

## Searching for the Right Stuff

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, at ...the 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 it

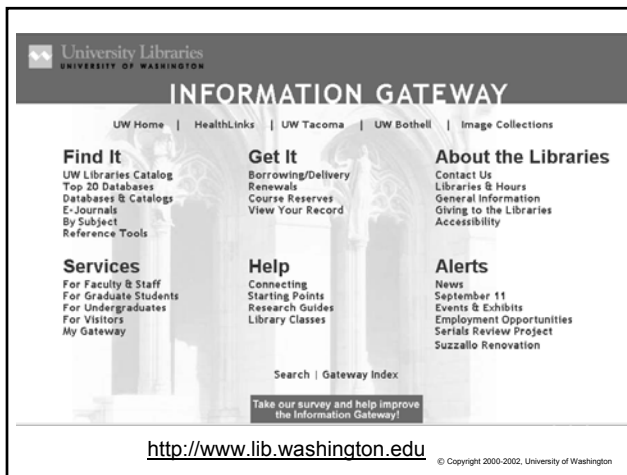
## Organizing Electronic Things

- ❖ 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

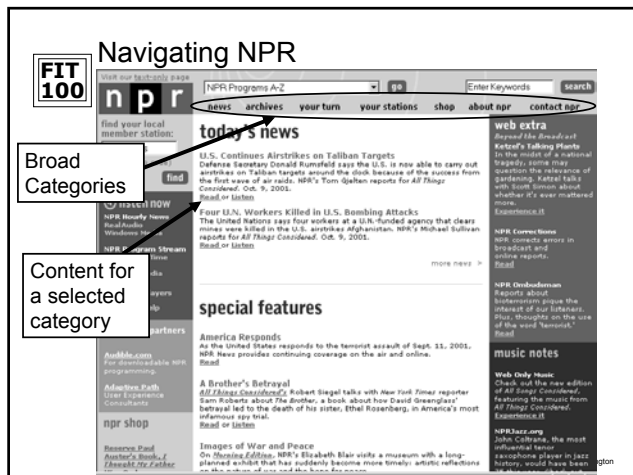
**FIT 100** A Well Organized Site

- ❖ 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...

© Copyright 2000-2002, University of Washington



© Copyright 2000-2002, University of Washington



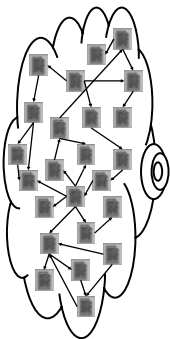
**FIT 100**

But what if you JUST DON'T KNOW?

Then the next step may be to go to do a search of information on the Web

© Copyright 2000-2002, University of Washington

## **FIT 100** 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 page B from A (but not vice versa!)
- ❖ However, searching for information by going from page to page, one document at a time, would take the rest of eternity

© Copyright 2000-2002, University of Washington

## **FIT 100** Directories: Yahoo!

- ❖ Not a search engine!
- ❖ A directory
  - ❑ Web pages have been organized into a hierarchical structure based on broad categories
  - ❑ Humans do the organizing, usually

<b>Arts &amp; Humanities</b> <a href="#">Literature, Photography...</a>	<b>News &amp; Media</b> <a href="#">Full Coverage, Newspapers, TV...</a>
<b>Business &amp; Economy</b> <a href="#">E2E, Finance, Shopping, Jobs...</a>	<b>Recreation &amp; Sports</b> <a href="#">Sports, Travel, Autos, Outdoors...</a>
<b>Computers &amp; Internet</b> <a href="#">Internet, WWW, Software, Games...</a>	<b>Reference</b> <a href="#">Libraries, Dictionaries, Quotations...</a>
<b>Education</b> <a href="#">College and University, K-12...</a>	<b>Regional</b> <a href="#">Countries, Regions, US States...</a>
<b>Entertainment</b> <a href="#">Picks, Movies, Humor, Music...</a>	<b>Science</b> <a href="#">Animals, Astronomy, Engineering...</a>
<b>Government</b> <a href="#">Elections, Military, Law, Taxes...</a>	<b>Social Science</b> <a href="#">Archaeology, Economics, Languages...</a>
<b>Health</b> <a href="#">Medicine, Diseases, Drugs, Fitness...</a>	<b>Society &amp; Culture</b> <a href="#">People, Environment, Religion...</a>

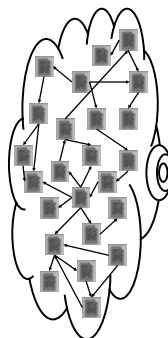
© Copyright 2000-2002, University of Washington

## **FIT 100** What is a Search Engine?

- ❖ A collection of programs designed to assist users in finding information
- ❖ Consists of four things (the book gives the 2 main ones):
  - ❑ A crawler (aka spider, robot)
  - ❑ A query processor
  - ❑ A user interface
  - ❑ A database
- ❖ 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

© Copyright 2000-2002, University of Washington

## **FIT 100** 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
  - ❑ Pages with multiple links

© Copyright 2000-2002, University of Washington

## **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 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.

© Copyright 2000-2002, University of Washington

## **FIT 100** Search Technique

- ❖ Search engines use basic logic to determine how to answer your question
- ❖ Boolean logic, in the form of Boolean operators, are the foundation of search logic:
  - ❑ AND
  - ❑ OR
  - ❑ NOT
- ❖ Many search engines now use "search math" instead of boolean terms
  - ❑ +, -
- ❖ Search for exact phrases, like titles, with quotation marks
  - ❑ "Fellowship of the Ring"

© Copyright 2000-2002, University of Washington

## **FIT 100** Search Technique

- ❖ Successful search takes more than just Boolean logic
  - ❑ There is too much out there
- ❖ Other search mechanisms
  - ❑ Ranking of search terms
  - ❑ Weighted terms
  - ❑ Proximity
- ❖ Some search engines have incorporated other ways to return better results

© Copyright 2000-2002, University of Washington

## **FIT 100** A bit about Google

The image shows a screenshot of the Google search page from around 2000. At the top, the word "Google" is displayed in its signature font. Below it, the text "Search 1,345,966,000 web pages" is visible. There are two main buttons: "Google Search" and "I'm Feeling Lucky". To the right of these buttons is a link for "Advanced Search" with a small arrow pointing to it. Below the search area, there is a link for "Google Web Directory" with the tagline "the web organized by topic". At the bottom of the page, there are several links: "Cool Jobs", "Add Google to Your Site", "Advertise with Us", "Google Toolbar", and "All About Google".

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

© Copyright 2000-2002, University of Washington

**FIT 100** Advanced Search

The screenshot shows the Google Advanced Search page. Callouts point to the search options:
 

- AND constraints:** points to the 'with all of the words' option.
- OR constraints:** points to the 'with any of the words' option.
- NOT constraints:** points to the 'without the words' option.

 Other visible options include 'with the exact phrase', 'Language', 'Occurrences', 'Domains', and 'SafeSearch'. A note at the bottom states: 'Google uses an extremely simple idea to create a very sophisticated search engine'.

**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

The diagram illustrates how Google indexes anchor text. It shows three pages:
 

- Molly's Home Page:** contains the anchor text 'Check out my favorite Thai place'.
- The Cuisine of Siam:** contains the anchor text 'LEMON GRASS:'.
- Springfield Restaurant Review:** contains the anchor text 'ETHNIC - Thai Lemon Grass Restaurant (lemongrass.com) is a well-appointed eatery run by Chate Koolerang. Their speciality..'.

 Arrows from the anchor text boxes point to a central box containing the indexed words: 'favorite, Thai, place,...' and 'restaurant, Thai, Lemon, Grass, Chate, ...'.

© Copyright 2000-2002, University of Washington

**FIT 100** 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.

© Copyright 2000-2002, University of Washington

**FIT 100** 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, -

© Copyright 2000-2002, University of Washington