

Introduction to Information Retrieval and Boolean Query

Lecture 1

CS 510

Information Retrieval on the Internet

IR 2010

1

Information retrieval

“Information retrieval (IR) deals with the representation, storage, organization of, and access to information items.”

- Baeza-Yates and Berthier Ribeiro-Neto in [Modern Information Retrieval](#), p. 1

“*Information retrieval* is often regarded as being synonymous with *document retrieval* and, nowadays, with *text retrieval*, implying that the task of an IR system is to retrieve documents or texts with information content that is *relevant* to a user’s information *need* ... the approaches that have been developed for this purpose are also applicable to a whole family of related information processing tasks that lie between, on the one hand, *data retrieval* and, on the other, *fact or knowledge retrieval*.”

- Sparck Jones and Willett in [Readings in Information Retrieval](#), p. 1

IR 2010

2

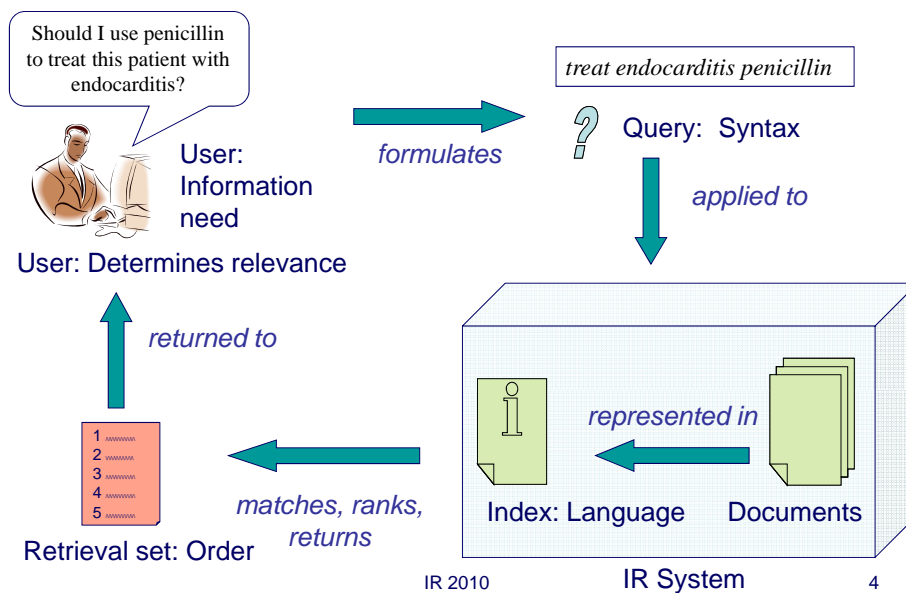
Information Retrieval

- History
 - Early work in IR in the 50's and 60's
 - Roots in library science much older, e.g.
 - Dewey Decimal system 1870s
 - Library of Congress Classification 1890s
- Important Components:
 - Indexing
 - Searching
 - User-system interaction

IR 2010

3

Typical information retrieval (IR) system



4

IR systems vs. DBMS

	DBMS	IR System
Target	Data (structured)	Information (text, multimedia)
Query/ Query language	Precise specification SQL	Imprecise specification Controlled vocabulary?
Matching	Exact match	“Similarity”
Results	Set of records	Set or ranked list of documents

IR 2010

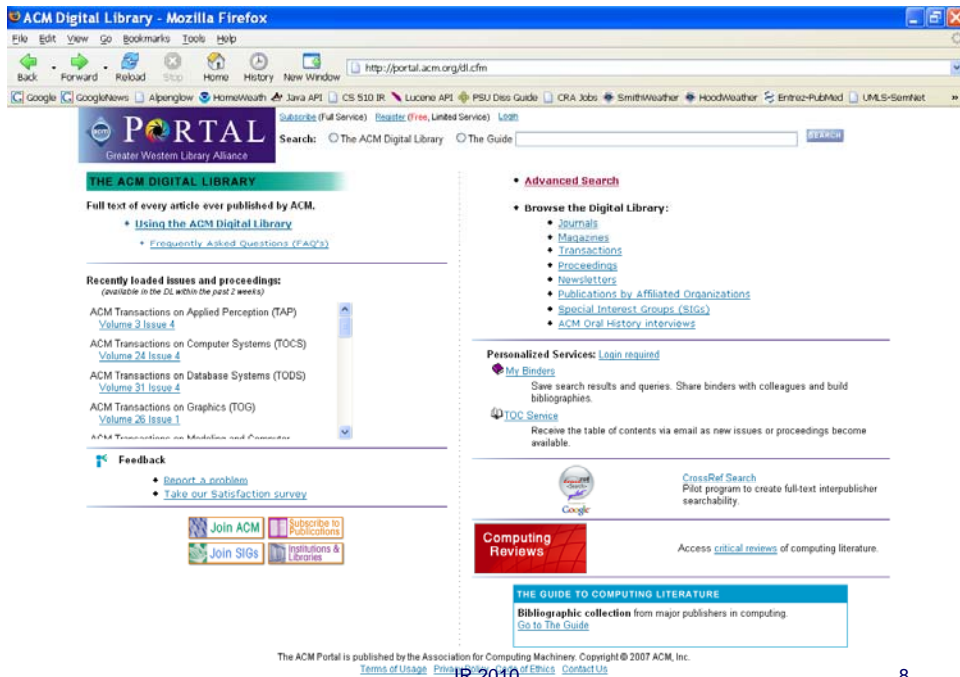
5

General types of IR systems

- Web
- Full text documents
- Bibliographic
- Distributed variations
 - Metasearch
 - Virtual document collections

IR 2010

6




www.ncstrl.org - Mozilla Firefox

File Edit View go bookmarks Tools Help

Back Forward Reload Stop Home History New Window

http://www.ncstrl.org/

Google GoogleNews Alpbanglow HomeWeather Java API CS SIO IR Lucene API PSU Diss Guide CRA Jobs SmithWeather HoodWeather Entrez-PubMed UMLS-SemNet

 **Networked Computer Science Technical Reference Library**

Simple Search Advanced Search Browse Register Submit to CoRR About NCSTR OAI Help

SEARCH RESULTS GROUPED BY ARCHIVE

This is page 1 of total 14 pages, hits (1-10) of total 122 hits.


Results Pages: 1 2 3 4 5 6 7 8 9 10 11 12 13 14

ARCHIVE .OGI

Title	<i>Multi-Stage Programming: Its Theory and Applications</i>
Authors	Walid Taha
Archive	OGI
Discovery Date	1999-11-01
Document ID	oaincstr@og_cse.ncstrl.org_cse/CSE-99-TH-002
Title	<i>Understanding and Adapting to Speaker Variability Using Correlation-based Principal Component Analysis</i>
Authors	Zhaohong Hu
Archive	OGI
Discovery Date	1999-10-01
Document ID	oaincstr@og_cse.ncstrl.org_cse/CSE-99-TH-001
Title	<i>The Efficiency of Multimodal Interaction for a Map-based Task</i>
Authors	Choben, Philip R. Mc-Gee, David Clow, Josh
Abstract	This paper compares the efficiency of using a standard direct-manipulation graphical user interface (GUI) with that of using the CheckIt voice/multimodal interface for navigation.

Now 3.0 users, please enable third party cookies

This is a collaborative project involving [NASA Langley](#), [Old Dominion University](#), [University of Virginia](#) and [Virginia Tech](#).

Powered by  Old Dominion University

IR 2010 9

Done

Dogpile Web Search Home Page - Mozilla Firefox

File Edit View go bookmarks Tools Help

Back Forward Reload Stop Home History New Window

http://www.dogpile.com/

Google GoogleNews Alpbanglow HomeWeather Java API CS SIO IR Lucene API PSU Diss Guide CRA Jobs SmithWeather HoodWeather Entrez-PubMed UMLS-SemNet

DOGPILE
All the best search engines piled into one.

Google YAHOO!SEARCH msn Search MSN

Web Images Audio Video News Yellow Pages White Pages

Go Fetch! Advanced

Review the Hottest Searches of 2006!

 **SPEAK!**

What Do You Think?
[Give your feedback on Dogpile.](#)

Download Toolbar | Joke of the Day | Search4py | Maps | Weather | Horoscopes | Add Dogpile Search to Your Site

Favorite Fetches: Part of Miami | Iraq Policy | Tax Software | Travel | 46664 | Two Headed Calf

Submit Your Site | Tell a Friend | [About Dogpile](#) | [Contact Us](#) | [Tools & Tips](#) | [Privacy Policy](#) | [Terms of Use](#)

InfoSpace  [About](#) | [Mobile](#) | [Search & Directory](#) | [Careers](#) | [Press](#) | [Investor Relations](#)

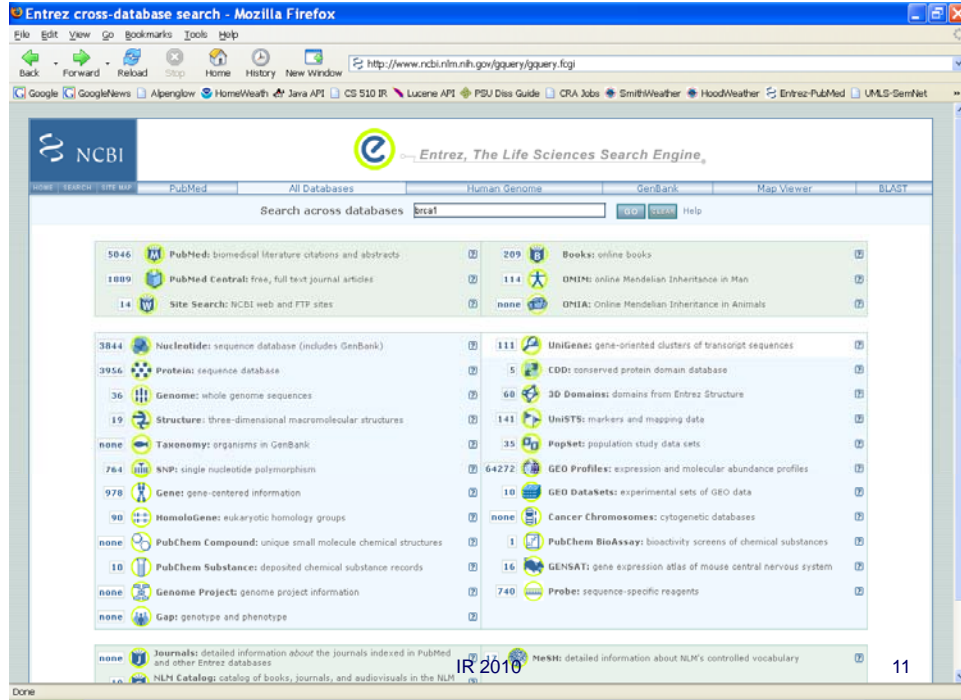
InfoSpace Network: Dogpile | WebCrawler | MetaCrawler | InfoSpace | Movers

InfoSpace Resources: Search Engines | Yellow Pages | White Pages | Town Directories | Business Categories | Maps | Ringlines

© 2007 InfoSpace, Inc. All Rights Reserved.

IR 2010 10

Done



Facets of IR systems: Scope

- Web
- Limited in Focus, Kind or Location
 - PubMed
 - Patents, News Articles
 - Company documents
- Local
 - Files
 - Email

Facets of IR systems: Documents

- Format
 - HTML
 - PDF
 - Word processed
 - Scanned OCR
- Type
 - Text
 - Multimedia
 - Semistructured (XML)
- Dynamic

IR 2010

13

Facets of IR systems: Engine

- Supported tasks: search, browse
- Basic model: Boolean vs. ranked results
- Indexing language: controlled keywords vs. natural language
- Indexing target: bibliographic data vs. full text
- Search syntax; available operations

IR 2010

14

Fact of IR Systems: Content vs. Context

- Content only
 - Boolean Query – presence or absence of term
- Extended content
 - Meta-data, anchor text
- Collection statistics
 - Word frequency: “snake” vs. “anaconda”
- Collection structure
 - In-links

IR 2010

15

User Tasks

- Mode
 - Search (retrieve)
 - Browse
- Purpose
 - Overview
 - Question answering/fact finding
 - Comprehensive research
 - Finding known item (document, page, or site)
 - Transaction (e.g. buy a book, download a file)

IR 2010

16

Document Representation: Logical

- Goal:
 - Represent the content
 - Represent other aspects (sometimes)
- Methods:
 - Assign descriptors (usually selected from a predetermined list)
 - Extract features (usually words or phrases if text document)
 - Descriptors provided by others: taggin

IR 2010

17

Document Representation: Logical

- History
 - Early systems mostly bibliographic
 - Systems contained brief surrogates for documents
 - Title, author, abstract, location of full document
 - Remember card catalogs??
 - Provided reference to the full document
 - Limited by available storage and processing power
 - Affected the way document content was represented

IR 2010

18

Document Representation: Logical

- **Methods:**
 - Manual
 - usually assign terms from a controlled vocabulary
 - Automatic
 - usually extract terms from the document
- **Considerations:**
 - Size of the representation
 - Improve likelihood of appropriate matching to queries

IR 2010

19

Document Representation: Logical

Abstract: The objectives of this study were to determine if (1) children with migraine experience greater sleep disturbances than their siblings, (2) those with more severe migraine have greater levels of sleep disruption, and (3) these sleep disturbances lead to greater behavioral problems and more missed school. Children aged 6 to 18 years with a diagnosis of migraine for > 6 months, who had at least one sibling without migraine in the same age range, were identified through our neurology clinic database or at the time of the clinic visit. Parents completed the (1) demographic, general health, and migraine information questionnaire; (2) Child Sleep Habits Questionnaire; and (3) Behavior Assessment System for Children: Second Edition (BASC-2) Parent Rating Scales for each child. Cases with migraine had higher total sleep ($P < .02$), sleep delay ($P < .03$), and daytime sleepiness scores ($P < .001$) than controls. Cases with more severe migraines had higher total sleep ($P < .01$) and sleep duration scores ($P < .03$) than those with milder headaches. In cases, higher total sleep...

Assigned Indexing Terms:

Activities of Daily Living
Adolescent
Case-Control Studies
Child
Child Behavior
Disorders/*complications/psychology
Circadian Rhythm/physiology
Female
Humans
Male
Migraine
Disorders/*complications/psychology
Severity of Illness Index
Sleep
Disorders/*complications/psychology
Sleep Stages/physiology
Wakefulness/*physiology

IR 2010

20

Document Representation: Logical

- Techniques to improve representation
 - Remove stopwords
 - Stemming
 - Exploit document structure (title, subject:)
- Possible text processing steps
 - Determine *what* to index (frames? page title? metadata?)
 - Strip formatting (e.g. html tags, w.p. instr.)
 - Recognize structure
 - Recognize words (or phrases)

IR 2010

21

Document Representation: Physical

- Associate document identifier with various descriptors
- If descriptors are extracted terms, may also include
 - Frequency of terms
 - Position of terms
- Considerations
 - Efficient storage
 - Representations that can be searched to allow a fast response to user requests

IR 2010

22

Queries

- User has an information need
- Information need communicated to IR system as a request
- Request encoded in a query expression
- Query expression must be interpretable by the system (query language)
 - Encoded in the interface
 - Terms and operators
 - Pattern syntax

IR 2010

23

Retrieval

- Match information request representation (query) to document representations (stored in an index)
- Use an algorithm to compute the matching
 - Yes/no match (Boolean expressions)
 - Similarity match
- Order and Organize results
 - Ranked by similarity (or other)
 - Clustered
 - Organized by type (images, news articles)

IR 2010

24

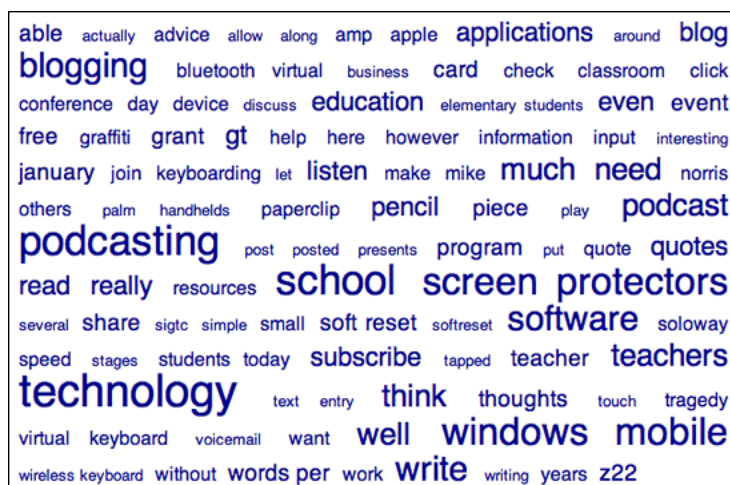
Unit of Retrieval

- Single page or file
- Site or cluster of documents
- Document fragment
- Summary or important terms (word cloud)
- Structured facts

IR 2010

25

Word Cloud



IR 2010

26

Other IR Tasks

- Classification
- Clustering
- Synoptic view
- Question answering
- Segmentation
- Summarization
- Information extraction

IR 2010

27

Themescape



IR 2010

28

System Issues

- Index organization and construction
- Scaling
- Collection building
- “Publicizing”

IR 2010

29

Course Organization

<http://web.cecs.pdx.edu/~maier/cs510iri/index.htm>

Please print off your own lecture notes for
subsequent classes

IR 2010

30

Boolean Search

Example of retrieval process

- Model: Document as set of words
- Query: Boolean combination of words
- Index: Inverted file
- Retrieval: Merging of lists

IR 2010

31

Boolean Search: Model

Document: Unit of retrieval

Term: Basic constituent of documents and queries

- Words: snake, cobra, python
- Stems: mak-, determin-
- Phrases: "garter snake"
- Special Chars: 1978, s70, help@ask.com

Document is a set of terms

- "to be or not to be" → {to,be,or,not}
- Term is in a document or not (not count, not position)

IR 2010

32

Boolean Search: Query

Terms connected with AND, OR, NOT

cobra AND python

(cobra OR python) NOT language

Could be entered through a form

Get all documents that satisfy condition,
none that don't

IR 2010

33

Boolean Search: Indexing

Inverted file: for each term, as list of
documents (IDs) containing the term

Usually sorted, elements called "postings"

cobra: 1,6,7,9,12

python: 1,5,8,9,11,12

language: 2,3,5,7,8,10,13,...

Terms organized in a dictionary, pointing to
appropriate list

IR 2010

34

Boolean Search: Retrieval

Union (OR), Intersection(AND), Different (NOT) of lists

```
cobra AND python
cobra: 1,6,7,9,12
python: 1,5,8,9,11,12
```

```
(cobra OR python) NOT language
```

```
t1 = cobra OR python
cobra: 1,6,7,9,12
python: 1,5,8,9,11,12
```

```
t1 NOT language
language: 2,3,5,7,8,10,13,...
```

IR 2010

35

Optimization

Look at sizes of lists to decide merge order

```
snake AND cobra AND python
11,345      280      156
```

Rewrite for smaller intermediate results

```
(snake OR language) AND python
```

```
11,345      23,667      156
```

```
(snake AND python) OR
(language AND python)
```

IR 2010

36

Extensions

Proximity of words

python NEAR language

python NEAR(3) language

python /P language

(same paragraph)

How can we process such queries?