

Korpora und Korpus-Technologie im Programmbereich „Mündliche Korpora“ am IDS

Im Programmbereich „Mündliche Korpora“ der Abteilung Pragmatik des Instituts für Deutsche Sprache Mannheim (IDS) werden Korporagesprochener Sprache erstellt und archiviert (FOLK und AGD) sowie Werkzeuge zur Erstellung und Bearbeitung von Korpora entwickelt (FOLKER und OrthoNormal). Der Publikation von Korpora und der Recherche darin dient die Datenbank für Gesprochenes Deutsch (DGD2). Für die Dokumentation von „Best Practice“-Verfahren zur Korpuserstellung und für ein Handbuch zur Arbeit mit mündlichen Sprachkorpora steht das Gesprächsanalytische Informationssystem (GAIS) zur Verfügung.

Archiv für Gesprochenes Deutsch (AGD)

Das Archiv für Gesprochenes Deutsch (AGD) ist die zentrale Sammelstelle für Korpora des gesprochenen Deutsch. Die Korpora werden im AGD aufbereitet und für Forschung und Lehre bereitgestellt. Im AGD werden Tonaufnahmen digitalisiert und maskiert, Dokumentationsdaten nach einem einheitlichen Schema erfasst und Transkripte in nachhaltigen, nutzbaren Formaten konvertiert. Aus dem Gesamtbestand von über 40 Korpora sind derzeit 17 Bestandskorpora (mit ca. 9000 dokumentierten Ereignissen) für die DGD2 fertiggestellt. Aktuell werden im AGD z. B. Korpora zu „Jugendsprache“ und „Deutsch in Australien“ bearbeitet. Unter der Adresse agd@ids-mannheim.de beantwortet das AGD Serviceanfragen zu seinen Beständen und zur Korpus-Technologie.

Web: <http://agd.ids-mannheim.de>

Kontakt: agd@ids-mannheim.de



Forschungs- und Lehrkorpus Gesprochenes Deutsch (FOLK)

Mit dem Forschungs- und Lehrkorpus Gesprochenes Deutsch (FOLK) baut das IDS ein kontinuierlich wachsendes Korpus auf, um Gesprächsdaten aus unterschiedlichen Bereichen des gesellschaftlichen Lebens (Arbeit, Freizeit, Bildung, Medien) im deutschen Sprachraum via Internet zugänglich zu machen. Die Datenakquise umfasst eigene Aufnahmen, eine Übernahme aus anderen IDS-Projekten (z. B. „MapTask“-Aufnahmen und Interviews aus „Deutsch heute“), Datenspenden aus externen Projekten (u. a. von „Gesprochene Wissenschaftssprache kontrastiv“ (GeWiss) und „Sprachvariation in Norddeutschland“ (SiN)), sowie Rundfunk-Mitschnitte. Bei der Datenaufbereitung werden Aufnahmen und Transkripte maskiert, die Metadaten nach DGD-Schema eingegeben, Transkriptionen nach GAT-Konventionen für das Minimaltranskript mit dem Editor FOLKER erstellt und die Transkripte mit OrthoNormal orthographisch normalisiert sowie mit dem TreeTagger lemmatisiert. 99 Ereignisse (ca. 70 Aufnahmestunden) sind für die DGD2 aufbereitet worden; weitere 110 Aufnahmestunden liegen vor und werden bearbeitet.

Web: <http://agd.ids-mannheim.de/folk.shtml>



FOLKER und OrthoNormal

FOLKER, Akronym für FOLK-Editor, ist als Transkriptionseditor optimiert für die Arbeit in FOLK, mit GAT und mit FOLK-Daten. FOLKER basiert auf EXMARaLDA und ist hiermit sowie mit Praat, ELAN etc. interoperabel. FOLKER steht zum kostenlosen Download bereit; seit März 2009 gibt es über 2500 Registrierungen. FOLKER wird in der Version 1.2 eine Syntaxkontrolle für GAT-Minimal- und Basistranskripte anbieten; dafür wird derzeit schon ein Preview angeboten. OrthoNormal ist ein Tool zur orthographischen Normalisierung in FOLK, allgemeiner: zum manuellen Annotieren auf Tokenebene.

Web: <http://agd.ids-mannheim.de/folker.shtml>



Datenbank für Gesprochenes Deutsch (DGD2)

Die Datenbank für Gesprochenes Deutsch (DGD2) präsentiert eine Auswahl aus den AGD-Korpora und den gesamten FOLK-Bestand, ermöglicht ein Browsen in Metadaten, Aufnahmen, Transkripten, Zusatzmaterial sowie mehrere Recherchemöglichkeiten: eine Volltextrecherche in Metadaten und Transkripten sowie eine tokenbasierte Recherche in Transkripten, jeweils mit Rückgriff auf Transkripte und Aufnahmen. Die DGD bietet zudem einen Download ausgewählter kompletter Datensätze.

Web: <http://dgd.ids-mannheim.de>