

Web Technologies
(Technologien für das Internet I)
Foundations of Information Retrieval
<http://www2.kbs.uni-hannover.de/internet1.html>

Introduction

Prof. Wolfgang Nejdl, Elena Demidova

Institut für Verteilte Systeme
L3S Research Center
Leibniz Universität Hannover

13 October 2011

Plan for today

- Organizational issues
- Course overview

Outline

1 Organizational issues

2 Course overview

Information for the ITIS students

Transmission will be available to the Internet Technologies and Information Systems (ITIS) students located outside Hannover upon request at

https://webconf.vc.dfn.de/foundations_of_ir/. If this applies to you, please ask for a password via email from Elena Demidova, demidova (at) L3S.de

Lecture and exam dates

- We meet every Thursday 11:30 - (ca.) 14:00. The exercise sessions follow the lectures (do not be late!).
- The lectures start next week (October 20).
- The exercise sessions start the week after (October 27).
- StudIP: please register for the exercise sessions (mailing list).
- Exam: 90 minutes written exam on the March 5, 2012.
- Prerequisites for the ITIS-students will be clarified later.

Outline

1 Organizational issues

2 Course overview

Literature

Selected Chapters of the IIR Book: Christopher D. Manning, Prabhakar Raghavan and Hinrich Schtze, Introduction to Information Retrieval, Cambridge University Press. 2008.
<http://www-nlp.stanford.edu/IR-book/>

Outlook for the lectures

- We will look at the algorithms and technologies used in the modern search engines to satisfy informational needs of the users from within large document collections (usually stored on computers).

Outlook for the lectures

- We will look at the algorithms and technologies used in the **modern search engines** to satisfy informational needs of the users from within large document collections (usually stored on computers).

Outlook for the lectures

- We will look at the algorithms and technologies used in the modern search engines to satisfy **informational needs** of the users from within large document collections (usually stored on computers).

Outlook for the lectures

- We will look at the algorithms and technologies used in the modern search engines to satisfy informational needs of the users from within **large document collections** (usually stored on computers).

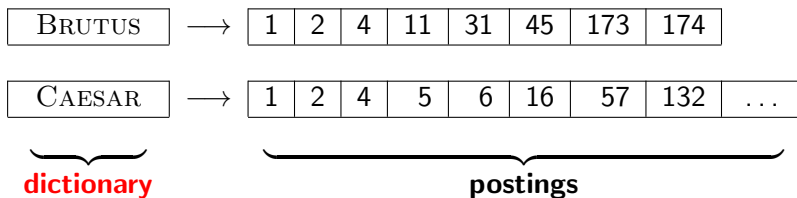
Outlook for the lectures

- We will look at the algorithms and technologies used in the **modern search engines** to satisfy **informational needs** of the users from within **large document collections** (usually stored on computers).

IIR 01: Boolean retrieval

- Design and data structures of a simple information retrieval system
- Queries are Boolean expressions, e.g., CAESAR AND BRUTUS
- The search engine returns all documents that satisfy the Boolean expression.

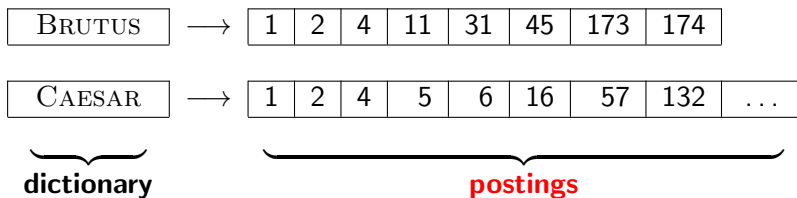
Inverted index:



IIR 01: Boolean retrieval

- Design and data structures of a simple information retrieval system
- Queries are Boolean expressions, e.g., CAESAR AND BRUTUS
- The search engine returns all documents that satisfy the Boolean expression.

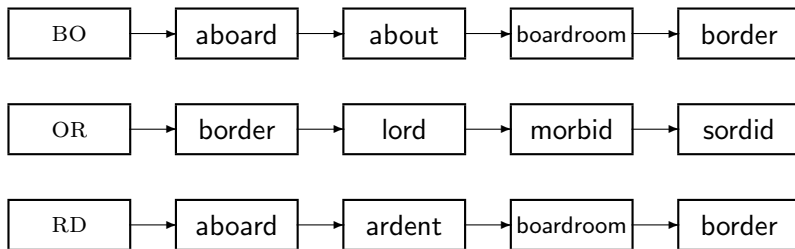
Inverted index:



IIR 02: The term vocabulary and postings lists

- Phrase queries: Stanford University
- Proximity: find Gates near Microsoft
- We need an index that captures position information for phrase queries and proximity queries.

IIR 03: Dictionaries and tolerant retrieval



IIR 06: Scoring, term weighting and the vector space model










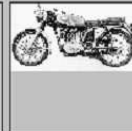


- Ranking search results
- Boolean queries only give inclusion or exclusion of documents.
- For ranked retrieval, we measure the proximity from query to each document.
- One formalism for doing this: the vector space model

IIR 08: Evaluation and dynamic summaries

- Benchmarks (e.g. TREC = Text Retrieval Conference)
- Measures (Precision, Recall, etc.)

IIR 09: Relevance feedback and query expansion

[Browse](#)
[Search](#)
[Prev](#)
[Next](#)
[Random](#)

					
(144538, 523493) 0.54182 0.231944 0.309876	(144538, 523835) 0.56319296 0.267304 0.295889	(144538, 523529) 0.584279 0.280881 0.303398	(144456, 253569) 0.64501 0.351395 0.293615	(144456, 253568) 0.650275 0.411745 0.23853	(144538, 523799) 0.66709197 0.358033 0.309039
					
(144473, 16249) 0.6721 0.393922 0.278178	(144456, 249634) 0.675018 0.4639 0.211118	(144456, 253693) 0.676901 0.47645 0.200451	(144473, 16328) 0.700339 0.309002 0.391337	(144483, 265264) 0.70170796 0.36176 0.339948	(144478, 512410) 0.70297 0.469111 0.233839

IIR 10: XML retrieval

- Semi-structured / structured documents vs. unstructured documents
- Can we utilize the structure of the data in IR systems?
- Databases support search for numerical range and exact match, e.g., salary < 60,000 and manager=Smith.
- If your data is structured and you only need precise queries like this (numerical, exact match etc), do not use an IR system.

IIR 13: Text classification and Naive Bayes

- Text classification = assigning documents automatically to predefined classes
- Examples:
 - a. Language (English vs. French)
 - b. Location

IIR 16: Flat clustering



the Web

[Advanced Search](#)
[Search](#)
[Help](#)

Clustered Results

Top 208 results of at least 20,373,974 retrieved for the query **jaguar** ([Details](#))

- ▶ [jaguar](#) (208)
- ▶ [Cars](#) (74)
- ▶ [Club](#) (34)
- ▶ [Cat](#) (23)
- ▶ [Animal](#) (13)
- ▶ [Restoration](#) (10)
- ▶ [Mac OS X](#) (8)
- ▶ [Jaguar Model](#) (8)
- ▶ [Request](#) (5)
- ▶ [Mark Webber](#) (6)
- ▶ [Maya](#) (5)
- ▼ [More](#)

Find in clusters:

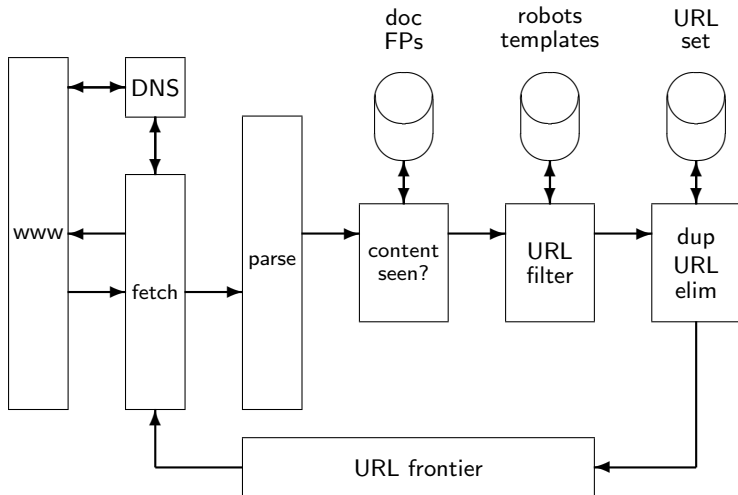


1. [Jag-lovers - THE source for all Jaguar information](#) [new window] [frame] [cache] [preview] [clusters]
 ... Internet! Serving Enthusiasts since 1993 The Jag-lovers Web Currently with 40661 members The Premier **Jaguar** Cars web resource for all enthusiasts Lists and Forums Jag-lovers originally evolved around its ...
[www.jag-lovers.org](#) - Open Directory 2, Wisenut 8, Ask Jeeves 8, MSN 9, Looksmart 12, MSN Search 18
2. [Jaguar Cars](#) [new window] [frame] [cache] [preview] [clusters]
 [...] redirected to [www.jaguar.com](#)
[www.jaguarcars.com](#) - Looksmart 1, MSN 2, Lycos 3, Wisenut 6, MSN Search 9, MSN 29
3. <http://www.jaguar.com/> [new window] [frame] [preview] [clusters]
[www.jaguar.com](#) - MSN 1, Ask Jeeves 1, MSN Search 3, Lycos 9
4. [Apple - Mac OS X](#) [new window] [frame] [preview] [clusters]
 Learn about the new OS X Server, designed for the Internet, digital media and workgroup management. Download a technical factsheet.
[www.apple.com/macosx](#) - Wisenut 1, MSN 3, Looksmart 26

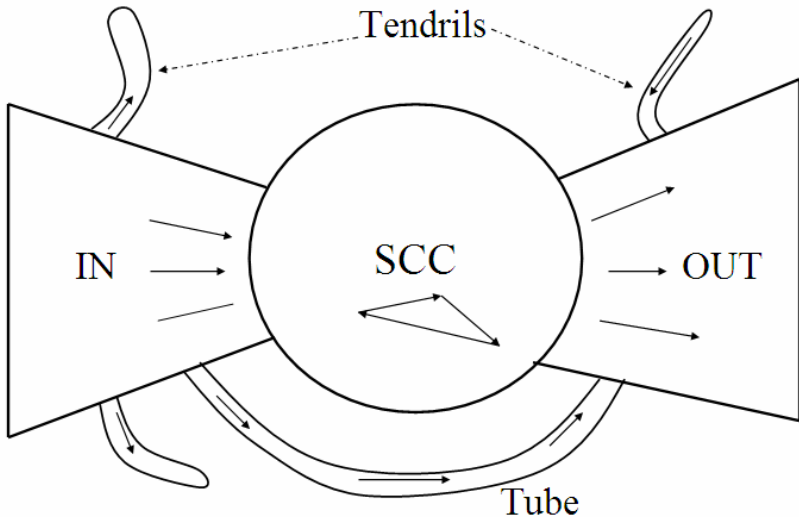
IIR 19: Web search

- Unusual and diverse documents
- Unusual and diverse users and information needs
- Beyond terms and text: exploit link analysis, user data
- How do web search engines work?
- How can we make them better?

IIR 20: Crawling



IIR 21: Link analysis / PageRank



Questions?

Thank you for your attention!