

Language Data: Curation Projects and Integration Scenarios



The aim of the curation projects was to identify key language data and resources and then to sustainably integrate those data into the CLARIN-D infrastructure. In the area of German studies three curation projects made it possible to add three types of data into the CLARIN-D centres: **historical texts** of the 15th-19th centuries, **spoken language data**, and **computer mediated communication (CMC)**.

In addition, we developed **integration scenarios** that can be used for further **integration use cases** in a sustainable way.

Curation Project 1 / UC 1 Integration of Historical Texts

Integration and Enhancement of Historical Textual Resources of the 15th-19th centuries

Integration und Aufwertung historischer Textressourcen des 15.-19. Jahrhunderts

Curation Project 2 / UC 2 Integration of Spoken Language Data

GeWiss: A Comparable Corpus of Spoken Scientific Language Gesprochene Wissenschaftssprache (GeWiss-Korpus)

The curation project aimed to integrate all existing published and unpublished resources of the GeWiss project and to make them available to the academic community at large, using a format compatible with CLARIN-standards. GeWiss provides the academic community with the first freely available corpus resource for the comparative analysis of spoken academic language.

The curated resources are made available through the existing web portal of the GeWiss project (free registration). In addition, the Institut für Deutsche Sprache (IDS) will make the data available through its CLARIN service centre infrastructure in the near future.

The curation project has also delivered an integration scheme for spoken language data that can be used in a sustainable way.

GeWiss 00:07 00:08

Gesprochene Wissenschaftssprache

Projekt Korpus Recherche Publikationen Kontakt

PO_DE_046

Nichtige Hinweise Important notes... Wähle westslawisch...

Um das Audio abzuspielen nutzen Sie bitte die Pfeile, oder klicken Sie auf die kleinen Play-Tasten im Transkript, um einzelne Segmente abzuspielen. N.B.: AUS Datensatzgründen wurden Sprecherbezeichnungen pseudonymisiert und in den Aufnahmen maskiert (verrauscht).

Navigation
Korpusübersicht
Liste der Kommunikationen
Download PDF
Sprachauswahl
Anmerkungen (x) Verbalspur (v) S

Rollen
Beisitzer B
Prüfer P
Profiling D

Transkript

[1]

JS_0215 [v] (wah...) ihnen in bezug auf das problem aussprache beziehungsweise artikulation 'h' (.) so bieten (0.4) also

[2]

JS_0215 [v] (0.2) ren sie die so vergleichen (.) also jetzt hm erkenlos das sys them (.) ja

[3]

JS_0215 [v] sondern (0.2) was ah bieten die beiden (0.3) modellie (0.3) für (0.2) ah(.) für das thema

[4]

JS_0215 [v] artikulation oder aussprache ((schmatzt)) okay (.) also zum ersten das (0.3) levvelmodell oder levvel Jah

[5]

JS_0215 [v] JS_0252 [x] wie auch immer es genannt wird 'h' ahm (.) also das modeld wird ja sehr oft Jah -rezipiert einfach auch in

Curation Project 3 / UC 3 Integration of CMC Data

ChatCorpus2CLARIN: Integration of the Dortmund Chat Corpus into CLARIN-D

ChatCorpus2CLARIN: Integration des Dortmunder Chat-Korpus in CLARIN-D

The screenshot shows a search result page for the term "smile" within the "Dortmunder Chat-Korpus". The results are sorted by date from 1998 to 2006. The first result is a link to a page titled "Ich stelle uns goldenen Schallplatten her ... smile ..." from the year 2000-01-12. The second result is a link to a page titled "... das sieht heute wieder einmal genial aus - grosses kompliment ! meine frage : welches sind dein markenzeichen ? sind es deine schuhe strafmessen ? ... bitte lass sie dran - bitte nicht wegschneiden - ok ? smile und küsses ; [MALE-PARTICIPANT-A20_]" from the year 2000-01-25.

Bibliography and Webliography

Publications CP1 / UC 1

Geyken, Alexander/Thomas Gloning (2015): "A living text archive of 15th-19th century German. Corpus strategies, technology, organization." In: Jost Gippert (eds.): Historical corpora: challenges and perspectives. Tübingen: Narr, 165-180.
Pfundt, Anna (2017): "Frauenwahlrecht? Oder Damenvahlrecht? Oder doch ein allgemeines Wahlrecht? – Zum Wortgebrauch in der Diskussion um das Frauenwahlrecht um 1900." In: Im Zentrum Sprache, 2.11.2017,
https://sprache.hypotheses.org/542 (last retrieved May 28, 2018).

Thomas, Christian/Frank Wiegand (2015): "Making great work even better. Appraisal and digital curation of widely dispersed electronic textual resources (c. 15th-19th centuries) in CLARIN-D." In: Jost Gippert (eds.): Historical corpora: challenges and perspectives. Tübingen: Narr, 181-196. [Online Version 31.10.2012 <https://edoc.bbaw.de/opus4-bbaw/frontdoor/index/index/docId/2005> [last retrieved November 30, 2017]].

Publications CP2 / UC2
Endruch, Christian/Cordula Meißner/Adriana Slavcheva (2012): "The ColWise Corpus: Comparing Spoken Academic German, English and Polish." In: Thomas Schmidt/Kai Wörner (eds.): *Multilingual corpora and multilingual corpus analysis*.

Fandrych, Christian/Co
Amsterdam: Benja

Fandrych, Christian/Cordula Meißner/Franziska Wallner (eds.) (2017): *Gesprochene Wissenschaftssprache – digital. Verfahren zur Annotation und Analyse mündlicher Korpora*. Tübingen: Stauffenburg.

Publications CP3 / IUC 3

Beißwenger, Michael/Harald Lüngen/Jan Schallaböck/John H. Weitzmann/Axel Herold/Pawel Kamocki/Angelika Storrer/Julia Wildgans (2017): "Rechtliche Bedingungen für die Bereitstellung eines Chat-Korpus in CLARIN-D: Ergebnisse eines Rechtsgutachtens." In: Michael Beißwenger (eds.): Empirische Erforschung internetbasierter Kommunikation. Berlin/Boston: De Gruyter, 7–46. (doi: <https://doi.org/10.1515/9783110567786-002>) (last retrieved November 30, 2017).