University Library Scholarship                                    University Library

5-20-2009

# "Beyond Google: searching the deep web."

Louise Buckley
*University of New Hampshire*, lousie.buckley@unh.edu

# Beyond Google:
# Searching the Deep Web

Louise Buckley

Reference Librarian

UNH Dimond Library

lbuckley@unh.edu or 603-862-1435

---

# Have You …

- Made a plane or hotel reservation online
- Gotten directions from MapQuest or Google Maps
- Searched for articles in an EBSCOhost, Gale or ProQuest database
- Searched an online catalog

If yes,
then you have searched the Deep Web

Visible Web = Surface Web = Shallow Web

Findable by general purpose search engines

Deep Web = Invisible Web = Hidden Web = Dark Web = Cloaked Web

Not findable by general purpose search engines

Timeline of Events Related to the Deep Web

http://papergirls.wordpress.com/2008/10/07/timeline-deep-web/

May 20, 2009                          NHLA Spring Conference                          3

---

# Deep Web larger than Surface Web



Image Source: http://www.mkbergman.com/?p=458

2007 Estimates:
http://netforbeginners.about.com/cs/secondaryweb1/a/secondaryweb.htm

May 20, 2009                          NHLA Spring Conference                          4

# How Visible Web is Created

# Deep Web Content Categories

- Databases
  - Academic Search, NewsBank, InfoTrac, ERIC, NASA Image eXchange (NIX), etc.
  - Some proprietary, some free
- Sites requiring forms to be completed
  - Generate dynamic information, often single use
    - Traveling direction sites, job hunting sites, phone books

# Deep Web Content Categories

- Deep (extensive) web sites
  - may be only partially indexed
    - Largest Deep Web™ Sites
    - http://aip.completeplanet.com/aip-engines/help/largest_engines.jsp

- Formats – new or older
  - http://www.google.com/advanced_search?hl=en
  - http://images.google.com/advanced_image_search?hl=en
  - http://video.google.com/videoadvancedsearch?q=butterflies%20fi letype%3Ajpg&hl=en

May 20, 2009        NHLA Spring Conference        7

# Deep Web Content Categories

- Very current or real-time content
- Sites requiring registration or password access
- Sites excluded by their owners or that are not linked to indexed pages
- Sites the general search engine chooses to exclude

Another image of Web content:
http://netforbeginners.about.com/library/diagrams/n4layers.htm

May 20, 2009        NHLA Spring Conference        8

# Why Use the Deep Web

- Quantity of information

- Quality of information

- Greater depth of information in specific topics

- Often structured information

# Searching

- Searching Surface Web More Effectively

- Article Databases

- Split-level Searching

- Subject Directories

- General Deep Web Search Engines

- Vertical Searching

- Digital Collections

# More Effective Surface Web Searching

**ranking.thumbshots.com**

Try more than one general purpose search engine

- Learn strengths/weaknesses/coverage

- "Different Engines, Different Results" http://www.infospaceinc.com/online prod/Overlap-DifferentEnginesDifferentResults.pdf

[Article on Thumbshots ranking tool]

**THUMBSH⊙TS**

Thumbshots Ranking

May 20, 2009                    NHLA Spring Conference                    11

---

# More Effective Surface Web Searching

- Use advanced search features
    - Domain/site restriction - .edu or .gov or .org
    - Other refinements – file format, page date, where keywords appear, language, usage rights
    - Advanced search features available for image, news, video searches
        - http://www.google.com/advanced_search?hl=en
        - http://search.yahoo.com/
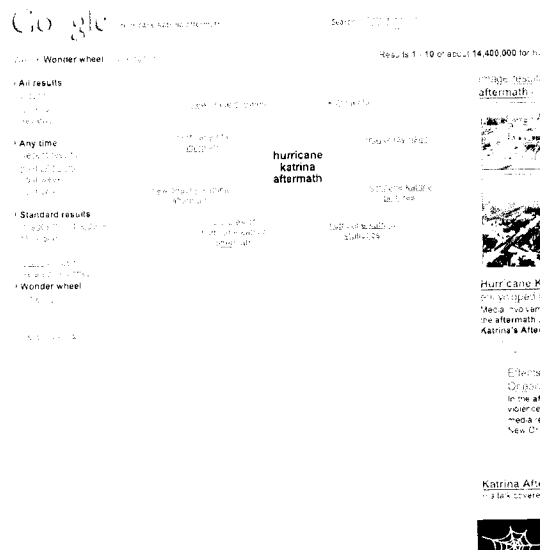
May 20, 2009                    NHLA Spring Conference                    12

**Searching the Surface Web
More Effectively**

Explore "Show Options" on Google
results page:
> All results
    Videos
    Forums
    Reviews
> Any time
    Recent results
    Past 24 hours
    Past week
    Past year
> Standard results
    Images from the page
    More text
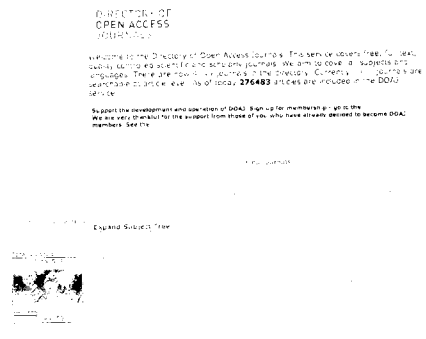> Standard view
    Related searches
    Wonder Wheel
    Timeline

Sample Search: <u>hurricane katrina</u>

May 20, 2009            NHLA Spring Conference        13

---

# Article Databases

- **Subscription databases**
  - EBSCOhost, Newsbank, etc.
- **Google Scholar**
- **Free databases**
  - Directory of Open Access Journals (DOAJ) - <u>http://www.doaj.org/</u>
  - HighWire Press - <u>http://highwire.stanford.edu/lists/freeart.dtl</u>
  - PubMed Central - <u>http://www.pubmedcentral.nih.gov/</u>

May 20, 2009            NHLA Spring Conference        14

# Split-Level Searching

- First Level
  - Use general purpose search engines to locate directories, specialized databases or search engines
- Second Level
  - Search the directory, database, search engine, etc.

May 20, 2009  NHLA Spring Conference  15

First level search strategy:
include "database" with topic term(s) in search

**nuclear explosions**          **nuclear explosions database**



May 20, 2009  NHLA Spring Conference  16

## Second level: search the database or portal

Nuclear Explosion DataBase
(NEDB)
http://www.rdss.info/database/nedb/nedb_ent.html

The NEDB consists of source information (date/time, location, yield, seismic magnitude, burial depth, etc.) drawn from a variety of official and unofficial sources, and a collection of related digital waveform data for all nuclear explosions conducted from 1945-2006. The database includes 2154 announced or presumed nuclear explosions for 837 of which over 120,000 digital waveforms have been compiled and archived.

A range of options are provided for accessing/downloading these data resources including menu-based and map-based alternatives. For a complete description of the NEDB data resources and access tools, the NEDB User Manual provides an orientation and guide including a summary of what data are available, functionality of various web-based access and display options, and step-by-step examples for several typical NEDB data queries.

Some features of the NEDB website are still being developed or supplemented.

---

# Subject Directories

- **About.com**
  - http://www.about.com/

- **Yahoo Directory**
  - http://dir.yahoo.com/

# Subject Directories

*INFOMINE*

- Librarians' Internet Index
  - http://www.lii.org

- INFOMINE
  - http://infomine.ucr.edu/

- Intute
  - http://www.intute.ac.uk/

# Subject Directories

- Open Directory Project (DMOZ)
  - http://www.dmoz.org

- World Wide Web Virtual Library
  - http://vlib.org/

# General Deep Web Search Engines

- CompletePlanet (70,000+ sites last harvested in 2004)

  http://www.completeplanet.com/

- Incy Wincy (content appears similar to DMOZ)
  - http://www.incywincy.com/

May 20, 2009                    NHLA Spring Conference                    21

# Vertical Searching

Scirus - http://www.scirus.com/



May 20, 2009                    NHLA Spring Conference                    22

# Vertical Searching

http://science.gov/ 

http://scienceresearch.com/



May 20, 2009     NHLA Spring Conference     23

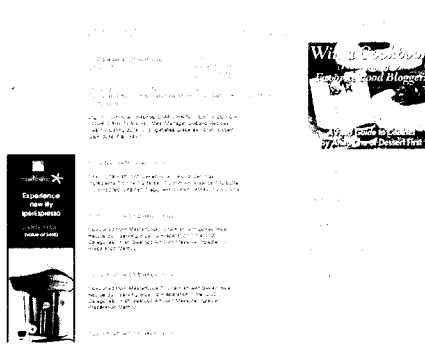# Vertical Searching

deepdyve.com

biznar.com/



May 20, 2009     NHLA Spring Conference     24

# Vertical Searching

**Foodieview.com** – a recipe search engine

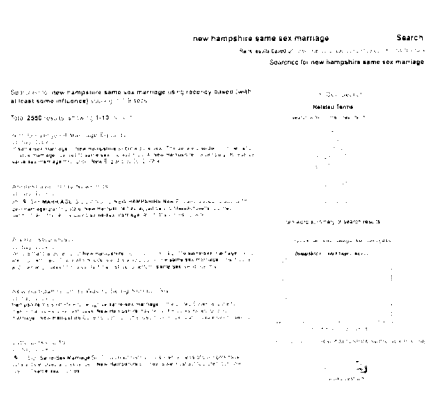**Kayak.com** – a travel search engine

# Vertical Searching

## http://www.blogscope.net/      http://technorati.com/

new hampshire same sex marriage     Search

# Vertical Searching

- Music – not just <u>iTunes</u>
  - <u>allmusic.com</u>
  - <u>Songza.com</u>
  - <u>iLike.com</u>
- Videos – not just <u>YouTube</u>
  - <u>blinkx.com</u>
  - <u>hulu.com</u>

- Additional Suggestions from Session Attendees:
  - <u>search.twitter.com</u>
  - <u>Pandora.com</u> (Internet radio service based on The Music Genome Project)

# Digital Collections

- May contain
  - digitized (i.e., scanned) books and articles
  - born-digital texts
  - audio files (e.g., wav, mp3)
  - images (e.g., tiff, gif)
  - movies (e.g., mp4, quicktime)
  - datasets (e.g., downloadable statistics files)
- May be freely accessible, restricted, or a combination

- A Compendium of Digital Collections *(A project of the UNH Library Digital Collections Initiative)* - <u>http://digitalcollections.wordpress.com/</u>

- Library of Congress Digital Collections & Services - <u>http://www.loc.gov/library/libarch-digital.html</u>

- Collaborative Digitization Programs in the United States – <u>http://frank.mtsu.edu/~kmiddlet/stateportals.html</u>

- <u>PolicyArchive.org</u>

# Digital Collections

- OAI – <u>Open Archives Initiative</u>

- OAIster - <u>http://www.oaister.org/</u>
  - Union catalog of digital resources
  - 20,928,590 records from 1,112 contributors

- Goals
  - One-stop "shopping" for academically-oriented digital resources
  - Access to digital materials in addition to searching of records for these resources

# OAIster Search

# OAIster Search

# Internet Archive

- Wayback Machine –
  http://www.archive.org/web/web.php
- Open Library Project
- Open Content Alliance – digitization of out-of-copyright books and accessible through a permanent archive
  - Text archive - http://www.archive.org/details/texts

# Internet Archive

Sample text search:
   durham nh report

Refine search by Related Creator: Durham (N.H.: Town)

Sort by Date

Can view PDF, B&W PDF, full text, Flip Book

May 20, 2009       NHLA Spring Conference       33

# UNH Digital Collections Initiative Database

- http://www.library.unh.edu/diglib/
- Simple search searches metadata only
- Advanced search can search within full text
- Only format is PDF

May 20, 2009       NHLA Spring Conference       34

# The Future

- Ever-shifting boundaries
- Google's efforts
- Different search & presentation methods
  - Kosmix.com
  - DeepPeep.org
- Semantic Deep Web?

# Keeping Up

- Search Engine Land
  - http://searchengineland.com/

- LibWorm - http://www.libworm.com/

- Google Alerts
  - http://www.google.com/alerts?hl=en&gl=us

# If You Want to Know More

Michael K. Bergman, 2001. "The 'Deep' Web: Surfacing Hidden Value" (A BrightPlanet White Paper). http://www.brightplanet.com/white-papers/119.html?task=view
[Also in *Journal of Electronic Publishing*, 7(1), online, August 2001].

Michael K. Bergman's Blogasbord : "AI[3]"– Deep Web Category
http://www.mkbergman.com/?cat=6

Anne Blecksmith, 2008. "Visual Resources Online: Digital Images of Primary Materials on Public Web Sites." *C&RL News*, 69 (5), pp. 275-278.
http://www.ala.org/ala/mgrps/divs/acrl/publications/crlnews/2008/may/visualresources.cfm
[This article is the basis of Sarah Perez's May 14, 2008 posting on ReadWriteWeb: "Digital Image Resources on the Deep Web." Comments to this posting include a number of other suggested resources. http://www.readwriteweb.com/archives/digital_image_resources_on_the_deep_web.php]

May 20, 2009      NHLA Spring Conference      37

# If You Want to Know More

Wendy Boswell, 2007. *The About.com Guide to Online Research*, Avon, MA: Adams Media. [Has chapters on the Invisible Web, searching the blogosphere, multimedia, social networks, and more.]

Jane Devine & Francine Egger-Sider, 2009. *Going Beyond Google: The Invisible Web in Learning and Teaching*. New York: Neal-Schuman Publishers.
[Good information & examples for teaching students about the Deep Web]

Jane Devine & Francine Egger-Sider, 2009. "Beyond Google: The Invisible Web" [Materials from their workshop]
http://library.laguardia.edu/invisibleweb

Bin He, Mitesh Patel, Zhen Zhang, & Kevin Chen-Chuan Chang, 2007. "Accessing the Deep Web." *Communications of the ACM*, 50 (5), pp. 95-101, May 2007.

May 20, 2009      NHLA Spring Conference      38

# If You Want to Know More

Alisa Miller, 2009. "100 Useful Tips and Tools to Research the Deep Web." Online College Blog. http://www.online-college-blog.com/index.php/features/100-useful-tips-and-tools-to-research-the-deep-web/

Chris Sherman and Gary Price, 2001. *The Invisible Web: Uncovering Information Sources Search Engines Can't See.* Medford, N.J. : CyberAge Books, Information Today. [The first published book about the Invisible Web. Chris Sherman is now the Executive Editor of SearchEngineLand.com. Gary Price, a librarian and online information consultant, is the founder and chief editor/compiler of ResourceShelf, a daily newsletter with resources of interest for online searchers.]

Alex Wright, 2008. "Searching the Deep Web," *Communications of the ACM,* 51 (10), pp. 14-15, October 2008. [Subscription only access at http://mags.acm.org/communications/200810/?CFID=5461527&CFTOKEN=11076271]

May 20, 2009         NHLA Spring Conference       39

# If You Want to Know More

Alex Wright, 2009. "Exploring a 'Deep web' That Google Can't Grasp" *The New York Times,* February 23, 2009, p. B4. http://www.nytimes.com/2009/02/23/technology/internet/23search.html?_r=1&scp=1&sq=deep%20web&st=cse

Marcus P. Zillman, 2009. Deep Web Research Resources and Sites. http://deepwebresearch.info/. [This is a Subject Tracer™ Information Blog developed and created by the Virtual Private Library™.]

May 20, 2009         NHLA Spring Conference       40

# Illustration Credits

- Spider Web that appears on every slide – Courtesy of Designed to a T

    - http://www.designedtoat.com/clipart2/hal11.gif

- Fish Trawler – Image Source: http://www.mkbergman.com/?p=458

- Spider on slide #5 – Courtesy of Webweaver - http://www.webweaver.nu/

- Other illustrations were screen captures by the presenter.


All links in this presentation were accessible on May 17, 2009.

May 20, 2009                          NHLA Spring Conference                          41