

# Forschungsdaten in den Digital Humanities - Der Verbund nestor und Datenzentren am Institut für Deutsche Sprache und an der Universität zu Köln

Andreas Witt



# Gliederung

- nestor – Verbund
  - nestor: Partner und Arbeitsgruppen
  - die AG Zertifizierung
- Forschungsdaten in den Digital Humanities
  - Institut für Deutsche Sprache Mannheim
  - Das Data Center for the Humanities der Universität zu Köln
- Was können Infrastruktureinrichtungen leisten
  - Herausforderungen
  - Langfristigkeit
- Wie helfen Verbünde, wie nestor, CLARIN, DARIAH

VERANSTALTUNGEN  
nestor-Praktikertag  
Workshops  
Vorträge

DIN-Normen



nestor-Schools  
nestor-Handbuch

nestor

INTERNATIONALE  
KOOPERATION

DPC  
NCDD  
NDIIPP  
EU-Projekte

PUBLIKATIONEN  
Kriterienkatalog  
Handbuch  
Ratgeber  
Expertisen  
Infosheets  
nestor edition  
Newsletter

INFORMATIONSSANGEBOTE  
News  
Veranstungskalender  
Personen  
Projekte



ARBEITSGRUPPEN z.B.  
AG Elektronische Akte  
AG Emulation  
AG Formaterkennung  
AG Forschungsdaten  
AG Media  
AG Zertifizierung

# nestor - Partner



# nestor - Arbeitsgruppen

- ..■ AG Digitale Bestandserhaltung
- ..■ AG Elektronische Akte
- ..■ AG Emulation / GI FG Langzeitarchivierung
- ..■ AG Formaterkennung
- ..■ AG Forschungsdaten
- ..■ AG Kooperation der Archive
- ..■ AG A-V Medien
- ..■ AG OAIS Review
- ..■ AG Personal Digital Archiving
- ..■ AG Zertifizierung

# nestor - AG Zertifizierung

- Auf der Basis der DIN 31644: "Kriterien für vertrauenswürdige digitale Langzeitarchive", hat die AG ein Zertifizierungsverfahren entworfen und betreut die Durchführung
- Aufgaben
  - ..■ Zertifizierung digitaler Archive
  - ..■ Fortentwicklung des nestor-Siegels

# Prämissen zur Vertrauenswürdigkeit

- Ein *digitales Langzeitarchiv* ist **vertrauenswürdig**:
  - wenn es gemäß seinen Zielen und Spezifikationen zum Informationserhalt über lange Zeiträume hinweg operiert
  - seine Nutzer, Produzenten, Betreiber, Partner ihm dieses zutrauen
- Digitale Archive sind derzeit (noch) nicht per se vertrauenswürdig
- Nutzer können nur in Einzelfällen die Vertrauenswürdigkeit überprüfen

# Nationale Entwicklung

- 2004/5 Gründung der nestor-AG „Vertrauenswürdige Archive“
- 2006: Kriterienkatalog vertrauenswürdige digitale Langzeitarchive, Version 1
- 2008: Kriterienkatalog vertrauenswürdige digitale Langzeitarchive, Version 2
- 2008: Kriterienkatalog wird im NABD 15 des DIN zu einem Normungsprojekt
- 2012: Veröffentlichung der DIN 31644





# nestor-Siegel: Verfahren



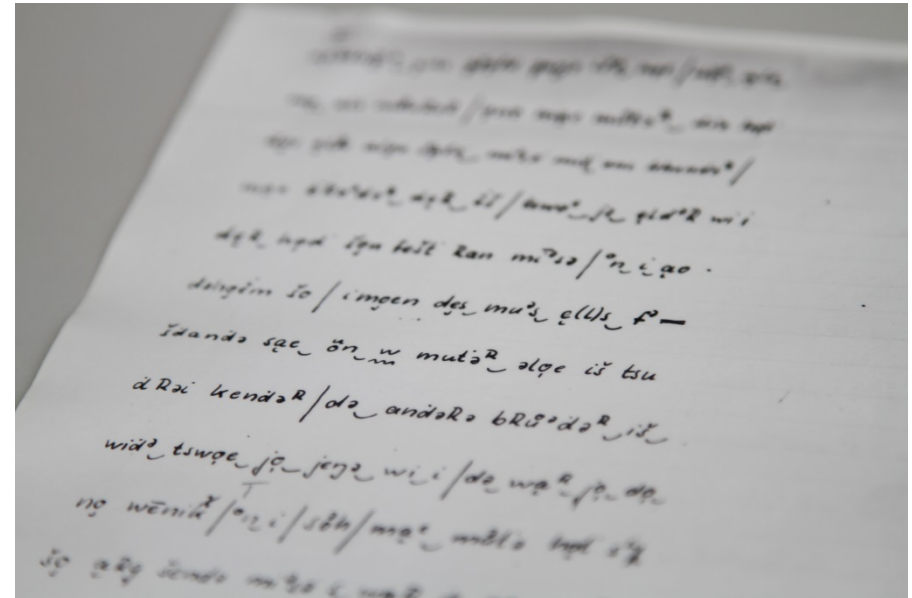
- Orientierung an DIN 31644
- Kein akkreditiertes Zertifizierungsverfahren
- „erweiterte Selbstevaluierung“ durch die nestor-AG Zertifizierung (2 Gutachter)
- Eingebettet in internationales MoU mit Data Seal of Approval und ISO 16363:  
<http://www.trusteddigitalrepository.eu/Trusted%20Digital%20Repository.html>

# Das Institut für Deutsche Sprache

- Sitz in Mannheim
- Gründung 1964
- Außeruniversitäre Forschungseinrichtung
- Mitglied der Leibniz Gemeinschaft
- Vernetzung
  - Betrieb eines CLARIN-B Zentrums
  - Mitglied in nestor
  - Sitz des Rats für Deutsche Rechtschreibung

# Beispiele aus der AGD am IDS

- Transkriptionen auf Papier



# Beispiele aus der AGD am IDS

- Tonbänder und Kassetten

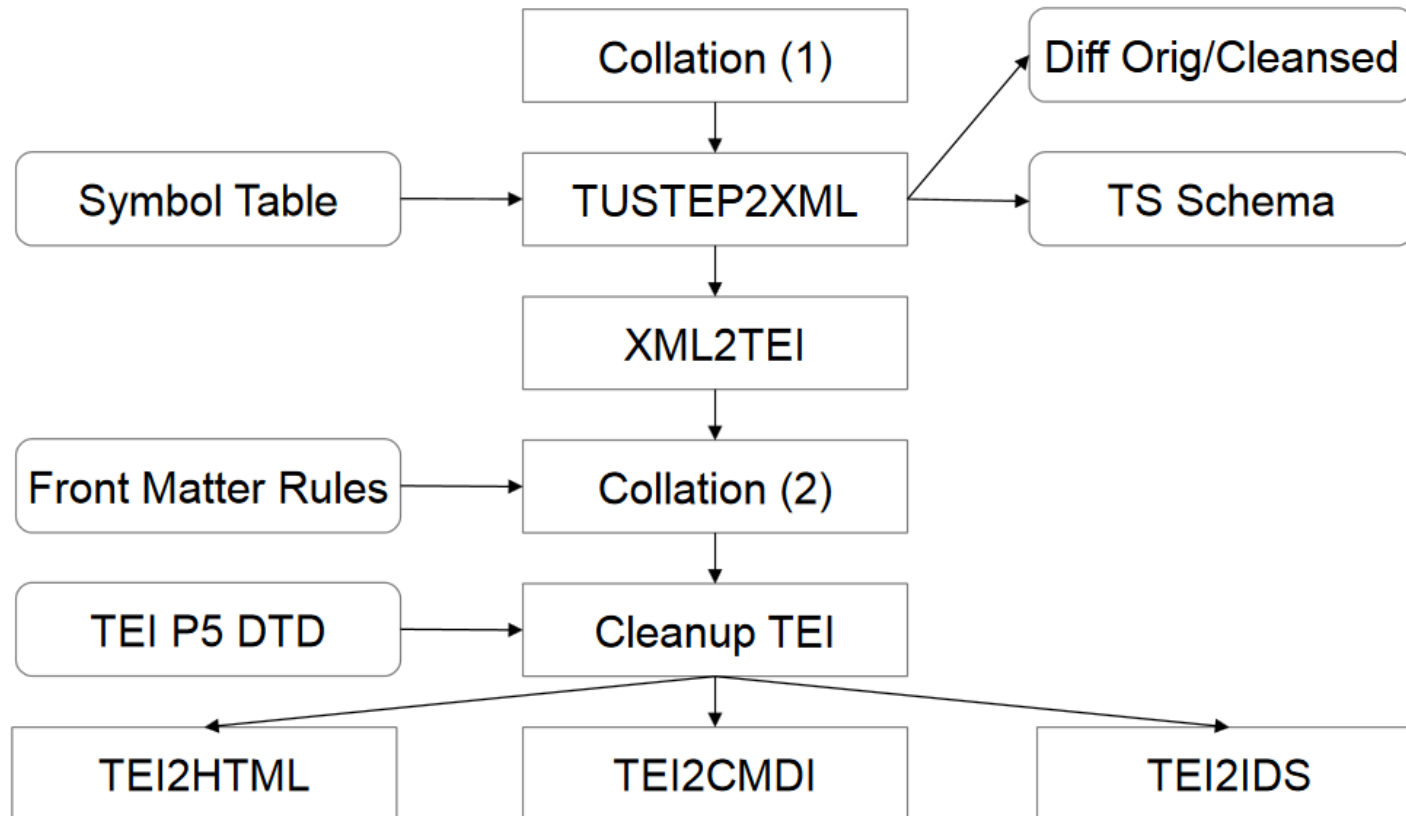


# Zwei Beispiele für Kurationsprojekte am IDS

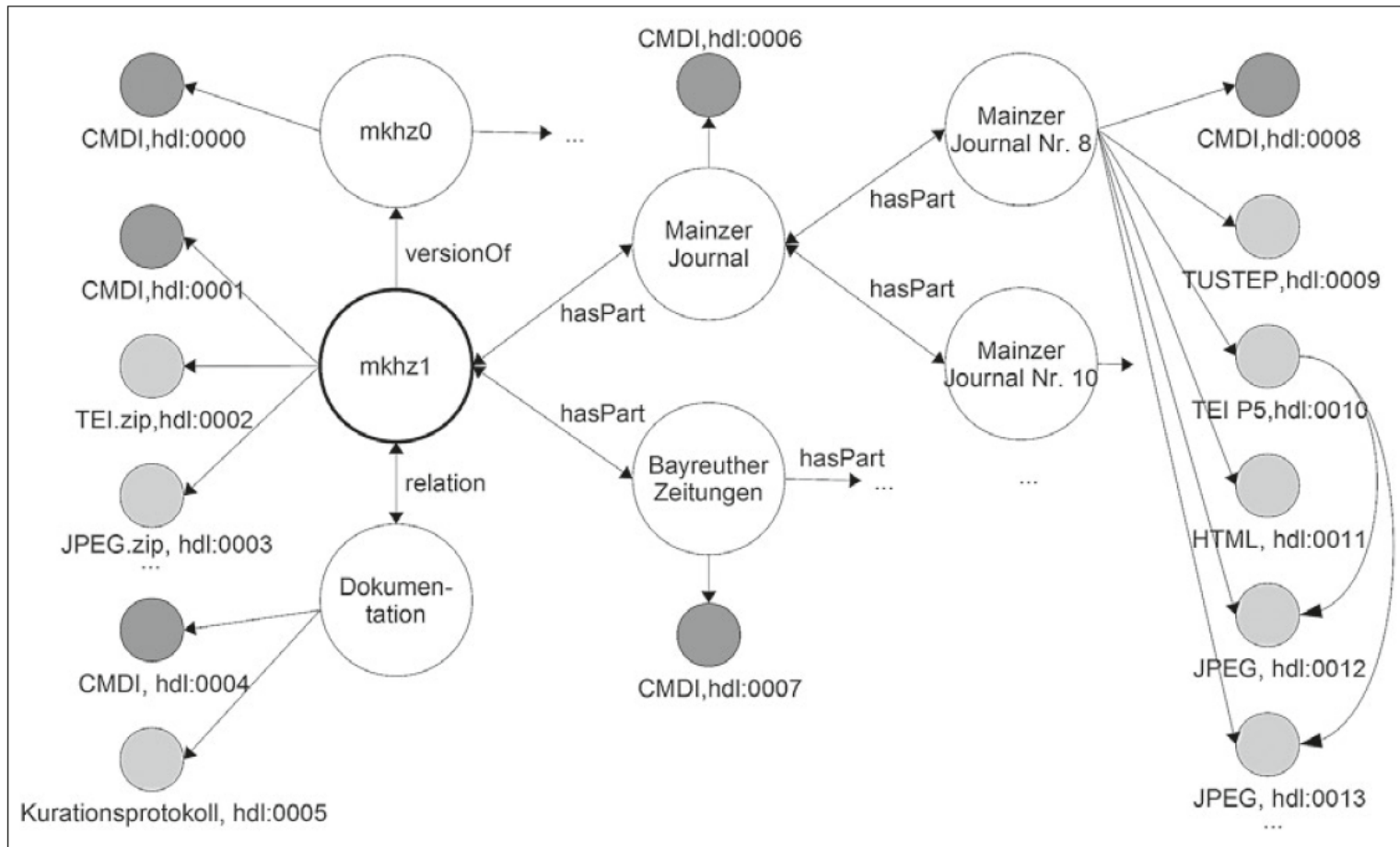
- Mannheimer Korpus Historischer Zeitungen und Zeitschriften (mkhz)
  - Konvertierung von TUSTEP nach TEI
  - COSMAS II
- Monash Corpus of Australien German
  - Daten von Michael Clyne
  - DGD
  - LZA



# Mannheimer Korpus Historischer Zeitungen und Zeitschriften



# Mannheimer Korpus Historischer Zeitungen und Zeitschriften





# Mannheimer Korpus Historischer Zeitungen und Zeitschriften

- Integration in COSMAS II
- Langzeitarchivierung
- Auffindbar im Virtual Language Observatory CLARIN
  
- Übergabe der Daten in das Deutsche Text Archiv an der Berlin-Brandenburgischen Akademie

(s. a. Fankhauser, P., Fiedler, N., & Witt, A. (2013).  
Forschungsdatenmanagement in den Geisteswissenschaften  
am Beispiel der germanistischen Linguistik. *Zeitschrift für  
Bibliothekswesen und Bibliographie*, 60(6), 296-306)



Monash University

CLAYTON VICTORIA AUSTRALIA

TELEPHONE: 544 0611 TELEGRAMS: Monashuni Melbourne

DEPARTMENT OF MODERN LANGUAGES

German Section

1. Oktober 1965

MC/AEK

Professor Dr. Dr. Zwirner  
Direktor des Deutschen Spracharchivs  
44 Münster, Westdeutschland  
Steinfurter Straße 107

6. Okt. 1965

DURCH LUFTPOST

Sehr geehrter Herr Professor Zwirner!

Ihr Schreiben vom Februar 1965 mit Bezug auf Tonbandaufnahmen von der Sprache deutscher Siedlungen in Australien ist mir aus unverständlichen Gründen erst vor kurzem in die Hand gekommen.

Wir sind in unserem germanistischen Seminar an dem Deutschen von Zweisprachigen in Australien besonders interessiert. Bisher haben wir uns hauptsächlich mit der "sprachlichen Assimilierung" deutscher Nachkriegseinwanderer befaßt. Wir hoffen, uns aber in nächster Zeit auch mit den Dialekten von deutschen Siedlungen in der Provinz Victoria zu beschäftigen.

Einige unserer "Postgraduate Students" und ich beabsichtigen, im Laufe der kommenden Ferien (Weihnachten bis Anfang Februar) einige Tonbandaufnahmen von den wenigen noch vorhandenen zweisprachigen Einwohnern der ehemaligen deutschen Siedlungen Westvictorias zu machen. Tonbänder und eventuelle Ergebnisse dieser kleinen Forschung würden wir Ihnen gerne zur Verfügung stellen. Sollten Sie irgendwelche Vorschläge oder Wünsche über die Durchführung dieser Untersuchung haben, so würde ich mich freuen, sie zu hören.

Mit freundlichen Grüßen,

hochachtungsvoll

I h r

*Michael Clyne*

(Michael Clyne)

Lecturer in German

Monash  
Corpus of  
Australian  
German

# Monash Corpus of Australian German

- 1966/67 rund 200 Aufnahmen aus drei Distrikten (Barossa Valley, Wimmera, Western District) und mit „Prewar“ Sprechern
- Die Daten wurden zur Forschung von u.a. Michael Clyne, Sandra Kipp, Brigitte Lambert, Doris Schüpach, Claudia Riehl, ... benutzt.
- 2009 – Michael Clyne und der Direktor des IDS einigen sich auf die Aufnahme in die AGD

# Michael Clyne

- \*1939 †2010
- Australischer Linguist mit deutschen Wurzeln
- Professor für Deutsche Linguistik an der Monash University, Melbourne
- Expert für Soziolinguistik, Zwei- und Mehrsprachigkeit, interkulturelle Kommunikation
- Gastwissenschaftler am SFB 538 in Hamburg
- Zusammenarbeit mit dem IDS seit den 1960s



# Monash Corpus of Australian German

- 2012 – Eine Kiste mit CDs, ohne Transkriptionen und ohne Metadaten wird ans IDS geliefert



- Claudia Riehl kann mit Namen und Emailadressen von drei Doktoranden von Michael Clyne in Melbourne aushelfen
- Alle haben (Word) Transkriptionen (und ein paar Metadaten) von verschiedenen Teilen des Korpus
- Die Daten auf den CDs scheinen nicht komplett zu sein.
- Sandra Kipp findet (Kopien?) der originalen Tonbänder und schickt sie ans IDS
- Sandra Kipp stellt Kontakte mit der Barossa Valley historical society und den Lutheran Archives in Adelaide her

# Monash Corpus of Australian German



- Kuration der Audiodaten
  - Digitalisierung der Bänder
  - Vergleich mit CDs
  - Editierung der besseren Version (meistens von den Bändern)
    - Normalisierung, Rauschunterdrückung
    - Schneiden, Zusammenfügen der Aufgezeichneten Interaktionen

# Monash Corpus of Australian German



- Metadaten Dokumentation
  - Einige Hinweise auf den Covern, Bändern, CDs
  - Einige Hinweise in Publikationen
  - Einige Hinweise in Transkriptionen
  - Biographische Informationen in den Aufnahmen selber
  - Immer noch große Lücken

# Monash Corpus of Australian German



- Transkriptionen
  - Konvertierung DOC > TXT, Vorverarbeitung mit Skripten
  - Import in EXMARaLDA mittels Simple EXMARaLDA
  - Manuales Alignment via Easy Alignmentfunktion
  - Manuale Vervollständigung (Interviewer meistens nicht transkribiert)
  - Pseudonymisierung der Transkripte / Anonymisierung im Audiosignal
  - Automatische Konvertierung, Orthographische Normalisierung, Lemmatisierung, Part-Of-Speech Tagging



# Monash Corpus of Australian German

- Etwa neun Stunden Aufnahmen des ersten Teilkorpus benötigten etwa 250 Personenstunden für die Kuration
- Übernahme in die DGD
- Langzeitarchivierung
- Auffindbar im Virtual Language Observatory  
CLARIN
  
- Gorisch, Jan/ Schmidt, Thomas/ Stift, Ulf-Michael (i.V.): Data of German Speech Minorities in the Archive for Spoken German. Eingereicht für: Boas, Hans C. (Hrsg.): Comparative language island research: Data, methods, and goals. Amsterdam: Benjamins.

# Data Center for the Humanities

- Forschungsdatenaktivitäten an der UzK
- Aktivitäten des DCH
  - Fallbeispiele
  - Beratungsworkflow
- Die Rolle des DCH im bei der Implementierung eines Forschungsdatenmanagements an der Universität zu Köln

# Forschungsdatenaktivitäten an der Philosophischen Fakultät der UzK

- Mitwirkung in der Pilotphase von DARIAH-DE
- Beteiligung am Verbund CLARIN
  - Koordination der F-AG zu Sprachtypologie
  - Etablierung von K-Zentrum auf der Ebene der CLARIN ERIC
  - Innerhalb des nationalen Verbunds CLARIN-D soll noch 2017 ein assoziiertes B-Zentrum zertifiziert werden
- KA<sup>3</sup> - BMBF-Zentrumsausschreibung
- Die UzK ist ein zentraler Partner der AWK-NRW und Akademienunion

# Fallbeispiele



## **Beratungsangebot Projektanträge**

Das DCH hilft Ihnen für Ihren Projektantrag ein Konzept für Datenhaltung, Archivierung und Bereitstellung Ihrer Forschungsdaten zu entwickeln und umzusetzen.



## **Projektentwicklung**

Das DCH führt in Kooperation mit Fachwissenschaftlern und unseren lokalen Partnern Entwicklungs- und Kurationsprojekte durch.



## **Digitale Ressourcen erhalten**

Laufend entstehen digitale Ressourcen, Datenbanken, Webseiten und Forschungsdatenbestände im Rahmen von zeitlich befristeten Projekten. Was geschieht mit diesen, nach dem Ende des Projekts?

# [1] Beispiel Anfrage

„Wir schreiben [...] gerade einen Antrag im Rahmen des Programms H2020 Reflective Societies, den wir in der kommenden Woche einreichen wollen.“

Dimension: Wann im Projektverlauf?

**vor Antragstellung – laufendes Projekt (anlaufend/auslaufend) – nach Projektende – legacy Data**

## [2] Beispiel Anfrage

„Wir starten in Kürze ein DFG-Projekt [...] und sind an der Einrichtung eines Coworking Spaces und einer Plattform für den Datenaustausch interessiert.“

Dimension: Welche Bedarfe gibt es?

## [3a] Beispiel Anfrage

„Im DFG-Projekt [...] nähert sich nach drei Jahren das Ende der Förderung (Februar 2018) und wir überlegen gerade, wie wir die angefallenen Forschungsdaten am besten aufbewahren können. Im Antrag hatten wir die Formulierung "Archivierung der Daten und Unterlagen" aufgenommen, wobei es sich in diesem Projekt um typische geisteswissenschaftliche qualitative Daten handelt.“

# [3b] Beispiel Anfrage

„Leider haben wir zu Beginn des Projekts versäumt, Standards festzulegen, deshalb habe ich drei Fragen: [...]

- Liegt nach dem Projektende das Urheberrecht - zum Beispiel für Präsentationen - beim Projekt oder bei mir als individueller Wissenschaftlerin?
- Wie sieht es mit dem Nutzungsrecht von digitalen Fotos aus, die ich in Archiven nach Registrierung eines persönlichen Accounts für den eigenen wissenschaftlichen Gebrauch gemacht habe, sowie bei Digitalisaten, die ich als Wissenschaftlerin für bestimmte Publikationen mit Projektgeldern erworben habe?
- Wie lässt sich nach dem Ende des Projekts die betreffende LitLink-Datenbank am besten archivieren?“



# Bedarfsgruppen

- Beratung/Support:
  - Datenmanagement im Projekt, wie?
  - Modellierung, Tools, Standards, Programme, Formate im Projekt, was/wie?
  - Verwertungsrecht/Wahrung von Persönlichkeitsrechten, was/wie?
- Daten
  - Workspace
  - Archivierung
  - Publikation
- Erhaltung von Systemen
  - Datenpräsentation (online-Portale, Repositories)
  - Begleitende Publikationen (begleitende Websites)
  - ...

# Beratungsworkflow

## Anfragen

## Interview und Planung

## Empfehlungen

## Akteure und Lösungen

bei Antragstellung  
(ab ovo)

laufende Projekte  
(in vita)

abgeschlossene Projekte  
(post mortem)

legacy data

Forschungsfrage verstehen  
Ausgangslage analysieren

Interview

Bedarfsbeschreibung

Datenmanagementplan

Metadaten/Kuratierung

Datenmanagement im Projekt

Archivierung/Speicher

Bereitstellung und Publikation

PID-Vergabe

projektspezifische  
Anwendungen

DCH-Lösungen



Language Archive  
Cologne

lokale Partner



Infrastruktureinrichtungen



Fachdatenzentren, z.B.



Beantragung von Drittmitteln für  
Entwicklungsprojekte



# DCH generischer FDM-Aktivitäten

- Schnittstelle zwischen generische Services der USB/RRZK und Forschungsprojekten/Forschenden aus den Disziplinen
- Blaupause für Implementierung FDM-Konzepts in den Fakultäten
  - Beratungsworkflow
  - Datenkuratierung
  - Datenmanagementpläne
  - Juristische Aspekte (von Verwertungs- bis Persönlichkeitsrechten)
  - Metadaten (von DublinCore über TEI bis CMDI)
  - PID-Vergabe (von ISO:24619 bis DOI)
  - OAIS konforme Architekturen
  - Repositorysysteme
  - Kodifizierte Übernahmeprozesse für legacy Data
  - dynamischen Forschungsdatensysteme

# DCH im Kontext von FDM-Aktivitäten

- Einbindung des auf Landesebene zentralen GW-Datenzentrums in die Aktivitäten der Universität
- Vertretung des Akteurs UzK auf nationale Ebene im Bereich gw-Forschungsdaten
  - Vertretung der AWK in der Akademienunion
  - CLARIN ERIC K-Centre (ZAS Berlin, SOAS London, HZSK Hamburg)
  - CLARIN-D assoziiertes B-Zentrum
  - Verbund der Datenzentren im Verband DHd (AG-Datenzentren)

# Infrastruktureinrichtungen

- Langfristige, sichere Aufbewahrung der Forschungsdaten bei einer vertrauenswürdigen Institution
  - Bereitstellen von Speicherplatz
  - Backups
    - Sicherung gegen Hardwareausfälle, menschliche Fehler, kleinere (Überspannung, Brände, Überflutung) und größere Katastrophen
      - Verschiedene Räume/ Gebäudeteile
      - Verschiedene Gebäude
      - Verschiedene Städte
    - Institut für Deutsche Sprache: Sicherung in verschiedenen Gebäuden in Mannheim

# Infrastruktureinrichtungen

- Zertifizierung
  - Data Seal of Approval (DSA)
    - Initiiert von DANS (Data Archiving and Networked Services)
    - Übergang in Core Trust Seal
  - nestor (Zertifizierung gemäß DIN-Norm 31644 „Kriterien für vertrauenswürdige digitale Langzeitarchive“)
  - ISO 16363
- Aufnahme in Infrastrukturverbänden
  - Verbund CLARIN (centre assessment)
  - Verbund DARIAH
  - ...

# Infrastruktureinrichtungen

- Zugangskontrolle
  - Wer hat Zugang?
    - Jeder, bestimmte Gruppen
    - Individuelle Regelungen
  - Wie werden Berechtigungen kontrolliert?
    - Shibboleth für CLARIN Anwendung am IDS
  - Herausforderungen:
    - Datensicherheit
    - Individuelle Regelungen

# Infrastruktureinrichtungen

- Zitierfähigkeit

Persistente Identifikatoren (PID)

- DOI
- Archival Resource Key
- Handle



# Infrastruktureinrichtungen

- Auffindbarkeit
- Metadaten
  - DC
  - CMDI
  - Herausforderung: Optimale Beschreibung der Daten mittels Metadaten
- Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH)
  - Beispiel: CLARIN Virtual Language Observatory (VLO)
  - Derzeit 37285 Metadatensätze vom IDS

# Infrastruktureinrichtungen

- Einfache Nutzbarkeit der Daten
  - Gute Beschreibung der Daten in den Metadaten
  - Sinnvolle Datenformate
  - Konsistenz der Daten und Datenformate

# Lösungsansätze für die Kuratierung von Daten aus den Digital Humanities

- Problemanalyse durch Personen mit einem Hintergrund in ‚data science‘ und der jeweiligen Fachdisziplin, z.B. Archäologie, Linguistik, Lusitanistik, Theologie, Zeitgeschichte
- Erstellung der Datenmanagementpläne
- Juristische Analyse (Möglichkeiten der Datenfreigabe)
- Daten- und Metadatenformatfestlegungen (Standards)
- Untersuchung von disziplinspezifischen und generischen Services

# Lösungsansätze für die Kuratierung von Daten aus den Digital Humanities

Archivierung und Bereitstellung der Forschungsdaten über zertifizierte, lokale Datenzentren, da

- Nähe zu den Datenerzeuger/innen
  - Akzeptanz
  - Interaktion
- Einbindung in lokale Infrastrukturen, wie Bibliotheken, Rechenzentren
  - Nutzung generischer Dienste
  - Interdisziplinäre Nutzung der Ressourcen in der Heimateinrichtung
- Einbindung in Verbünde, z.B. in CLARIN und DARIAH
  - Authentifizierungs- und Autorisierungsinfrastruktur
  - Nachweis in Katalogen

Perspektive: Einbindung in eine NFDI

# Vielen Dank an:

- Sabrina Kistner (nestor, DNB)
- Thomas Schmidt (IDS)
- Jan Gorisch (IDS)
- Peter Fankhauser (IDS)
- Denis Arnold (IDS)
- Jonathan Blumtritt (DCH an der UzK)
- Felix Rau (UzK)
- Patrick Sahle (UzK)
- Brigitte Mathiak (UzK, Gesis)

# Forschungsdaten in den Digital Humanities - Der Verbund nestor und Datenzentren am Institut für Deutsche Sprache und an der Universität zu Köln

Andreas Witt

