Dr.-Ing. Gjergji Kasneci

*gjergji.kasneci@hpi.uni-potsdam.de*

HPI Potsdam, Winter Term 2012-13

# INTRODUCTION TO INFORMATION RETRIEVAL

# Organization

- ➢ Timetable
  - ➢ Lectures: Mondays 13:30 – 15:00, room H-E.51
  - ➢ Exercises: biweekly, Tuesdays 15:30 – 17:00, room H-E.52

- ➢ Exam
  - ➢ Condition for admission: oral presentation of at least two solutions during the tutorials
  - ➢ Form of exam: oral exam at the end of the term

# What is this lecture about?

"Know where to find the information and how
   to use it – That's the secret of success."

**Albert Einstein**

➢ Subject

Automated techniques for finding information relevant to an information need from a large corpus of data.

- ➢ Content & structure analysis
- ➢ Indexing, search, ranking
- ➢ Classification, clustering
- ➢ Annotation, extraction, summarization, visualization
- ➢ Personalization & collaboration
- ➢ Interaction with knowledge bases

➢ **Web search**
  ➢ Web pages (or multimedia) about people, products, locations, events, news, problem solutions, …
  ➢ Main problems: scalability, spam, format diversity, deep (hidden) web, …

➢ **Enterprise search**
  ➢ Intranet documents/pages (or multimedia) about people/experts, products, divisions, problem solutions, …
  ➢ Main problems: low information redundancy, few links, security issues, …

➢ Indexing and retrieval for knowledge discovery & data mining tasks
  ➢ Bioinformatics: finding promising medical therapies, retrieving salient interaction paths in biochemical networks, …
  ➢ Business intelligence: risk identification, composition of equity baskets, …
  ➢ Scientific observations: identifying typical or noticeable patterns
  ➢ Socio-political observations: analyzing social concerns, developments, political issues, …

# Differences to database systems

| | Information retrieval systems | Databases |
|---|---|---|
| **Corpus** | Unstructured, semi-structured information (text, multimedia, etc.) | Structured, semi-structured data, e.g., database table(s), XML-document(s) |
| **Query** | Imprecise, keyword-based, typically with "and'ish" search semantics | Precise, structured query language, e.g., SQL, Xpath, Datalog, SPARQL … |
| **Retrieval** | Approximate, frequency-, similarity-based | Exact match of attribute values |
| **Results** | Ranked by relevance | Unranked, or ranked by specific attribute(s) (from the order-by clause of the query) |

# Historic background

➢ 1945: far-reaching ideas about personal & collective information gathering, indexing, aligning, associating, sharing, reasoning, … mentioned in "As we may think" by Vannevar Bush



Vannevar Bush



Memex extends human memory and reasoning capabilities

# Roots in search over library catalogs

➤ Early work on IR, in the 50's and 60's

  ➤ Boolean search (i.e., combine search terms with Boolean operators: **and, or, not**) unsatisfactory

  ➤ Need for approximate search with "vaguely" Boolean operators needed (e.g., "**and'ish**" search semantics)

  ➤ Need for ranked results

# Early research

- ➢ **Hans Peter Luhn** (IBM researcher)
  - ➢ 1957: published "A Statistical Approach to Mechanized Encoding and Searching of Literary Information"

- ➢ **Gerard Salton** (professor at Cornell)
  - ➢ **1960's**: supervised the development of **SMART** (System for the Mechanical Analysis and Retrieval of Text)
  - ➢ 1969: published "*Automatic Information Organization and Retrieval*"
  - ➢ Introduces important concepts, e.g., **vector space model**, **relevance feedback, Rocchio classififcation**, …

- ➢ Other influenatial scientists
  - ➢ **Karen Spärck Jones** introduced **inverse-document-frequency** in 1972
  - ➢ **Steven Robertson** introduced **Okapi BM25** ranking model together with Karen Spärck Jones in 1976

~ 630 Mio. active sites

Numbers from:
http://news.netcraft.com/archives/2012/

Web 2.0 era

Collaborating

(real-time)

~ 315 Mio. active sites

Searching

Publishing

| 1969 ARPANET | 1989 WWW "Mesh" proposal by Berners Lee | 1998 Google | 2000/01 Wikipedia | 2004/05 Social networks Web 2.0 | 2009 Cloud, real-time, mobile |

# General infrastructure for web IR



Source: Modern Information Retrieval,
Baeza-Yates, Ribeiro-Neto

# Crawling

# Indexing

**Document analyzer**
(Structure, language, detect spam, duplicated or dynamic content …)

**Cleanser**
(Remove HTML tags and scripts, tokenize, stem, remove stop words …)

**Score allocator**
(Term frequency, inverted document frequency, …)

**Thesaurus**
(with synonyms, hyponyms, hypernyms, meronyms, …)

**Inverted Index (B+ tree)**
camera          iphone

URLs & scores          URLs & scores

# Query answering

query (expressing search need)

**User interface**

**Query processor**
(parsing, cleansing, spell correction, entity recognition, query expansion, reformulation, …)

**Result provider**
(top-k retrieval & ranking by relevance, authority, freshnes, ….)

**Thesaurus**
(with synonyms, hyponyms, hypernyms, meronyms, …)

**Inverted Index (B+ tree)**
camera          iphone

URLs & scores    URLs & scores

# Merging results through metasearch



query    response

Metasearch engine

User interface

Query distributor

Result aggregator

Search engine 1

Search engine 2

Search engine 3

Web

# Metasearch example

# Taxonomy of search on the web

➢ Typical semantics: **and'ish** search, e.g., "cheap flights Rome"

➢ Categorization of query types (by A. Broder, "A taxonomy of web search", 2002) through query log analysis and user study:

  ➢ **~50% informational queries**: "General relativity", "Capital of Taiwan", "Obama campaign", …

  ➢ **~20% navigational queries**: "facebook", "amazon", "ebay", …

  ➢ **~30% transactional/resource-oriented queries**: "hotels Seattle", "flights Berlin Rome", "eminem mp3", "skype download", …

➢ Different numbers presented in study by Rose & Levinson in 2004 (~62% informational, ~13% navigational, ~25% transactional queries)

# Advanced query types

➢ Regular expression queries (e.g., on Google):

-inurl:(htm|html|php) intitle:"index of" +"last modified" +"parent directory" +description +size +(wma|mp3) "Nirvana"

➢ Conversion queries:

10 Euro in USD

➢ Calculation queries:

110*(654/8+3)

➢ Natural language questions:

What is the time in Hong Kong

# Search results

18

# Ranking results

➢ **By relevance**

 ➢ Binary (relevant or non-relevant).

 ➢ Query/information-need dependent (i.e., similarity estimation between query terms and each document)

➢ **By importance**

 ➢ Authority/popularity of result page (e.g., likelihood of being useful to a majority of users)

   ➢ Links pointing to that page (PageRank, HITS)

   ➢ Estimation of click-through-rate

   ➢ …

 ➢ Interestingness & context (e.g., likelihood of being useful for specific user in specific context → personalization)

 ➢ Freshness (e.g., most recent news).

Modern search engines combine relevance and importance criteria for ranking!

Source: http://yippy.com/

Grouping topical categories

# Visualizing result clusters



Source: http://search.carrotsearch.com/

# Human edited categories in DMOZ



dmoz open directory project

about dmoz

[Search] | the entire directory ▼ |

**Top**: **Regional**: **Europe**: **Germany** *(3,986)*

- **Localities** *(1)*
- **States** *(3,273)*

- **Arts and Entertainment** *(34)*
- **Business and Economy** *(227)*
- **Education** *(10)*
- **Employment**@ *(5)*
- **Government** *(103)*
- **Guides and Directories** *(11)*
- **Health** *(15)*
- **Maps and Views** *(5)*
- **News and Media** *(21)*

- **Real Estate**@ *(7)*
- **Recreation and S**
- **Science and Envi**
- **Shopping**@ *(23)*
- **Society and Cultu**
- **Transportation** *(9*
- **Travel and Touri**
- **Weather** *(7)*

dmoz open directory project

about dmoz

[Search]

**Top**: **Regional**: **Europe**: **Germany**: **Arts and Entertainment**

- **Berlin**@ *(86)*
- **Bremen**@ *(4)*
- **Hamburg**@ *(2)*
- **Saarland**@ *(4)*

- **Architecture** *(3)*
- **Bodyart Studios**@ *(1)*
- **Clubs and Venues** *(0)*
- **Dance** *(9)*
- **Film Festivals**@ *(11)*
- **Film Societies**@ *(0)*
- **Galleries**@ *(31)*

- **Literature**@ *(10)*
- **Magicians**@ *(5)*
- **Modeling Agencies**@ *(13)*
- **Museums** *(0)*
- **Music** *(17)*
- **Theatre** *(2)*

23

# Semantic search: question answering

# Computational question answering



WolframAlpha™ computational... knowledge engine

Population development in China and India

**Input interpretation:** *Mathematica* form

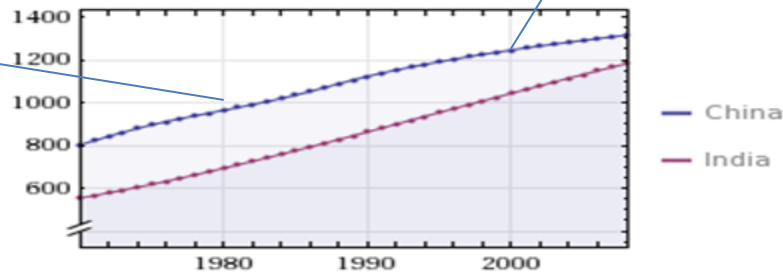| China / India | population | time series |

**Results:**

| China | mean | 577.5 million people |
| | lowest | 138 million people (1700) |
| | highest | 1.314 billion people (2008) |
| India | mean | 507 million people |
| | lowest | 113.3 million people (1600) |
| | highest | 1.181 billion people (2008) |

Population and Family Planning Law

One-Child Policy

Recent population history: Show projections | Log scale

— China
— India

27

# Community-based question answering

# Mobile web search



- ➢ Multimodal (videos, images, sounds, …)

- ➢ Personalized

- ➢ Context-aware

- ➢ Natural language

- ➢ Exact (answers)

- ➢ Social

- ➢ …

Source: http://googlemobile.blogspot.de/

In a nutshell, Goggles lets users search for objects using images rather than words. Simply take a picture with your phone's camera, and if we recognize the item, Goggles returns relevant search results. Right now Goggles identifies landmarks, works of art, and products (among other things), and in all cases its ability to "see further" is rooted in powerful computing, pervasive connectivity, and the cloud:

- We first send the user's image to Google's datacenters
- We then create signatures of objects in the image using computer vision algorithms
- We then compare signatures against all other known items in our image recognition databases; and
- We then figure out how many matches exist; and
- We then return one or more search results, based on available meta data and ranking signals; and
- We do all of this in just a few seconds

# Social annotation for image search



> ➢ Game with a purpose (exploits "wisdom of the crowds")
> ➢ Project initiated by Luis von Ahn et al. In 2004
> ➢ Estimation: 5000 people could tag all images on the web in 4 weeks

➢ **Social collaboration**
  ➢ Tagging, answering, computing, rating, editing, recommending, discussing, …
  ➢ Examples: Wikipedia, flickr, del.icio.us, librarything, Amazon, Netflix, answers.com, ESP, foldit, …

# Google Advertising Programs

## For Advertisers: Google AdWords

Reach people when they are actively looking for information about your products and services online, and send targeted visitors directly to what you are offering. With AdWords cost-per-click pricing, it's easy to control costs—and you only pay when people click on your ad.
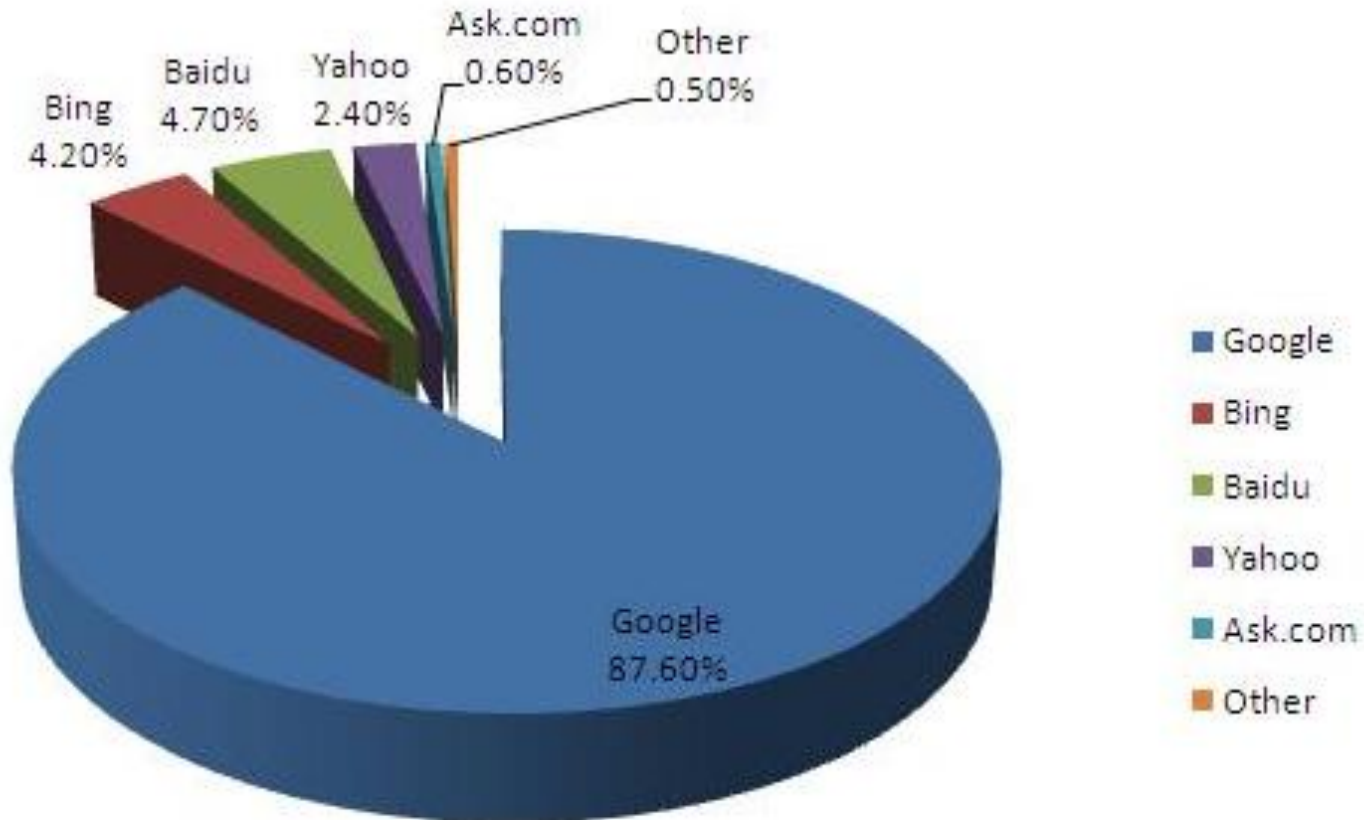
Search Term

AdWords Ads

Search Results

## For Web Publishers: Google AdSense

Earn more revenue from your website, while providing visitors with a more rewarding online experience. Google AdSense™ delivers ads targeted to your content pages and, when you add Google WebSearch to your site, AdSense delivers targeted ads to your search results pages too. With AdSense you earn more ad revenue with minimal effort—and no additional cost.
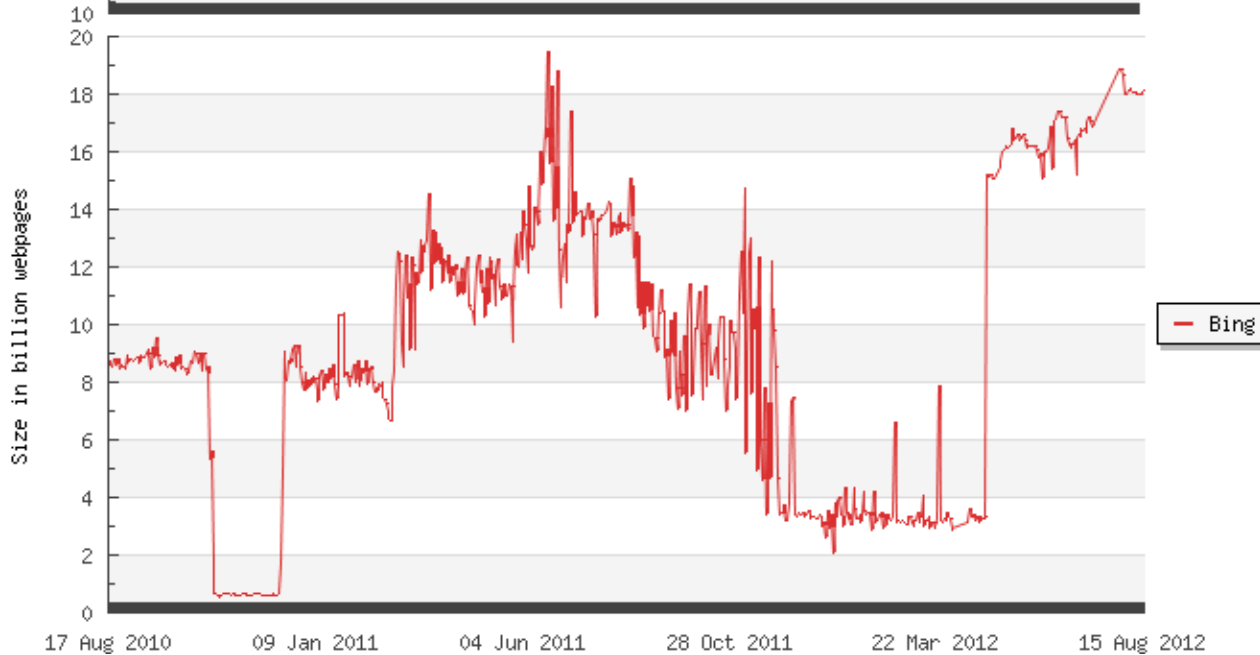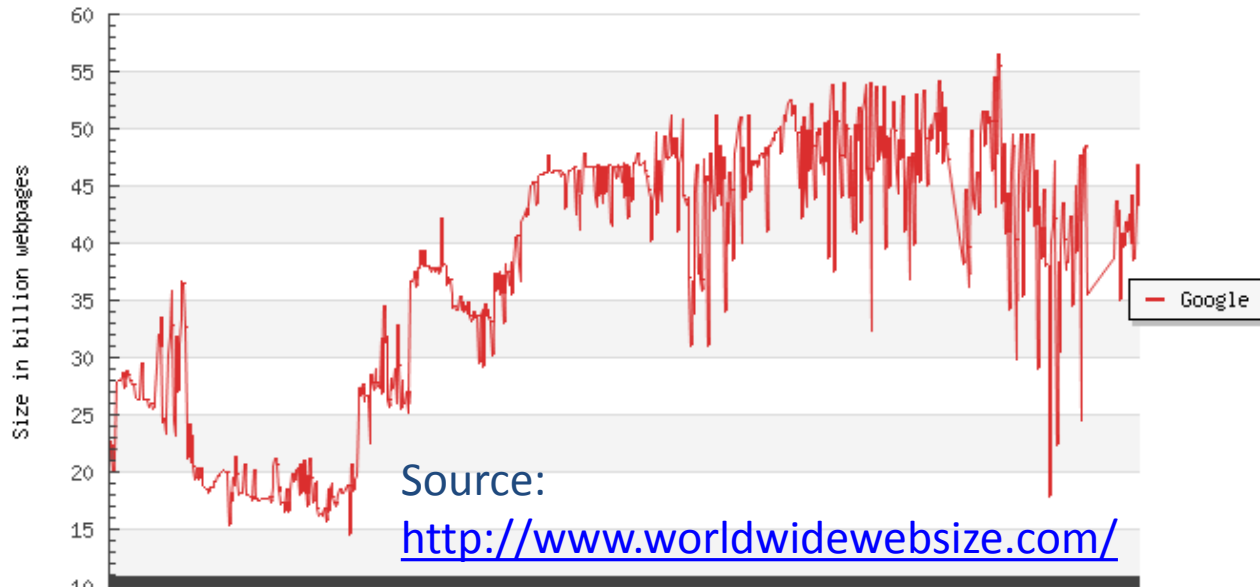
34

# Worldwide market share (as of May 2012)



Source :   http://www.karmasnack.com/about/search-engine-market-share/

# Indexed web pages



Source:
http://www.worldwidewebsize.com/

# Deep web

- Content of the web that
  - resides in backend databases and
  - is not accessible through hyperlinks
  - can only be surfaced by submitting queries to query forms

- Examples
  - Yellow pages
  - Scientific databases (bio-medical-, geo-, weather-related data, publications, measurements, …)
  - Census data (i.e., population analysis, statistics, etc.)

- Estimation: several times larger than the visible web
- High-quality & structured information

# Deep web: example (1)



Source: http://aip.completeplanet.com

# Deep web: example (2)



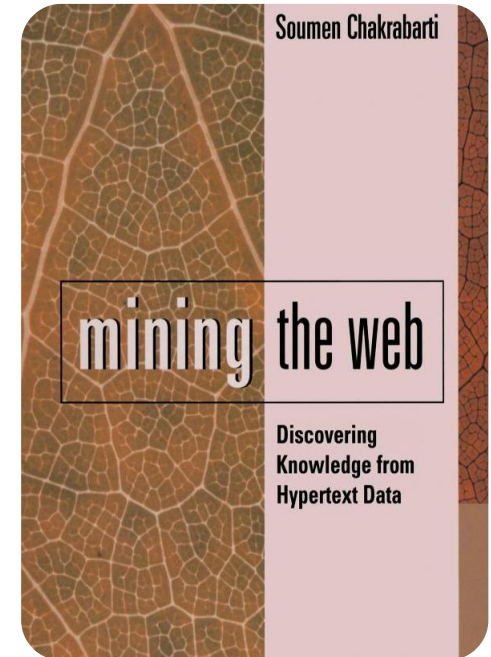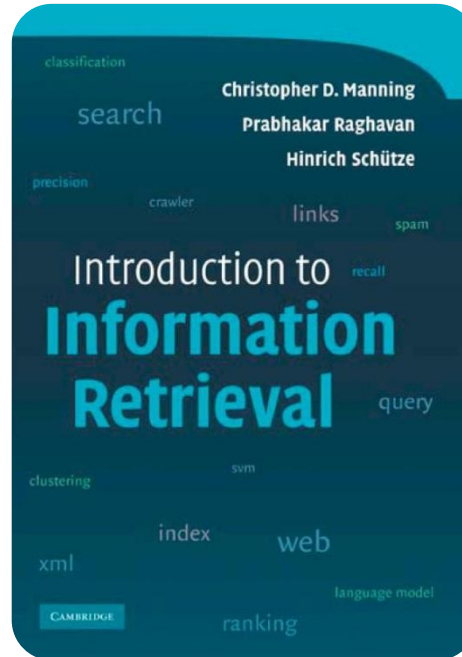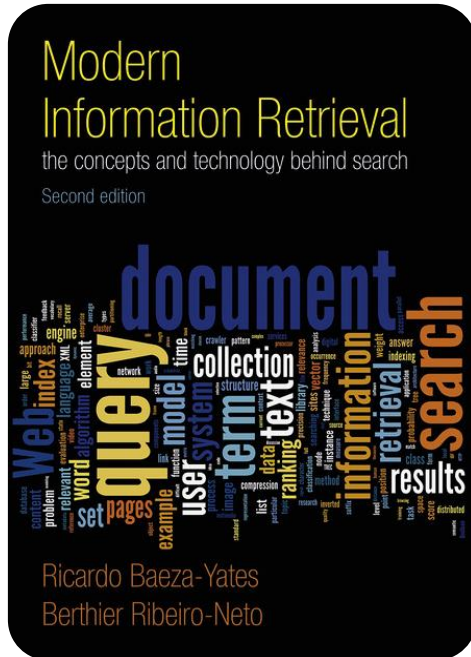Source: http://ned.ipac.caltech.edu

# Sufacing the deep web

➢ **Strategy 1** (used by Google, see [Madhavan et al. 2006](#)):

  ➢ Probe query boxes in the form with random keywords

  ➢ Use words in the results as new query terms (bootstrapping)

➢ **Strategy 2** (much more challenging -- semantic analysis):

  ➢ Infer the data type (or semantic concept) behind a query box by analyzing scripts (for required values), exploiting surrounding terms, or probing

  ➢ Surface information by querying with appropriate values (e.g. from other databases / knowledge bases)

  ➢ Analyze associations in surfaced data to infer underlying schema

➢ **Problems**

  ➢ Map queries to relevant interfaces

  ➢ Integrating schemas of sources from similar domains

  ➢ Merging results into global ranking

# Related literature

➢ **Books**



➢ **Important conferences**

    ➢ SIGIR, TREC, WSDM, WWW, WebSci, CIKM, ECIR, ICML, ECML, KDD, JCDL, …

➢ **IR toolkits**

    ➢ Lucene (http://lucene.apache.org/core/)

    ➢ The Lemur Project (http://www.lemurproject.org/)

Source: funnyjunk.com