# Search Engines & Information Retrieval
## 406.424 Internet Applications

**Jonghun Park**

[jonghun@snu.ac.kr](mailto:jonghun@snu.ac.kr)

**Dept. of Industrial Eng.**

**Seoul National University**

**9/1/2010**

# search and information retrieval

- search on the web is a daily activity for many people throughout the world
- **search** and **communication** are most popular uses of the computer
- applications involving search are everywhere
- the field that is most involved with R&D for search is **information retrieval (IR)**

# information retrieval

- *"Information retrieval is a field concerned with the structure, analysis, organization, storage, searching, and retrieval of information."* (Salton, 1968)
- originated from "library science" community
- primary focus of IR since the 50s has been on **text** and **documents**

# what is a document?

- examples:
  - web pages, email, books, news stories, scholarly papers, text messages, Word, Powerpoint, PDF, forum postings, patents, IM sessions, etc.

- common properties
  - significant text content
  - some structure (e.g., title, author, date for papers; subject, sender, destination for email)

# Documents vs. Database Records

- Database records (or tuples in relational databases) are typically made up of well-defined fields (or attributes)
  - e.g., bank records with account numbers, balances, names, addresses, social security numbers, dates of birth, etc.
- Easy to compare fields with well-defined semantics to queries in order to find matches
- Text is relatively **unstructured**: more difficult
- Example bank database query
  - Find records with balance > $50,000 in branches located in Amherst, MA.
  - Matches easily found by comparison with field values of records
- Example search engine query
  - bank scandals in western mass
  - This text must be compared to the text of entire news stories

# Comparing Text

- Comparing the query text to the document text and determining **what is a good match** is the core issue of information retrieval

- **Exact matching of words is not enough**
  - Many different ways to write the same thing in a "natural language" like English
  - e.g., does a news story containing the text "bank director in Amherst steals funds" match the query?
  - Some stories will be better matches than others

# Dimensions of IR

- Any application that involves a collection of text or other unstructured information will need to organize and search that information

- IR is more than just text, and more than just web search
  - although these are central

- New applications increasingly involve **new media**
  - e.g., video, photos, music, speech

- Like text, media content is difficult to describe and compare
  - text may be used to represent them (e.g. tags)

# IR Tasks

- Ad-hoc search: Find relevant documents for an arbitrary text query
- Filtering: Identify relevant user profiles for a new document
- Classification: Identify relevant labels for documents
- Question answering: Give a specific answer to a question

| Content | Applications | Tasks |
|---|---|---|
| Text | Web search | Ad hoc search |
| Images | Vertical search | Filtering |
| Video | Enterprise search | Classification |
| Scanned docs | Desktop search | Question answering |
| Audio | Forum search | |
| Music | P2P search | |
| | Literature search | |

# web search engine market



Search Engine Rankings — May 2007 by Hitwise

- Google 65.1%
- Yahoo! 20.9%
- msn 8.4%
- Ask 3.9%
- All the Rest 1.7%

Chart by www.internetworldstats.com, Copyright © 2007



BRUCE CLAY, INC
Internet Business Consultants
866-517-1900 • www.bruceclay.com

SEARCH ENGINE RELATIONSHIP CHART
United States Edition
http://www.bruceclay.com/serc.htm

LEGEND

SUPPLIES ⟶ RECEIVES PRIMARY SEARCH RESULTS
SUPPLIES ⟶ RECEIVES SECONDARY SEARCH RESULTS
SUPPLIES ⟶ RECEIVES PAID RESULTS

CLICK ON A LOGO FOR SEARCH ENGINE INFORMATION

CLICK HERE TO SELECT A DIFFERENT CHART

# niche search engines: by domain

# niche search engines: by type

# presentation of search results

# enhancing search engine usability

# big issues in IR

- **relevance**
  - simple definition: a relevant document contains the information that a person was looking for when s/he submitted a query to the search engine
  - many factors influence a person's decision about what is relevant: e.g., task, context, novelty, style
  - **topical relevance** (same topic) vs. **user relevance** (everything else)

# big issues in IR

- relevance
  - **retrieval models** define a view of relevance
    - retrieval model: a formal representation of the process of matching a query and a document
  - ranking algorithms used in search engines are based on retrieval models
  - most models describe **statistical properties** of text rather than linguistic
    - i.e. counting simple text features such as words instead of parsing and analyzing the sentences
    - **statistical approach** to text processing started with Luhn in the 50s
      - e.g., TF
    - linguistic features can be part of a statistical model

# big Issues in IR

- evaluation
    - experimental procedures and measures for comparing IR system output with user expectations
        - originated in Cranfield experiments in the 60s
    - IR evaluation methods now used in many fields
    - typically use test collection of documents, queries, and relevance judgments
        - most commonly used are **TREC collections**
    - recall and precision are two examples of effectiveness measures
        - precision: proportion of retrieved documents that are relevant
        - recall: proportion of relevant documents that are retrieved
        - assumption: all the relevant documents for a given query are known

# big issues in IR

- precision = TP / (TP + FP)
- recall = TP / (TP + FN)

# big issues in IR

- users and information needs
  - search evaluation is **user-centered**
  - keyword queries are often poor descriptions of actual information needs
    - e.g., "apple"
  - interaction and context are important for understanding user intent
    - analysis of clickthrough data
  - query refinement techniques such as **query expansion**, **query suggestion, relevance feedback** improve ranking

# PRF example

# IR and search engines

- a search engine is the practical application of information retrieval techniques to **large scale** text collections

- web search engines are best-known examples, but many others
  - Google web search engine must be able to crawl many **terabytes of data**, and then provide **subsecond response times** to millions of queries submitted everyday from around the world
  - open source search engines are important for R & D
    - e.g., Lucene, Lemur/Indri, Galago

# IR and search engines

information retrieval

- relevance

  *- effective ranking*

- evaluation

  *- testing and measuring*

- information needs

  *- user interaction*

search engines

- performance

  *- response time, query throughput, indexing speed*

- incorporating new data

  *- coverage and freshness*

- scalability

  *- growing with data and users, distributed processing*

- adaptability

  *- tuning for applications*

- specific problems

- *- e.g. spam*

# search engine issues

- **indexes** are data structures designed to improve search efficiency
  - designing and implementing them are major issues for search engines
- **dynamic** data
  - "collection" for most real applications is **constantly changing** in terms of updates, additions, deletions
    - e.g., web pages
  - acquiring or "crawling" the documents is a major task
    - typical measures are **coverage** (how much has been indexed) and **freshness** (how recently was it indexed)
  - updating the indexes while processing queries is also a design issue

# spam

- for web search, spam in all its forms is one of the major issues

- affects the efficiency of search engines and, more seriously, the effectiveness of the results

- many types of spam
  - e.g. spamdexing or term spam, link spam, "optimization"

- new subfield called **adversarial IR**, since spammers are "adversaries" with different goals
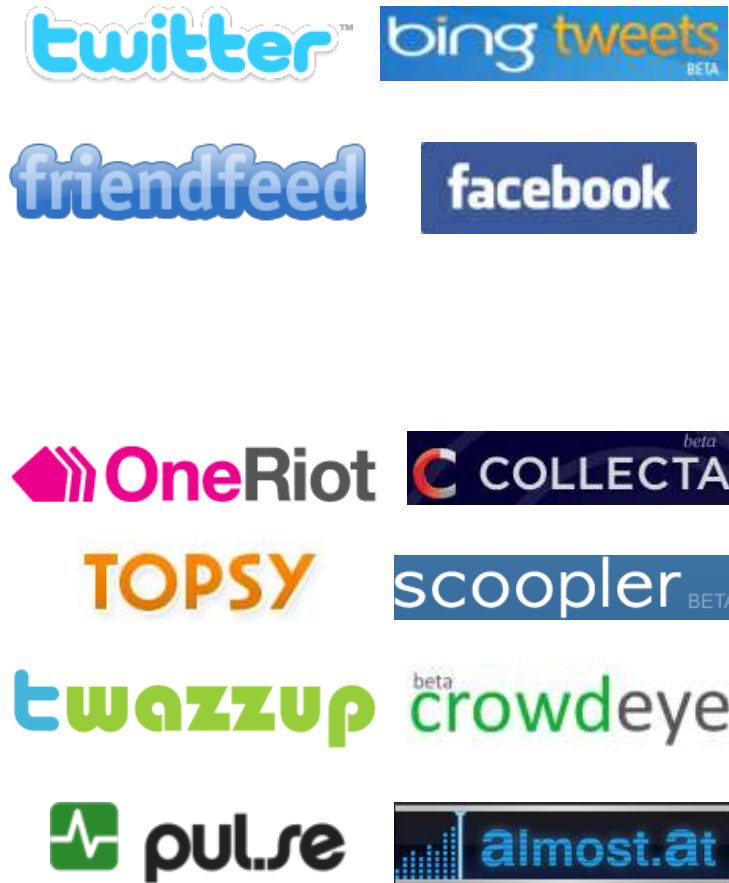
# search engine mega trend



Relevance Features

Time

Real-Time

Social

Semantic

Quantity/Quality

# real-time search examples

# real-time search engines

# history of search services in Korea

2001   **NAVER** 통합검색

2002   **NAVER** 지식검색

2003   **empas** 지식검색

2004   Google 한국어 검색

2005   **Daum** 신지식 검색 **empas** 열린검색

2006   블로그, 까페, 게시판 검색

2007   동영상검색   **Daum** 자체검색엔진 개발

2008   **empas** **NATE** 합병

2009   **NATE** 시맨틱 검색 **Daum** 소셜 검색

2010   **LIVE K** 실시간 검색
살아있는검색

포털 통합검색 점유율

구글 2.45    파란 0.29
야후 3.59

네이트 9.81

다음 20.15

(단위:%)

네이버 63.72

〈자료:코리안클릭〉

# IR related conferences

- SIGIR: ACM Special Interest Group on IR
- WWW: World Wide Web Conference
- ECIR: European Conference on IR
- CIKM: Conference on Information and Knowledge Management
- WSDM: Web Search and Data Mining Conference
- WISE: Web Information System Engineering
- ICWE: International Conference on Web Engineering
- HT: Hypertext and Hypermedia