

Information Retrieval and Data Mining

**Winter Semester 2005/06
Saarland University, Saarbrücken**

*Prof. Dr. Gerhard Weikum
weikum@mpi-inf.mpg.de*

http://www.mpi-inf.mpg.de/departments/d5/teaching/ws05_06/irdm/

Organization

- **Lectures:**

Tuesday 14-16 and Thursday 14-16 in 45/001

Office hours Prof. Weikum: appointment by e-mail

- **Assignments / Tutoring Groups:**

Friday 9-11, 11-13, or 14-16

Monday 9-11, 11-13, or 13-15

Paper assignments given out in Tuesday lecture, to be solved until next Tuesday

First paper assignment given out on Tuesday, Oct 25

First meetings of tutoring groups: Friday, Nov 4, and Monday, Nov 7

- **Requirements for obtaining 9 credit points:**

- will be announced in second week

Outline

1. Overview and System Architectures
2. Basics from Probability Theory and Statistics (1)
3. Basics from Probability Theory and Statistics (2)
4. Top-k Query Processing and Indexing
5. Advanced IR Models (1)
6. Advanced IR Models (2)
7. Advanced Link Analysis
8. Ontologies and Query Expansion
9. XML Search with Ranking
10. Peer-to-Peer Search

11. Automatic Classification
12. Clustering and Graph Mining
13. Information Extraction (1)
14. Information Extraction (2)
15. Rule Mining

***Part I:
Introduction &
Foundations***

***Part II:
Information
Search***

***Part III:
Information
Organization***

General Literature (1)

Information Retrieval:

- **Soumen Chakrabarti: Mining the Web: Analysis of Hypertext and Semi Structured Data, Morgan Kaufmann, 2002. see also <http://www.cse.iitb.ac.in/~soumen/mining-the-web/>**
- **David A. Grossman, Ophir Frieder: Information Retrieval: Algorithms and Heuristics, Springer, 2004.**
- **Christopher D. Manning, Hinrich Schütze: Foundations of Statistical Natural Language Processing, MIT Press, 1999.**
- **Ian H. Witten: Managing Gigabytes: Compressing and Indexing Documents and Images, Morgan Kaufmann, 1999.**
- **Ricardo Baeza-Yates, Berthier Ribeiro-Neto: Modern Information Retrieval, Addison-Wesley, 1999.**
- **Norbert Fuhr: Information Retrieval, Vorlesung im SS 2005, Universität Duisburg, http://www.is.informatik.uni-duisburg.de/courses/ir_ss05/index.html**
- **Christopher Manning, Prabhakar Raghavan, Hinrich Schütze: Introduction to Information Retrieval, Cambridge University Press, 2007, <http://www-csli.stanford.edu/~schuetze/information-retrieval-book.html> ; see also: <http://www.stanford.edu/class/cs276/cs276-2005-syllabus.html> <http://www.stanford.edu/class/cs276a/syllabus2004.html> <http://www.stanford.edu/class/cs276b/syllabus.html> <http://www.ims.uni-stuttgart.de/~schuetze/ws2004ir/>**
- **Pierre Baldi, Paolo Frasconi, Padhraic Smyth: Modeling the Internet and the Web - Probabilistic Methods and Algorithms, Wiley & Sons, 2003.**
- **Max-Planck Institute for Informatics, ADFOCS Summer School 2004, <http://www.mpi-inf.mpg.de/conferences/adfocs-04/program.html>**
- **Berkeley School of Information Management and Systems: Search Engines: Technology, Society, and Business, <http://www.sims.berkeley.edu/courses/is141/f05/schedule.html>**

General Literature (2)

Data Mining:

- **Jiawei Han, Micheline Kamber: Data Mining: Concepts and Techniques, Morgan Kaufmann, 2000.**
- **Ian H. Witten, Eibe Frank: Data Mining: Practical Machine Learning Tools and Techniques, Morgan Kaufmann, 2005.**
see also <http://www.cs.waikato.ac.nz/ml/weka/>
- **Margaret H. Dunham: Data Mining, Pearson Education, 2003**
- **David J. Hand, Heikki Mannila, Padhraic Smyth: Principles of Data Mining, MIT Press, 2001.**
- **Andrew Moore: Statistical Data Mining Tutorials, CMU,**
<http://www.autonlab.org/tutorials/>
- **Tobias Scheffer, Steffen Bickel: Maschinelles Lernen und Data Mining, Vorlesung SS 2004, Humboldt-Universität, Berlin,**
http://www.informatik.hu-berlin.de/Forschung_Lehre/wm/index_e.html

Foundations from Statistical Machine Learning:

- **Richard O. Duda, Peter E. Hart, David G. Stork: Pattern Classification, Wiley&Sons, 2000.**
- **Trevor Hastie, Robert Tibshirani, Jerome H. Friedman: Elements of Statistical Learning, Springer, 2001.**
- **Tom M. Mitchell: Machine Learning, McGraw-Hill, 1997.**

General Literature (3)

Foundations from Stochastics:

- **Larry Wasserman: All of Statistics, Springer, 2004.**
<http://www.stat.cmu.edu/~larry/all-of-statistics/index.html>
- **George Casella, Roger L. Berger: Statistical Inference, Duxbury, 2002.**
<http://www.stat.ufl.edu/~casella/>
- **Arnold Allen: Probability, Statistics, and Queueing Theory with Computer Science Applications, Academic Press, 1990.**

Practical Tools and Programming:

- **Ian H. Witten, Eibe Frank: Data Mining: Practical Machine Learning Tools and Techniques, Morgan Kaufmann, 2005.**
see also <http://www.cs.waikato.ac.nz/ml/weka/>
- **Erik Hatcher, Otis Gospodnetic: Lucene in Action, Manning Publications, 2004.**
see also <http://lucene.apache.org/>
- **Tony Loton: Web Content Mining with Java, John Wiley & Sons, 2002.**

Chapter 1: IRDM Applications and System Architectures

- 1.1 Overview of IRDM Technologies and Applications**
- 1.2 Web Search Engines**
- 1.3 Towards Semantic Search Engines**
- 1.4 Deep Web Search**
- 1.5 Intranet and Enterprise Search**
- 1.6 Personalized Search and Personal Info Management**
- 1.7 Peer-to-Peer Search and Collaboration**
- 1.8 Multimedia and NLP Search**

1.1 Overview of IRDM Applications and Technologies

Objective: Satisfy information demand & curiosity of human users – and eliminate the (expensive) bottleneck of human time !

Information Retrieval (IR):

- document content & structure analysis
- indexing, search, relevance ranking
- classification, grouping, segmentation
- interaction with knowledge bases
- annotation, summarization, visualization
- personalized interaction & collaboration

application areas:

- Web & Deep Web search
- intranet & enterprise search
- XML & text integration
- personalized filtering
- P2P search & collaboration
- multimedia search

Data Mining (DM):

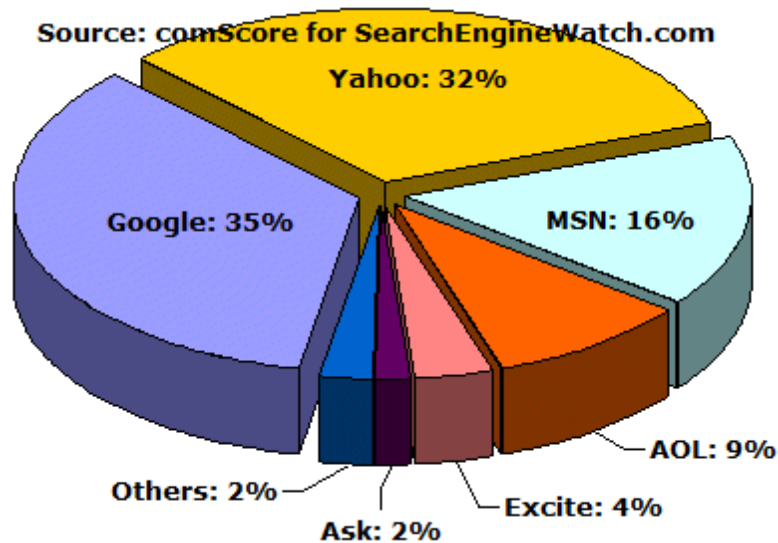
- learning predictive models from data
- pattern, rule, trend, outlier detection
- classification, grouping, segmentation
- knowledge discovery in data collections
- information extraction from text & Web
- graph mining (e.g. on Web graph)

application areas:

- bioinformatics, e.g.: protein folding, medical therapies, gene co-regulation
- business intelligence, e.g.: market baskets, CRM, loan or insurance risks
- scientific observatories, e.g.: astrophysics, Internet traffic (incl. fraud, spam, DoS)
- Web mining & ontology construction

connected to natural language processing (NLP) and statistical machine learning (ML)

1.2 Web Search Engines

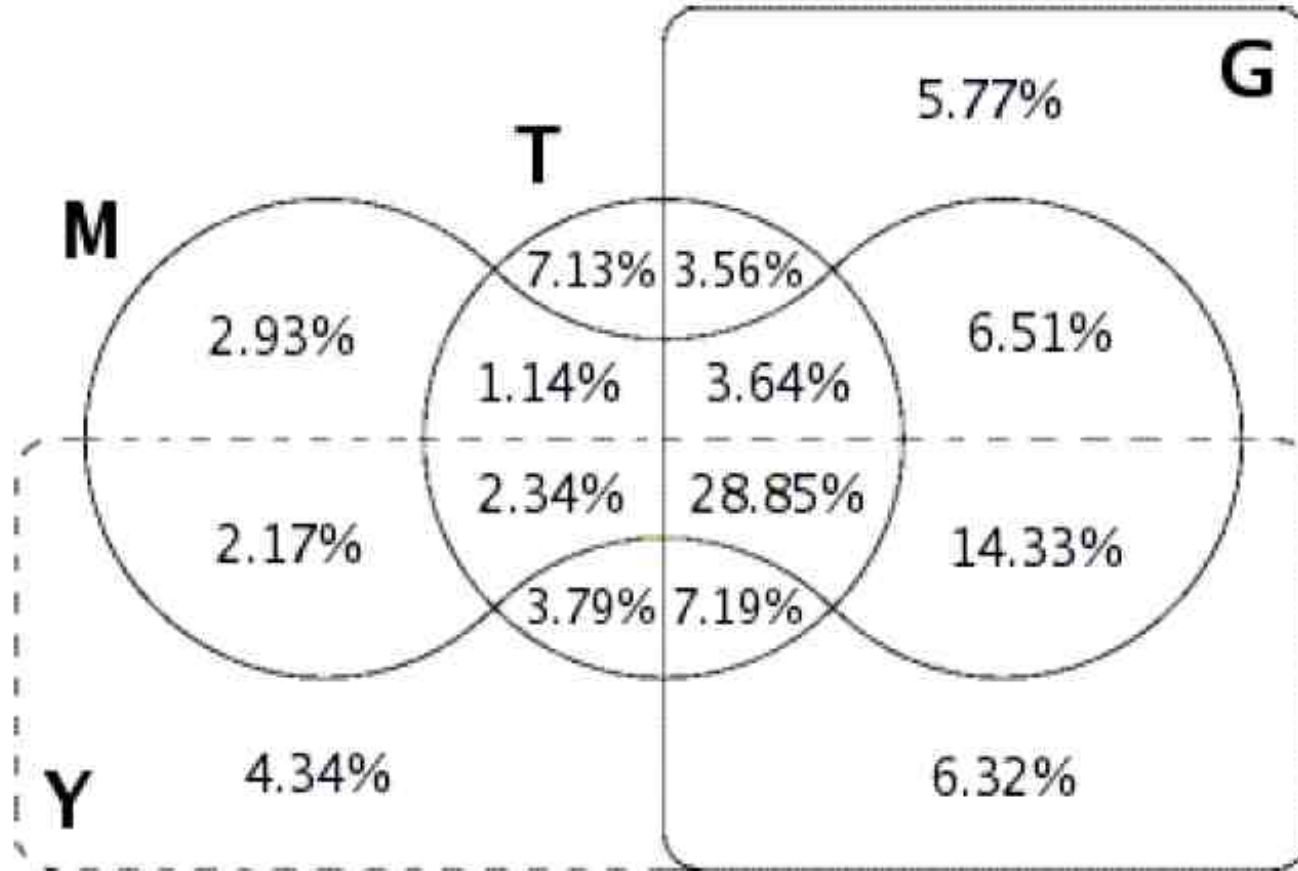


- > 11 Billion pages ($11 * 10^9$)
- > 450 Million daily queries
- > 8 Billion US \$ annual revenue

Outline:

- Web IR basics
- System architecture
- Usage patterns & quality assessment
- Limitations

Web Size and Web Coverage



Source:
A. Gulli, A. Signorini,
WWW 2005

G = Google
M = Msn Beta
T = Ask/Teoma
Y = Yahoo!

Google > 8 Bio., MSN > 5 Bio., Yahoo! > 4 Bio., Ask/Teoma > 2 Bio.
overlap statistics → *(surface) Web > 11.5 Bio. pages (> 40 TBytes)*

Deep Web (Hidden Web) estimated to have 500 Bio. units (> 10 PBytes)

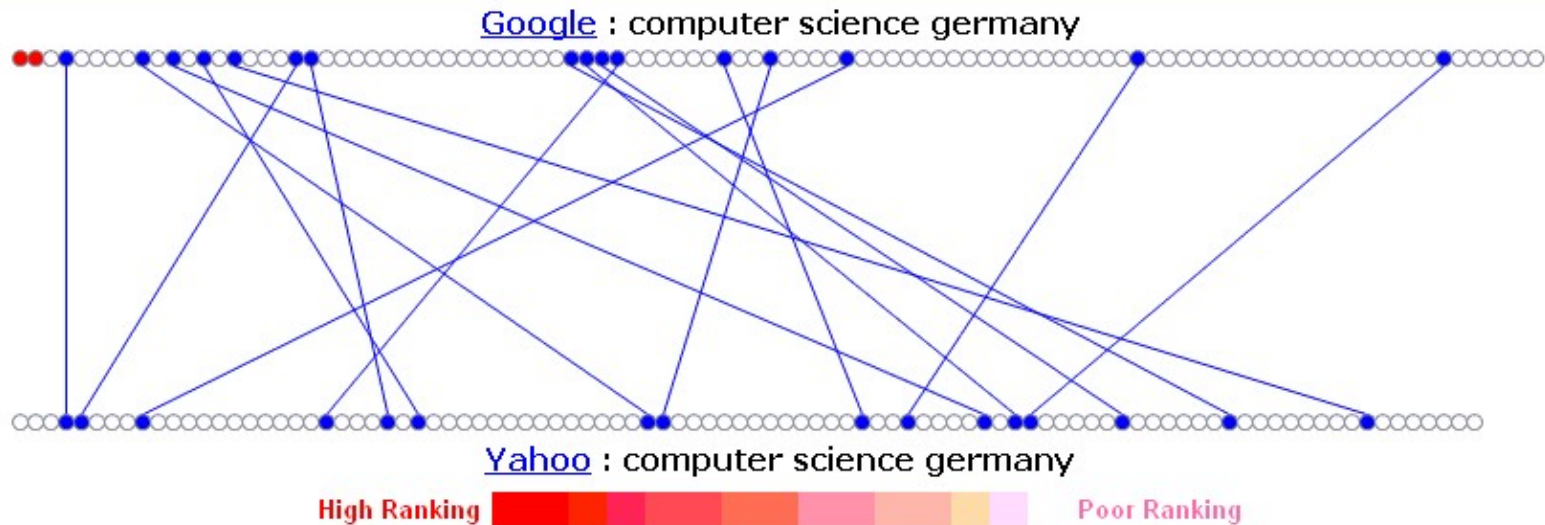
Web Coverage by Search Engines

THUMBSHOTS RANKING by Smartdevil

Search: versus

Highlight Site: (optional eg. mysite.com) *This service requires IE5+*

Due to high traffic, results may not appear. Please try again later.
Support Thumbshots Ranking by [donating or sponsoring](#) on Thumbshots.org.



Google (computer science germany)	Yahoo (computer science germany)
Overlapping Links: 16 (16 %)	Overlapping Links: 16 (17 %)
Unique Links: 84 (84 %)	Unique Links: 80 (83 %)
Total Links: 100	Total Links: 96

<http://ranking.thumbshots.com/>
<http://rankcomparison.di.unipi.it/>

Web Archiving

Enter Web Address: All [Adv. Search](#) [Compare Archive Pages](#)

Searched for <http://www.mpi-sb.mpg.de>

373 Results

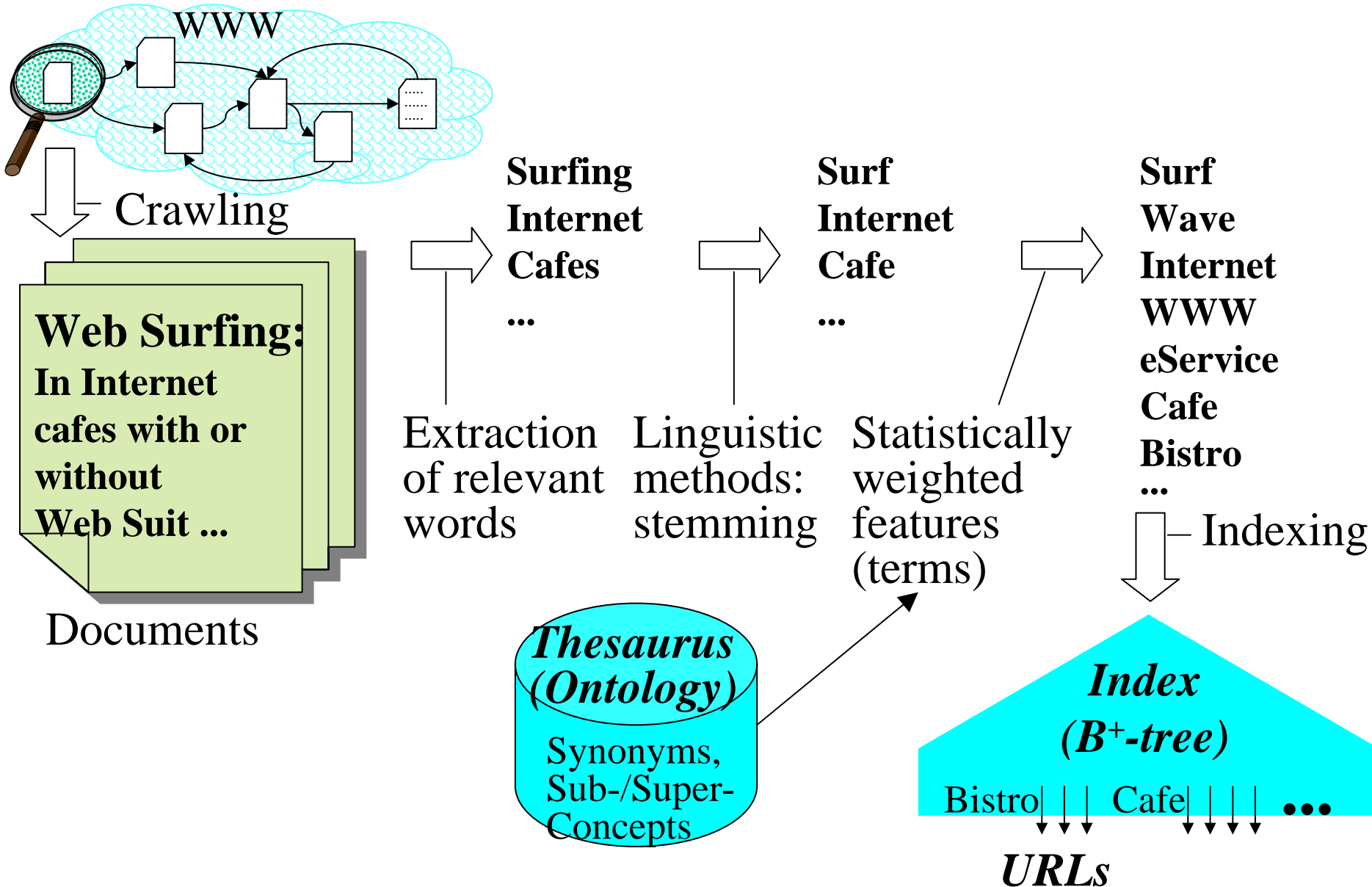
Note some duplicates are not shown. [See all.](#)
* denotes when site was updated.

Search Results for Jan 01, 1996 - Oct 10, 2005

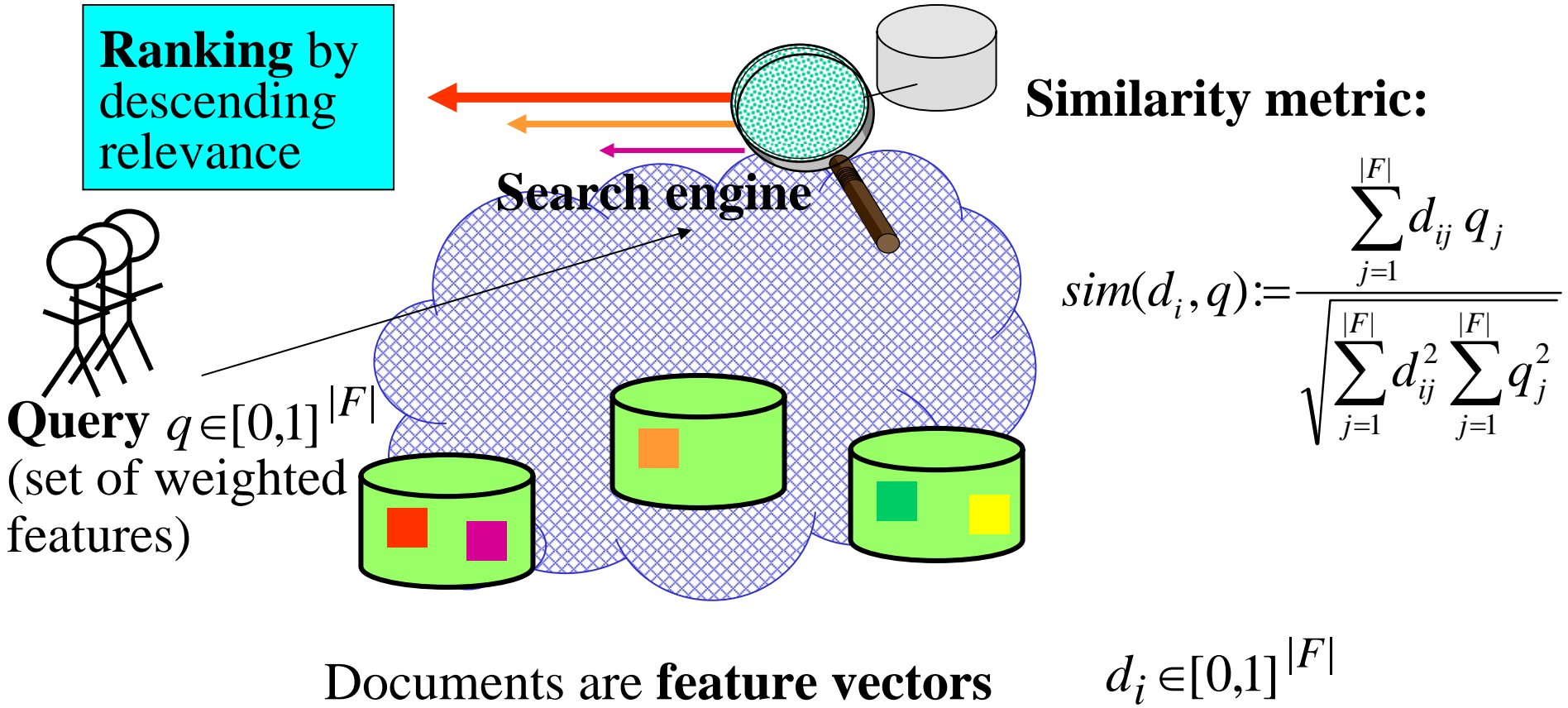
1996	1997	1998	1999	2000	2001	2002	2003	2004	2005
2 pages	5 pages	4 pages	8 pages	12 pages	39 pages	11 pages	16 pages	36 pages	0 pages
Nov 11, 1996 * Dec 27, 1996 *	Feb 17, 1997 * Feb 18, 1997 Mar 05, 1997 * Apr 28, 1997 * Aug 14, 1997 *	Jan 25, 1998 * Jul 03, 1998 * Dec 02, 1998 * Dec 12, 1998	Jan 17, 1999 Jan 25, 1999 Jan 27, 1999 Feb 03, 1999 Apr 17, 1999 * Apr 23, 1999 * Oct 03, 1999 * Nov 03, 1999 *	Mar 04, 2000 * Apr 08, 2000 May 11, 2000 * May 19, 2000 May 20, 2000 Jun 19, 2000 * Jun 21, 2000 Aug 17, 2000 * Oct 18, 2000 * Oct 19, 2000 Oct 22, 2000 Dec 04, 2000 *	Feb 02, 2001 * Feb 26, 2001 * Mar 01, 2001 * Mar 02, 2001 Mar 09, 2001 Mar 31, 2001 Apr 03, 2001 Apr 04, 2001 Apr 05, 2001 Apr 06, 2001 Apr 07, 2001 Apr 10, 2001 Apr 11, 2001 Apr 12, 2001 Apr 13, 2001 Apr 14, 2001 Apr 17, 2001 Apr 18, 2001 Apr 19, 2001 Apr 20, 2001 Apr 21, 2001 Apr 22, 2001 Apr 23, 2001	Feb 10, 2002 * May 29, 2002 * May 30, 2002 Jun 01, 2002 Jul 22, 2002 * Aug 02, 2002 Sep 28, 2002 * Oct 13, 2002 Nov 26, 2002 * Nov 28, 2002 Dec 04, 2002	Feb 01, 2003 * Feb 03, 2003 Feb 28, 2003 * Mar 27, 2003 * Apr 19, 2003 * Apr 22, 2003 Apr 24, 2003 May 26, 2003 * Jun 11, 2003 Jul 29, 2003 * Aug 08, 2003 Sep 30, 2003 * Oct 26, 2003 Dec 05, 2003 * Dec 13, 2003 Dec 21, 2003	Feb 01, 2004 * Apr 02, 2004 * May 11, 2004 May 22, 2004 * May 25, 2004 Jun 06, 2004 * Jun 14, 2004 * Jun 15, 2004 Jun 16, 2004 * Jun 18, 2004 Jun 24, 2004 Jun 26, 2004 Jun 28, 2004 Jul 03, 2004 Jul 11, 2004 Jul 15, 2004 Jul 16, 2004 Jul 18, 2004 Jul 25, 2004 Aug 11, 2004 Aug 13, 2004 * Sep 21, 2004 * Sep 29, 2004 *	

40 Billion URLs archived every 2 months since 1996 → 500 TBytes
<http://www.archive.org>

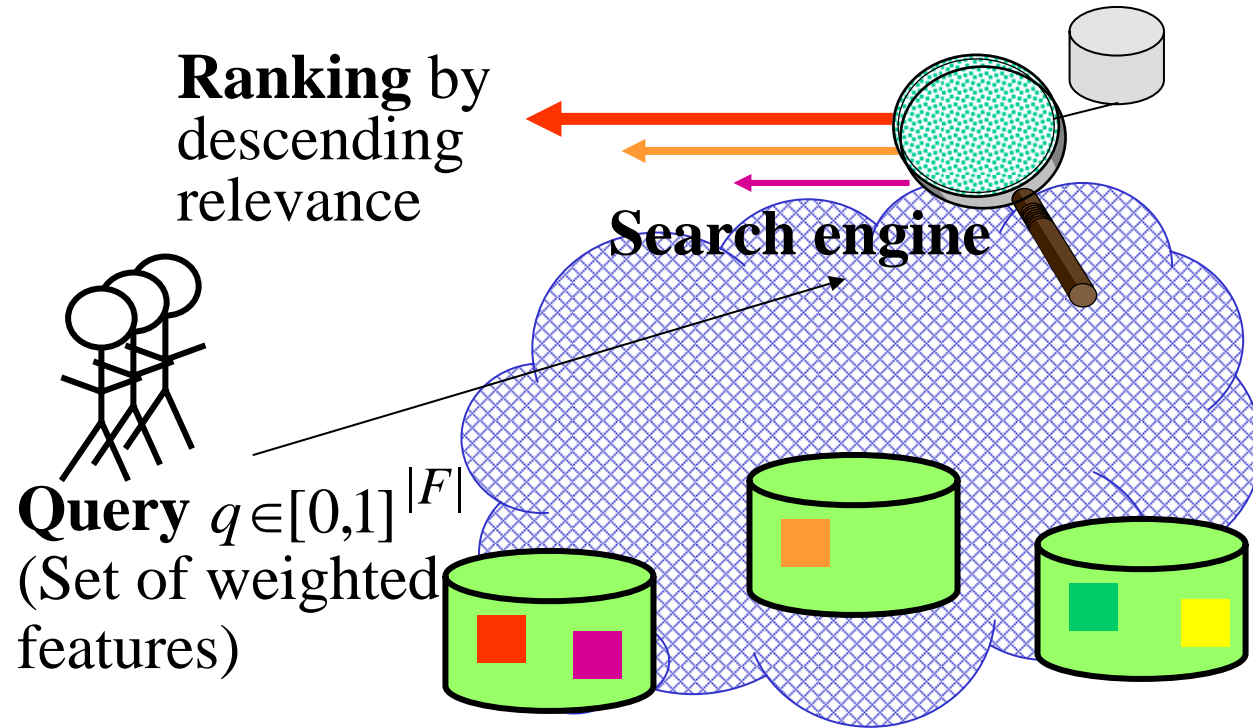
Web Content Gathering and Indexing



Vector Space Model for Content Relevance Ranking



Vector Space Model for Content Relevance Ranking



Similarity metric:

$$sim(d_i, q) := \frac{\sum_{j=1}^{|F|} d_{ij} q_j}{\sqrt{\sum_{j=1}^{|F|} d_{ij}^2 \sum_{j=1}^{|F|} q_j^2}}$$

e.g., using:

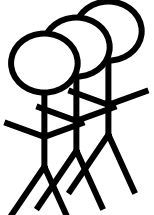
$$d_{ij} := w_{ij} / \sqrt{\sum_k w_{ik}^2}$$

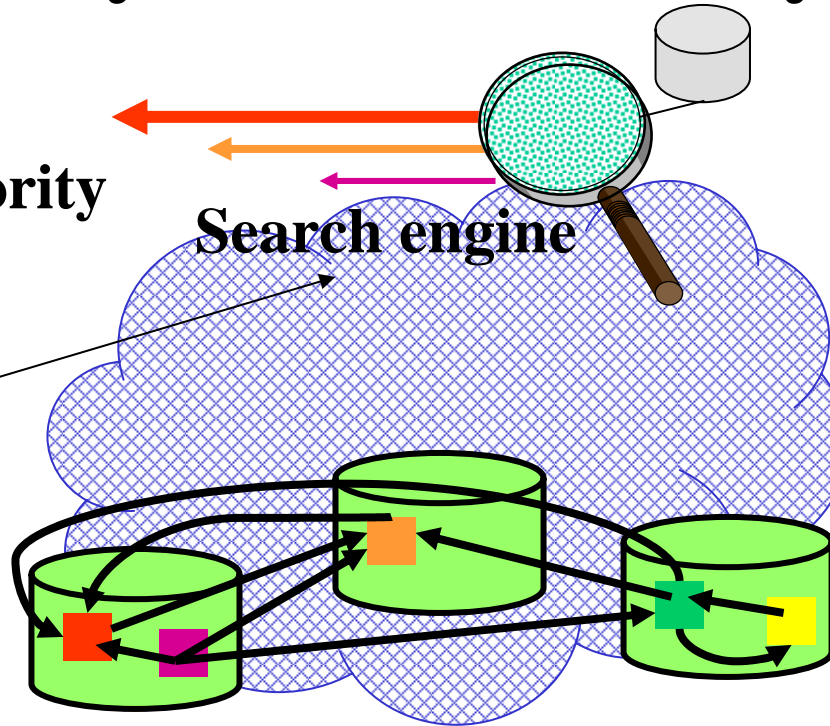
$$w_{ij} := \frac{freq(f_j, d_i)}{\max_k freq(f_k, d_i)} \log \frac{\# docs}{\# docs \text{ with } f_i}$$

**tf*idf
formula**

Link Analysis for Authority Ranking

Ranking by
descending
relevance & authority


Query $q \in [0,1]^{|F|}$
(Set of weighted
features)



+ Consider in-degree and out-degree of Web nodes:

Authority Rank (d_i) :=

Stationary visit probability [d_i]

in random walk on the Web

Reconciliation of relevance and authority by ad hoc weighting

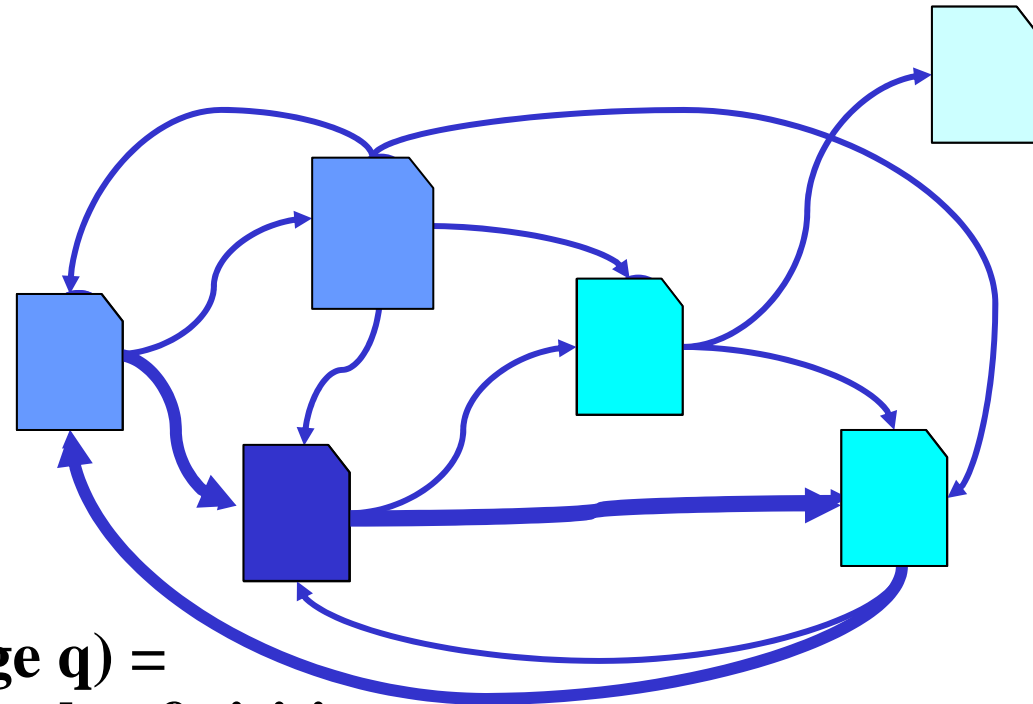
Google's PageRank in a Nutshell

random walk on the Web graph:

uniformly random choice of **links** + random jumps

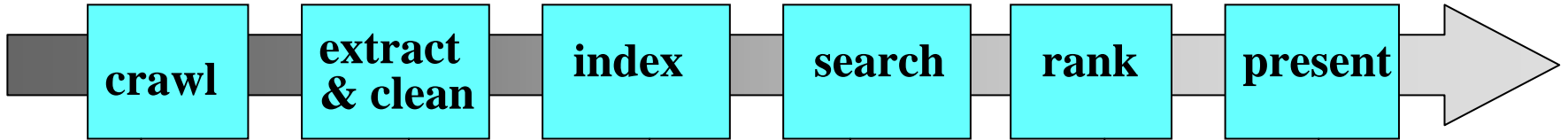
$$PR(q) = \varepsilon \cdot j(q) + (1 - \varepsilon) \cdot$$

$$\sum_{p \in IN(q)} PR(p) \cdot t(p, q)$$



**Authority (page q) =
stationary prob. of visiting q**

System Architecture of a Web Search Engine



strategies for crawl schedule and priority queue for crawl frontier

handle dynamic pages, detect duplicates, detect spam

build and analyze Web graph, index all tokens or word stems

fast top-k queries, query logging and auto-completion

scoring function over many data and context criteria

GUI, user guidance, personalization

special file system for high-performance storage management

index-entry and query-result caching for fast search

server farm with > 10 000 nodes, distributed/replicated data

Search Engine Users

Zeitgeist This Month

Popular Celebrities August 2005	Music-Related Queries August 2005	Popular References August 2005	Travel-Related Queries August 2005
1. madonna	1. lyrics	1. dictionary	1. expedia
2. jessica simpson	2. my chemical romance	2. maps	2. travelocity
3. pamela anderson	3. beyonce	3. weather	3. orbitz
4. paris hilton	4. mariah carey	4. white pages	4. southwest airlines
5. jessica alba	5. green day	5. yellow pages	5. american airlines

Google News Queries

Katrina-Related Queries August 2005	Popular Sports Queries August 2005	Popular Newsmakers August 2005
1. hurricane katrina	1. real madrid	1. natalee holloway
2. new orleans	2. arsenal	2. cindy sheehan
3. hurricane katrina photos	3. cricket	3. peter jennings
4. slidell	4. nhl	4. lance armstrong
5. french quarter	5. nba	5. tiger woods

Frequent Web queries:

<http://www.google.com/press/zeitgeist.html>

Web-Search Usage Patterns

classification of queries [Rose/Levinson: WWW 2004]:

- **navigational**: find specific homepage with unknown URL, e.g. Cirrus Airlines
- **informational**: learn about topic
 - focused, e.g. Chernoff bounds, soccer world championship qualification
 - unfocused, e.g. undergraduate statistics, dark matter, Internet spam
 - seeking advice, e.g. help losing weight, low-fat food, marathon training tips
 - locating service, e.g. 6M pixel digital camera, taxi service Saarbrücken
 - exhaustive, e.g. Dutch universities, hotel reviews Crete, MP3 players
- **transactional**: find specific resource, e.g. download Lucene source code, Sony Cybershot DSC-W5, Mars surface images, hotel beach south Crete August
- **embedded in business workflow** (e.g. CRM, business intelligence) or **personal agent** (in cell phone, MP3 player, or ambient intelligence at home) **with automatically generated queries**
- **natural-language question answering (QA)**:
 - **factoids**, e.g. when was Johnny Depp born, where is the Louvre, who is the CEO of Google, what kind of particles are quarks, etc.
 - **list queries**, e.g. in which movies did Johnny Depp play

Search Result Organization (1)



[company](#) | [products](#) | [solutions](#) | [customers](#) | [demos](#) | [press](#)

[Advanced Search](#)

[Help](#)

NEW search the [Wikipedia](#) at [Clusty.com](#)

Clustered Results

- ▶ [java lava](#) (115)
- ⊕ ▶ [Stone](#) (14)
- ⊕ ▶ [Coffee](#) (12)
- ⊕ ▶ [Volcanic](#) (9)
- ⊕ ▶ [Design](#) (9)
- ▶ [Lava Lamp](#) (6)
- ⊕ ▶ [Spa, Treatment](#) (8)
- ⊕ ▶ [Java Lava Trading](#) (5)
- ▶ [College, Pierce](#) (6)
- ▶ [Kona, Lava Java Kailua](#) (4)
- ▶ [Cafe](#) (4)
- ▼ [More](#)

Find in clusters:



Top 115 results of at least 272,100 retrieved for the query **java lava** ([Details](#))

1. [Lava](#) [new window] [frame] [cache] [preview] [clusters]
Lava Lava has been discontinued and is no longer sold or supported. The closest alternate is Lava3 Core . You may enter the archive to find what you are looking for, but please anticipate that some of ...
[sharkysoft.com/software/java/lava](#) - Lycos 2, Ask Jeeves 2, Open Directory 3, MSN Search 12
2. [Lava Java- coffee house and Bistro in Charlotte, NC, Coffee shop ...](#) [new window] [frame] [cache] [preview] [clusters]
... house & bistro,a cafe restaurant & bistro serving Charlotte,NC,Coffee shop Links 704-567-4577 **Lava Java** Coffee shop & bistro in Charlotte, NC is run ...
[www.global-espresso.com](#) - Open Directory 2, Wisenut 4, MSN Search 24, Ask Jeeves 37
3. [Java Lava Trading Company - Steelco Industries, Inc.](#) [new window] [frame] [cache] [preview] [clusters]
Costa Ricans have been enjoying Imperial beer since 1924. Imperial is created from a balanced formula combining malts, grains and hops, without a pronounced overtone in its taste. It is what is known ...
[www.javalavatrading.com](#) - Wisenut 1, MSN Search 4, Ask Jeeves 30
4. [Java Lava PT](#) [new window] [frame] [preview] [clusters]
Lava stone products including architectural, garden and decorative tiles and products.
[www.javalava.biz](#) - Open Directory 1, Lycos 3, Ask Jeeves 3
5. [LAVA LAMP](#) [new window] [frame] [cache] [preview] [clusters]
How I Created the **Lava** Lamp My **Lava** Lamp is featured in: TeamJava **Java** Links There were at least three steps in creating this animation applet. Create the animation frames. I was inspired by the wonderful ...
[smc.vnet.net/javalamp.html](#) - Lycos 1, Ask Jeeves 1, Wisenut 7, MSN Search 8
6. [Java Lava Stone: Stone Craft Collection](#) [new window] [frame] [cache] [preview] [clusters]
java lava stone, Java Lava Stone: Stone Craft Collection, Java Lava Stone, Company Profile Product

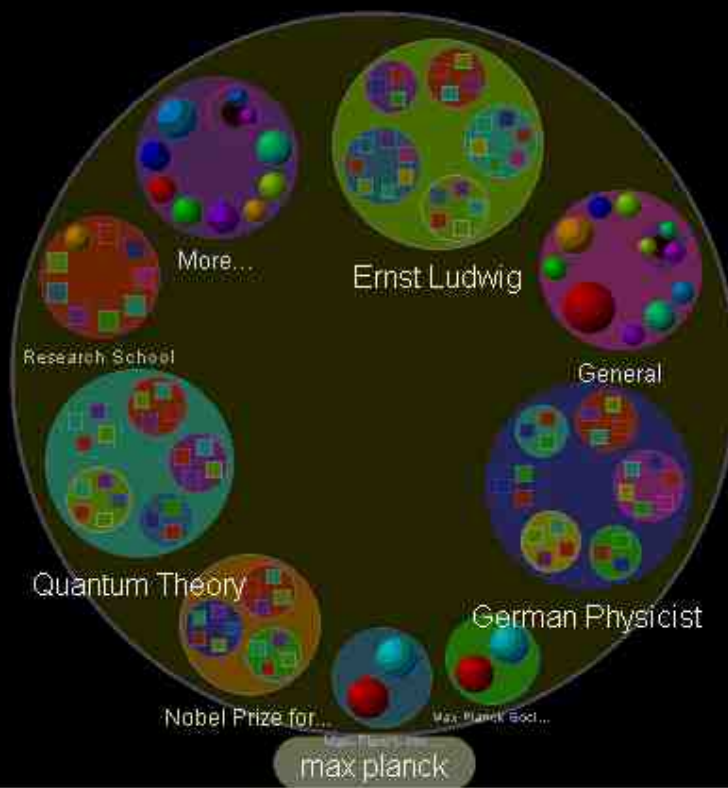
cluster search results into topic areas

<http://www.vivisimo.com/>

Search Result Organization (3)

I grok DEMO max planck GROK EMAIL YOUR GROKKER MAP! | HELP powered by YAHOO! SEARCH

← BACK | TOP



max planck

SHOW TOOLS

SPONSOR RESULTS

[MAX Conferencing System - ClearOne](#)
The first to launch the wireless conferencing phone. Products include MAX EX, MA...
www.clearone.com

[Max - Dragon Tales Party Supplies](#)
Party Poopers stocks a large inventory of party supplies, including tableware, i...
www.partypoopers.com

[The Max Digital Magnifier for TVs](#)
The Max Digital Magnifie is a powerful, portable, hand-held digital magnifier d...
www.firststreetonline.com

visualize cluster hierarchies for search results

<http://www.grokker.com/>

Search Result Organization (4)



max planck

[Search tips](#)

Search

- [Advanced Search](#)
- [Preferences](#)

Find this phrase

Sponsored Link

[Max Planck](#)

Neu und gebraucht Mitbieten oder Sofort-Kaufen!
www.ebay.de

Results

Relevant web pages

Showing 1-10 of about 1,259,000:

[Max Planck - Biography](#)

Max Planck – Biography **Max** Karl Ernst Ludwig **Planck** was born in Kiel, Germany, on April 23, 1858, the son of Julius Wilhelm and Emma (née...
www.nobel.se/physics/laureates/1918/planck...

[\[Related Pages\]](#)

[Biography of Max Planck](#)

Max Planck (1858-1947)

wwwchem.csustan.edu/chem3070/Raul1.htm | [Cached](#)

[Max-Planck-Gesellschaft - Website der MPG](#)

Max-Planck-Institute betreiben Grundlagenforschung in den Natur-, Bio- und Geisteswissenschaften im Dienste der Allgemeinheit. Insbesondere greift...

www.mpg.de/ | [Cached](#)

[Max Planck Society - Max-Planck-Portal](#)

Refine

Suggestions to narrow your search

[Max Planck INSTITUTE](#)

[Max Planck Research](#)

[Karl Ernst Ludwig](#)

[Corporate Governing](#)

[Science Odyssey](#)

[Homesite Helpaboutsearch](#)

[\[Show All Refinements\]](#)

Resources

Link collections from experts and enthusiasts

[Max-Planck-Gesellschaft - Sonstige Ausstattungen](#)

[www.planck.de/...](http://www.planck.de/)

[SurfWax: News, Reviews and Articles On Max Planck](#)

point out related queries
<http://www.teoma.com>

Search Result Organization (5)



[Web](#) [Images](#) [Groups](#) [News](#) [Froogle](#) [Local](#) [Desktop](#) [more »](#)

As you type, Google suggests:

max pl	
max planck	2,240,000 results
max planck institute	1,220,000 results
max plank	216,000 results
max plank institute	61,400 results
max planck institut	1,110,000 results
max plugins	665,000 results
max plus	15,000,000 results
max planck biography	35,200 results
max planck society	710,000 results
max plus ii	11,100,000 results

[Advanced Search](#)
[Preferences](#)
[Language Tools](#)

[Learn more](#)

©2005 Google

auto-complete queries

<http://labs.google.com/suggest/>
<http://www.mpi-inf.mpg.de>

Search Result Organization (6)

The screenshot displays the exalead search engine interface. At the top, the search bar contains 'max planck' and shows '275,327 results in 0.96 s'. Below the search bar, there are navigation options for 'audio' and 'video', and links for 'preferences', 'advanced search', 'feedback', and 'help'. The main content area is divided into several sections:

- RELATED TERMS:** A list of related terms including 'Max Planck Institute', 'Max-Planck-Institut', 'Max-Planck-Gesellschaft', etc.
- RELATED CATEGORIES:** Categories such as 'Science and Environment', 'Computers', 'Reference', and 'Kids and Teens'.
- WEB SITE LOCATION:** Geographic locations like 'Europe', 'North America', and 'Asia'.
- DOCUMENT TYPE:** File formats like 'PDF', 'TXT', 'DOC', and 'XLS'.
- MATCHING DOCUMENTS:** A list of search results, each with a title, a brief description, and a URL. The results include:
 - Max-Planck-Gesellschaft - Website der MPG:** Über die Max-Planck-Gesellschaft Forschungsgebiete der Max-Planck-Gesellschaft Forschungsergebnisse der Max-Planck-Gesellschaft Wissenschaftliche Ressourcen und Kooperationen [...]
 - Max-Planck-Institut für ausl. öffentliches Recht ...:** Direktoren: Prof. Dr. Armin von Bogdandy Prof. Dr. Dr. h.c. Rüdiger Wolfrum Aktuelles Über das Institut Profil, Arbeitsbereiche, Adresse und Anfahrt Forschung Mitarbeiter Bibliothek OPAC Virtuelles Institut © Max-Planck-Institut für ausländisches öffentliches Recht und Völkerrecht, Heidelberg
 - Homepage: Max-Planck-Institut für Informatik:** ... homepage Departments Location People Services Research School Max Planck Center Computer Science Cluster Sitemap Search [...]
 - Max-Planck Institut für Wissenschaftsgeschichte:** www.mpiwg-berlin.mpg.de/ - 13k - 08 Jun 2005
 - Max Planck Institut fuer Radioastronomie Bonn:** [english] Aktuell Das Institut Forschung Mitarbeiter Öffentlichkeit Intranet webmaster@mpifr-bonn.mpg.de [english]

On the right side of the interface, there are several thumbnail images of website pages.

show broader context
of search results

<http://www.exalead.com/>

Evaluation of Search Result Quality: Basic Measures

ideal measure is user satisfaction

heuristically approximated by benchmarking measures

(on test corpora with query suite and relevance assessment by experts)

Capability to return **only** relevant documents:

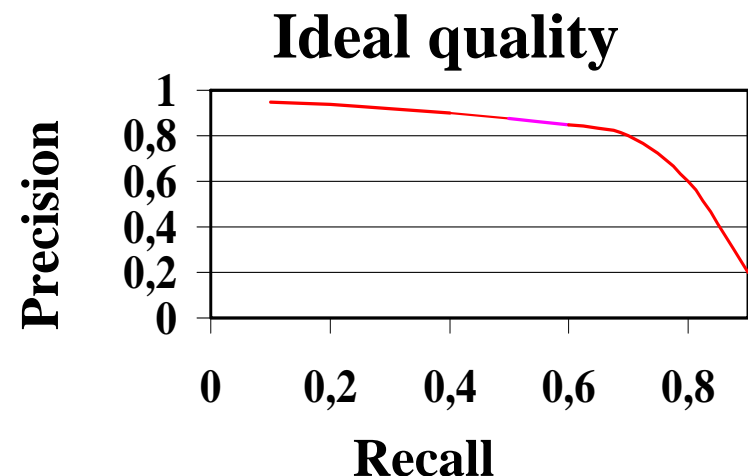
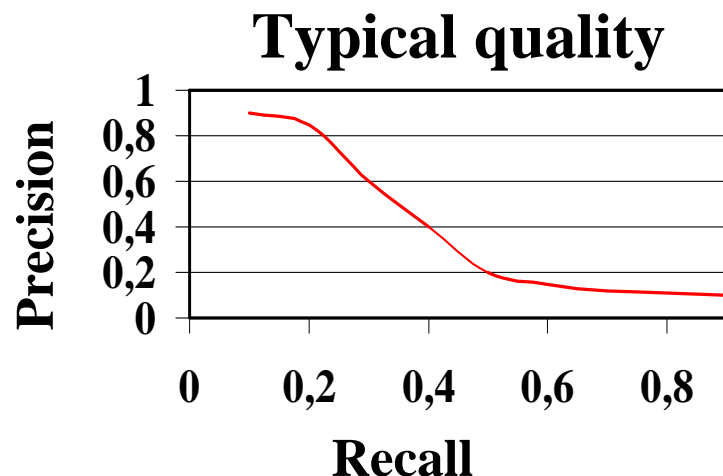
$$\textit{Precision (Prazision)} = \frac{\# \textit{ relevant docs among top } r}{r}$$

typically for
 $r = 10, 100, 1000$

Capability to return **all** relevant documents:

$$\textit{Recall (Ausbeute)} = \frac{\# \textit{ relevant docs among top } r}{\# \textit{ relevant docs}}$$

typically for
 $r = \text{corpus size}$



Evaluation of Search Result Quality: Aggregated Measures

Combining precision and recall into **F measure**

(e.g. with $\alpha=0.5$:

harmonic mean **F1**):

$$F = \frac{1}{\alpha \frac{1}{\textit{precision}} + (1 - \alpha) \frac{1}{\textit{recall}}}$$

Precision-recall breakeven point of query q :

point on precision-recall curve $p = f(r)$ with $p = r$

for a set of n queries q_1, \dots, q_n (e.g. TREC benchmark)

$$\begin{array}{l} \textit{Macro evaluation} \\ \textit{(user-oriented)} \\ \textit{of precision} \end{array} = \frac{1}{n} \sum_{i=1}^n \textit{precision}(q_i)$$

$$\begin{array}{l} \textit{Micro evaluation} \\ \textit{(system-oriented)} \\ \textit{of precision} \end{array} = \frac{\sum_{i=1}^n \# \textit{ relevant \& found docs for } q_i}{\sum_{i=1}^n \# \textit{ found docs for } q_i}$$

analogous
for recall
and F1

Evaluation of Search Result Quality: Integrated Measures

- **Interpolated average precision** of query q
with precision $p(x)$ at recall x
and step width Δ (e.g. 0.1):

$$\frac{1}{1/\Delta} \sum_{i=1}^{1/\Delta} p(i\Delta)$$

area under precision-recall curve

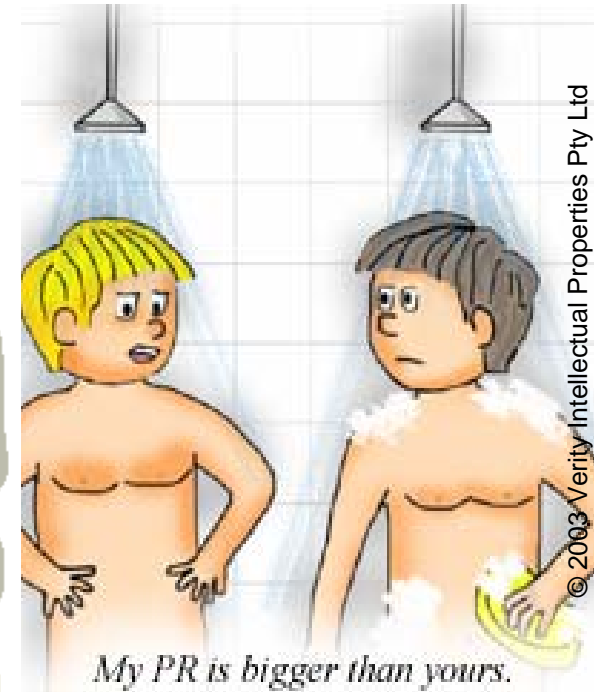
- **Uninterpolated average precision** of query q
with top- m search result rank list d_1, \dots, d_m ,
relevant results d_{i_1}, \dots, d_{i_k} ($k \leq m, i_j \leq i_{j+1} \leq m$):

$$\frac{1}{k} \sum_{j=1}^k \frac{j}{i_j}$$

- **Mean average precision (MAP)** of query benchmark suite
macro-average of per-query interpolated average precision
for top- m results (usually with recall width 0.01)

more measures in the literature

Google & Co: Where Do We Stand Today?



- ★ great for e-shopping, school kids, scientists, doctors, etc.
- ★ high-precision results for simple queries
- ★ superb scalability (now >8 Bio. docs, >1000 queries/sec)
- ★ continuously enhanced: Froogle, Google Scholar, alerts, multilingual for >100 languages, query auto-completion, etc.

What Google Can't Do

Killer queries (disregarding QA, multilingual, multimedia):

- *professors from Saarbruecken who teach DB or IR and have projects on XML*
- *drama with three women making a prophecy to a British nobleman that he will become king*
- *the woman from Paris whom I met at the PC meeting chaired by Jennifer Widom*
- *best and latest insights on percolation theory*
- *pros and cons of dark energy hypothesis*
- *market impact of XML standards in 2002 vs. 2004*
- *experienced NLP experts who may be recruited for IT staff*

The following words are very common and were not included in your search: with a. [details]

Web Results 1 - 10 of about 85,800 for **drama with three women making a prophecy** . 0.29 seconds)

Quick Tips for Meeting, Dating, and Attracting Women.

... Don't Lose Your Nerve. The 3-Day Test. Do You Have a Phone? ... The Pity Kiss. An Obvious Talent - Palm Reading. Learn to Play Golf. ... Don't Pursue **Women**. Attract Them. ...
www.sosuave.com/quick/default.htm - 49k - [Cached](#) - [Similar pages](#)

An eyewitness to Shakespeare's plays

... of his death (though he may have committed suicide to make his **prophecy** ... More than a play. ... riding through a wood, there stood before them **three women** fairies or ...
ise.uvic.ca/Library/SLTnoframes/life/forman.html - 6k - [Cached](#) - [Similar pages](#)

The Guide to World Drama - Plays L

... Hammersmith in 1989, this is a careful examination of the role of pornography in our society and the way it affects **three** young **women** in ... Play. ... 4 men, 3 **women**. ...
www.4-wall.com/plays/plays_l/lovelyhappy.htm - 12k - [Cached](#) - [Similar pages](#)

CliffsNotes::Oedipus Trilogy - The Oedipus Trilogy: Study Help ...

... 3. In Antigone, who is the real main character ... and Ismene in their views of **women** in society. ... the following statement: Antigone is primarily a **drama** of politics ...



Web Images Groups ^{New!} News Froogle more »

drama three women prophecy british nobleman | Search Advanced Search Preferences

Web Results 1 - 10 of about 647 for **drama three women prophecy british nobleman king.** (0.22 seconds)

A survey course in British literature

... plays: Cardenio, Two Noble Kinsmen, Sir Thomas More. **Three** pages in ...
Poems: Venus and Adonis, 1592-3; The Rape of Lucrece ... An historic **drama**
(jointly with Coleridge ...

www.unibuc.ro/eBooks/filologie/tupan/indexofauthors.htm - 84k -
[Cached](#) - [Similar pages](#)

Sponsored Links

British Drama
Research **British drama** at
the world's largest online library.
www.questia.com

"KING ARTHUR" ON THE STAGE

... of Richard II., and Margaret's curse in Richard III (i., 3). There is ... and a vision of
the mystic barge and the **three** queens. ... A **Drama** in a Prologue and Four Acts ...

www.lib.rochester.edu/camelot/carrbond.htm - 54k - [Cached](#) - [Similar pages](#)

Metroactive Movies | Reviews La-Lm

... Liberty: 3 Stories About Life and Death Worthy of ... The comedy/**drama** may be a
vehicle for comedians ... adorable son (Giorgio Cantarini) when the **three** are
deported ...

www.metroactive.com/movies/capsule-la.html - 101k - [Cached](#) - [Similar pages](#)

ERIC STOLTZ VIDEOS & DVD'S AT HOLLYWOOD TEEN MOVIES

... sides of a romantic triangle between **three** best friends ... a true story, this powerful



Web Images Groups ^{New!} News Froogle more »

drama woman prophecy scottish nobleman Search Advanced Search Preferences

Web Results 1 - 10 of about 456 for **drama woman prophecy scottish nobleman** (0.25 seconds)

Macbeth

... Lady Macbeth hears of the witches' **prophecy**, Duncan's ... Lady Macbeth and the three witches are extremely wicked ... Macbeth is the focus of the **drama's** moral ...

www.sparknotes.com/shakespeare/macbeth/section1.html - 38k - [Cached](#) - [Similar pages](#)

Macbeth

... to Orlando theater, is a strongly sensual Lady Macbeth, whose ... eyes at the witches' first apparently ridiculous **prophecy**. ... You see Macbeth as play-actor, the man ...

www.shakespearefest.org/macbeth_99.htm - 27k - [Cached](#) - [Similar pages](#)

[PDF] Macbeth

File Format: PDF/Adobe Acrobat - [View as HTML](#)

... the new titles and appears afraid of the **prophecy**. ... womb and so "not born of woman." Macbeth conquers ... different treatments in historical source than in **drama**. ...

www.openstage.com/productions/macbethguid.pdf - [Similar pages](#)

[PDF] Macbeth Study Guide 2004.indd

File Format: PDF/Adobe Acrobat - [View as HTML](#)

... fate and that they truly **prophecy**, leaving Macbeth ... son about her husband being gone: Lady Macduff: "How ... speech used by common people in Shakespearean **drama**. ...

[Amazon.com: Books: Macbeth \(Dover Thrift Editions\) \[UNABRIDGED\]](#)

... witches's prophecies are deceptively clear: no man born of **woman** may harm ... Thus, the nature of **prophecy** becomes an integral part of the play's dynamic. ...

www.amazon.com/exec/obidos/tg/detail/-/0486278026?v=glance&st=* - 90k - [Cached](#) - [Similar pages](#)

[\[PDF\] Microsoft PowerPoint - weikum-er2004](#)

File Format: PDF/Adobe Acrobat - [view as HTML](#)

... (on Web / Deep Web / Intranet / Personal Info) Which **drama** has a scene in which a **woman** makes a **prophecy** to a **Scottish nobleman** that he will become king? ...

www.cs.fudan.edu.cn/er2004/news/ppt/Towards%20a%20Statistically%20Semantic%20Web.pdf - [Similar pages](#)

[\[PDF\] The Web of the Future](#)

File Format: PDF/Adobe Acrobat - [view as HTML](#)

... large deviation theory? Which **drama** has a scene in which a **woman** makes a **prophecy** to a **Scottish nobleman** that he will become king? ...

depend.cs.uni-sb.de/fileadmin/user_upload/depend/teaching/WS04/pers/weikum.pdf - [Similar pages](#)

[\[PPT\] The Web of the Future](#)

File Format: Microsoft Powerpoint 97 - [view as HTML](#)

... Which **drama** has a scene in which a **woman** makes a **prophecy**. to a **Scottish nobleman** that he will become king? Which professors from Saarbruecken (SB). ...

depend.cs.uni-sb.de/fileadmin/user_upload/depend/teaching/WS04/pers/weikum.ppt - [Similar pages](#)

[Character Directory](#)

... none of **woman** born / Shall harm Macbeth' (4.1 ... The tensions of the play tighten with this episode, the first of Macbeth's rise, in the Witches' prophecy of 1.3

1.3 Towards Semantic Search Engines

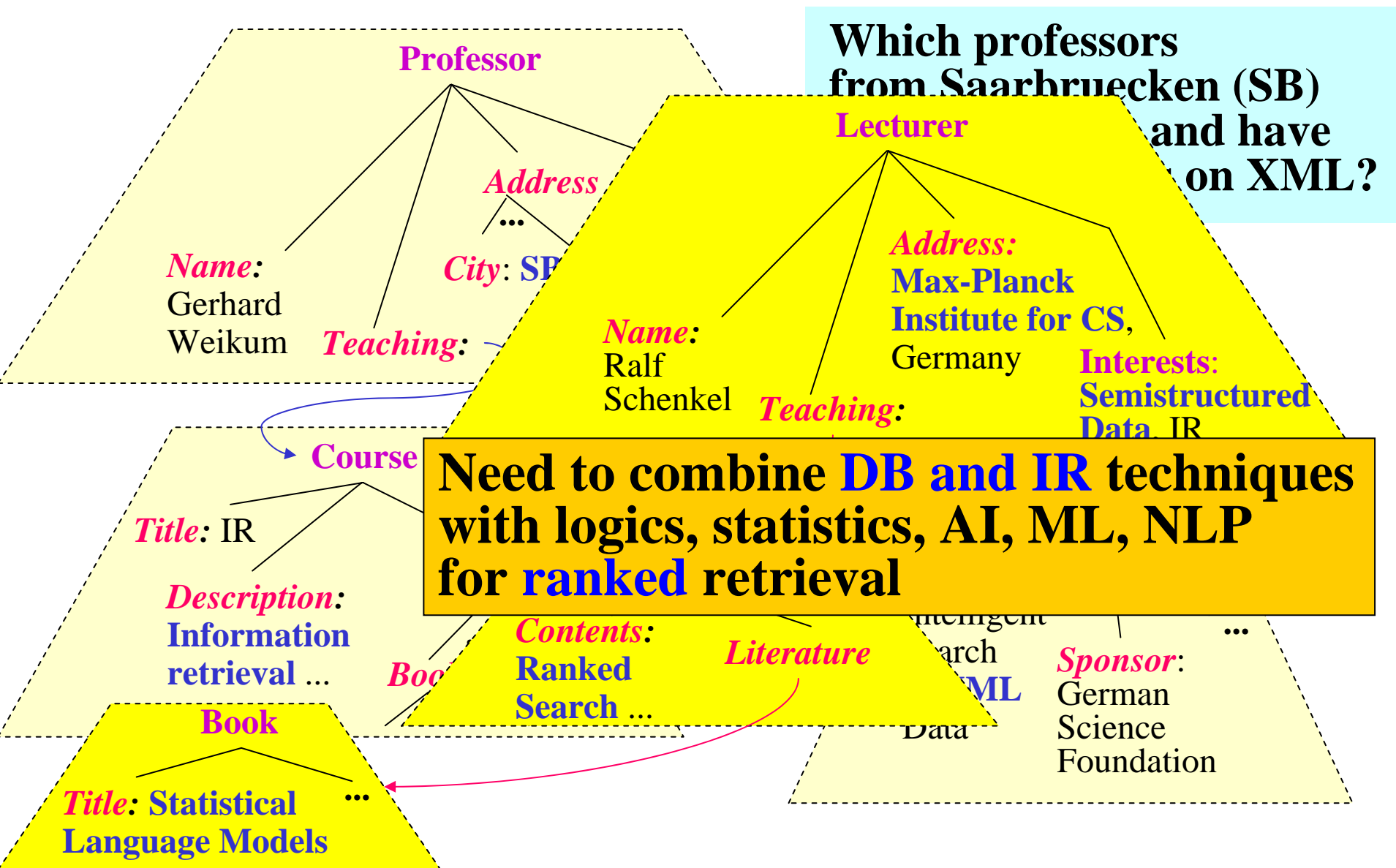
better search result quality (with low human effort) needs:

- **richer content representation** and
- **context awareness** (user profile, place&time, info authors, etc.)

strategic research directions:

- ***background knowledge***
 - ontologies & thesauri, statistics
- ***(semi-)structured and „semantic“ data***
 - metadata, XML, info extraction, annotation & classification
- ***personalization***
 - geo & time, user behavior
- ***humans in the loop***
 - collaboration, recommendation, P2P networks

Rich Content Representation in XML



Which professors from Saarbruecken (SB) and have on XML?

„Semantic Search“ with TopX Engine

User query: $\sim c = \sim t1 \dots \sim tm$

Example:

$\sim professor$ and ($\sim course = „\sim IR“$)

// $professor$ [// $place = „SB“$]// $course = „IR“$

Thesaurus/Ontology:

concepts, relationships, glosses
from WordNet, Gazetteers,
Web forms & tables, Wikipedia

Term2Concept with WSD

Query expansion

$\exp(ti) = \{w \mid \text{sim}(ti, w) \geq \theta\}$

Weighted expanded query

Example:

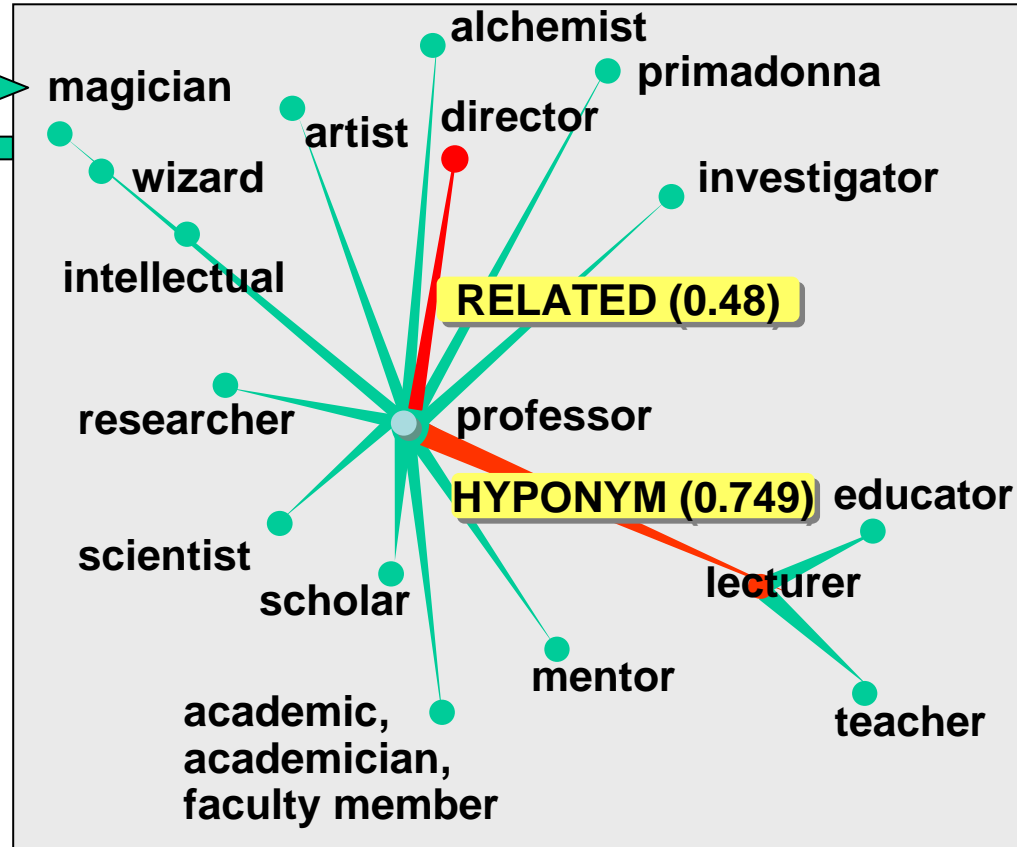
($professor$ $lecturer$ (0.749) $scholar$ (0.71) ...)

and (($course$ $class$ (1.0) $seminar$ (0.84) ...)

= („ $IR“$ „ $Web search“$ (0.653) ...))

Efficient top-k search
with dynamic expansion

better recall, better mean
precision for hard queries



Towards a Statistically Semantic Web

Isaac Newton

From Wikipedia, the free encyclopedia.

<Person>

Sir Isaac Newton (25 December 1642 – 20 March 1727 by the Julian calendar in use in England at the time; or 4 January 1643 – 31 March 1727 by the Gregorian calendar) was an English physicist, mathematician, astronomer, philosopher, and alchemist; who wrote the *Philosophiæ Naturalis Principia Mathematica* (published 5 July 1687)¹, where he described **universal gravitation** and, via his laws of motion, laid the groundwork for classical mechanics. Newton also shares credit with **Gottfried Wilhelm Leibniz** for the development of differential calculus. However, their work was not a collaboration; they developed calculus separately but nearly contemporaneously.

<Person>

Information extraction yields:
(via reg. expr., lexicon, HMM, MRF, etc.)

Person	TimePeriod	...
Sir Isaac Newton ... Leibniz ... Kneller	4 Jan 1643 - ...	

Publication	Topic
Philosophiæ Naturalis	... gravitation

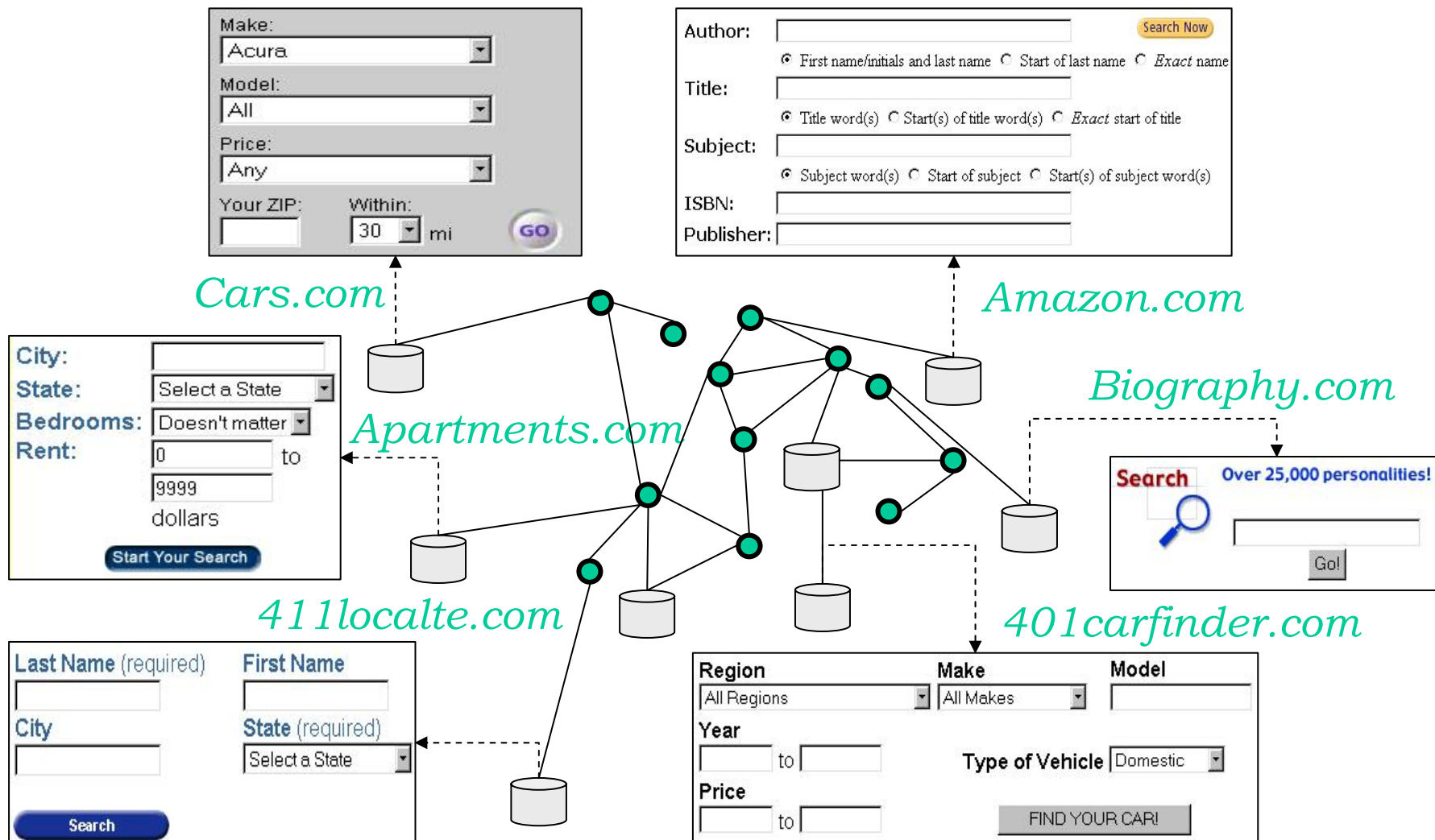
Author	Publication
... Newton	Philosophia ...

Scientist
Sir Isaac Newton ... Leibniz

but with confidence < 1

- Semantic-Web database with uncertainty !
- ranked XML/DB retrieval !

1.4 Deep Web Search



Source: Kevin Chen-Chuan Chang, CIDR 2005

Deep Web Sources

Data accessible only through query interfaces
(HTML forms, WSDL web services)

Study by B. He, M. Patel, Z. Zhang, K. Chang, CACM 2006:
> 300 000 sites with > 450 000 databases and > 1 200 000 interfaces
coverage in directories (e.g. dmoz.org) is < 15%,
total data volume estimated 10-100 PBytes

Examples of Deep Web sources:

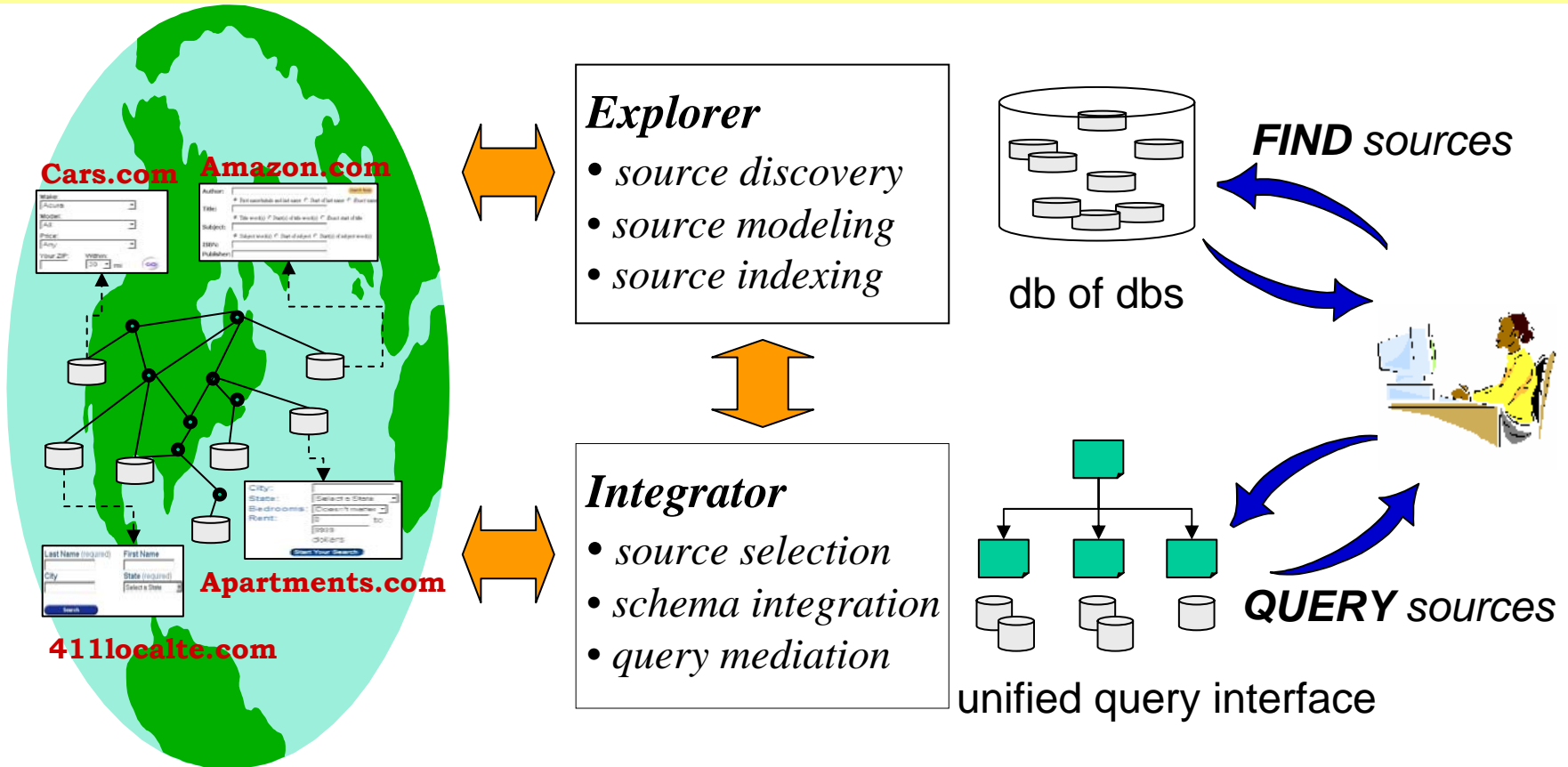
e-business and entertainment: amazon.com, ebay.com, realtor.com, cars.com,
imdb.com, reviews-zdnet.com, epinions.com

news, libraries, society: cnn.com, yahoo.com, spiegel.de, deutschland.de,
uspto.gov, loc.gov, dip.bundestag.de, destatis.de, ddb.de, bnf.fr, kb.nl, kb.se,
weatherimages.org, TerraServer.com, lonelyplanet.com

e-science: NCBI, SRS, SwissProt, PubMed, SkyServer, GriPhyN

Deep Web Research Issues

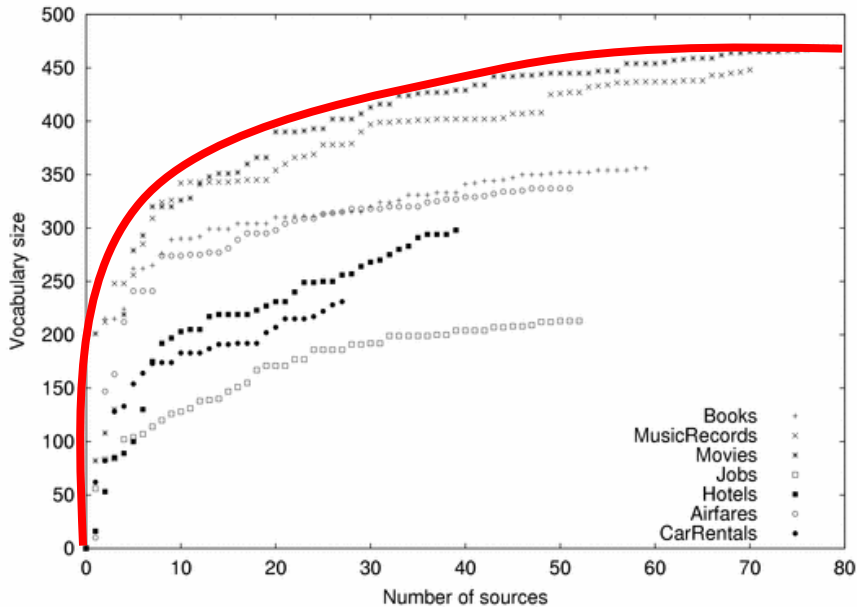
- find important/relevant sources in the Deep Web (aka. Hidden Web)
- map user queries onto source-specific interfaces
(metasearch engines are a simple special case)
- merge results into global ranking



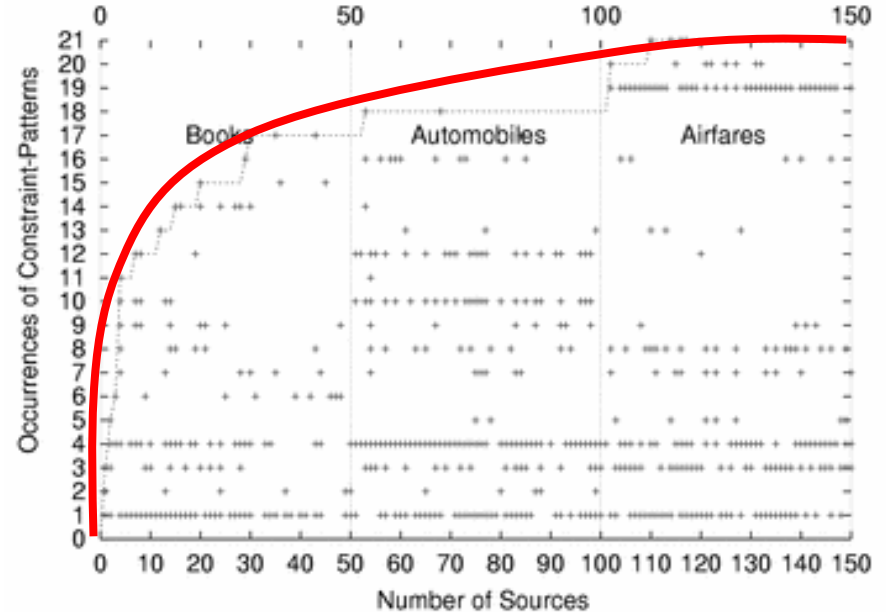
Source: Kevin Chen-Chuan Chang, CIDR 2005

Statistics on Diversity of Names

Attributes converge
in a domain!



Condition patterns converge
even across domains!



Source: Kevin Chen-Chuan Chang, CIDR 2005

→ can use statistical learning to derive mappings
among different sources within same domain

Query-to-Source Mapping

attribute *operator* *value*

Author:

First name, initials and last name Start of last name Exact name

Title:

Title word(s) Start(s) of title word(s) Exact start of title

Subject:

Subject word(s) Start of subject Start(s) of subject word(s)

ISBN:

Publisher:

Author: Last Name:
First Name:

Title:

Subject:

ISBN: ISBN:

Publisher: Category:

Artist:

Title:

Label:

Media:

Format: CD Cassette DVD Audio Vinyl

Album: Exact Phrase

Used only:

Source: Kevin Chen-Chuan Chang, CIDR 2005

Deep Web search with MetaQuerier
<http://metaquerier.cs.uiuc.edu/formext/>

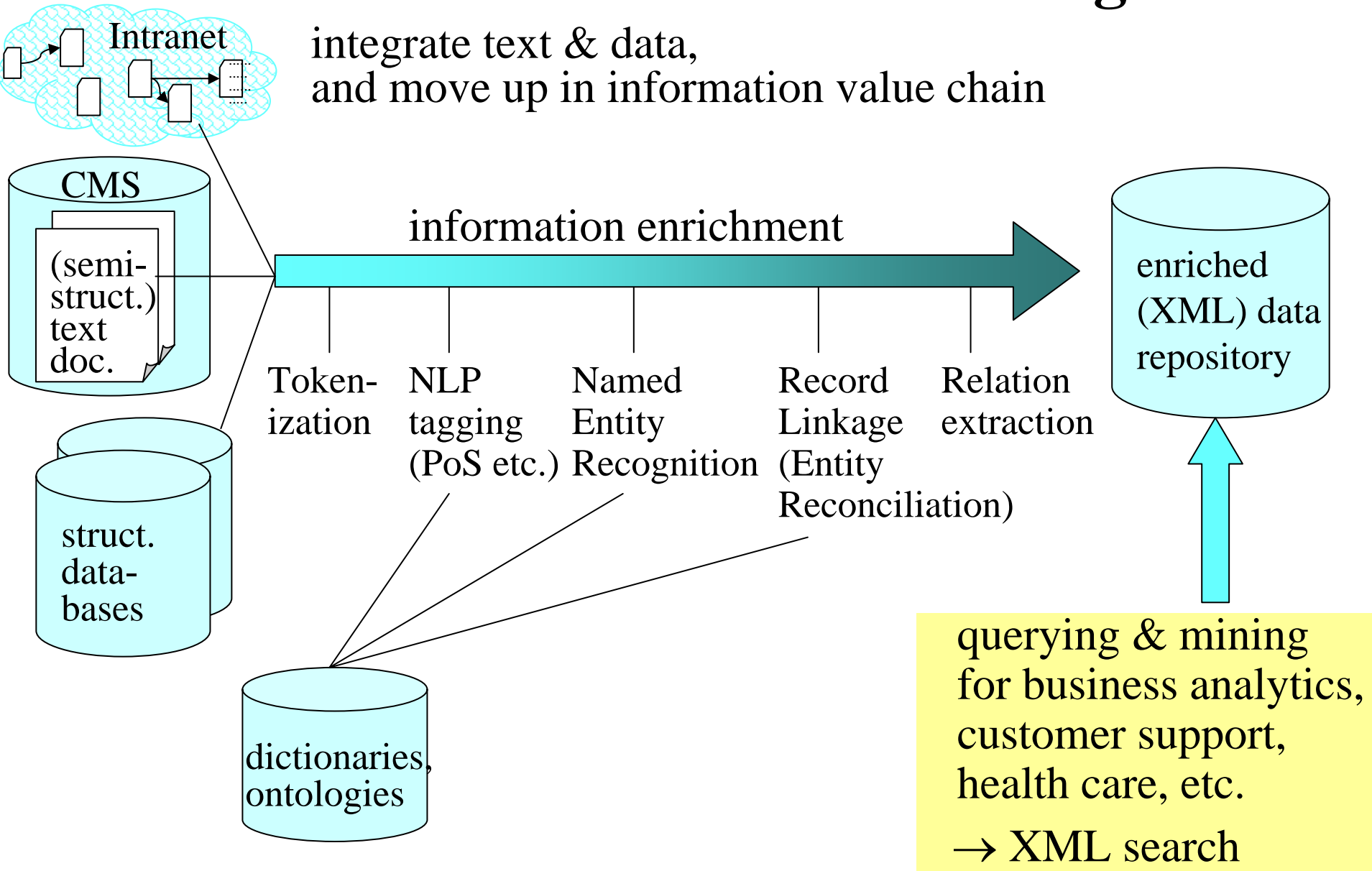
1.5 Intranet and Enterprise Search

Important differences to Web search:

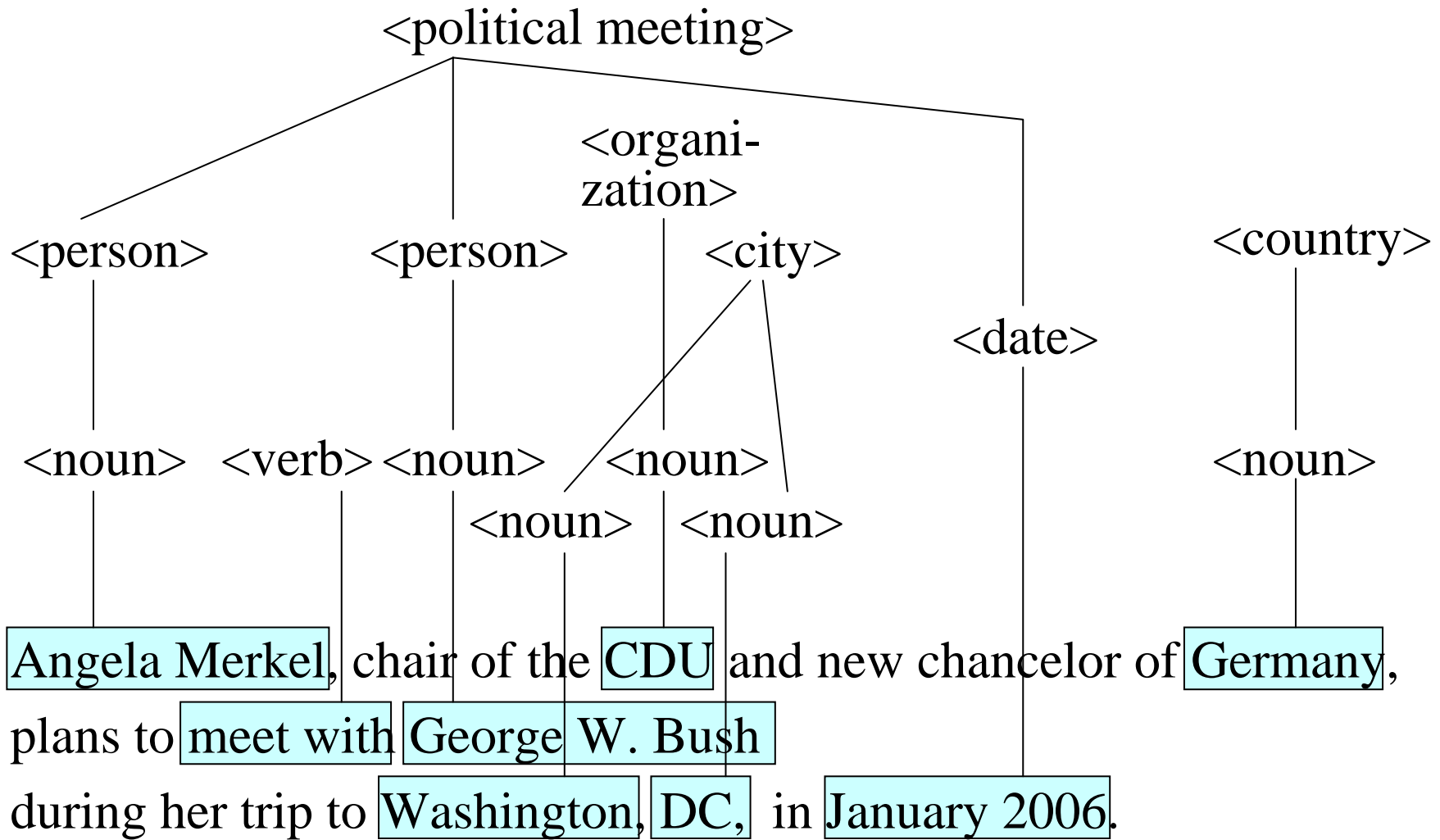
- more professional query topics, more directed queries
- higher user skills, but also higher expectations
- fewer spam and low-quality pages
- fewer hyperlinks (endorsements)
- more meaningful metadata (e.g. author and date of document)
- more context knowledge about user and data
(e.g. organizational units, location)
- opportunity for user (group) profiling

non-difference: intranet/enterprise data can also be huge !

System Architecture for Unstructured and Semistructured Information Management



Information Enrichment Example



Information Extraction from Web Pages

ANNIE Output for http://en.wikipedia.org/wiki/Che_Guevarra

Annotation Key:

Person **Location** **Organization** **Date** **Address** **Money** **Percent**

>>/**/>/**/

Che Guevara

From Wikipedia, the free encyclopedia.

(Redirected from **Che Guevarra**)

Jump to: [navigation](#) , [search](#)



Che Guevara ☐

Ernesto Rafael Guevara de la Serna (**June 14, 1928** ^[1]? **October 9, 1967**), commonly known as **Che Guevara** or **el Che**, was an **Argentine** -born **Marxist revolutionary** and **Cuban guerrilla** leader. **Guevara** was a member of **Fidel Castro** 's " **26th of July Movement** " that seized power in **Cuba** in **1959** . After serving in various important posts in the new government, **Guevara** left **Cuba** in **1965** with the hope of fomenting revolutions in other countries, first in the Congo-Kinshasa (currently the **Democratic Republic of the Congo**) and later in **Bolivia** , where he was captured in a **CIA** -organized military operation. It is believed by some that the **CIA** wished to keep **Guevara** alive for **interrogation** but, after his capture in the Yuro ravine, he died at the hands of the **Bolivian Army** in **La Higuera** near **Vallegrande** on **October 9, 1967** . Testimony by various individuals who were participants in, or

Leading open-source tool: GATE/ANNIE

<http://www.gate.ac.uk/annie/>

1.6 Personalized Search and PIM

Personalization:

- query interpretation depends on personal interests and bias
- need to learn user-specific weights for multi-criteria ranking (relevance, authority, freshness, etc.)
- can exploit user behavior (feedback, bookmarks, query logs, click streams, etc.)



or



Personal Information Management (PIM):

- manage, annotate, organize, and search all your personal data
 - on desktop (mail, files, calendar, etc.)
 - at home (photos, videos, music, parties, invoices, tax filing, etc.) and in smart home with ambient intelligence

Query-Log and Click-Stream Sharing in Communities



Communities

About I-Spy

max planck

Search

PRIVATE SEARCH

You are currently in the computer science community

Recent Queries

- [padprints](#)
- [joachims webwatcher](#)
- [seligmann live web stat...](#)
- [insyder content-based r...](#)
- [mann visualization www](#)

[VIEW ALL](#)

Recent Webpages

- [PadPrints: Graphical Mu...](#)
- [Citations: Webwatcher ...](#)
- [Dying Link](#)
- [Bell Labs presents new ...](#)
- [Bell Labs: Bell Labs Pr...](#)

[VIEW ALL](#)

Recent Communities

- [fourth year cs lab](#)
- [first year cs lab](#)
- [second year cs lab](#)
- [online shopping](#)
- [hdip/msc cs lab](#)

[VIEW ALL](#)

Popular Queries

- [help with java](#)
- [chr ei](#)

Popular Webpages

- [Java Technology](#)
- [Welcome to the UCD](#)

Popular Communities

collect user queries, derive community profiles,
and learn/adjust weights in multi-criteria ranking

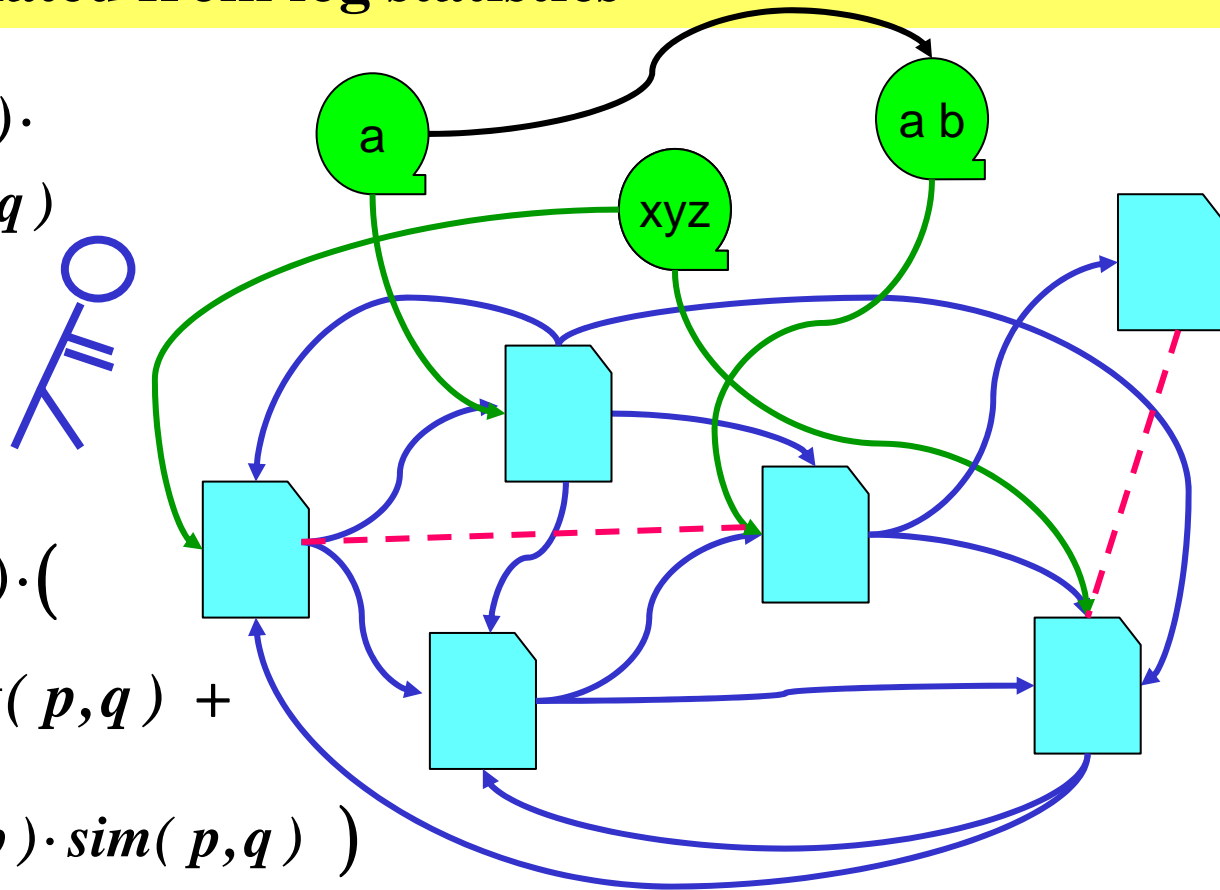
<http://ispy.ucd.ie/>

Exploiting Query Logs and Click Streams

from PageRank: uniformly random choice of **links** + random jumps
to QRank: + **query-doc transitions** + query-query transitions
+ **doc-doc transitions** on implicit links (w/ thesaurus)
with probabilities estimated from log statistics

$$PR(q) = \varepsilon \cdot j(q) + (1 - \varepsilon) \cdot \sum_{p \in IN(q)} PR(p) \cdot t(p, q)$$

$$QR(q) = \varepsilon \cdot j(q) + (1 - \varepsilon) \cdot \left(\alpha \sum_{p \in explicitIN(q)} PR(p) \cdot t(p, q) + (1 - \alpha) \sum_{p \in implicitIN(q)} PR(p) \cdot sim(p, q) \right)$$



Preliminary Experiments

Setup:

70 000 Wikipedia docs, 18 volunteers posing Trivial-Pursuit queries
ca. 500 queries, ca. 300 refinements, ca. 1000 positive clicks
ca. 15 000 implicit links based on doc-doc similarity

Results (assessment by blind-test users):

- QRank top-10 result preferred over PageRank in 81% of all cases
- QRank has 50.3% precision@10, PageRank has 33.9%

Untrained example query „philosophy“:

<u>PageRank</u>	<u>QRank</u>
1. Philosophy	Philosophy
2. GNU free doc. license	GNU free doc. license
3. Free software foundation	Early modern philosophy
4. Richard Stallman	Mysticism
5. Debian	Aristotle

1.7 Peer-to-Peer (P2P) Search

P2P systems: decentralized, self-organizing, highly dynamic networks of loosely coupled, autonomous computers

Applications:

- Large-scale distributed computation (SETI, PrimeNumbers, etc.)
- **File sharing** (Napster, Gnutella, KaZaA, BitTorrent, etc.)
- IP telephony (Skype)
- Publish-Subscribe Information Sharing (Auctions, Blogs, etc.)
with continuous queries (subscriptions) and alerting on updates
- Collaborative Work (Games, etc.)
- Collaborative Data Mining
- **(Collaborative) Web Search** (much harder than file search)

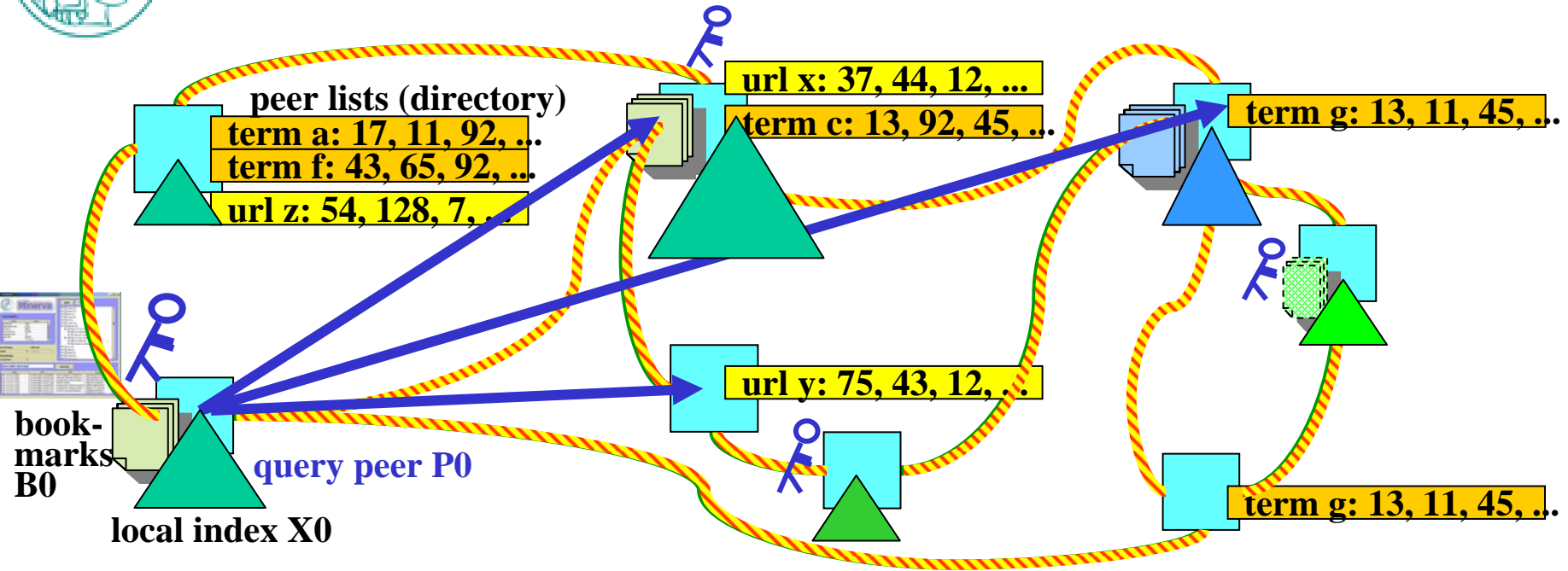
Why P2P Web Search

Objective: Self-organizing P2P Web Search Engines
with Google-or-better functionality

- **Scalable & Self-Organizing** Data Structures and Algorithms
(DHTs, Semantic Overlay Networks, Epidemic Spreading, Distr. Link Analysis, etc.)
- Better Search Result **Quality** (Precision, Recall, etc.)
 - Powerful Search Methods for Each Peer
(Concept-based Search, Query Expansion, Personalization, etc.)
 - Leverage Intellectual Input at Each Peer
(Bookmarks, Feedback, Query Logs, Click Streams, Evolving Web, etc.)
 - Collaboration among Peers
(Query Routing, Incentives, Fairness, Anonymity, etc.)
- Small-World Phenomenon
Breaking Information Monopolies

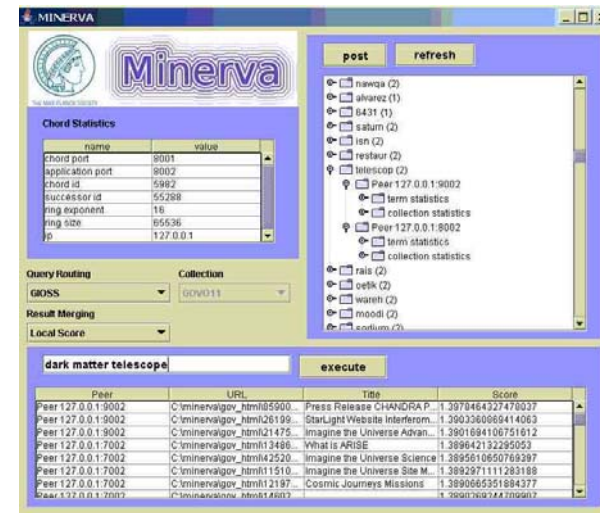


Minerva System Architecture



Query routing aims to optimize benefit/cost driven by distributed statistics on peers' content similarity, content overlap, freshness, authority, trust, performability etc.

Dynamically precompute „good peers“ to maintain a **Semantic Overlay Network** using random but biased graphs



1.8 Multimedia and NLP Search

search for images, speech, audio files, videos, etc.:

- based on **signal-level content features**
(color distribution, contours, textures, video shot sequence, pitch change patterns, harmonic and rhythmic features, etc. etc.)
- complement signal-level features with **annotations** from context
(e.g. adjacent text in Web page, GPS coordinates from digital camera)
- **query by example**: similarity search w.r.t. given object(s)
plus relevance feedback

question answering (QA) in natural language:

- express query as NL question: Who ..., When ..., Where ..., What ...
- provide short NL passages as query result(s), not entire documents

Content-based Image Retrieval by Example (1)

SIMPLICITY

Semantics-sensitive Integrated Matching for Picture Libraries

Option 1 --> Image ID or URL
similar images:

Option 2 --> **Random**

Option 3 --> Click an image to find



<http://wang.ist.psu.edu/IMAGE/>

Content-based Image Retrieval by Example (2)

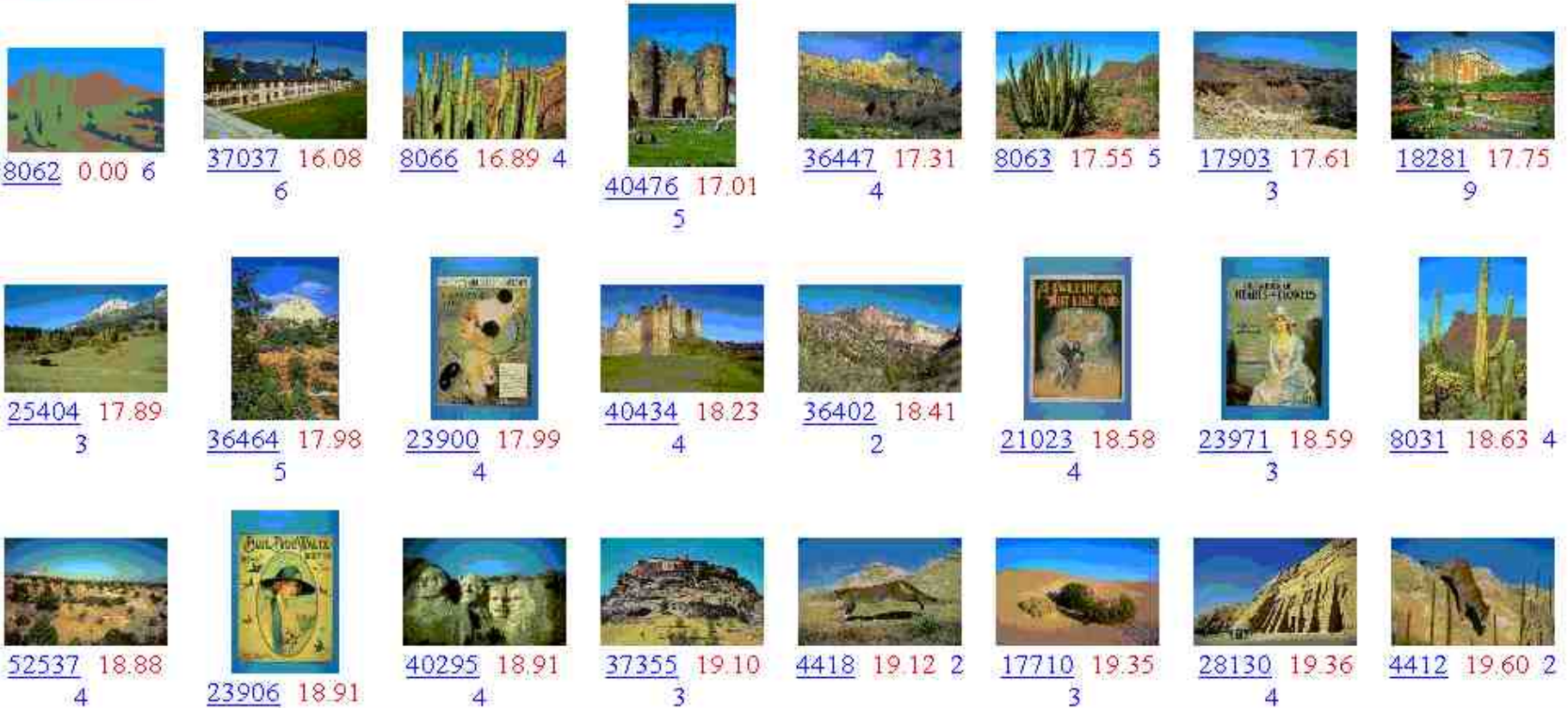
S I M P L I C I T Y

Semantics-sensitive Integrated Matching for Picture Libraries

Option 1 --> Image ID or URL

Option 2 --> **Random**

Option 3 --> Click an image to find similar images



Automatic Annotation of Images

S-I-M-P-L-I-c-i-t-y *a-LIP*

Automatic Linguistic Indexing of Pictures (wang.ist.psu.edu/IMAGE)

Random <- Click **random** to see more examples randomly selected from 60,000 images.
Or, click on a keyword below to search the computer-annotated image database by the selected keyword.

All keywords for the top 5 categories are shown. Words in **bold** are those top picks, selected by computer based on their statistical significances, for annotation. The computer is selecting keywords from **a dictionary of 600 automatically learned concepts**. Click [here](#) for more information about the project.



Computer Predictions-
grass

Grand_canyon

lion lizard animal mushroom
wild_cat waterfall cloud lake
rock mountain

Manual Category Annotation-
lizard animal rock



Computer Predictions-
horse grass

barn_yard animal
nature car rural

plant flower man-made landscape

Manual Category Annotation-
Kenya Africa animal landscape
people



Computer Predictions-
balloon sky ski city

man-made snow mountain
cloth building

historical_building people landscape

Manual Category Annotation-
city historical_building building



Computer Predictions-
man-made sport
car pill dining



Computer Predictions-
Toronto aviation rural
scene grass plane
England landscape city



Computer Predictions-
sun sky cloud dawn
drink dusk orange food
plant indoor

<http://wang.ist.psu.edu/IMAGE/>

Natural-Language Question Answering

AnswerBus

who invented quantum theory

Ask

Type in your question in English, French, Spanish, German, Italian or Portuguese:

Question:

who invented quantum theory

Possible answers: [XML](#) [TXT](#)

1. [Shortly after quantum field theory was invented, people started trying to invent a quantum field theory of gravity.](#)
2. [Quantum field theory was invented to deal simultaneously with special relativity and quantum mechanics, the two greatest discoveries of early twentieth-century physics, but it has become increasingly important to many areas of physics.](#)
3. [Quantum field theory was invented to reconcile quantum mechanics with special relativity.](#)
4. [Planck himself in \[7\] explains how despite having invented quantum theory he did not understand it himself at first: - I tried immediately to weld the elementary quantum of action somehow in the framework of classical theory.](#)
5. [Meanwhile, there is no doubt that quantum mechanics is the most successful theory of physical phenomena yet invented by the human mind.](#)
6. [It has not been reached?not by quantum theory, not by special or general relativity, not by anything invented since.](#)
7. [Austrian physicist Erwin Schrödinger, who, like Albert Einstein, never really believed in quantum theory, invented the story of a cat, now named after him, to illustrate how absurd the situation is.](#)

find compact text passages for question answering
<http://answerbus.coli.uni-sb.de/index.shtml>

Additional Literature (1)

important conferences on IR and DM

(see DBLP bibliography for full detail, <http://www.informatik.uni-trier.de/~ley/db/>)

SIGIR, ECIR, CIKM, TREC, WWW, KDD, ICDM, ICML, ECML

performance evaluation initiatives:

- Text Retrieval Conference (TREC), <http://trec.nist.gov>
- Cross-Language Evaluation Forum (CLEF), www.clef-campaign.org
- Initiative for the Evaluation of XML Retrieval (INEX),
<http://inex.is.informatik.uni-duisburg.de/>
- KDD Cup, <http://www.kdnuggets.com/datasets/kddcup.html>
and <http://kdd05.lac.uic.edu/kddcup.html>
- Language-Independent Named-Entity Recognition,
www.cnts.ua.ac.be/conll2003/ner/

Additional Literature (2)

Crawling, storage, and server management:

- S. Brin, L. Page: The Anatomy of a Large-Scale Hypertextual Web Search Engine, WWW 1998
- S. Ghemawat, H. Gobioff, S.-T. Leung: The Google File System, SOSP 2003
- P. Boldi, B. Codenotti, M. Santini, S. Vigna. Ubicrawler: A Scalable Fully Distributed Web Crawler, Software: Practice & Experience, 34(8):711-726, 2004.
- A. Heydon, M. Najork: Mercator: A Scalable, Extensible Web Crawler, WWW 1999.
- V. Shkapenyuk, T. Suel: Design and Implementation of a High-Performance Distributed Web Crawler, ICDE 2002
- X. Long, T. Suel: Three-level caching for efficient query processing in large Web search engines, WWW 2005

Web structure, size, dynamics:

- D.E. Rose, D. Levinson: Understanding User Goals in Web Search, WWW 2004
- A. Gulli, A. Signorini: The Indexable Web is More Than 11.5 Billion Pages, WWW 2005
- A. Ntoulas, J. Cho, C. Olston: What's New on the Web? The Evolution of the Web from a Search Engine Perspective, WWW 2004
- D. Fetterly, M. Manasse, M. Najork: Spam, Damn Spam, and Statistics, WebDB 2004
- D. Fetterly, M. Manasse, M. Najork, J. Wiener: A Large-Scale Study of the Evolution of Web Pages, WWW 2003
- A.Z. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins,
- J.L. Wiener: Graph structure in the Web, WWW 2000
- D. Donato et al.: Mining the Inner Structure of the Web Graph, WebDB 2005
Web graph data and tools, <http://webgraph.dsi.unimi.it/>
- Center for Complex Network Research, <http://www.nd.edu/~networks/>
- Search Engine Watch, <http://searchenginewatch.com/>

Additional Literature (3)

Geo- and time-aware search:

- A. Markowetz, Y.-Y. Chen, T. Suel, X. Long, B. Seeger: Design and Implementation of a Geographic Search Engine, WebDB 2005
- D. Ancona, J.Frew, G.Jané, and D.Valentine: Accessing the Alexandria Digital Library from Geographic Information Systems, JCDL 2004
- J. Ding, L. Gravano, N. Shivakumar: Computing Geographical Scopes of Web Resources, VLDB 2000
- K. Berberich, M. Vazirgiannis, G. Weikum: Time-aware Authority Ranking, to appear in Internet Mathematics Journal

Deep Web search:

- Kevin Chen-Chuan Chang, Bin He, Zhen Zhang: Toward Large Scale Integration: Building a MetaQuerier over Databases on the Web, CIDR 2005
- B. He, M. Patel, Z. Zhang, K. C.-C. Chang : Accessing the Deep Web: a Survey, Communications of the ACM, 2006
- Luciano Barbosa, Juliana Freire: Searching for Hidden-Web Databases, WebDB 2005

Intranet and enterprise search:

- IBM Systems Journal 43(3), 2004, Special Issue on Unstructured Information Management, <http://www.research.ibm.com/journal/sj43-3.html>
- R. Fagin, R. Kumar, K.S. McCurley, J. Novak, D. Sivakumar, J.A. Tomlin, D.P. Williamson: Searching the Workplace Web, WWW 2003
- Bjorn Olstad: Why Search Engines are used increasingly to Offload Queries from Databases, Keynote, VLDB 2005, <http://www.vldb2005.org/program/slides/tue/s1-olstad.ppt>
- Aleksander Ohrn: Contextual Insight in Search: Enabling Technologies and Applications, Tutorial, VLDB 2005, <http://www.vldb2005.org/program/slides/wed/s1366-ohrn.ppt>

Additional Literature (4)

Personalized search and PIM:

- J. Luxenburger, G. Weikum: Query-Log Based Authority Analysis for Web Information Search, WISE 2004
- E. Balfe, B. Smyth: An Analysis of Query Similarity in Collaborative Web Search, ECIR 2005
- E. Balfe, B. Smyth: Improving Web Search through Collaborative Query Recommendation, ECAI 2004: 268-272
- T. Mitchell: Computer Workstations as Intelligent Agents, Keynote, SIGMOD 2005, <http://www.cs.cmu.edu/~tom/>
- G. Bell: MyLifeBits: a Memex-Inspired Personal Store; Another TP Databa, Keynote, SIGMOD 2005, <http://research.microsoft.com/users/GBell/>
- X. Dong, A. Halevy: A Platform for Personal Information Management and Integration, CIDR 05
- X. Dong, A. Halevy, Jayant Madhavan: Reference Reconciliation in Complex Information Spaces, SIGMOD 2005

P2P search and collaboration:

- M. Bender, S. Michel, P. Triantafillou, G. Weikum, C. Zimmer: Improving Collection Selection with Overlap Awareness in P2P Search Engines, SIGIR 2005
- J.X. Parreira, G. Weikum: JXP: Global Authority Scores in a P2P Network, WebDB 2005
- J. Zhang, T. Suel: Efficient Query Evaluation on Large Textual Collections in a Peer-to-Peer Environment, Int.Conf. on Peer-to-Peer Computing, 2005
- F. M. Cuenca-Acuna, C. Peery, R. P. Martin, T. D. Nguyen: PlanetP: Using Gossiping to Build Content Addressable Peer-to-Peer Information Sharing Communities, HPDC 2003
- Christos Tryfonopoulos, Stratos Idreos, Manolis Koubarakis: Publish/Subscribe Functionality in IR Environments using Structured Overlay Networks, SIGIR 2005

Additional Literature (5)

Multimedia and NLP search:

- J.Z. Wang, J. Li, G. Wiederhold: SIMPLIcity: Semantics-sensitive Integrated Matching for Picture Libraries, IEEE Trans. on Pattern Analysis and Machine Intelligence 23(9), 2001
- J. Li, J.Z. Wang: Automatic linguistic indexing of pictures by a statistical modeling approach, IEEE Transactions on Pattern Analysis and Machine Intelligence 25(9), 2003
- M. Ortega-Binderberger, S. Mehrotra: Relevance feedback techniques in the MARS image retrieval system. Multimedia Syst. 9(6), 2004
- J. Fauqueur, N. Boujemaa: Region-based image retrieval: fast coarse segmentation and fine color description. J. Vis. Lang. Comput. 15(1), 2004
- A. Natsev, R. Rastogi, K. Shim: WALRUS: A Similarity Retrieval Algorithm for Image Databases, IEEE Trans. Knowl. Data Eng. 16(3), 2004
- Y. Zhu, D. Shasha: Warping Indexes with Envelope Transforms for Query by Humming, SIGMOD 2003
- E. Agichtein, S. Lawrence, L. Gravano: Learning to find answers to questions on the Web. ACM TOIT 4(2): 129-162 (2004)
- Ganesh Ramakrishna, Soumen Chakrabarti, Deepa Paranjpe, Pushpak Bhattacharya: Is question answering an acquired skill? WWW 2004
- Zhiping Zheng: Question Answering Using Web News as Knowledge Base, EACL 2003