



# Recherchieren zu einem Wissenschaftlichen (Informatik-)Thema

Ulf Leser

# Inhalt

- Übersicht
- Wissenschaftliche Publikationen
- Bewertung von Publikationen
- Ressourcen und Techniken
- Literatursammlung erstellen und managen

# Wissenschaftliches Thema

- In (Informatik-)Seminaren werden Themen meist gestellt
  - Ein Problem, eine Methode, ein Vergleich
  - Meist Paper und / oder kurze Erklärung
- **Granularität**: Themen können sehr breit oder sehr eng sein
  - Breit: Übersichtsarbeiten, geringe technische Tiefe
  - Eng: Spezialarbeiten, hohe technische Tiefe
- Themen stehen meist nicht genau fest
  - Recherche ergibt neue Aspekte – Themenanpassung
  - Hier können **persönliche Präferenzen** eingebracht werden
- Anders: Recherchieren zur Themenfindung
  - Explorativ versus zielgerichtet

# Ziel einer Recherche

- Im Ergebnis der Recherche kennt man
  - Verschiedene **Aspekte des Problems**
  - Die verschiedenen Ansätze zur Lösung
  - Die **wichtigste Literatur** zum Thema
  - **Eigener Schwerpunkt** und welche Literatur man verwenden wird
- Vollständigkeit
  - IdR ist es unmöglich, alle Literatur zum Thema zu sichten
  - Bewertung der **Relevanz** und Auswahl notwendig
  - Teilweise macht man das selbst (lesen), teilweise vertraut man anderen (Wikipedia, Buch, Übersichtsartikeln, Datenbanken, ...)
  - Wichtig ist **gestufte Bewertung**: Titel/Autoren, Zitate, Erscheinungsort, Abstract, querlesen ...
  - Literaturliste immer mit dem Betreuer besprechen

# Allgemeines Vorgehen

- **Vor** der genauen Recherche muss Thema (intuitiv) verstanden sein
  - ... sonst kann man **Relevanz nicht beurteilen**
  - Um was geht es in dem Problem?
  - Für was wird eine Methode verwendet?
- Prozesssicht
  - **Recherche verläuft iterativ**
    - Man sucht, liest, sucht mit anderen Begriffen, liest, folgt Links, ...
  - Management des Prozesses wichtig
    - Verwaltung und **Systematisierung** des Gefundenen
  - Beim späteren Schreiben ergeben sich neue Aspekte
    - Ergänzende Recherchen
  - Ggf. auch **Anpassung des Themas** möglich oder notwendig

# Eng gefasste Themen

- Manches Seminarthema besteht aus einem Paper
- Das heisst nicht, dass man nicht auch andere Paper suchen und lesen muss!
  - Andere Ansätze kennenlernen
  - Übersicht gewinnen
  - Zugeteiltes Paper besser einschätzen (was ist besser / neu / ...)

# Beispiel

- Breit gefasst: Multidimensionale Indexstrukturen
  - Etabliertes Thema; es gibt Bücher und Übersichtsartikel; es gibt 1000+ Originalarbeiten; es gibt andere Seminararbeiten dazu im Web; es gibt Wikipedia Artikel; es gibt kommerzielle Produkte; ...
  - Übersicht gefragt: Was wird indiziert, welche Suchen werden unterstützt, welche grossen Klassen von Ansätzen gibt es, was sind deren Vor- und Nachteile, wo wurden sie schon mal verglichen, ...
- Eng gefasst: kdb-Trees
  - Sehr spezielles Thema; vielleicht eine (kurze) Erwähnung in einem Buch; vielleicht 10-20 wirklich relevante Arbeiten, die man vor allem in Spezialdatenbanken findet; kaum Erwähnung im Web; ...
  - Details sind gefragt: Wie funktioniert Löschen, welche Komplexität haben alle Operationen, wo gibt es empirische Untersuchungen, sind die Abhängig von den Daten, gibt es aktuelle Weiterentwicklungen, ...
- Themenanpassung (MDI)
  - Punkte / Flächen/Körper? Exakte / Ähnlichkeitssuche? Memory / Disk?
  - Was interessiert Sie mehr?

# Inhalt

- Übersicht
- Wissenschaftliche Publikationen
- Bewertung von Publikationen
- Ressourcen und Techniken
- Literatursammlung erstellen und managen



# Arten von Publikationen

- Zentrales Merkmal: **Peer-review or not**
- Peer-reviewed: **Journale, Konferenzen**, Workshops
  - In anderen Fächern ist das anders!
  - Informatikspezifisch: Konferenzen fast wichtiger als Journale
- Nicht: Bücher, Buchkapitel, (Technical) Reports
  - Reports: Früher pro Institut, heute eher internationale, **fachspezifische Repositorien** (arXiv, corr)
- Nicht wissenschaftlich: Blogs, White Paper, Wikipedia
- Sonderfälle
  - Dissertationen, Diplom- und Seminararbeiten (Reports)
  - Poster, Meeting-Abstracts (in der Informatik kaum existent)
  - Vortragsfolien, Vorlesungen, auch Video

# Peer-Review

- „Gute“ oder „neue“ Wissenschaft nicht objektiv definierbar
- Entscheidung soll durch Community getroffen werden
- **Peer-Review**: Paper werden von 2-3 ExpertInnen bewertet
  - Ablehnung, Annahme, Revision
  - Einschätzung der **Relevanz, Neuheit, Qualität** des Textes
  - Unterschiedliche Anforderungen je nach Erscheinungsort
    - Nature: Extrem hoch, insb. Neuheit und Allgemeinrelevanz
    - Kleine Workshops: Eher niedrig, oft sehr spezielle Ergebnisse
- Vielfältige Kritik
  - ExpertInnen die keine sind; schlampiges Lesen; Unterdrückung von Konkurrenz (single-blind); Trend zu Modethemen und inkrementellen Ergebnissen; idR **keine Ergebnisvalidierung**; ...
- Aber bestes bekanntes Verfahren

# WikiPedia

- Zur frühen Recherche sehr nützlich
- Erfahrungsgemäß (in Informatik) meistens korrekt, aber nur **geringe Themenabdeckung oder Tiefe**
- Als Referenz allgemein **nicht** akzeptiert; zum Finden von Referenzen schon

# Übersichtsartikel / Benchmarks

- Zu nahezu allen Themen erscheinen regelmäßig Übersichtsartikel (surveys, reviews)
  - Vor allem als Buchkapitel / in Journalen
- Unterschiedliche Qualität und Tiefe
- **Gute und aktuelle Surveys** sind ein Segen und die halbe Miete für eine Seminararbeit (oder eine BA)
- Meistens sehr lange und nützliche Literaturlisten
- Für praktische Themen sehr relevant: Empirische Übersichts- und **Vergleichsarbeiten** (Benchmarks)

# Lexika

- Lange Zeit bedeutungslos
- In den letzten Jahren aber wieder etwas populärer
  - Encyclopedia of ... Database Systems (Springer): „Comprehensive reference to about 1,400 entries, covering key concepts and terms in the broad field of database systems.“
  - Synthesis Lectures on ... Data Management (Morgan Claypool)
  - ...

# Englisch, Deutsch, ...

- Alles wichtige, aktuelle ist Englisch
- Deutsche Lehrbücher als Einstieg

# Inhalt

- Übersicht
- Wissenschaftliche Publikationen
- **Bewertung von Publikationen**
- Ressourcen und Techniken
- Literatursammlung erstellen und managen

# Impact Factor, Zitationen

- Qualität ist nicht binär
  - Erhebliche Unterschiede auch zwischen peer-reviewed Papern
- Wichtiges Indiz: Erscheinungsort und Autor
  - Welches Journal, welche Konferenz
  - Was keinen eindeutigen Autor hat, zählt überhaupt nicht
- Typische Qualitätsmaße
  - **Impact Factor**: Durchschnittliche Zahl Zitierungen eines Papers in diesem Journal nach 5 Jahren; Wertebereich 0-32
  - Kritik: **Nicht alle Paper eines Journals gleich gut**; Themen mit kleinere Communities haben weniger Zitate; **Fächerunterschiede**; ...
  - Alternative: **Zitationszahlen pro Paper**
    - Google Scholar, Microsoft Academic Search, CiteSeer, Web-of-Science
  - Früher: Verlag



# Wem trauen?

- Zentral: **Zahl der Zitate**
  - Gewichtet nach Jahr – 10 pro Jahr ist ganz ordentlich
  - Mit etwas Exploration ein Gefühl für das Thema bekommen
    - Viele Paper mit vielen Zitierungen?
    - Spezialgebiet mit insgesamt wenig Zitierungen?
  - Bringt einen Bias: „rich get richer“
- Auch wichtig: Erscheinungsort
  - Suche nach „Erschienen in“ und sehen, wie oft Paper typischerweise zitiert werden
  - Vorsicht vor Abkürzungen und Schreibweisen
- Auch wichtig: Autor
  - Verlangt viel Erfahrung, starker Bias gegen Newcomer

# Welchen Zahlen treuen?

The image displays four overlapping browser windows, illustrating a research workflow. The leftmost window shows a Google Scholar profile for Ulf Leser. The second window shows a CiteSeerX search result for 'Querying Distributed...'. The third window shows a PubMed search for 'leser-u[au]' with a list of 71 results. The rightmost window shows a PubMed search for 'leser-u[au]' with a list of 71 results.

**Ulf Leser - Google Scholar**

Meistbesucht Frequent WBI Lehre

Web Images

Querying Distributed...  
B Quilitz, U Leser  
Europe

Quality-driven Integr...  
F Naumann, U Leser  
Humboldt

Fast and practical in...  
S Trißl, U Leser  
Proceedings

The st...  
A Alexa, U Leser  
The VL

Feder...  
S Buss, U Leser  
Forschung

Comp...  
F Naumann, U Leser  
Inform

AliBab...  
C Plake, U Leser  
Bioinform

Inform

**CiteSeerX — Search Results**

citeseerx.ist.psu.edu/search?q=author%3Aleser&sort=cite&t=doc

Meistbesucht Frequent WBI Lehre

Documents Authors

**CiteSeerX**

Results 1 - 10 of 206

[Querying Distributed...](#)  
by Bastian Quilitz, Ulf Leser  
"... Abstract. Integrated access to distributed data is a challenge. As a reaction to this challenge, we have developed a system. However, the current standard is not sufficient. Abstract - Cited by 129 (0 se

[Quality-driven Integr...](#)  
by Felix Naumann, Ulf Leser  
"... Integrated access to information and commercial domains. While the important ... Abstract - Cited by 101 (16 se

[Fast and practical in...](#)  
by Silke Trißl, Ulf Leser - In S...  
"... Many applications work well. We present the queries in graphs. GRIPP req Abstract - Cited by 76 (2 se

[A comprehensive b...](#)  
by Domonkos Tikk, Peter Pal...  
"... The most important way of (PPs) reported in scientific p been prop ... Abstract - Cited by 47 (8 se

[Ultrastructural distri...](#)  
by Stanislav Fakan, George...  
"... ABSTRACT The ultrastru or rat liver, embedded in Low means of a protein ... Abstract - Cited by 44 (2 se

[Federated Informati...](#)  
by Susanne Busse, Ralf-Det...  
context of the DDI Extraction 2

**PubMed.gov**

US National Library of Medicine  
National Institutes of Health

NCBI Resources How To

PubMed

Create RSS Create alert Advanced

Article types  
Clinical Trial  
Review  
Customize ...

Text availability  
Abstract  
Free full text  
Full text

PubMed Commons  
Reader comments  
Trending articles

Publication dates  
5 years  
10 years  
Custom range...

Species  
Humans  
Other Animals

Clear all  
Show additional filters

Format: Summary Sort by: Most Recent Per page: 20 Send to

**Search results**

Items: 1 to 20 of 71

<< First < Prev Page 1 of 4 Next > Last >>

1. [Robust In-Silico identification of cancer cell lines based on next generation sequencing.](#)  
Otto R, Sers C, **Leser U**.  
Oncotarget. 2017 Mar 10. doi: 10.18632/oncotarget.16110. [Epub ahead of print]  
PMID: 28415721 [Free Article](#)  
[Similar articles](#)

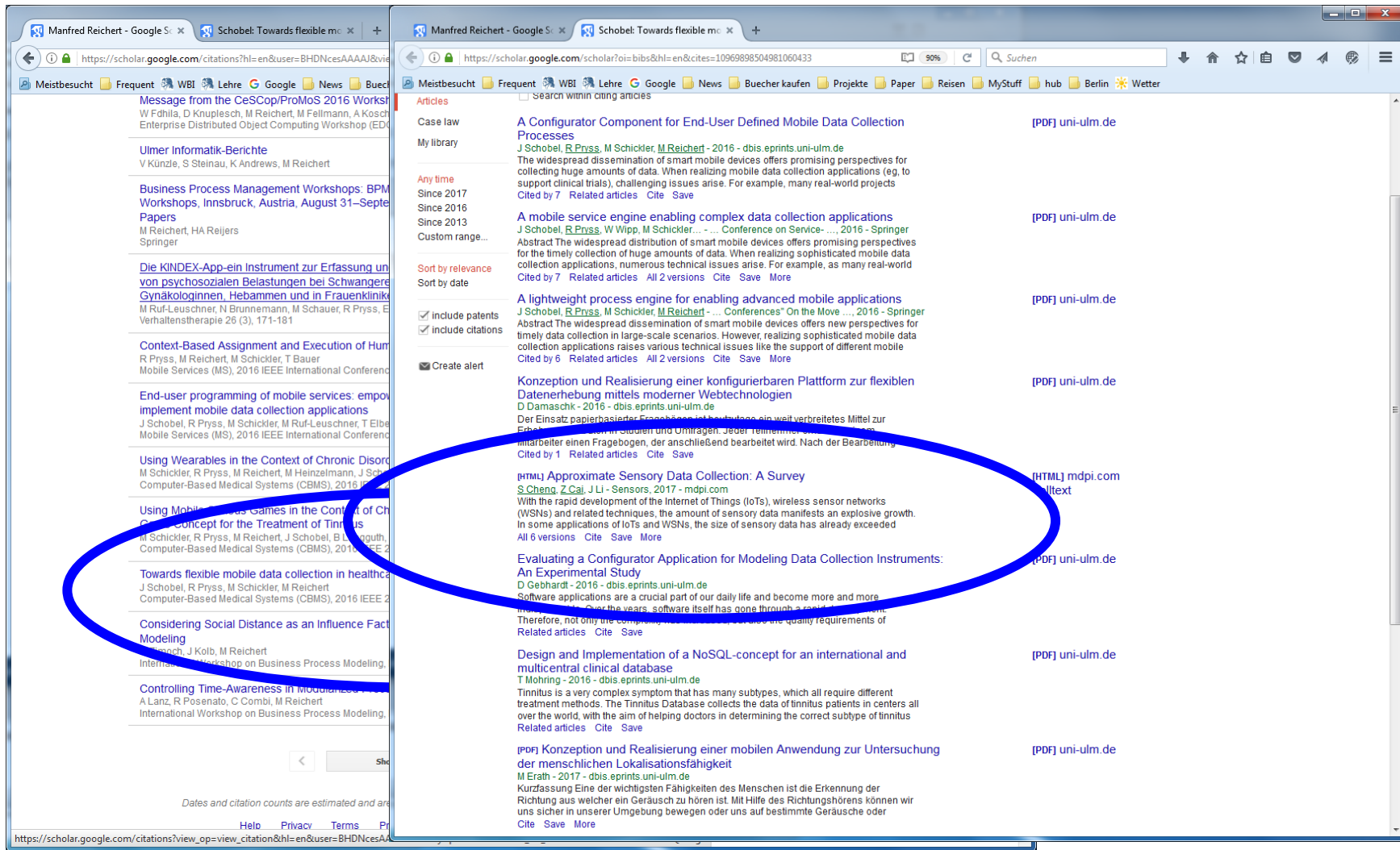
2. [Estimating genome-wide regulatory activity from multi-omics data sets using mathematical optimization.](#)  
Trescher S, Münchmeyer J, **Leser U**.  
BMC Syst Biol. 2017 Mar 27;11(1):41. doi: 10.1186/s12918-017-0419-z.  
PMID: 28347313 [Free PMC Article](#)  
[Similar articles](#)

3. [DNA copy number changes define spatial patterns of heterogeneity in colorectal cancer.](#)  
Mamlouk S, Childs LH, Aust D, Heim D, Melching F, Oliveira C, Wolf T, Durek P, Schumacher D, Blaker H, von Winterfeld M, Gastl B, Mohr K, Menne A, Zeugner S, Redmer T, Lenze D, Tierling S, Möbs M, Weichert W, Folprecht G, Blanc E, Beule D, Schäfer R, Morkel M, Klauschen F, **Leser U**, Sers C.  
Nat Commun. 2017 Jan 25;8:14093. doi: 10.1038/ncomms14093.  
PMID: 28120820 [Free PMC Article](#)  
[Similar articles](#)

4. [Recognizing chemicals in patents: a comparative analysis.](#)  
Habibi M, Wiegand DL, Schmedding F, **Leser U**.  
J Cheminform. 2016 Oct 28;8:59. eCollection 2016.  
PMID: 27843493 [Free PMC Article](#)  
[Similar articles](#)

5. [Comparative assessment of differential network analysis methods.](#)  
Lichtblau Y, Zimmermann K, Haldemann B, Lenze D, Hummel M, **Leser U**.  
Brief Bioinform. 2016 Jul 29. pii: bbw061. [Epub ahead of print]  
PMID: 27473063  
[Similar articles](#)

# Welchen Zahle trauen?



# Spezifika eines Fachs

	<b>Informatik</b>	<b>Life Sciences</b>	<b>Humanities</b>
Wo wird publiziert	Konferenzen, Workshops, Journals, TRs	Nur Journals	Bücher
Wie wird bewertet	Zitierungen, Ort	Impact Factor	Autor
Wo wird gesucht	Google Scholar (ResearchGate?)	Web of Science	Bibliothek
Was ist „viel“	100 pro Jahr	1000 pro Jahr	10 pro Jahr
Autorenlisten	2-5	Oft >20	1
Struktur von Papern	Frei	Fest (Intro, Methods, Results, Diskussion)	Frei

# Open Access

- Idee: Artikel lesen umsonst, publizieren kostet Geld
  - Traditionell: Publizieren ist umsonst, Lesen kostet Geld
- Das ist **unabhängig von Peer-Review**
  - Die meisten Open Access Journale sind peer-reviewed
- Das ist unabhängig von kommerziellen Interessen
  - Die meisten Open Access Journale sind kommerziell orientiert
- Erheblicher politischer Druck zu Open Access
- Informatik: Bisher keine relevanten OA-Journale
  - Autoren stellen traditionell viele Artikel frei zur Verfügung
  - Workshop/Konferenz-Beiträge sind fast immer frei verfügbar
  - Bedeutung von **TR-Repositorien**

# Alter – Wann ist Literatur nicht mehr aktuell?

- So alt ist die Informatik nicht
  - Viele grundlegende Arbeiten aus den 60ziger – 80ziger
- **Theoretische Resultate** (Beweise) halten „ewig“
- Methodische Resultate halten lange
  - Geschäftsmodelle, Best Practise, Algorithmen, Modellierungsansätze, ...
- **Empirische Resultate** können schnell veralten
  - Besser nicht älter als 10 Jahre
  - Hauptspeicher - andere Indexverfahren
  - Cachelines – Scan statt Index
  - Breitbandinternet – Berechnungen verteilen
  - ...

# Inhalt

- Übersicht
- Wissenschaftliche Publikationen
- Bewertung von Publikationen
- Ressourcen und Techniken
- Literatursammlung erstellen und managen

# Recherche-Datenbanken

- Google Scholar / (Microsoft Academic Search) / CiteSeer
  - Volltextindexierung wissenschaftlicher Veröffentlichungen
  - GS: Ziemlich **vollständig** und aktuell
  - Wichtiges Feature: **Zitierende Arbeiten**
  - Voll automatisch; viele Qualitätsmängel
  - Zugriff auf Volltexte über Links
    - Innerhalb der HU weit mehr Content als außerhalb
  - Bibliographische Informationen oft unvollständig
    - Selber nach-recherchieren
- DBLP
  - Manuell gepflegt, Fokus auf DB+LP
  - Unvollständig, keine Volltexte, Links auf PDFs
  - Sehr gut bzgl. **bibliographischer Daten**



# Weitere Datenbanken

- ResearchGate
  - Closed, eventuell mehr PDF, Zugriffsregeln unklar, Zitierungszahlen unzuverlässig, manuell kuriiert (Autoren), dubiose „Scores“ für Autoren
- PubAnnotator, Mendely, ...
- ACM-DL, IEEE-Explore, Springer Link, ...
- Web-of Science, Scopus, ...
- „First pirate ...“ [sci-hub.io](http://sci-hub.io)

# Bibliotheken

- Zur Recherche wenig nutzbringend
- Sehr gut als [Arbeitsort](#)
- Sehr wichtig zur Organisation des Zugriffs auf [lizensierten Content](#)
  - eBooks, online-subscriptions, ...

# Recherche Techniken

- Suchfeature benutzen
  - Keywords, Phrasen, Negation
  - Suche nach Autoren
  - Suche nach zitierenden Arbeiten
  - Suche nach Jahren
  - [Suche nach Erscheinungsort]
- Auch wichtig: Wichtige Referenzen in guten Papern

# Gestufte Bewertung

- Man muss **vieles ansehen**, aber das meisten **nur kurz**
- Erste Einschätzung: Titel, Ort, Zitierungen (C)
  - Journal oder Konferenz; bei Zahl Zitierungen das Jahr beachten
- Zweite Einschätzung: **Abstract** (C)
  - Ggf schnell aufhören, wenn erkennbar irrelevant
- Dritte Einschätzung: **Paper querlesen** (B)
  - Welches Problem, welcher Ansatz, welche Daten, welche Resultate?
  - Algorithmen, Beweise, Formeln überspringen
  - Paper mit Stichworten **in Liste aufnehmen** („Cite“ – kopieren)
- Vierte Einschätzung: **Paper genau lesen** (A)
  - Eigene Zusammenfassung **in Liste aufnehmen**

# Iterativer Prozess

- Explorationsphase: Liste B-Paper erstellen
  - Keywordsuche, Surveys suchen, normierte Zitationszahlen beachten
  - Man findet die „seminal paper“ und erhält Übersicht
  - Paper eher schnell beurteilen, erstmal sammeln
- Komplettierungsphase: Die A-Paper finden
  - Wichtigste Paper aus der B-Gruppe
  - Auch gezielte weitere Suchen
    - Aktuelle Arbeiten, die A/B-Paper zitieren
    - Wichtige Arbeiten aus Referenzlisten von A/B-Papern
  - Führt zu weiteren A und B Papern (ständig bewerten)
- Verwendungsphase
  - Kann zu gezielten weiteren Recherchen führen
  - Aber nicht verlieren!  $|A|=10$  ist viel; ggf. Betreuer fragen

# Life-Beispiel: SIMD in Datenbanken

- Autorensuche
- Erscheinungsjahr
- Zitierende Arbeiten
- Versionen
- Bibliographische Informationen
- [Related Articles, Web of Science]

---

## Data broadcasting in SIMD computers

D Nassimi, S Sahni - Computers, IEEE Transactions on, 1981 - [ieeexplore.ieee.org](http://ieeexplore.ieee.org)

... The PE's are **indexed** 0 through TV — 1 and may be referenced as PE ...  $2^k$ -block independent of the remaining  $2^k$ -blocks (ie, as if each  $2^k$ -block defined a separate SIMD computer ... Complexity of Rank: 1) MCCs with row-major **indexing**: Let  $TV = n \cdot k = 2P$  be the number of PE's in a ...

Cited by 324 [Related articles](#) All 3 versions [Web of Science](#): 134 [Cite](#) [Save](#) [More](#)

[\[PDF\] from computer.org](#)  
[Volltext](#)

---

## Implementing database operations using SIMD instructions

J Zhou, KA Ross - Proceedings of the 2002 ACM SIGMOD international ..., 2002 - [dl.acm.org](http://dl.acm.org)

... When one returns all matches, we write the results to an array called result, **indexed** by a ... resident **indexes**, as well as to the layout of a disk page within disk-based **indexes**. In this section, we study techniques for employing SIMD instructions to make index traversal more efficient ...

Cited by 142 [Related articles](#) All 11 versions [Cite](#) [Save](#) [More](#)

[\[PDF\] from columbia.edu](#)

---

## SIMD-scan: ultra fast in-memory table scan using on-chip vector processing units

T Willhalm, N Popovici, Y Boshmaf, H Plattner... - Proceedings of the ..., 2009 - [dl.acm.org](http://dl.acm.org)

... In this paper, we assume 64-bit little-endian architecture with **indexing** starting at and increasing by ... For example, if the compressed data represents an array of ascending integers (**indexes**) starting from "0d" till ... the index array would be "1d, 2d, 3d, 4d, 5d" for a 0-**indexed** array. ...

Cited by 99 [Related articles](#) All 5 versions [Cite](#) [Save](#)

[\[PDF\] from ubc.ca](#)  
[Volltext](#)

# Inhalt

- Übersicht
- Wissenschaftliche Publikationen
- Bewertung von Publikationen
- Ressourcen und Techniken
- [Literaturliste managen](#)

# Aufbereiten einer Textsammlung

- **Notizen und Zusammenfassungen** zu allen A und B Papern machen und **verwalten**
- Bibliographische Informationen gleich in hoher Qualität sammeln
- **Strukturieren**: Nach Problemvarianten, nach verwendeten Datensätzen, nach Architektur, nach Denkschule, ...
- Individuelle **Schlagworthierarchie** sinnvoll
  - Und notwendig bei Dissertationen
- Spezielle Software benutzen (später)



# Exzerpte erstellen

- **Eigene Worte** verwenden, nicht Abstract kopieren
- Texte drucken und **markieren**, dann zusammenfassen
- **Systematik erstellen** (Farben) – Kernaussage, Kritikpunkt, experimentelle Ergebnisse, ...
- A Paper ausführlich beschreiben, B Paper nur kurz
- Gute Zusammenfassungen kann man ggf. in die Seminararbeit übernehmen
  - Meist sind die aber zu kurz

# Zum Zitieren

- Alle: Autoren, Titel, Jahr
- Konferenzen/Workshops: Name der Konferenz, Ort
- Journal: Ausgabe (volume, issue), Seitenzahl
- Report (auch Dissertation, DA ...): Institution, Nummer
  
- Möglich: DOI
- Nicht in Referenzliste: URL
- Nützlich zur Suche: Abstract
- Nützlich zum Nachsehen: PDF

# Management von Literatur

- Spezielle Software zum
  - Datenbank für Literatur (Bibliographie, Link, Verschlagwortung, ...)
  - Ggf mit PDF Verwaltung und Volltextindexen
  - Ggf. automatische Extraktion bibliographischer Daten aus Webseiten
  - Formatierung von Referenzen in Texten nach Vorlagen
    - Word Plug-In, BibRef
  - Sortieren, suchen, filtern
  - Zugriff auf Datenbanken (PubMed)
- Endnote (Campuslizenz), Bibliographix, Citavi, jab-ref, ...
- Online-Systeme: Mendely, zetero, ...

# Live-Beispiel

- Suchen
- Zugriff PubMed
- Word-Formatierung
- Referenzstile
- Formatierter Export
- BibTex-Export

