open.michigan

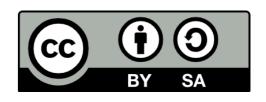
Unless otherwise noted, the content of this course material is licensed under a Creative Commons BY-SA 3.0 License. http://creativecommons.org/licenses/by-sa/3.0/

Copyright © 2009, Robert Frost.

You assume all responsibility for use and potential liability associated with any use of the material. Material contains copyrighted content, used in accordance with U.S. law. Copyright holders of content included in this material should contact open.michigan@umich.edu with any questions, corrections, or clarifications regarding the use of content. The Regents of the University of Michigan do not license the use of third party content posted to this site unless such a license is specifically granted in connection with particular content. Users of content are responsible for their compliance with applicable law. Mention of specific products in this material solely represents the opinion of the speaker and does not represent an endorsement by the University of Michigan. For more information about how to cite these materials visit http://michigan.educommons.net/about/terms-of-use.

Any medical information in this material is intended to inform and educate and is not a tool for self-diagnosis or a replacement for medical evaluation, advice, diagnosis or treatment by a healthcare professional. You should speak to your physician or make an appointment to be seen if you have questions or concerns about this information or your medical condition. Viewer discretion is advised: Material may contain medical images that may be disturbing to some viewers.





Finding Information

 All about how we now search for and (hopefully) find information and determine (or not!) its veracity. Theories of classif cation and the revolt against old hierarchical systems. In passing, issues of how information finding systems morph into knowledge navigation and knowledge architecture frameworks. Metadata as the new magic bullet. Tagging as a solution, but a laborious one. The advantages and shortcomings of Web search engines. What we [think] we know about how people try to f nd information.

Ordering and Finding Information

- Traditions of cataloguing and classif cation
 - Collections, relics, and cabinets of curiousities
 - Diderot and the Tree of Knowledge; Linneas
 - Hierarchies of knowledge
 - Social & cultural hierarchies
 - Melvil Dewey and his system: a single subject... +
 - How many places right of the dot?
 - Problems of scalability
 - LC headings: Dewey made scalable
 - Dilemmas of LC headings: reproducing the cataloguer's thinking
 - Now: f at structures, relations, parsing, semantics
 - Records and the rule of provenance

The Key Problem in "Information Retrieval": Increasing the Signal-to-Noise Ratio

- A consequence of "information overload"
 - Most searches yield oceans of spurious "hits"
- Strategies to attain greater accuracy, better results:
 - Rethink the use of hierarchies
 - Be aware of the Bradford distribution (the 80/20 rule)
 - Shift toward faceted systems, relational aspects
 - Structure with rich metadata structures/f elds, then "populate it"
 - Problems in tagging information objects
 - —but who will do the tagging? Our "objects dB" project...
 - is automated tagging a credible option?

P roblems of Hierarchies

- Taxonomies, nomenclatures, and controlled vocabularies—who def nes these??
 - Buffon, Linneas, and the rage for order: phenetics vs. cladistics
 - Standardizing language: nomenclatures, thesauri, etc.
 - Controlled vocabularies and authority lists
 - LC headings
 - Authority lists as def ning proper nouns
 - Controlled vocabularies are broader
- Nota bene: Taxonomies do not yield new knowledge—they can become conceptual prisons as they order things
- New issues: loadable lists & DTDs via XML + friends

From Hierarchies to Flatness

- The birth of relational databases
- Hierarchic and relational databases compared
 - Parallels to org charts and peer-to-peer
 - Consider EndNote: can we conf dently share our f les?
- Computation power and the ease of the fat search
 - Illusory "hits"
 - Second-guessing non-standard vocabularies and keywords

Booleans for simple fliering Searches

- Standardizing records
 - Standardized vocabularies
 - OPACs and MARC format; Z39.50
 - EAD and quasi-searchable archival f nding aids; CIDOC-CRM
- Data descriptors def ned: form vs. content
 - Page languages, SGML, Postscript
 - Metadata
 - From HTML to XML; semantics and auto-tagging
 - Endless need to "metatag" digital objects
 - ls there a politics of tagging? (we know there's a politics of

Slide No. 7

The Compromise:

Flat/Relational Systems vs. Hierarchies

- Hierarchies yield better-focused results, and less processing, but require knowledge of the classif cation system
- Primitive f at systems vastly ease access—no need to "second guess" how hierarchy is ordered, but often yield too many [useless] results
- Faceted searching (Bates): use tagging (read: XML) and metadata to define several "faces" of information objects, then access them with f at searching

Searching for Meanings and Ideas, Not Just Words

- First step: "do you mean...?" —Google
- Web caching: meanings def ned by frequencies & contexts of queries, and by analysis of hyperlinks
- Semantics
 - Def ned
 - From word-search algorithms to deriving meanings
 - The role of statistics

How Well Do "Finding Aids" Filter Information?

- First, how Google works...
 - Step 1: "spidering," to harvest what's publicly available on the Web—a constant process (and sites can implement "do not spider" rules)
 - Step 2: indexing with word searches, then building rankings based in large part on frequency of links to specif c pages
 - Consequences: a. no validity tests, so "mind-share" equates to "truth;" b. massive parts of the Web don't appear in Google; c. system is more a gathering and ranking than a f Iter
- Contrast this to for-pay and specialized databases (OCLC, RLIN-Eureka, Ovid, ProQuest, Dialog, MedLine...)
 - Minimal ranking, as the ordering algorithm is the traditional author/title/subject/[keyword] framework
 - Better general veracity, as content has already been vetted
 - At the same time, access rules and interfaces stink

More on Databases and Search Engines...

- Keep in mind that Google and Yahoo do also maintain canonical (or near-canonical) subject heading indexes...
- They also work closely with the older, more "mainstream" info aggregators such as OCLC and LC's National Union Catalog effort
- Many database services (ProQuest, MedLine, Nexis/Lexis, Gale, etc.) are actually "data aggregators;" they thus inherit the foibles and oddnesses of the originals
- Both Web search engines and higher-end databases redefine "browsing" in ways that are probably not as fruitful as old-fashioned "stack browsing"
- Note well the emergence of "grass-roots" indexing with Wikipedia, del.icio.us: do folksonomies consolidate or diffuse f ndability frameworks? Do they undermine the utility of thesauri and controlled vocabularies?

An Important Insight... with startling implications.

- In terms of the development of personally-held knowledge (what and how people learn), research has often indicated that discovery by unintentional means, from browsing to thinking by metaphor, is often quite effective and profound. Consider the implications:
 - Search-engine results may often be too narrow, precluding serendipity
 - The best search-behavior model is recursive: people learn and adapt in an extended, often less-directed fashion
 - "Stack browsing" is often more useful than we would otherwise believe
 - Maybe the best interaction with a reference librarian is "Socratic"...
 - [This is what "Information Searching Behavior" is all about]—current preferred model is "berry-picking"