PROJECT REPORT (Final Year Project 2007-2008)

Hybrid Search Engine

Project Supervisor Mrs. Shikha Mehta

INTRODUCTION

Search Engines

Definition:

A search engine is an information retrieval system designed to help find information stored on a computer system, such as on the World Wide Web, inside a corporate or proprietary network, or in a personal computer.

Examples:

Various search engines are available on the internet e.g. Google, Alta Vista, Ask.com, Yahoo, Lycos, Alltheweb, Myspace, etc.

The popularity of search engines can be estimated by the fact that approximately 112 * 10⁶ searches are made in a single day from one search engine alone.

How do Search Engines work?

There are differences in the ways various search engines work, but they all perform three basic tasks:

- They search the Internet -- or select pieces of the web -- based on important words. [CRAWLER]
- They keep an index of the words they find, and where they find them. [INDEXER]
- They allow users to look for words or combinations of words found in that index. [SEARCHER]



PROBLEM STATEMENT

On the basis of recent studies made on the structure and dynamics of the web itself, it has been analyzed that the web is still growing at a high pace, and the dynamics of the web is shifting. More and more dynamic and real-time information is made available on the web.

Our aim is to design a search engine that meets the challenges of web growth and update dynamics.



Our Proposed Design



Data Flow



Snapshots



INPUT Initial set of URL's taken for the sample :

Url	Depth
http://www.google.com/dirhp/	1
http://www.sites-internationaux.com/	1
http://www.diroot.com/	1
http://www.urlz.net/	1
http://www.dmoz.org/	1
http://www.hotvsnot.com/	1
http://www.webworldindex.com/	1

OUTPUT

As the crawler visits these URLs, it identifies all the hyperlinks in the page and adds them to the list of URLs to visit, called the **crawl frontier**. URLs from the frontier are recursively visited according to a set of policies.

Url	Depth
http://www.webworldindex.com/phtml/USA_Real_Estate/Alabama_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/California_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/Colorado_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/Florida_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/Indiana_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/Maryland_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/North_Carolina_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/South_Carolina_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/Texas_Real_Estate/	2
http://www.webworldindex.com/phtml/USA_Real_Estate/Virginia_Real_Estate/	2
http://www.webworldindex.com/phtml/Regional/	2
http://www.webworldindex.com/phtml/Regional/Countries/	2

HTML Parser

Programming Language

C#

Input

After the crawler crawls the web and store the pages in the repository, we need to extract the useful information from the web page like title, no. of forward links etc. of all the web pages.

Output

The extracted information is then stored in the database.

Output

Url	Nof	Bit	Uri	Title	Buffer
http://www.webworldindex.com/phtml/People/Families/	36	a04da0ae4c8a302390bdaabcb7e5cfab	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test1.txt	Families Directory of Families Websites	15872
http://www.webworldindex.com/phtml/Recreation_and_Sports/	55	53692463eb8a6f4a9147d38a20acebbc	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test2.txt	Recreation and Sports Directory of Recreation and Sports Websites	24269
http://www.webworldindex.com/phtml/Recreation_and_Sports/Gambling/	53	afd860964c743e2f344297f188542f56	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test3.txt	Gambling > Casino Directory > Web World Directory	21917
http://www.webworldindex.com/phtml/Recreation_and_Sports/Sports/	148	32a6a1afdcd331cd1e1e763f8600a3b6	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test4.txt	Sports Directory of Sports Websites	41075
http://www.webworldindex.com/phtml/Education/Reference/Libraries/	30	9cbb66e4a992f6bd7e130d1fe36eac99	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test5.txt	Libraries Directory of Libraries Websites	16487
http://www.webworldindex.com/phtml/Education/Reference/Thesauri/	20	41c2c65f8ff1e8aa557f535ba07473ca	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test6.txt	Thesauri Directory of Thesauri Websites	11742
http://www.webworldindex.com/phtml/Science/Agriculture/	61	9fdf3183a263cd12e8f729820fd297d4	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test7.txt	Agriculture Directory of Agriculture Websites	21747
http://www.webworldindex.com/phtml/Science/Biology/	68	3bfab2877017644a6b1ac87b22d25d4e	F:\STUDY MATERIAL\Final Year Project\SEARCH ENGINE\Database\test8.txt	Biology Directory of Biology Websites	23350

Compressor-Decompressor

Programming Language

C#

Input

After the crawler crawls the web, this module compress the pages and store them in the repository. We need to decompress all the web pages to search a keyword.

Output

The compressed pages are stored in the database.



Content Seen Tester

Programming Language

C#

Input

The content seen tester generates a bit sequence of all the web pages using MD5 algorithm.

Output

The bit sequence of every web page is stored in the database.



Indexer

Sorts the results found on the basis of a rank distribution algorithm.

Programming Language C#

Input

The links between all the web pages are fetched from the database.

Output

The rank of each web page is stored in the database.



🖼 file:///C:/Documents and Settings/Administrator/Desktop/EVALUATED_MODULES/Indexer/Ind 🗕 🗖 🗙
story/348269.html
ttn://www.cricinfo.com//feedback/ httn://www.cricinfo.com//
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//db/
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//db/RSS/
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//linktous/
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//feedback/
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//feedback/
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//ci/content/page/1560
56.html
ttp://www.cricinfo.com//feedback/ http://www.cricinfo.com//db/jobs/
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//link to database/MAN
GEMENT/PRIVACY POLICY.html
http://www.cricinfo.com//feedback/ http://www.cricinfo.com//link to database/MAN
GEMENT/TERMS_USE.html
http://www.cricinfo.com//feedback/ http://www.espn.com
http://www.cricinfo.com//feedback/ http://esun.go.com/
http://www.cricinfo.com//feedback/ http://esunsoccernet.com
http://www.cricinfo.com//feedback/ http://www.scrum.com
processing
rl>http://www.imdb.com//title/tt0499448/ rank>0.290530797439244
rl>http://www.imdb.com//title/tt0499448/taglines_rank>0.150791510185331
rl>http://www.imdh.com//title/tt0499448/plotsummary_rank>0.150791510185331
trl>http://www.imdb.com/sunopsis rank>1.24883192717323
rl>http://www.imdb.com//keuword/sword-and-sorcery/ rank>0.150791510185331
······································

Refresher

Updates the local database with fresh copies of web pages.

Programming Language

Input

C#

The cached pages from the database.

Output

The refreshed pages are stored in the repository.





User Interface

Programming Language ASP .NET

Input

The user enters a keyword or multiple keywords.

Output The results are fetched to the user.



