# SGG 3643

# Computer Programming III

## The Internet, URLs, Protocols

Ivin Amri Musliman

# Table of Contents

# Web References

- HyperText Markup Language (HTML)
  Home Page, http://www.w3.org/MarkUp/

- *Getting started with HTML* by Dave Raggett

- Quittner, Joshua. "Network Designer:Tim Berners-Lee". 1999, Time Magazine.
  http://www.time.com/time/time100/scientist/profile/bernerslee.html

# Web References II

- CSS Primer:
  http://www.moock.org/webdesign/css/index.html

- Elements of Style:
  http://www.webtechniques.com/archives/2001/03/desi/

- HTML with Style Tutorials:
  http://webreference.com/html/tutorials/

- Style Sheet Reference Guide:
  http://www.webreview.com/style/

# Related Work

- XML: the universal format for structured documents and data on the Web. It allows you to define your own mark-up formats when HTML is not a good fit. XML is being used increasingly for data; for instance, W3C's metadata format RDF.

- W3C's Cascading Style Sheets language (CSS) provides a simple means to style HTML pages, allowing you to control visual and aural characteristics; for instance, fonts, margins, line-spacing, borders, colors, layers and more. W3C is also working on a new style sheet language written in XML called XSL, which provides a means to transform XML documents into HTML.

# Related Work II

- <u>Document Object Model</u>
  - Provides ways for scripts to manipulate HTML using a set of methods and data types defined independently of particular programming languages or computer platforms. It forms the basis for dynamic effects in Web pages, but can also be exploited in HTML editors and other tools by extensions for manipulating HTML content.

# Definitions

- "an internet-wide distributed hypermedia information retrieval system" [Liu et al. 1994]

- "the World Wide Web is a global, seamless environment in which all information (text, images, audio, video, computational services) that is accessible from the Internet and can be accessed in a consistent and simple way by using a standard set of naming and access conventions" [WebMaster Magazine 1996]

- "the World Wide Web (known as "WWW', "Web" or "W3") is the universe of network-accessible information, the embodiment of human knowledge" [W3C 1999]
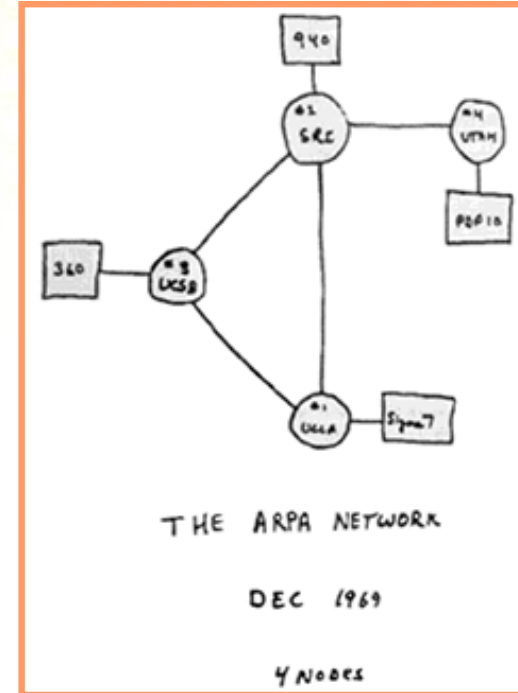
# The Web: Ideas and Objectives

- A uniform naming scheme for locating resources on the Web (e.g., URIs).
- Protocols, for access to named resources over the Web (e.g., HTTP).
- Hypertext, for visualization / presentation and easy navigation among resources (e.g. HTML, media objects).
- Availability on different computer systems
- Uniform access (reading and writing) using standard interface
- Integration of external information resources (i.e. database systems)

# History

- 1965 "Hypertext" and "Hypermedia" (Ted Nelson)

- 1969 ARPANET (4 nodes)

- 1974 TCP protocol

- 1983 "Internet"

- 1989 World Wide Web (Berners-Lee, Cailliau; Release 1991)

- 1993 Mosaic Browser (Web increases 341634%
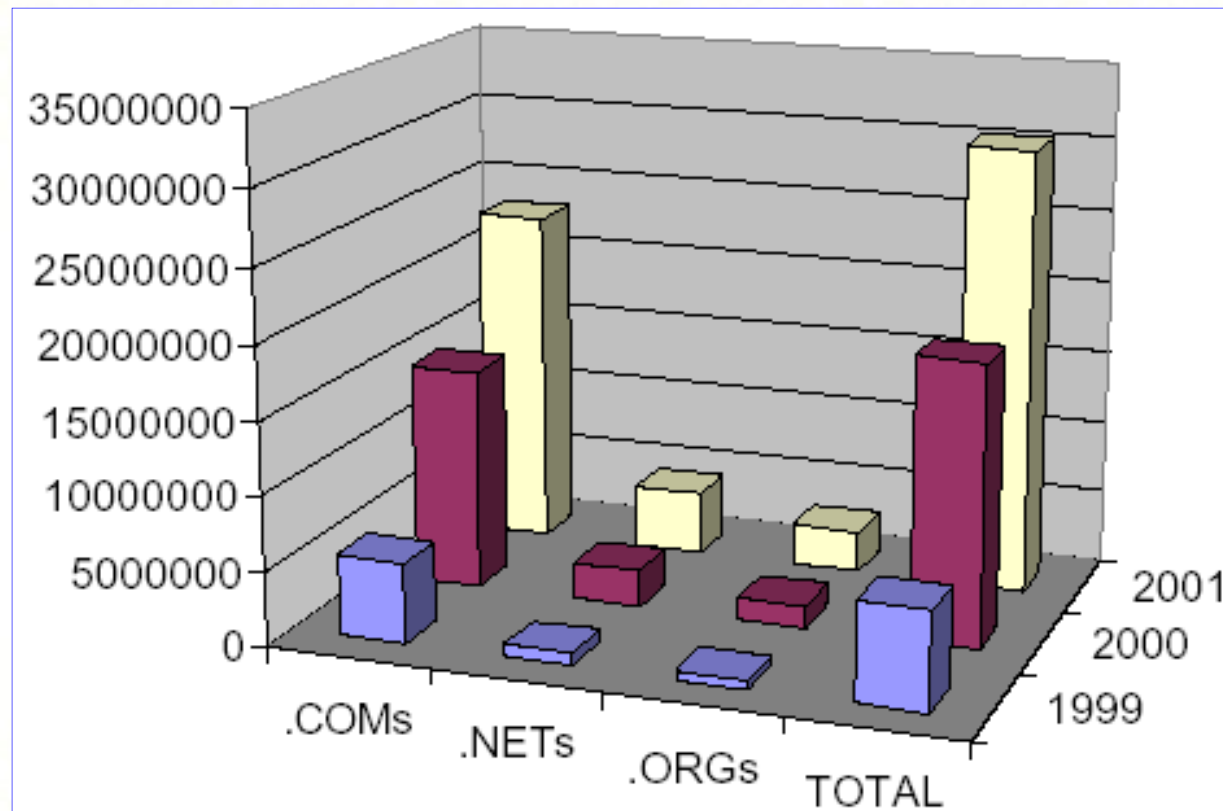
- per year)

- 1995 Web exceeds FTP in transfer volume

**See also: A Brief History of the Internet,**
**http://www.isoc.org/internet/history/brief.shtml**



THE ARPA NETWORK

DEC 1969

4 NODES

http://kartoweb.itc.nl/webcartography/webbook/ch06/ch06.htm

# Internet Growth



Source: http://www.netfactual.com

# Some statistics

- Total domains registered worldwide:

  | | |
  |---|---|
  | 31497437 | INTERNATIONAL (COM) |
  | 21783099 | INTERNATIONAL (EDU) |
  | 7429 | INTERNATIONAL (NET) |
  | 3658670 | INTERNATIONAL (ORG) |
  | 2408744 | United Kingdom (CO.UK) |

Source: http://www.domainstats.com/, Nov 23,2002

# Use and Users of the Web



**Internet users**

| | |
|---|---|
| · | 200 - 50,000 |
| · | 50,001 - 500,000 |
| ● | 500,001- 5,000,000 |

**Top 15 nations (x 1 million)**

| | | | |
|---|---|---|---|
| ● Spain 5.2 | ● France 9.0 | ● China 15.8 |
| ● Netherlands 5.4 | ● Brazil 10.6 | ● United Kingdom 17.9 |
| ● Taiwan 6.5 | ● Italy 11.6 | ● Germany 19.1 |
| ● Russia 6.6 | ● South Korea 14.8 | ● Japan 26.9 |
| ● Australia 8.1 | ● Canada 15.2 | ● United States 135.7 |

Corné P.J.M. van Elzakker: Use and Users of Maps on the Web, http://www.nacis.org/cp/cp37/CP37_34_50.pdf

# Main organisations



http://www.w3c.org/

http://www.ietf.org/

# Substantial ideas

- Document format
  - Hypertext Markup Language, HTML
  - Document Type Definition (DTD) Standardized General Markup Language (SGML)
- Transfer Protocol
  - Hypertext Transfer Protocol, HTTP (ASCII-coded Request-Reply protocol using TCP/IP)

# Substantial ideas II

- Uniform identification schema: Uniform Resource Identifier (URI)
- Information resources
    - Information resources must be identifiable (per name, per adress / location)
    - Each resources in the Internet should be identifiable
    - Identification schema (string!) must be expandable and complete.

# Uniform Resource Identifier

# Uniform Resource Identifier (URI)

- Syntax for all Identificators [RFC2396]

  <uri> ::= <scheme>":"<scheme-specific-part>

- An absolute URI contains the name of the scheme being used (<scheme>) followed by a colon (":") and then a string (the <scheme-specific- part>) whose interpretation depends on the scheme.

- <scheme> Name schema for this URI

- <scheme-specific-part> contains actual identification corresponding to this scheme

# Uniform Resource Identifier (URI) II

- URIs can be:

    - Locations / adresses: Uniform Resource Locator (URL)

    - Names: Uniform Resource Name (URN, Objective: Uniform names for all resources)

    - Metainformationen: Uniform Resource Characteristic (URC)

# URI examples

- ftp://ftp.is.co.za/rfc/rfc1808.txt - ftp scheme for File Transfer Protocol services

- http://www.math.uio.no/faq/compression-faq/part1.html- http scheme for Hypertext Transfer Protocol services

- mailto:mduerst@ifi.unizh.ch- mailto scheme for electronic mail addresses

- news:comp.infosystems.www.servers.unix- news scheme for USENET news groups and articles

- telnet://melvyl.ucop.edu/- telnet scheme for interactive services via the TELNET Protocol

# URN – experimental state

- URN [RFC 1737, RFC 2141, RFC 3061]
  (<scheme> ::= "urn")
  <urn> ::= "urn:" <nid> ":" <nss>

  – nid = Namespace Identifier

  – nss = Namespace Specific String

- Properties::

- Uniform Resource Names (URNs) are intended to serve as persistent, location-independent, resource identifiers.

  – globally unique (global scope and uniqueness)

  – persistent

  – Extensible, independent

# Uniform Resource Locator (URL)

- URL scheme Definitionen [RFC1738]:

  http, https, ftp, news, mailto, telnet and others

- scheme-specific-part has the general format:

  ["//"] [user [":"password] "@"] host [":"port] ["/"url-path]

- relative URLs possible [RFC 1808]

- General scheme-specific-part can be reduced to:
  <http_URL> = "http://" <host> [ ":" <port> ] [<abs_path>]

  Example: http://www.gis-news.de/links/xml.htm

# URL: ["//"] [user [":"password] "@"] host [":"port] ["/"url-path]

- user: An optional user name. Some schemes (e.g., ftp) allow the specification of a user name.

- password: An optional password. If present, it follows the user name separated from it by a colon.

- The user name (and password), if present, are followed by a at-sign "@". Within the user and password field, any ":", "@", or "/" must be encoded.

- host: The fully qualified domain name of a network host, or its IP address (RFC 1034).

- port: The port number to connect to; in most cases: default port number.

- url-path: supplies the details of how the specified resource can be accessed. Note that the "/" between the host (or port) and the url-path is NOT part of the url-path.

# Base URLs

- The Base URL of a resource is everything up to and including the last slash in its pathname.

| Absolute URL | Base URL |
|---|---|
| http://gis-news.de/ | http://gis-news.de/ |
| http://gis-news.de/xml/ | http://gis-news.de/xml/ |
| http://gis-news.de/format/about.html | http://gis-news.de/format/ |
| http://gis-news.de/map/map1.html?x=3500100.0 | http://gis-news.de/map/ |

# Absolute vs. relative URLs

- absolute URLs: identified a resource *independent* of their context.

- relative URLs: a way to identify a resource *relative* to their base URL.

# Relative URL Examples

- Example:
  Base URL: http://WebReference.com/html/

| Relative URI | Absolute URI |
|---|---|
| about.html | http://WebReference.com/html/about.html |
| tutorial1/ | http://WebReference.com/html/tutorial1/ |
| tutorial1/2.html | http://WebReference.com/html/tutorial1/2.html |
| / | http://WebReference.com/ |
| //www.internet.com/ | http://www.internet.com/ |
| /experts/ | http://WebReference.com/experts/ |
| ../ | http://WebReference.com/ |
| ../experts/ | http://WebReference.com/experts/ |
| ../../../ | http://WebReference.com/ |
| ./ | http://WebReference.com/html/ |
| ./about.html | http://WebReference.com/html/about.html |

# Rules for URI, URN and URL

**Several Rules are prerequisites for uniform resource identifier:**

- Character set: ISO Latin-1 (similar ASCII-character set)
- No blanks: Masqueraded by %20.
- Reserved characters are among others:
  - Escape character (**%**): Identification of characters.
  - /: Hierarchical separation of directory and file names.
  - Fragment Delimiter (**#**): Separates URI from a fragment or a partial field of a data object.
  - Query Delimiter (**?**): Separates query string of a ressource from the URI **à** see Common Gateway Interface (CGI)

http://www.utoronto.ca/webdocs/HTMLdocs/NewHTML/iso_table.html

# Example:

- http://localhost/cgi-bin/mapserv.exe?img.x=348&img.y=318&zoomdir=1&zoomsize=2&layer=St%E4dte&layer=motorway&layer=roads&layer=Fl%FCsse&layer=H%F6henstufen&layer=Europe&map=C%3A%2FProgramme%2FApache+Group%2FApache%2Fhtdocs%2Fdemo.map&imgext=4.296614+49.186302+7.284114+52.174938&imgxy=274.5+274.5

# IP-Addresses

# IP-Addresses

- Every computer that communicates over the Internet is assigned an IP address that uniquely identifies the device and distinguishes it from other computers on the Internet.

- An IP address consists of 32 bits, often shown as 4 octets of numbers from 0-255 represented in decimal form instead of binary form.

- Example:

  168.212.226.204

  10101000.11010100.11100010.11001100

# IP-Addresses

- An IP address consists of two parts, one identifying the network and one identifying the node, or host.

- The Class of the address determines which part belongs to the network address and which part belongs to the node address. All nodes on a given network share the same network prefix but must have a unique host number.

- -> Network classes A, B, C [, D, E]

| Netztyp | IP-Adressierung | Typische IP-Adresse |
|---|---|---|
| Klasse-A-Netz | **xxx**.xxx.xxx.xxx | **103**.234.123.87 |
| Klasse-B-Netz | **xxx.xxx**.xxx.xxx | **151.170**.102.15 |
| Klasse-C-Netz | **xxx.xxx.xxx**.xxx | **196.23.155**.113 |

Details: http://www.oreilly.com/catalog/coreprot/chapter/appb.html

# IP Classes

- **Class C**

  This is the most widely used class by small businesses. When you look at the IP address, you'll notice that class C networks start with a first number that's between 192 and 223
  Example: 205.161.74.x

  Network identity    Host in this network

- There can be up to 2,097,151 class C networks and each network can handle close to 254 computers.

# IP Classes B and A

- **Class B**

  IP addresses of this type starts with a number between 128 and 191. It's possible to have 16,384 of these networks and each class B network can handle up to 65,534 IP addresses or computers.

- **Class A**

  Starts with a number between 1 and 126. Only 126 of these networks are available, however each class A network can handle 16,777,214 IP addresses or computers.

# Class D and E Networks

- **Class D Network:**
  Binary addresses start with 1110, therefore the decimal number can be anywhere from 224 to 239. Class D networks are used to support multicasting.

- **Class E Network:**
  Binary addresses start with 1111, therefore the decimal number can be anywhere from 240 to 255. Class E networks are used for *experimentation*. They have never been documented or utilized in a standard way.

# DNS - Domain Name Service

- Mapping of host names to IP addresses

Example: utm.my

**u** Sub-Level-Domain

- – Can contain additional sub domains, i.e. java.sun.com
- – Name must respect trademarks etc.

**u** Top-Level-Domains, i.e.

- – my - Malaysia
- – com - (US-) Companies
- – edu - education
- – org – Organisations

Responsible: National authorities, like Suruhanjaya Komunikasi & Multimedia Malaysia (SKMM), and Telekom Malaysia.

# Domain Name: Definition

- A domain name is a unique, clear and descriptive addressing standard used to locate a specific destination on the Internet. They are primarily used to point to website pages or used in the creation of email addresses.

- A domain name often relates to the name of a business, organisation or service, and like any company name, it has to be registered.

# Length and characters of domain names

- Usually between 3 and 63 characters long (excluding the suffix such as .com),

- Alphanumeric characters and the hyphen (e.g., 0 to 9, a to z and the hyphen (-) ). Some domains can have less than 3 characters, or may not allow as many as 63, but as a general rule these are typical.

- A space cannot be used in a domain name, and should not begin or end with a hyphen (-).

- Recommendation: that the domain name be kept short for ease of use and to make it easier to remember.

# Top Level Domains

- **Top Level Domain**: The highest level domains.
- They are either
  - gTLD - Generic Top Level Domain (e.g..com, .net) or
  - ccTLD - Country Code Top Level Domain, e.g. .my, .sg
  - Country codes can be found at ISO
- **Second Level Domain** (SLD) - Domain names with two or more suffixes, e.g. co.uk.

# Internet Services

- World Wide Web (WWW)
- Newsgroups (News)
- E-Mail
- Telnet
- File Transfer (FTP)
- Chat (IRC)
- Gopher, ...

# Internet Services

- World-Wide Web

  An application that uses the Internet to transport hypertext/multimedia documents. These documents are viewed by a ***browser***

- Newsgroups

- e-mail

- telnet

- ftp

# Internet Services
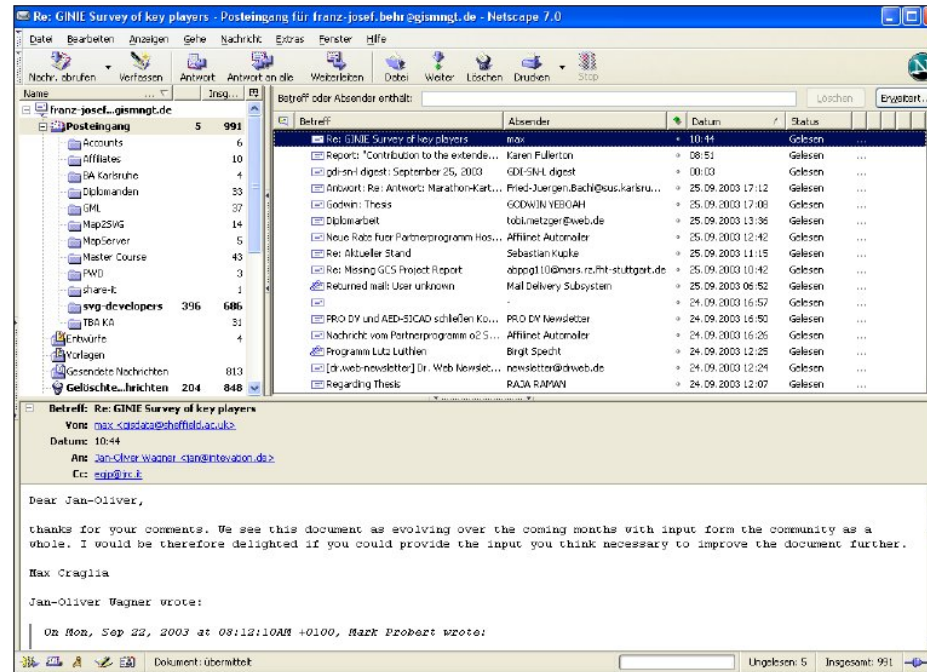
- World-Wide Web
- Newsgroups

  An on-line forum that allows users from all over the world to participate in a discussion about a specific topic

- e-mail
- telnet
- ftp

# Internet Services

- World-Wide Web

- Newsgroups

- e-mail

  Electronic mail

- telnet

- ftp

# Internet Services

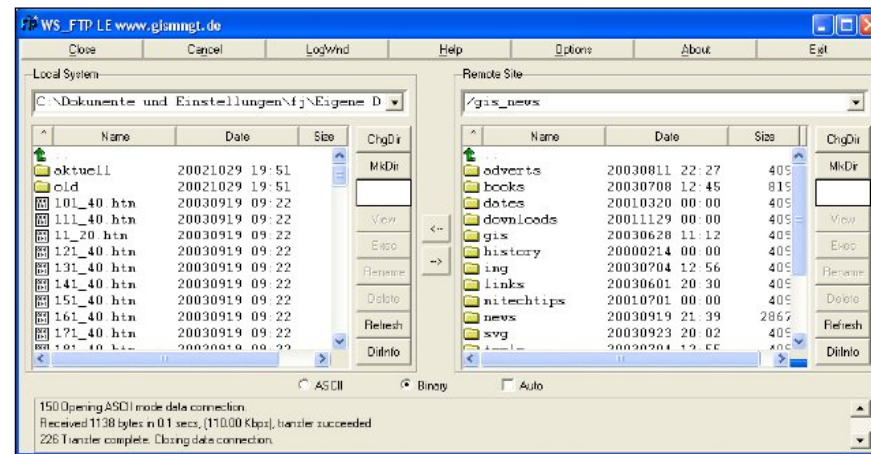- World-Wide Web

- Newsgroups

- e-mail

- telnet

    A program which permits a user on one computer to use that computer as a terminal to access another, perhaps distant, computer

- ftp

# Internet Services

- World-Wide Web

- Newsgroups

- e-mail

- telnet

- ftp



***File Transfer Protocol:*** protocol for transferring files between computers – important for uploading your files to a web site!

# MIME Types

- Multipurpose Internet Mail Extensions
- redefine the format of messages to allow for
  - textual message bodies in character sets other than US-ASCII,
  - an extensible set of different formats for non-textual message bodies (i.e. image media, audio media, video media, applications)
  - multi-part message bodies (data consisting of multiple entities of independent data types), and
  - textual header information in character sets other than US-ASCII.

# MIME Types

- General form: type and subtype

  text/plain

- always case-insensitive.

Content-type: text/plain; charset=iso-8859-1

- The "text" media type is intended for sending material which is principally textual in form.

- A "charset" parameter may be used to indicate the character set of the body text.

- Plain: text that does not contain any formatting commands or directives. Plain text is intended to be displayed "as-is", that is, no interpretation of embedded formatting commands, font attribute specifications, processing instructions, interpretation directives, or content markup should be necessary for proper display.

# Other content types

- text/html
- text/css
- image/jpeg
- application/excel

# Hypertext Transfer Protocol, HTTP

- Properties:
  - Based on TCP/IP
  - stateless, HTTP 1.1: persistent connections.
  - ASCII coded
  - Transfer of text and multimedia resources, i.e. graphic, images, audio, video
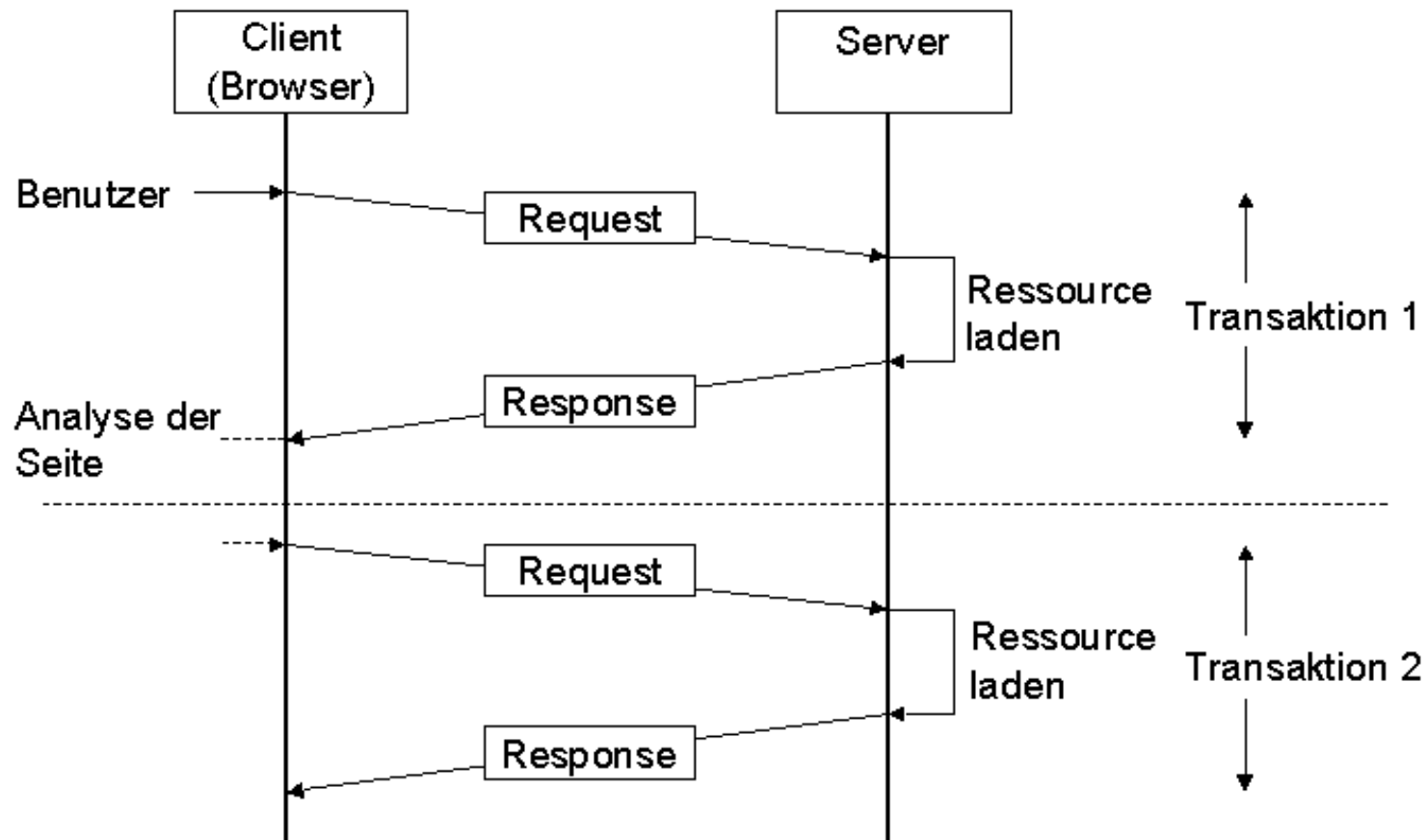  - based on MIME (multipurpose internet mail extensions)

More information:
http://www.mhonarc.org/~ehood/MIME/2046/rfc2046.html
http://www.w3.org/Protocols/

# HTTP Example for Transaction



Source: http://www.teco.uni-karlsruhe.de/lehre/webe/webev221099/sld046.htm

# HTML user clients (i.e. browser)

- HTTP: Client server oriented protocol
- "User Client": browser
  - parses the characters of an HTML document into data characters and markup.
  - robust, extensible (i.e. plug-ins, JAVA)
  - supports the `ISO-8859-1' character encoding scheme and processes each character in the ISO Latin Alphabet No. 1
  - allows the user to traverse hyperlinks from <a> elements in an HTML document.

# HTML user clients (i.e. browser)

| + Expand All | | | |
|---|---|---|---|
| **+ Windows IE** | **+ Mosaic** | **+ Netscape** | **+ Opera** |
| 1.0 Final Aug. 1995 | 1.0 Final Nov 1993 | 1.0 Final Dec 1994 | 2.1 Final Dec. 1996 |
| 2.0 Final Nov. 1995 | 2.0 Final Nov. 1995 | 1.1 Final Apr. 1995 | 3.0 Final Dec. 1997 |
| 3.0 Final Aug. 1996 | 2.1 Final Jan. 1996 | 2.0 Final Mar. 1996 | 3.5 Final Nov. 1998 |
| 4.0 Final Oct. 1997 | 3.0 Final Jan. 1997 | 3.0 Final Aug. 1996 | 4.0 Final Jun. 2000 |
| 5.0 Final Mar. 1999 | Mosaic Ends | 4.0 Final Jun. 1997 | 5.0 Final Dec. 2000 |
| 5.5 Final Jul. 2000 | | 4.5 Final Oct. 1998 | 6.0 Final Nov. 2001 |
| 6.0 Final Oct. 2001 | | 6.0 Final Nov. 2000 | 7.0 Final Jan. 2003 |
| **+ Macintosh IE** | | 7.0 Final Aug. 2002 | |
| 2.0 Final Apr. 1996 | | | |
| 2.1 Final Sep. 1996 | | | |
| 3.0 Final Jan. 1997 | | | |
| 4.0 Final Jan. 1998 | | | |
| 4.5 Final Jan. 1999 | | | |
| 5.0 Final Mar. 2000 | | | |

http://www.blooberry.com/indexdot/history/browsers.htm

# HTML user clients (i.e. browser)

- Today we have > 80% IE 5.0, 5.5 and IE6, less than 10% NN4.xx and approximately 3.5% IE4.

- Netscape 7, Opera and other browsers are fighting for the rest (http://www.upsdell.com/BrowserNews/stat.htm).

- 15% have browsers with Java disabled or unsupported.

# HTTP request headers

- Here is an example request:

```
GET /X/Y.html HTTP/1.0
Connection: Keep-Alive
User-Agent: Mozilla/4.61 [en] (Win95; I)
Host: site.org:8000
Accept: image/gif, image/x-xbitmap, image/jpeg,
image/pjpeg, image/png, */*
Accept-Encoding: gzip
Accept-Language: en
Accept-Charset: iso-8859-1,*,utf-8
Extension: Security/Remote-Passphrase
```

# HTTP response headers

- A possible response to that request:

Status Code

```
HTTP/1.1 200 OK
Date: Mon, 11 Oct 1999 18:55:17 GMT
Server: Apache/1.2.6 Red Hat
Last-Modified: Sun, 10 Oct 1999 22:45:13 GMT
Content-Length: 52
Content-Type: text/html

<html>
This is page Y.html in directory X.
</html>
```

Explanation of Status Code: ftp://ftp.isi.edu/in-notes/rfc2616.txt