



Autoren-Disambiguierung in bibliographischen Metadaten

Der Fall dblp

Marcel R. Ackermann

4. April 2016



Indexierung von Publikationen

home | browse | search | about



metadata

[+] Search dblp

> Home

Dagstuhl

[-] Venue results

Likely matches

- International Journal of **Metadata**, Semantics and Ontologies
- International Conference on **Metadata** and Semantics Research (MTSR)
- Metadata** Management in Grid and P2P Systems (MMGPS)

[-] Publication results

found 3,375 matches

2016

- Gregory Giuliani, Yaniss Guigoz, Pierre Lacroix, Nicolas Ray, Anthony Lehmann:
Facilitating the production of ISO-compliant metadata of geospatial datasets. Int. J. Applied Earth Observation and Geoinformation 44: 239-243 (2016)
- Evangelos Pafilis, Pier Luigi Buttigieg, Barbra Ferrell, Emiliano Pereira, Julia Schnetzer, Christos

[-] Refine list

refine by author

- Erik Duval (16)
- Stuart Weibel (15)
- Demetrios G. Sampson (15)
- Miguel-Ángel Sicilia (14)
- Nikos Manouselis (13)
- Wolfgang Nejdl (13)
- Rik Van de Walle (13)
- C. Lee Giles (12)
- Panos Balatsoukas (12)
- Yasushi Kiyoki (12)
- 7,359 more options

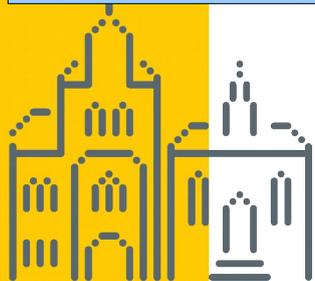
refine by venue

(54)

Was ist dblp?

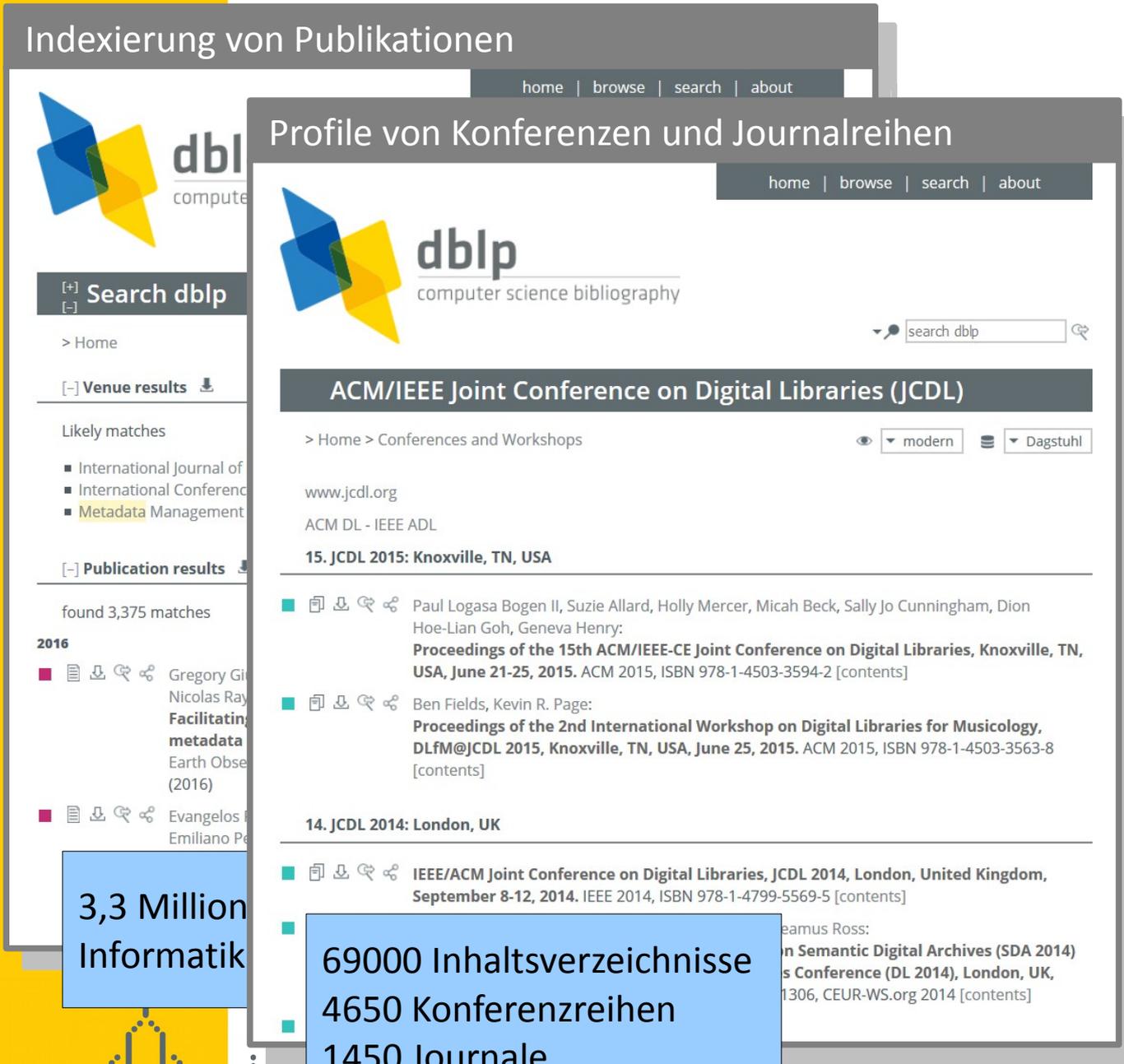


3,3 Millionen Publikationen aus der Informatik und Nachbardisziplinen



Indexierung von Publikationen

Was ist dblp?

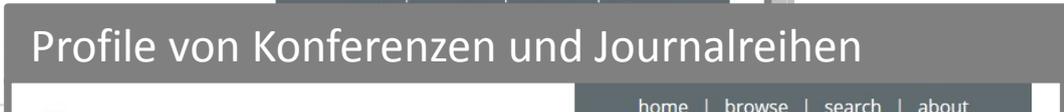


The screenshot shows the dblp website interface. At the top, there is a navigation bar with 'home | browse | search | about'. Below this is the dblp logo and the text 'computer science bibliography'. A search bar is visible with the text 'search dblp'. The main content area is titled 'Profile von Konferenzen und Journalreihen' and shows details for the 'ACM/IEEE Joint Conference on Digital Libraries (JCDL)'. It includes a breadcrumb trail '> Home > Conferences and Workshops', a search filter for 'modern', and a 'Dagstuhl' button. The conference details for '15. JCDL 2015: Knoxville, TN, USA' are listed, including authors like Paul Logasa Bogen II, Suzie Allard, Holly Mercer, Micah Beck, Sally Jo Cunningham, Dion Hoe-Lian Goh, and Geneva Henry. The proceedings title is 'Proceedings of the 15th ACM/IEEE-CE Joint Conference on Digital Libraries, Knoxville, TN, USA, June 21-25, 2015'. Below this, the details for '14. JCDL 2014: London, UK' are shown, including the title 'IEEE/ACM Joint Conference on Digital Libraries, JCDL 2014, London, United Kingdom, September 8-12, 2014' and authors like Thomas D. S. Rasmussen and Thomas D. S. Rasmussen.

3,3 Million
Informatik

69000 Inhaltsverzeichnisse
4650 Konferenzreihen
1450 Journale





[+] Search dblp
[-]

> Home

[-] Venue results

Likely matches

- International Journal of
- International Conferenc
- Metadata Management

[-] Publication results

found 3,375 matches

2016

- Gregory G...
Nicolas Ray...
Facilitating
metadata
Earth Obser
(2016)
- Evangelos I...
Emiliano Pe...

ACM/IEEE Joint

> Home > Conferences and J

www.jcdl.org

ACM DL - IEEE ADL

15. JCDL 2015: Knoxville, TN

[+] Kai Eckert 0001

[-]

> Home > Persons

by year | Dagstuhl

[-] Person information

- affiliation: Hochschule der Medien, Stuttgart, Germany
- affiliation: University of Mannheim, Institute of Computer Science and Business Informatics

[-] Other persons with the same name

- Kai Eckert
- Kai Eckert 0002 — Autonomous University of Barcelona, Theoretical Physics Group

[-] Other persons with a similar name

- Kai-Helmut Eckert

[-] 2010 - today

Refine list

2015

- [c21] Thomas Bosch, Erman Acar, Andreas Nolle, Kai Eckert:
The role of reasoning for RDF validation.
SEMANTICS 2015: 33-40
- [13] Thomas Bosch, Andreas Nolle, Erman Acar,

showing all 27 records

refine by search term

refine by type

- Journal Articles (only)
- Conference and Workshop Papers (only)
- Books and Book Chapters (only)

3,3 Million Informatik

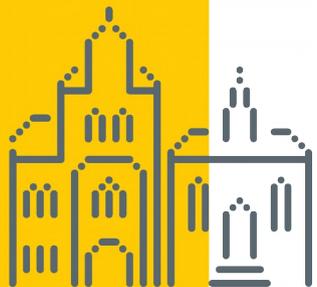
69000 Inhalte
4650 Konferenzen
1450 Journale

1,7 Millionen Autoren und Editoren
Homonym/Synonym-Disambiguierung



Grundsätze

- Offenheit
- Neutralität der Daten
- **Datenqualität**
- Semantische Anreicherung
- Ermöglichen von
Forschung
- Nutzerorientierung



Grundsätze

- Offenheit
- Neutralität der Daten
- **Datenqualität**
- Semantische Anreicherung
- Ermöglichen von Forschung
- Nutzerorientierung

Was dblp nicht ist ...

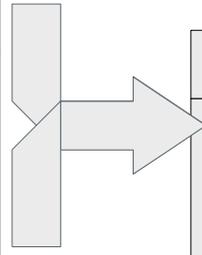
- keine Digital Library
- kein bibliothekarisches Projekt
- kein fachübergreifender Dienst
- „kein Geschäftsmodell“
- kein großes Team



Der dblp-Workflow

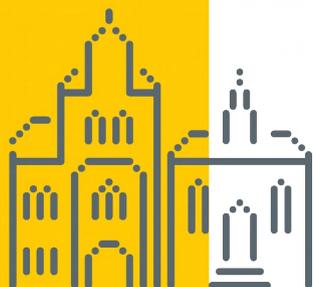
1. Datenakquise

- Webcrawler oder Datenlieferungen
- Normalisierung



Internes Datenformat (XML)

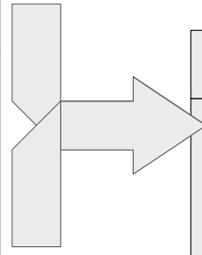
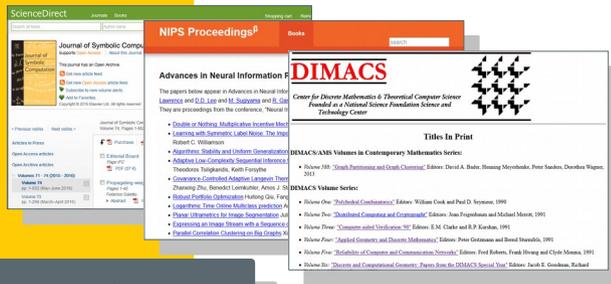
```
<article key="journals/jsc/Kemper16">
<author>Gregor Kemper</author>
<title>Using extended Derksen ideals in computational
invariant theory.</title>
<pages>161-181</pages>
<year>2016</year>
<volume>72</volume>
<journal>J. Symb. Comput.</journal>
<ee>http://dx.doi.org/10.1016/j.jsc.2015.02.004</ee>
<url>db/journals/jsc/jsc72.html#Kemper16</url>
</article>
```



Der dblp-Workflow

1. Datenakquise

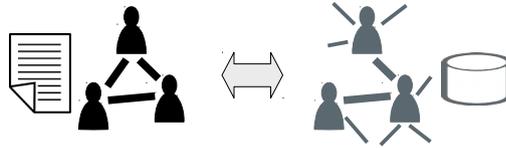
- Webcrawler oder Datenlieferungen
- Normalisierung



Internes D

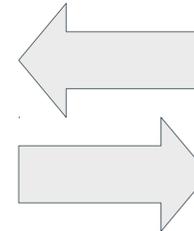
```
<article key="jour
<author>Gregor Kem
<title>Using exten
invariant theory.<
<pages>161-181</pa
<year>2016</year>
<volume>72</volume
<journal>J. Symb.
<ee>http://dx.doi.
<url>db/journals/j
</article>
```

2. Initiale Datensäuberung



- Zuordnen der Autorenschaft
- Homonym/Synonym-Erkennung
- Korrektur der Neudaten
- Vervollständigungen (zB: abgek. Namen, fehlende Namensteile)
- Anreicherung (Homepages, Affiliations, PIDs, zB ORCIDs, ...)

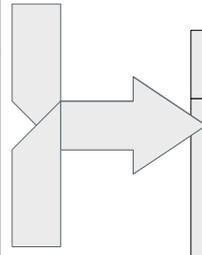
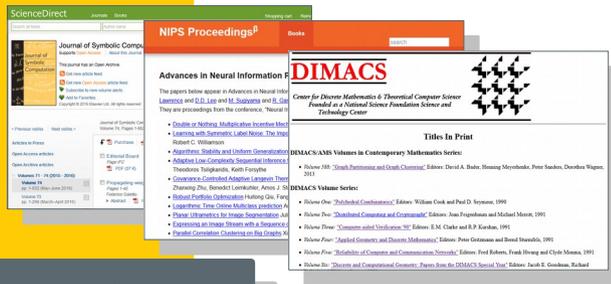
eigener Korpus als Authority



Der dblp-Workflow

1. Datenakquise

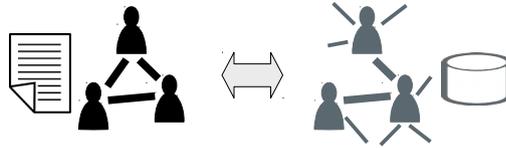
- Webcrawler oder Datenlieferungen
- Normalisierung



Internes D

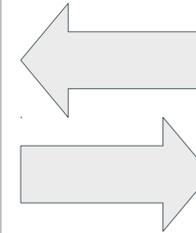
```
<article key="jour
<author>Gregor Kem
<title>Using exten
invariant theory.<
<pages>161-181</pa
<year>2016</year>
<volume>72</volume
<journal>J. Symb.
<ee>http://dx.doi.
<url>db/journals/j
</article>
```

2. Initiale Datensäuberung

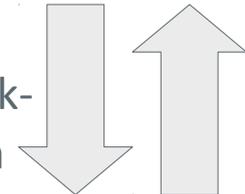


- Zuordnen der Autorenschaft
- Homonym/Synonym-Erkennung
- Korrektur der Neudaten
- Vervollständigungen (zB: abgek. Namen, fehlende Namensteile)
- Anreicherung (Homepages, Affiliations, PIDs, zB ORCIDs, ...)

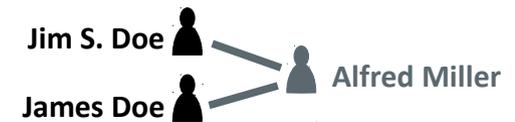
eigener Korpus als Authority



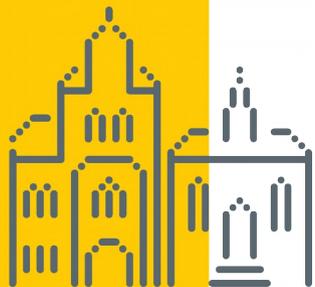
Netzwerk-
analysen



3. Kontinuierliche Qualitätskontrolle



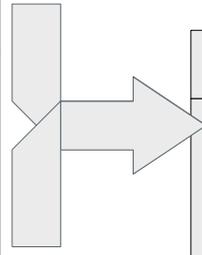
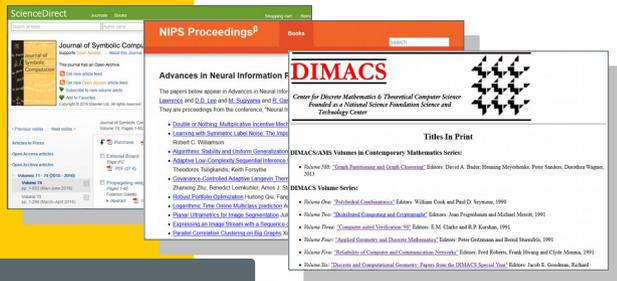
- Homonym/Synonym-Erkennung
- Korrekturen am bestehenden Corpus
- neue Informationen durch neuste Ergänzungen



Der dblp-Workflow

1. Datenakquise

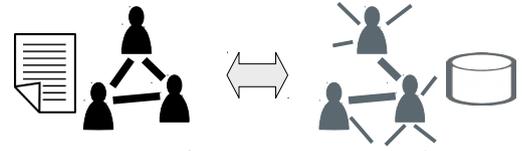
- Webcrawler oder Datenlieferungen
- Normalisierung



Internes D

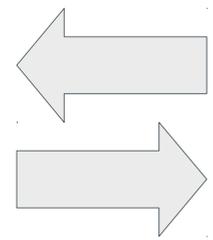
```
<article key="jour
<author>Gregor Kem
<title>Using exten
invariant theory.<
<pages>161-181</pa
<year>2016</year>
<volume>72</volume
<journal>J. Symb.
<ee>http://dx.doi.
<url>db/journals/j
</article>
```

2. Initiale Datensäuberung

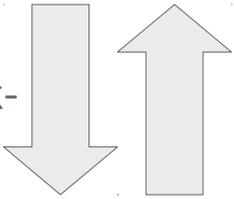


- Zuordnen der Autorenschaft
- Homonym/Synonym-Erkennung
- Korrektur der Neudaten
- Vervollständigungen (zB: abgek. Namen, fehlende Namensteile)
- Anreicherung (Homepages, Affiliations, PIDs, zB ORCIDs, ...)

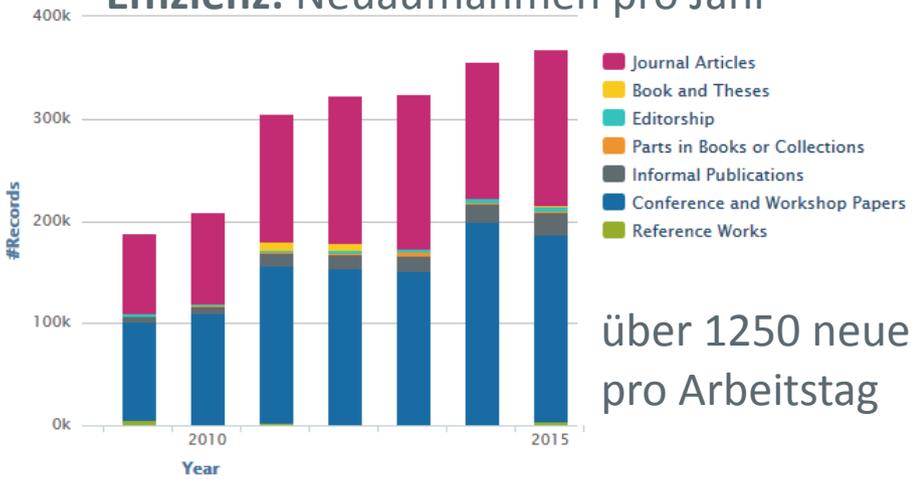
eigener Korpus als Authority



Netzwerk-
analysen

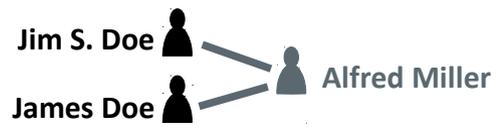


Effizienz: Neuaufnahmen pro Jahr



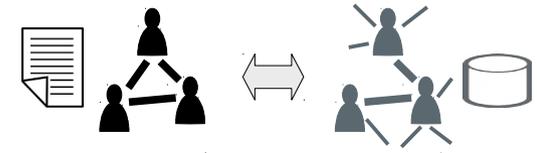
über 1250 neue Artikel pro Arbeitstag

3. Kontinuierliche Qualitätskontrolle



- Homonym/Synonym-Erkennung
- Korrekturen am bestehenden Corpus
- neue Informationen durch neuste Ergänzungen

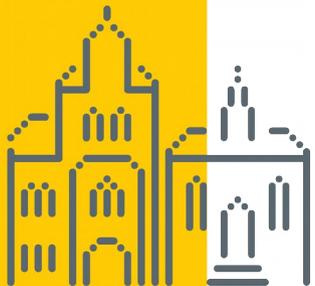




Disclaimer:

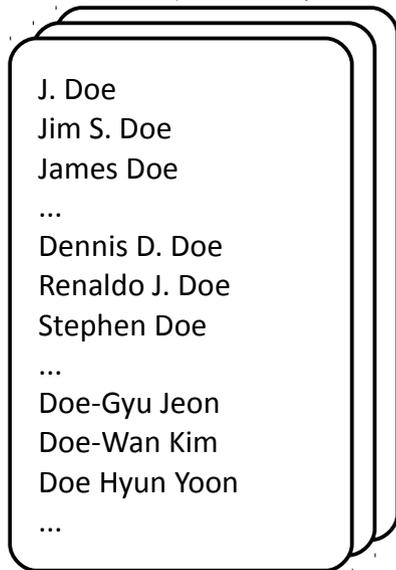
Your mileage may vary!

- optimiert für die Informatik
- basiert nur auf Kern-Metadaten



1. Initiale Kandidatenmenge erzeugen

- Suche auf gesamtem Korpus (rein syntaktisch)
- benötigt sehr schnelle Ähnlichkeitsmaße und Datenstrukturen

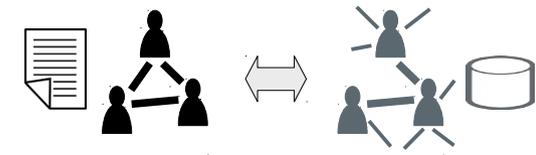


Simple Blocking-Strategien, z.B.:

- letzter Namensteil
- längster Namensteil
- seltenster Namensteil
- vorkommende Konsonanten

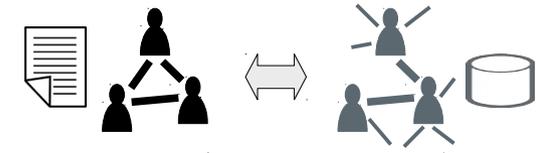
alle Vergleiche unter Unifikation
von diakritischen Zeichen

Algorithmus *dblpi*



1. Initiale Kandidatenmenge erzeugen

Algorithmus *dblpi*



2. Zu große Kandidatenmengen adaptiv filtern

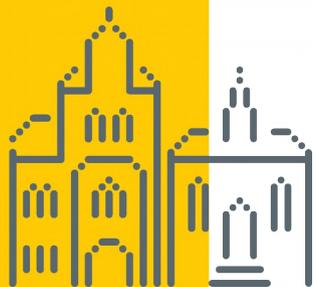
- zu groß: mehr als 10 Kandidaten (z.B. „W. Wang“ \approx 15000)
- Filter darf teurere Maße verwenden als initiale Erstellung



J. Doe
Jim S. Doe
James Doe
...
~~Dennis D. Doe~~
Renaldo J. Doe
~~Stephen Doe~~
...
Doe-Gyu Jeon
~~Doe Wan Kim~~
~~Doe Hyun Yoon~~
...

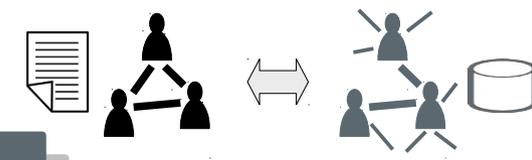
Mögliche Filterregel:

- alle Namensteile als Präfix enthalten
- in eine oder andere Richtung
- ggf. permutierte Reihenfolge



1. Initiale Kandidatenmenge erzeugen

Algorithmus *dblpi*



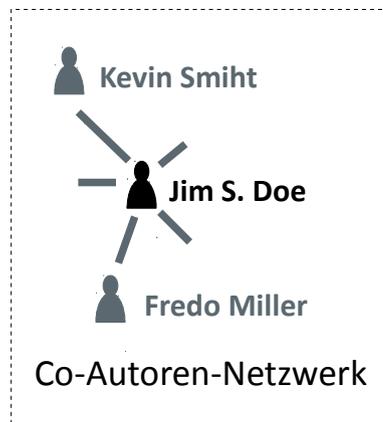
2. Zu große Kandidatenmengen adaptiv filtern

3. Weitere Kandidaten im Kontext finden

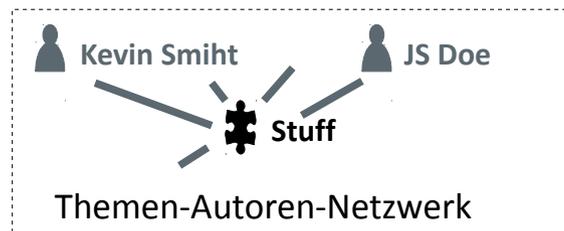
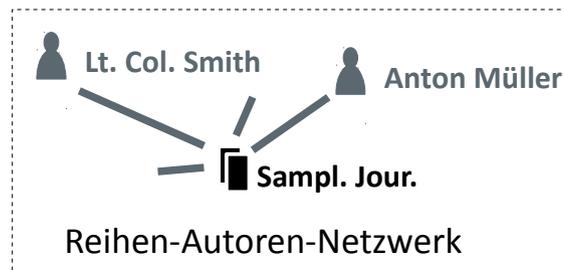
- Ausnutzen von semantischen Zusammenhängen
- darf teurere Ähnlichkeitsmaße verwenden

J. Doe, A. Miller, K. Smith: „Important Stuff.“ Sampl. Jour., 1975

Kontext der Kandidatenmenge:



weitere Kontexte, z.B.:



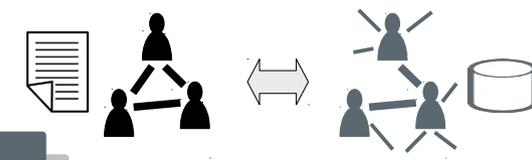
Mögliche Maße:

- Levenshtein
- Jaro-Winkler
- n-Gramme
- Soundex-Varianten mit Anpassungen an Personennamen



1. Initiale Kandidatenmenge erzeugen

Algorithmus *dblpi*



2. Zu große Kandidatenmengen adaptiv filtern

3. Weitere Kandidaten im Kontext finden

4. Kandidaten bewerten



Kontext
Kandidat



Co-A

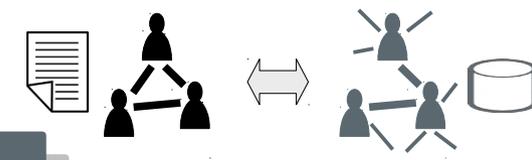
- Bewertung anhand mehrdimensionaler Features
- Werte je zwischen 0 und 1
- ähnlicher Name
- gemeinsame Publikation
- zeitnahe Publikationen
- gleiche Affiliation
- publiziert in gleicher Reihe
- publiziert in *ähnlicher* Reihe
- verwendet gleiche Themen
- ggf. weitere ...
- Ranking mittels gewichteter Summe

Themen-Autoren-Netzwerk



1. Initiale Kandidatenmenge erzeugen

Algorithmus *dblpi*



2. Zu große Kandidatenmengen adaptiv filtern

3. Weitere Kandidaten im Kontext finden

4. Kandidaten bewerten

5. Autoren zuordnen

- Ausnahmefälle
- darf nicht

- Bewertung
- Werte
- Ähnlichkeit
- Gemeinsamkeiten
- zeitliche Nähe
- gleiche
- Ranking
- Co-Autoren

Develop - dblpi interface

Load new data 5 / 18

©2019 De Koster, Tom Van Cutsem, Theo D'Hondt:
 Domains safe sharing among actors.
 11-22
 http://dx.doi.org/10.1145/2414639.2414644#view

authors: Ehsan khamespaneh **Ehsan Khamespanali** * 2a 3a 4a 5a
 Zeynab Sabahi-Kavian **Zeynab Sabahi-Kavian** * 1a 3a 4a 5a
 Ramtin Khorzavi **Ramtin Khorzavi** * 1a 2a 4a 5a
 Marjan Sirjani **Marjan Sirjani** * 1a 2a 3a 5a Behzod Sirjani * Abolfazl Sirjani Arash Shahabi-Sirjani
 Mohammad-Javad Izadi **Mohammad-Javad Izadi** * 1a 2a 3a 4a Mohammad Izadi * M. Hamed Izadi * Mohammad Hadi Izadi
 M. Izadi M. A. Izadi

Ask Yahoo: [Names] [Title] [Names & Title]
 Ask Google: [Names] [Title] [Names & Title] (e) | Google Scholar: [Names] [Title] [Names & Title] (e) | DeepDyve

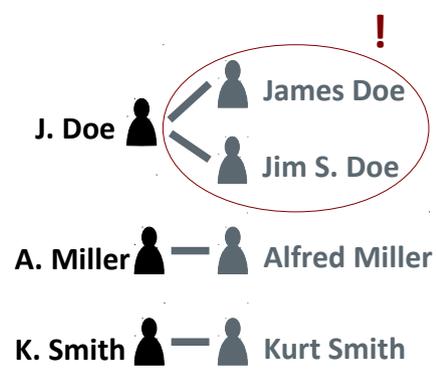
title: Titled-reduce schedulability and deadline-freedom analysis using floating-time transition system.

pages: 23-34

ee: <http://doi.acm.org/10.1145/2414639.2414645> visit

©2019 Thierry Renaux, Ludic Hoche, Stefan Mann, Wolfgang De Maess:
 Digital genome recognition with soft real-time guarantees.
 05-06
 http://dx.doi.org/10.1145/2414639.2414644#view

©2019 Ham W. Kurnia, Arif Poespoh-Herfan:
 A regional trace logic for simple hierarchical actor-based component systems.
 07-08
 http://dx.doi.org/10.1145/2414639.2414647#view

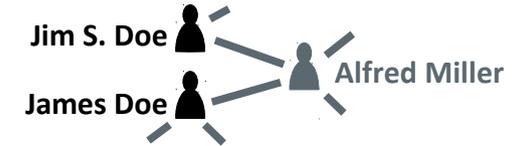
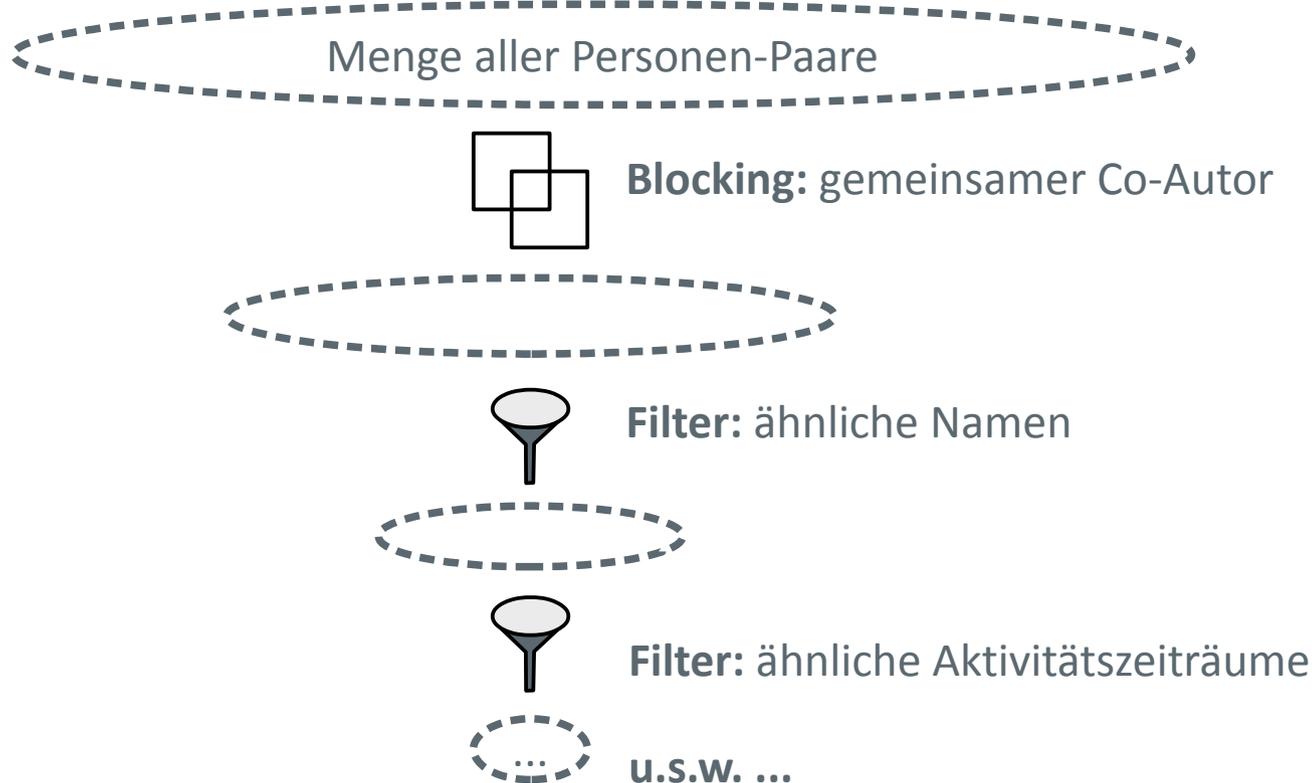


- manuelle Zuordnung anhand der Kandidatenlisten
- im Zweifel: Recherche mit externen Quellen
- Fehler im Korpus werden durch die Vorschläge oft sichtbar!



Kontinuierliche Qualitätskontrolle

Synonyme entdecken



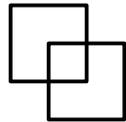
- Ranking der Vorschlagslisten & manuelle Kontrolle



Kontinuierliche Qualitätskontrolle

Synonyme entdecken

Menge aller Personen-Paare



Blöcke



Filter

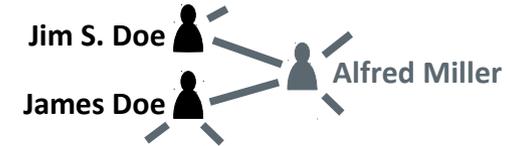


Filter



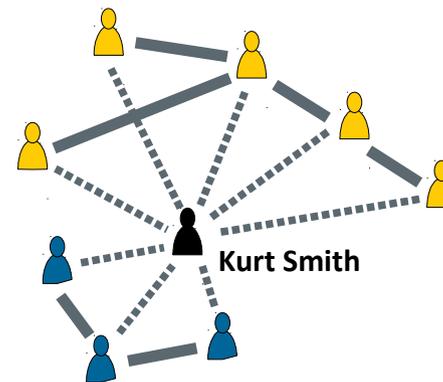
u.s.

- Ranking der Vorschlagslisten & ma



Homonyme erkennen

- Ansatz: Netzwerkanalyse von Co-Autor-Communities
- Auch andere semantische Netzwerke nutzbar



- Derzeit noch kein automatisiertes Verfahren im Einsatz
- Erforschung solcher Teil von laufendem Drittmittel-Projekt



Vielen Dank!

<http://dblp.dagstuhl.de>

<http://dblp.uni-trier.de>

<http://dblp.org>

✉ dblp@dagstuhl.de

🐦 [@dblp_org](https://twitter.com/dblp_org)

