

Searching for the Right Stuff [FIT chapter 5]

And boy, is there a lot of stuff to search through!

The best way to narrow down the amount of stuff you have to look through is by searching in the right area to begin with

Finding and Evaluating Information

- ❖ **Key principle:**
 - The best place to search for something is where it is likely to be found.
 - Question: Where is the best place to find tax information?
+ Used car information?
- ❖ What are the potential places that one would find the best information?
 - Library – not everything is on the Internet, and it never will be! – and the library specializes in organizing information.
 - Specialized databases – if you do search for information electronically, go to the places where the content has already been evaluated and authenticated. Often these electronic databases can be used, for free, atthe Library!
- ❖ So why all this focus on the library?

What do Humans do?

- ❖ We associate things with other things and we organize them (group them together)
- ❖ Humans have a natural tendency to organize (cluster) similar things together
 - When you were a child, did you put all of your Barbie dolls in one place, or all of your robot toys? How about the same color toys?
- ❖ Think about the subjects you study in school: Math, Social Sciences, Geography, Art, etc.
 - Each of these subjects is further divided, becoming more and more detailed, specialized
- ❖ Librarians tackle the problem of trying to organize and place information where people can most easily find it³

Organizing Electronic Things

- ❖ If we tend to group similar things together, then we also tend to look for things that are similar in the same places
 - Problem: Do any 2 people organize things in exactly the same way? Do any 2 groups of people?
- ❖ Now, go to the web where all of that stuff sits
 - NO organization
- ❖ But, individual web sites can try to organize their "stuff". A well-designed and well-organized site will help you find the "stuff" faster



A Well Organized Site

- ❖ A web site that wants to help users find relevant information fast is usually organized in a manner similar to a Library:
 - ❑ The site is divided into broad categories (subjects)
 - ❑ Each category in turn has sub categories
 - Familiar? Didn't we try to organize folders on Dante in a meaningful manner?
- ❖ If a user that visits a site (you) can quickly see how things are categorized then they will find what they need much quicker
 - ❑ This assumes they have started looking in the right spot...

1/27/2002

L05-5
© Copyright 2000-2002, University of Washington

1/27/2002

L05-7



Searching NPR

1/27/2002

L05-7
© Copyright 2000-2002, University of Washington

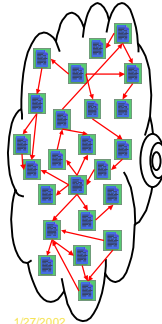
1/27/2002

L05-8
© Copyright 2000-2002, University of Washington

But what if you JUST DON'T KNOW?

Then the next step may be to go to a search engine

Structure of the Web

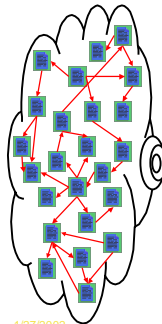


- ❖ The Web is just lots of documents on lots of machines
- ❖ Each document has a link to it from some other document
- ❖ If Page A links to Page B, you can get to document B from A (but not vice versa!)

What is a Search Engine?

- ❖ A collection of programs designed to assist users in finding information
- ❖ Consists of three things:
 - A crawler (aka spider, robot)
 - A query processor
 - An interface (you will see the GUI's in Assignment 1)
- ❖ A crawler does what the name implies:
 - "Crawls the Internet" building an index of URLs and key terms that are hopefully an indicator of the content of the page.
- ❖ A query processor takes a request from the user (search terms)
 - Retrieves the list of URLs associated with a given set of key word terms according to the index

A Web Spider: See how they crawl!



- ❖ Start at a document
- ❖ Repeat:
 - Store document
 - Extract all links
 - Index Terms
 - Follow every link
- ❖ Until: all links followed
- ❖ Redo when necessary
- ❖ Need to pick "top" page, and good pages

FIT 100

The Search, literally, never ends

- ❖ Crawling is an ongoing process. Why?
 - ❑ To keep up with the millions of new pages that are added to the Web on a weekly basis
 - ❑ To go back and revisit sites and make sure that links aren't broken
- ❖ Lists of URL's are created on the spot for users
 - ❑ But the query processor doesn't have a robot search the web each time
 - ❑ Goes to the database of searched pages and matches search terms with index terms
- ❖ The effectiveness of a search engine will depend in part on how much of the Internet the crawlers have seen and how discriminating the index is.

1/27/2002

L05-13
© Copyright 2000-2002, University of Washington

FIT 100

A bit about Google

Google uses an extremely simple idea to create a very sophisticated search engine

The "basic" search, seen here, is often effective just based on their rankings

Advanced search is usually better if you spend a moment thinking about it



1/27/2002

L05-14
© Copyright 2000-2002, University of Washington

FIT 100

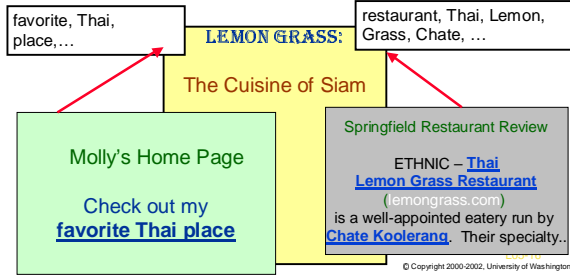
Advanced Search



FIT 100

Google's Sophistication

- ❖ Google is a search engine that indexes a page by the words used in the anchor tags of pages that link to it



© Copyright 2000-2002, University of Washington



More Google Sophistication

- ❖ Popularity is also a key to Google's rankings
- ❖ If page A links to page B, then that is considered a vote by page A for page B
- ❖ If page A is also a very popular site that many other sites link to, then page A's vote is worth more
- ❖ How can WE get a sense of how popular a web site is? Do a simple link search in Google to see how many sites link to the site that interests us.

1/27/2002

L05-17
© Copyright 2000-2002, University of Washington



Summary

- ❖ Search where you are most likely to find the information
- ❖ Good sites will have effective navigation that you can easily figure out
- ❖ Local searches in good sites can quickly find candidate pages
- ❖ Search Engines build indexes to assist in searching the web
- ❖ When doing a search of two or more words or phrases, specify whether
 - ❑ Both words MUST be present: AND, +
 - ❑ At least one of the words must be present: OR
 - ❑ The word(s) must NOT be present: NOT, -

1/27/2002

L05-18
© Copyright 2000-2002, University of Washington