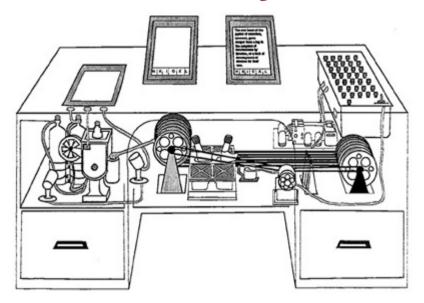
# 15-213 "The course that gives CMU its Zip!"

## Web Services November 10, 2009

#### **Topics**

- HTTP
- Serving static content
- Serving dynamic content

### **Web History**



"Consider a future device for individual use, which is a sort of mechanized private file and library. It needs a name, and to coin one at random, "memex" will do. A memex is a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility. It is an enlarged intimate supplement to his memory."

#### 1945:

- Vannevar Bush, "As we may think", Atlantic Monthly, July, 1945.
  - Describes the idea of a distributed hypertext system.
  - A "memex" that mimics the "web of trails" in our minds.

-2- 15-213, F'09

### **Web History**

#### 1989:

- Tim Berners-Lee (CERN) writes internal proposal to develop a distributed hypertext system.
  - Connects "a web of notes with links."
  - Intended to help CERN physicists in large projects share and manage information

#### 1990:

■ Tim BL writes a graphical browser for NeXT machines.

-3-

# Web History (cont)

#### 1992

- NCSA server released
- 26 WWW servers worldwide

#### 1993

- Marc Andreessen releases first version of NCSA Mosaic browser
- Mosaic version released for (Windows, Mac, Unix).
- Web (port 80) traffic at 1% of NSFNET backbone traffic.
- Over 200 WWW servers worldwide.

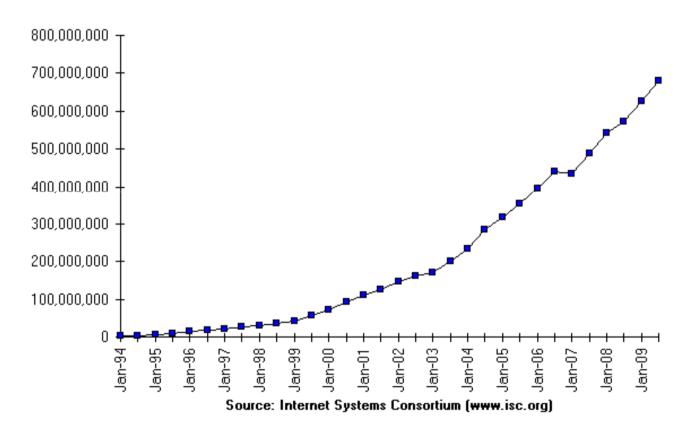
#### 1994

Andreessen and colleagues leave NCSA to form "Mosaic Communications Corp" (predecessor to Netscape).

- 4 - 15-213, F'09

### **Internet Hosts**

#### Internet Domain Survey Host Count



■ How many of the 2<sup>32</sup> IP addresses have registered names?

-5-15-213, F'09

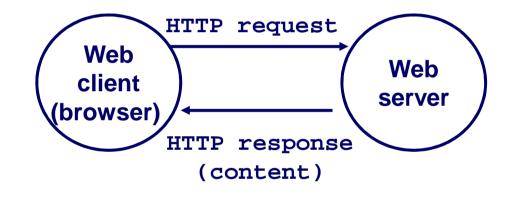
### **Web Servers**

Clients and servers
communicate using the
HyperText Transfer
Protocol (HTTP)

- Client and server establish TCP connection
- Client requests content
- Server responds with requested content
- Client and server close connection (usually)

#### **Current version is HTTP/1.1**

■ RFC 2616, June, 1999.



http://www.w3.org/Protocols/rfc2616/rfc2616.html

### **Web Content**

#### Web servers return *content* to clients

 content: a sequence of bytes with an associated MIME (Multipurpose Internet Mail Extensions) type

#### **Example MIME types**

■ text/html
HTML document

■ text/plain Unformatted text

application/postscript Postcript document

■ image/gif Binary image encoded in GIF format

■ image/jpeg Binary image encoded in JPEG

**format** 

### **MIME and Andrew**



-8-15-213, F'09

### **Static and Dynamic Content**

The content returned in HTTP responses can be either static or dynamic.

- Static content: content stored in files and retrieved in response to an HTTP request
  - Examples: HTML files, images, audio clips.
- Dynamic content: content produced on-the-fly in response to an HTTP request
  - Example: content produced by a program executed by the server on behalf of the client.

Bottom line: All Web content is associated with a "file" that is managed by the server.

- 9 - 15-213, F'09

### **URLs**

# Each file managed by a server has a unique name called a URL (Universal Resource Locator)

#### **URLs for static content:**

- http://www.cs.cmu.edu:80/index.html
- http://www.cs.cmu.edu/index.html
- http://www.cs.cmu.edu
  - Identifies a file called index.html, managed by a Web server at www.cs.cmu.edu that is listening on port 80.

#### **URLs for dynamic content:**

- http://www.cs.cmu.edu:8000/cgi-bin/adder?15000&213
  - Identifies an executable file called adder, managed by a Web server at www.cs.cmu.edu that is listening on port 8000, that should be called with two argument strings: 15000 and 213.

- 10 - 15-213, F'09

### **How Clients and Servers Use URLs**

Example URL: http://www.aol.com:80/index.html

Clients use *prefix* (http://www.aol.com:80) to infer:

- What kind of server to contact (Web server)
- Where the server is (www.aol.com)
- What port it is listening on (80)

#### Servers use suffix (/index.html) to:

- Determine if request is for static or dynamic content.
  - No hard and fast rules for this.
  - Convention: executables reside in cgi-bin directory
- **■** Find file on file system.
  - Initial "/" in suffix denotes home directory for requested content.
  - Minimal suffix is "/", which all servers expand to some default home page (e.g., index.html).

15-213, F'09

### **Anatomy of an HTTP Transaction**

```
Client: open connection to server
unix> telnet www.aol.com 80
                                        Telnet prints 3 lines to the terminal
Trying 205.188.146.23...
Connected to aol.com.
Escape character is '^]'.
GET / HTTP/1.1
                                        Client: request line
                                        Client: required HTTP/1.1 HOST header
host: www.aol.com
                                        Client: empty line terminates headers.
                                        Server: response line
HTTP/1.0 200 OK
MIME-Version: 1.0
                                        Server: followed by five response headers
Date: Mon, 08 Jan 2001 04:59:42 GMT
Server: NaviServer/2.0 AOLserver/2.3.3
                                        Server: expect HTML in the response body
Content-Type: text/html
                                        Server: expect 42,092 bytes in the resp body
Content-Length: 42092
                                        Server: empty line ("\r") terminates hdrs
                                        Server: first HTML line in response body
<html>
                                        Server: 766 lines of HTML not shown.
. . .
                                        Server: last HTML line in response body
</html>
Connection closed by foreign host. Server: closes connection
                                        Client: closes connection and terminates
unix>
```

- 12 - 15-213, F'09

### **HTTP Requests**

# HTTP request is a request line, followed by zero or more request headers

#### Request line: <method> <uri> <version>

- <version> is HTTP version of request (HTTP/1.0 or HTTP/1.1)
- <uri>is typically URL for proxies, URL suffix for servers.
  - A URL is a type of URI (Uniform Resource Identifier)
  - See http://www.ietf.org/rfc/rfc2396.txt
- method> is either GET, POST, OPTIONS, HEAD, PUT,
  DELETE, Or TRACE.

- 13 - 15-213, F'09

### **HTTP Requests (cont)**

#### **HTTP methods:**

- GET: Retrieve static or dynamic content
  - Arguments for dynamic content are in URI
  - Workhorse method (99% of requests)
- POST: Retrieve dynamic content
  - Arguments for dynamic content are in the request body
- OPTIONS: Get server or file attributes
- HEAD: Like GET but no data in response body
- PUT: Write a file to the server!
- DELETE: Delete a file on the server!
- **TRACE: Echo request in response body** 
  - Useful for debugging.

- 14 - 15-213, F'09

### **HTTP Requests (cont)**

(HTTP request is a request line, followed by zero or more request headers)

Request headers: <header name>: <header data>

■ Provide additional information to the server.

#### Major differences between HTTP/1.1 and HTTP/1.0

- HTTP/1.0 uses a new connection for each transaction.
- HTTP/1.1 also supports *persistent connections* 
  - multiple transactions over the same connection
  - Connection: Keep-Alive
- HTTP/1.1 requires ноят header
  - Host: kittyhawk.cmcl.cs.cmu.edu
- HTTP/1.1 supports chunked encoding (described later)
  - Transfer-Encoding: chunked
- HTTP/1.1 adds additional support for caching

### **HTTP Responses**

HTTP response is a *response line* followed by zero or more *response headers*.

#### Response line:

<version> <status code> <status msg>

- <version> is HTTP version of the response.
- <status code> is numeric status.
- <status msg> is corresponding English text.

• 200 OK Request was handled without error

• 403 Forbidden Server lacks permission to access file

• 404 Not found Server couldn't find the file.

#### Response headers: <header name>: <header data>

- Provide additional information about response
- Content-Type: MIME type of content in response body.
- Content-Length: Length of content in response body.

- 16 - 15-213, F'09

# **GET Request to Apache Server**From IE Browser

URI is just the suffix, not the entire URL

```
GET /test.html HTTP/1.1
Accept: */*
Accept-Language: en-us
Accept-Encoding: gzip, deflate
User-Agent: Mozilla/4.0 (compatible; MSIE 4.01; Windows 98)
Host: euro.ecom.cmu.edu
Connection: Keep-Alive
CRLF (\r\n)
```

- 17 - 15-213, F'09

### **GET Response From Apache Server**

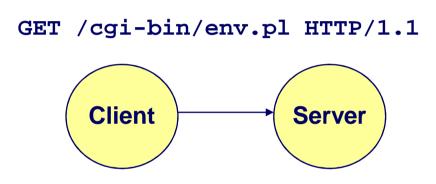
```
HTTP/1.1 200 OK
Date: Thu, 22 Jul 1999 04:02:15 GMT
Server: Apache/1.3.3 Ben-SSL/1.28 (Unix)
Last-Modified: Thu, 22 Jul 1999 03:33:21 GMT
ETag: "48bb2-4f-37969101"
Accept-Ranges: bytes
Content-Length: 79
Keep-Alive: timeout=15, max=100
Connection: Keep-Alive
Content-Type: text/html
CRLF
<html>
<head><title>Test page</title></head>
<body>
<h1>Test page</h1>
</html>
```

- 18 - 15-213, F'09

## **Serving Dynamic Content**

Client sends request to server.

If request URI contains the string "/cgi-bin", then the server assumes that the request is for dynamic content.

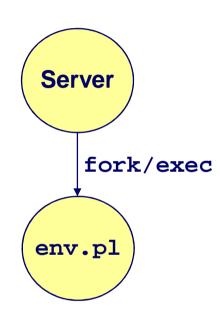


- 19 - 15-213, F'09

## **Serving Dynamic Content (cont)**

The server creates a child process and runs the program identified by the URI in that process



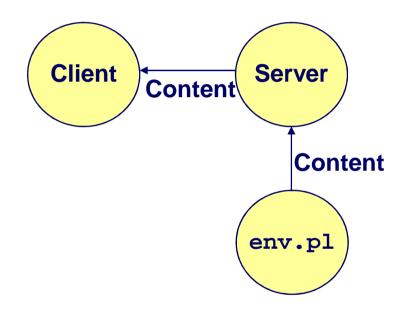


- 20 - 15-213, F'09

# **Serving Dynamic Content (cont)**

The child runs and generates the dynamic content.

The server captures the content of the child and forwards it without modification to the client



-21 - 15-213, F'09

### **Issues in Serving Dynamic Content**

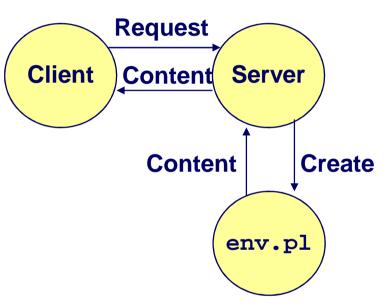
How does the client pass program arguments to the server?

How does the server pass these arguments to the child?

How does the server pass other info relevant to the request to the child?

How does the server capture the content produced by the child?

These issues are addressed by the Common Gateway Interface (CGI) specification.



- 22 - 15-213, F'09

### **CGI**

Because the children are written according to the CGI spec, they are often called CGI programs.

Because many CGI programs are written in Perl, they are often called *CGI scripts*.

However, CGI really defines a simple standard for transferring information between the client (browser), the server, and the child process.

- 23 - 15-213, F'09

# add.com: THE Internet addition portal!

Ever need to add two numbers together and you just can't find your calculator?

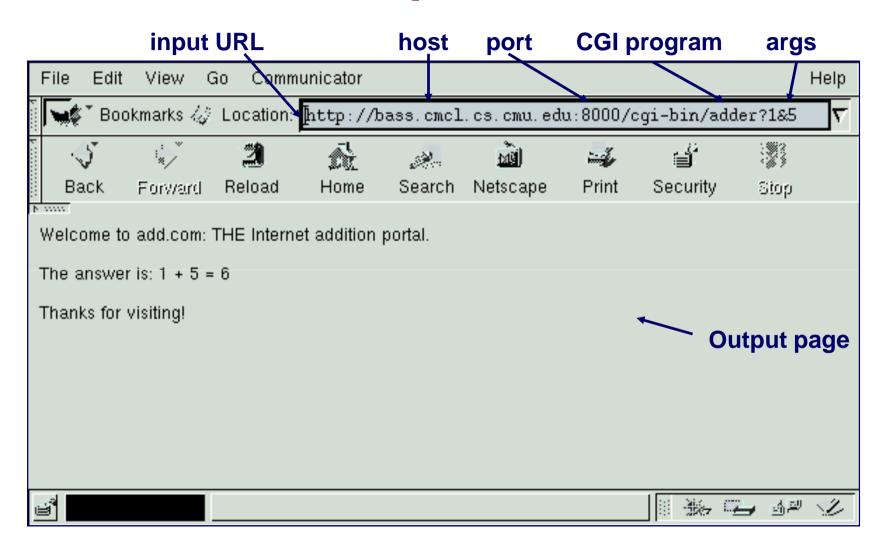
Try Dr. Dave's addition service at "add.com: THE Internet addition portal!"

- Takes as input the two numbers you want to add together.
- Returns their sum in a tasteful personalized message.

After the IPO we'll expand to multiplication!

- 24 - 15-213, F'09

### The add.com Experience



- 25 - 15-213, F'09

**Question:** How does the client pass arguments to the server?

**Answer: The arguments are appended to the URI** 

Can be encoded directly in a URL typed to a browser or a URL in an HTML link

- http://add.com/cgi-bin/adder?1&2
- adder is the CGI program on the server that will do the addition.
- argument list starts with "?"
- arguments separated by "&"
- spaces represented by "+" or "%20"

#### Can also be generated by an HTML form

#### **URL**:

■ http://add.com/cgi-bin/adder?1&2

#### Result displayed on browser:

Welcome to add.com: THE Internet addition portal.

The answer is: 1 + 2 = 3

Thanks for visiting!

# **Question:** How does the server pass these arguments to the child?

#### **Answer: In environment variable QUERY\_STRING**

- A single string containing everything after the "?"
- For add.com: QUERY\_STRING = "1&2"

```
/* child code that accesses the argument list */
if ((buf = getenv("QUERY_STRING")) == NULL) {
   exit(1);
}

/* extract arg1 and arg2 from buf and convert */
...
n1 = atoi(arg1);
n2 = atoi(arg2);
```

- 28 - 15-213, F'09

**Question:** How does the server pass other info relevant to the request to the child?

Answer: In a collection of environment variables defined by the CGI spec.

- 29 - 15-213, F'09

### Some CGI Environment Variables

#### General

- SERVER\_SOFTWARE
- SERVER\_NAME
- GATEWAY\_INTERFACE (CGI version)

#### Request-specific

- SERVER\_PORT
- REQUEST\_METHOD (GET, POST, etc)
- QUERY\_STRING (contains GET args)
- REMOTE\_HOST (domain name of client)
- REMOTE\_ADDR (IP address of client)
- CONTENT\_TYPE (for POST, type of data in message body, e.g., text/html)
- CONTENT\_LENGTH (length in bytes)

### Some CGI Environment Variables

In addition, the value of each header of type *type* received from the client is placed in environment variable HTTP\_type

- **Examples:** 
  - HTTP ACCEPT
  - HTTP HOST
  - HTTP\_USER\_AGENT (any "-" is changed to "\_")

- 31 - 15-213, F'09

**Question:** How does the server capture the content produced by the child?

Answer: The child generates its output on stdout. Server uses dup2 to redirect stdout to its connected socket.

■ Notice that only the child knows the type and size of the content. Thus the child (not the server) must generate the corresponding headers.

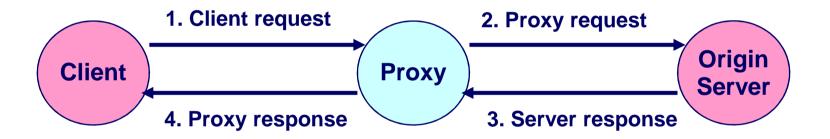
- 32 - 15-213, F'09

```
bass> ./tiny 8000
                                           HTTP request received by
   GET /cgi-bin/adder?1&2 HTTP/1.1
                                           Tiny Web server
   Host: bass.cmcl.cs.cmu.edu:8000
   <CRLF>
   kittyhawk> telnet bass 8000
   Trying 128.2.222.85...
   Connected to BASS.CMCL.CS.CMU.EDU.
   Escape character is '^]'.
   GET /cgi-bin/adder?1&2 HTTP/1.1
   Host: bass.cmcl.cs.cmu.edu:8000
                                           HTTP request sent by client
   <CRLF>
                                           HTTP response generated by
   HTTP/1.1 200 OK
   Server: Tiny Web Server
                                           the server
   Content-Tength: 102
                                           HTTP response generated by
   Content-type: text/html
                                           the CGI program
   <CRLF>
   Welcome to add.com: THE Internet addition portal.
   The answer is: 1 + 2 = 3
   Thanks for visiting!
   Connection closed by foreign host.
   kittyhawk>
-33-
                                                             15-213, F'09
```

### **Proxies**

# A proxy is an intermediary between a client and an origin server.

- To the client, the proxy acts like a server.
- To the server, the proxy acts like a client.

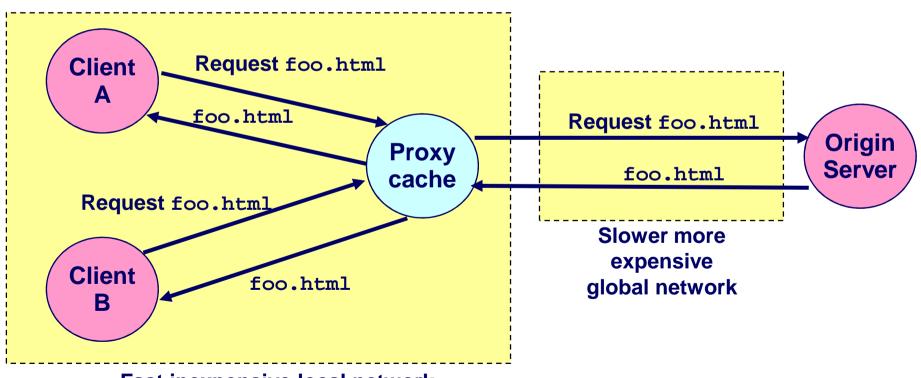


- 34 - 15-213, F'09

## Why Proxies?

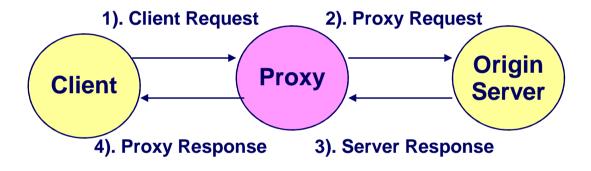
# Can perform useful functions as requests and responses pass by

Examples: Caching, logging, anonymization, filtering, transcoding



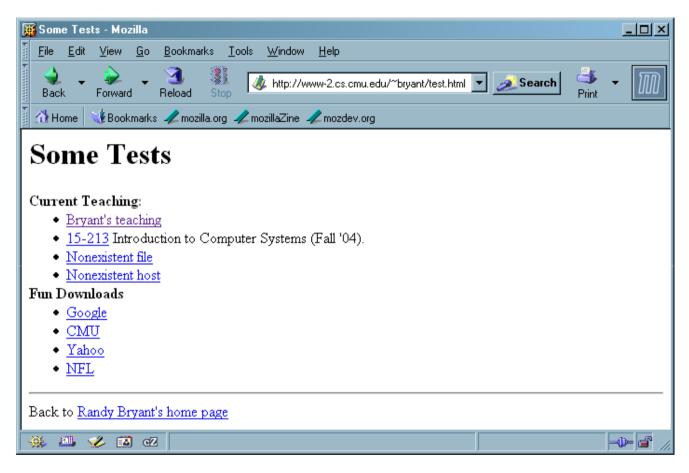
**Fast inexpensive local network** 

# Putting it Together: Web Proxy Demonstration



- 36 - 15-213, F'09

# Servicing Web Page Request



- 37 - 15-213, F'09

### Client → Proxy

### The browser sends a URI that is a complete URL

```
GET http://www-2.cs.cmu.edu/~bryant/test.html HTTP/1.1\r\n
Host: www-2.cs.cmu.edu\r\n
User-Agent: Mozilla/5.0 (Windows; U; Windows NT 5.1; en-US; rv:1.7.3)
    Gecko/20040910\r\n
Accept:
    text/xml,application/xml,application/xhtml+xml,text/html;q=0.9,tex
    t/plain;q=0.8,image/png,*/*;q=0.5\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,*;q=0.7\r\n
Keep-Alive: 300\r\n
Proxy-Connection: keep-alive\r\n
\r\n
```

- 38 - 15-213, F'09

### Proxy → Server

The proxy sends a URI that is a path

```
GET /~bryant/test.html HTTP/1.1\r\n
Host: www-2.cs.cmu.edu\r\n
User-Agent: Mozilla/5.0 (Windows; U; Windows NT 5.1; en-US; rv:1.7.3)
    Gecko/20040910\r\n
Accept:
    text/xml,application/xml,application/xhtml+xml,text/html;q=0.9,tex
    t/plain;q=0.8,image/png,*/*;q=0.5\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,*;q=0.7\r\n
Keep-Alive: 300\r\n
Connection: keep-alive\r\n
\r\n
```

- 39 - 15-213, F'09

### Server → Proxy → Client

```
HTTP/1.1 200 OK\r\n
Date: Mon, 29 Nov 2004 01:27:15 GMT\r\n
Server: Apache/1.3.27 (Unix) mod_ssl/2.8.12 OpenSSL/0.9.6
   mod pubcookie/a5/1.76-009\r\n
Transfer-Encoding: chunked\r\n
Content-Type: text/html\r\n
\r\n
```

#### **Chunked Transfer Encoding**

- Alternate way of specifying content length
- Each "chunk" prefixed with chunk length
- See http://www.w3.org/Protocols/rfc2616/rfc2616-sec3.html

- 40 - 15-213, F'09

# Server → Proxy → Client (cont)

```
First Chunk: 0x2ec = 748 bytes
2ec\r\n
<head><title>Some Tests</title></head>\n
<h1>Some Tests</h1>\n
<dl>\n
 <dt> <strong>Current Teaching: </strong>\n
 <u1>\n
 <a href="teaching.html">Bryant's teaching</a>\n</a>
  <a href="/afs/cs.cmu.edu/academic/class/15213-f04/www/">\n</a>
   15-213</a> Introduction to Computer Systems (Fall '04).\n
  <a href="http://www.cs.cmu.edu/nothing.html">Nonexistent file</a>\n
  <a href="http://nowhere.cmu.edu/nothing.html">Nonexistent host</a>\n

 <dt><strong>Fun Downloads</strong>\n
 <ul>\n
 <a href="http://www.google.com">Google</a>\n
 <a href="http://www.cmu.edu">CMU</a>\n
 <a href="http://www.yahoo.com">Yahoo</a>\n
 <a href="http://www.nfl.com">NFL</a>\n

</dl>

< hr > \ n
Back to <a href="index.html">Randy Bryant's home page</a>\n
n
\r\
         Second Chunk: 0 bytes (indicates last chunk)
0\r\n
r n
```

### For More Information

#### Study the Tiny Web server described in your text

- Tiny is a sequential Web server.
- Serves static and dynamic content to real browsers.
  - text files, HTML files, GIF and JPEG images.
- 220 lines of commented C code.
- Also comes with an implementation of the CGI script for the add.com addition portal.

#### See the HTTP/1.1 standard:

http://www.w3.org/Protocols/rfc2616/rfc2616.html

- 42 - 15-213, F'09