



Google, Deep Web und Fachdatenbanken

Dirk Lewandowski

dirk.lewandowski@haw-hamburg.de

Vision

Eine Datenbank

- die alle Themenbereiche abdeckt,
- eine optimale Recherche erlaubt,
- leicht zu bedienen ist,
- die Dokumente optimal erschließt,
- die Volltexte sofort (und möglichst kostenlos) zur Verfügung stellt.

Agenda

Google

Das Deep Web und seine Bedeutung für wissenschaftliche Inhalte

Überblick Wissenschaftssuchmaschinen

Entwicklungen bei Fachdatenbanken

Fazit

Agenda

Google

Das Deep Web und seine Bedeutung für wissenschaftliche Inhalte

Überblick Wissenschaftssuchmaschinen

Entwicklungen bei Fachdatenbanken

Fazit

Google Websuche

- **In Deutschland laufen mehr als 80 Prozent aller Web-Suchen über Google.**
- **Aber: Hoher Anteil von Suchanfragen, die anderswo in Suchfelder eingetragen werden.
-> Suche ist nicht gleichbedeutend mit Web-Suche.**
- **Google wird von den Nutzern als die beste Suchmaschine wahrgenommen (empirisch nicht nachgewiesen).**
- **Google hat die Standards für SM-Interfaces gesetzt.**
- **Erfolgsfaktoren von Google**
 - Bedienbarkeit
 - Geschwindigkeit
 - Indexgröße
 - Qualität des Rankings

Agenda

Google

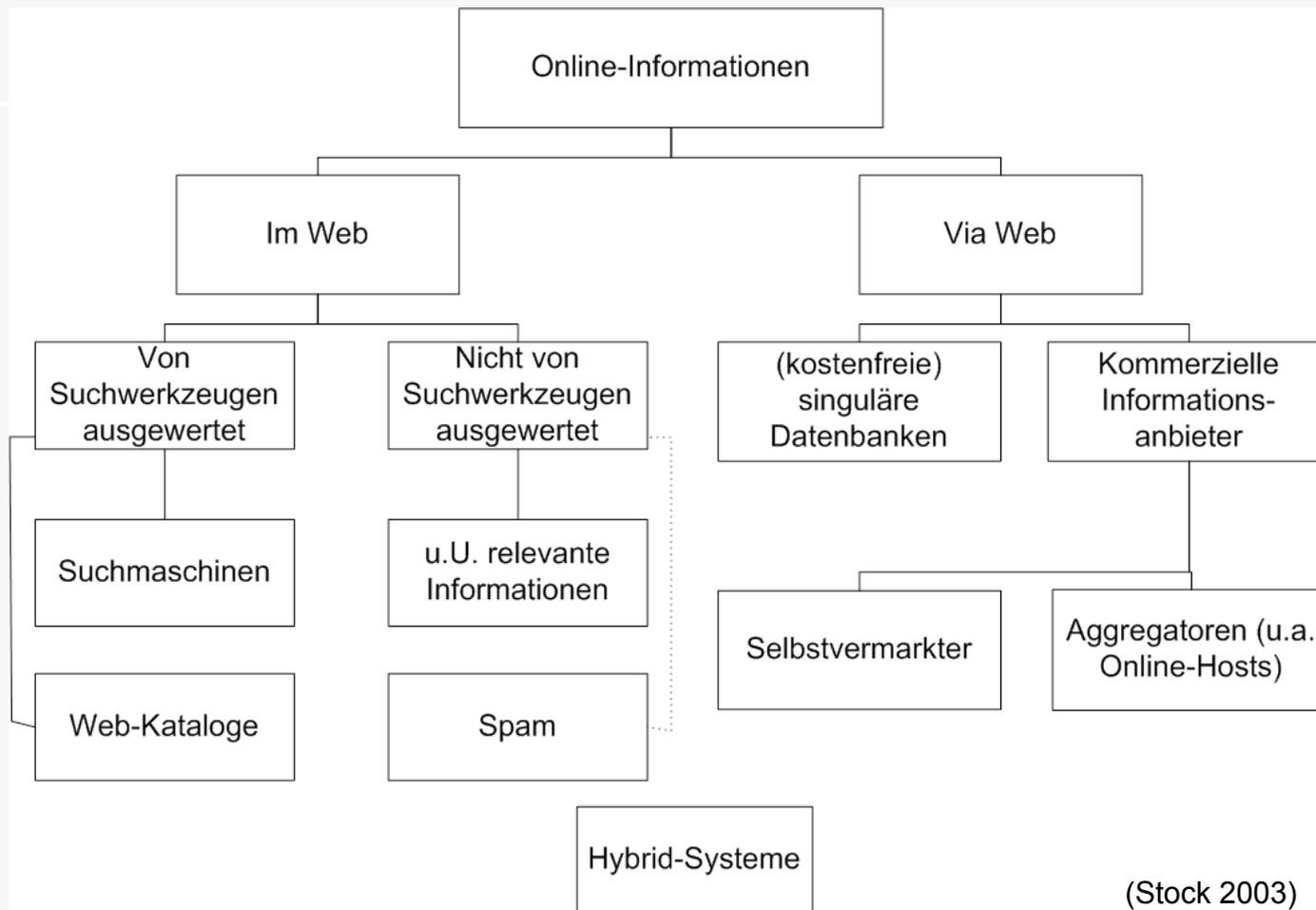
Das Deep Web und seine Bedeutung für wissenschaftliche Inhalte

Überblick Wissenschaftssuchmaschinen

Entwicklungen bei Fachdatenbanken

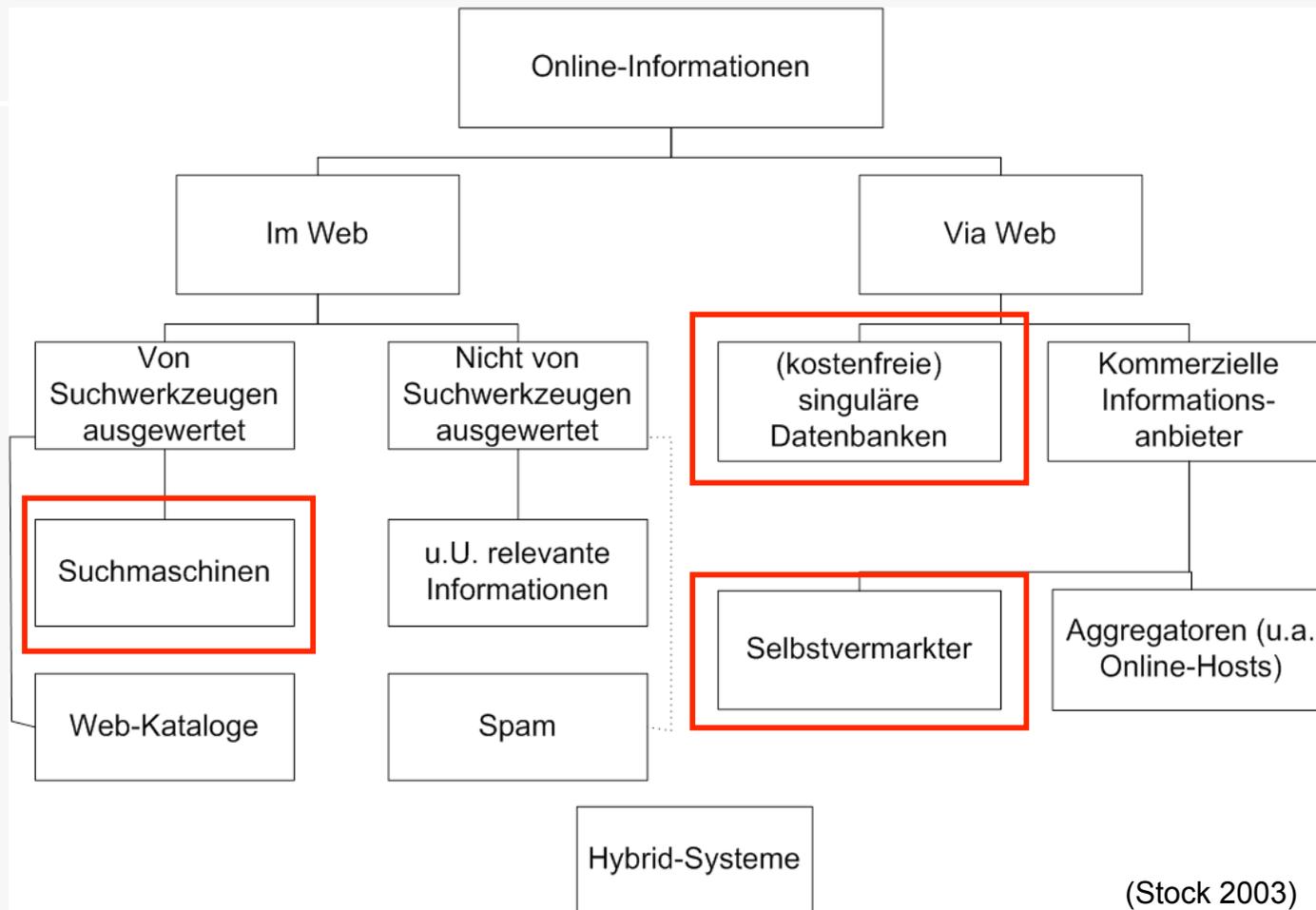
Fazit

Taxonomie der digitalen Online-Information



(Stock 2003)

Taxonomie der digitalen Online-Information



Surface Web vs. Invisible Web

Definitionen des Invisible/Deep Web

- “Text pages, files, or other often high-quality authoritative information available via the World Wide Web that general-purpose search engines cannot, due to technical limitations, or will not, due to deliberate choice, add to their indices of Web pages” (Sherman u. Price 2001).
- “The deep Web - those pages do not exist until they are created dynamically as the result of a specific search“ (Bergman 2001).
- **Surface Web:** Alle Inhalte, die von den allgemeinen Suchmaschinen erschlossen werden (können).
- **Invisible Web:** Alle Inhalte, die von den allgemeinen Suchmaschinen nicht erschlossen werden (können), vor allem die Inhalte von Datenbanken, die über das Web erreichbar sind.

Bereiche des (Academic) Web

- **Academic Surface Web**

- Wissenschaftliche Inhalte im Oberflächenweb.
- Alle Seiten von Unis, Forschungseinrichtungen, usw.
- Wissenschaftliche Texte.

- **Academic Invisible Web**

- Vor allem Inhalte aus wissenschaftlich relevanten Datenbanken.
- Bibliothekskataloge, Literaturdatenbanken, Bücher, Aufsätze, Forschungsdaten,
...

Bedeutung des Academic Invisible Web

- **Die Inhalte sind für den gesamten wissenschaftlichen Prozeß von Bedeutung.**
 - Literatur (Artikel, Dissertationen, Report, Bücher, usw.).
 - Forschungsdaten.
 - Reine Online-Inhalte (u.a. Open-Access-Inhalte).
- **Anbieter von IW-Inhalten**
 - Datenbank-Anbieter (Metadaten + intellektuelle Erschließung).
 - Bibliotheken (Bibliothekskataloge, Sammlungen + intellektuelle Erschließung).
 - Verlage (Volltexte + automatische/teile intellektuelle Erschließung).
 - Repositories.

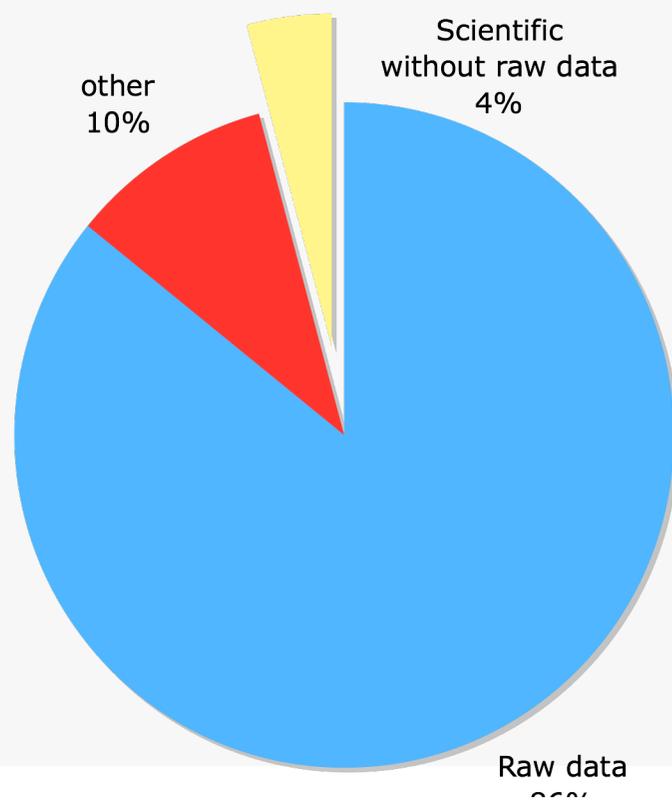
Größe des (Academic) Invisible Web

- **Größe des Invisible Web nach Bergman (2001): 550 Milliarden Dokumente**
 - Berechnung: Durchschnittliche Größe der bekannten (großen) IW-Datenbanken * geschätzte Gesamtzahl der IW-Datenbanken.
 - Problem: Verteilung der Datenbank-Größen stark linksschief (Median: 4.950 Dokumente je Datenbank).
 - Wenige Datenbanken enthalten viele Dokumente (>100 Millionen), viele Datenbanken nur einige Tausend.
 - Tatsächliche Größe des IW dürfte bei <100 Milliarden Dokumenten liegen (Lewandowski&Mayr, 2006).
- **Gesamtgröße aller Datenbanken im Gale Directory of Databases: 18,92 Milliarden Dokumente.**
 - Verzeichnis von ca. 16.000 Datenbanken.
 - Manche der in Bergmans Liste aufgeführten Datenbanken fehlen.

Inhalte des Academic Invisible Web

Basis: Top60 größte IW-Datenbanken aus Bergman (2001)
Größenanteile auf Basis der Dateigrößen; nicht Zahl der Dokumente!

Contents of Bergman's Top 60



Zugang zum Academic Invisible Web - verschiedene Ansätze

- **Kommerzielle Suchmaschinen**
 - Google Scholar
 - Windows Live Academic
 - Scirus
- **Bibliotheken und Datenbank-Anbieter**
 - BASE (Bielefeld Academic Search Engine)
 - HBZ-Suchmaschine
 - Vascoda (Integration von Bibliotheks- und Datenbank-Inhalten)
- **Open Access Repositories**
 - Citebase
 - OpenROAR

Agenda

Google

Das Deep Web und seine Bedeutung für wissenschaftliche Inhalte

Überblick Wissenschaftssuchmaschinen

Entwicklungen bei Fachdatenbanken

Fazit

Wissenschaftssuchmaschinen werden das Academic Web erschließen.

Die Inhalte der Wissenschaftssuchmaschinen

- **Verlagsinhalte**
 - Bücher
 - Aufsätze
- **Graue Literatur aus dem Web**
 - Reports
 - Manuskripte
- **Open Access**
 - Zeitschriften
 - Repositorien
- **Inhalte aus Datenbanken**
- **(Forschungsdaten)**

Google Scholar: Inhalte

Inhalte von Google Scholar

- **Fächer**
 - prinzipiell alle Fächer, Schwerpunkt bei STM.
- **Quellen**
 - freies Web
 - Verlage und Fachgesellschaften
 - Open-Access-Archive und -Zeitschriften
 - **Kein** Quellenverzeichnis; Umfang der Quellen unklar
- **Dateiformate**
 - PDF, PS

Google Scholar bietet keine Qualitätskontrolle.

Arten von Inhalten in Google Scholar

- **Texttypen**
 - Zeitschriftenaufsätze (peer review), Konferenzbeiträge
 - Preprints, Postprints
 - Reports
 - Seminararbeiten
 - ...



Dokumenttypen

- Direkter Link auf den Volltext
- „Normale Literaturangabe“
- Zitation
- Buch

[BUCH] **Wegweiser im Netz: Qualität und Nutzung von Suchmaschinen**
M Machill... - 2003 - Verl. Bertelsmann-Stiftung
[Zitiert durch: 34](#) - [Ähnliche Artikel](#) - [Websuche](#) - [Library Search](#)

[Abfragesprachen und erweiterte Funktionen von WWW-Suchmaschinen](#)
D Lewandowski - IWP-Information: Wissenschaft und Praxis, 2004 - [www-public.rz.uni-due](#)
... **WWW-Suchmaschinen*** ... Zusammenfassung Der vorliegende Artikel stellt die erweiterten Suchmöglichkeiten in den wichtigsten **Suchmaschinen** vor. ...
[Zitiert durch: 11](#) - [Ähnliche Artikel](#) - [HTML-Version](#) - [Websuche](#)

[Suchmaschinen und Anfragen im World Wide Web](#)
U Masermann, G Vossen - Informatik-Spektrum, 1998 - Springer
Page 1. U.Masermann,G.Vossen: **Suchmaschinen** und Anfragen im World Wide Web
9 ... Hauptbeitrag **Suchmaschinen** und Anfragen im World Wide Web ...
[Zitiert durch: 13](#) - [Ähnliche Artikel](#) - [Websuche](#)

[BUCH] **Suchmaschinen im Internet**
M Glögler - 2003 - Springer
[Zitiert durch: 13](#) - [Ähnliche Artikel](#) - [Websuche](#) - [Library Search](#)

[Deutsche Suchmaschinen im Vergleich: AltaVista.de, Fireball.de, Google.de](#)
J Griesbaum, M Rittberger, B Bekavac - Information und Mobilität. Optimierung und Vermittlung
... S. 201 – 223 Deutsche **Suchmaschinen** im Vergleich: AltaVista.de, ... Es zeigen sich Vorteile für Google.de gegenüber den anderen **Suchmaschinen**. ...
[Zitiert durch: 10](#) - [Ähnliche Artikel](#) - [HTML-Version](#) - [Websuche](#)

[ZITATION] **Suchmaschinen. Metamedien im Internet?**
H Winkler - Virtualisierung des Sozialen
[Zitiert durch: 10](#) - [Ähnliche Artikel](#) - [Websuche](#)

[BUCH] **Reisen ohne Karte: Wie funktionieren Suchmaschinen?**

Instanzen eines Artikels werden zusammengeführt.

[The freshness of web search engine databases](#) - [Gruppe von 9 »](#)

[D Lewandowski, H Wahlig, G Meyer-Bautor - Journal of Information Science, 2006](#) - [jiss.sagepub.com](#)

Page 1. The **freshness of web search engine** databases 131 Journal of Information Science, 32 (2) 2006, pp. 131–148 © CILIP, DOI: 10.1177/0165551506062326 ...

[Zitiert durch: 8](#) - [Ähnliche Artikel](#) - [Websuche](#)

Instanzen eines Artikels werden zusammengeführt

- **Verlag**

[The freshness of web search engine databases](#)

D Lewandowski, H Wahlig, G Meyer-Bautor - Journal of Information Science, 2006 - jis.sagepub.com
Journal of Information Science, 32 (2) 2006, pp. 131–148 © CILIP, DOI:
10.1177/0165551506062326 ... Dirk Lewandowski, Henry Wahlig and Gunnar
Meyer-Bautor ... Department of Information Science, ...
[Zitiert durch: 8](#) - [Ähnliche Artikel](#) - [Websuche](#)

- **Portal**

[The freshness of web search engine databases](#)

D Lewandowski, H Wahlig, G Meyer-Bautor - Journal of Information Science, 2006 - portal.acm.org
This study measures the frequency with which search engines update their
indices. Therefore, 38 websites that are updated on a daily basis were analysed
within a time-span of six weeks. The analysed search engines were Google, ...
[Websuche](#)

- **Private Homepage**

[The freshness of web search engine databases](#)

D LEWANDOWSKI, H WAHLIG, G MEYER-BAUTOR - Journal of information science, 2006 - cat.inist
This study measures the frequency with which search engines update their
indices. Therefore, 38 websites that are updated on a daily basis were analysed
within a time-span of six weeks. The analysed search engines were Google, ...
[Websuche](#)

- **Open-Access-Archiv**

[The Freshness of Web search engines' databases](#)

D Lewandowski, H Wahlig, G Meyer-Bautor - durchdenken.de
The Freshness of Web search engines' databases ... Dirk Lewandowski, Henry
Wahlig and Gunnar Meyer-Bautor ... Department of Information Science,
Heinrich-Heine-University Düsseldorf, Germany ... Correspondence to: Dirk ...
[HTML-Version](#) - [Websuche](#)

[The Freshness of Web search engines' databases](#)

D Lewandowski, H Wahlig, G Meyer-Bautor - citebase.eprints.org
This study measures the frequency in which search engines update their indices.
Therefore, 38 websites that are updated on a daily basis were analysed within a
time-span of six weeks. The analysed search engines were Google, Yahoo and ...
[Im Cache](#) - [Websuche](#)

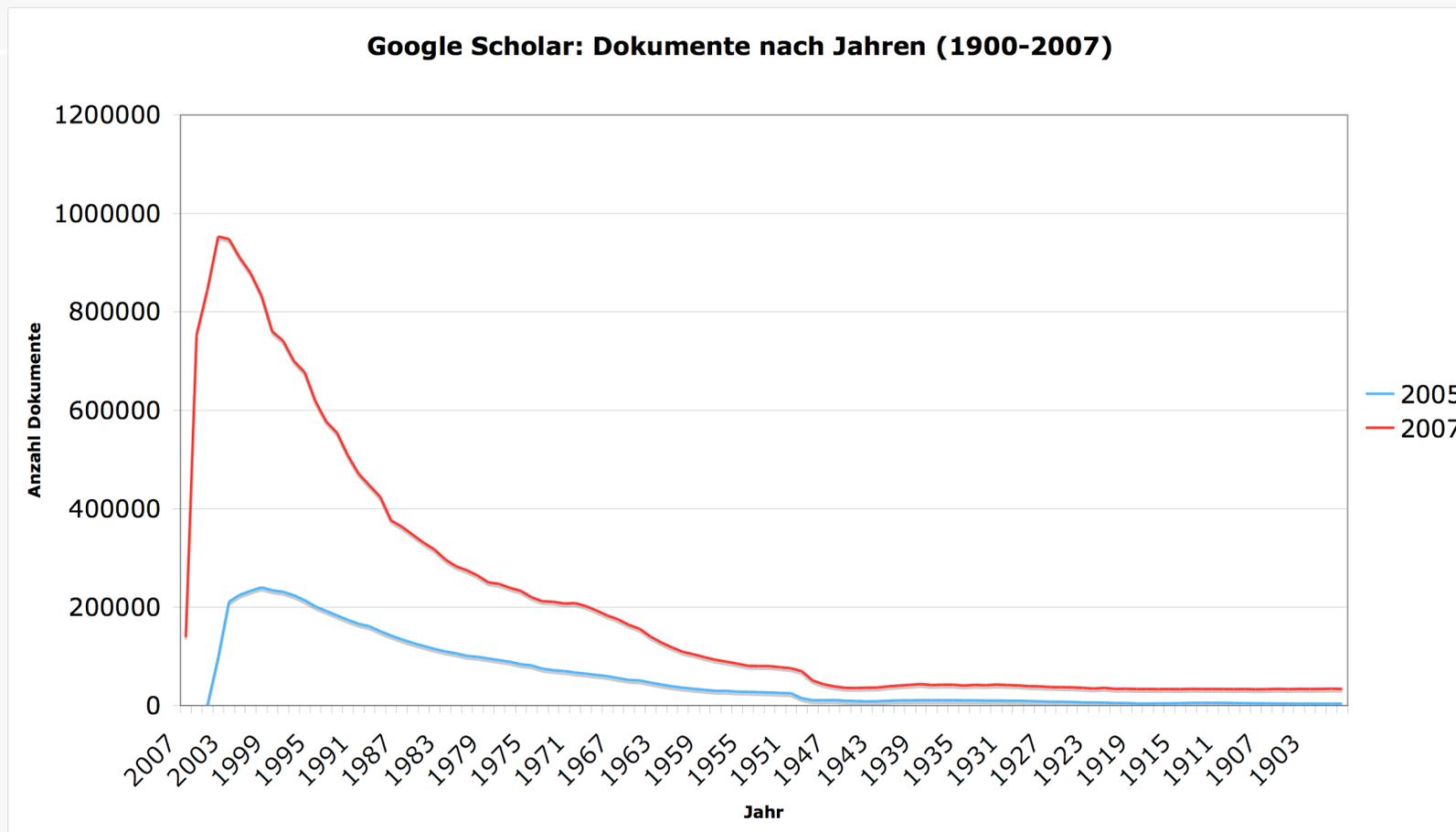
Die Suchmöglichkeiten sind sehr beschränkt.

Suchfunktionen

- **Suchformulare**
 - Einfache und erweiterte Suche analog zu Google.com
- **Einschränkungsmöglichkeiten**
 - Boolesche Operatoren (mit den von Google bekannten Einschränkungen)
 - Phrasensuche
 - Beschränkung auf Titel, Autor, Publikation, Datum. Nur eingeschränkt funktionstüchtig.
- **Browsing**
 - entlang der Zitationen.

Der Datenbestand ist in den letzten Jahren gewachsen

Gesamtbestand 2005: 2-7 Millionen Dok; 2007: ca. 25 Mio. Dok.



Windows Live Academic

- **Suche in Verlagsarchiven (Aufsätze).**
- **Verwendung der Verlags-Metadaten.**
- **Bereits in die MSN-Suche eingebunden.**
- **Frühes Entwicklungsstadium.**



Windows Live Academic

The screenshot shows the Windows Live Academic search interface. At the top, there is a search bar with the query "freshness of web search engines" and a "Sign In" link. Below the search bar, there are navigation tabs for "Web", "Images", "News", "Maps", "MSN", and "More", with "Academic Beta" selected. A "Sort by:" dropdown menu is set to "Relevance", with other options including "Date - Oldest", "Date - Newest", "Author", "Journal", and "Conference".

The search results are displayed in a list. The first result is titled "The freshness of web search engine databases - Published Version (2006)" by Lewandowski, Dirk | Wahlig, Henry | Meyer-Bautor, Gunnar. The abstract snippet reads: "... indices varies and more than one engine should be used when searching for current content. Key Words: search engines online information retrieval world wide web index quality Index freshness Search Web Hide Abstract".

The second result is "A Weighted Freshness Metric for Maintaining Search Engine Local ... - Published Version" by Cercone, N. | Han, Jianchao | Hu, Xiaohua. The abstract snippet reads: "Web search engines create and maintain a local repository, a local copy of a portion of ... Campbell, Internet search engine freshness by web server help, Proc. of the Symposium on the Internet ... Search Web Hide Abstract".

The third result is "Internet Search Engine Freshness by Web Server Help (2000)" by Gupta, Vijay | Campbell, Roy. The abstract snippet reads: "... web servers themselves track the changes happening to their content files for propagating updates to search engines. We propose an algorithm which uses both freshness and popularity of data at the web ... Search Web CiteSeer: 4 Hide Abstract".

The fourth result is "Internet search engine freshness by Web server help - Published Version".

On the right side, a detailed view of the first result is shown. It includes the title "The freshness of web search engine databases", the journal name "Journal of Information Science", and a full abstract: "This is a preprint of an article published in the Journal of Information Science Vol. 32, No. 2, 131-148 (2006). This study measures the frequency in which search engines update their indices. Therefore, 38 websites that are updated on a daily basis were analysed within a time-span of six weeks. The analysed search engines were Google, Yahoo and MSN. We find that Google performs best overall with the most pages updated on a daily basis, but only MSN is able to update all pages within a time-span of less than 20 days. Both other engines have outliers that are quite older. In terms of indexing patterns, we find different approaches at the different engines: While MSN shows clear update patterns, Google shows some outliers and the update process of the Yahoo index seems to be quite chaotic. Implications are that the quality of different search engine indices varies and not only one engine should be used when searching for current content." The authors listed are Lewandowski, Dirk | Wahlig, Henry | Meyer-Bautor, Gunnar, and the volume is 32.

At the bottom of the page, there is a copyright notice: "© 2007 Microsoft Trademarks | Privacy | Legal | For Site Owners" and a "Help Central | Account | Feedback" link.

Andere

Andere

- **Scirus**
 - Academic Surface Web (keine Beschränkung auf Literatur)
 - Teile des Academic Invisible Web
 - Elsevier-Content
- **Forschungsportal.net**
 - Websites der in Deutschland öffentlich geförderten Forschungseinrichtungen
 - Online-Dissertationen DDB
- **Nicht zu vergessen:**
 - Interdisziplinäre Literaturdatenbanken (Web of Science, Scopus)
 - Verlagsangebote (Springerlink, Science direct).
 - Google Buchsuche, Open Content Alliance, Amazon.

Die Wissenschaftssuchmaschinen starten mit einem hohen Anspruch.

Vorteile der kommerziellen Wissenschaftssuchmaschinen

- **Inhalte**
 - Alle Aufsätze, alle Bücher.
 - (teils) andere Inhalte des Academic Web.
- **Erschließung**
 - Volltexterschließung.
 - Anreicherung durch Volltexte (bzw. Ausschnitte), Rezensionen, „tags“
 - Empfehlungssysteme.
- **Suche**
 - Schnelle und einfache Suche.
 - Suchinterfaces wie bei allgemeinen Suchmaschinen.

Wissenschaftssuchmaschinen zeigen wenig Transparenz, die Inhalte sind meist schlecht erschlossen.

Nachteile der kommerziellen Wissenschaftssuchmaschinen

- **Unklare Quellenlage**
 - Ausgewertete Quellen werden nicht angegeben.
 - Quellen werden oft nicht vollständig erschlossen.
- **Mangelhafte Erschließung**
 - Hohe Fehlerrate (Autorennamen, Zeitschriftentitel)
 - Keine Erschließung mit Klassifikation, Schlagwörtern, usw.
- **Zu allgemeine Community**
 - Rezensionen bei Amazon werden von *irgendwelchen* Nutzern geschrieben.

Agenda

Google

Das Deep Web und seine Bedeutung für wissenschaftliche Inhalte

Überblick Wissenschaftssuchmaschinen

Entwicklungen bei Fachdatenbanken

Fazit

Klassische Stärken der Fachdatenbanken

- **Quellen**
 - Quellentransparenz.
 - Vollständige Abdeckung eines Fachgebiets.
- **Erschließung**
 - Kontrolliertes Vokabular
- **Zusatzdienste**
 - Alerting-Dienste
 - Gespeicherte Suchen
 - TOC-Alerts

Probleme der Fachdatenbanken und Reaktionen darauf

Lösung auf technischer Ebene durch Suchmaschinen-Technologie

- **Suchmöglichkeiten auf Profis ausgerichtet**
 - Endnutzer-Interfaces wie bei Web-Suchmaschinen.
 - Ranking.
- **Beschränkung auf ein Fachgebiet**
 - Zusammenführen von Fachdatenbanken unter einer Oberfläche.
- **Fehlende Navigation**
 - Drill-Down mit Menüs/Clustering
- **Bruch vom Literaturnachweis zum Dokument**
 - Link Resolver
- **Fehlende Volltexterschließung**



IBLK Fachportal Internationale Beziehungen und Länderkunde deutsch | english

Startseite
Suche in allen Datenbanken
Nur World Affairs Online
Nur Abkommen
Schlagwörter
Online Contents IBLK

Wir über uns
Aktuelles
Hilfe
FAQ
Kontakt
Impressum
Sitemap

Ihre Suche **schnell** **erweitert** ?

Ihre letzte Suche war:
Tite/Abstract: unruhen (alle Wörter)

Filter **Datenbasen** **Publikationstypen** **Sprache** ?

Filter der letzten Suche waren:
Datenbasen: Alle | Publikationstypen: Alle | Sprache: Alle

Suchergebnisse **Treffer (390)** ?

Optionen:

Sortieren nach: Relevanz

Markierte Treffer

alle markieren

1	<input type="checkbox"/> Zeitschriftenaufsatz: Beirut - Konflikt und Koexistenz in einer geteilten Stadt <input type="button" value="v"/> Autor/Hrsg.: Hanf, Theodor; 1985 In: Geographische Rundschau Publikationstyp: Zeitschriftenaufsatz <input type="checkbox"/> @ <input type="checkbox"/>	WAO
2	<input type="checkbox"/> Buchaufsatz: South Africa's military relations with its neighbours <input type="button" value="v"/> Autor/Hrsg.: Spence, Jack E.; 1986 Publikationstyp: Buchaufsatz <input type="checkbox"/>	WAO

Schränken Sie Ihr Ergebnis ein: ?

Person

- > Reuter, Jens (10)
- > Weggel, Oskar (5)
- > Peters, Ralph-Michael (3)
- > Oschlies, Wolf (3)
- > Toennes, Bernhard (2)

Datenbank

- > WAO (389)
- > PAIS (1)

Schlagwort

- > Politische Unruhen (119)
- > Innenpolitischer Konflikt (113)
- > Innenpolitische Lage/Entwicklung (71)
- > Politische Partei (43)
- > Innenpolitik (40)

Agenda

Google

Das Deep Web und seine Bedeutung für wissenschaftliche Inhalte

Überblick Wissenschaftssuchmaschinen

Entwicklungen bei Fachdatenbanken

Fazit

Fazit

- **Fachdatenbanken werden zunehmend zu Wissenschaftssuchmaschinen.**
- **Wissenschaftssuchmaschinen aus der Fach-Community können durch gute Erschließung punkten.**
- **Volltexte werden für die Ausnutzung der Stärken der Suchmaschinen-Technologie benötigt.**
- **Zusammenführung von (auch kommerziellen) Quellen muss vorangetrieben werden.**
- **Zugang zu den Volltexten muss vereinfacht werden.**

Fazit: Strategien für die Literaturrecherche

- **Primärquelle zur Recherche in einem Fachgebiet bleiben die Fachdatenbanken.**
- **Kommerzielle Wissenschaftssuchmaschinen eignen sich zur ergänzenden (interdisziplinären) Recherche.**
- **Crawlende Wissenschaftssuchmaschinen (Google Scholar) eignen sich in vielen Fällen zur schnellen Beschaffung von Volltexten.**



Vielen Dank für Ihre
Aufmerksamkeit.

Prof. Dr. Dirk Lewandowski
www.durchdenken.de/lewandowski
E-Mail: dirk.lewandowski@haw-hamburg.de