

UNIVERSIDADE DA CORUÑA
Departamento de Tecnoloxías da Información
e as Comunicacións

LABORATORIO DE RC:
PRÁCTICA 1: IMPLEMENTACIÓN DE UN SERVIDOR WEB



PRÁCTICA 3: Implementaci3n de un Servidor Web

Esta pr3ctica consiste en el desarrollo de un servidor Web multi-thread, capaz de procesar m3ltiples solicitudes de servicio simult3neas en paralelo. El objetivo final ser3 obtener un servidor Web capaz de interactuar con un cliente para navegar a trav3s de un sitio Web.

Para la realizaci3n del servidor Web se utilizar3 el protocolo HTTP versi3n 1.0 (definido en el RFC 1945), en donde se generan solicitudes HTTP independientes para cada componente de la p3gina Web. El servidor procesar3 m3ltiples solicitudes en paralelo, utilizando para ello m3ltiples hilos de ejecuci3n.

En el hilo principal, el servidor permanece escuchando en un puerto fijo y cuando se recibe una petici3n de conexi3n TCP, se establece la conexi3n a trav3s de otro puerto y se resuelve el servicio en un hilo de ejecuci3n separado.

Para la realizaci3n de la pr3ctica se recomienda una aproximaci3n en dos etapas. En una primera etapa, el servidor multi-thread 3nicamente mostrar3 por pantalla el contenido de las solicitudes HTTP recibidos, para a continuaci3n (una vez que el programa se ejecuta convenientemente) incluir el c3digo que genere la respuesta adecuada. Para comprobar el correcto funcionamiento del servidor se recomienda utilizar el comando `nc` para conectarse al puerto de escucha, e introducir directamente por l3nea de comandos las 3rdenes HTTP.

El servidor deber3 lanzarse en un puerto no reservado (i.e. por encima del puerto 1024), y en el cliente se le deber3 indicar el fichero que se est3 solicitando. A modo de ejemplo, una URL que se podr3 solicitar desde el navegador tendr3 el siguiente formato: <http://localhost:1111/index.html>. El servidor Web localizar3 el fichero solicitado y ser3 devuelto al cliente en el formato adecuado. En caso de que se produzca alg3n error (por ejemplo, el fichero solicitado no se encuentra), el servidor deber3 responder con el c3digo HTML adecuado para cada situaci3n err3nea.

La versi3n b3sica de esta pr3ctica debe incluir:

- Implementaci3n de un servidor Web multi-thread.
- Implementaci3n de los m3todos GET y HEAD.
- Soporte de los formatos b3sicos HTML, texto plano, GIF y JPEG.

Adem3s, dentro de la pr3ctica se incluyen las siguientes opciones:

- Implementaci3n de la opci3n “If-modified-since” para el m3todo GET, que se deber3 activar a trav3s del navegador Web mediante la recarga de la p3gina (bot3n “Actualizar”). De esta manera, el navegador solicitar3 al servidor Web la p3gina Web con la opci3n “If-modified-since”, descarg3ndola si ha habido cambios. Para verificar el correcto



- funcionamiento de esta opci3n se recomienda utilizar el complemento “*Tamper Data*” de Firefox.
- Definici3n de dos ficheros de log: para accesos y para errores.
 - o El fichero de log para accesos tendr3 el siguiente formato:
 - Direcci3n IP del cliente que realiza la solicitud.
 - Fecha y hora en la que se recibió la petici3n: [día/mes/año hora:minuto:segundo zona_horaria]
 - Línea de petici3n del cliente (entre comillas dobles).
 - C3digo de estado que el servidor envía como respuesta al cliente.
 - Tamaño (en bytes) del objeto enviado al cliente.
 - o El fichero de log para errores tendr3 el siguiente formato:
 - Fecha y hora en la que se produjo el error (mismo formato que para el fichero de log de accesos).
 - Direcci3n IP del cliente que gener3 el error.
 - Mensaje de error.
 - Petici3n que provoc3 el error.
 - Configuraci3n b3sica del servidor Web: el servidor Web dispondr3 de un **fichero de configuraci3n** a trav3s del cual se podr3n indicar, al menos, los siguientes par3metros (se recomienda utilizar la clase `java.util.Properties` para el desarrollo de esta opci3n):
 - o Par3metro PORT: Puerto de inicio
 - o Par3metro DIRECTORY_INDEX: Nombre del fichero a utilizar por defecto, en caso de no incluirse en la petici3n (por ejemplo, <http://localhost:1111/>)
 - o Par3metro DIRECTORY: Directorio raíz del servidor Web, en donde se ubicar3n las p3ginas Web. Se incluir3 la directiva ALLOW, que funcionar3 de la siguiente manera:
 - Si existe el fichero por defecto, se mostrar3 el fichero por defecto.
 - Si est3 activada y no existe fichero por defecto en ese directorio, se mostrar3 el listado de todos los archivos/directorios que lo componen, con un enlace que permitir3 abrir cada fichero/directorio.
 - Si no est3 activada y no existe el fichero por defecto, se mostrar3 el error 403 (Access forbidden).



- Ejecuci3n de CGIs (*Common Gateway Interface*): el servidor Web deber3 crear un subproceso y ejecutar un programa especificado en la petici3n, recibiendo como par3metros de entrada los datos enviados por el navegador a trav3s del m3todo GET. Ser3 necesario capturar la salida est3ndar del proceso y el resultado ser3 enviado directamente al navegador. Para comprobar el correcto funcionamiento de esta opci3n se utilizar3 un **formulario** cuyos datos ser3n procesados por un CGI que mostrar3 el resultado por pantalla. Para esta opci3n se dispondr3 de un `saludo.html` y de un `saludoCGI.java`, si bien tambi3n se puede optar por utilizar un formulario y un CGI propios.

Ser3 necesario definir un directorio espec3fico para la localizaci3n de los CGIs con permisos de ejecuci3n.

Adem3s, la pr3ctica deber3 ser realizada utilizando las funcionalidades propias de los sockets en Java, sin utilizar ninguna clase que implemente total o parcialmente las caracter3sticas del protocolo HTTP (como por ejemplo, la clase `URLConnection`).

Introducci3n a HTTP

El protocolo HTTP (HyperText Transfer Protocol) est3 especificado en el RFC 1945 (versi3n 1.0) y en el RFC 2616 (versi3n 1.1). Este protocolo define c3mo los clientes (navegadores) solicitan p3ginas Web a los servidores Web, y c3mo 3stos realizan la transferencia de estas p3ginas.

El protocolo HTTP se basa en el protocolo TCP (que ofrece un servicio orientado a conexi3n y fiable), lo que garantiza que cada mensaje HTTP emitido por el cliente o el servidor es recibido en el otro extremo sin sufrir modificaciones. Adem3s, HTTP es un protocolo sin estado, esto es, el servidor HTTP no almacena ning3n tipo de informaci3n sobre sus clientes. Cada petici3n recibida por el servidor se trata independientemente de las peticiones anteriormente recibidas de ese u otros clientes.

El protocolo HTTP 1.0 utiliza conexiones no persistentes, ya que es necesario establecer una conexi3n TCP independiente para cada uno de los componentes de una p3gina Web. De manera esquem3tica, el procedimiento para la petici3n de una URL (p.e. <http://www.tic.udc.es/index.html>) ser3 el siguiente:

1. El cliente HTTP inicia la conexi3n TCP con el servidor `www.tic.udc.es` en el puerto 80.
2. El cliente HTTP env3a al servidor el mensaje de petici3n solicitando el objeto `/index.html`
3. El servidor HTTP recibe la petici3n, lee el objeto, lo encapsula en el mensaje HTTP de respuesta y lo env3a.
4. El servidor finaliza la conexi3n TCP.

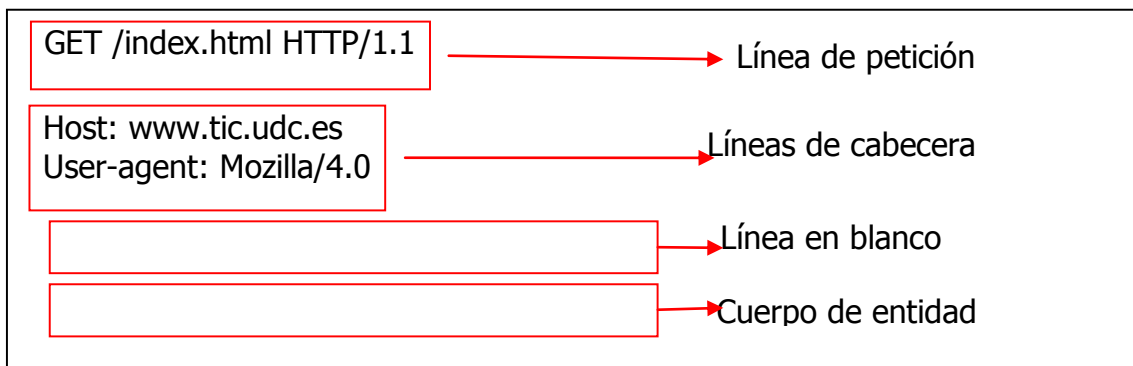


5. El cliente HTTP recibe la respuesta y finaliza la conexi3n TCP.
6. El cliente extrae el archivo del mensaje de respuesta, examina el archivo HTML y encuentra referencias a otros objetos HTML (p.e. imágenes)
7. Volver al paso 1, para cada uno de los nuevos objetos HTML.

El principal inconveniente de esta aproximaci3n se centra en el retardo introducido, ya que es necesario esperar dos veces el RTT (Round Trip Time), uno para el establecimiento de la conexi3n y otro para la petici3n y respuesta del objeto solicitado. Además, el hecho de utilizar una conexi3n para cada objeto solicitado implica una mayor utilizaci3n de recursos (buffers, variables, timeouts, ...), tanto en el cliente como en el servidor.

Estos problemas se resuelven mediante el protocolo HTTP 1.1 que utiliza conexiones persistentes, en donde el servidor deja abierta las conexiones TCP establecidas en espera de nuevas peticiones. Tras un per3odo de inactividad estas conexiones se cierran por parte del servidor.

El protocolo HTTP sigue un sencillo modelo de peticiones y respuestas. El formato de una petici3n se puede observar en la siguiente figura. Los únicos campos obligatorios son la l3nea de petici3n y la l3nea en blanco.



La l3nea de petici3n especifica el tipo de petici3n que se est3 realizando, y est3 formada por tres campos:

- M3todo
- URL: objeto al que se hace referencia.
- Versi3n: del protocolo HTTP utilizada por el navegador.

Los m3todos b3sicos definidos en el protocolo HTTP son los siguientes:

- GET: solicitud de un objeto por parte del cliente.
- HEAD: se solicita al servidor que responda con un mensaje HTTP, aunque sin incluir el objeto en la respuesta.
- POST: permite incluir datos en el cuerpo de entidad.
- PUT: permite a un cliente cargar un archivo en la ruta especificada (s3lo en HTTP 1.1).



- DELETE: permite a un cliente borrar un archivo de un servidor (sólo en HTTP 1.1).

Dentro las líneas de cabecera, un navegador típico puede incluir múltiples opciones, siendo las más relevantes:

- Host: especifica el host en el que reside el objeto solicitado.
- User-agent: especifica el tipo de navegador que está haciendo la petici3n.
- If-modified-since: utilizada junto con el método GET, especifica una fecha para que el servidor, si el objeto no ha sido modificado con posterioridad a esa fecha, no lo envíe de nuevo.

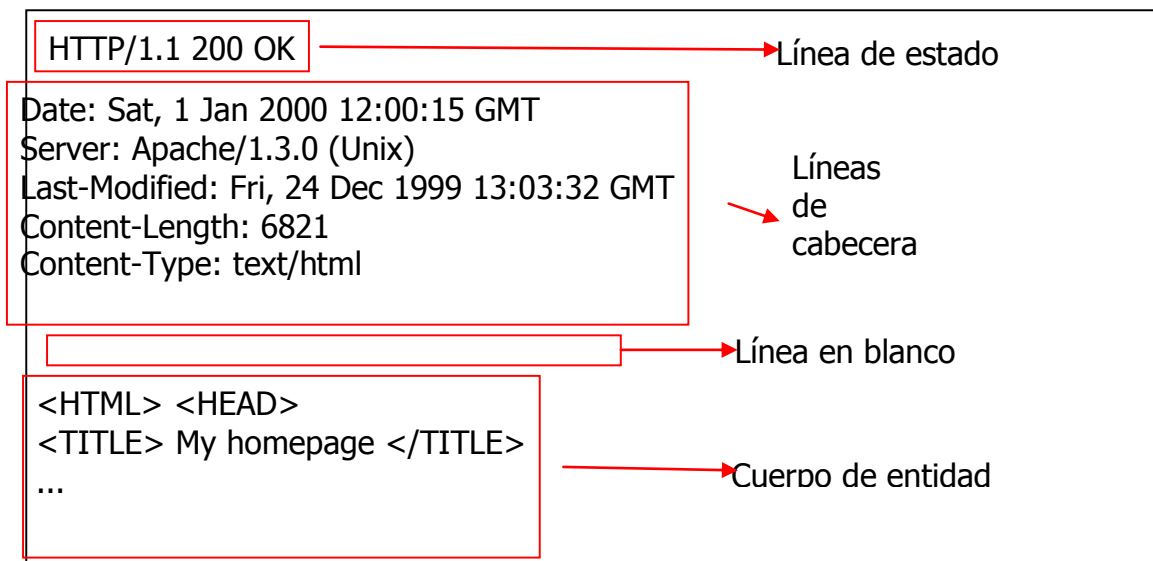
La opci3n “*If-modified-since*” es utilizada en las implementaciones de las cachés en el cliente. El formato de la petici3n sería el siguiente:

```
GET /images/udc.gif HTTP/1.1
User-agent: Mozilla/4.0
If-modified-since: Fri, 24 Dec 1999 13:03:32 GMT
```

Mientras que en el servidor, si el objeto no ha sido modificado, la respuesta sería la siguiente:

```
HTTP/1.1 304 Not Modified
Date: Wed, 5 Jan 2000 20:30:43 GMT
Server: Apache/1.3.0 (Unix)
```

El formato de las respuestas HTTP se puede observar en la siguiente figura.





La lnea de estado incluye tres campos:

- Versi3n: versi3n utilizada por el servidor Web.
- C3digo de estado: c3digo num3rico que representa si la respuesta ante la petici3n es satisfactoria o si ha habido alg3n error.
- Frase: asociado a cada c3digo num3rico existe una frase que informa sobre la naturaleza del c3digo. En el RFC 1945 se pueden consultar todos los c3digos y frases del protocolo, destacando los siguientes:
 - o 200 OK
 - o 400 Bad Request: petici3n no comprendida por el servidor.
 - o 403 Forbidden: petici3n comprendida por el servidor, pero que rechaza satisfacer.
 - o 404 Not Found: el objeto solicitado no existe en el servidor.

Respecto a las lneas de cabecera, existen m3ltiples par3metros que el servidor puede especificar, destacando los siguientes:

- Date: fecha y hora en la que se cre3 y envi3 la respuesta HTTP.
- Server: especifica el tipo de servidor Web que ha atendido la petici3n.
- Last-Modified: indica la fecha y hora en que el objeto fue creado o modificado por 3ltima vez.
- Content-Length: indica el n3mero de bytes del objeto enviado.
- Content-Type: indica el tipo de objeto incluido en el cuerpo de entidad. Este campo es necesario, ya que la extensi3n del archivo no especifica (formalmente) el tipo de objeto asociado. Los tipos m3s com3nmente utilizados son:
 - o text/html: indica que la respuesta est3 en formato HTML.
 - o text/plain: indica que la respuesta est3 en texto plano.
 - o image/gif: indica que se trata de una imagen en formato gif.
 - o image/jpeg: indica que se trata de una imagen en formato jpeg.
 - o application/octet-stream: utilizado cuando no se identifica el formato del archivo.

FECHA DE ENTREGA

La pr3ctica 3 se iniciar3 la semana del 18 al 22 de Octubre de 2010. La defensa de la pr3ctica se realizar3, **como m3ximo**, la semana del 15 al 19 de Noviembre de 2010. Para ello, cada alumno estar3 apuntado en uno de los grupos de pr3cticas de la asignatura y la pr3ctica ser3 presentada y defendida en el d3a de



prácticas habitual. En caso de hacerlo en un grupo posterior se podrá considerar que la práctica ha sido entregada fuera de plazo.

Para la presentación de esta práctica, se mostrará su funcionamiento en el laboratorio de prácticas y el alumno deberá ser capaz de explicar el funcionamiento de cualquier componente de la misma.

Además, el código fuente de la práctica deberá ser copiado al repositorio de entrega de prácticas una vez finalizada su defensa. En caso contrario, se considerará que la práctica NO ha sido presentada.