



## Forschungsdaten

Anke Lüdeling

Sprach- und literaturwissenschaftliche Fakultät

- 14:00 - 15:00 Uhr:  
Kurze Einführungen von Anke Lüdeling und Malte Dreyer (CMS)
- 15:00 - 16:30 Uhr:  
Postersession
- ab 16:30 Uhr:  
Diskussion

## 1 Forschungsdaten

## 2 Warum soll man Daten veröffentlichen?

- Forschung
- Replizierbarkeit
- Wiederverwendbarkeit der Daten
- Anforderung der Geldgeber und Universitäten

## 3 Technisches zur Veröffentlichung

- Formate und Architekturen
- Repositorien

**Editionen** Quellen (Bilder) und Transliterationen, kritische Apparate etc.

**Korpora** Textsammlungen (gesprochene oder geschriebene Daten, Videos), mit Metadaten und ggf. Annotationen

**Lexika** Einträge mit bestimmten Informationen, dazu ggf. Verweise auf Belege

**Daten aus Erhebungen und Experimenten** Fragebögen, Stimulusmaterial, Messdaten (Reaktionszeiten, Potentiale, etc.)

**Programme und Skripte** Software zur Suche in Daten, zur Verarbeitung & Analyse von Daten

¶ Item wer der beyfuß wurczeln  
vber die thor des hauses legt oder  
hencket / dē haufz mag nichcz v̄bels  
c̄ v̄ngeheürigkait zūgefügt wer  
den ¶ Der hoch geleert maister Ga  
lienus spricht. dz baide beyfuß rot  
vnd weiß gūt sey den frawen ge  
nūctz weñ es in not sey. Vñ auch  
fast wol bekome den die den stain  
habē in den lenden. ¶ Der maister

¶ Item wer der beyfuß wurczeln  
vber die thor des hauses legt oder  
hencket / dē haufz mag nichcz v̄bels  
od̄ v̄ngeheürigkait zūgefügt wer  
den ¶ Der hoch geleert maister Ga  
lienus spricht + dz baide beyfuß rot  
vnd weiß gūt sey den frawen ge  
nūctz weñ es in not sey + Vñ auch  
fast wol bekome den die den stain  
habē in den lenden + ¶ Der maister  
(Gart der Gesundheit, 1487)

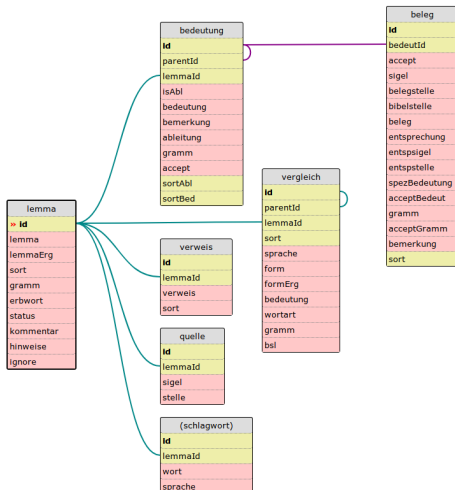
Erst	spielen	die	Dallgower	Gemeindevertreter	so	statisch	und	verzagt	wie	die
erst	spielen	der	Dallgower	Gemeindevertreter	so	statisch	und	verzagt	wie	der
--	3.Pl.Pres.Ind	Nom.Pl.Masc	Pos.*.*	Nom.Pl.Masc	--	Pos	--	Pos	--	Nom.Sg.Fem
ADV	WFIN	ART	ADJA	NN	ADV	ADJD	KON	ADJD	KOKOM	ART



Steilpass Wunder gibt es immer wieder ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie die deutsche Abwehr von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitspieler ihn aufnehmen können . Ein Befreiungsschlag ist es allerdings .

# Lexika/strukturierte Dateien (hier ALEW)

**advērija** (1) sf. 'Türpfosten, Türrahmen': BrB<sub>III</sub> [189]<sub>v11</sub> (Spr 8,34) g.pl. *idant lauktų ušu stulpų* [Glpfosten *Adwieriju*] *mana wartų* '(das er warte an den pfosten meiner thür)'; SzD<sup>1</sup> 128d<sub>3</sub> *adweria aukštine* 'podowy & podwoie', 'poftis, limen superius, superliminare'; **atvērija, atvērija** (1) sf. 'Türpfosten, Türrahmen' BrB<sub>I</sub> [242]<sub>r27</sub> (Dtn 6,9) g.pl. *rafchik ios ant aukšctieių slenkšnių* [Gl wirschutinių *atveriu*] *tawa Namų* '(solt sie vber deins Hauses pfosten schreiben)'; **advērnikas** (1) sm. 'Türhüter' DaP 449<sub>29</sub> *adwėrnikas daģuieģis* '(odźwierny Niebieŝki)'; **advērnīkē** (1) sf. 'Türhütererin' ChB<sub>I</sub> [96]<sub>a29</sub> (Joh 18,16) i. sg. *kalbejo fu ta kurij fergiejo duriu adwėrnīkie* '(sprack met de deurwaerŝter)'.  
Alit. *advērija* ist aus dem Ostslav. entlehnt, vgl. aosl. \**odvernje* snf. (aruss. *odverse*, *odverie*, russ. dial. *odvēre* 'Türrahmen'); die Variante *atverija* ist nach Sommer (1914: 25) volksetymologisch nach *atvērti* 'öffnen' aus *advērija* umgebildet. Alit. *advėrnīkas* hingegen ist hybride Lehnbildung nach apoln. *odźwierny*, *odwierzny*, *odwierzny* sm. 'Torwächter'. rf □ LEW 1.2; SEJL 3; SLA 25; 36.



	Trial	h	beding	code	RT	group
1	1	27	130	170	508	EAHP
6554	254	167	130	170	311	EAHP
13108	288	78	131	180	383	Kontrolle
19662	192	223	140	170	412	LAHP
26216	96	270	141	181	806	LALP
32770	350	288	143	180	343	LALP

**Table:** Reaktionszeitdaten (Dank an Felix Golcher und Juliane Domke)



Forschungsdaten in den Geisteswissenschaften können sehr unterschiedlich sein.

**Aufbereitungszustand:** Primärtext, Annotationen, Verweise auf externe Ressourcen | abgeschlossen oder erweiterbar

**Zweck:** Forschung, Lehre | Datengrundlage, Verarbeitung

**Quelle:** historisch, modern | handschriftlich, gedruckt, digital | Bild, Text, Audiodatei, Videodatei

**Größe:** ein Text, Internetkorpora

**Sprache:** viele – monolingual, multilingual | L1, L2

...

# Warum soll man Forschungsdaten veröffentlichen?

Oder auch: Reicht es nicht, die Forschungsergebnisse zu veröffentlichen?

- Erstellung von Ressourcen ist Forschung
- Nachvollziehbarkeit der Ergebnisse (Replizierbarkeit)
- Wiederverwendbarkeit der Daten
- Anforderungen der Geldgeber

In die Erstellung von qualitativ hochwertigen Ressourcen fließt viel Forschung und gute Ressourcen ermöglichen weitere Forschung.

Gut entworfene und nutzbare Ressourcen sind veröffentlichungswürdig - es gibt Konferenzen & Zeitschriften für die Veröffentlichung von Ressourcen.

Ressourcenveröffentlichungen werden viel zitiert!

# Transparenz und Replizierbarkeit

Bei der Erstellung *jeder* digitalen Ressource wird interpretiert.

**Erstellung der Primärebene:** Entscheidungen über Spatien, Entscheidungen über Wiedergabe von Zeichen, Entscheidungen über Wiedergabe von Audiodaten etc.

**Erstellung von Annotationen (Kategorisierungen):** Entscheidungen über Kriterien, Entscheidungen über Werkzeuge, Evaluationen etc.

**Erstellung von Verweisen:** Entscheidungen über die externen Quellen

**Analyse:** Entscheidungen über Analyseverfahren, statistische Verfahren etc.

...

Das bedeutet, dass Ergebnisse nur nachvollziehbar (und ggf. replizierbar) sind, wenn alle Daten veröffentlicht werden.

Die Erstellung von digitalen Ressourcen ist oft teuer und aufwändig. Die Daten können in vielen Fällen für weitere Forschungsfragen genutzt werden.

- innerhalb eines Faches: z.B. können mehrere (unabhängige) linguistische Forschungsfragen auf derselben Textgrundlage behandelt werden
- fächerübergreifend: z.B. kann ein historischer Text, der in der Geschichtswissenschaft digitalisiert wird, für literaturwissenschaftliche oder linguistische Zwecke ausgewertet werden

”Primärdaten als Grundlagen für Veröffentlichungen sollen auf **haltbaren und gesicherten Trägern** in der Institution, wo sie entstanden sind, für **zehn Jahre** aufbewahrt werden.”

*Vorschläge zur Sicherung guter wissenschaftlicher Praxis:  
Empfehlungen der Kommission „Selbstkontrolle in der Wissenschaft“;*

*Denkschrift DFG, 1998*

”Primärdaten als Grundlagen für Veröffentlichungen sollen auf **haltbaren und gesicherten Trägern** in der Institution, wo sie entstanden sind, für **zehn Jahre** aufbewahrt werden.”

*Vorschläge zur Sicherung guter wissenschaftlicher Praxis:  
Empfehlungen der Kommission „Selbstkontrolle in der Wissenschaft“;*

*Denkschrift DFG, 1998*

”Wenn aus Projektmitteln systematisch Forschungsdaten oder Informationen gewonnen werden, die für die **Nachnutzung** durch andere Wissenschaftlerinnen und Wissenschaftler geeignet sind, legen Sie bitte dar, ob und auf welche Weise diese für andere zur Verfügung gestellt werden. Bitte berücksichtigen Sie dabei auch - sofern vorhanden - die in Ihrer Fachdisziplin existierenden Standards und die Angebote existierender **Datenrepositorien oder Archive.**”

*DFG-Vordruck 54.01; Leitfaden für die Antragstellung Projektanträge [06/14]*

Die Pflicht zur Datenveröffentlichung kann zu einer differenzierten methodischen Diskussion in einem Fach (und als Konsequenz zur Veränderung der Fachkultur) führen.

Gutes Beispiel: Digitale Editionen (Hinweise zu Mindeststandards und Veröffentlichung seit 2015).

Kooperative, offene Kulturen entstehen/werden verstärkt - Transparenz, Möglichkeiten zur Beteiligung, Einbeziehung von Nachwuchs etc.

Gutes Beispiel: Computerlinguistik



- Abklärung der **Rechte** (Copyright und Persönlichkeitsrechte), z. B. durch CLARINs Legal Help Desk  
<http://de.clarin.eu/en/training-helpdesk/legal-helpdesk.html>
- genaue **Dokumentation** der Daten erforderlich (dies betrifft die Erstellung und alle Verarbeitungsschritte, insbesondere auch Evaluationen)
- **Versionierung** des Korpus, Bezugnahme auf die verwendete Version, Aufbewahrung alter Versionen
- freie Veröffentlichung des Korpus unter einer passenden **Lizenz**:  
Daten und Programme sollten immer unter einer möglichst offenen Lizenz abgelegt werden, sonst können sie nicht (gefahrlos) weiter genutzt werden. Für Daten sind oft die Creative Commons-Lizenzen nützlich, für Programme beispielsweise die Gnu Public License oder Apache-Lizenzen (siehe <http://opensource.org/licenses>).

- Frei** wie “für umsonst”. Nicht jedeR hat ausreichend Geld oder den Zugang zu Campuslizenzen.
- Frei** wie “offen und transparent”. EinE WissenschaftlerIn sollte die Möglichkeit haben, nachzuvollziehen, wie eine Ressource aufgebaut wurde oder was ein Programm tut (Ablage auf dem GitHub-Server)
- Frei** wie “nicht fest”: Freie Software und freie Ressourcen sind durch andere ForscherInnen modular erweiterbar.
- für die Statistiksoftware R (am 14.02.2018 12149 Packages)
  - unser Lernerkorpus Falko wird von vielen Studierenden für Abschlussarbeiten genutzt und dabei ständig erweitert

**Formate** Die verwendeten Formate sollten gut beschrieben und weit verbreitet sein. Es ist sinnvoll, die Daten in mehreren Formaten abzulegen - eines davon möglichst ein XML-Format (für gut beschriebene Formate gibt es oft Konvertierer, zum Beispiel SaltNPepper <http://korpling.german.hu-berlin.de/saltnpepper/>). Für viele Ressourcentypen gibt es etablierte (zertifizierte oder de facto) Standards (beispielweise von der Text Encoding Initiative <http://www.tei-c.org/index.xml> oder von der ISO).

**Architekturen** Die Daten sollten in Architekturen abgelegt werden, die offen und erweiterbar sind. Auch hier ist es sinnvoll, bereits existierende Lösungen zu nutzen und ggf. zu erweitern.

In Repositorien werden Daten aufbewahrt und - durch die Angabe von **Metadaten** - suchbar gemacht. Es gibt ein Kontinuum von sehr allgemeinen Repositorien, in denen alle möglichen Daten abgelegt werden (wie bspw. das Virtual Language Observatory <https://www.clarin.eu/content/virtual-language-observatory>), bis hin zu sehr spezifischen Repositorien, in denen nur Daten eines bestimmten Typs abgelegt werden (wie bspw. das LAUDATIO-Repositorium <http://www.laudatio-repository.org/repository/> für historische Texte).

Problem: Es gibt bisher für geisteswissenschaftliche Forschungsdaten keine **langfristig gesicherten** Repositorien, die mehr versprechen als die reine 'Aufbewahrung':

- Daten können nicht verändert werden
- Programme, mit denen die Ressourcen gelesen und analysiert werden können, werden nicht notwendigerweise gewartet

Im Moment wird eine nationale Forschungsdateninfrastruktur diskutiert (angestoßen vom Rat für Informationsinfrastrukturen) – es ist noch unklar, wie diese aussehen wird; es wird wahrscheinlich viel Geld in die NFDI fließen und das könnte bedeuten, dass es dann keine (weniger) Ressourcen für andere Infrastrukturmodelle geben wird.

Workshop *Wissenschaftsgeleitete Forschungsinfrastrukturen für die Geistes- und Kulturwissenschaften in Deutschland* gestern, weitere Workshops

- Forschungsdaten** An der Fakultät gibt es viele digital vorliegende Forschungsdaten in ganz unterschiedlichen Formaten.
- Veröffentlichung** Es gibt gute Gründe, digitale Ressourcen frei zu veröffentlichen. Für die technischen Aspekte (Formate, Rechte) kann man Lösungen finden.
- Langfristige Speicherung & Zugänglichkeit** Wir brauchen die Unterstützung der Universität bei der nachhaltigen Speicherung und der Sicherung der Zugänglichkeit.

“Die in einem System integrierte Verarbeitung und Darstellung heterogener Datenmengen aus ganz unterschiedlichen Quellen, erfasst mit verschiedenen Instrumenten, erfordert einen langen Prozess an Transformation, Speicherung und Übermittlung. Dieser Prozess muss bewusst und nachvollziehbar gestaltet werden, damit die erzeugten Daten ihre wissenschaftliche Aussagekraft behalten und für die Auswertung zugänglich bleiben. Das ist die Aufgabe des Forschungsdatenmanagements.

Es findet in dem Bewusstsein statt, dass lokale Lösungen Bestandteil einer übergreifenden Forschungsdateninfrastruktur sein müssen. Das Forschungsdatenmanagement muss so gestaltet werden, dass Datenzugriff und -auswertung unabhängig vom Datenerzeuger möglich wird und bleibt.”

*Büttner, St., Hobohm, H. & Müller, L. (Hg.) Handbuch Forschungsdatenmanagement. Bock und Herchen Verlag. Bad Honnef. S.29*

**Vernetzung** innerhalb der Fakultät und gerne auch mit KollegInnen aus anderen Fakultäten, die mit ähnlichen Daten arbeiten (Geschichte, Musikwissenschaft, Kulturwissenschaften, Psychologie etc.).

**Forschungsdatenstrategien** der Fakultät entwickeln und Forschungsdatenstrategien der Universität beeinflussen (Forschungsdaten gibt es nicht nur in den Naturwissenschaften)

→ Wir haben eine Forschungsdatengruppe an der Fakultät (eingesetzt vom FR). Ein Ergebnis ist die Ausstattung mit virtuellen Servern. → Wir planen die Einrichtung einer/eines Forschungsdatenbeauftragten an der Fakultät, Kofinanzierung mit dem CMS



# Vielen Dank!

Roland Baumgarten, Markus Egg, Roland Meyer, Katja Münster,  
Muriel Norde, Erika Thomalla, Anne Wolfsgruber