

Índice

- 1 Tema 1: Internet
- 2 Tema 2: La web
- 3 Tema 3: Principios de Recuperación de Información
- 4 Tema 4: Búsqueda de información en la web
- 5 Tema 5: Búsqueda avanzada en la web
- 6 Tema 6: Integración de conocimiento lingüístico
- 7 Tema 7: Más allá de la búsqueda textual**
- 8 Pruebas de evaluación

Tipos de archivos en Internet

- Archivos ejecutables (*software*):
 - Archivos creados para funcionar por sí mismos (programas fuente, programas compilados).
 - Archivos que sirven de apoyo a algún archivo ejecutable (librerías, extensiones de aplicaciones, controladores de dispositivos, etc.).
- Archivos no ejecutables:
 - Archivos de datos (documentos, bases de datos, etc.).
 - Archivos multimedia (audio, imágenes, vídeo, etc.).

Licencias de uso y distribución

- Una *licencia de software* es la autorización que concede el titular del derecho de autor de un programa informático al usuario del mismo, para utilizar éste conforme a unas condiciones convenidas.
- La licencia puede ser gratuita o de pago. Debe señalar los derechos de uso, modificación y distribución concedidos a la persona autorizada. Además puede incluir un plazo de duración, un territorio de aplicación, y otra serie de cláusulas que el titular del derecho de autor establezca.
- Ejemplo: licencia GPL
 - Original en inglés: <http://www.gnu.org/licenses/gpl.html>
 - Traducción: <http://www.viti.es/gnu/licenses/gpl.html>
- Clasificación de las licencias de software:
 - Según su destinatario:
 - Licencia de usuario final.
 - Licencia de distribuidor.

Licencias de uso y distribución (cont.)

- Clasificación de las licencias de software (cont.):
 - Según los derechos que cada autor se reserve sobre su obra:
 - BSD, Apache, LGPL: uso gratis, permite copia, proporciona código fuente, permite modificación sin protección heredada (se puede redistribuir con una licencia distinta y más restrictiva).
 - GPL, Mozilla: uso gratis, permite copia, proporciona código fuente, permite modificación con protección heredada ("*copyleft*": se puede redistribuir, pero con la misma licencia).
 - Shared source (Microsoft, Apple): no permite uso, no permite copia, sí proporciona código fuente, no permite redistribución.
 - Freeware: uso gratis, permite copia, no proporciona código fuente, no permite redistribución.
 - Shareware: uso gratis parcial (hay partes gratis y partes de pago), permite copia parcial, no proporciona código fuente, no permite redistribución.
 - Software propietario: uso pagando, no permite copia, no proporciona código fuente, no permite redistribución.

Servidores FTP

- FTP (*File Transfer Protocol*) es un protocolo de transferencia de archivos entre sistemas conectados a una red TCP, que utiliza los puertos de red 20 y 21, y que está basado en una arquitectura cliente-servidor, de manera que desde un equipo cliente nos podemos conectar a un servidor para descargar archivos desde él, o para enviarle nuestros propios archivos, independientemente del sistema operativo utilizado en cada equipo.
- Pensado para ofrecer la mayor velocidad en la conexión, pero no la mayor seguridad: todo el intercambio de datos (incl. login/passwd) va en texto plano sin cifrar
 - Solución: aplicaciones tipo `sftp` o `scp`, que sí cifran el tráfico

Servidores FTP (cont.)

- Modos de acceso del cliente FTP:
 - Acceso anónimo (*anonymous*): acceso público pero restringido (p.ej. sólo lectura)
 - Acceso de usuario (*login* y *password*): acceso privado y total
 - Acceso de invitado (*guest*): acceso privado pero restringido
- Tipos de transferencia de archivos en FTP:
 - Tipo *ascii*: para archivos que sólo contengan caracteres imprimibles.
 - Tipo *binary*: para archivos comprimidos, ejecutables, imágenes, audio, video, etc.
- Con la llegada de Internet y de los navegadores modernos, ya no es necesario conocer los complejos comandos FTP. Este protocolo se puede utilizar escribiendo la URL del servidor al que queramos conectar el navegador web, pero con el prefijo `ftp://` o `sftp://`, en lugar de `http://`.

Redes *peer-to-peer* o P2P

- Una red informática *peer-to-peer* o P2P (de punto a punto, o entre pares/iguales), es una red que no tiene clientes ni servidores fijos, sino una serie de nodos que se comportan simultáneamente como clientes y como servidores de los demás nodos de la red.
- Las redes P2P aprovechan, administran y optimizan el uso de banda ancha que acumulan todos los usuarios de la red, obteniendo en algunos casos mejores rendimientos en las conexiones y transferencias que algunos métodos centralizados convencionales.
 - La eficacia de cada nodo depende de su configuración y de su conexión local.

Redes *peer-to-peer* o P2P (cont.)

- Suelen estar basadas en una filosofía meritocrática:
 - Todos los usuarios deben compartir recursos
 - "El que más comparta, más privilegios y mejor acceso tiene a los contenidos"
- Aplicaciones de las redes P2P: en general en aplicaciones y servicios que requieran de una ingente cantidad de recursos (almacenamiento, ancho de banda, capacidad de procesamiento, etc.)
 - Intercambio de archivos entre usuarios (texto, datos, imágenes, audio, vídeo). P.ej. *Napster*, *eDonkey* (sobre la que trabaja *eMule*), *BitTorrent*.
 - Mejorar la distribución el tráfico de la telefonía basada en *voz por IP* (*Skype*).
 - Alternativa a la distribución convencional de películas y programas de televisión (desde 2006, Warner Bros. y BBC mediante *BitTorrent*).
 - Cálculos científicos que procesen enormes cantidades de datos (procesamiento de señales, aplicaciones bioinformáticas, ...).

Redes *peer-to-peer* o P2P (cont.)

- Características deseables:
 - Escalabilidad: para que cuantos más nodos se conecten y compartan, mejor sea el funcionamiento de la red pues los recursos totales aumentan.
 - Robustez ante fallo: por su naturaleza distribuida y redundante nadie es imprescindible
 - Descentralización: no existen nodos con funciones especiales ni imprescindibles
 - Algunas redes P2P no lo cumplen totalmente. P.ej. eDonkey2000, donde existen *servidores* para la gestión de los recursos compartidos
 - Reparto de costes: se proporcionan recursos a cambio de recursos
 - Seguridad: actualmente poco implementada
 - Cifrado, gestión de derechos de autor, comunicaciones seguros, protección contra virus ...

Aspectos legales de las redes P2P: controversia legal

- Buena parte de los archivos compartidos que se pueden descargar en estas redes son archivos de música y vídeo con contenidos sujetos a propiedad intelectual y derechos de autor, perjudicando a los depositarios de los mismos.
 - Problema: falta de regulación legal
- Sin embargo, también se intercambia gran cantidad de contenidos privados o no sujetos a derechos de autor, así como obras cuyos autores no han prohibido dichos intercambios (por ejemplo, distribuciones Linux sujetas a la licencia GPL).
- De hecho, existen aplicaciones específicas de redes P2P directamente orientadas al intercambio de este tipo de contenidos y obras, como por ejemplo Skype (voz por IP) o Picasa (álbumes de fotos personales), etc.

Búsqueda de personas. Páginas Blancas.

- Internet no es un club que posea una lista de socios. Poder conectarse a Internet no implica que haya que estar incluido en ninguna lista mundial de clientes usuarios.
- No obstante, existen en Internet unos lugares denominados “páginas blancas” en las que muchas personas se inscriben para que otros usuarios de Internet puedan localizarlas.
- El número de personas en Internet con e-mail y página web es muy superior al número de personas que aparecen en estos índices, pero buscar en ellos puede ser útil.
- Algunas páginas blancas útiles:
 - WhoWhere (<http://www.whowhere.com>)
 - Internet Address Finder (<http://www.iaf.net>)
 - Páginas blancas de Telefónica (<http://www.paginasblancas.es>)
- En la actualidad las denominadas “**redes sociales**” (*Facebook, Tuenti*) y las “**redes profesionales**” (*LinkedIn*) actúan como puntos de encuentro que, gracias a su progresiva implantación, permiten cubrir también esta funcionalidad

Páginas amarillas

- Y en Internet es posible también localizar información sobre empresas y entidades.
- Algunas “páginas amarillas” útiles:
 - Páginas amarillas de Telefónica (<http://www.paginasamarillas.es>)
 - Páginas amarillas internacionales (<http://www.paginasamarillas.com>)
 - Índice de las guías fax del mundo (<http://www.infobel.com/World>)

Búsquedas multimedia: audio, imágenes, video

- Audio:
 - Canciones, mp3, radio.
- Imágenes:
 - Búsqueda de imágenes en Google (<http://images.google.es>)
 - Galerías de fotos compartidas (<http://www.flickr.com>)
 - Fotos geolocalizadas (<http://www.panoramio.com>, ahora integrado en Google Maps)
- Vídeo:
 - Búsqueda de vídeos en Google (<http://video.google.com>)
 - Motores de búsqueda de vídeos (<http://www.purevideo.com>, <http://www.blinkx.com>)
 - Compartición de vídeos (<http://www.youtube.com>, <http://www.dailymotion.com>)
- Wikimedia Commons (<http://commons.wikimedia.org>):
repositorio de imágenes, vídeos y audios de libre uso y distribución

Sitios de mapas

- Necesitamos saber:
 - ¿Cómo ir en coche desde Carballo a Santiago?
 - ¿Cómo buscar restaurantes japoneses en La Coruña?
 - ¿Cómo conseguir un plano detallado de Ferrol?
- Los sitios de mapas más populares son:
 - MapQuest (<http://www.mapquest.com>). Más antiguo. Interfaz básica. Ofrece mapas y direcciones, pero no mucho más.
 - Google Maps (<http://maps.google.com>, <http://maps.google.es>). Proporciona información local sobre los mapas (restaurantes, museos, farmacias, etc.). Interfaz interactiva que permite zoom y combinación con imágenes de satélite. Planifica y calcula rutas entre diferentes puntos de origen y destino.
 - "*Street View*": nuevo servicio que permite recorrer virtualmente una ciudad a pie de calle

Sitios de mapas (cont.)

- Los sitios de mapas no crean sus propios mapas, ni la información subyacente a ellos. Obtienen todo esto a partir de determinadas agencias gubernamentales, y sobre todo a partir de proveedores comerciales de información sobre mapas (eg. Navteq, más de 600 investigadores, a lo largo de 23 países). Estos proveedores venden sus bases de datos no sólo a estos sitios, sino también a empresas particulares que necesiten este tipo de información.

Google Earth

- Google ofrece una herramienta más sofisticada que los simples mapas: Google Earth.
- Google Earth permite “volar” de forma virtual a cualquier lugar de la tierra, utilizando fotografías y animaciones de alta resolución.
- Fue inicialmente creado por la empresa Keyhole (comprada y rebautizada por Google).
- Las imágenes se obtienen de muchas formas diferentes utilizando satélites, cámaras giroscópicas instaladas en aviones volando a alturas entre 5000 y 9000 metros, y en algunos casos incluso mediante globos y cometas.
- Google Earth añade sobre las fotografías las capas de información sobre fronteras, calles, parques, escuelas, museos, restaurantes, hospitales, etc.
- Para utilizar Google Earth es necesario descargar e instalar previamente el software específico de Google Earth.

Google Desktop

- A veces, la información que necesitamos está en nuestro propio escritorio. Google Desktop (<http://desktop.google.es>) realiza búsquedas en nuestro equipo tan fácilmente como lo hace Google en la web.
- Se trata de una aplicación de búsqueda en el escritorio que permite localizar documentos de texto, hojas de cálculo, música, fotografías, texto en correos electrónicos, conversaciones chat, páginas web visitadas, archivos suprimidos, etc.
- Es capaz de generar vistas previas de los resultados, sin necesidad de esperar a que una aplicación los abra, simplemente para verificar que se trata del archivo que estábamos buscando.
- Es capaz también de sincronizarse con la web, con el fin de recopilar y organizar nueva información entrante, mediante los Google Gadgets: conjunto de miniaplicaciones interactivas que avisan de la llegada de correo, de las previsiones metereológicas, de la publicación de noticias personalizadas, etc.

Google Desktop (cont.)

- Incorpora un sistema de indexación inteligente: En cuanto se instala el programa, Google Desktop empieza a indexar los archivos, los mensajes de correo y los historiales web. Aunque esta indexación inicial puede consumir horas, se realiza cuando el equipo está inactivo durante más de 30 segundos, por lo que no debería ralentizar su rendimiento. Posteriormente, el sistema se asegura de que nuestro índice se mantenga actualizado con la nueva información que nos va llegando o que nosotros mismos vamos generando.
- En resumen, se trata de un auténtico mini-sistema de recuperación de información en nuestro propio ordenador.