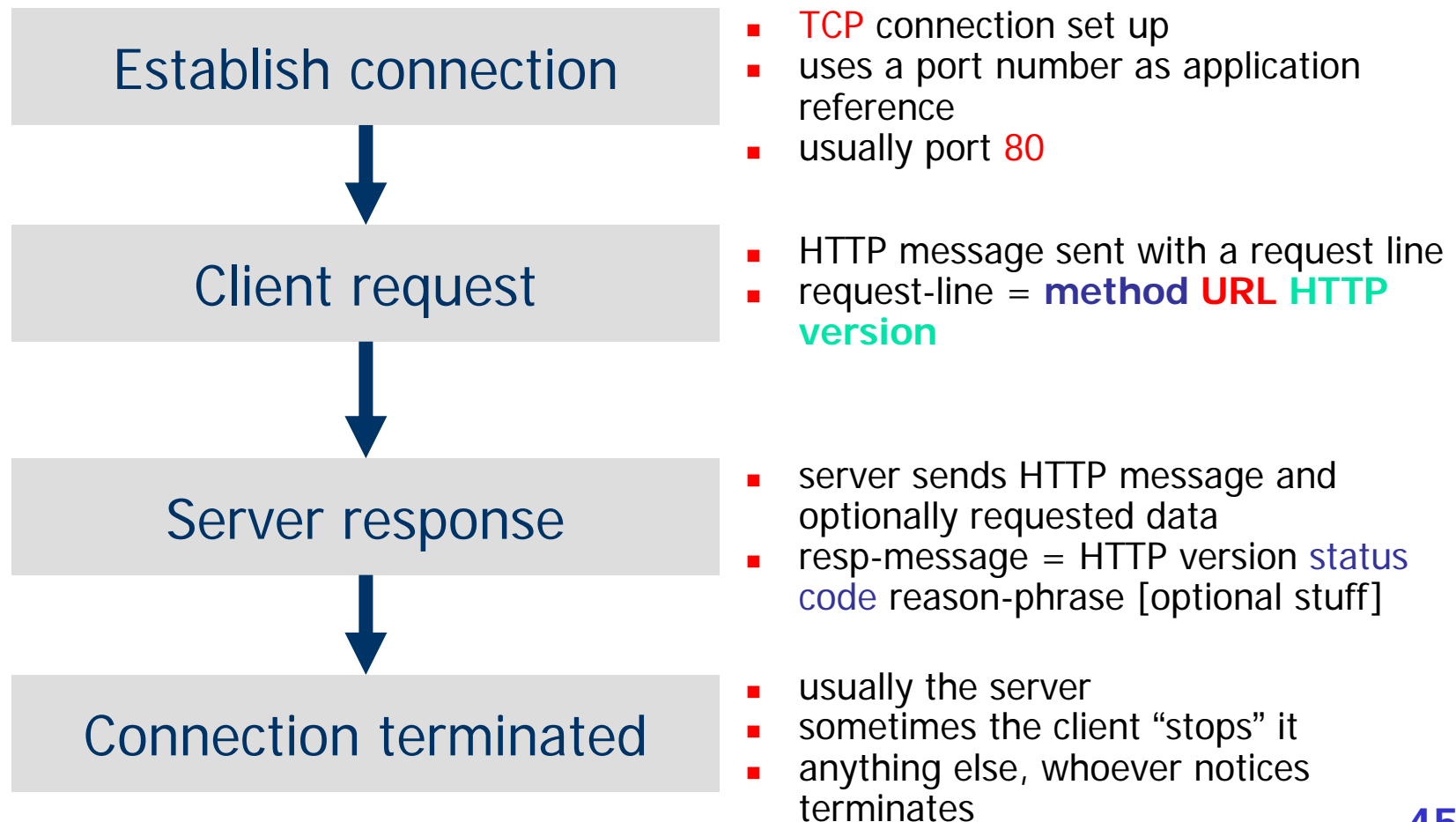


HTTP – Basics

- The heart of the Web
- Features
 - Application layer protocol for client/server communication
 - Request/response based
 - Stateless
 - Bi-directional transfer
 - Capability negotiation
 - Support for cache
 - Support for intermediaries: HTTP proxy



HTTP – HTTP Transaction





HTTP – Status Codes

- 1xx – for information only
- 2xx – action successful
- 3xx – further action needed (redirect)
- 4xx – client request error
- 5xx – server error



HTTP – Getting Remote Web Pages

- The browser determines the URL
- Browser asks DNS for the IP address of web-page being referred to
- DNS returns the IP address to the browser
- The browser makes a TCP connection to port 80 at the web-page IP address
- The browser sends a get request, eg.
 - GET /dir/FileName.html/HTTP/1.0
- The remote server sends the file FileName.html
- The TCP connection is released
- The browser displays all the text in FileName.html
- The browser fetches and displays all the images in FileName.html



HTTP – HTTP Methods

Method	Description
GET	retrieve document specified by URL
PUT	store specified document under given URL
HEAD	identical to GET except that the server MUST NOT return a message-body in the response
OPTIONS	retrieve information about available options
POST	give information (eg. annotation) to the server
DELETE	remove document specified by URL
TRACE	loopback request message
CONNECT	reserved for use with a proxy

HTTP – An ASCII/MIME protocol

- Because HTTP is an **ASCII / MIME protocol**, it is simple for a user at a terminal to communicate directly to a Web server
 - ASCII: defined in RFC 2822
 - MIME: Multipurpose Internet Mail Extension
- Each interaction consists of one **ASCII request**, followed by one **RFC2822 / MIME-like response**
 - e.g. Content-type: text/html
 - Data type/subtype
 - text/html
 - text/plain
 - image/gif
 - video/mpeg
 - application/msword
 - etc.



HTTP – An ASCII/MIME protocol

■ Content types and subtypes defined by MIME

Type	Subtype	Description
Text	Plain	Unformatted text
	Richtext	Text including simple formatting commands
Image	Gif	Still picture in GIF format
	Jpeg	Still picture in JPEG format
Audio	Basic	Audible sound
Video	Mpeg	Movie in MPEG format
Application	octet-stream	An uninterpreted byte sequence
	Postscript	A printable document in PostScript
Message	RFC2822	A MIME RFC 2822 message
	Partial	Message has been split for transmission
	External-body	Message itself must be fetched over the net
Multipart	Mixed	Independent parts in the specified order
	Alternative	Same message in different formats
	Parallel	Parts must be viewed simultaneously
	Digest	Each part is a complete RFC 2822 message



HTTP/1.1 Performance Enhancements

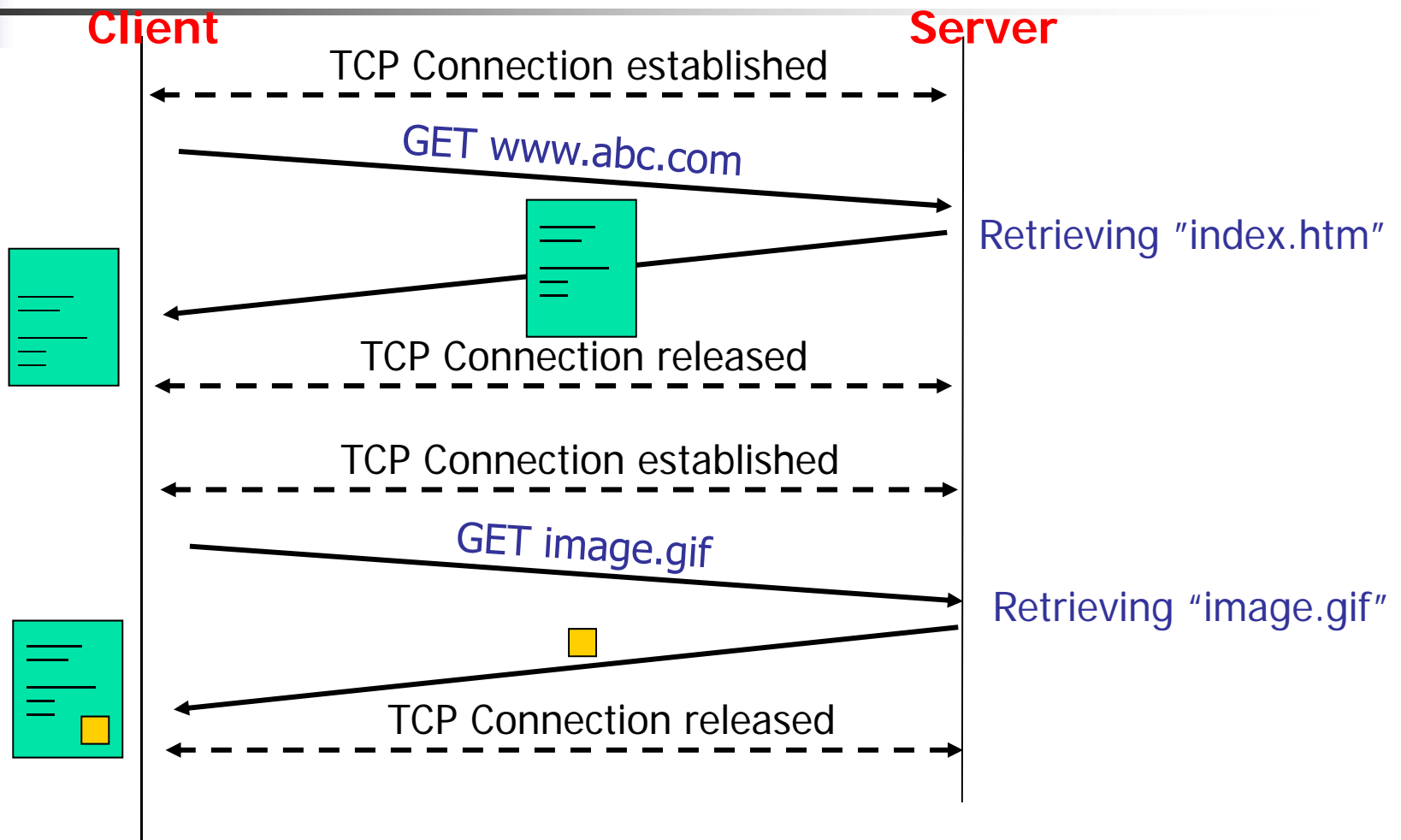
- HTTP/1.0 is a “stop and wait” protocol
 - Separate TCP connection for each file
 - Connect setup and tear down is incurred for each file
 - Inefficient use of packets
 - Server must maintain many connections
- HTTP/1.1 specification focus on performance enhancements
 - Persistent connections
 - Pipelining
 - Enhanced caching options
 - Support for compression



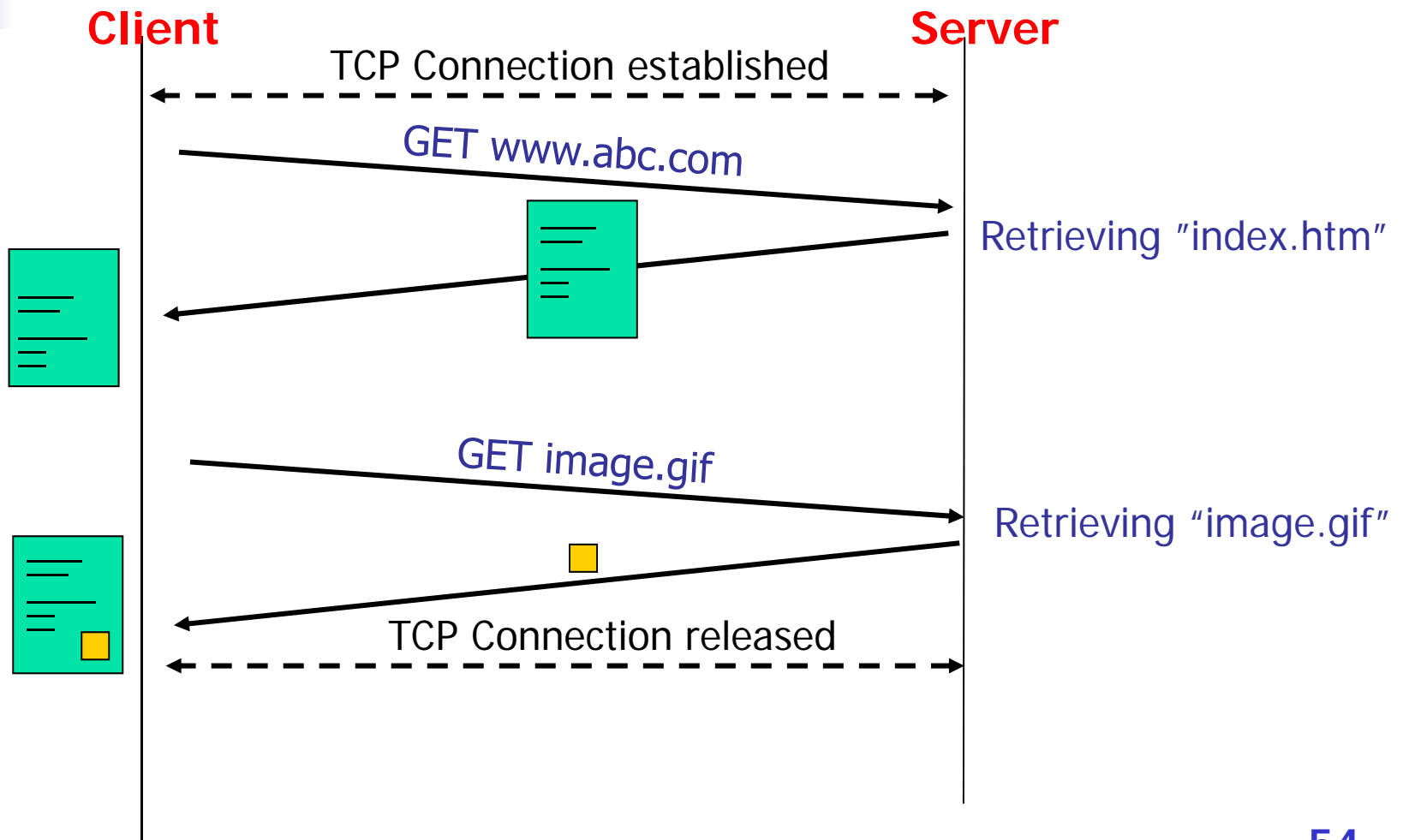
HTTP/1.1 Persistent Connections and Pipelining

- Persistent connections
 - Use the same TCP connection(s) for transfer of multiple files
 - Reduces packet traffic significantly
- Pipelining
 - Multiple HTTP requests can be written out to a socket together without waiting for the corresponding responses
 - Pack several HTTP requests into one TCP/IP packet

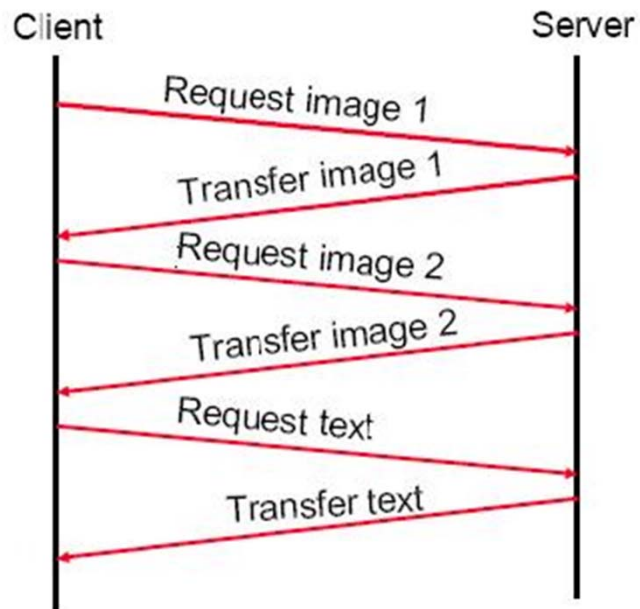
Example of Non-persistent Connections



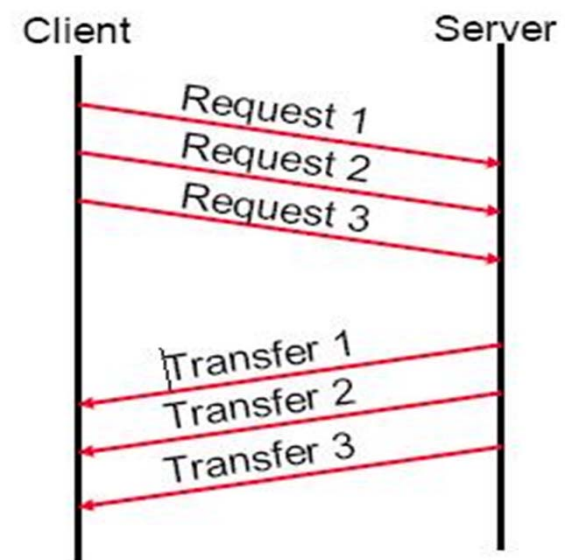
Example of Persistent Connections



Example of Pipelining



Non-pipelining



Pipelining



User-server state: cookies

Many major Web sites use cookies

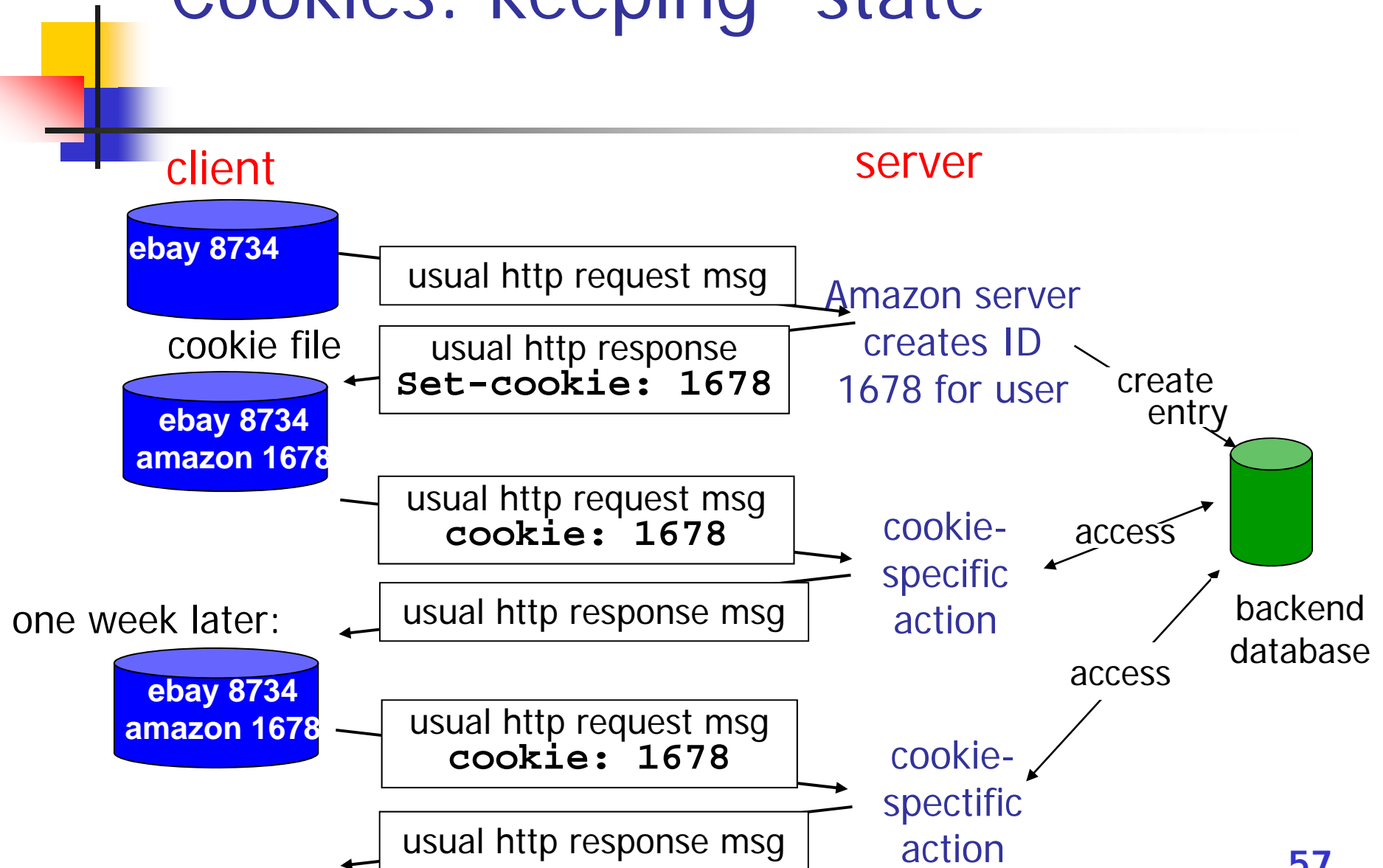
Four components:

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

Example:

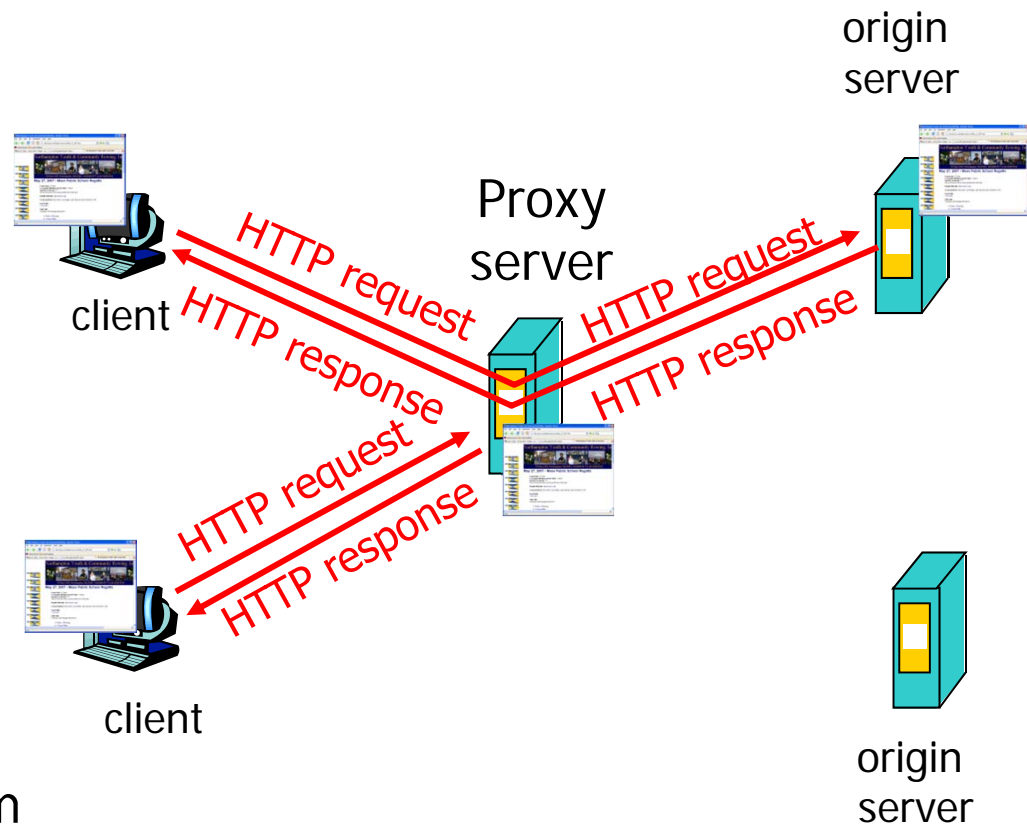
- Susan always access Internet always from PC
- visits specific e-commerce site for first time
- when initial HTTP requests arrives at site, site creates:
 - unique ID
 - entry in backend database for ID

Cookies: keeping "state"

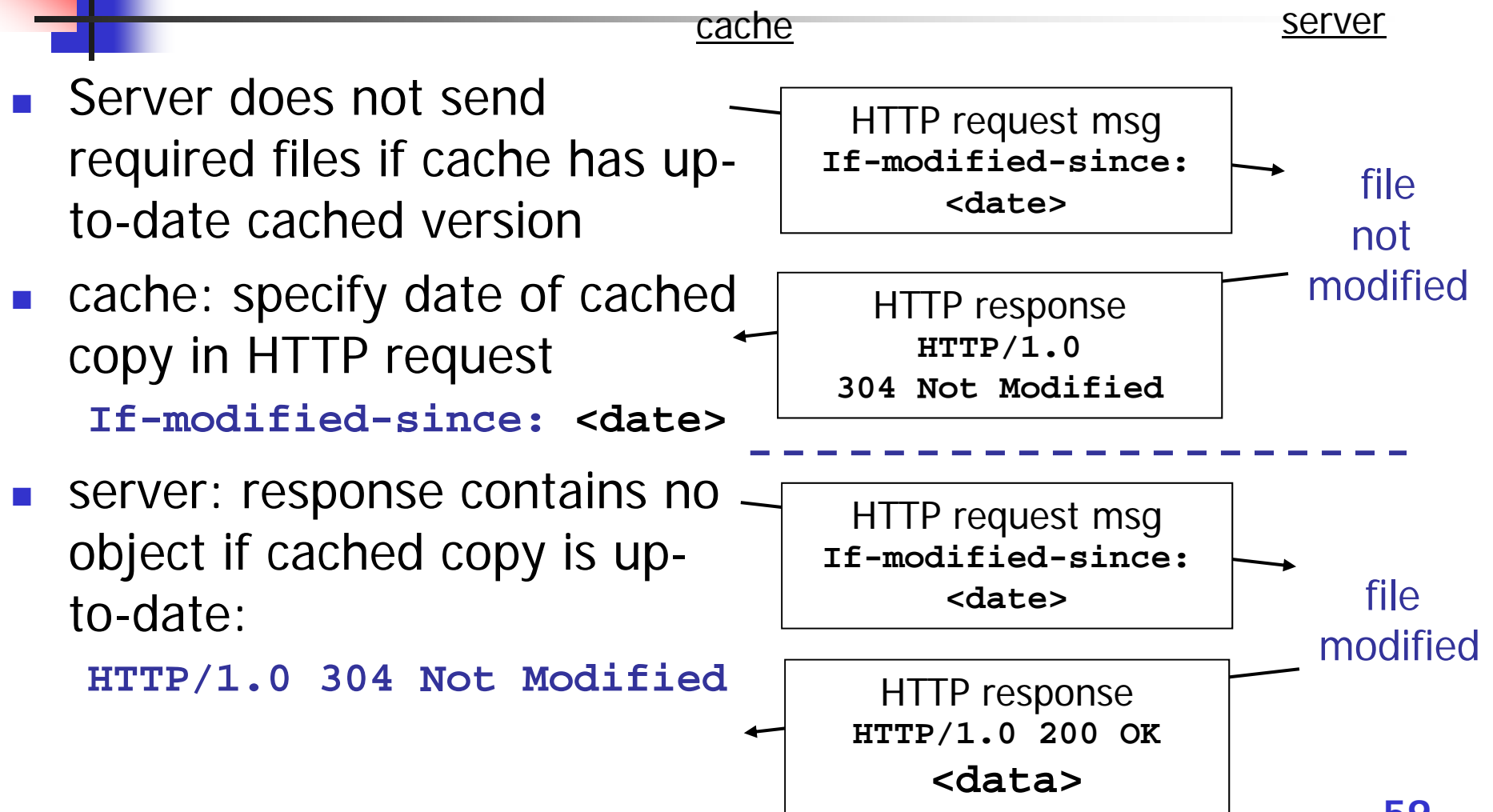


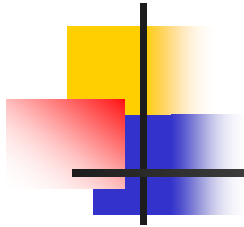
Web Caches(Proxy Server)

- Motivation: satisfy client request without involving origin server
- User sets browser: Web accesses via a **proxy server**
- Browser sends all HTTP requests to proxy server
 - If requested file in cache: proxy server returns file
 - else proxy requests file from origin server, then forwards to client



Conditional get





Web Applications

Web Applications

- Navigating the Web
- Information search
- Information download
- Advertisement and dissemination
- Remote education, remote diagnosis
- ...



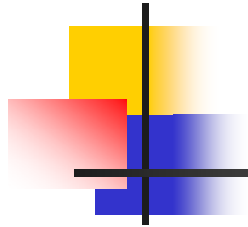
无线信号处理与网络实验室

中文站

Wireless Signal Processing and Networks Laboratory

English





Summary



Summary

- Terminologies
 - WWW, the Web, W3
 - URL
 - HTML
 - HTTP
- WWW components
 - Client/browser
 - Web server
- URL
 - Structure
 - Used for different services
- HTML
 - Basic web page structure
 - Basic tags
 - Static vs. dynamic
 - CGI
- HTTP
 - Features
 - Transaction
 - Methods and responses
 - Performance enhancement of HTTP 1.1



Questions

- Does a web address equal to a domain name?
- What are the disadvantages of the stateless feature of HTTP?
- How is the procedure of client-based dynamic pages?
- What is the cookies?
- How does the HTTP proxy work?



Useful URLs

- W3C
 - <http://www.w3.org/>
- The Web
 - <http://www.learnthenet.com/ENGLISH>
- HTML
 - <http://www.w3.org/MarkUp/Guide/>
 - <http://www.w3.org/MarkUp/Guide/Advanced.html>
 - <http://www.w3.org/MarkUp/2004/xhtml1-faq>
 - <http://www.jmarshall.com/easy/html/>
 - <http://www.dreamdu.com/>
- HTTP
 - <http://www.jmarshall.com/easy/http/>
- A detailed description of Internet history
 - <http://www.zakon.org/robert/internet/timeline/>