

Buscadores Web

Referencias:

<http://searchenginewatch.com/>

<http://searchengineshowdown.com/>

Curso: Recuperación de Información – 2002/1

-- ChaTo

Introducción

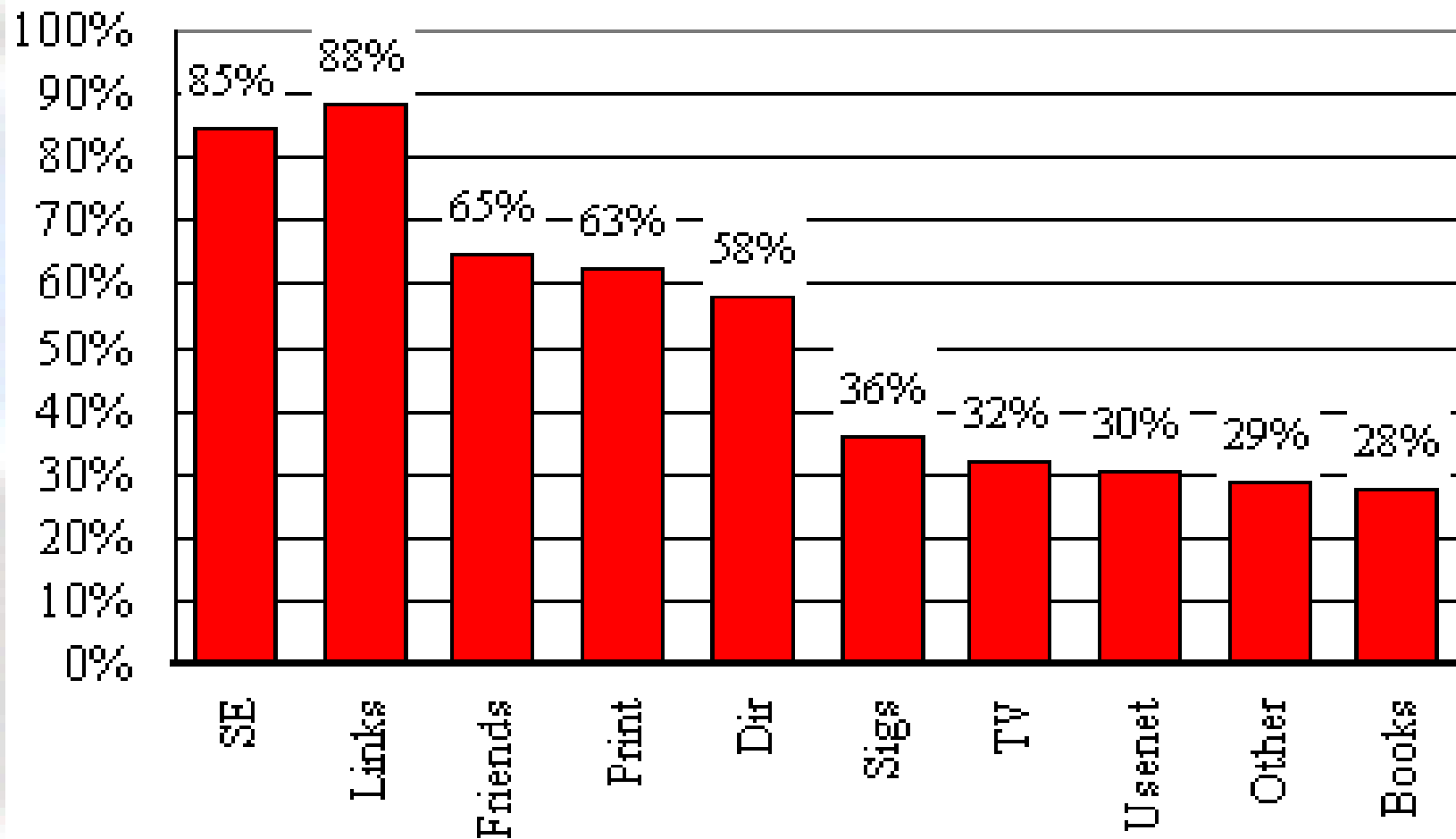
- Sistemas para localizar documentos
 - Buscadores basados en crawlers
 - Aprox. 1M páginas por hora
 - Directorios creados por humanos
 - DMOZ recibe 250 sitios nuevos por hora

Importancia Buscadores

- Webmaster: cómo atraer tráfico al sitio
- Generan 7 a 8% del tráfico
- 97% sitios corporativos (Fortune 100) tienen problemas estructurales para máquinas de búsqueda.
- Compras en línea
 - 25% buscar
 - 5% directorio
 - 2% banner

Importancia Usuario

How Do You Find New Web Pages/Sites?



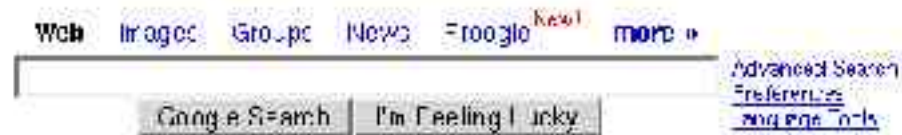
Importancia Usuario (2)

- 1/3 sesiones de usuario involucran buscador
- Escribir directamente la URL
 - 2001: 46%
 - 2002: 52%
- Buscar
 - 2001: 46%
 - 2002: 41%
- Buscador: descubrir nuevos sitios

Máquinas de búsqueda

- Google – google.com
 - Pagerank-based
- Yahoo – search.yahoo.com
 - Directory + Search engine
- AskJeeves – askjeeves.com
 - Natural language

Google™



- 1998, Backrub search engine.
- Múltiples buscadores especializados.
 - Catálogos, Imágenes, Noticias, Grupos, etc.
- Spellchecker, traducción, definiciones, etc.
- Adwords



- 1998, Primeros con lenguaje natural
 - Activo ahora en ajkids.com
- Máquina de búsqueda por Teoma

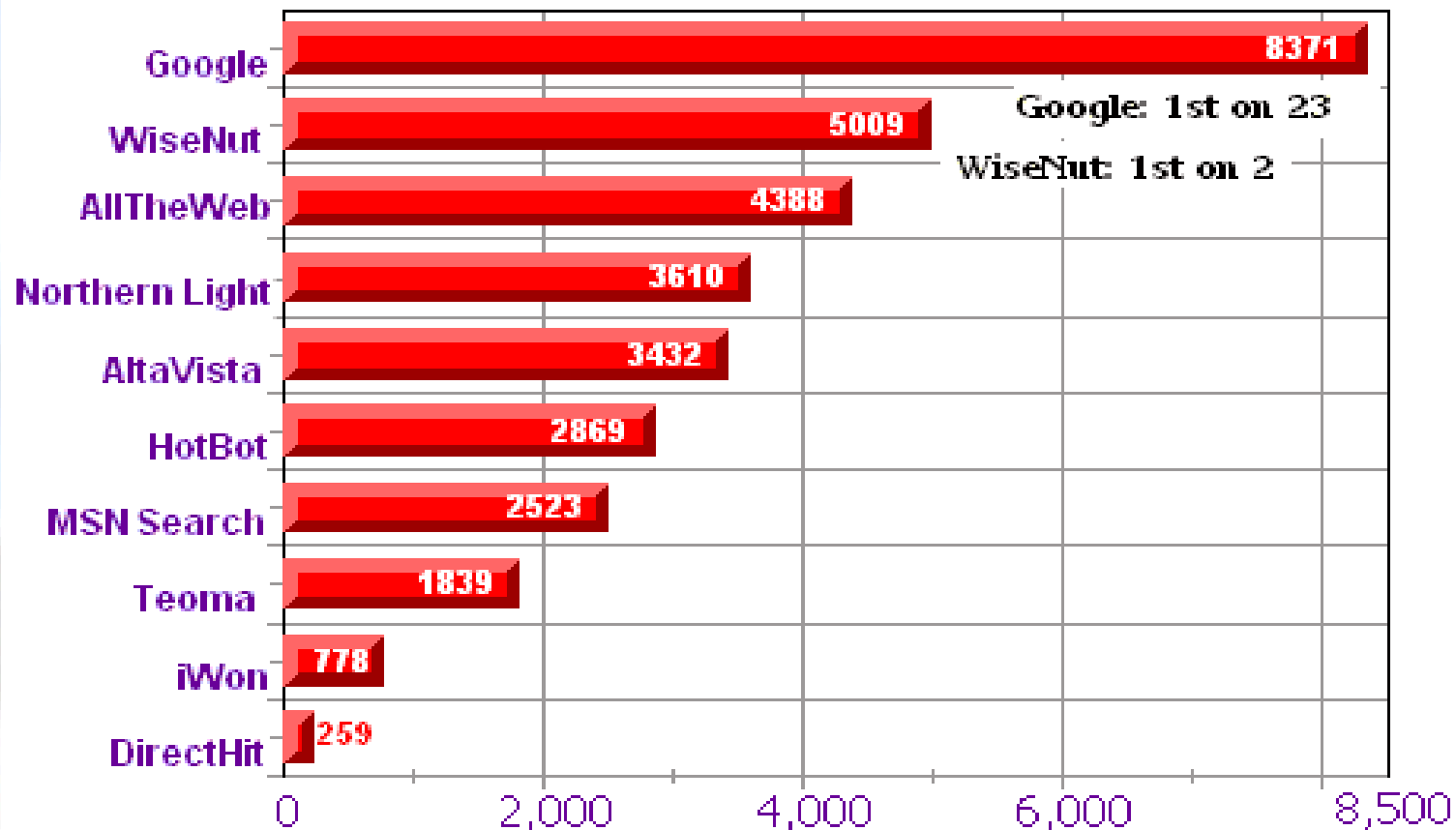
Otros buscadores

- AllTheWeb: incluye MP3 y FTP search
- Altavista: desde 1998.
- LookSmart: directorio.
- Lycos: uno de los primeros crawlers.

Comparación: Cobertura

Total Hits from 25 Searches March 4-6, 2002

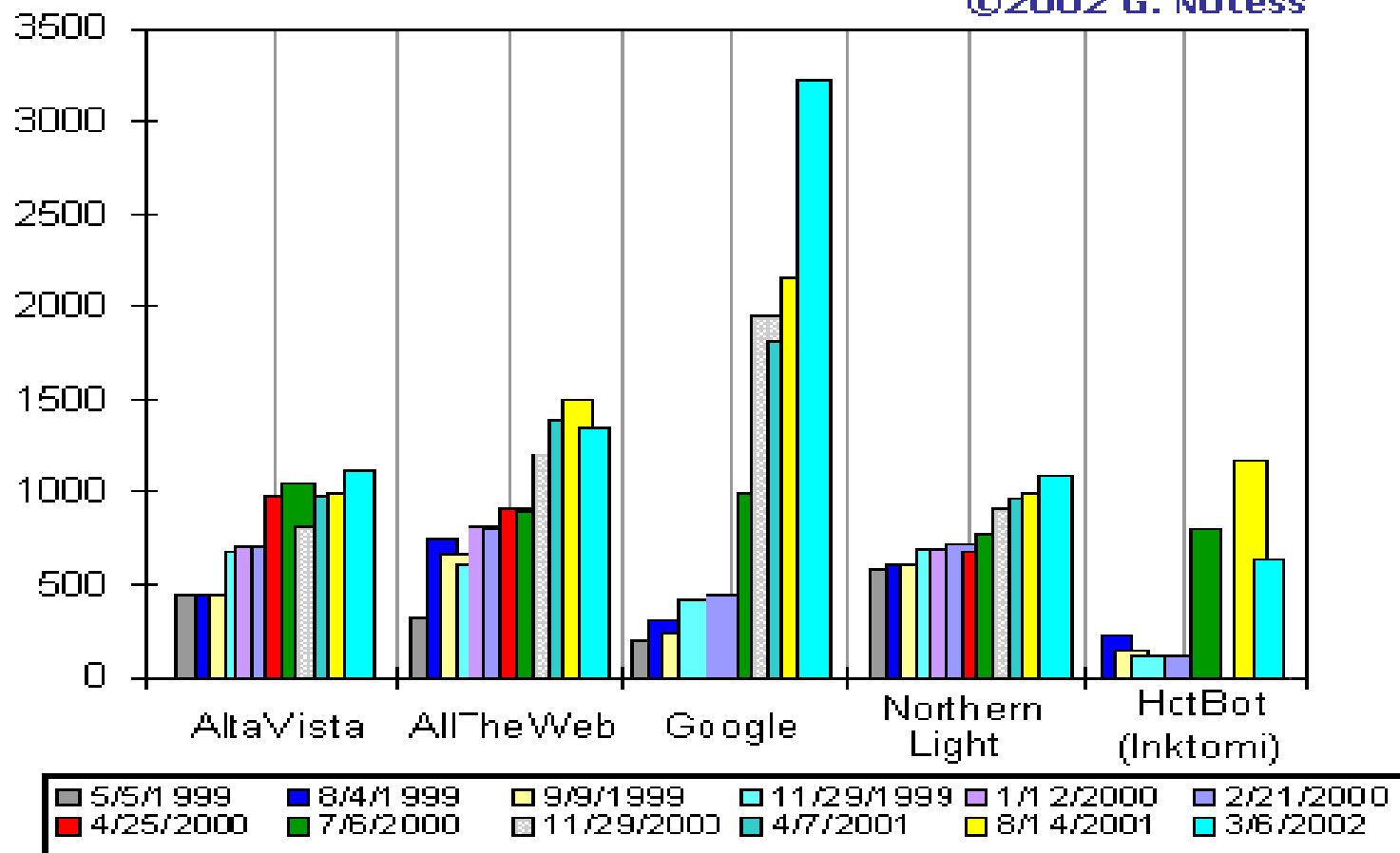
©2002 G. Notess



Comparación: Crecimiento

Top 5: Size Changes 5/99-3/2002 Results from same 8 searches

©2002 G. Notess



Comparación: errores

Search Engine	% Dead	% 400 errors only
AltaVista	13.7%	9.3%
Excite	8.7%	5.7%
Northern Light	5.7%	2.0%
Google!	4.3%	3.3%
HotBot	2.3%	2.0%
Fast	2.3%	1.8%
MSN Inklomi	1.7%	1.0%
Answers	1.3%	0.7%

Pero (Precisión @ 5)

- Dog (Google)
 - dogpile, dog.com, hotdog, explodingdog, dogplay
- Dog (Altavista)
 - dogofday, dog.com, yellowdog, dogracing, stuffdog
- Dog (Wisenut)
 - Ilovedogs, dog.com, hotdog, dogfriendly, yellowdog
- Dog (Yahoo)
 - dog.com, dogpile, dogplay, d.o.g., explodingdog

Google (antiguo)

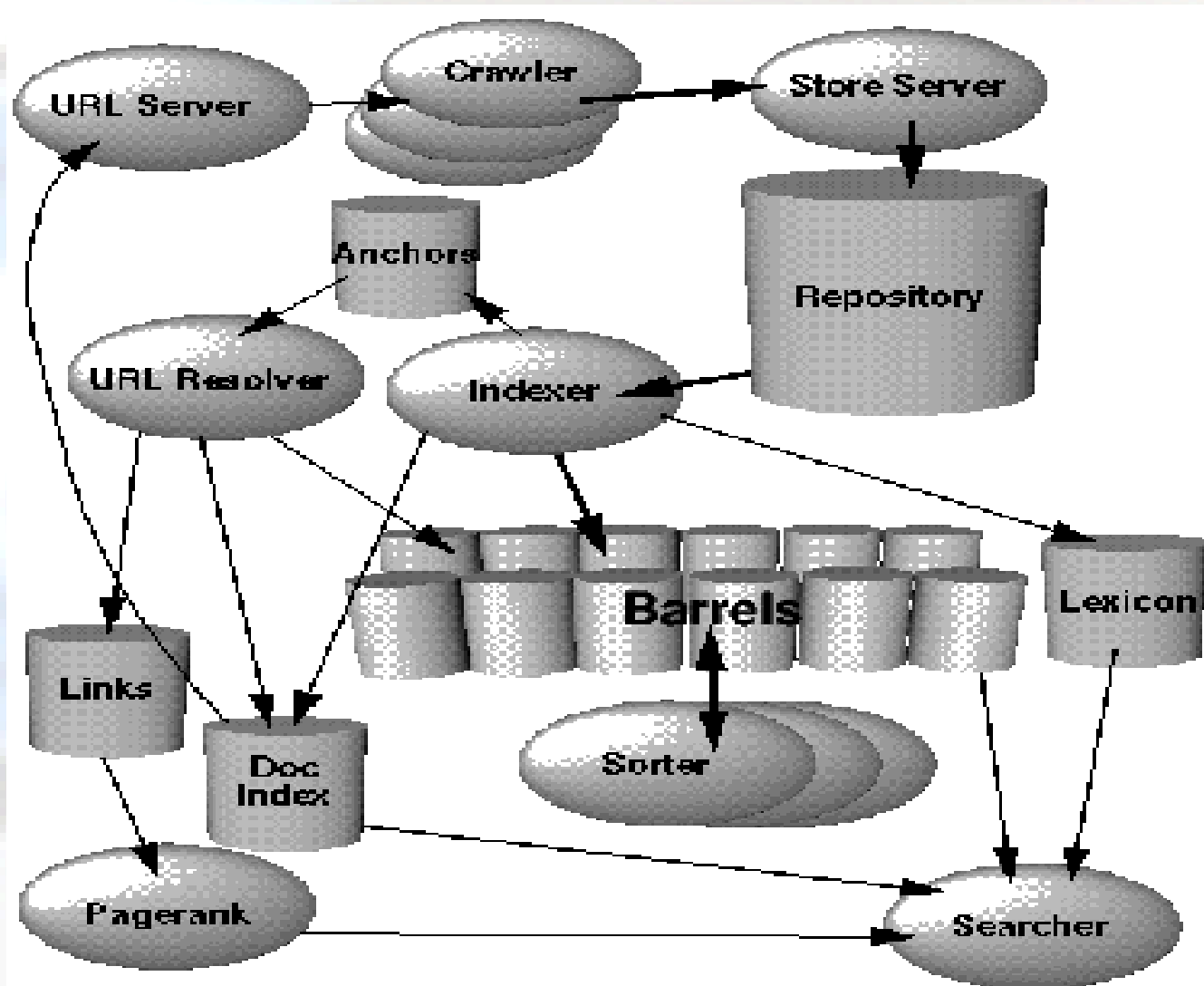
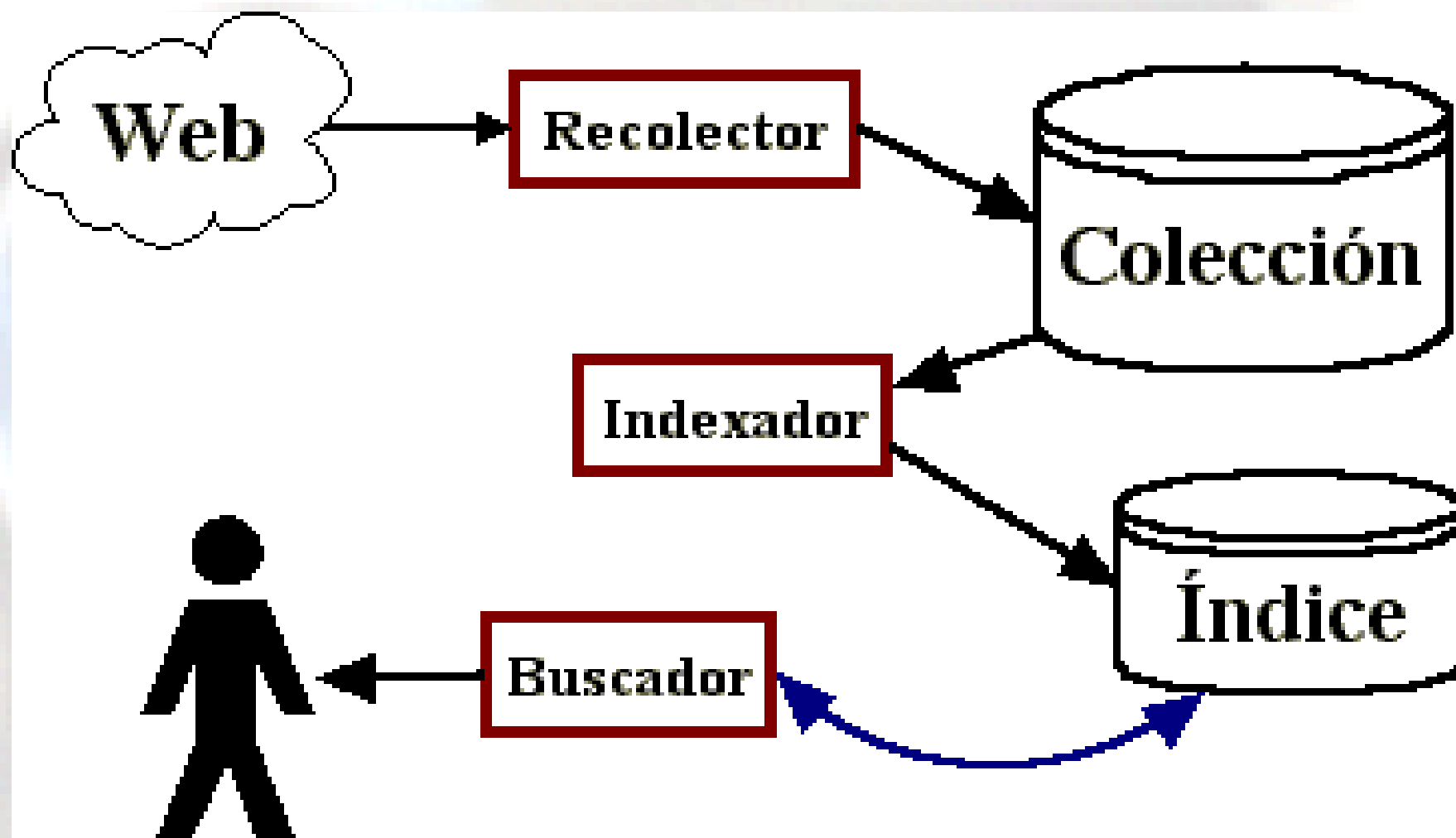
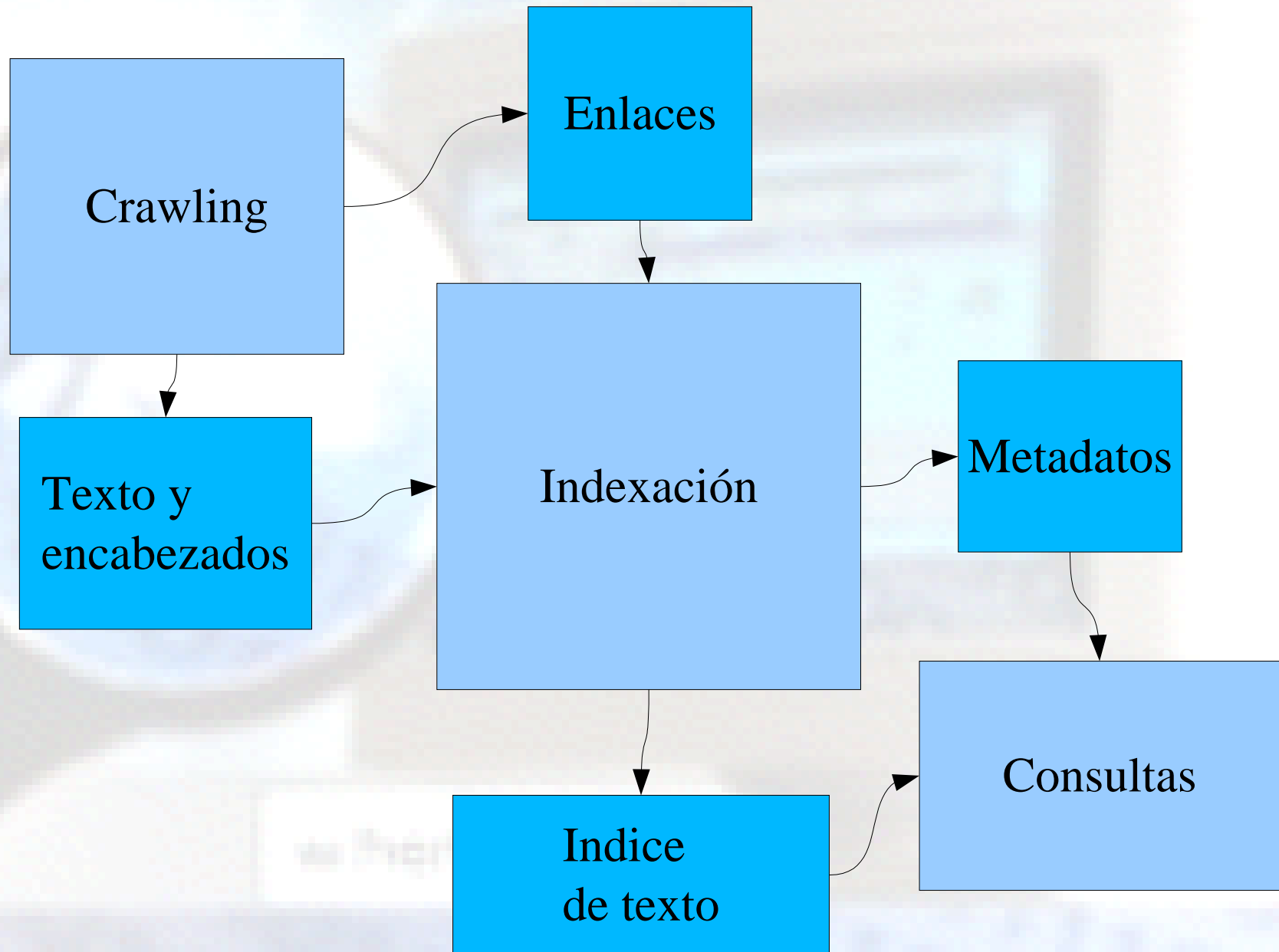


Diagrama Simple



Arquitectura general



Crawling

- Algoritmo para recorrer el grafo
- Métrica parada, Web infinita
- Parsing de los datos
 - HTML mal codificado (comillas, largo de los tags)
 - Binarios servidos como HTML
- Soporte frames
- Soporte de metatags, imagemaps
- Exclusión robots: robots.txt, meta robots
- Implementación de HTTP, etc.

Indexador

- Archivos de texto que no son HTML (PDF/PPT/etc.)
 - Toman tiempo de parsear
 - No se reducen siempre a términos (ej.: imágenes)
 - ¿Cuánto del texto indexar?
- ¿Almacenar o botar el texto completo?
 - Índice aprox. 1/3 del texto
- Indexar números
- Operaciones sobre texto: stemming
 - Eliminación de stopwords

Indexador (cont.)

- Eliminación de duplicados
 - Duplicados exactos: hashes
 - Duplicados cercanos: shingles
- Scrappers
 - Indexar definiciones
 - Indexar productos
 - Indexar imágenes/multimedios (sin bajarlos)
- Paralelización de índice
 - Por términos – Por documentos

Indexador (cont.)

- Re-indexación parcial
 - Crawling con 2 o más períodos distintos
 - Combinar rankings
- Charsets
- Idiomas en que el concepto de “palabra” no está tan claro

Buscador

- Operadores de búsqueda
 - Booleanos?
- Palabras “trigger”
 - Google: define, calculator
 - Yahoo: facts (encyclopedia), synonym, hotels/traffic
- Lenguaje natural


Buscador (triggers)

Yahoo! My Yahoo! Mail Welcome, **chato** [Sign Out, My Account]

YAHOO! search Yahoo! Search [Advanced Preferences](#)

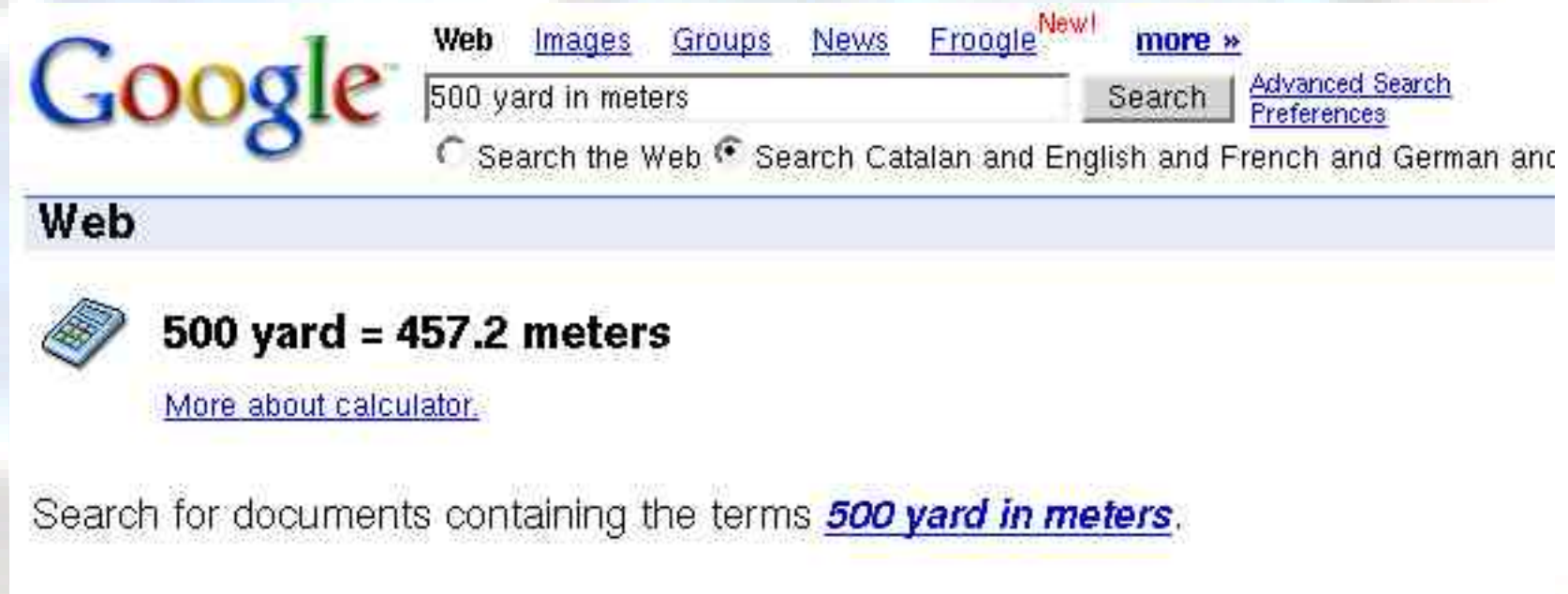
[Web](#) [Images](#) [Directory](#) [Yellow Pages](#) [News](#) [Products](#)

YAHOO! SHORTCUT ([What's this?](#))

 **brain:** the supervisory center of the nervous system in all vertebrates. It also serves as the site of emotions, memory and thought... [More](#)

View all results for **brain** in the: [Columbia Encyclopedia](#)

Buscador (especiales)



The screenshot shows the Google search interface. At the top left is the Google logo. To its right are navigation links: [Web](#), [Images](#), [Groups](#), [News](#), [Froogle](#) (with a 'New!' tag), and [more »](#). Below these is a search input field containing the text '500 yard in meters' and a 'Search' button. To the right of the search bar are links for [Advanced Search](#) and [Preferences](#). Below the search bar, there are radio buttons for 'Search the Web' (which is selected) and 'Search Catalan and English and French and German and...'. A horizontal bar labeled 'Web' is positioned below the search options. Underneath this bar, a calculator icon is shown next to the text '500 yard = 457.2 meters'. Below this result is a link for [More about calculator.](#). At the bottom of the search results, a line of text reads: 'Search for documents containing the terms 500 yard in meters'.

Buscador (lenguaje natural)



Ad Zone

[Privacy Policy](#) · [Home](#) · [About](#) · [Help](#) · [Parents](#) · [AJKids' Books](#)

You asked: 

Jeeves knows these answers:

 [Where do babies come from?](#)

 [Where can teens get advice online?](#)

 [Where can I learn about the human reproduction term](#)  [?](#)

Buscador (corregir)

- Sugerencias de corrección
 - Muchas palabras: cuáles borrar
 - Pocas palabras: cuáles agregar
- Errores ortográficos
 - “Galiello” Yahoo: Galileo Google: Galileo
 - “Galielxo” Yahoo: - Google: Galileo
 - “Aglilloe” Yahoo: - Google: Galileo (3 errores!)
 - “Xalielo” Yahoo: - Google: Galileo
 - “Cihlxe” Yahoo: - Google: Chile

Comparación: características

Search Engine	Boolean	Default	Proximity	Truncation	Case	Fields	Limits	Stop	Sorting
Google Review	-, OR	and	Phrase	No	No	title, url, more	Language, filetype, date	Yes + searches	Relevance, site
WebCrawler Review	< OR >	and	Phrase	No	No	No	Language	Yes + searches	Relevance, site
AllTheWeb Review	+,-, OR, WITH ()	and	Phrase	No	No	title, URL, link, more	Language, domain	No	Relevance, site
Lycos Review	+,-	and	Phrase	No	No	title, URL, link, more	Language, domain	No	Relevance
Northern Light Review	and, or, not, (), +,-	and	Phrase	Yes * %, auto plurals	No	title, URL, more	Doc type, date, more	No	Relevance, site, date, folders
AllCrawler Simple Review	+,-, AND, OR, AND NOT, ()	and usually	Phrase, NEAR	Yes " < 3 ** unlimited	No	title, URL, link, more	Language	No	Relevance, site
AllCrawler Adv Review	and, or and not, ()	phrase	Phrase, near, within, <, >	Yes " < 3 ** unlimited	Yes	title, URL, link, more	Language, date	No	Relevance, if used
HotBot Review	and, or, not, (), +,-	and	Phrase	Yes "	Yes	title, more	Language, date, more	Yes	Relevance, site
MSN Search Review	AND, OR, NOT, (), +,-	and	Phrase	No	Yes	title, link	Language, date, more	Yes	Relevance
Teoma Review	= only	and	Phrase	No	No	No	No	Yes + searches	Folders, Relevance, metafiles

Desafíos Generales

- Gran volumen de datos
 - Buscadores verticales
- Datos distribuidos
 - Problemas de red
- Datos volátiles
 - 404 Not Found
- Datos redundantes
 - 10% copias

Desafíos Generales (2)

- Datos no estructurados
 - No adhieren a estándares
- Datos de baja calidad
 - Información no confiable
- Datos heterogéneos
 - Formatos
 - Información que se desestructuro (ej.: BD->HTML)

Desafíos Específicos

- Crawler
 - DNS cuello botella
 - Velocidad variable de sitios
- Indexador
 - Conversión formatos distintos de texto
 - HTML no es respetado
- Buscador
 - Ranking: encontrar fácil, rankear difícil
 - Metabuscaadores

Spamming

- URL con sessionID (oculta)
 - Diferente URL misma página
 - Variaciones ligeramente distintas
- Spamming keywords
 - Texto pequeño o transparente
 - Páginas para crawlers y páginas para humanos
- Defensa: ranking enlaces sin contar links internos
 - Ataque: varios dominios mismo dueño

Spamming (práctica)

- Texto pequeño o invisible (color similar fondo)
- Metadatos que no reflejan el contenido de la página
- Páginas que tienen redirects automáticos vía Javascript
- “Link farms”, páginas que son esencialmente links a otras páginas
- Page-swapping (página para el crawler, página para los humanos)

Spamming (práctica)



merchandise is received, please call Bob Burt at 1-321-784-6129 for a full refund or replacement.

American Flag Cases
Bob Burt
731 Java RD
Cocoa Beach, FL 32931
Phone: 1-321-784-6129
usflag@qnc.net

3231

merchandise is received, please call Bob Burt at 1-321-784-6129 for a full refund or replacement.

American Flag Cases
Bob Burt
731 Java RD
Cocoa Beach, FL 32931
Phone: 1-321-784-6129
usflag@qnc.net

3231

flag cases sells wooden flag cases,
veteran burials require flag cases from American flag cases,
visit American flag cases and save on memorial flag cases
American Honor Flag Display Cases: Memorial Woodcrafts, American Honor Flag Display Case
Hand-crafted, Armed Forces
Custom display cases featuring memorial flag cases, sailor cases, medals cases, collectible cases,
antiques, solid oak construction, wholesale and retail.

Diseño de sitios

- Títulos y meta-tags
 - Principalmente títulos adecuados
- Links encontrables
 - Etiquetas de texto cuando sea posible
 - Evitar imágenes como links
 - No ocultar links tras javascript (ej.: falsos sessionid)
- Tener buen uptime del sitio
 - DNS
 - Servidor Web

Optimización de sitios

- Sitios densamente conectados
- Todos los enlaces visibles por máquinas de búsqueda
- Todas las palabras relevantes mencionadas en cada página, o en los meta-tags
- Tener enlaces desde buenos sitios
- Search engine optimization ? Depende de cuánto prometan...

 Quality Traffic From The Industry Leaders	 The Best SEO Software: Free Download!
 Search Engine Optimization Tools and Services	 Search Engine Marketing Services Free Report