

Silke Reineke/Arnulf Deppermann/Thomas Schmidt  
(Mannheim/Basel)

# Das Forschungs- und Lehrkorpus für Gesprochenes Deutsch (FOLK)

Zum Nutzen eines großen annotierten Korpus gesprochener  
Sprache für interaktionslinguistische Fragestellungen

**Abstract:** Der Beitrag illustriert die Nutzung des Forschungs- und Lehrkorpus Gesprochenes Deutsch (FOLK) für interaktionslinguistische Fragestellungen anhand einer exemplarischen Studie. Zunächst werden die Stratifikation (Datenkomposition) des Korpus, das zugrundeliegende Datenmodell und dessen Annotations Ebenen sowie Typen von Untersuchungsinteressen vorgestellt, für die das Korpus nutzbar ist. Im Hauptteil wird Schritt für Schritt anhand einer Studie zur Verwendung des Formats *was heißt X* in der sozialen Interaktion gezeigt, wie mit FOLK relevante Daten gefunden und analysiert werden können. Abschließend weisen wir auf einige Vorsichtsmaßnahmen bei der Benutzung des Korpus hin.

## 1 Das Forschungs- und Lehrkorpus Gesprochenes Deutsch im Überblick

### 1.1 Aufbau

Während in den letzten drei Jahrzehnten national wie international zunehmend eine große Menge an schriftlichen Korpora verschiedenster Art für die sprachwissenschaftliche Forschung zur Verfügung steht, sind wissenschaftsöffentlich zugängliche mündliche Korpora Mangelware. Dies gilt umso mehr für Audio- und Videoaufnahmen natürlicher Interaktion. Vorhandene Korpora sind meist nur für die Angehörigen einer bestimmten Institution zugänglich, sie sind nicht maschinell erschließbar, beinhalten keine Videodaten und benutzen keine interoperablen Datenformate. Darüber hinaus sind sie für eine breitere wissenschaftliche Nutzung auch nicht autorisiert. Vor diesem Hintergrund wurde am Leibniz-Institut für Deutsche Sprache (IDS) im Jahre 2006 die Entscheidung getroffen, ein großes, wissenschaftsöffentlich verfügbares Korpus verbaler Interaktionen aufzubauen, das aktuellen korpustechnologischen Standards entspricht. Dies ist das

<https://doi.org/10.1515/9783111085708-005>

Forschungs- und Lehrkorpus Gesprochenes Deutsch (FOLK), das auch auf internationaler Ebene, mit Ausnahme weniger Korpora,<sup>1</sup> eine einzigartige, wissenschaftsöffentliche Ressource für konversationsanalytische und interaktionslinguistische Untersuchungen darstellt. Dieser Beitrag illustriert die Nutzung von FOLK und der im Korpus enthaltenen, annotierten Daten gesprochener Sprache für die Bearbeitung interaktionslinguistischer Fragestellungen am Beispiel einer konkreten Untersuchung. Wir charakterisieren zunächst die allgemeine Anlage von FOLK (Abschn. 1.1), gehen dann auf die Stratifikation (Datenkomposition) des Korpus (Abschn. 1.2) und auf das Datenmodell von FOLK ein (Abschn. 1.3). In Abschnitt 2 diskutieren wir unterschiedliche Typen von interaktionslinguistischen Fragestellungen, für die FOLK geeignet ist. Der Hauptteil des Aufsatzes, Abschnitt 3, ist dann der Darstellung des konkreten Vorgehens der Korpusnutzung anhand einer exemplarischen Studie, der Untersuchung der Verwendung des Formats *was heißt X* in der sozialen Interaktion, gewidmet. Wir zeigen hier Schritt für Schritt, wie mit Hilfe von FOLK relevante Daten gefunden und analysiert werden können. Im abschließenden Fazit weisen wir auf einige Vorsichtsmaßnahmen, die bei der Benutzung des Korpus zu beachten sind, hin und wir resümieren den Nutzen der Verwendung von Daten aus FOLK im Vergleich zur (alleinigen) Benutzung von selbst erhobenen Daten für eine interaktionslinguistische Untersuchung (Abschn. 4).

Das Forschungs- und Lehrkorpus für Gesprochenes Deutsch wird seit 2008 am Leibniz-Institut für Deutsche Sprache aufgebaut. Das Ziel des FOLK-Korpus ist die wissenschaftsöffentliche Bereitstellung einer großen, nach aktuellen Standards erschlossenen und breit diversifizierten Datenbasis zur Untersuchung gesprochener Sprache in natürlicher Interaktion. Zielgruppe des Korpus sind Forschende, Lehrende und Studierende aus Gesprächsforschung bzw. Konversationsanalyse und Interaktionaler Linguistik sowie aus Korpuslinguistik und angrenzenden Fachgebieten. Das Korpus bildet ‚natürliche‘ Interaktionen ab, d. h. solche Interaktionen, die nicht durch Forschende elizitiert wurden, also auch ohne deren Zutun stattgefunden hätten. Es ist das größte und korpustechnologisch avancierteste Korpus mit den meisten Erschließungsmöglichkeiten unter den im Archiv für Gesprochenes Deutsch (AGD) verfügbaren Gesprächskorpora. FOLK wird der wissenschaftlichen Öffentlichkeit (d. h. Forschenden, Leh-

---

1 Auf internationaler Ebene vergleichbar sind das französischsprachige Corpus CLAPI (<http://clapi.ish-lyon.cnrs.fr/>, Stand: 29.8.2022) sowie für das amerikanische Englisch das Santa Barbara Corpus of Spoken American English (<https://www.linguistics.ucsb.edu/research/santa-barbara-corpus>, Stand: 29.8.2022) und für das australische Englisch das Griffith Corpus of Spoken Australian English (<https://www.ausnc.org.au/corpora/gcscouse>, Stand: 29.8.2022).

renden und Studierenden) über die Datenbank für Gesprochenes Deutsch (DGD) via <dgd.ids-mannheim.de> zur Verfügung gestellt.

Seit Beginn seines Aufbaus im Jahre 2008 ist das FOLK-Korpus stetig gewachsen. Es umfasst in der Version vom Mai 2021 insgesamt 374 Gesprächsaufnahmen (mit rund 314 Stunden Audio-Aufnahmen, davon 140 Stunden auch als Video-Aufnahmen) sowie vollständige Transkriptionen aller Aufnahmen mit knapp 3 Mio. Tokens. Mit dem nächsten Release 2022 wird das Korpus um weitere 26 Ereignisse mit einem Umfang von knapp 22 Stunden wachsen.

Alle Gesprächsaufnahmen werden nach zeitgemäßen Standards erschlossen (d. h., sie werden vollständig transkribiert, linguistisch annotiert und mit Metadaten zu Gespräch und den Beteiligten dokumentiert, vgl. Schmidt 2017), bevor sie in der DGD zur Verfügung gestellt werden. Dort können Nutzerinnen und Nutzer gezielt über die insgesamt 2.9 Millionen transkribierten Wörter auf 4 Annotationsebenen (Transkription, Normalisierung, Lemmatisierung, Part-of-Speech-Tagging (POS)) recherchieren. Darüber hinaus können auch über die Metadaten von Gesprächen und Gesprächsbeteiligten gezielt Gespräche mit bestimmten Merkmalen ausgesucht und gefiltert werden (z. B. institutionelle Gespräche oder Gespräche mit Beteiligung bestimmter Altersgruppen u. v. m.). Schließlich ist auch freies Explorieren der Daten möglich, indem man sich einzelne Gespräche über die Funktionalitäten der Audio-, Video- und Transkriptionanzeige ansieht bzw. anhört.<sup>2</sup>

Im Gegensatz zu Gesprächskorpora, die für Projekte mit spezifischen Forschungsfragen erhoben werden und deren Erhebung zeitlich begrenzt ist, wird das FOLK-Korpus kontinuierlich ausgebaut. Es speist sich dabei insbesondere aus zwei Quellen: Zum einen werden Aufnahmen im FOLK-Projekt bzw. vom IDS ausgehend erhoben, z. B. im Rahmen von Seminaren an der Universität Mannheim, die von FOLK-Mitarbeiter/-innen geleitet werden, oder durch direkte Aufrufe zur Teilnahme an Erhebungsaktionen. Zum anderen stammen die Aufnahmen in FOLK aus Erhebungen, die Interaktionsforschende für ihre Projekte durchgeführt haben und von denen sie Aufnahmen für das FOLK-Korpus und damit für die Community spenden. FOLK ist daher ein Korpus von der wissenschaftlichen Fachgemeinschaft für die wissenschaftliche Fachgemeinschaft. Aus diesem Grund ruft das FOLK-Korpus auch kontinuierlich zu ‚Datenspenden‘ für das Korpus auf. Die Projektmitarbeitenden beraten im Gegenzug gerne schon

---

<sup>2</sup> Für Hilfestellung in der Benutzung des FOLK-Korpus in der DGD sei hier auch verwiesen auf die Hilfe-Seiten der DGD, dort finden sich unter „Hilfe > Materialien“ Hinweise und Links zu Videotutorials, Handreichungen und Publikationen zu den Funktionalitäten der DGD (siehe [https://dgd.ids-mannheim.de/dgd/pragdb.dgd\\_extern.help](https://dgd.ids-mannheim.de/dgd/pragdb.dgd_extern.help), Stand: 29.8.2022).

in frühen Planungsphasen einer Erhebung zu allen Aspekten der Erhebung von Audio-/Videodaten und Metadaten. Im Rahmen einer verbindlichen Kooperation können dann auch Daten, die für das FOLK-Korpus gespendet werden, im Projekt zeitnah transkribiert und nach zeitgemäßen Standards der Audio- und Videodokumentation sowie Annotation aufbereitet werden. Die aufbereiteten Daten werden dann dem erhebenden Forschungsprojekt zur Verfügung gestellt und erst nach einer Sperrfrist (sog. ‚Embargo‘) im FOLK-Korpus in der DGD veröffentlicht.

## 1.2 Stratifikation

Die Stratifikation des Korpus folgt dem langfristigen Ziel des Ausbaus von FOLK hin zu einem Referenzkorpus des Gesprochenen Deutsch. Im Unterschied zu anderen Korpora, die zumeist nur einen (und oftmals wissenschaftlich elizitierten) Gesprächstyp beinhalten, stand für FOLK die Leitvorstellung Pate, die volle Breite des kommunikativen Haushalts (Luckmann 1986) der deutschen Gegenwartsgesellschaft im Korpus abzubilden (siehe Deppermann/Hartung 2011). FOLK soll daher möglichst breit diversifiziert sein (vgl. dazu ausführlich Kaiser 2018). Für die Stratifikation von FOLK ist dabei das Konzept des Interaktions- bzw. Gesprächstyps zentral. Interaktionsdomäne, Lebensbereich und Aktivität sind die wesentlichen Eigenschaften von Gesprächen, die möglichst breit abgedeckt werden sollen. Die sekundären Stratifikationsparameter sind soziodemographische Merkmale der Sprecher/-innen wie Geschlecht, Alter, regionale Herkunft und Bildungsstand (siehe Abb. 1).<sup>3</sup> Das Ideal wäre eine ausgewogene bzw. noch besser: eine repräsentative Abdeckung sämtlicher möglicher Parameter-Kombinationen. Dies implizierte z. B. Daten verschiedener Dienstleistungsgespräche jeweils mit Beteiligten aller Herkunftsregionen, innerhalb jeder Region auch von allen Altersgruppen jeweils jeden Geschlechts und aller Bildungsniveaus – und dies genauso für alle anderen Gesprächstypen.

Es liegt auf der Hand, dass dieses maximale Ziel utopisch und aus verschiedenen (aufwandsbezogenen, datenschutzrechtlichen u. a.) Gründen nicht realisierbar ist. Das vorrangige Ziel ist für uns daher zunächst die Abdeckung der ein-

---

<sup>3</sup> Die allermeisten Sprecher/-innen in FOLK haben Deutsch als Erstsprache, eine geringe Anzahl von Sprecher/-innen hat Deutsch als Zweitsprache. Da sich FOLK als Korpus des usuellen gesprochenen Deutsch versteht, werden zwar L2-Sprecher/-innen als Teil der gesellschaftlichen Realität mit erfasst; sie bilden aber keine Sprechergruppe, auf die die Stratifikation von FOLK abzielt.

zelen Parameter-Ausprägungen, also z. B. die Abdeckung möglichst jeder Altersgruppe und jeder regionalen Herkunft, erst einmal ungeachtet des jeweiligen Gesprächstyps. Umgekehrt streben wir eine möglichst große Variation von Gesprächstypen innerhalb der Interaktionsdomänen und Lebensbereiche an, hier zunächst ungeachtet der genauen Sprecher/-innen-Verteilung darin. Sowohl private als auch institutionelle und öffentliche Typen mündlicher Interaktion sollen in möglichst großer Breite im Korpus enthalten sein (siehe Abb. 2).

INTERAKTIONSTYP	INTERAKTIONS-DOMÄNEN		Institutionell				Öffentlich		Anderes
	Privat		Bildung	Verwaltung	Interprofessionelle Kommunikation	Vereinsleben	Politik	Unterhaltung	[Anderes]
	[Privat]		Religion/Kirche	Kultur (Unterhaltung, Kunst, Sport)	Dienstleistungen	Medizin	Wissenschaft	Wirtschaft	
AKTIVITÄTEN	Nicht aktivitätsgeleitet	Renovieren, Urlaubsplanung, ...	Meeting, Fahrstunde; ...			Mediation; Panel-Diskussion; ...		Experimentelles Spiel; Interview; ...	

Primäre Parameter: Interaktion



Geschlecht	männlich			weiblich		anderes
Alter	0-18		19-39	40-65	66-99	
Region	nord-west	mittel-west	süd-west	nord-ost	mittel-os	süd-ost
Bildung	hoch		mittel	niedrig		

Sekundäre Parameter: Sprecher

Grafik 11: Schema Stratifikation

Abb. 1: Schema der Stratifikation von FOLK (Kaiser 2018, S. 543)

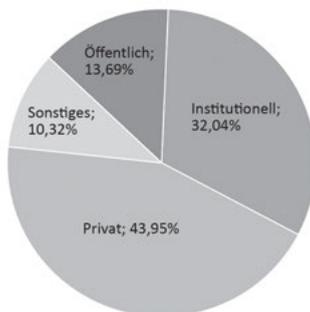
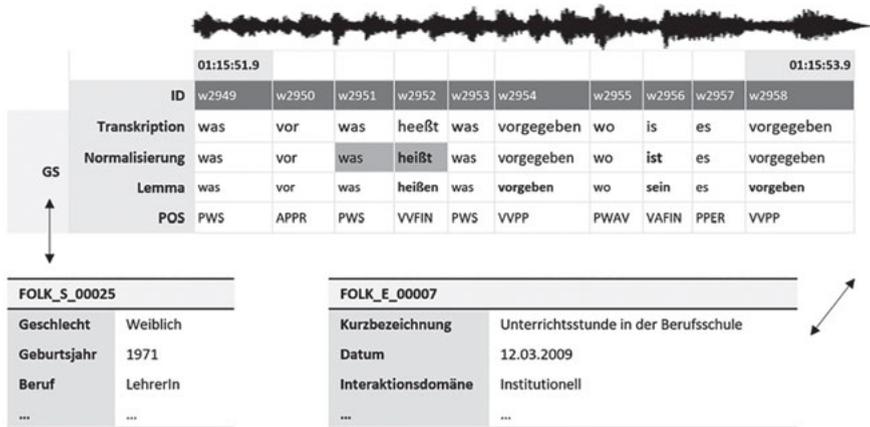


Abb. 2: Verteilung der Daten in FOLK (Version 2.16, 2021) nach Interaktionsdomänen

### 1.3 Datenmodell und Interoperabilität – Das FOLK-Korpus in der DGD

Wenn man in der DGD eine Belegstelle gefunden hat – zum Beispiel nach einer strukturierten Suche (siehe Abschn. 3.1.1) oder durch Exploration eines Transkriptes – gibt die Datenbank Zugriff auf sehr reichhaltige Informationen, die mit dieser Belegstelle im Zusammenhang stehen. Dies wird ermöglicht durch ein Datenmodell, das die Audio- und Videodaten, den transkribierten Text, Annotationen auf Token-Ebene sowie die Metadaten in einer feinteiligen Struktur miteinander verknüpft (Abb. 3).



**Abb. 3:** Zugrundeliegendes Datenmodell der für die DGD aufbereiteten FOLK-Daten

Dazu gehören als Erstes Zeitmarken, die im Abstand von etwa 2 bis 5 Sekunden vom Transkript in die zugrundeliegende Aufnahme zeigen. Damit lässt sich zu jeder Belegstelle präzise die zugehörige Stelle der zugehörigen Audio-/Videoaufnahme ansteuern. Nach einer Suchanfrage ist für jede Belegstelle in der *Keyword-in-Context*-Ansicht (KWIC) immer der linke und rechte Kontext sichtbar. Darüber hinaus kann auch immer der weitere sequenzielle Transkriptkontext in variabler Ausdehnung angezeigt werden, also die vorhergehenden und folgenden Äußerungen des betreffenden Sprechers oder der weiteren Interaktionsteilnehmerinnen (Abb. 4).



Abb. 4: Anzeige von KWIC und Transkript in der DGD

Auf Token-Ebene werden lexikalische Tokens, also Wörter, unterschieden von Tokens, die spezifisch für das Mündliche sind. Letztere umfassen Pausen, Ein- und Ausatmen sowie Beschreibungen non-verbalen Verhaltens (wie z. B. Räuspern). Jedem Wort-Token werden drei Annotationen zugeordnet: die orthographische Normalisierung führt abweichende literarisch transkribierte Formen auf eine orthographische Normalform zurück, im Beispiel etwa dialektales „heeßt“ auf „heißt“.<sup>4</sup> Auf Grundlage der Normalisierung wird dann zum einen mit dem TreeTagger (Schmid 1995) und der in Westpfahl (2020) für mündliche Daten trainierten Parameterdatei automatisch eine Lemmatisierung vorgenommen, die z. B. der flektierten Form „heißt“ den Infinitiv „heißen“ zuordnet. Im selben Annotationsvorgang werden zum anderen die Wort-Tokens mit einem Part-Of-Speech-Tag, z. B. VVFIN für finites Vollverb, versehen. Wir verwenden hierzu das Stuttgart-Tübingen-Tagset (STTS) in der Version 2.0 (Westpfahl et al. 2017), die Erweiterungen und Anpassungen für gesprochene Sprache enthält.

Für alle Belegstellen sind die zugehörigen Metadaten abrufbar – dies betrifft sowohl Informationen zum betreffenden Gespräch, etwa dessen Interaktionsdomäne, als auch Informationen zu den beteiligten Sprecher/-innen, z. B. deren Alter, sprachliche Herkunft, formalen Bildungsstand etc.

All diese Bestandteile sind mit geeigneten IDs versehen. Diese ermöglichen beispielsweise ein präzises Adressieren eines spezifischen Transkriptausschnitts, oder auch Standoff-Annotationsmethoden, also die Möglichkeit, dem Transkriptdokument zusätzliche analytische Information mittels Verweisen auf annotierte Elemente hinzuzufügen.

<sup>4</sup> Dies erfolgt im FOLK-Aufbereitungsworkflow nach der abgeschlossenen Transkription. Hierzu werden die Transkripte zunächst automatisch mit dem Tool OrthoNormal (Teil des EXMARALDA-Pakets) normalisiert. In einer anschließenden händischen Durchsicht wird diese initiale automatische Normalisierung korrigiert und verbessert.

Dieses Datenmodell ist in mehrerlei Hinsicht interoperabel mit anderen technischen Lösungen, die für mündliche Daten gebräuchlich sind. Textdaten und Metadaten werden als XML-Daten serialisiert, deren formale Korrektheit mit entsprechenden Schemata (also Grammatiken, die zulässige XML-Strukturen beschreiben) überprüft werden kann. Die Transkripte sind in dieser Form kompatibel mit den wichtigsten Transkriptions- und Annotationstools wie ELAN, Praat und EXMARaLDA. Sie orientieren sich außerdem am ISO-Standard „Transcriptions of Spoken Language“ (ISO 2016).

Praktisch wird so zum einen eine größere Flexibilität im Workflow für den Korpusaufbau ermöglicht, da wir z. B. ohne größeren Aufwand auch Transkripte integrieren können, die ursprünglich in ELAN, Praat oder EXMARaLDA angefertigt wurden. Zum anderen können umgekehrt auch Transkripte oder Ausschnitte aus der DGD direkt in diesen Formaten heruntergeladen werden. Dies unterstützt Arbeitsformen, in denen Forscher/-innen Transkriptausschnitte lokal weiterbearbeiten möchten, also z. B. zusätzliche Annotationen der Videodaten in ELAN oder akustische Analysen in Praat vornehmen. Weiterhin bietet die DGD auch die Möglichkeit, Transkripttextausschnitte von Belegstellen in einer Form (z. B. als HTML-Datei) zu exportieren, mit der sie dann in Standard-Office-Programmen wie Word oder Excel eingelesen und dort weiter bearbeitet werden können. Dies ist sowohl für weitere Auswertungen und Aufarbeitungen der Daten (z. B. Erstellung multimodaler Transkripte) als auch für Publikationen überaus relevant.

## 2 Nutzung des Korpus für verschiedene Typen von Fragestellungen

Das FOLK-Korpus wird breit genutzt. Etwa zwei Drittel der Datenbank-Anfragen der rund 14.500 registrierten DGD-Nutzer/-innen, die insbesondere aus der gesprächsanalytischen und korpuslinguistischen Forschung und angrenzenden Fachgebieten stammen, beziehen sich auf FOLK. Das FOLK-Korpus wird oft als alleinige Untersuchungsgrundlage genutzt, aber ebenso auch als Referenz- oder Vergleichskorpus zu selbst erhobenen Daten. Unter <https://www.ids-mannheim.de/prag/muendlichekorpora/bibliographie-folk/> stellt das FOLK-Projekt eine Bibliografie von Arbeiten, die auf FOLK-Daten basieren, online zur Verfügung. Die Bibliografie wird regelmäßig erweitert um die von Forschenden gemeldeten Neuerscheinungen entsprechender Arbeiten.

Je nach Forschungsfrage und Ziel und auch nach Offenheit im eigenen Forschungsparadigma kann man das Korpus auf unterschiedliche Weise erschlie-

ßen. Die grundlegenden Wege der Erschließung der Korpusdaten entsprechen den verschiedenen methodischen Ansätzen für unterschiedliche Fragestellungen im Bereich der Interaktionalen Linguistik. Ausgangspunkte können sein:

- **Ein formbasierter Zugang:** Hier bildet die Suche nach dem Lemma oder der orthographischen Transkription einer sprachlichen Form oder einer Kombination von sprachlichen Formen, eventuell auch zusammen mit ihrer Position im Turn, den Ausgangspunkt. Dieser Ausgangspunkt ist der am häufigsten gewählte. Bei diesem Vorgehen kann man am intensivsten die verschiedenen Datenerschließungsinstrumente des Korpus nutzen. Beispiele sind: *komm* (Proske 2014), *(das) stimmt* (Betz 2015), *irgendwie* (Günthner/König 2015), Intonation und Bedeutung (Moroni 2015), *ich weiß nicht* (Helmer/Reineke/Deppermann 2016), Progressivkonstruktion (Katelhön 2016), *hesitation markers* (Wieling et al. 2016), *weiß nich, keine Ahnung* (Bergmann 2017), Argumentrealisierung (Deppermann/Proske/Zeschel 2017), *machen* (Kress 2017), direkte Rede (Katelhön/Moroni 2018), *sehr sehr* (Staffeldt 2018), *halt, eben* (Torres Cajo 2019), *ich dachte* (Deppermann/Reineke 2020), Adressierung (Droste/Günthner 2020), *oder* (König 2020), *adjective intensifiers* (Stratton 2020), Interrogative (Gubina 2021).
- **Ein sequenz- oder handlungsbasierter Zugang:** Hier stehen interpretative Größen wie bestimmte Handlungen, Sequenztypen oder Interaktionsaufgaben und -probleme am Anfang der Korpuserschließung. Im Gegensatz zum formbasierten Zugang kann nach solchen Phänomenen nicht durch gezielte maschinelle Anfragen gesucht werden, sondern die Gesprächsereignisse im Korpus müssen händisch gesichtet werden. Als Heuristik kann hier dienen, Gesprächstypen zu sichten, in denen die entsprechende Handlung bzw. Sequenz erwartungsgemäß häufiger vorkommen wird, z. B. Instruktionen in pädagogischen Interaktionen, Klatsch-Sequenzen in privaten Konversationen oder Zeigeaktivitäten im gemeinsamen praktischen Handeln oder bei Museumsführungen. Beispiele sind: Ironie (Moroni 2016), *suspended assessments* (Aldrup et al. 2021), Interpretationen (Zinken/Küttner 2022).
- **Ein Zugang ausgehend von einem spezifischen Interaktionstyp:** Hier werden Ereignisse im Korpus eines bestimmten Interaktionstyps (z. B. Dienstleistungsgespräche) ausgewählt und analysiert. Beispiele sind: Rettungsübungen (Deppermann 2014), Fahrschule (Deppermann 2018), WG-Casting (Bies 2020), Schlichtung Stuttgart 21 (Helmer/Deppermann 2022; Reineke 2016).
- **Ein Zugang ausgehend von Metadaten:** Hier beginnt man mit einem Filter, z. B. nach spezifischeren Gesprächsmerkmalen oder der Erhebungsform als Suchheuristik. Man grenzt die Suche z. B. nach vorhandenen Videoaufnah-

men ein, wie es für einen videoanalytischen Zugang (z. B. Deppermann/Gubina 2021; Deppermann/Schmidt 2021) notwendig ist.

Die Anwendungsbeispiele in der obigen Liste sind nicht exhaustiv, aber illustrativ für die Verteilung typischer Zugangswege zum Korpus.

### **3 Ein Beispiel für die Nutzung von Korpusfunktionalitäten für eine interaktionslinguistische Studie: Die Untersuchung von *was heißt X***

#### **3.1 Vorgehen und Korpusfunktionalitäten**

Ausgangspunkt unserer hier vorzustellenden Studie ist ein formbasierter Zugang. Im Rahmen von Untersuchungen zu Praktiken der Bedeutungskonstitution (vgl. Deppermann 2020, im Dr. a und b) haben uns Verwendung und Funktionen des Formats *was heißt X* interessiert. Das Format fiel unter anderem im Rahmen von Untersuchungen zu Definitionspraktiken auf (vgl. Helmer 2020), denn es wird häufig als Format zur Initiierung von Definitionen verwendet.

Wir schildern in den folgenden Abschnitten, wie wir diese Fragestellung mithilfe der Daten des FOLK-Korpus und der Funktionen in der DGD zur Erschließung der Korpusdaten bearbeitet haben. Dabei gehen wir insbesondere auf die methodischen Schritte ein, die eng mit der Nutzung von FOLK und den Funktionalitäten der DGD verbunden sind. In einer kurzen Ergebnisübersicht über Formen und Verwendungen des Formats (vgl. Abschn. 3.2) werden wir diese vorstellen. Sequenzanalytische Details und die zugehörigen detaillierten Einzelfallanalysen werden wir hier nicht darstellen, da sie nicht im engeren Sinne auf die Korpusfunktionalitäten zugreifen. Diese Punkte werden in Deppermann (im Dr. a) in Bezug auf unterschiedliche Praktiken der Verwendung von *was heißt X* genauer ausgeführt; in Deppermann/Reineke (in Vorb.) gehen wir vertieft auf die Relevanz von und den Umgang mit Metadaten in der Analyse ein.

Der Forschungsprozess einer formbasierten Untersuchung von FOLK-Daten mithilfe der DGD lässt sich schematisch wie in Abbildung 5 darstellen:



**Abb. 5:** Forschungsprozess einer formbasierten Untersuchung von FOLK-Daten mithilfe der Funktionalitäten der DGD

Anhand dieser Schritte im Forschungsprozess werden wir unser Vorgehen nun schildern.

### 3.1.1 Suchen und Filtern

Für die initiale Suchanfrage muss man entscheiden, wie man nach der interessierenden linguistischen Struktur sucht. Wir haben in diesem Fall über die Menüpunkte *Recherche* > *Tokens* in der DGD im Feld *Normalisiert* nach „heißt“ gesucht (Abb. 6).

**Abb. 6:** Suche nach „heißt“ im Feld *Normalisiert*

Man könnte natürlich auch mit einer Suche nach „was“ beginnen. In diesem Fall bekäme man aber mehr Treffer als in der Datenbank aktuell ausgegeben werden können, da eine Grenze von 10.000 Treffern überschritten wird. In solchen Fällen bietet die DGD dann nur die Möglichkeit, mit einer Zufallsstichprobe weiterzuarbeiten. Wenn nur die häufige Form selbst interessiert, kann das unproblematisch sein. In unserem Falle würden aber von der erwartungsgemäß viel geringeren Anzahl der Belege des Formats *was heißt X* durch die eingeschränkte Zufallsstichprobe vorab Treffer ausgeschlossen. Daher empfiehlt es sich für die initiale Suche ggfs. über Ausprobieren verschiedener Suchanfragen aus der interessierenden Phrase den Suchbegriff zu wählen, der voraussichtlich am seltensten ist.<sup>5</sup> Bei

<sup>5</sup> In unserem Fall konnten wir durch die Suche nach „heißt“ trotz der technologischen Einschränkungen zielführend ohne Verlust zu allen relevanten Belegstellen gelangen. Für spezifi-

Ausdrücken, die in mehreren syntaktischen Kategorien vorkommen (z. B. *was, eben, bitte*), empfiehlt es sich, den entsprechenden POS-Tag zur Suche zu benutzen. Allerdings muss man beachten, dass die POS-Annotation nicht immer völlig zuverlässig ist. Daher empfiehlt es sich bei diesem Vorgehen, stichprobenartig Belege des gleichen Lemmas, die aber mit einem anderen POS-Tag versehen sind, zu sichten, um die Zuverlässigkeit der Annotation abzuschätzen. Suchanfragen sollten einerseits möglichst präzise, andererseits aber offen genug formuliert sein, dass auch unerwartete Varianten der interessierenden Struktur gefunden werden können. Z. B. besteht die Gefahr, bei der Suche nach transkribierten Formen relevante phonetische Varianten nicht zu finden; bei der Suche nach normalisierten Formen werden keine anderen morphosyntaktischen Varianten des Lemmas gefunden.

Durch unsere Suche nach „heißt“ im Feld *Normalisiert* wurden 3.026 Belege gefunden (Abb. 7). Bevor wir die Suche weiter eingrenzen, speichern wir das Ergebnis mit einem Klick auf das Disketten-Symbol in der Funktionsleiste über der KWIC-Ansicht ab. Grundsätzlich ist zu empfehlen, bei der Arbeit mit der DGD die jeweiligen Suchergebnisse immer durch Notizen zu dokumentieren und zusätzlich für alle wesentlichen Zwischenschritte in der DGD zu speichern. So kann im weiteren Bearbeitungsprozess der eigenen Studie immer wieder auf frühere Zwischenstände zurückgegriffen werden und händische Filter-Ergebnisse und ursprüngliche Suchanfrage-Ergebnisse können eingesehen und nachvollzogen werden.

---

sche Suchanfragen, in denen diese maximale Treffergrenze nicht zu umgehen ist oder für die mehrschrittiges und sehr komplexes Filtern der Ergebnismenge notwendig ist, empfiehlt es sich, über das Tool „Zu-Recht“ des Projektes und der Anwendung „ZuMult“ eine CQP-Anfrage zu stellen (vgl. zu Anwendungsmöglichkeiten von ZuMult auch Fandrych/Wallner (in diesem Band) sowie Frick/Wallner/Helmer (in Vorb.) zu Nutzungsmöglichkeiten von ZuRecht). Dies ist einfach durch Verlinkungen von DGD und der Plattform ZuMult möglich. Der Zugang zu ZuMult erfolgt mit den Anmelde Daten zur DGD.

ÜBER DIE DGD    BROWSING    RECHERCHE    DOWNLOAD    MEINE DGD    HILFE    ABMELDEN

POSITION    **TOKEN**    KONTEXT    METADATEN    ANZEIGE

Transkribiert: z.B. 'kannst'    Normalisiert: heißt  
 Lemma: z.B. 'können'    POS: z.B. 'VMFIN'

Reguläre Ausdrücke    Suche starten    CQP

*Recherche 0 Tokens*

Suchergebnis gespeichert.

Ergebnisse 1 bis 20 von 3026 ( 3026 / 0 ans-abgewählt)    Seite 1 von 152

	Sprechereignis	Sprecher	Treffer
<input checked="" type="checkbox"/>	1	FOLK_00001_01 LB	wie heißt n des
<input checked="" type="checkbox"/>	2	FOLK_00001_01 LB	das heißt
<input checked="" type="checkbox"/>	3	FOLK_00001_01 LB	noch ne schutzschaltung gegen irgendwelche kurzschlüsse ... heißt
<input checked="" type="checkbox"/>	4	FOLK_00001_01 LB	das heißt ein mess widerstand
<input checked="" type="checkbox"/>	5	FOLK_00001_01 LB	denk des kennen sie die schaltung des heißt
<input checked="" type="checkbox"/>	6	FOLK_00001_01 LB	genau das heißt unser magnetfeld bricht zusammen in der primärspule
<input checked="" type="checkbox"/>	7	FOLK_00001_01 LB	das heißt die selektive prüfung des steuergerätes
<input checked="" type="checkbox"/>	8	FOLK_00001_01 LB	das heißt diese spulen würden überlasten
<input checked="" type="checkbox"/>	9	FOLK_00001_01 LB	des problem begegnet uns immer wieder des heißt wir
<input checked="" type="checkbox"/>	10	FOLK_00001_01 LB	das heißt sie haben problem mit der einspritzung
<input checked="" type="checkbox"/>	11	FOLK_00001_01 LB	das heißt
<input checked="" type="checkbox"/>	12	FOLK_00001_01 LB	des heißt der normale monteur macht des natürlich oft nicht
<input checked="" type="checkbox"/>	13	FOLK_00001_01 LB	das spannungssignal des halbgewers das heißt gemessen mit einem voltmeter
<input checked="" type="checkbox"/>	14	FOLK_00001_01 LB	des heißt sich beküme ein zündimpuls
<input checked="" type="checkbox"/>	15	FOLK_00001_01 LB	des heißt hier steht entsprechend alles
<input checked="" type="checkbox"/>	16	FOLK_00001_01 LB	das heißt jetzt kommt die selektive funktion des
<input checked="" type="checkbox"/>	17	FOLK_00001_01 LB	net unbedingt will s selektiv prüfen das heißt selektiv heißt der halbgewer
<input checked="" type="checkbox"/>	18	FOLK_00001_01 LB	will s selektiv prüfen das heißt selektiv heißt der halbgewer
<input checked="" type="checkbox"/>	19	FOLK_00001_01 LB	das heißt was müsse mer mache immer wenn mer prüfen hier
<input checked="" type="checkbox"/>	20	FOLK_00001_01 LB	ja das heißt die haltschicht ist aktiv

Ergebnisse 1 bis 20 von 3026 ( 3026 / 0 ans-abgewählt)    Seite 1 von 152

**Abb. 7:** KWIC-Ansicht nach der initialen Suche im Reiter *Token* nach „heißt“ im Feld *Normalisiert*

Um nun nur noch Fälle unseres Zielformats zu bekommen, haben wir weiter gefiltert mithilfe des Reiters *Kontext* (Abb. 8):

RECHERCHE    DOWNLOAD    MEINE DGD    HILFE    ABM

**KONTEXT**    METADATEN    ANZEIGE

Normalisiert: was    Kontext: 2 Tokens    links

POS: z.B. 'VMFIN'    Skopus: Betrag

Reguläre Ausdrücke    Kontext filtern

**Abb. 8:** Filtern nach linkem Kontext im Reiter *Kontext* mit Suche nach „was“ im Feld *Normalisiert*

Dort tragen wir im Feld *Normalisiert* „was“ ein und wählen „2 Tokens“ „links“ mit dem Skopus „Beitrag“ aus. In den meisten Fällen wird „was“ wohl nur 1 Token links von „heißt“ stehen. Ein etwas weiterer Skopus ist aber anzuraten, um falsche Negative zu vermeiden und eventuell sogar unbekannte Varianten zu entdecken. Denn so kann man Belege mit Phänomenen gewinnen, die man normgrammatisch nicht erwarten würde. Gesprochensprachlich kann auch bspw. ein „äh“ oder ein Wortabbruch innerhalb des Syntagmas erscheinen. Umgekehrt ist der Skopus weiterhin relativ eng gewählt, um zu vermeiden, dass nicht zu viele falsch-positive Belege händisch ausgefiltert werden müssen. Auch hier ist es immer hilfreich, über verschiedene Suchanfragen auszuprobieren, ob sich die Belegliste verändert, und jeweils kursorisch zu prüfen, ob man unnötig viele falsch-positive Belege mit einem bestimmten Skopus produziert.

Durch diesen Kontext-Filter werden dann die nicht passenden Ergebnisse durchgestrichen (siehe ausgefilterte und durchgestrichene Ergebnisse am Beispiel des nächsten Schrittes in Abb. 9). Hier sollte man ebenfalls kursorisch gegenprüfen, ob die richtigen Belege ausgefiltert wurden. Man kann dann die ausgefilterten Ergebnisse mit einem Klick auf den Papierkorb in der Funktionsleiste über der KWIC löschen.

Im nächsten Schritt blieben in unserer Untersuchung 296 Belege übrig. Wir gehen diese händisch auf die Form hin durch, hauptsächlich transkriptbasiert. Wir schauen aber immer wieder in den Kontext und hören die Daten an. Das geht ganz schnell, in dem man ein Beispiel ausklappt, es abspielt und dann behält oder abwählt. Hier kann man nach Bedarf den jeweils angezeigten Transkriptkontext in der KWIC durch Klicken auf das Lupen-Symbol erweitern (vgl. Beleg Nr. 17 im Beispiel in Abb. 9) oder bei Bedarf die weiteren Funktionen der DGD nutzen und z. B. einen Ausschnitt zusammen mit dem gegebenenfalls zugehörigen Video aufrufen oder ein vollständiges Transkript (z. B. in einem anderen Tab) anzeigen. Man sieht hier z. B. in Abbildung 9, dass auch falsch-positive Ergebnisse mit Verbletzstellung automatisch in unsere Suche eingeschlossen waren. Diese haben wir dann händisch getilgt, so z. B. „was es heißt“ hier in Zeile 1.<sup>6</sup> Insgesamt blieben nach der händischen Durchsicht noch 275 Belege übrig.

---

<sup>6</sup> Dieses Ergebnis ist ein zu erwartender falsch-positiver Beleg durch den Skopus von 2 Tokens links.

The screenshot shows the 'Recherche' (Search) interface for the FOLK corpus. At the top, there is a search bar containing the query 'norm\_heißt\_was\_2\_l' and a 'Transkriptausschnitt berechnet' (Transcript excerpt calculated) checkbox. Below the search bar, the results are displayed in a table with columns for 'Sprechereignis' (Speech event), 'Sprecher' (Speaker), and 'Treffer' (Hit). The results are numbered 1 through 20, with checkboxes for selection. A detailed view of result 0074 is shown below the table, displaying the original sentence and the KWIC extraction.

Sprechereignis	Sprecher	Treffer
1	FOLK_00004_01 TE	steht.auch.dunter.was.es heißt
2	FOLK_00004_01 GS	was heißt das dann für t seinen zukünftichn erfolg
3	FOLK_00004_01 AB	was heißt n kognitiv
4	FOLK_00004_01 GS	genau danke was heißt kognitiv
5	FOLK_00007_01 GS	was heißt eigentlich kognitiv
6	FOLK_00007_01 GS	was vor was heißt was vorgegeben wo is es vorgegeben
7	FOLK_00013_01 CJ	was heißt das
8	FOLK_00014_01 CJ	heißt un willst du wissen was ängstlich heißt auf türkisch
9	FOLK_00014_01 CJ	doch ne vom papa und was heißt zu hause auf türkisch weißt du es
10	FOLK_00014_01 CJ	des kennstsch du genau und was heißt spielen auf türkisch weißt du des
11	FOLK_00015_01 CH	zuerst mal ah die frage stellen was heißt
12	FOLK_00015_01 CH	was heißt sprachmelodie was heißt auch spr zeit verlauf was dass wir
13	FOLK_00015_01 CH	was heißt sprachmelodie was heißt auch spr zeit verlauf was dass wir da genauer sehen
14	FOLK_00015_01 CH	methode f für ah erstsprachenweb wa was heißt das metho methode für die datengewinnung
15	FOLK_00015_01 CH	kön können sie mir zeigen was das heißt wie ich zu meinen daten komme
16	FOLK_00020_01 EM	was heißt des
17	FOLK_00021_01 JZ	was heißt nur

Below the table, a detailed view of result 0074 is shown:

- 0074 (0.34)
- 0075 SK sag mehr
- 0076 NI [dann sag mehr (,)] dann is es deiner
- 0077 JZ [was heißt nur]
- 0078 (0.35)
- 0079 XM1 ((Lachensatz))
- 0080 (0.66)

At the bottom of the interface, there are more search results (18-20) and a footer indicating 'Ergebnisse 1 bis 20 von 295 ( 294 / 2 aus- / abgewählt) Seite 1 von 15'.

**Abb. 9:** KWIC-Ansicht der händischen Durchsicht der Belege nach Abwahl falsch-positiver Belege

### 3.1.2 Export der Suchergebnisse

Dies waren die wichtigsten Schritte der initialen Suchanfrage und Belegsichtigung in der DGD. Es sind innerhalb der DGD noch vielfältige Filter und Sortierschritte möglich, die zur strukturierten Suche und Ergebnissicherung dienen können, die aber auch einen ersten Eindruck von formalen Varianten, Kookkurrenzen und Verteilungen im Korpus geben können. So kann man beispielsweise nach rechtem Kontext filtern (in diesem Fall beispielsweise nach „das“ oder, offener, nach ‚Pronomen 1 Token rechts‘ von ‚heißt‘). Man kann auch über den Reiter *Metadaten* nur Belege eines bestimmten Gesprächstyps oder bestimmter Teilnehmerkonstellationen anzeigen lassen etc. Da knapp 300 Belege gut einzeln unsortiert durchgesehen werden können, haben wir unsere übrigen Belege in diesem Fall aber als nächstes erneut als KWIC-Suchergebnis in der DGD gespeichert und als XML-Tabelle mit den transkribierten Belegstellen exportiert. Dieser Export ist über das Excel-Symbol in der Funktionsleiste oberhalb der KWIC möglich (siehe Abb. 9 in der oberen Leiste).

Es wird dann eine Tabelle ausgegeben, die man z. B. mit Excel öffnen kann und in der man die einzelnen Belege ansehen kann (Abb. 10). Die Tabelle enthält zusätzlich basale Metadaten-Informationen der jeweiligen Belegstellen: Transkript-ID und Sprecher-ID, den linken Kontext, das gesuchte Keyword und den rechten Kontext jeder Belegstelle sowie einen Link, der direkt zum Datum in der DGD führt.

	A	B	C	D	E	F
1	transcript-id	speaker-id	left-context	match	right-context	dgd-link
2	FOLK_E_00004_SE_01_T_FOLK_S_00025	was	heißt	das dann für t seinen zukünftigen	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
3	FOLK_E_00004_SE_01_T_FOLK_S_00014	was	heißt	n kognitiv	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
4	FOLK_E_00004_SE_01_T_FOLK_S_00025	genau danke was	heißt	kognitiv	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
5	FOLK_E_00007_SE_01_T_FOLK_S_00025	was	heißt	eigentlich kognitiv	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
6	FOLK_E_00007_SE_01_T_FOLK_S_00025	was vor was	heißt	was vorgegeben wo is es vorgegeb	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
7	FOLK_E_00013_SE_01_T_FOLK_S_00030	was	heißt	das	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
8	FOLK_E_00014_SE_01_T_FOLK_S_00030	heißt un willst du wisser	heißt	auf türkisch	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
9	FOLK_E_00014_SE_01_T_FOLK_S_00030	doch ne vom papa und w	heißt	zu hause auf türkisch weißt du des	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
10	FOLK_E_00014_SE_01_T_FOLK_S_00030	des kennst du genau u	heißt	spielen auf türkisch weißt du des	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
11	FOLK_E_00015_SE_01_T_FOLK_S_00034	zuerst mal äh die frage s	heißt		<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
12	FOLK_E_00015_SE_01_T_FOLK_S_00034	was	heißt	sprachmelodie was heißt auch spr	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
13	FOLK_E_00015_SE_01_T_FOLK_S_00034	was heißt sprachmelodi	heißt	auch spr zeit verlauf was dass wir d	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
14	FOLK_E_00015_SE_01_T_FOLK_S_00034	methode f für äh erstspr	heißt	das metho methode für die dateng	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
15	FOLK_E_00020_SE_01_T_FOLK_S_00034	was	heißt	des	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
16	FOLK_E_00021_SE_01_T_FOLK_S_00039	was	heißt	nur	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
17	FOLK_E_00021_SE_01_T_???	was	heißt	hier schon	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
18	FOLK_E_00021_SE_01_T_FOLK_S_00043	hao halt ho ho was	heißt	hier du kaufst	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
19	FOLK_E_00024_SE_01_T_FOLK_S_00048	ja was	heißt	jetzt nett	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
20	FOLK_E_00026_SE_01_T_FOLK_S_00047	halt in die a und e was	heißt	n des e is erziehungshilfe odder wa	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
21	FOLK_E_00026_SE_01_T_FOLK_S_00049	was	heißt	n des	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
22	FOLK_E_00026_SE_01_T_FOLK_S_00048	was	heißt	hau hat man des in der grundschul	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	
23	FOLK_E_00027_SE_01_T_FOLK_S_00182	was	heißt	ich muss will	<a href="http://dgd.ids-mannheim.de/">http://dgd.ids-mannheim.de/</a>	

Abb. 10: Ausschnitt aus der Belegliste der exportierten XML-Tabelle der Suchergebnisse

### 3.1.3 Qualitative Analyse und Kodieren

Diese Excel-Tabelle war die Grundlage für die Organisation der Belegstellen und die Dokumentation unserer weiteren Analyse und der Kodierung der Belege nach Funktionen. Im Rahmen der Kodierung wird jeder Beleg nochmals auf seine Zugehörigkeit zur Kollektion geprüft. Nach Tilgen weiterer falsch-positiver Belege und uninterpretierbarer Abbrüche verblieben 250 Fälle für unsere Untersuchung.

Das Arbeiten mit der Belegsammlung in Excel ermöglicht es, beliebig viele Kodier-Kategorien in zusätzlichen Spalten einzufügen und jeden Beleg entsprechend zu kodieren. Der Kodierung voraus geht jedoch eine extensive qualitative Analyse von Einzelfällen, die wir hier nicht nachzeichnen. Wir setzen hier vielmehr an dem Punkt an, an dem wir in den qualitativen Analysen bestimmte, robust erscheinende Analysekategorien und ihre Varianten gefunden haben. Grundsätzlich ergibt sich jedoch während jeder Kodierung in der Auseinander-

setzung mit den Daten oft die Notwendigkeit, zuvor zugewiesene Codes eines Beleges abzuändern oder die Codes selbst zu modifizieren. Auch Beobachtungen, die zunächst offen notiert werden, können Anlass zu neuen Codes geben. Die exportierte Tabelle ermöglicht durch die Verlinkung jeden Belegs den schnellen Zugriff auf sämtliche zugehörige Originaldaten in der DGD, was den Analyseprozess unterstützt. Dies ist für die interaktionslinguistische Arbeit auch essenziell, denn wir analysieren jeden Beleg stets im sequenziellen Kontext und mit Video- bzw. Audioausschnitten. Dieses Vorgehen ermöglicht zugleich die Datenbank-unabhängige Dokumentation der eigenen Analyse-Schritte und Kodierungen, ohne auf den Zugriff auf die Originaldaten verzichten zu müssen. Durch die Hyperlinks, die von der XML-Tabelle direkt die jeweilige Belegstelle in einem Transkript ansteuern, ist dies für Nutzende ohne Aufwand möglich. Durch das Ansteuern des Belegs kann man dann alle dem Beleg zugrundeliegenden Originaldaten (Audiodatum, ggf. Videodatum, Transkript und weitere Annotationen) anzeigen und nach Bedarf erweitern und so den Beleg im sequenziellen Kontext der Methodik entsprechend in die Analyse miteinbeziehen. Zusätzlich zu den Anzeigemöglichkeiten von bild-/tonaligniertem Transkript und Metadaten in der DGD ist es ebenfalls möglich, den Beleg über die Oberfläche von „ZuViel“ anzusehen (wiederum durch Verlinkung) und so weitere Anzeige-Möglichkeiten und Darstellungsmöglichkeiten zu nutzen. Dies ist für die Korpora FOLK und GeWiss („Gesprochene Wissenschaftssprache kontrastiv“, Fandrych/Meißner/Slavcheva 2012) möglich; siehe zu den Anwendungen der Plattform „ZuMult“ auch Fandrych/Wallner (in diesem Band).

### 3.2 Analyse/Ergebnispräsentation

An einem kleineren Datensatz hatte Susanne Günthner schon 2015 responsive (fremd- und selbstresponsive) Verwendungen von *was heißt X* untersucht und dort problematisierende und klarifizierende Verwendungen gefunden (49 Fälle: fremdresponsive: Problematisierung, Klarifizierung, Mischformen; selbstresponsive; siehe auch De Stefani (im Dr.) zum analogen italienischen Format (*Che cosa vuol dire X*, für das er die gleichen Verwendungen feststellt).

Im Rahmen einer qualitativen Vorstudie (N = 90) haben wir feinere Unterscheidungen getroffen, da uns spezifische Praktiken der Bedeutungskonstitution interessiert haben. Dabei haben wir folgende Verwendungen identifiziert:

- Definition (allgemeine Bedeutung, z. B. *was heißt theistisch*),
- Spezifikation (Bedeutung im lokalen Kontext, z. B. *was heißt fast*),
- Übersetzung (z. B. *was heißt n des auf deutsch*),
- Konsequenz (z. B. *was heißt das für den finanzierungsrahmen*)

- Adäquatheit (sprachliche bzw. sachliche Angemessenheit von X, z. B. *was heißt in der Diskussion*).

Diese Verwendungen bestätigten sich in der Hauptstudie. Ihre Zuordnungs- und Abgrenzungskriterien wurden aber anhand der größeren Datenbasis verfeinert. Dazu kamen weitere Unterscheidungen für einzelne Verwendungen, die z. B. die Frage betreffen, wer eine mit *was heißt X* formulierte Reparatur-Initiierung produziert (selbst vs. fremd) und wer dann die folgende Reparatur ausführt (selbst vs. fremd). Bei der Verwendungsform ‚Spezifikation‘ wird somit z. B. unterschieden, ob die eigene Äußerung oder die eines anderen spezifiziert werden soll.

Im Folgenden stellen wir die Verwendungsformen an Beispielen vor. Die häufigsten Verwendungsformen sind Adäquatheit und Spezifikation (Tab. 1).

**Tab. 1:** Häufigkeiten der Verwendungsformen von *was heißt X* (N = 250)

Verwendungsform	Häufigkeit	Häufigkeit in Prozent
Definition	30	12%
Spezifikation	79	32%
Übersetzung	26	10%
Konsequenz	12	5%
Adäquatheit	103	41%

Wir beginnen mit einem typischen Beispiel für die Verwendung von *was heißt X* zur Elizitierung einer Definition, also einer verallgemeinerbaren Bedeutungsangabe (30 Fälle von 250, 12%). Der Ausschnitt stammt aus einem Expertenvortrag des Sprechers WW im Rahmen der Schlichtungsgespräche zu Stuttgart 21. WW wird vom Schlichter HG um eine Definition des Ausdrucks *Bemessungsquelldruck* gebeten:

- (1) FOLK\_E\_00069\_SE\_01\_T\_01\_c436<sup>7</sup>
- 01 WW: °hh und diese (.)  
 ABdichtungsbauwe[rke sind, hh° ]  
 -> 02 HG: [was **was heisst beMES**]Sungs

<sup>7</sup> Die Transkriptüberschrift unserer Beispiele enthält jeweils den Korpusnamen (hier „FOLK“), die Ereignis-ID (hier „E\_00069“), die Sprechereignis-Nummer (hier „SE\_01“) die Transkript-Nummer (hier „T01“) sowie die ursprüngliche Contribution-Nummer der fokalen Phrase aus der DGD (hier „c436“).

- (.) QUELL (.) druck;  
 03 (0.4)  
 04 HG: [(sagen sie) was so-]  
 05 WW: [bemessungsQUE ]LLduck- h°  
 06 HG: was IS des;  
 07 WW: °h das is DER druck, h°  
 08 (.) gegen DEN,  
 09 (.) °h das BAUwerk,  
 10 (.) °h beMESSen wird; h

Der Ausdruck, dessen Bedeutung hier von Heiner Geißler erfragt wird, war nicht mündlich vorerwähnt. In diesem Fall stand er auf einer Vortragsfolie, die wir in diesem Fall auch im Videodatum in FOLK sehen können (Abb. 12):

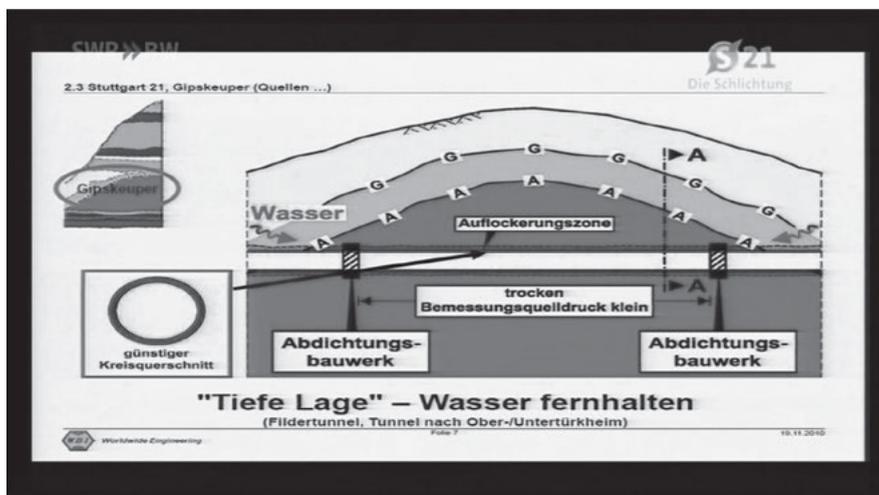


Abb. 11: Folie aus dem Expertenvortrag von WW<sup>8</sup>

Die von HG erbetene Bedeutungsklärung wird von WW in den Zeilen 07–10 als allgemein gültige Definition der Bedeutung eines Fachbegriffs formuliert. Die Reformulierung durch HG in Zeile 06 mit „was ist des“ desambiguiert seine

<sup>8</sup> Der Ausschnitt (1) beginnt bei 00:20:07. Die Vortragsfolie ist in der Aufnahme während des Vortrages sichtbar zu den Zeitpunkten 00:18:29–00:19:37, 00:19:50–00:20:06 sowie 00:20:11–00:20:18.

was-*heißt*-X-Äußerung nachträglich als Frage nach einer Definition. Dies könnte dazu dienen, andere, alternative Lesarten (wie z. B. Adäquatheit, siehe unten) auszuschließen.

### Exkurs: Transkripte in der DGD und Analysestandards

Eine Bemerkung zu den Transkripten, wie wir sie hier in der Ergebnispräsentation verwenden: Die Transkripte in der DGD sind erstellt nach den Konventionen von cGAT, den „Konventionen für das computergestützte Transkribieren in Anlehnung an das Gesprächsanalytische Transkriptionssystem 2 (GAT2)“ (Schmidt/Schütte/Winterscheid 2015). Diese Konventionen sind ausgelegt auf die Transkription großer Mengen von Audio-/Videodaten für maschinell durchsuchbare Korpora und angelehnt an die Konventionen für GAT2-Minimaltranskriptniveau. Das erfüllt den Zweck der Durchsuchbarkeit und weiteren Annotation – auf basaler Ebene. Die in FOLK bzw. der DGD veröffentlichten Transkripte sind jedoch keine Analyse- und Publikationstranskripte. Für die Analyse ist der Mindeststandard, auf Basis-Niveau von GAT2 (Selting et al. 2009) nachzutranskribieren und – je nach Analyseinteresse – auch Konventionen des Feintranskripts oder zusätzliche multimodale Annotationen zu benutzen. Abbildung 13 illustriert, wie der Unterschied zwischen einem Ausschnitt aus der DGD auf cGAT-Niveau und dem überarbeiteten GAT2-Basis-Transkript aussehen kann. Dies zeigt sich einerseits an der zusätzlichen prosodischen Notation (hier: Fokusakzente und Endintonationsnotation), es kann darüber hinaus aber auch zu Abweichungen in der Segmentierung und damit zu veränderten Zeilenzählungen oder modifiziertem Wortlaut führen. Zur Zitation empfiehlt es sich deshalb, eine feste Referenz zum Ausschnitt in der DGD zu zitieren. Dazu eignet sich die Beitrags-Nummer (Contribution-Number, siehe „c“ im Transkriptkopf). Man kann hierzu entweder die Beitrags-Nummer des fokalen Beitrags verwenden oder die erste Zeile des in einer Publikation zitierten Ausschnitts. Entsprechend können dann von der DGD abweichende Zeilennummerierungen verwendet werden, ohne dass die Referenz zum DGD-Beleg dadurch nicht mehr nachzuvollziehen wäre.

<pre>(1) FOLK_E_00069_SE_01_T_01_DF_01_c436  0434 WW [(schmatzt)] "hh und diese abdichtungsbauwe[rke sind hh"] 0435 HG [was was heißt bemes]sungs (.) quell (.) druck 0436 (0.38) 0437 WW [bemessungs]quell druck h° 0438 HG [sagen sie was soll] 0439 HG was is des 0440 WW "h das is der druck h° (.) gegen den "h das bauwerk (.) "h       bemessen wird h°</pre>	<pre>(1) FOLK_E_00069_SE_01_T_01_DF_01_c436  01 WW: °hh und diese (.)       Abdichtungsbauwe[rke sind, hh° ] -&gt; 02 HG: [was was heißt bemes]sungs       (.) QUELL (.) druck; 03 (0.4) 04 HG: [(sagen sie) was so-] 05 WW: [bemessungs]QUE ]LLdruck- h° 06 HG: was IS des; 07 WW: °h das is DER druck, h° 08 (.) gegen DEN, 09 (.) °h das BAUwerk, 10 (.) °h beMESSen wird; h</pre>
<p>cGAT Minimaltranskript Schmidt/Schütte/Winterscheid 2015 (hier: HTML-Ausgabe aus DGD)</p>	<p>GAT2 Basistranskript Selting et al. 2009</p>

Abb. 12: Vergleich von cGAT-Minimaltranskript und GAT2-Basis-Transkript

Eine zweite Verwendungsform ist *was heißt X* zur Elizitierung einer Spezifikation (79 Fälle von 250, 32%). Hier geht es um die Bedeutung von X in einem spezifischen Verwendungskontext, also nicht um eine allgemeingültige Definition. In den meisten Fällen handelt es sich dabei um die vom Sprecher im konkreten Kontext intendierte Bedeutung. Im folgenden Beispiel aus einem Pärchengespräch geht es um eine solche im lokalen Kontext gültige Bedeutung:

- (2) FOLK\_E\_00030\_SE\_01\_T01\_c901
- 01 AM: °h aber (.) GIB des doch nom ma ein?  
 02 es kann nich sein dass die nur so WEnige::-  
 03 (.) ähm ((schnalzt)) **so wenige AUFlistungen haben;**  
 04 (0.6)  
 05 PB: wie,  
 06 (3.0)  
 -> 07 PB: **was HEISST nur wenige auf[(li-) ]**  
 08 AM: [ja wir ham **nur**  
**zwei EINträge gefunden für die-**  
 09 (0.4)

Nach AMs Frage „was heißt nur wenige aufli“ (Z. 07) klärt PB sowohl die lokale Bedeutung von ‚wenige‘ (= „zwei“) als auch von ‚auflistungen‘ (= „einträge“).

*Was heißt X* kann drittens zur Elizitierung von Übersetzungen in eine andere Sprache (26 Fälle von 250, 10%) verwendet werden:

- (3) FOLK\_E\_00331\_SE\_01\_T\_03\_DF\_01\_c207
- 01 (2.0)  
 -> 02 RA: **was heißt denn SMOOTH?**  
 03 (0.6)  
 04 CA: **wei:ch;**  
 05 RA: weich.  
 06 (2.0)  
 07 RA: °hh

Der Bezugsausdruck wurde hier bereits 38 Sekunden früher (c177) von CA produziert, die aus einem englischen Rezept vorliest: „mix till smooth“.

*Was heißt X* wird viertens dazu verwendet, um eine Konsequenz zu formulieren (12 Fälle von 250, 5%). Hier geht es um die Frage, welche Auswirkungen X auf Y hat bzw. welche Handlungskonsequenzen angesichts von X zu ergreifen sind.

Spezifisch für diese Verwendung ist, dass sie an die formale Variante *was heißt X für Y* gebunden ist. Das folgende Beispiel stammt aus einer Unterrichtsinteraktion in einer Berufsschule, in der angehende Ausbilder auf ihre Rolle vorbereitet werden. Es geht hier um die Frage, welche Konsequenz eine mangelnde Ausbildung in Methodenkompetenzen für Auszubildende haben kann.

(4) FOLK\_E\_00004\_SE\_01\_T01\_c1097

01 GS: wenn (0.21) sie (.) wie geFORDert,  
 02 (.) diese neuen (.) lehrmethoden an:: WENden-  
 03 wie SELBSTgesteuertes lernen,  
 (0.34)  
 04 fördern sie ja damit AUDOmatisch die  
 methodenkompetenz.=ne,  
 (0.77)  
 05 wenn sie ihm so was NICHT anbieten,  
 (.) wird er (HIER),  
 (0.79)  
 06 gegen NULL gehen;  
 -> 07 **was HEISST das dann für t seinen zukünftichen erfolg?**  
 08 (0.2)  
 09 AB: **null.**  
 10 (0.34)  
 11 AB: (**nichts.**)  
 12 GS: **sieht a NIT so gut aus; (.) ja?**

Die letzte und häufigste Verwendungsform von *was heißt X* ist die Infragestellung der Adäquatheit des Ausdrucks X im vorangehenden Kontext (103 Fälle von 250, 41%). Hier wird mit *was heißt X* angezeigt, dass die Verwendung eines Ausdrucks sachlich, sprachlich oder stilistisch inadäquat war. So kann X zum Beispiel zu extrem oder zu unpräzise sein oder falsche Inferenzen nahelegen. Mit *was heißt X* wird in diesen Fällen eine Reparatur eingeleitet und projiziert, dass der Ausdruck durch einen passenden zu ersetzen ist. Die Einschätzung von Inadäquatheit wird von den Interaktionsteilnehmer/-innen mit der Konstruktion selbst vorgenommen, es ist keine Einschätzung aus Analytiker/-innen-Sicht.

Hier ein Beispiel einer Selbstreformulierung nach einer *was heißt X*-Äußerung. In einem ethnographischen Interview zu deutsch-türkischer Migration hatte die Befragte AB den Ausdruck „abgefunden“ (Z. 04) benutzt, um die Haltung ihrer Schwester bezüglich ihrer Entscheidung, nicht mehr in die Türkei zurückzukehren, zu beschreiben.

## (5) FOLK\_E\_00258\_SE\_01\_T\_02\_DF\_01\_c450

- 01 AB: die hat da jetzt nich mehr vor in die türkei  
zuRÜCKzukehrn.
- 02 also die HAT sich da- °hh
- 03 (0.3)
- 04 AB: mit ABgefunden glaub ich [mal;]**
- 05 ZY: [ hmh]m (.) hmhm.
- 06 AB: oder sie HAT,
- 07 (0.3)
- 08 AB: ja,
- > **09 °hwas heißt ABgefunden;**
- 10 also wie gesagt ich hab keine negativen erfahrungen**  
**geMACHT aber-**
- 11 ich !FÜHL! mich in der türKEI viel?**
- 12 (0.82)**
- 13 AB: ((schmatzt))**
- 14 (0.52)**
- 15 AB: Eher(.) zu hause als in DEUTSCHland;**

Im weiteren Verlauf markiert sie diesen Ausdruck aber als unpassend mit „was heißt abgefunden“ (Z. 04) und projiziert so eine Reparatur dieses Ausdrucks. Der Ausdruck wird aber nicht ersetzt, sondern es erfolgt eine mit „also“ eingeleitete, längere Erläuterung ihrer eigenen Erfahrungen und Einstellungen zur Türkei gegenüber Deutschland (Z. 10–15). Die Erläuterung lässt aber schließen, dass sie die negativen Konnotationen von „abgefunden“ abzuschwächen sucht.

Solche Adäquatheitsreparaturen finden wir besonders oft in den Interview-Gesprächen im FOLK-Korpus. Dieser Zusammenhang könnte darin begründet sein, dass sich in dieser Gattung die Gesprächsteilnehmer/-innen offenbar an erhöhten Präzisionskriterien des Ausdrucks orientieren.<sup>9</sup>

Neben diesen Selbstreparaturen/-reformulierungen haben wir auch Fälle, in denen die Verwendung von X durch andere Sprecher als inadäquat markiert und repariert wird. In diesen Fällen werden Zuschreibungen mit einem bestimmten Ausdruck zurückgewiesen. Das kann oft in einer spaßhaften Modalität sein oder – wie im folgenden Ausschnitt aus einem sprachbiographischen Interview – sich auf die vorangegangene Formulierung eines Sachverhalts durch den Ge-

<sup>9</sup> Es bestehen vermutlich weitere Zusammenhänge zwischen der Häufung von Adäquatheitsfällen in Interviewdaten, auf die wir in Deppermann/Reineke (in Vorb.) vertieft eingehen.

sprächspartner beziehen, für den der korrigierende Sprecher selbst die Wissensautorität hat (ein sogenanntes „B-Event“, Labov/Fanshel 1977).

## (6) FOLK\_E\_00179\_SE\_01\_T\_02\_DF\_01\_c631

- 01 NL: **WAS is da die lust an diesem verfremden?**  
 02     bist du da der EINzige? =oder;  
 03 ZI: **ach NEE na ja-**  
 -> 04     **was heeßt verFREMden man,**  
 05     (0.4)  
 06 ZI: **m:: schreibt halt SO wie man manchmal spricht**  
       **u[nd,]**  
 07 NL: [hm ]hm-  
 08 ZI: °h **manchmal spricht man halt so wie man sich fühlt**  
       und, ((Lachansatz))  
 09     °h (.) das is e-  
 10     °h beDINGT halt;  
 11     (0.5)  
 12 ZI: ähm wie man SCHREIBT wie man spricht und so.  
 13     (.) °h und sonst das is jetzt nur so een  
       ABSichtlicher SPASS dass man sagt-  
 14     och ich SCHREIB heut ma so und morgen wieder so.

Die Kategorisierung einer zuvor vom Befragten ZI geschilderten kreativen Sprachverwendungspraxis als „verfremden“ durch den Interviewer NL (Z. 01) wird hier zurückgewiesen. Der interviewte Schüler ZI formuliert im Folgenden die Umstände, aufgrund derer er diese Praxis nicht als „verfremden“ bezeichnen würde (Z. 06–14).

Wie sich auch in den bisherigen Beispielen gezeigt hat, ist X üblicherweise vorerwähnt – das kann im Gespräch selbst sein, im direkten lokalen oder etwas weiter vorangegangenen sequenziellen Kontext oder aber andere Entitäten im lokalen Wahrnehmungsraum wie Vortragsfolien o.ä. betreffen. Eine besonders interessante Verwendungsweise von *was heißt X* zur Anzeige der Inadäquatheit eines Ausdrucks sind demgegenüber solche Fälle, in denen ein Sprecher oder eine Sprecherin sich auf eigene vorangehende Äußerungen bezieht, X aber nicht vorerwähnt war. Im folgenden Beispiel aus einem ethnographischen Interview findet sich eine solche Verwendung (Z. 07/08).

## (7) FOLK\_E\_00148\_SE\_01\_T\_01\_DF\_01\_c749

- 01 HF: °h oder ma MUSS eben-  
 02     °hh die die die ÖFFnungszeiten

03       (.) Ändern oder ANpassen oder sonst wat; =eh,  
 04       (0.2)  
 05 TS: ja.  
 06 HF: °hh zum beispiel is äh kann ich ja Sagen hier;  
 07       ham wir is bei **uns in der**,  
 -> 08       **wat heißt in der diskUSION**,  
 09       aber **is auffällig?**  
 10       °h (.) SAMStach sin die öffnungszeiten immer von  
       zehn bis ei:ns?  
 11 TS: hm\_hm. h°  
 12 HF: °h und DAT is sind (.) sagen wir mal (.)  
       ÖFFnungszeiten,  
 13       °hh die (.) den (.) geWOHNheiten von so\_m;;  
 14       äh äh (.) wochenend (.) äh verHALten?  
 15       (0.7)  
 16 HF: nich entSPRECHen.

Wir haben es hier mit einer spezifischen Art von Reparaturen zu tun, die Stoltenburg (2012) „Präparatur“ nennt: X wurde (noch) gar nicht erwähnt, sondern erscheint nur im Format *was heißt X*. Dadurch erhält dieses eine ‚neue‘ Funktion: X wird als nicht passender Ausdruck eingeführt, dann aber sofort durch einen nachfolgenden ersetzt. Damit wird aber dennoch die potenzielle Relevanz von X nahegelegt, der Sprecher negiert aber zugleich, sich auf das Zutreffen von X zu verpflichten.

Die Beispiele geben einen Einblick in die basalen Verwendungsweisen des von uns anhand der FOLK-Daten untersuchten Formats. Wie es in qualitativen Untersuchungen oft der Fall ist, haben wir in den Fallanalysen Eindrücke gesammelt, die auf systematische Zusammenhänge zwischen dem untersuchten Format und Kontextparametern hindeuten. In diesem Fall waren dies vor allem Gesprächstypabhängigkeiten. So schien/schienen z. B.

- *was heißt X* zur Initiierung einer Reparatur inadäquater Ausdrücke besonders häufig in Interviews verwendet zu werden (siehe hierzu auch Deppermann/Reineke i. Vorb.),
- Definitionen v. a. in Unterrichtsinteraktionen und in Prüfungsgesprächen abgefragt zu werden und
- das Format bevorzugt in bestimmten Sprecher/-innen-Rollenkonstellationen verwendet zu werden.

Um solche Eindrücke und Hypothesen systematisch prüfen zu können, haben wir unser gespeichertes Suchresultat in der DGD um weitere relevante Kategorien im

Reiter „Metadaten“ in unserer KWIC ergänzt. In Abbildung 13 sehen wir beispielhaft die Metadaten für Art des Gesprächs und Rolle der Sprecher und Sprecherin für jeden Beleg aus unserem Suchresultat.

The screenshot shows the 'Metadaten' tab in the DGD interface. At the top, there are navigation tabs: 'POSITION', 'TOKEN', 'KONTEXT', 'METADATEN', and 'ANZEIGE'. Below these, there are search filters: 'Deskriptor: SE: Kurzbezeichnung (\"Art\")' and 'Deskriptor: SES: Rolle'. A search bar contains the text 'nom\_heißt\_was\_2\_k\_gef'. Below the search bar, there is a table with the following columns: 'Sprecher\*in', 'Sprecher', 'Träger', 'Art', and 'Rolle'. The table contains 20 rows of data, each representing a different instance of the word 'heißt' in various contexts and roles.

Sprecher*in	Sprecher	Träger	Art	Rolle
1	FOLK_00004_01 GS	was heißt das dann für 1 seinen zusätzlichen erlog	Unterrichtsba...	Schülerin
2	FOLK_00004_01 AB	was heißt n kognitiv	Unterrichtsba...	Schülerin
3	FOLK_00004_01 GS	genau danke was heißt kognitiv	Unterrichtsba...	Schülerin
4	FOLK_00007_01 GS	was heißt eigentlich kognitiv	Unterrichtsba...	Lehrerin
5	FOLK_00007_01 GS	was vor was heißt was vorgegeben so es vorgegeben	Unterrichtsba...	Lehrerin
6	FOLK_00013_01 CJ	was heißt das	Vorlesen für K...	Familienmitglied ; Vort...
7	FOLK_00014_01 CJ	heißt un willst du wissen was ängstlich heißt auf türkisch	Vorlesen für K...	Familienmitglied ; Vort...
8	FOLK_00014_01 CJ	doch ne vom papa und was heißt zu hause auf türkisch weißt du des	Vorlesen für K...	Familienmitglied ; Vort...
9	FOLK_00014_01 CJ	des kennst du genau und was heißt spielen auf türkisch weißt du des	Vorlesen für K...	Familienmitglied ; Vort...
10	FOLK_00019_01 CH	zuerst mal ah die frage stellen was heißt	Prüfungsgegp...	Prüferin
11	FOLK_00019_01 CH	was heißt sprachmelodie was heißt auch spr zeit verlauf was dass wir	Prüfungsgegp...	Prüferin
12	FOLK_00019_01 CH	was heißt sprachmelodie was heißt auch spr zeit verlauf was dass wir da genauer sehen	Prüfungsgegp...	Prüferin
13	FOLK_00019_01 CH	methode für eh ersprachenerw na wa heißt das metho methode für die datengewinnung	Prüfungsgegp...	Prüferin
14	FOLK_00020_01 EM	was heißt dies	Tischgespräch	Familienmitglied
15	FOLK_00021_01 JZ	was heißt mir	Spätereraktio...	Mitspielerin
16	FOLK_00021_01 XM1	was heißt hier schon	Spätereraktio...	---
17	FOLK_00021_01 PL	hao hatt ho ho was heißt hier du kauft	Spätereraktio...	Mitspielerin
18	FOLK_00024_01 SZ	ja was heißt jetzt nett	Meeting in ein...	Mitarbeiterin ; Meetin...
19	FOLK_00026_01 HM	hatt in die a und e was heißt n des e is erziehungshilfe oder was	Meeting in ein...	Chefin ; Gruppenleit...
20	FOLK_00026_01 MS	was heißt n des	Meeting in ein...	Mitarbeiterin ; Meetin...

Abb. 13: Um Metadatenangaben erweiterte KWIC-Ansicht in der DGD

Diese Daten haben wir dann erneut zur Bearbeitung in Excel ausgegeben und in unsere Tabelle mit den kodierten Belegen ergänzt. So können wir auf Belegebene auswerten, ob und wie bestimmte Verwendungen der untersuchten Form mit Metadatenparametern systematisch zusammenhängen; das können wir hier nicht im Einzelnen zeigen, aber mit gängigen Anwendungen, z. B. über Pivot-Tabellen-Funktionen in Excel können wir uns mit ein paar Klicks Verwendungsverteilungen ausgeben lassen und auswerten. Für die Auswertung ist es auch hier wesentlich, sich mit der Systematik der Metadatenwerte im FOLK-Korpus auseinanderzusetzen, bevor man sie zur Interpretation von Verteilungen nutzt (vgl. auch Deppermann/Reineke i. Vorb.). Gegebenenfalls müssen auch einzelne Parameter zu für die jeweils eigene Analyse relevanten Kategorien neu gruppiert oder zusammengefasst werden.

Wir haben in diesem Aufsatz versucht, an einem Untersuchungsbeispiel zu zeigen, welche großen Potenziale das FOLK-Korpus und die in ihm enthaltenen annotierten Daten für interaktionslinguistische Fragestellungen haben. Die sachgerechte Arbeit mit dem Korpus erfordert es aber, einige Dinge adäquat in Rech-

nung zu stellen, was unserer Erfahrung nach leider nicht in jeder Untersuchung getan wird. Dazu gehören vor allem folgende Punkte:

- Der Nutzung des Korpus für eine Untersuchung sollte immer eine intensive Auseinandersetzung mit der Architektur des Korpus, seinen Annotationsebenen sowie den Funktionsweisen der DGD vorausgehen. So kann verhindert werden, dass mangelnde Kenntnisse von Transkriptionsmodell, Datenbestand und Funktionsweise der Suchen zu Fehlschlüssen führen. So sind bspw. Transkriptsegmente nicht mit Äußerungen, Turns oder Turnkonstruktionseinheiten gleichzusetzen; der Aufnahmeort ist nicht mit einer für die Sprachproduktion relevanten Sprachregion gleichzusetzen; phonetische Variation ist in den transkribierten Wortformen nicht immer konsistent repräsentiert etc.
- Automatische Suchergebnisse ersetzen nicht die Einzelfallanalyse der zugrundeliegenden Daten und die Prüfung der Korrektheit ihrer Annotation und Transkription. Konversationsanalytische Kriterien der detaillierten Sequenzanalyse und der Kollektionsanalyse dürfen nicht außer Acht gelassen werden.
- Eine etwaige Relevanz von Metadatenannotationen für die eigene Untersuchung ist in der Datenanalyse zu erweisen und kann nicht *a priori* angenommen werden. Die Metadaten in FOLK werden auf Sprecher- und Gesprächsebene erhoben und im Projekt systematisiert; einige Werte werden gemäß den Stratifikationsparametern vom FOLK-Team zugewiesen. Es kann daher sein, dass sich bspw. das Verständnis des Gesprächstyps, das die Teilnehmenden selbst von der gegenwärtigen Aktivität haben, oder die lokal für die Interaktion relevanten Identitäten von dem im Korpus global für das gesamte Gespräch annotierten Wert unterscheiden. Außerdem finden wir häufig Gesprächsphasen, die nicht dem übergeordneten Gesprächstyp entsprechen, z. B. Smalltalk in der Fahrstunde, das Gespräch über gemeinsame Freunde im WG-Casting oder die Planung der Mittagessenstermine in den Schlichtungsgesprächen zu Stuttgart21.

## 4 Fazit

In diesem Aufsatz haben wir an einer Untersuchung exemplarisch die Nutzung des Korpus FOLK für interaktionslinguistische Fragestellungen vorgehensbezogen gezeigt. Abschließend wollen wir die wesentlichen Potenziale resümieren, die sich aus unserer Sicht mit der Nutzung von FOLK verbinden und die dazu führen, dass die Nutzung für FOLK für viele interaktionslinguistische Untersu-

chungen aus unserer Sicht die Datengrundlage der Wahl sein kann, da das Korpus erhebliche Vorteile gegenüber der Arbeit mit neu zu erhebenden oder anderen, bestehenden Korpora aufweist.

Der größte Vorzug von FOLK, welcher von Beginn an den Aufbau des Korpus motiviert hat, ist die Tatsache, dass es den wissenschaftsöffentlichen Zugang zu einer großen Anzahl an Daten aus einer großen Spannweite von Interaktionstypen und Sprecher/-innengruppen bietet, die in dieser Form weder in anderen Korpora bereitstehen noch je von einzelnen Projektgruppen selbst erhoben werden könnten. Dies erlaubt vor allem für relativ häufige Phänomene eine belastbarere, differenziertere und generalisierbarere Untersuchung als sie mit der Nutzung der üblicherweise viel kleineren und auf einen oder wenige soziale Kontexte beschränkten Projektkorpora möglich ist. Das Angebot von FOLK ermöglicht eine enorme Effizienzsteigerung von Untersuchungen: Die eigene Erhebung und die basale Transkription der Daten entfallen für viele Arten von Untersuchungen. Das heißt auch, dass Erhebungen, die einmal für FOLK gemacht worden sind, wissenschaftsgeschichtlich kumulativ werden und nachgenutzt werden können. Die wissenschaftsöffentliche Verfügbarkeit der Daten und die Sicherstellung zeitüberdauernder, zitierbarer Referenzierbarkeit und Zugänglichkeit der Daten ermöglicht es, mit FOLK vorgenommene Untersuchungen zu prüfen und zu replizieren, wie dies mit anderen Untersuchungen in der Gesprächsforschung nicht möglich ist.

FOLK liefert dabei insbesondere nach interaktionslinguistischen und konversationsanalytischen Maßstäben hochwertige Daten. Die Daten sind natürlich, sie sind an *best-practice* Standards orientiert, die die Datenqualität in vielerlei Hinsicht absichern und auf langjähriger Erprobung und Optimierung beruhen. Dies gilt für Aufnahmetechnik, Transkription, Dokumentation und auch den Datenschutz (informierte Einwilligung), der in einem wissenschaftsöffentlichen Korpus naturgemäß besonders streng gehandhabt wird.

Das Korpus eignet sich besonders gut für formbasierte Untersuchungen, denn hier entfalten die automatischen Suchmöglichkeiten ihr Potenzial. Die Fülle der Daten in FOLK, die Erschließungsmöglichkeiten und die Möglichkeiten der Treffervisualisierung in der DGD bieten Möglichkeiten der induktiven Entdeckung von Phänomenen und Mustern von Phänomenen durch die Inspektion von Suchergebnissen. In einer langfristigen Perspektive kann man von FOLK nicht mehr nur von einem synchronen Korpus sprechen, sondern es wird mit der Zeit zu einem diachronen Korpus.

Mit der stetig wachsenden Größe des Korpus bieten sich zunehmend mehr Möglichkeiten zur Erstellung virtueller Korpora nach ausgewählten Metadatenmerkmalen. Dies ist heute schon individuell nach spezifischen Forschungsinteressen möglich, in Zukunft werden virtuelle Korpora auch in FOLK selbst nach

spezifischen Stratifikationsparametern (z. B. virtuelle Korpora vergleichbarer Erhebungszeiträume oder Gesprächstypen) kuratiert werden.

## Literatur

- Aldrup, Marit/Küttner, Uwe-A./Lechler, Constanze/Reinhardt, Susanne (2021): Suspended assessments in German talk-in-interaction. In: Kupetz, Maxi/Kern, Friederike (Hg.): Prosodie und Multimodalität. (= OraLingua 18). Heidelberg: Winter, S. 31–66.
- Bergmann, Pia (2017): Gebrauchsprofile von *weiß nich* und *keine Ahnung* im Gespräch. Ein Blick auf nicht-responsive Vorkommen. In: Blühdorn, Hardarik/Deppermann, Arnulf/Helmer, Henrike/Spranz-Fogasy, Thomas (Hg.): Diskursmarker im Deutschen. Reflexionen und Analysen. Göttingen: Verlag für Gesprächsforschung, S. 157–182.
- Betz, Emma (2015): Indexing epistemic access through different confirmation formats. Uses of responsive (*das stimmt*) in German interaction. In: Journal of Pragmatics 87, S. 251–266.
- Bies, Andrea (2020): WG-Castings im DaF-Unterricht. In: Deutsch als Fremdsprache 57, 2, S. 88–101.
- De Stefani, Elwys (im Dr.): A prima facie resource for problematizing meaning. Taking a stance with (Che) cosa vuol dire X? („What does X mean“). In: Interactional Linguistics.
- Deppermann, Arnulf (2014): Multimodal participation in simultaneous joint projects. Interpersonal and intrapersonal coordination in paramedic emergency drill. In: Haddington, Pentti/Keisanen, Tiina/Mondada, Lorenza/Nevile, Maurice (Hg.): Multiactivity in social interaction. Beyond multitasking. Amsterdam/Philadelphia: Benjamins, S. 247–282.
- Deppermann, Arnulf (2018): Changes in turn-design over interactional histories. The case of instructions in driving school lessons. In: Deppermann, Arnulf/Streeck, Jürgen (Hg.): Time in embodied interaction. Synchronicity and sequentiality of multimodal resources. (= Pragmatics & Beyond New Series 293). Amsterdam/Philadelphia: Benjamins, S. 293–324.
- Deppermann, Arnulf (2020): Interaktionale Semantik. In: Hagemann, Jörg/Staffeldt, Sven (Hg.): Semantiktheorien II. Analysen von Wort- und Satzbedeutungen im Vergleich. (= Stauffenburg Einführungen 36). Tübingen: Stauffenburg, S. 235–278.
- Deppermann, Arnulf (im Dr. a): An exercise in interactional semantics. Definitions and specifications provided in response to was heißt x? („What does X mean?“). In: Interactional Linguistics.
- Deppermann, Arnulf (im Dr. b): „What do you understand by X?“. Semantics in interactional linguistics. In: Selting, Margret/Barth-Weingarten, Dagmar (Hg.): New perspectives in interactional linguistic research. Amsterdam/Philadelphia: Benjamins.
- Deppermann, Arnulf/Gubina, Alexandra (2021): When the body belies the words. Embodied agency with *darf/kann ich?* („May/Can I?“) in German. In: Frontiers in Communication 6, S. 1–16.
- Deppermann, Arnulf/Hartung, Martin (2011): Was gehört in ein nationales Gesprächskorpus? Kriterien, Probleme und Prioritäten der Stratifikation des „Forschungs- und Lehrkorpus Gesprochenes Deutsch“ (FOLK) am Institut für Deutsche Sprache (Mannheim). In: Felder,

- Ekkehard/Müller, Markus/Vogel, Friedemann (Hg.): Korpuspragmatik. Thematische Korpora als Basis diskurslinguistischer Analysen. (= Linguistik – Impulse und Tendenzen 44). Berlin/Boston: De Gruyter, S. 414–450.
- Deppermann, Arnulf/Reineke, Silke (2020): Practices of indexing discrepant assumptions with German *ich dachte* (‘I thought’) in talk-in-interaction. In: *Functions of Language* 27, 2, S. 113–142.
- Deppermann, Arnulf/Reineke, Silke (in Vorb.): Zur Verwendung von Metadaten in der konversationsanalytischen Arbeit mit Korpora – am Beispiel einer Untersuchung anhand des Korpus FOLK.
- Deppermann, Arnulf/Schmidt, Axel (2021): Micro-sequential coordination in early responses. In: *Discourse Processes* 58, 4, S. 372–396.
- Deppermann, Arnulf/Proske, Nadine/Zeschel, Arne (Hg.) (2017): Verben im interaktiven Kontext. Bewegungsverben und mentale Verben im gesprochenen Deutsch. (= Studien zur Deutschen Sprache 74). Tübingen: Narr.
- Droste, Pepe/Günthner, Susanne (2020): „das machst du bestimmt AUCH du;“. Zum Zusammenspiel syntaktischer, prosodischer und sequenzieller Aspekte syntaktisch desintegrierter *du*-Formate. In: Imo/Lanwer (Hg.), S. 75–109.
- Fandrych, Christian/Meißner, Cordula/Slavcheva, Adriana (2012): The GeWiss corpus. Comparing spoken academic German, English and Polish. In: Schmidt, Thomas/Wörner, Kai (Hg.): *Multilingual corpora and multilingual corpus analysis*. (= Hamburg Studies on Multilingualism 14). Amsterdam/Philadelphia: Benjamins, S. 319–338.
- Frick, Elena/Wallner, Franziska/Helmer, Henrike (in Vorb.): ZuRecht: Neue Recherchemöglichkeiten in Korpora gesprochener Sprache für Gesprächsanalyse und Deutsch als Fremd- und Zweitsprache. In: KorDaF, Themenheft „Zugänge zu multimodalen Korpora gesprochener Sprache“.
- Gubina, Alexandra (2021): Availability, grammar, and action formation. On simple and modal interrogative request formats in spoken German. In: *Gesprächsforschung – Online-Zeitschrift zur verbalen Interaktion* 2 (Themenheft: ‚How to get things done‘ – Anforderungen und Instruktionen in der multimodalen Interaktion), S. 272–303. <http://www.gespraechsforschung-online.de> (Stand: 29.8.2022).
- Günthner, Susanne (2015): Grammatische Konstruktionen im Kontext sequenzieller Praktiken – „was heißt x“-Konstruktionen im gesprochenen Deutsch. In: Bücker, Jörg/Günthner, Susanne/Imo, Wolfgang (Hg.): *Konstruktionen im Spannungsfeld von sequenziellen Mustern, kommunikativen Gattungen und Textsorten*. (= Konstruktionsgrammatik 5). Tübingen: Stauffenburg, S. 187–219.
- Günthner, Susanne/König, Katharina (2015): Temporalität und Dialogizität als interaktive Faktoren der Nachfeldpositionierung – ‚irgendwie‘ im gesprochenen Deutsch. In: Vinckel-Roisin, Héléne (Hg.): *Das Nachfeld im Deutschen. Theorie und Empirie*. (= Reihe Germanistische Linguistik 303). Berlin/Boston: De Gruyter, S. 255–278.
- Helmer, Henrike (2020): How do speakers define the meaning of expressions? The case of German *x heißt y* (‘x means y’). In: *Discourse Processes* 57, 3, S. 278–299.
- Helmer, Henrike/Deppermann, Arnulf (2022): Verständlichkeit und Partizipation in den Schlichtungsgesprächen zu Stuttgart 21. In: Kämper, Heidrun/Plewnia, Albrecht (Hg.): *Sprache in Politik und Gesellschaft. Perspektiven und Zugänge*. (= Jahrbuch des Instituts für Deutsche Sprache 2021). Berlin/Boston: De Gruyter, S. 263–294.

- Helmer, Henrike/Reineke, Silke/Deppermann, Arnulf (2016): A range of uses of negative epistemic constructions in German. *ich weiß nicht* as a resource for dispreferred actions. In: *Journal of Pragmatics* 106, S. 97–114.
- Imo, Wolfgang/Lanwer, Jens P. (2020): *Prosodie und Konstruktionsgrammatik*. Berlin/Boston: De Gruyter.
- ISO (2016): ISO 24624:2016 Language resource management – Transcription of spoken language. [www.iso.org/standard/37338.html](http://www.iso.org/standard/37338.html) (Stand: 29.8.2022).
- Kaiser, Julia (2018): Zur Stratifikation des FOLK-Korpus. Konzeption und Strategien. *Gesprächsforschung – Online-Zeitschrift zur verbalen Interaktion* 19, S. 515–552. <http://www.gespraechsforschung-online.de> (Stand: 29.8.2022).
- Katelhön, Peggy (2016): Verbale Progressivkonstruktionen im Deutschen und Italienischen. Ein korpusbasierter Sprachvergeich. In: Selig, Maria/Morlicchio, Elda/Dittmar, Norbert (Hg.): *Gesprächsanalyse zwischen Syntax und Pragmatik. Deutsche und italienische Konstruktionen*. (= *Stauffenburg Linguistik* 78). Tübingen: Stauffenburg, S. 169–188.
- Katelhön, Peggy/Moroni, Manuela C. (2018): Inszenierungen direkter Rede in mündlichen Interaktionen. In: *Quaderni dell'ALG* 1, S. 179–208.
- König, Katharina (2020): Prosodie und *epistemic stance*. Konstruktionen mit finalem *oder*. In: Imo/Lanwer (Hg.), S. 167–199.
- Kress, Karoline (2017): Das Verb *machen* im gesprochenen Deutsch. Bedeutungskonstitution und interaktionale Funktionen. (= *Studien zur Deutschen Sprache* 78). Tübingen: Narr.
- Labov, William/Fanshel, David (1977): *Therapeutic discourse. Psychotherapy as conversation*. New York: Academic Press.
- Luckmann, Thomas (1986): Grundformen der gesellschaftlichen Vermittlung des Wissens. Kommunikative Gattungen. In: *Kölner Zeitschrift für Soziologie und Sozialpsychologie* (Sonderheft: Kultur und Gesellschaft 27), S. 191–211.
- Moroni, Manuela C. (2015): Intonation und Bedeutung. Die nuklear steigend-fallende Intonationskontur in einer deutschen und einer italienischen Varietät. In: *Deutsche Sprache* 43, S. 255–286.
- Moroni, Manuela C. (2016): Ironie und Intonation im privaten Gespräch. In: Amann, Klaus/Hackl, Wolfgang (Hg.): *Satire – Ironie – Parodie. Aspekte des Komischen in der deutschen Sprache und Literatur*. (= *Innsbrucker Beiträge zur Kulturwissenschaft. Germanistische Reihe* 85). Innsbruck: Innsbruck University Press, S. 167–185.
- Proske, Nadine (2014): ‚Oh ach KOMM; hör AUF mit dem kIEInkram‘. Die Partikel *komm* zwischen Interjektion und Diskursmarker. In: *Gesprächsforschung – Online-Zeitschrift zur verbalen Interaktion* 15, S. 121–160. <http://www.gespraechsforschung-online.de> (Stand: 29.8.2022).
- Reineke, Silke (2016): Wissenszuschreibungen in der Interaktion. Eine gesprächsanalytische Untersuchung impliziter und expliziter Formen der Zuschreibung von Wissen. (= *OraLingua* 12). Heidelberg: Winter.
- Schmid, Helmut (1995): Improvements in part-of-speech tagging with an application to German. In: Tzoukermann, Evelyne (Hg.): *Proceedings of the ACL SIGDAT-Workshop, Dublin*. <https://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/data/tree-tagger2.pdf>.
- Schmidt, Thomas (2017): Construction and dissemination of a corpus of spoken interaction. Tools and workflows in the FOLK project. In: *Journal for Language Technology and Computational Linguistics (JLCL)* 31, 1, S. 127–154.

- Schmidt, Thomas/Schütte, Wilfried/Winterscheid, Jenny (2015): cGat. Konventionen für das computergestützte Transkribieren in Anlehnung an das Gesprächsanalytische Transkriptionssystem 2 (GAT2). Mannheim: Institut für Deutsche Sprache.
- Selting, Margret/Auer, Peter/Barth-Weingarten, Dagmar/Bergmann, Jörg/Bergmann, Pia/Birkner, Karin/Couper-Kuhlen, Elizabeth/Deppermann, Arnulf/Gilles, Peter/Günthner, Susanne/Hartung, Martin/Kern, Friederike/Mertzluff, Christine/Meyer, Christian/Morek, Miriam/Oberzaucher, Frank/Peters, Jörg/Quasthoff, Uta/Schütte, Wilfried/Stukenbrock, Anja/Uhmann, Susanne (2009): Gesprächsanalytisches Transkriptionssystem 2 (GAT 2). In: Gesprächsforschung – Online-Zeitschrift zur verbalen Interaktion 10, S. 53–402. <http://www.gespraechsforschung-online.de> (Stand: 29.8.2022).
- Staffeldt, Sven (2018): Über *sehr sehr*. Beobachtungen zum Vorkommen einer totalen Reduplikation im gesprochenen Deutsch. In: Filatkina, Natalia/Stumpf, Sören (Hg.): Konventionalisierung und Variation. Phraseologische und konstruktionsgrammatische Perspektiven. (= Sprache – System und Tätigkeit 71). Berlin: Lang, S. 179–200.
- Stoltenburg, Benjamin (2012): „ich will jetzt nicht sagen Reparaturen, aber...“ – Eine Gesprächsstrategie zur Indizierung von Problemstellen. In: gidi Arbeitspapier 47, 10, <https://arbeitspapiere.sprache-interaktion.de/arbeitspapiere/arbeitspapier47.pdf> (Stand: 25.8.2022).
- Stratton, James M. (2020): Adjective intensifiers in German. In: Journal of Germanic Linguistics 32, 2, S. 183–215.
- Torres Cajo, Sarah (2019): Zwischen Strukturierung, Wissensmanagement und Argumentation im Gespräch. Interaktionale Verwendungsweisen der Modalpartikeln *halt* und *eben* im gesprochenen Deutsch. In: Deutsche Sprache 47, S. 289–310.
- Westpfahl, Swantje (2020): POS-Tagging für Transkripte gesprochener Sprache. Entwicklung einer automatisierten Wortarten-Annotation am Beispiel des Forschungs- und Lehrkorpus Gesprochenes Deutsch (FOLK). (= Studien zur Deutschen Sprache 83). Tübingen: Narr.
- Westpfahl, Swantje/Schmidt, Thomas/Jonietz, Jasmin/Borlinghaus, Anton (2017): STTS 2.0. Guidelines für die Annotation von POS-Tags für Transkripte gesprochener Sprache in Anlehnung an das Stuttgart Tübingen Tagset (STTS). Version 1.1. Mannheim: Leibniz-Institut für Deutsche Sprache. URN: urn:nbn:de:bsz:mh39-60634.
- Wieling, Martijn/Grieve, Jack/Bouma, Gosse/Fruehwald, Josef/Coleman, John/Liberman, Mark (2016): Variation and change in the use of hesitation markers in Germanic languages. In: Language Dynamics and Change 6, 2, S. 199–234.
- Zinken, Jörg/Küttner, Uwe-A. (2022): Offering an interpretation of prior talk in everyday interaction. A semantic map approach. In: Discourse Processes 59, 4, S. 298–325.