



# Forschungsdaten in den Geisteswissenschaften – die germanistische Sprachwissenschaft

Andreas Witt

Institut für Deutsche Sprache, Mannheim

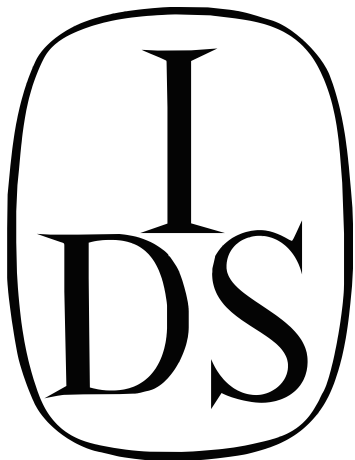
Workshop Forschungsdaten

WGL Geschäftsstelle Berlin

2012-05-10

# Institut für Deutsche Sprache

- gegründet 1964
- primäres außeruniversitäres  
Forschungsinstitut zur Erforschung  
und Dokumentation der deutschen  
Gegenwartssprache



# Linguistische Primärdaten am IDS

- Lexikographische Informationssysteme
  - z.B. OWID
- Geschriebensprachliche Korpora
  - z.B. DeReKo
- Gesprochensprachliche Korpora
  - z.B. FOLK

# Online-Wortschatz- Informationssystem Deutsch (OWID)

- Plattform für ein lexikografisches Informationssystem
- Bereitstellung lexikografischer und lexikologischer Arbeiten in computerlexikografisch angemessener Weise
- Erforschung der gemeinsamen Präsentierbarkeit unterschiedlicher lexikografischer Ressourcen, der Vernetzbarkeit von Ressourcen und des Umgangs der Benutzer mit neuen Formen eines Wortschatzinformationssystems

**Stichwortsuche:**

- Groß-/Kleinschreibung beachten

Suche in ...

- elexiko
- Neologismen
- Feste Wortverbindungen
- Diskurswörterbuch 1945-55

Optionen

- in Gesamtstichwortliste
- in allen lexikografisch ausgearbeiteten Artikeln

## Willkommen in OWID, dem Online-Wortschatz-Informationssystem Deutsch des Instituts für Deutsche Sprache, Mannheim

In OWID – einem Portal für wissenschaftliche, korpusbasierte Lexikografie – finden Sie Wörterbücher zum Deutschen mit unterschiedlichen inhaltlichen Schwerpunkten (mehr unter [Projekt OWID](#)).

Sie können über die allgemeine Stichwortsuche auf alle Wörterbücher dieses Informationssystems zugreifen. Über die „Suche in“-Liste unter dem Suchfenster können Sie festlegen, ob alle in OWID enthaltenen Wörterbücher in die Suche einbezogen werden sollen oder nicht. Darüber hinaus können Sie über „Optionen“ auswählen, ob Sie die Suche über die gesamte Stichwortliste von OWID laufen lassen wollen (diese besteht aus etwa 300.000 Stichwörtern, zu denen es in der Mehrzahl Informationen zur Rechtschreibung, in Kürze mit automatisch gewonnenen Belegen gibt), oder ob Sie die Suche auf lexikografisch ausgearbeitete Artikel mit tiefergehenden Informationen aus den einzelnen Wörterbüchern einschränken wollen (mehr [Informationen zur Suche](#)).

Wenn Sie sich für eines der Wörterbücher in OWID speziell interessieren, klicken Sie bitte auf die entsprechende Schaltfläche.

[elexiko](#)

[Neologismenwörterbuch](#)

[Feste Wortverbindungen](#)

[Diskurswörterbuch 1945-55](#)

# Deutsches Referenzkorpus (DeReKo)

- weltweit größte Sammlung deutschsprachiger Korpora (ca. 5 Milliarden Wörter, umfangreiche Annotationen)
- Texte aus Gegenwart und jüngerer Vergangenheit als Dokumentation des tatsächlichen Gebrauchs der deutschen Sprache in der geschriebenen Modalität
- empirische Basis für nationale und internationale Germanisten, sowie für interdisziplinäre Forschungen
- überwiegend urheberrechtlich geschütztes Material, daher:
- kein Zugriff auf vollständige Korpus-texte, sondern nur auf begrenzte Kontexte



# COSMAS II

Ergebnisse von 'derartig\*' in Korpus: mk - Mannheimer Korpora 1+2

Ansicht Andocken Extras Hilfe

Treffer	Texte	Dokument
1	1	MK1/LBC Heinrich Böll: Ansichten eines Clowns, [Roman], 1963
8	1	MK1/LFH Max Frisch: Homo faber, [Roman], 1957
1	1	MK1/LGB Günter Grass: Die Blechtrommel, [Roman], 1962
1	1	MK1/LSO Erwin Strittmatter: Ole Bienkopp, [Roman], 1963
4	1	MK1/TPM Heinz Pinkwart: Mord ist schlecht für hohen Blutdruck,
1	1	MK1/WBO Peter Bamm: Ex ovo, [Sachbuch], 1956
3	1	MK1/WGW Otto Willi Gail ; W. Petri: Weltraumfahrt, [Sachbuch],
1	1	MK1/WHK Hermann Heimpel: Kapitulation vor der Geschichte? [Sachb
2	1	MK1/WJA Karl Jaspers: Die Atombombe und die Zukunft des Menschen
20	1	MK1/WPE Rudolf Pörtner: Die Erben Roms, [Sachbuch], 1964
24	1	MK1/WUB Karl Ullrich: Wehr Dich, Bürger! [Sachbuch], 1960
8	3	MK1/ZBW Bild der Wissenschaft, [Fachzshr.], Jan. - März 1967
1	1	MK1/ZB4 Bildzeitung, [Tagesztg.], April 1967 (ausgewählte Artike
2	2	MK1/ZB6 Bildzeitung, [Tagesztg.], Juni 1967 (ausgewählte Artike
21	17	MK1/ZFA Frankfurter Allgemeine, [Tagesztg.], Dez. 1965 - Febr.
9	1	MK1/ZSG Studium generale, [Fachzshr.], Dez. 1966, (ausgewählte
5	2	MK1/ZUR Urania, [Fachzshr.], 1966-1967, (ausgewählte Artikel)
3	3	MK1/ZWE Die Welt, [Tagesztg.], Jan. - Febr. 1966, (ausgewählte
3	1	MK2/HAB Anweisungen, (1961, 1966, 1972)
1	1	MK2/PRO Informationsmaterial, (1960-1974)
1	1	MK2/TRI Trivialliteratur, (1960-1974)
1	1	MK2/WF1 Forschungsberichte des Instituts für Deutsche Sprache,
9	7	MK2/ZTG Zeitungen und Zeitschriften, (1960, 1969, 1973)
<b>130</b>	<b>51</b>	<b>23 Dokumente</b>

Textsuche in Korpus: hist - historische Korpora

Korpus Palette Optionen Fenster Hilfe

\$ SUCHWORT

\$ X NICHT Y

\$ LEMMA

\$ X ABSTAND Y

Bsp. Mauerspecht

STR()  
"beispiellos\*"

beispiellos\*

Neu

Suchen

KWIC (chronologisch sortiert)

Aktive Zeilen: 1-9072 setzen

[1 - 100] von 9072

Options... Hilfe ?

KWIC	Volltext	Text
1	GOE/AGF 02286	ist und daß er ur <b>Entschuldigung</b> , daß im vorhergehenden beim
2	GOE/AGF 02286	entspringende Hindernisse sind, zur <b>Entschuldigung</b> angeführt werden können, von
3	GOE/AGI 00000	die reinste Würdigung oder vielmehr <b>Abwürdigung</b> der irdischen Dinge, den
4	GOE/AGI 00000	tun und zuletzt eine unschätzbare <b>Befriedigung</b> hoffen, man trifft Spuren
5	GOE/AGI 00000	, daß ältere Künstler die <b>Verköndigung</b> Mania also vorstellen, daß
6	GOE/AGI 00000	den 25. März 1787. <b>Verköndigung</b> Mania. ob ich gleich
7	GOE/AGI 00000	sorgenlos dahin leben, augenblickliche <b>Befriedigung</b> , mäßiger Genuß, vorübergehend
8	GOE/AGI 00000	Gewalt des römischen Zauberkreises zur <b>Entschuldigung</b> dienen, als wir nach
9	GOE/AGI 00000	klarsten Menschenverstand, die reinste <b>Würdigung</b> oder vielmehr <b>Abwürdigung</b> der irdi
10	GOE/AGI 00000	Anstrengung hinübersetzen kann, an <b>Verteidigung</b> ist hier nicht zu denken
11	GOE/AGI 00000	weitere Reise als Gründe der <b>Entschuldigung</b> möchte gelten lassen, er
12	GOE/AGI 00000	einer lebenswürdigen jungen Dame, <b>Huldigung</b> , anzunehmen gewohnt und geneigt.
13	THM/AM1.000...	unmittelbare Erheiterung, Erwärmung, <b>Befriedigung</b> erweckt, die ich bei
14	THM/AM1.003...	Motiven versehen habe, die <b>Entschuldigung</b> ist so häufig wie der
15	THM/AMN.00...	Treulosigkeit in der Treue. Ents <b>1813</b> ung seiner Gefühle durch die Vorstel
16	THM/AMN.00...	intelligente Scheu vor Menschen- und <b>Sachbeschädigung</b> , Sinn für das Schicksliche
17	THM/AMN.06...	des Kontinents auf Grund der <b>Verständigung</b> und und Freundschaft zwischen
18	THM/AM3.084...	Ermanglung direkter Beweise, die <b>Beschuldigung</b> , Hatvany habe in Zeitungsartike

# Linguistische Primärdaten zur gesprochenen Sprache

- Archiv für Gesprochenes Deutsch
  - 44 Korpora
  - 6.000 h Audio-Aufzeichnungen
  - 6.000.000 Token transkribiert
- Datenbank für Gesprochenes Deutsch
  - aufbereitete Daten
  - web-basierter Zugang
  - Relaunch 2012 als DGD 2.0
  - Nationales Referenzkorpus



# Archiv für Gesprochenes Deutsch (AGD) / Datenbank für Gesprochenes Deutsch (DGD)

- Sammlung von Korpora des gesprochenen Deutsch
- Audioaufnahmen, Videoaufnahmen, Transkriptionen
- Webschnittstelle (DGD)

# DGD

DATENBANK FÜR  
GESPROCHENES  
DEUTSCH

[ÜBER DIE DGD](#)

[KORPORA](#)

[RECHERCHE](#)

[HILFE](#)

[FAQ](#)

[ABMELDEN](#)

## Herzlich willkommen in der Datenbank für Gesprochenes Deutsch (DGD 2.0 - extern)

### Schnelleinstieg

Verschaffen Sie sich über den Menüpunkt "[Über die DGD > Bestand](#)" einen **Überblick** über die verfügbaren Daten. Detaillierte **Beschreibungen** der Korpora (Korpusmetadaten) finden Sie unter "[Korpora > Korpusbeschreibungen](#)".

Zum **Browsen** der Ereignis- oder Sprecherdokumentationen sowie von Transkripten und Zusatzmaterialien bzw. zum **Anhören** von Audioaufnahmen eines ausgewählten Korpus wählen Sie die entsprechenden Unterpunkte im Menü "[Korpora](#)".

Zum **Durchsuchen** von Ereignis- oder Sprecherdokumentationen sowie von Transkripten wählen Sie die entsprechenden Unterpunkte im Menü "[Recherche > Volltext](#)".

Einen direkten Zugang zu den **FOLK-Daten** erhalten Sie zum Beispiel, indem Sie die zugehörigen Transkripte über "[Korpora > Transkripte](#)" browsen oder über "[Recherche > Volltext](#)" durchsuchen.



02.05.2012

DATENSCHUTZ IMPRESSUM

# Forschungsdatenmanagement

1. Eingabe (Ingest Model)
2. Ausgabe und Zugriff (Access Model)
3. Technische Basis (Repositoryum)
4. Nachhaltiger Betrieb (Maintenance Model)

# Ingest Model (Input)

- Datenaufbereitung
- Definition von Standards
- lebendes Archiv vs. Langzeitarchiv

# Access Model (Output)

- Identifizierbarkeit
- Referenzierbarkeit
- Zugriffsberechtigungen

# Repository (Technik)

- Phys. Speicherstrategien
- Implementierung Input & Output
  - Ingest Model
  - Access Model
- Langzeitarchivierung

# Maintenance Model

- technische Gewährleistung der Nachhaltigkeit
- strukturelle Gewährleistung der Nachhaltigkeit
- Implementierung von Weiterentwicklungen
- Etablierung von Workflows bei Nutzern

# Arbeitsabläufe beim Umgang mit Primärdaten

## Bedarf:

- Suche in den vorhandenen Daten
- Referenzierung von gefundenen Daten
- Definition von eigenen Zusammenstellungen
- Zugriff auf eigene Definitionen
- Entdeckung von unbekanntem Daten

## Lösungen

- Nutzung oder Entwicklung von Standards
- Definition von Workflows



# Standardisierungsbeispiel 1: Metadaten

## Bedarf:

- Anforderungen wie für Primärdaten
- umfassender bzw. universeller Katalog (nicht möglich)
- dynamische Metadaten-Schemata

## Lösungen

- Component MetaData Infrastructure (CMDI)
  - komponentenbasierter Aufbau mit Registrierung
  - konform zu ISOcat (ISO 12620, Data Category Registry)
- fester, allgemeingültiger „Kern-Katalog“
- flexible, bedarfsspezifische Ergänzungen
- Harvesting: OAI-PMH-Schnittstelle

# Standardisierungsbeispiel 2: Identifizierung von Ressourcen

- Stabile Referenzen durch Persistent Identifier
  - Dienst zum Registrieren und Auflösen von PIDs
  - Verwendung von Handles
  - unabhängig vom Freiheitsgrad der Verfügbarkeit einer Ressource
- Dies erlaubt u.a. Querverweise zwischen Primärdaten und Publikationen







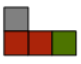

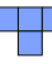

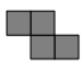


# Technisches Bewusstsein

- Softwarekonzept als tragende Säule
- Nachhaltigkeit durch modulare Architekturen
- Kontinuierliche Analyse und Hinterfragung der bestehenden Softwarelösungen

# Vernetzung

- Neben der Mitwirkung in der WGL-AG zu Forschungsdaten ist das IDS vielfältig vernetzt:
  - Fachspezifische Vernetzung
  - Fächerübergreifende Vernetzung
- Netzwerke:
  - Verbundprojekte
  - Offizielle Gremien
  - Interessenvereinigungen

# Projekte und Mitgliedschaften


<p><b>FI1</b></p> 	<p>Aufbau eines Zentrums „Digitale Forschungsressourcen für die germanistische Sprachwissenschaft“</p> <p>2009-2010 – gefördert vom BMBF</p>	<p><b>TG</b></p> 	<p><b>DIN</b></p> 	<p>TextGrid – Vernetzte Forschungs-umgebung in den eHumanities</p> <p>2009-2012 – gefördert vom BMBF</p> <p>Deutsches Institut für Normung, DIN</p> <p>seit 2004 – ohne Förderung</p>
<p><b>FI2</b></p> 	<p>Forschungsinfrastrukturen in wissenschaftlichen Einrichtungen: Implementierung eines Prototyps am Institut für Deutsche Sprache</p> <p>2011-2013 – gefördert vom BMBF</p>	<p><b>WG</b></p> 	<p><b>ISO</b></p> 	<p>WissGrid – Grid für die Wissenschaft</p> <p>2009-2012 – gefördert vom BMBF</p> <p>Internationale Organisation für Normung, ISO</p> <p>seit 2004 – ohne Förderung</p>
<p><b>DSPIN</b></p> 	<p>Deutsche Sprachressourcen-Infrastruktur D-SPIN</p> <p>2008-2011 – gefördert vom BMBF</p>	<p><b>FP1</b></p> 	<p><b>TEI</b></p> 	<p>Zentrum für germanistische Forschungsprimärdaten</p> <p>2011-2012 – gefördert von der DFG</p> <p>Text Encoding Initiative, TEI</p> <p>seit 2011 – ohne Förderung</p>
<p><b>CL/D</b></p> 	<p>Common Language Resources and Technology Infrastructure CLARIN in Deutschland</p> <p>2008-2011 – gefördert vom BMBF</p>	<p><b>VG</b></p> 	<p><b>nestor</b></p> 	<p>Verwertung Geist – Analyse und Konzepterstellung zur systemat. Verwertung geistes- und sozial-wissensch. Forschungsergebnisse</p> <p>2011-2013 – gefördert vom BMBF</p> <p>Network of Expertise in long-term Storage and availability of digital Resources in Germany, nestor</p> <p>seit 2010 – ohne Förderung</p>
<p><b>CL/EU</b></p> 	<p>Common Language Resources and Technology Infrastructure CLARIN in Europe</p> <p>2011-2014 – gefördert vom BMBF</p>			

**FI2**



## **Forschungsinfrastrukturen in wissenschaftlichen Einrichtungen: Implementierung eines Prototyps am IDS**

- Prototypische Einrichtung einer Organisationseinheit „Forschungsinfrastrukturen“ am IDS
- Entwicklung einer Art infrastrukturellen Werkzeugkastens
- Verwendung der Arbeitsergebnisse als mögliche Blaupause für andere Forschungsinstitute, die institutionelle Lösungen für den Umgang mit Forschungsdaten suchen

DSPIN	CL/D	CL/EU
		

## Deutsche Sprachressourcen-Infrastruktur; Common Language Resources and Technology Infrastructure in Deutschland / Europe

- Paneuropäische Initiative zum koordinierten Umgang mit Sprachressourcen und –technologien (ERIC)
- Organisation, Ausbau, Vernetzung und Vereinheitlichung existierender nationaler Infrastrukturen, EU-geförderter Projekte und ihrer Ressourcen
- Aufbau von dedizierten Zentren für eine funktionsfähige Infrastruktur
- Ressourcenanbieter-Föderation innerhalb des DFN
- Darstellung der rechtlichen Rahmenbedingungen
- Bereitstellung von Ressourcen, Daten und Werkzeugen über State-of-the-Art-Registraturen und Web-Services

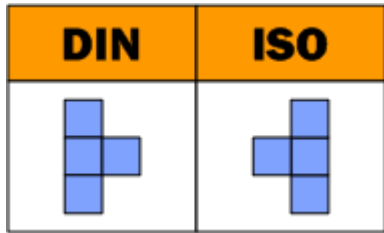
**FP1**



## Zentrum für germanistische Forschungsprimärdaten

- Etablierung einer einheitlichen Infrastruktur für die nachhaltige Bereitstellung von Forschungsprimärdaten aus der germanistischen Linguistik
- Entwicklung datentechnischer, dokumentatorischer und rechtlicher Standards sowie Best-Practice-Richtlinien für die nachhaltige Archivierung und Verfügbarmachung von sprach-wissenschaftlichen Primärdaten
- Perspektive: Öffnung des Zentrums für germanistische Forschungsdaten fremder Institutionen





## Deutsches Institut für Normung, DIN; Internationale Organisation für Normung, ISO

- nationale und internationale Vereinigung von Normungsorganisationen
- Mitarbeit im ISO-Komitee ISO/TC37/SC4 (Terminology and Other Language and Content Resources: Language Resources Management) und in der nationalen Spiegelgruppe in DIN





## Network of Expertise in long-term Storage and availability of digital Resources in Germany, nestor

- Kompetenznetzwerk Langzeitarchivierung und Langzeitverfügbarkeit digitaler Ressourcen in Deutschland
- IDS ist Partner seit 2010

# Programmbereich Forschungsinfrastrukturen

Bibliothek	Langzeit- archivierung	Technische Infrastrukturen	Recht	Infrastruktur- forschung



# Forschungsdaten in den Geisteswissenschaften – die germanistische Sprachwissenschaft

Andreas Witt

Institut für Deutsche Sprache, Mannheim

Workshop Forschungsdaten

WGL Geschäftsstelle Berlin

2012-05-10