

Automatisierte Verfahren zur Datenanreicherung

Angela Vorndran

Team Datenmanagement, Deutsche Nationalbibliothek

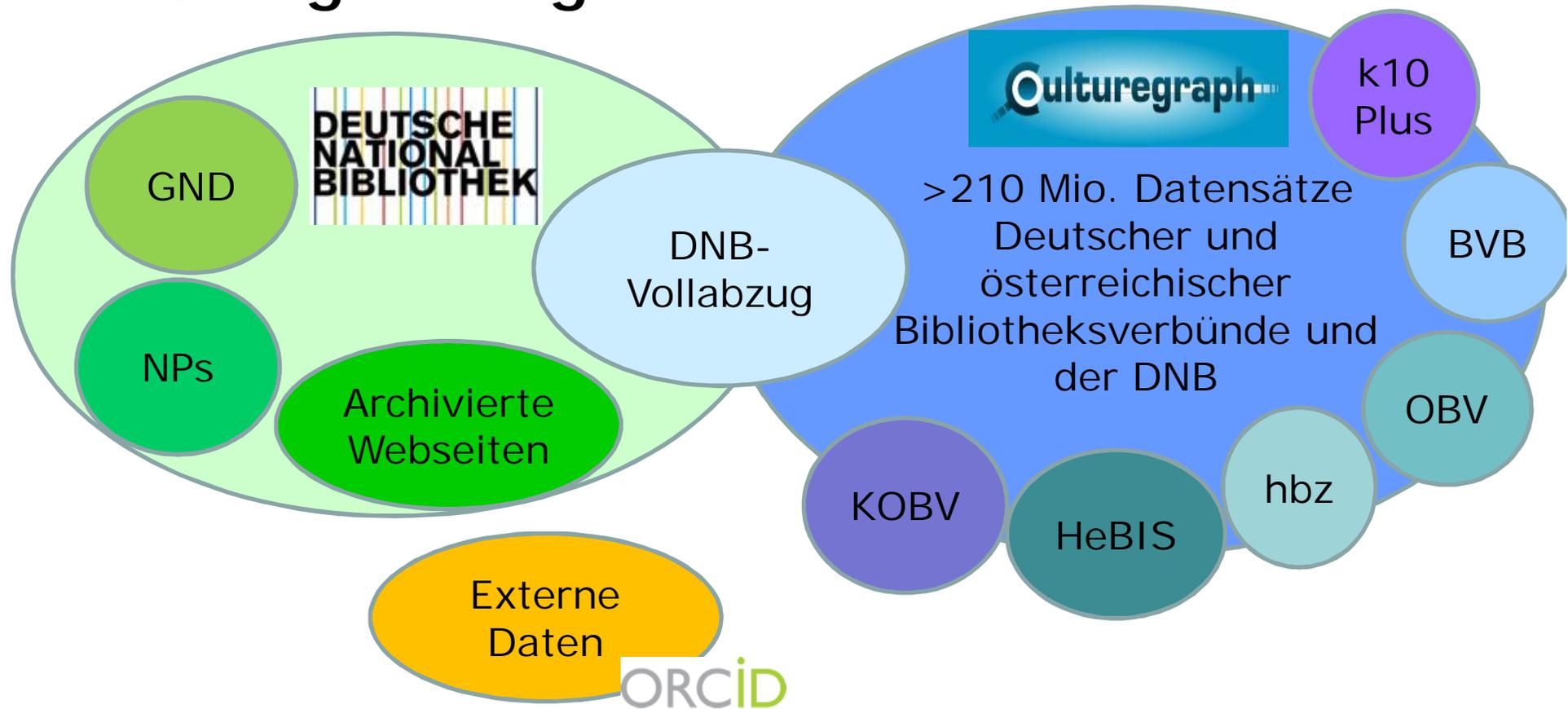
Hintergrund

- Im Bestand der DNB befinden sich über 43 Millionen Medien und mittlerweile mehr als 10 Millionen Netzpublikationen
- Im Schnitt kommen jeden Tag 5.800 neue NPs dazu
- Maschinelle Lösungen sind unvermeidbar, um die Metadaten v.a. der Netzpublikationen möglichst auch mit GND-Verknüpfungen zu erschließen und ihre Qualität damit zu verbessern
- Ein Lösungsansatz: Fremddatenübernahme

Gliederung

1. Werkbündelung
2. Personen mit Normdatenverknüpfungen anreichern
3. Datenabgleich mit externen Informationsquellen
zur Datenanreicherung, hier: ORCID

Datengrundlagen



Werkbündelung



Datengrundlage

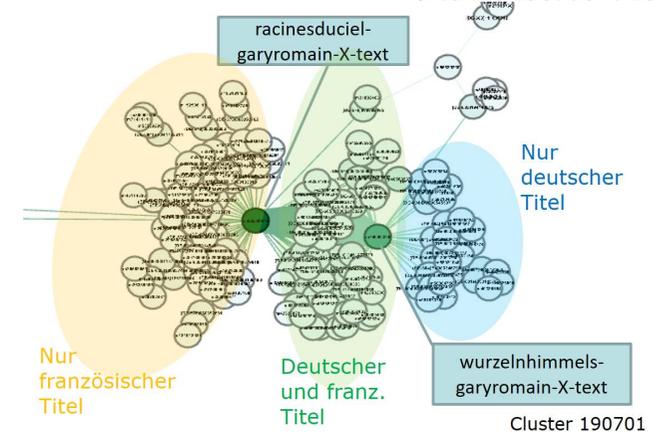
Für jeden Datensatz
Schlüssel für den
Abgleich erstellen



oclc961837268-X-meeralsversprechen-text

9783525253243-X-meeralsversprechen-text

meeralsversprechenbedeutungfunktionvonseehersch
aftbeithukydidess-koppfans-X-text



Breitensuche vereinigt Datensätze
mit mind. einem identischen
Schlüsseln in einem Bündel

Titel im Culturegraph-Bestand abgleichen

- Abgleich der Publikationen durch Schlüssel
- Nach definierten Regeln werden relevante Datenelemente aus dem MARC-Format extrahiert und kombiniert
- Datensätze mit gleichem Schlüssel werden zu Bündeln zusammengefasst
- Potenzial zur Übertragung von
 - Sacherschließungsdaten
 - Anreicherung durch Normdatenverknüpfungenvon einem Bündelmitglied auf ein anderes, das darüber nicht verfügt

Hauptsachtitel	Das Meer Als Versprechen
Zusatz	Bedeutung Und Funktion Von Seeherrschaft Bei Thukydides
Inhalt	<u>Inhaltsverzeichnis</u> <u>Inhaltstext</u> <u>Inhaltstext</u>
Person	aut Kopp, Hans 1127785095
Körperschaft	pbl Vandenhoeck & Ruprecht
Umfang	303 Seiten
Erscheinungsjahr	2017
Material	book
Erscheinungsort	Göttingen
Herausgeber	Vandenhoeck & Ruprecht
Schlagwort	655 Hochschulschrift 4113937-9 650 Seeherrschaft 4077302-4 600 Thucydides 4138063-0
Klassifikation	DDC 355 DDC 880 DDC 930 936.05
Standard-Identifizier	OCLC 961837268 EKI DNB1115738097
Verlags-Identifizier	GTIN14 9783525253243 ISBN13 9783525253243

Titel+Titelzusatz-ErstellerIn-Band-Publikationstyp:
meeralsversprechenbedeutungfunktionvonseeherrschaftbeithukydidess-kopp-hans-X-text



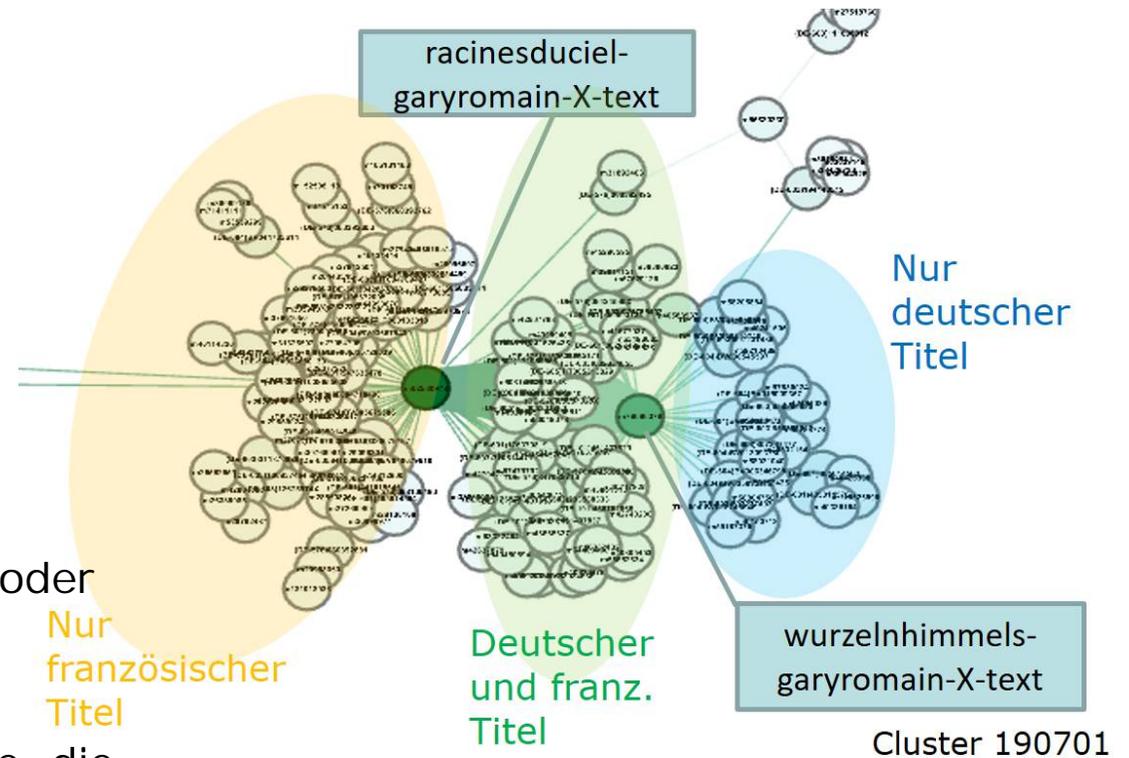
Kontrollnummer-Band-Titel-Publikationstyp:
oclc961837268-X-meeralsversprechen-text



Kontrollnummer-Band-Titel-Publikationstyp:
DNB1115738097-X-meeralsversprechen-text

ISBN-Band-Titel-Publikationstyp:
9783525253243-X-meeralsversprechen-text

Werkbündel



- Für jeden Datensatz liegen ein oder mehrere Schlüssel vor
- Bündel vereinen alle Datensätze, die mindestens einen gemeinsamen Schlüssel mit einem anderen Mitglied besitzen

Statistik Werkbündel (Stand 6/2022)

- 132.880.199 Datensätze wurden verarbeitet (nicht behandelt: unselbständige Teile von mehrbändigen Werken, fortlaufende Ressourcen etc.)
- 21.295.838 Bündel mit mehr als einem Mitglied
- Darin enthalten 101.948.360 Datensätze
- Durchschnittliche Bündelgröße (wenn mehr als ein Mitglied):
4,8

Überblick	
Titel	Nachlaß Und Erbe Im Steuerrecht
Zusatz	Handbuch Zur Steuerlichen Abwicklung Des Erbfalls
Person	aut GND Max Troll
Umfang	XVI, 370 S.
Erscheinungsjahr	1978
Schlagwort	Inheritance and transfer tax GND Nachlass GND Steuerrecht GND Erbe
Klassifikation	DDC 343/.43/053 RVK PP 5345 RVK QL 500

Überblick	
Titel	NACHLASS UND ERBE IM STEUERRECHT. HANDBUCH ZUR STEUERLICHEN ABWICKLUNG DES ERBFALLS. NACHTRAG NACH D. STANDE VOM 1. JANUAR 1970. VON MAX TROLL. MUENCHEN (U.A.): BECK 1967-1970. 27, XII, 364 S.
Person	aut GND Max Troll
Erscheinungsjahr	1967

Überblick	
Titel	Nachlass Und Erbe Im Steuerrecht
Person	aut GND Max Troll
Erscheinungsjahr	1967

Überblick	
Titel	Nachlass Und Erbe Im Steuerrecht
Zusatz	Handbuch Zur Steuerl. Abwicklung D. Erbfalls
Person	aut GND Max Troll
Umfang	XVI, 370 S.
Erscheinungsjahr	1978
Schlagwort	Erbfall Steuerrecht Steuer Nachlaß Erbrecht

Mögliche Anreicherungen

Link zu diesem Datensatz	http://d-nb.info/820108170
Titel	Tage und Nächte steigen aus dem Strom : e. Donaufahrt Georg Lentz
Person(en)	Buchheim, Lothar-Günther (Verfasser)
Ausgabe	1. Aufl.
Verlag	[München] : Goldmann
Zeitliche Einordnung	Erscheinungsdatum: 1981
Umfang/Format	269 S. ; 18 cm
ISBN/Einband/Preis	978-3-442-06343-7 kart. : DM 6.80 3-442-06343-4 kart. : DM 6.80
Beziehungen	Ein Goldmann-Taschenbuch ; 6343
Anmerkungen	Lizenz d. Langen-Müller-Verl., München, Wien
Schlagwörter	Donau
Sachgruppe(n)	61 Geographie, Heimat- und Länderkunde, Reisen

Schlagwortfolge	Donau Bootsfahrt Reisebericht 1938 Buchheim, Lothar-Günther Autobiographie
Einzelschlagwörter (GND)	4049275-8 Reise 4015701-5 Europa 4263903-7 Flusswandern
Schlagwörter aus anderen Thesauri	Donau Reisebeschreibungen Danube River Description and travel Autobiographie 1941 Donau – Bootsfahrt – Reisebericht Buchheim, Lothar-Günther
Klassifikation	BKL 21.30 BKL 74.05 SFB ERD 242 RVK GN 9999 BKL 18.10 BKL 17.97

IE-Übernahme aus Werkbündeln

- Ziel: Werkbündel nutzen, um inhalts- und formalerschließende Metadaten von anderen Bündelmitgliedern zu übernehmen
- Schlagwörter und Klassifikation retrospektiv in Fremddatenfelder übernehmen wenn keine intellektuelle Erschließung vorhanden ist
- Ausgewählte Thesauri und Klassifikationssysteme
- Für Neuzugänge alles übernehmen

Gliederung

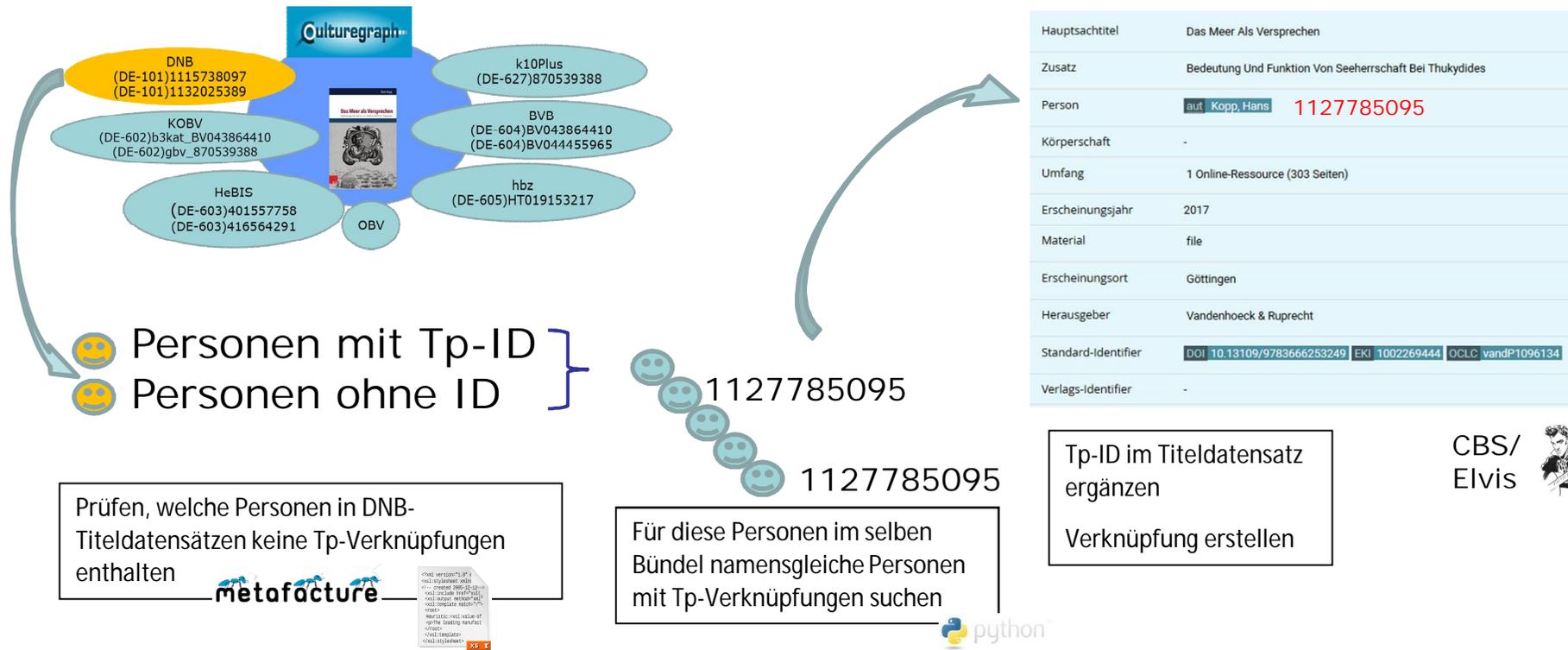
1. Werkbündelung
2. Personen mit Normdatenverknüpfungen anreichern
3. Datenabgleich mit externen Informationsquellen zur Datenanreicherung, hier: ORCID

Übernahme von Verknüpfungen aus Culturegraph-Werkbündeln

Hauptsachtitel	Das Meer Als Versprechen	Hauptsachtitel	Das Meer Als Versprechen
Zusatz	Bedeutung Und Funktion Von Seeherrschaft I	Zusatz	Bedeutung Und Funktion Von Seeherrschaft Bei Thukydides
Person	aut Kopp, Hans	Inhaltsverzeichnis	Inhaltstext
Körperschaft	-	Person	aut Kopp, Hans 1127785095
Umfang	1 Online-Ressource (303 Seiten)	Körperschaft	-
Erscheinungsjahr	2017	Umfang	303 Seiten
Material	file	Erscheinungsjahr	2017
Erscheinungsort	Göttingen	Material	book
Herausgeber	Vandenhoeck & Ruprecht	Erscheinungsort	Bristol, CT
Standard-Identifizier	DOI 10.13109/9783666253249 EKI 1002269444 OCLC vandP1096134	Herausgeber	Vandenhoeck & Ruprecht
Verlags-Identifizier	-	Schlagwort	655 Hochschulschrift 4113937-9 650 Seeherrschaft 4077302-4 600 Thucydides 4138063-0
		Klassifikation	DDC 938.05 PVK FH 26175 PVK NH 5850

Anreicherung

Ablauf beim Abgleich von Personen



Erfolgte Normdatenanreicherung

Ausgabe

1021053295, (DE-101)962759643, Petersen§ Holger, 100
Tp-ID, IDN Titeldatensatz, Name, Feld

- 7,5 Mio. DNB-Personen in Werkbündeln ohne Tp-Verknüpfung
- Ca. 1,4 Mio. Tp-Anreicherungen in 1,2 Mio. Datensätzen ermittelt (19 %)
- Im Sommer 2021 Anreicherungen vorgenommen

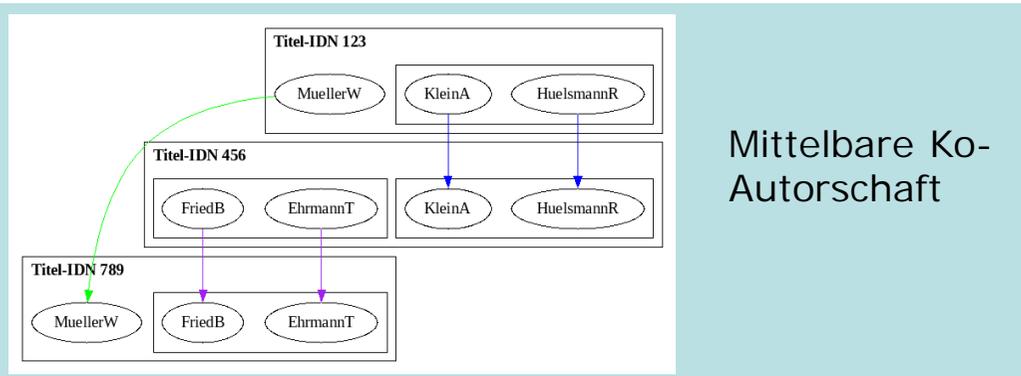
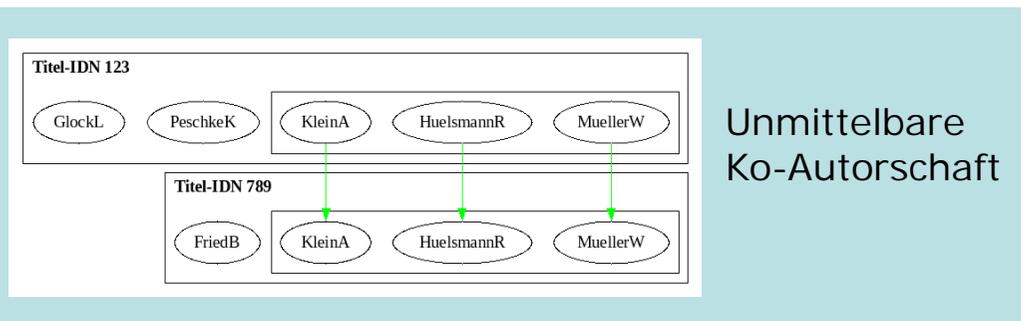
Manifestationsbündelung

- Testen der Übernahme von Metadatenelementen von anderen Manifestationen (=gleiche Ausgabe eines Werkes)
- Exaktere Übereinstimmung einer Publikation notwendig als bei Werken
- Übernahme von Autor*innen und ggf. beteiligten Personen
- Auffinden im Bestand fehlender Publikationen

Nutzung von Ko-Autor*innen-Beziehungen

- Wissenschaftskommunikation erfolgt in den meisten Fällen kollaborativ
- Forschende aus derselben Domäne (oder auch interdisziplinär), aber u.U. aus unterschiedlichen Institutionen arbeiten über längere Zeiträume immer wieder zusammen und veröffentlichen gemeinsam

Ko-Autorschaften auswerten



- Bei selben Namen können vorhandene GND-Verknüpfungen übertragen werden
- Bei 2 identischen Ko-Autoren Anreicherungen bei 276.000 Personen im DNB-Bestand möglich

Gliederung

1. Werkbündelung
2. Personen mit Normdatenverknüpfungen anreichern
3. Datenabgleich mit externen Informationsquellen zur Datenanreicherung, hier: ORCID

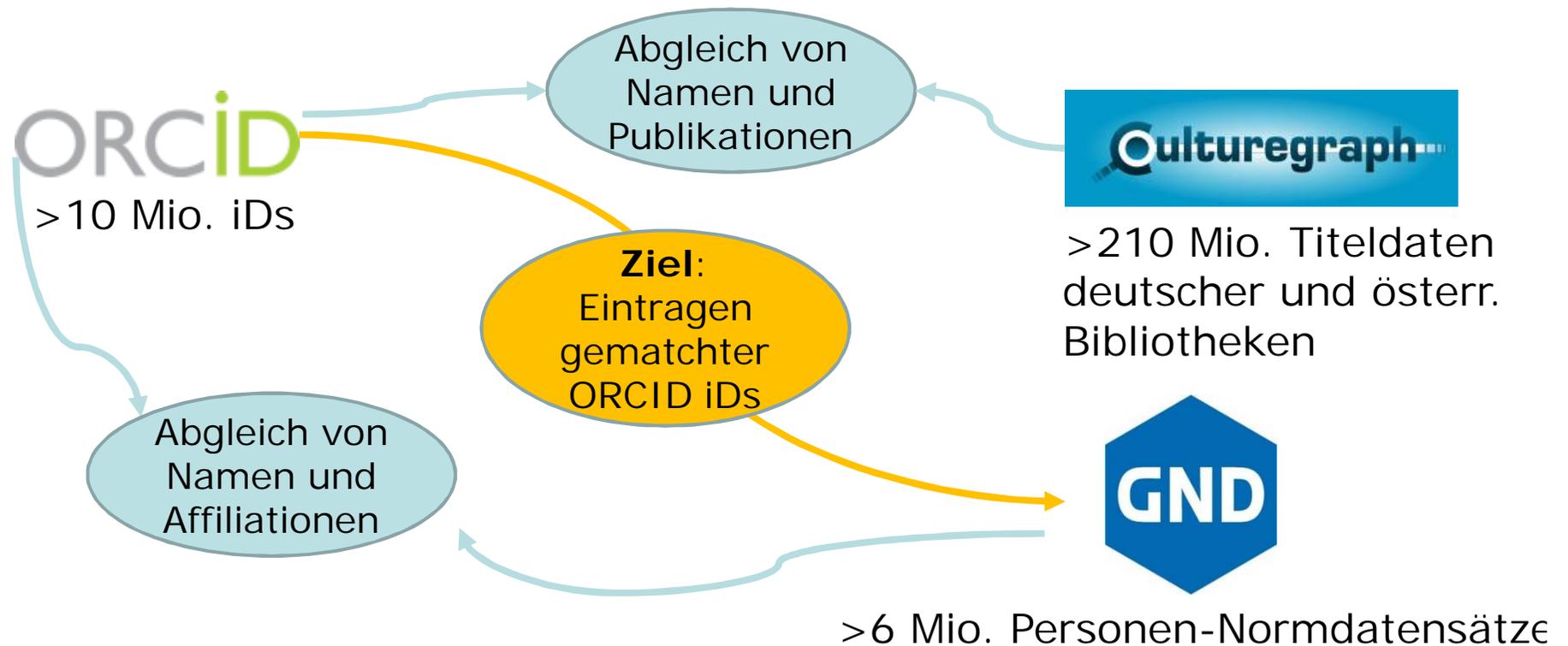


ORCID



- Open Researcher and Contributor ID
- Eindeutige Kennung zur Identifikation einer Person
- Wissenschaftler*innen erstellen Record mit der Möglichkeit, Informationen zu Werdegang, Publikationen, u.a. zu hinterlegen

Verfügbare Ressourcen



Abgleich von ORCID und GND-Sätzen

005 Tp3
006 http://d-nb.info/gnd/123157269
008 piz
011 f
012 v
024 orcid: 0000-0003-1905-6461\$ orcid
035 gnd/123157269
039 pnd/123157269\$vzg
039 pnd/17421569X
040 \$erda
043 XA-DE:XB-CN
100 Gruschke, Andreas
375 m
510 !101594133!Universität Leipzig\$bOrientalisches Institut [Tb1]\$4affi
548 1960\$b2018\$4datl
548 16.04.1960\$4datx
550 Prof. Dr.\$4akad
550 !040202143!Geograf ITs1!\$4berc

Immer vollständige
Übereinstimmung des Namens
oder alternativer Namensformen

Affiliation kann bis
zu einem Matchwert
von 0.89 variieren,
Universität und
Institut werden
berücksichtigt

The screenshot shows the ORCID profile for Andreas Gruschke. The profile includes the name 'Andreas Gruschke', an ORCID iD, and a list of employment records. The employment records are: 1. Sichuan University: Chengdu, Sichuan (2012-11-01 to 2018-01-30), Professor (Institute for Social Development and Western China Development Studies), marked as a preferred source. 2. Universität Leipzig: Leipzig, Sachsen (2011-01-01 to 2012-06-31), Post-doc researcher (Orientalisches Institut), marked as a preferred source. 3. Helmholtz-Zentrum für Umweltforschung UFZ: Leipzig, Sachsen (2010-08-01 to 2010-12-31), Post-doc researcher (Department Head Department Ecological Modelling (OESA)), marked as a preferred source.



ORCID
Connecting Research and Researchers

EDIT YOUR RECORD ABOUT ORCID CONTACT US HELP

Eldad Davidov

ORCID ID
<https://orcid.org/0000-0002-3396-969X>

Print view

Country
Germany, Switzerland

Keywords
Structural equation modeling, comparative empirical research, values (the Schwartz value theory), attitudes toward immigrants

Websites
Publications

Other IDs
ISNI: 0000000089475125

Employment (2)

- University of Cologne: Cologne, Germany
2017-01 to present | Prof. Dr. (Institute of Sociology and Social Psychology)
Source: Eldad Davidov
- University of Zurich: Zurich, Switzerland
2009-09-01 to present | Prof. Dr.
Source: Eldad Davidov

Works (1 of 1)

- Methods, theories, and empirical applications in the social sciences : Festschrift for Peter Schmidt
edited-book
DOI: 10.1007/978-3-531-17130-2
Part of ISBN: 978-3-531-17130-2
Source: Eldad Davidov Preferred source

Schlüssel:
davidoveldad*methodstheoriesand
empiricalapplicationsinsoc

Schlüssel:
davidoveldad*methodstheoriesand
empiricalapplicationsinsoc

Culturegraph

Titel Methods, Theories, And Empirical Applications In The Social Sciences

Zusatz Festschrift For Peter Schmidt

Person [GND Samuel Salzborn](#) [GND Eldad Davidov](#)
[GND Peter Reinecke](#) [GND Peter Schmidt](#)

Umfang 351 S.

Erscheinungsjahr 2012

Schlagwort [GND Rosenstock-Huessy, Eugen](#)
[GND Empirische Sozialforschung](#) [GND Methodologie](#)
[GND Politische Wissenschaft](#) [GND Sozialpsychologie](#)
[GND Empirische Sozialforschung](#) [Aufsatzsammlung](#)
[Aufsatzsammlung](#)

Klassifikation [DDC 300](#) [DDC 301.0723](#)

GND-ID:
1022834061

Wenn GND-ID
vorhanden,
Eintrag in Tp-Satz

005 Tp3
006 <http://d-nb.info/gnd/1022834061>\$z<http://d-nb.info/gnd/129032956>
008 piz
011 f
012 v
024 [orcid: 0000-0002-3396-969X](https://orcid.org/0000-0002-3396-969X)\$vHerkunft: orcid
035 [gnd/1022834061](http://d-nb.info/gnd/1022834061)
039 [gnd/129032956](http://d-nb.info/gnd/129032956)
039 [pnd/129032956](http://d-nb.info/gnd/129032956)\$vzg
043 XA-DE;XA-CH
100 Davidov, Eldad
510 [!004815750!](http://nbn-resolving.org/urn:nbn:de:hbz:5:1-63884-p0011-9)Universität zu Köln [Tb 1]\$4affi
510 [!000361909!](http://nbn-resolving.org/urn:nbn:de:hbz:5:1-63884-p0011-9)Universität Zürich [Tb 1]\$4affi
548 1971\$4datl
550 [!040550300!](http://nbn-resolving.org/urn:nbn:de:hbz:5:1-63884-p0011-9)Soziologie [Tb 1]\$4here



Standardidentifizier

- Herstellen von Verknüpfungen zur GND via gelieferte ORCiDs, ISNIs, etc.
- Einbringen von weiteren ORCiDs in Titel- und GND-Datensätze aus anderen Quellen (DNB-Claiming-Logs, BASE-Claiming-Logs, ORCiD-Werke-Dump)
- Jede Übernahme/Verknüpfung wird zusätzlich geprüft (z.B. stimmen Name im Titel- und GND-Satz überein, passt der geclaimte Name zum ORCiD-Record)

Personen-Vorschlagssätze mit Informationen aus ORCID

- Automatisierte Erstellung von Vorschlägen für neue Personennormdatensätze
- Informationen aus dem Titeldatensatz und dem ORCID-Record
 - Andere Standardnummern für Personen (z.Zt. ISNI)
 - Ländercode
 - Abweichende Namensformen
 - Affiliation
 - Keywords
 - Wirkungsdaten
 - Akademischer Grad
- Vorschläge können perspektivisch aus verschiedenen Quellen stammen

Beispiel Vorschlagsdatensatz Person

Titel

Eingabe: 1140:05-09-17 Änderung: 9999:20-08-20 12:35:34 Status: 1140:05-09-17
0500 Cof
0501 Text**Sb**bt
0502 Computermedien**Sbc**
0503 Online-Ressource**Sbcr**
0550 Webformular
0551 **Sb**sm
0598 NPL004
0600 ro,rb
1100 2016
1101 cr
1500 /1ger
1700 /1XA-DE-NW
2050 urn:nbn:de:101:1-2017090535307
2105 17,010
2150 NPSozialer Zusammenhalt in Bremen: Anlage 1 Codebuch
2198 535658540
2240 DNB:1139306340
2242 |o|1005511653
3000 Arant, Regina**Sb**Verfasser**S4**aut**S**Ea**S**Hnpi**S**D2020-06-12
3010 Larsen, Mandi**Sb**Verfasser**S4**aut**S**(orcid)0000-0001-5057-0085**S**Ea**S**Horcid**S**D2020-08-20
3010 Boehnke, Klaus**Sb**Verfasser**S4**aut**S**Ea**S**Hnpi**S**D2020-06-12
3119 Bertelsmann Stiftung**Sb**Sonstige**S4**oth
4000 Sozialer Zusammenhalt in Bremen : Anlage 1 Codebuch / Regina Arant, Mandi Larsen, Klaus Boertelsmann Stiftung
4030 Gütersloh : Bertelsmann Stiftung
4060 Online-Ressource
4083 =A \$
4085 "HTTP"=q pdf=U
http://www.bertelsmann-stiftung.de/de/publikationen/publikation/did/codebuch-sozialer-zusammenhalt-H=z LF
4233 Saaa**S**DE-101
5050 300**S**Em**S**Haeps**S**K0,966**S**D2017-09-06
5050 300**S**Ea**S**Hnpi**S**D2017-09-05
5051 3KK_A8_03_20160930_08
5052 3380**S**0,952**S**390**S**0,839**S**D2017-09-06
5450 [noScheme]nbbre
5550 [ckw]040081354|Bremen [Tg]t**S**Ea**S**Hstwg**S**K1**S**D2019-04-25

Vorschlagsdatensatz

Eingabe: 9999:04-02-22 Änderung: 9999:04-02-22 17:00:52 Status: 9999:04-02-22
005 TX
008 piz
011 tf
024 orcid: 0000-0001-5057-0085**S**vHerkunft: orcid
043 DE:US
100 Larsen, Mandi
400 Larsen, Mandi M. **S**vName aus Titel
510 Jacobs University Bremen gGmbH**S**gBremen, Bremen, DE**S4**am
510 Universitätsklinikum Hamburg-Eppendorf**S**gHamburg, Hamburg, DE**S4**affi
510 Safe Horizon**S**gNew York, New York, US**S4**affi
510 Columbia University **S**gNew York, NY, US**S4**affi
548 2004**S4**datw
550 Dr. **S4**akad
550 social cohesion**S4**them
550 intimate partner violence**S4**them
550 health inequalities**S4**them
667 Vorschlags-Ranking: 7**S5**DE-101
667 Quelle: Auswertung DNB-ORCID-Claiming-Logs; zugehöriger Titel: 1108361051113565-101
667 Quelle: Auswertung DNB-ORCID-Claiming-Logs; zugehöriger Titel: 1108361051113565-101
667 Quelle: Auswertung DNB-ORCID-Claiming-Logs; zugehöriger Titel: 11139306340**S5**DE-101
667 Quelle: Auswertung DNB-ORCID-Claiming-Logs; zugehöriger Titel: 11140167820**S5**DE-101
667 Quelle: Auswertung DNB-ORCID-Claiming-Logs; zugehöriger Titel: 11148691148**S5**DE-101
667 Quelle: Auswertung DNB-ORCID-Claiming-Logs; zugehöriger Titel: 11175694444**S5**DE-101
667 Quelle: Auswertung DNB-ORCID-Claiming-Logs; zugehöriger Titel: 1117569505X**S5**DE-101
672 Health Inequities Related to Intimate Partner Violence Against Women**Sb**The Role of Social Policy in the States, Germany, and Norway**S**I2016**S**w(DE-101)1080799184
672 Health Inequities Related to Intimate Partner Violence Against Women**Sb**The Role of Social Policy in the States, Germany, and Norway**S**I2016**S**w(DE-101)1080799184
672 Sozialer Zusammenhalt in Bremen**Sb**Anlage 1 Codebuch**S**I2016**S**w(DE-101)1139306340**S**0(urn)urn:nbn:de:101:1-2017090535307
672 Sozialer Zusammenhalt in Bremen**S**I2016**S**w(DE-101)1140167820**S**0(urn)urn:nbn:de:101:1-201709211672
672 Sozialer Zusammenhalt in Bremen**Sb**Anlage 2 Methodenbericht**S**I2016**S**w(DE-101)1148691148**S**0(urn)urn:nbn:de:101:1-2017121418157
672 What Holds Asian Societies Together?**S**bInsights from the Social Cohesion Radar; Codebook**S**I2017**S**w(DE-101)1175694444**S**0(urn)urn:nbn:de:101:1-2019011815404429194545
672 What Holds Asian Societies Together?**S**bInsights from the Social Cohesion Radar; Methods

ORCID record

The screenshot shows the ORCID record for Mandi Larsen. Key elements include:

- Name:** Mandi Larsen
- Other IDs:** Scopus Author ID: 7201731343, Loop profile: 1199434, Scopus Author ID: 57202743638
- Keywords:** social cohesion, intimate partner violence, health inequalities
- Countries:** Germany
- Activities:** Employment (5) is expanded to show:
 - Jacobs University Bremen gGmbH: Bremen, Bremen, DE** (2016-10-01 to present) | University Lecturer / Methods Center Coordinator (Bremen International Graduate School for Social Sciences) | Employment | Source: Mandi Larsen
 - Jacobs University Bremen gGmbH: Bremen, Bremen, DE** (2015-06-01 to 2016-09-30) | Postdoctoral Researcher (Diversity Focus Area) | Employment | Source: Mandi Larsen
 - Universitätsklinikum Hamburg-Eppendorf: Hamburg, Hamburg, DE** (2008-10-01 to 2010-05-31) | German Chancellor Fellow (Institute of Forensic Medicine) | Employment | Source: Mandi Larsen
 - Safe Horizon: New York, New York, US** (2004-11-01 to 2008-05-31) | Research and Evaluation Associate (Research and Evaluation Department) | Employment | Source: Mandi Larsen

Vernetzung von Normdaten stärken



- Grundlage dieser Verfahren: Bereits vorhandene Normdaten und ihre Verknüpfungen sollen nachgenutzt werden
- Notwendig: Mehr (verknüpfte) Personennormdatensätze, z.B. auch für Autor*innen von neuen Netzpublikationen, die nicht mehr intellektuell bearbeitet werden

Danke für Ihre Aufmerksamkeit!



Dr. Angela Vorndran
Team Datenmanagement
Deutsche Nationalbibliothek
a.vorndran@dnb.de