



Institut für Deutsche Sprache

Sylvia Dickgießer

unter Mitarbeit von Caren Brinckmann und Joachim Gasch

Metadaten-Schemata in der Datenbank für Gesprochenes Deutsch (DGD 2.0)

Stand: 2009-09-01

© Institut für Deutsche Sprache, Mannheim

Inhaltsverzeichnis

1.	Einleitung	4
2.	Externe Metadatenschemata	4
3.	IDS-Datenmodell für die Dokumentation von Korpora der gesprochenen Sprache	5
4.	Generische Schemata und projektspezifische Subschemata	6
5.	Generisches Schema für die Dokumentation von Korpusbestandteilen auf der Ereignisebene	6
5.1.	Ereignis	7
5.1.1.	Basisdaten	8
5.1.2.	Rundfunksendung	9
5.1.3.	Projekt	10
5.1.4.	Quellaufnahme	10
5.1.4.1.	Basisdaten	11
5.1.4.2.	Aufnahmetechnik	12
5.1.4.3.	Technische Fassungen	12
5.1.4.4.	Archivierung und Distribution	15
5.1.5.	Zusatzmaterial	16
5.1.5.1.	Basisdaten	17
5.1.5.2.	Technische Fassungen	17
5.1.5.3.	Archivierung und Distribution	19
5.1.6.	Sprechereignis	19
5.1.6.1.	Basisdaten	20
5.1.6.2.	Beschreibung	20
5.1.6.3.	Sprecher	21
5.1.6.3.1.	Basisdaten	22
5.1.6.3.2.	Sprachdaten	22
5.1.6.4.	Sprechereignisspezifische Aufnahme	23
5.1.6.4.1.	Basisdaten	23
5.1.6.4.2.	Technische Fassungen	24
5.1.6.4.3.	Archivierung und Distribution	26
5.1.6.5.	Transkript	26
5.1.6.5.1.	Basisdaten	27
5.1.6.5.2.	Annotation	27
5.1.6.5.2.1.	Basisdaten	28
5.1.6.5.2.2.	Erstellung	28
5.1.6.5.2.3.	Alignment	29
5.1.6.5.3.	Technische Fassungen	30
5.1.6.5.4.	Archivierung und Distribution	31
5.1.6.6.	Zusatzmaterial	32
5.1.6.6.1.	Basisdaten	32
5.1.6.6.2.	Technische Fassungen	33
5.1.6.6.3.	Archivierung und Distribution	34
5.1.7.	Dokumentationsgeschichte	34

6.	Generisches Schema für die Dokumentation allgemeiner Sprecherdaten	35
6.1.	Sprecher	36
6.1.1.	Basisdaten	36
6.1.2.	Ortsdaten	37
6.1.3.	Sprachdaten	38
6.1.3.1.	Sprachkenntnisse	38
6.1.3.2.	Sprachproduktion	39
6.1.3.3.	Sprachgebrauch	40
6.1.4.	Beziehungen zu anderen Sprechern	41
6.1.5.	Sonstige Bezugspersonen	41
6.1.5.1.	Bezugspersonen kompakt	41
6.1.5.2.	Einzelne Bezugsperson	42
6.1.5.2.1.	Basisdaten	43
6.1.5.2.2.	Ortsdaten	43
6.1.5.2.3.	Sprachdaten	44
6.1.5.2.3.1.	Sprachkenntnisse	44
6.1.5.2.3.2.	Sprachgebrauch	44
6.1.6.	Rechteverwaltung	45
6.1.7.	Zusatzmaterial	46
6.1.7.1.	Basisdaten	46
6.1.7.2.	Technische Fassungen	46
6.1.7.3.	Archivierung und Distribution	48
6.1.8.	Dokumentationsgeschichte	49
7.	Generisches Schema für die Dokumentation von Zusatzmaterial auf Korpusebene	49
7.1.	Basisdaten	50
7.2.	Technische Fassungen	50
7.3.	Archivierung und Distribution	52
7.4.	Dokumentationsgeschichte	53
8.	Generisches Schema für die Korpusbeschreibung	53
8.1.	Erstellungsprojekt	54
8.2.	Korpusarbeiten	55
8.3.	Aufzeichnungsobjekte	55
8.4.	Korpusbestandteile	57
8.4.1.	Quellaufnahmen	57
8.4.1.1.	Basisdaten	57
8.4.1.2.	Aufnahmetechnik	58
8.4.1.3.	Technische Fassungen	58
8.4.1.4.	Archivierung und Distribution	60
8.4.2.	Sprechereignisspezifische Aufnahmen	61
8.4.2.1.	Basisdaten	61
8.4.2.2.	Transkribierte SE-Aufnahmen	62
8.4.2.3.	Technische Fassungen	62
8.4.2.4.	Archivierung und Distribution	63
8.4.3.	Transkripte	64
8.4.3.1.	Basisdaten	64
8.4.3.2.	Annotation	64
8.4.3.2.1.	Basisdaten	65
8.4.3.2.2.	Erstellung	65
8.4.3.2.3.	Alignment	66
8.4.3.3.	Technische Fassungen	66
8.4.3.4.	Archivierung und Distribution	67
8.4.4.	Zusatzmaterial	68

8.4.4.1.	Basisdaten	68
8.4.4.2.	Technische Fassungen	69
8.4.4.3.	Archivierung und Distribution	69
8.5.	Dokumentationsgeschichte	70
9.	Abschließende Bemerkungen	71
10.	Anmerkungen	72

1. Einleitung

Die Datenbank für Gesprochenes Deutsch (DGD 2.0) enthält eine Metadatenkomponente, die auf einem neuen Modell für die Dokumentation von Korpora der gesprochenen Sprache beruht und vier darauf aufbauende (XML-)Schemata umfasst.

Die Entwicklung dieser Komponente orientierte sich an folgenden Richtlinien:

- Unabhängigkeit von spezifischen Forschungsansätzen
- Vermittlung zwischen projektübergreifenden und projektspezifischen Anforderungen
- detaillierte Datenstruktur
- kalkulierte Redundanz
- validierbare Datenerfassung
- konsistente zentrale Datenspeicherung
- variable benutzerfreundliche Darstellung
- effektives korpusübergreifendes Retrieval
- datenschutz- und datensicherheitsgerechte Benutzerverwaltung

Die Unabhängigkeit von spezifischen Forschungsansätzen soll es uns ermöglichen, Daten aus verschiedenen Bereichen in übergeordnete Strukturen integrieren zu können.

2. Externe Metadatenschemata

Die Anzahl empfohlener Metadatenschemata ist beachtlich. Die bekanntesten sind Dublin Core (DC) [1], die Systematik des Open Language Archives (OLAC) [2], die der Text Encoding Initiative (TEI) [3], MPEG 7 [4] und schließlich die Schemata der ISLE Meta Data Initiative (IMDI) [5].

Bei der Auswahl und Bewertung möglicher Vorgaben für unsere Arbeit stützten wir uns v.a. auf eine Studie von Thorsten Trippel und Tanja Baumann, die auf der Suche nach einem geeigneten Metadatenstandard für „multimodale Korpora“ („linguistische Korpora“) die Schemata von DC, OLAC, TEI und IMDI miteinander verglichen und auf ihre Nützlichkeit für die Dokumentation solcher Korpora hin überprüft haben. Die Autoren kamen zu dem Ergebnis:

„Für die Archivierung von Ressourcen sind verschiedene Standards definiert und betrachtet worden: Dublin Core: kleinster gemeinsamer Nenner von Metadaten [...], wobei der Schwerpunkt auf der Katalogisierung von Ressourcen liegt. OLAC: DC Erweiterung für mehrsprachige und vor allem auch in anderen Medien vorliegenden Ressourcen. TEI: Struktur für Metadaten für gedruckte und textuelle Medien [...], wobei weder Mehrsprachigkeit noch andere Medien vorgesehen sind. IMDI: geeignetster Standard, da er die anderen Standards konzeptuell einschließt und gleichzeitig Probleme mehrsprachiger Ressourcen und verschiedener Medien berücksichtigt. Einzig die fehlenden Datenkategorien auf Annotationsebene stellen ein Problem dar, wobei aber auch in anderen Standards hierfür keine Kategorien bekannt sind.“ [6]

IMDI hat bislang zwei Schemata für die Dokumentation linguistischer Ressourcen bereitgestellt: Ein Schema für eine Korpusbeschreibung (catalogue descriptions) und eines für die Dokumentation von Korpusbestandteilen (session descriptions). Diese Schemata wurden für das Retrieval (mit Hilfe spezieller Browser) und für die Publikation von Metadaten konzipiert. Sie werden u.a. im Max-Planck-Institut für Psycholinguistik in Nijmegen verwendet.

„Session“ wird folgendermaßen erläutert: „The session concept bundles all information about the circumstances and conditions of the linguistic event, groups the resources belonging to this linguistic event, records the administrative information of the event and describes the content of the event. Since version 3.0. also written resources other than annotations can be included in a session. For written resources the definition of session is extended to include all documents that pertain to the creation, analysis and commentary of a document.“ [7]

Wir haben das IMDI-Session-Schema geprüft und sind zunächst auf zwei für uns problematische Punkte aufmerksam geworden:

a) Um eine Vergleichbarkeit der Daten gewährleisten und zugleich sehr heterogenen Datenbeständen gerecht werden zu können, enthält das Schema nur eine relativ kleine Anzahl verbindlicher Informationselemente. Daneben sind in allen Abschnitten optionale „description elements“, in denen unstrukturierte Texte abgelegt werden können, sowie optionale „keys“ vorgesehen, die es verschiedenen Forschungsgruppen ermöglichen, gruppenspezifische Strukturen in das Schema zu integrieren. Dadurch wird eine große Flexibilität gewährleistet, aber auch eine Beliebigkeit gefördert, die für unsere Zwecke nicht sinnvoll ist.

b) Das Session-Konzept bezieht sich auf „linguistic events“ und speichert alle Daten über die sozialen Kontexte („circumstances and conditions“) dieser „linguistic events“ sowie alle Daten für die an einer „Session“ Beteiligten („actors“), in einem Schema. Das führt zu Redundanzen in der Datenbasis, wenn mehrere „linguistic events“ in einem sozialen Kontext zu beschreiben sind und wenn einzelne Personen an mehreren dieser „linguistic events“ beteiligt waren.

In Anbetracht dieser Besonderheiten entschieden wir uns für die Entwicklung eines eigenen Modells unter Berücksichtigung vorhandener Schemata (wie IMDI) und Empfehlungen für die Dokumentation von Korpora der gesprochenen Sprache, wie z.B. die des Bayerischen Archivs für Sprachsignale (BAS) [8].

3. IDS-Datenmodell für die Dokumentation von Korpora der gesprochenen Sprache

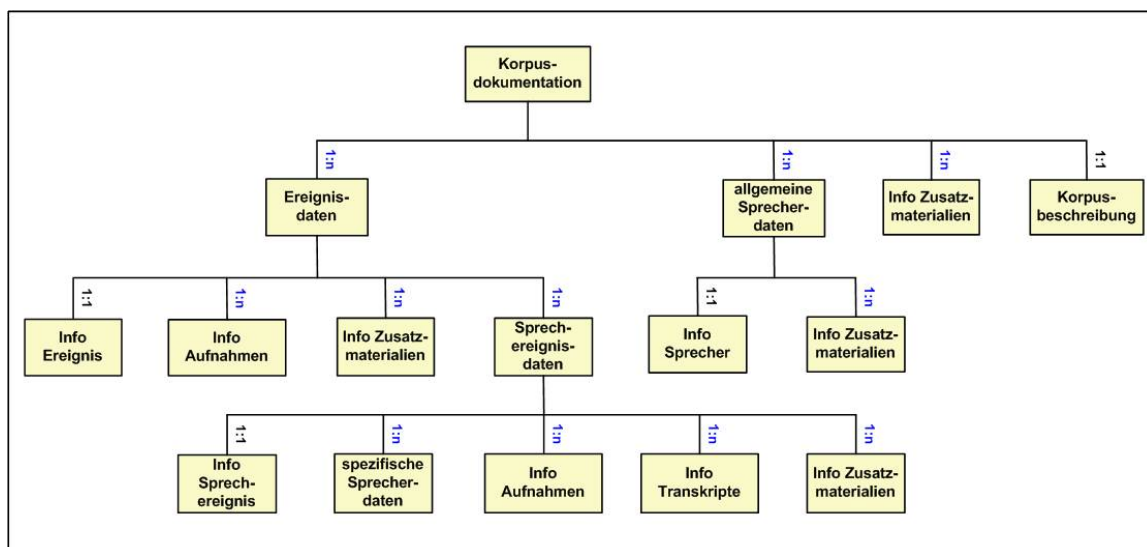


Abb. 1, Datenmodell für die Korpusdokumentation

Abb. 1 zeigt unser Datenmodell für die Dokumentation von Korpora der gesprochenen Sprache, das vier Bereiche vorsieht, die mithilfe von (XML-)Schemata strukturiert werden: einen Bereich für Ereignisdaten, einen Bereich für ereignisübergreifende, allgemeine Sprecherdaten, einen Block für Informationen über Zusatzmaterialien auf Korpusebene (z.B. Transkriptionskonventionen oder Texte, die von allen Informanten vorgelesen wurden) und eine Korpusbeschreibung.

Dieses Modell hebt sich von dem Modell, das der IMDI-Session-Beschreibung zugrundeliegt, v.a. dadurch ab, dass unterschieden wird zwischen Ereignis und Sprecherereignis und dass Sprecherdaten in zwei Bereichen abgelegt werden. Damit möchten wir übermäßige Redundanzen in der Datenbasis vermeiden.

Im Laufe der Arbeiten an den Schemata haben wir ein neues System von Kennungen entwickelt. Kennungen sind systematische, eindeutige Kurzbezeichnungen für die Bezugsobjekte der Dokumentation: Ereignis, Sprechereignis, Sprecher, verschiedene Korpusbestandteile und deren technische Fassungen. Diese aufeinander abgestimmten Kurzbezeichnungen werden im Zusammenhang mit den Schemata in den Abschnitten 5. bis 7. vorgestellt.

4. Generische Schemata und projektspezifische Subschemas

Um zwischen projektübergreifenden und projektspezifischen Anforderungen vermitteln zu können, wurden zunächst umfangreiche Sammlungen von Informationselementen für die Bereiche Ereignisdaten und Sprecherdaten zusammengestellt, strukturiert und in (XML-)Schemata übertragen.

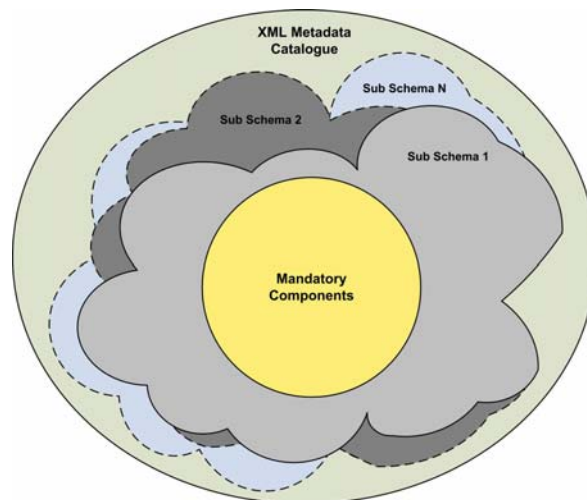


Abb. 2, Generisches Schema (Catalogue) und projektspezifische Subschemas

Mit diesen generischen Schemata werden Standards gesetzt. Sie enthalten obligatorische und fakultative Komponenten, Felddefinitionen und Standardwerte und bilden die Grundlage für projektspezifische Subschemas. Bei der Ableitung eines Subschemas aus dem generischen Schema müssen zunächst alle obligatorischen Komponenten übernommen werden. Diese werden ergänzt durch eine Auswahl fakultativer Komponenten, die für das auswählende Projekt verbindlich werden. Darüber hinaus können einzelne, in den generischen Schemata vorgegebene Werte an Projektbedürfnisse angepasst werden. Diese Anpassung geschieht durch die Spezifikation projektspezifischer Muster, mit denen die eingegebenen Werte schon bei der Erfassung verglichen werden, und eine Vorbelegung von Feldern mit projektspezifischen Werten, u.a. in Form von Auswahllisten, wobei die Vorgaben verschiedener Projekte koordiniert werden sollten.

5. Generisches Schema für die Dokumentation von Korpusbestandteilen auf der Ereignisebene

Die Kategorie „Ereignis“ dient als Startknoten eines generischen (XML-)Schemas, das folgende Informationen vorsieht: Angaben über Aufzeichnungsobjekte (Ereignis, Sprechereignis, Sprecher), Angaben über Korpusbestandteile (Audioaufnahmen, Videoaufnahmen, Transkripte, Zusatzmaterialien auf Ereignis- und Sprechereignisebene) sowie eine Dokumentationsgeschichte.

Das Schema enthält obligatorische und fakultative Komponenten. Obligatorische Komponenten sind in allen projektspezifischen Subschemas zu berücksichtigen, fakultative Komponenten

stehen zur Wahl. Wenn Sie verwendet werden, müssen alle Kennungsfelder und alle Felder, die ein Fragezeichen (?) enthalten, bearbeitet werden.

Eingaben für fehlende Daten in Feldern mit Fragezeichen (?) sind standardisiert: „Nicht dokumentiert“ bedeutet: Es kann ein Datum geben, das bei der Datenerfassung jedoch nicht bekannt ist. Ein Beispiel dafür wäre: „Sonstige_Bezeichnung: Nicht dokumentiert“ - zu lesen als: „Im Projekt kann eine andere Kurzbezeichnung als die Ereigniskennung vergeben worden sein, die jedoch nicht bekannt ist.“ „Nicht vorhanden“ bedeutet: Es gibt kein Datum. Ein Beispiel dafür wäre: „An_E_teilnehmende_Techniker: Nicht vorhanden“ - zu lesen als: „An diesem Ereignis hat kein Techniker teilgenommen.“ In einem Feld („Datenrate“) kann der Wert „Nicht relevant“ verwendet werden.

Das an vielen Stellen vorgesehene Feld „Anmerkungen“ ist für Anmerkungen zu Angaben in anderen Feldern und für nicht kategorisierte Angaben vorgesehen. Das Feld kann leer bleiben.

Einzelne Komponenten des Schemas wurden als iterativ gekennzeichnet, d.h. dass sie bei der Datenerfassung vervielfältigt werden können.

Die nachfolgenden Abbildungen stammen aus einem projektneutralen Erfassungsformular, das nur zu Demonstrationszwecken angelegt wurde.

5.1. Ereignis

Unter „Ereignis“ (E) verstehen wir eine Phase des sozialen Geschehens, die von Beteiligten bzw. Korpusproduzenten als abgrenzbare Einheit wahrgenommen und aufgezeichnet wird. Diese Definition ist aus arbeitspraktischen Gründen bewusst sehr allgemein gehalten. Wir stellen lediglich ein für die Dokumentation von Korpusbestandteilen relevantes Konzept bereit, keine linguistischen Segmentierungskriterien. Zur Veranschaulichung unseres Ereigniskonzeptes nennen wir im Folgenden drei Beispiele:

Im Korpusprojekt „Deutsch heute“ gelten mehrstündige Aufnahmesitzungen in Schulen und Volkshochschulen, die von Projektmitarbeitern geleitet wurden, als zu dokumentierende Ereignisse. Das IDS-Korpus „Stadtssprache: Mannheim“ enthält u.a. Aufzeichnungen von Treffen sozialer Gruppen in bestimmten Stadtteilen. Jedes dieser Gruppentreffen kann als Ereignis dokumentiert werden. Ein im IDS-Korpus „Biographische und Reiseerzählungen“ aufgezeichnetes Ereignis wurde folgendermaßen beschrieben: „Gemeinsames Kaffeetrinken während einer Seminarpause. Das Treffen zwischen Studentinnen und Dozentinnen wurde organisiert, um Reiseerzählungen aufzuzeichnen.“

IDS-Schema für die Dokumentation von Korpusbestandteilen auf der Ereignisebene

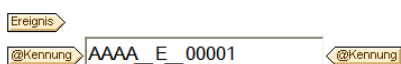


Abb. 3, Ereignis - Kennung

An erster Stelle der Ereignisdaten wird eine Kennung eingetragen, die eine vierstellige Korpuskennung, den Kennbuchstaben E (für „Ereignis“) und eine fünfstelligen laufende Nummer umfasst. Ein Beispiel für eine Ereigniskennung finden Sie in Abb. 3. Die Kennung eines zweiten Ereignisses müsste im o.g. Beispiel AAAA_E_00002 lauten.

5.1.1. Basisdaten

Die Ereignis-Basisdaten werden mit dem Feld „Sonstige_Bezeichnung“ eröffnet. Damit sind projektinterne Kurzbezeichnungen des Ereignisses gemeint. Im IDS-Projekt „Deutsch heute“ z.B. wurden dreistellige Ortskürzel wie „ODF“ (für „Oberstdorf“) benutzt.

Im Anschluss daran wird eine kurze Charakterisierung des aufgezeichneten Ereignisses erwartet. Hier sollte man nach Möglichkeit auch angeben, ob das Ereignis geplant oder nicht geplant war sowie ob und wann die Beteiligten über die Aufnahmen informiert wurden. Unsere erste Beispielbeschreibung stammt aus der Dokumentation des Korpus „Deutsch heute, die zweite aus der Dokumentation des Korpus Grundstrukturen: Freiburger Korpus: 1.) „Geplante Aufnahmeaktion im Rahmen des Spracherhebungsprojekts Deutsch heute, wobei von jedem Sprecher das gleiche Material erhoben wird (Lesesprache, Interview, Maptask). Die Sprecher waren im Vorfeld über die Aufzeichnungen informiert worden.“ 2.) „Lesung von Günther Grass aus "Katz und Maus" und anschließende Diskussion..“

Im Komplex „Ort“ sind Angaben über den Ort zusammengefasst, an dem das jeweilige Ereignis stattfand. Für Werte im Feld „Land“ gibt es eine ISO-Liste. Das Feld „Region“ ist für die amtliche Bezeichnung eines Bundeslandes, eines Kantons oder einer Provinz vorgesehen. In die Felder „Kreis“ und „Ortsname“ sollen ebenfalls amtliche Bezeichnungen eingetragen werden. Die Komponente „Koordinaten“ ist fakultativ. Wenn sie von einem Projekt gewählt wird, ist entweder der Geocode oder das Planquadrat (Kategorie im DSAV-Katalog [9]) des Ortes zu verzeichnen. Im Feld „Ortsteil“ kann der Name des Ortsteils, in dem das Ereignis stattfand, notiert werden. Weitere Informationen über den Ort des Ereignisses kann man im Feld „Ortsbeschreibung“ festhalten.

The image shows a data entry form with two main sections: 'Basisdaten' and 'Ort'. The 'Basisdaten' section includes fields for 'Sonstige_Bezeichnung' (with a question mark), 'Beschreibung' (with a question mark), and 'Einwohnerzahl' (with the value '0'). The 'Ort' section includes a dropdown for 'Land', text boxes for 'Region', 'Kreis', and 'Ortsname', a text box for 'Geografische_Breite' (with '000.0'), a text box for 'Geografische_Länge' (with '000.0'), a dropdown for 'Geocode', a text box for 'Planquadrat' (with a question mark), a text box for 'Ortsbeschreibung' (with a question mark), and a text box for 'Ortsteil' (with a question mark). There are also several 'Anmerkungen' (notes) fields. The form is styled with a light background and orange accents.

Abb. 4, Ereignis - Basisdaten (1)

Institution	?	Institution
Räumlichkeiten	?	Räumlichkeiten
Datum		
YYYY-MM-DD	9999-01-01	YYYY-MM-DD
Anmerkungen		Anmerkungen
Datum		
Dauer	?	Dauer
Zeitraum	?	Zeitraum
Aufnahmebedingungen	?	Aufnahmebedingungen
Anmerkungen		Anmerkungen
Basisdaten		

Abb. 5, Ereignis - Basisdaten (2)

Die Bezeichnung „Institution“ verstehen wir im Sinne von „Organisation“. Im DH-spezifischen Erfassungsformular z.B. wird an dieser Stelle der Name der Schule bzw. Volkshochschule, in der Aufnahmen gemacht wurden, verzeichnet. In der Dokumentation des Korpus „Schlichtungs- und Gerichtsverhandlungen“ z.B. gibt es an dieser Stelle u.a. die Werte „Vergleichsbehörde“, „Schlichtungsstelle der Handwerkskammer“ oder „Amtsgericht“. Im Feld „Räumlichkeiten“ kann man Angaben zur räumlichen Umgebung des Ereignisses ablegen.

Das Modul „Datum“ ist für Angaben über das Datum, an dem das Ereignis stattfand, vorgesehen. Wenn sich ein Ereignis über mehrere Tage erstreckte, sollte nur das Anfangsdatum aufgenommen werden. Das zugehörige Feld „Anmerkungen“ bietet die Möglichkeit, auf ungenaue Daten hinzuweisen. Die Dauer des aufgezeichneten Ereignisses ist im gleichnamigen Feld zu erfassen. An dieser Stelle sind auch ungenaue Angaben wie z.B. „Ca. 6 Stunden“ möglich. Über den Zeitraum, innerhalb dessen ein Ereignis stattfand, kann im gleichnamigen Feld z.B. folgendermaßen informiert werden: „Von 16 bis 17 Uhr“, „Vormittags“, „Zwei Tage“. Ereignisse, die zu völlig unterschiedlichen Zeiten stattfanden, gelten als unterschiedliche Ereignisse, die in verschiedenen Dokumenten zu beschreiben sind.

Unter dem Stichwort „Aufnahmebedingungen“ sollten Angaben über besondere Umstände der Aufzeichnungsaktion, z.B. Zeitdruck oder problematische akustische Verhältnisse, erfasst werden.

5.1.2. Rundfunksendung

Rundfunksendung		
@Rundfunktyp	?	@Rundfunktyp
Rundfunkanstalt	?	Rundfunkanstalt
Organisationsform	?	Organisationsform
Programm	?	Programm
Titel_Sendung	?	Titel_Sendung
Themen	?	Themen
Sendedatum		
YYYY-MM-DD	9999-01-01	YYYY-MM-DD
Anmerkungen		Anmerkungen
Sendedatum		
Sendezeit	?	Sendezeit
Sendeform	?	Sendeform
Sendart	?	Sendart
Beteiligte	?	Beteiligte
Anmerkungen		Anmerkungen
Rundfunksendung		

Abb. 6, Rundfunksendung

Übertragungen eines Ereignisses im Rundfunk können im fakultativen Komplex „Rundfunksendung“ dokumentiert werden. Für den Fall, dass Mitschnitte von mehreren Übertragungen zu dokumentieren sind, wurde der Komplex im Schema als iterativ gekennzeichnet.

An erster Stelle ist der Rundfunktyp - Hörfunk oder Fernsehen - anzugeben. Es folgen Felder für den Namen der Rundfunkanstalt, die eine Sendung verantwortet (z.B. „Bayerischer Rundfunk“, „Hessischer Rundfunk“, „Österreichischer Rundfunk“, „ZDF“) und für eine Information über deren Organisationsform (z.B. „Öffentlich-rechtlich“ oder „Privat“). Im Feld „Programm“ soll der „Kanal“ verzeichnet werden. Mögliche Werte sind z.B. „Kinderkanal“, „SWR-1“, „WDR-1-regional“, „RBB-Berlin“. [10]

Für den gesamten Titelkomplex einer Sendung steht das Feld „Titel_Sendung“ bereit. Hier können Haupt- und Untertitel von einzelnen Sendungen und Sendereihen eingetragen werden. Soweit Angaben über das Thema einer Sendung nicht im Titel enthalten sind, kann man sie im Feld „Themen“ verzeichnen.

Die Komponente „Sendedatum“ ist für Angaben über das Datum, an dem das Ereignis übertragen wurde, vorgesehen. Es folgt ein Feld für die Sendezeit, wo die Anfangs- und die Endzeit eingetragen werden sollten. Im Feld „Sendeform“ können Bezeichnungen wie „Bericht“, „Diskussion“, „Interview“, „Magazin“, „Nachrichten“, „Show“ etc. erfasst werden. Das Feld „Sendertyp“ ist für Hinweise wie z.B. „Erstausstrahlung“ oder „1. Wiederholung“ etc. gedacht.

Für die Namen von Produzenten (dazu zählen u.a. Regisseure und Moderatoren) und sonstigen Mitwirkenden (z.B. Gäste bei einer Talkshow) wurde das Feld „Beteiligte“ eingerichtet.

5.1.3. Projekt

In der Komponente „Projekt“ werden Angaben über das Projekt, in dem das Ereignis aufgezeichnet wurde, zusammengefasst. Um Kooperationen zwischen mehreren Projekten gerecht werden zu können, wurde die Komponente als iterativ gekennzeichnet.

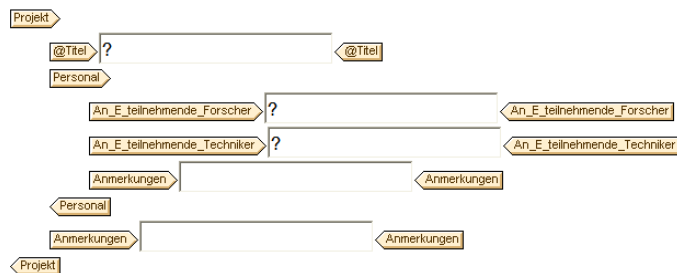


Abb. 8, Ereignis - Projekt

Wir nehmen an, dass die in Abb. 8 dargestellten Kategorien keiner Erläuterungen bedürfen.

5.1.4. Quellaufnahme

Unter „Quellaufnahmen“ verstehen wir Rohdaten, Originalaufnahmen von Ereignissen oder Aufnahmen, die für die dokumentierende Stelle Originalcharakter haben. Diese Aufnahmen können Quellen für sprechereignisspezifische Kopien sein. Da ein Korpus nicht unbedingt Originalaufnahmen umfassen muss, ist dieser Komplex fakultativ. Um mehrere Quellaufnahmen eines Ereignisses dokumentieren zu können, haben wir ihn im Schema als iterativ gekennzeichnet.

Quellaufnahme

@Kennung > AAAA_E_00001_A_01 < @Kennung

Abb. 9, Quellaufnahme - Kennung

An erster Stelle der Dokumentation einer Quellaufnahme steht eine Kennung, die sich zusammensetzt aus der Kennung des Ereignisses, dem Kennbuchstaben A (für „Aufnahme“) und einer zweistelligen laufenden Nummer. Ein Beispiel für die Kennung einer ersten Quellaufnahme finden Sie in Abb. 9. Die Kennung der zweiten Quellaufnahme im o.g. Beispiel lautet: AAAA_E_00001_A_02.

5.1.4.1. Basisdaten

Im ersten Feld der Quellaufnahme-Basisdaten werden „Sonstige_Bezeichnungungen“ abgefragt. Damit sind eventuell im Projekt vergebene Kurzbezeichnungen gemeint. Für eine Quellaufnahme aus dem Korpus „Stadtsprache: Mannheim“ z.B. könnte man an dieser Stelle eine im Projekt gewählte „Diskursnummer“ (z.B. „2036.50“) eintragen.

Quellaufnahmen können unterschiedlichen Typs sein: Audioaufnahme, Videoaufnahme und ggf. auch Tonkopie von Videoaufnahme. Für Angaben über die Dauer der jeweiligen Quellaufnahme wurde ein Zeitfeld mit dem Format hh:mm:ss vorbereitet.

Quellaufnahmen können Daten enthalten, die nach dem Willen der Urheber und aus datenschutzrechtlichen Gründen Außenstehenden nicht kenntlich werden dürfen, wie z.B. persönliche Sprecherdaten. Für Informationen über solche Daten ist das Feld „Schutzbedürftige_Daten“ vorgesehen.

Basisdaten

Sonstige_Bezeichnungungen ? < Sonstige_Bezeichnungungen

Typ ? < Typ

Dauer 00:00:00 < Dauer

Schutzbedürftige_Daten ? < Schutzbedürftige_Daten

Qualität

Aufnahmeablauf ? < Aufnahmeablauf

Sprachlich ? < Sprachlich

Anmerkungen < Anmerkungen

Relation_zu_E

Vollständigkeit ? < Vollständigkeit

Zeitabschnitt ? < Zeitabschnitt

Anmerkungen < Anmerkungen

Relation_zu_E

Relation_zu_anderer_Quellaufnahme

@Kennung_anderer_Quellaufnahme > AAAA_E_00001_A_02 < @Kennung_anderer_Quellaufnahme

Art ? < Art

Anmerkungen < Anmerkungen

Relation_zu_anderer_Quellaufnahme

Anmerkungen < Anmerkungen

Basisdaten

Abb. 10, Quellaufnahme - Basisdaten

Die Komponente „Qualität“ ist fakultativ. Wir haben die Felder „Aufnahmeablauf“ und „Sprachlich“ v.a. im Hinblick auf eine mögliche Übernahme der DSAV-Katalogdaten in die Struktur eingesetzt. Auf Seite 5 des DSAV-Gesamtkatalogs [9] sind Kriterien zusammengestellt, die zu einer schlechten Beurteilung katalogisierter Aufnahmen führten: a) Sprachliche Qualität - „Zahnschäden, Nuscheln, Lispeln, Mund- oder Prothesengeräusche, Stottern, scharfe, zischende s oder z, Asthma, Heiserkeit, [...] starke Befangenheit oder Erregung, zu leise oder zu laute

Sprache, [...]“; b) Aufnahmeablauf - „Pausen, Zwischenfragen, Absätze, stockendes Erzählen, Dazwischenreden.“

Unter der Überschrift „Relation_zu_E“ werden Informationen über das Verhältnis der jeweiligen Quellaufnahme zum Ereignis erfasst. Im Feld „Vollständigkeit“ wird eingetragen, ob eine vollständige oder eine unvollständige Aufnahme eines Ereignisses dokumentiert wird. Die Angabe „Unvollständig“ sollte im Feld „Zeitabschnitt“ präzisiert werden. „Zeitabschnitt“ meint den in der jeweiligen Quellaufnahme aufgezeichneten Zeitabschnitt des Ereignisses. Wenn die genaue Zeit nicht zu ermitteln ist, kann hier darüber informiert werden, welcher Abschnitt des Ereignisses in der jeweiligen Aufnahme festgehalten ist, z.B. „1. Abschnitt“, „2. Abschnitt“ usw. Bei vollständigen Aufnahmen wird der Wert „Vollständig“ erwartet.

Die Komponente „Relation_zu_anderer_Quellaufnahme“ ist fakultativ und iterativ. Wenn es nur eine Quellaufnahme gibt, ist sie nicht relevant. An dieser Stelle kann man z.B. auf Überlappungen von Quellaufnahmen hinweisen. Benötigt werden die Kennung der anderen Quellaufnahme und eine Information über die Art der Beziehung.

5.1.4.2. Aufnahmetechnik

Unter dem Stichwort „Aufnahmetechnik“ werden Informationen über die Aufnahmeapparatur (Aufnahmegerät, Mikrofone), eine ggf. eingesetzte Aufzeichnungssoftware, die Aufnahmege-
schwindigkeit (bei Spulentonbandaufnahmen relevante Angabe in cm/s) und Rauschunterdrückungsverfahren (z.B. Dolby B) zusammengefasst.

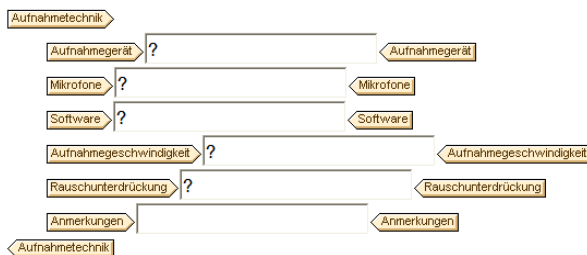


Abb. 11, Quellaufnahme - Aufnahmetechnik

5.1.4.3. Technische Fassungen

Quellaufnahmen liegen in bestimmten technischen Fassungen vor. Das können analoge und / oder digitale Fassungen sein. Für jeden Typ gibt es einen eigenen Abschnitt im Schema. Beide Abschnitte sind fakultativ und iterativ, wenigstens ein Abschnitt muss bei der Erstellung projektspezifischer Schemata übernommen werden.

Für jede technische Fassung wird eine Kennung generiert. Diese Kennung ist zusammengesetzt aus der Kennung der jeweiligen Quellaufnahme, dem Kürzel AF (für „Analoge_Fassung“) bzw. DF (für „Digitale_Fassung“) und einer zweistelligen laufenden Nummer.

Im Anschluss an die Kennungen werden Typen analoger und digitaler Fassungen benannt. Grundlage für eine Typisierung analoger Fassungen sind Datenschutz und Kanäle, für die Typisierung digitaler Fassungen sind außerdem noch technische Daten (z.B. das Dateiformat) relevant. Als Typenbezeichnungen dienen die Kürzel AFT (analoge Fassung) und DFT (digitale Fassung) in Verbindung mit einer zweistelligen Nummer.

Beispiele für Kennungen und Typenbezeichnungen finden Sie in den Abb. 12 und 14.

Abb. 12, Quellaufnahme - Analoge Fassung

Das Modul „Datum“ ist für Angaben über das Erstellungsdatum der technischen Fassung vorgesehen. „Datenschutz“ meint technische Maßnahmen zum Datenschutz, wie z.B. die Überlagerung von Personennamen in Aufnahmen. Das Feld „Kanäle“ steht für Angaben wie „Mono“ oder „Stereo“ bereit.

Die Komponente „Bewertung“ im Abschnitt „Qualität“ ist fakultativ. Die Bewertung unterschiedlicher Aspekte der Aufnahmequalität soll anhand einer Skala erfolgen, die in Abb. 13 dargestellt wird.

Abb. 13, Aufnahmequalität - Bewertungsskala

Das Feld „Allgemein“ soll eine Bewertung des Gesamteindrucks aufnehmen. Da uns für die anderen Aspekte kein neuer Kriterienkatalog vorliegt, verwenden wir zur Veranschaulichung ältere Hinweise aus dem DSAv-Gesamtkatalog [9] (S. 5).

Die akustische Aufnahmequalität können lt. Katalog beeinträchtigen: „halliger Aufnahmeraum, Unruhe bzw. Störgeräusche im Sprecherraum oder von außen, Sprecher zu nahe am Mikrofon (harte p, t, k), wechselnder bzw. ungünstiger Mikrofonabstand [sic!] zum Sprecher“. Folgende Erscheinungen führen lt. Katalog zu einer negativen Bewertung der technischen Aufnahmequa-

lität: „Netzbrummen, Leitungston, Kopiereffekt, verzerrte Modulation durch elektrische Fehler an den Geräten, Bandfehler“. In dieser älteren Aufzählung fehlt das bekannte Phänomen „Verzerrung“, das durch Übersteuerung bei analogen Aufnahmen (und in neuerer Zeit auch durch Fehler bei der Digitalisierung) entstehen kann.

Das Feld „Optisch“ ist für eine Bewertung der optischen Qualität von Videoaufnahmen vorgesehen. Im Feld „Probleme“ können konkrete Qualitätsprobleme benannt werden. Mögliche Werte sind z.B.: „Lauter Verkehrslärm“ ; „Leises Brummen“ ; „Die Aufnahme ist stellenweise verzerrt“ ; „Die Aufnahme ist an vielen Stellen unverständlich, da oft durcheinander geredet wird“.

Da eine technische Fassung auf verschiedenen Datenträgern gespeichert sein kann, haben wir das entsprechende Modul im Schema als iterativ gekennzeichnet. An erster Stelle wird eine eindeutige Inventarnummer des zu dokumentierenden Datenträgers erwartet. Im Feld „Sonstige_Bezeichnung“ können in einem Korpusprojekt möglicherweise generierte Ordnungskennzeichen (Name, Ordnungsnummer, etc.) erfasst werden. Im nächsten Schritt sollte man über den Typ des Datenträgers (z.B. Kompaktkassette oder Tonband) informieren.

Für den Abschnitt „Digitale_Fassung“ sind außer den bisher beschriebenen Komponenten die Felder „Digitalisierungssoftware“ und „Elektronische Speicherort“ sowie das Modul „Technische_Daten“ relevant.

Im Feld „Digitalisierungssoftware“ sollte notiert werden, welche Software ggf. bei der Digitalisierung einer analogen Fassung verwendet wurde. [11] „Elektronischer_Speicherort“ verlangt einen Pfadnamen oder eine URL.

The image shows a complex form titled "Digitale_Fassung" with multiple sections and fields. The fields are as follows:

- @Kennung:** A text field containing "AAAA_E_00001_A_01_DF_01".
- Basisdaten:**
 - Typ:** A dropdown menu with "DFT_01" selected.
 - Datum:** A date field with "9999-01-01" entered.
 - Anmerkungen:** Two empty text input fields.
- Digitalisierungssoftware:** A dropdown menu with "?" selected.
- Datenschutz:** A dropdown menu with "?" selected.
- Kanäle:** A dropdown menu with "?" selected.
- Qualität:**
 - Bewertung:** Four dropdown menus labeled "Allgemein", "Akustisch", "Optisch", and "Technisch", each with "?" selected.
 - Probleme:** A text input field with "?".
 - Anmerkungen:** A text input field.
- Datenträger:**
 - @Inventarnummer:** A dropdown menu with "?".
 - Sonstige_Bezeichnung:** A text input field with "?".
 - Typ:** A dropdown menu with "?".
 - Anmerkungen:** A text input field.

Abb. 14, Quellaufnahme - Digitale Fassung (1)

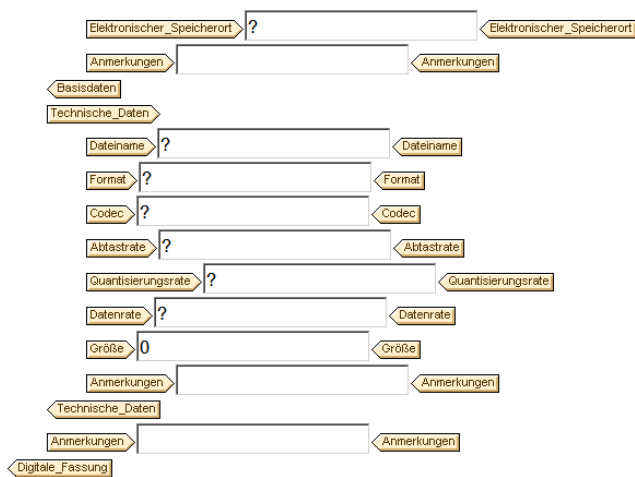


Abb. 15, Quellaufnahme - Digitale Fassung (2)

Im Modul „Technische_Daten“ ist zunächst der Dateiname anzugeben. Für die folgenden Kategorien bieten wir hier lediglich kurze Erläuterungen an. Weitere Informationen finden Sie z.B. im Gesprächsanalytischen Informationssystem (GAIS) <http://gais.ids-mannheim.de/> und unter den u.g. Adressen.

Im Feld „Format“ werden Angaben über das Dateiformat und das Digitalisierungsverfahren bzw. Audioformat erfasst. Ein möglicher Feldwert wäre: „WAVE (Linear PCM)“. Informationen über relevante Audioformate finden Sie unter der Adresse <http://de.wikipedia.org/wiki/Audioformat>.

Als „Codec“ bezeichnet man ein Verfahren bzw. Programm, das Daten oder Signale digital kodiert und dekodiert. Unter der Adresse <http://de.wikipedia.org/wiki/Codec> finden Sie eine Liste mit Namen gängiger Codecs.

„Abtastung“ (engl. sampling) bezeichnet die Registrierung von Messwerten zu diskreten, meist äquidistanten Zeitpunkten. Aus einem zeitkontinuierlichen Signal wird so ein zeitdiskretes Signal gewonnen. Die Anzahl der Abtastungen pro Zeiteinheit wird Abtastrate genannt und meist in Hertz (Hz = Anzahl pro Sekunde) angegeben. Mögliche Werte sind z.B. „44100“ oder „48000“.

Nach der Abtastung erfolgt die Quantisierung des zeitdiskreten, aber noch wertkontinuierlichen Signals. Dadurch entsteht ein zeit- und wertdiskretes Signal. Die Quantisierungsrate (auch Samplingtiefe oder Bittiefe) gibt die Anzahl der Bits an, die bei der Quantisierung pro Abtastwert verwendet werden. Typische Quantisierungsraten sind 8, 16 und 24 Bit.

Bei komprimierten Daten wird die Datenrate relevant - die Anzahl der Informationseinheiten, die pro Zeiteinheit gespeichert werden. Sie wird in kBit/s angegeben. Die Größe einer digitalen Fassung soll in Bytes erfasst werden.

5.1.4.4. Archivierung und Distribution

In den Abb. 16 und 17 werden die Bausteine „Archivierung“ und „Distribution“ vorgestellt, die an verschiedenen Stellen des Ereignisschemas vorkommen. Sie sollen Informationen über rechtliche und organisatorische Aspekte der Korpusbestandteile aufnehmen.

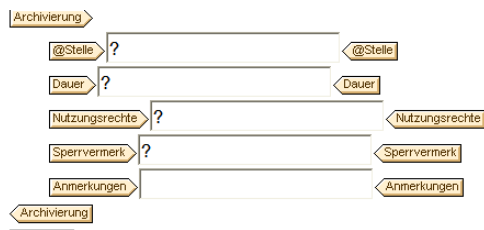


Abb. 16, Quellaufnahme - Archivierung

Das Modul „Archivierung“ wurde im Schema als iterativ gekennzeichnet. Zunächst soll der Name der archivierenden Stelle vermerkt werden. Es folgt ein Feld für Informationen über die vorgesehene Archivierungsdauer. Hier könnte z.B. „Bis 2018“ oder „Langfristig“ stehen. Die Nutzungsrechte der archivierenden Stelle können von der ausschließlichen wissenschaftlichen Auswertung durch den Aufnahmeleiter bis hin zur Veröffentlichung einer Aufnahme im Internet reichen. Sperrvermerke, wie z.B. „Bis 2010 für Externe gesperrt“, können die Nutzungsmöglichkeiten einschränken.

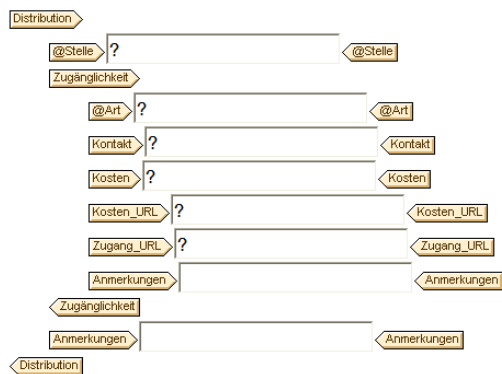


Abb. 17, Quellaufnahme - Distribution

Das Modul „Distribution“ ist ebenfalls iterativ und umfasst neben einem Feld für den Namen der für die Distribution zuständigen Stelle die iterative Komponente „Zugänglichkeit“. In dieser Komponente sollen folgende Angaben verzeichnet werden: Art der Zugänglichkeit, E-Mail-Kontaktadresse, Angaben über die Kosten, ggf. eine URL dieser Angaben sowie ggf. eine URL, die einen direkten Zugang zum jeweiligen Korpusbestandteil ermöglicht.

5.1.5. Zusatzmaterial

Unter „Zusatzmaterial“ auf der Ereignissebene verstehen wir Dokumente, die zusätzlich zu Quellaufnahmen vorhanden sein können. Das können z.B. Reiseberichte der Aufnahmeleiter sein, Protokolle von Aufnahmesitzungen, Fotos von Aufnahmeorten, Notizen zu einer Sitzordnung etc. Der gesamte Komplex „Zusatzmaterial“ ist fakultativ. Für den Fall, dass mehrere Dokumente zu einem Ereignis zu beschreiben sind, wurde er im Schema als iterativ gekennzeichnet.

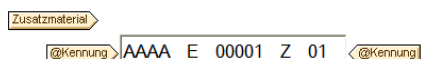


Abb. 18, Zusatzmaterial - Kennung

Die Kennung für Zusatzmaterial auf der Ereignissebene setzt sich zusammen aus der Ereigniskennung, dem Kennbuchstaben Z (für „Zusatzmaterial“) und einer zweistelligen laufenden Nummer. Ein Beispiel finden Sie in Abb. 18. Die Kennung eines zweiten zu dokumentierenden Dokuments müsste im o.g. Beispiel AAAA_E_00001_Z_02 lauten.

5.1.5.1. Basisdaten

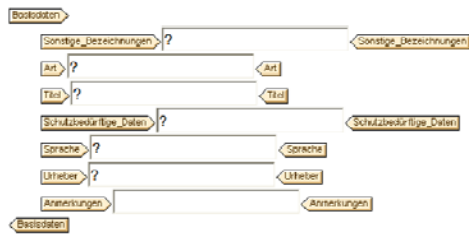


Abb. 19, Ereignis - Zusatzmaterial - Basisdaten

An erster Stelle der Basisdaten für Zusatzmaterial steht das Feld „Sonstige_Bezeichnungungen“, wo eine eventuell im Erstellungsprojekt für ein Dokument vergebene Kurzbezeichnung abgelegt werden kann. Im Feld „Art“ wird eine Angabe über die Art des Dokuments (z.B. „Skizze der Sitzordnung“) erwartet. Zusatzmaterialien können Daten enthalten, die nach dem Willen der Urheber und aus datenschutzrechtlichen Gründen Außenstehenden nicht kenntlich werden dürfen, wie z.B. persönliche Sprecherdaten, über die man im Feld „Schutzbedürftige_Daten“ informieren kann. Die Sprache, in der ein Textdokument abgefasst ist, sollte ebenso vermerkt werden wie der Urheber eines Dokuments, wobei „Urheber“ für Autoren, Grafiker, Fotografen etc. verwendet wird.

5.1.5.2. Technische Fassungen

Da zusätzliche Dokumente in verschiedenen technischen Fassungen vorliegen können, haben wir auch an dieser Stelle die fakultativen und iterativen Komponenten „Analoge_Fassung“ und „Digitale_Fassung“ eingefügt. Wenigstens eine Komponente muss bei der Erstellung eines projektspezifischen Schemas gewählt werden.

Zunächst wird eine Kennung abgefragt, die aus der Kennung des Zusatzmaterials, dem Kürzel AF (für „Analoge_Fassung“) bzw. DF (für „Digitale_Fassung“) und einer zweistelligen laufenden Nummer besteht. Beispiele finden Sie in den Abb. 20 und 21.

Im Anschluss an die Kennungen werden Typen analoger und digitaler Fassungen benannt. Grundlage für eine Typisierung analoger Fassungen ist der Datenschutz, für die Typisierung digitaler Fassungen sind außerdem noch technische Daten (z.B. das Dateiformat) relevant. Als Typenbezeichnungen dienen die Kürzel AFT (analoge Fassung) und DFT (digitale Fassung) in Verbindung mit einer zweistelligen Nummer.

Beispiele für Kennungen und Typenbezeichnungen finden Sie in den Abb. 20 und 21.

Das Modul „Datum“ ist für Angaben über das Erstellungsdatum der technischen Fassung vorgesehen. „Datenschutz“ meint technische Maßnahmen zum Datenschutz, wie z.B. die Maskierung von Personennamen in Texten.

Da eine technische Fassung auf mehreren Datenträgern gespeichert sein kann, wurde dieser Abschnitt im Schema als iterativ gekennzeichnet. An erster Stelle dieses Abschnitts wird eine eindeutige Inventarnummer des zu dokumentierenden Datenträgers erwartet. Im Feld „Sonstige_Bezeichnungungen“ können weitere Ordnungskennzeichen (Name, Ordnungsnummer, etc.) erfasst werden. Im nächsten Schritt ist über den Typ des Datenträgers (z.B. Papier oder Mikrofilm) zu informieren.

Analoge_Fassung

@Kennung > AAAA_E_00001_Z_01_AF_01 <@Kennung

Typ > AFT_01 <Typ

Datum > YYYY-MM-DD 9999-01-01 <YYYY-MM-DD

Anmerkungen > <Anmerkungen

Datum > <Datum

Datenschutz > ? <Datenschutz

Datenträger > @Inventarnummer ? <@Inventarnummer

Sonstige_Bezeichnungungen > ? <Sonstige_Bezeichnungungen

Typ > ? <Typ

Datenträger > <Datum

Anmerkungen > <Anmerkungen

Analoge_Fassung

Abb. 20, Ereignis - Zusatzmaterial - Analoge Fassung

Digitale_Fassung

@Kennung > AAAA_E_00001_Z_01_DF_01 <@Kennung

Basisdaten > Typ > DFT_01 <Typ

Datum > YYYY-MM-DD 9999-01-01 <YYYY-MM-DD

Anmerkungen > <Anmerkungen

Datum > <Datum

Digitalisierungssoftware > ? <Digitalisierungssoftware

Datenschutz > ? <Datenschutz

Datenträger > @Inventarnummer ? <@Inventarnummer

Sonstige_Bezeichnungungen > ? <Sonstige_Bezeichnungungen

Typ > ? <Typ

Anmerkungen > <Anmerkungen

Datenträger > <Datum

Abb. 21, Ereignis - Zusatzmaterial - Digitale Fassung (1)

Elektronischer_Speicherort > ? <Elektronischer_Speicherort

Anmerkungen > <Anmerkungen

Basisdaten > Technische_Daten > Dateiname > ? <Dateiname

Format > ? <Format

Character_Encoding > ? <Character_Encoding

Größe > 0 <Größe

Anmerkungen > <Anmerkungen

Technische_Daten > Anmerkungen > <Anmerkungen

Digitale_Fassung

Abb. 22, Ereignis - Zusatzmaterial - Digitale Fassung (2)

Nur für digitale Fassungen relevant sind die Felder „Digitalisierungssoftware“ und „Elektronischer_Speicherort“ sowie die Komponente „Technische_Daten“. Im Feld „Digitalisierungssoftware“ soll das Programm, mit dem eine analoge Fassung digitalisiert wurde, genannt werden. Im Feld „Elektronischer Speicherort“ wird eine URL oder ein Pfadname erwartet.

Das erste Feld der Komponente „Technische_Daten“ ist für den Dateinamen vorgesehen, an den sich eine Information über das Dateiformat anschließen sollte. „Character_Encoding“ steht

für die Zeichencodierung in einer Textdatei (z.B. ASCII oder UTF-16BE). Im Feld „Größe“ ist die Dateigröße (Anzahl von Bytes) anzugeben.

5.1.5.3. Archivierung und Distribution

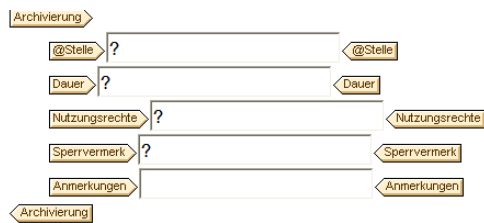


Abb. 23, Ereignis - Zusatzmaterial - Archivierung

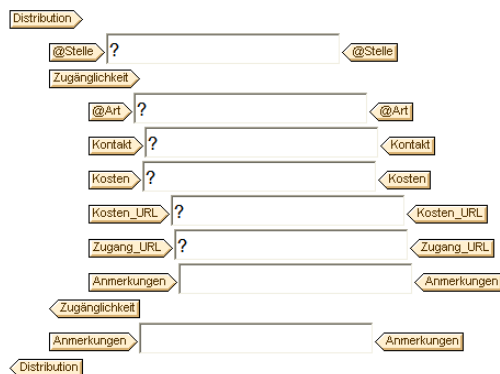


Abb. 24, Ereignis - Zusatzmaterial - Distribution

Die Module „Archivierung“ und „Distribution“, die Informationen über rechtliche und organisatorische Aspekte der Korpusbestandteile aufnehmen sollen, wurden bereits vorgestellt. Die Erläuterungen in Abschnitt 4.1.4.4. gelten auch für Zusatzmaterial.

5.1.6. Sprechereignis

Unter „Sprechereignis“ (SE) verstehen wir den aufgezeichneten sprachlichen / kommunikativen Gehalt eines Ereignisses bzw. Segmente dieses Gehalts. Hier wie in Abschnitt 5.1. gilt: Diese Definition ist aus arbeitspraktischen Gründen bewusst sehr allgemein gehalten. Wir stellen lediglich ein für die Dokumentation von Korpusbestandteilen relevantes Konzept bereit, keine linguistischen Segmentierungskriterien. Zur Veranschaulichung unseres Sprechereigniskonzeptes nennen wir im Folgenden einige Beispiele:

Im Korpusprojekt „Deutsch heute“ gilt jede Aufgabe, die im Rahmen einer mehrstündigen Aufnahmesitzung bearbeitet wurde, als Sprechereignis. Zu diesen Aufgaben gehören u.a. Bildbenennung, Verlesen einer Wortliste, Übersetzung und Interview. Im IDS-Korpus „Stadtsprache: Mannheim“ sind Aufnahmen von Gruppentreffen enthalten, bei denen z.B. Witze erzählt, Klatsch ausgetauscht und gemeinsame Unternehmungen geplant wurden. Solche kommunikativen Sequenzen können als einzelne Sprechereignisse dokumentiert werden. Das IDS-Korpus „Elizitierte Konfliktgespräche“ enthält Aufzeichnungen von Settings, in denen jeweils eine Mutter-Tochter-Dyade zwei Konfliktgespräche führte. Das Thema des ersten Gesprächs wurde von der Mutter eingebracht, das Thema des zweiten von der Tochter. Wir betrachten jedes dieser Gespräche als ein Sprechereignis.

Damit mehrere Sprechereignisse pro Ereignis dokumentiert werden können, wurde dieser Bereich im Schema als iterativ gekennzeichnet.

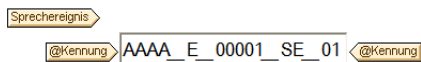


Abb. 25, Sprechereignis - Kennung

An erster Stelle der SE-Daten steht eine Kennung, die die jeweilige Ereigniskennung, das Kürzel SE (für „Sprechereignis“) und eine zweistellige laufende Nummer umfasst. Ein Beispiel finden Sie in Abb. 25. Die Kennung eines zweiten Sprechereignisses lautet im o.g. Beispiel: AAAA_E_00001_SE_02 .

5.1.6.1. Basisdaten

Zu Beginn der Basisdaten steht das Feld „Sonstige_Bezeichnung“, wo eventuell im Projekt vergebene Kurzbezeichnungen abgelegt werden können. Im Projekt kann auch ein Titel für das zu dokumentierende Sprechereignis vergeben worden sein, der im gleichnamigen Feld zu erfassen ist.

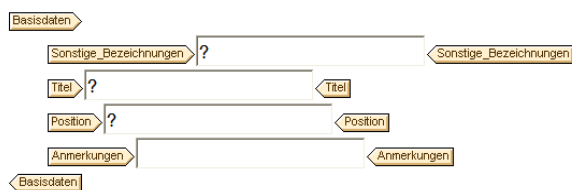


Abb. 26, Sprechereignis - Basisdaten

Die Position eines Sprechereignisses im Ereignis kann relevant sein, wenn Segmente des aufgezeichneten sprachlichen / kommunikativen Gehalts eines Ereignisses betrachtet werden. In solchen Fällen kann man hier die Zusammenhänge beschreiben. Eine mögliche Positionsbeschreibung wäre: „Beginnt unmittelbar nach der Begrüßung der Beteiligten und endet vor der ersten längeren Pause“.

5.1.6.2. Beschreibung

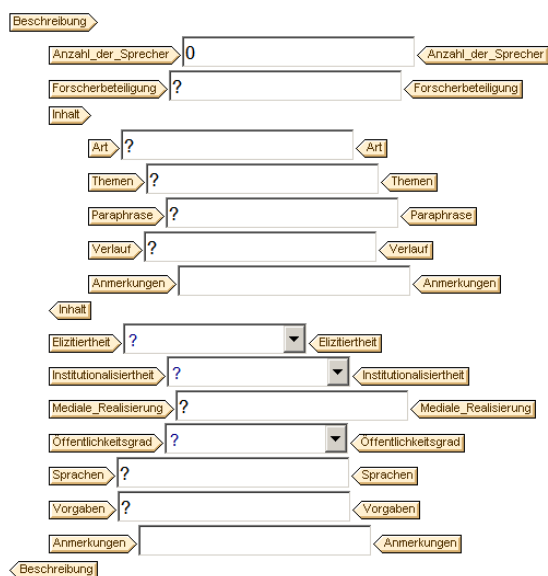


Abb. 27, Sprechereignis - Beschreibung

Im Komplex „Beschreibung“ wollen wir wesentliche systematische Aspekte eines Sprechereignisses darstellen.

Zunächst wird die Zahl der Sprecher notiert, wobei verbal beteiligte Forscher / Aufnahmeleiter mitgezählt werden sollten, was man dann im Feld „Forscherbeteiligung“ verdeutlichen kann. Für „Forscherbeteiligung“ haben wir die Werte „Verbal beteiligt“, „Nicht verbal beteiligt“ und „Nicht vorhanden“ (für „Forscher nicht anwesend“) vorgesehen. Bei Ereignissen, an denen mehrere Forscher teilgenommen haben, kann es nötig werden, den oder die beteiligten Forscher zu benennen. In solchen Fällen können die Namen den Werten „Verbal beteiligt“ bzw. „Nicht verbal beteiligt“ vorangestellt werden.

Das Modul „Inhalt“ umfasst die Elemente „Art“, „Themen“, „Paraphrase“ und „Verlauf“.

Wir verwenden „Art“ anstelle von Kategorien wie „Textsorte“, „Texttyp“, „Interaktionstyp“, „Gesprächstyp“, „Diskurstyp“, „Genre“, „Gattung“, die aus verschiedenen Forschungsansätzen stammen, um Daten aus allen Bereichen aufnehmen zu können. Nach unserer Vorstellung können in dieses Feld mehrere Angaben eingetragen werden. Wir denken dabei an Werte wie „Erzählung“, „Rede“, „Anleitung“, „Beschreibung“, „Benennung“, „Übersetzung“, „Interview“, „Beratung“, „Diskussion“, „Begrüßung“ etc. Mit diesen Beispielwerten wollen wir keine Vorentscheidung über eine im Einzelfall anzuwendende Systematik treffen.

Themenangaben sollten stichwortartig sein (z.B. „Politik“, „Recht“, „Studium“, „Lebenslauf“). Im Feld „Paraphrase“ kann man eine an die Werte in den Feldern „Art“ und „Themen“ anschließende Darstellung notieren. Wir nennen im Folgenden zwei fiktive Beispiele: Art: „Erzählung“ - Paraphrase: „Großvater erzählt seinem Enkel ein altes türkisches Märchen mit dem Titel xyz, das er in seiner Jugend in seinem türkischen Heimatdorf gehört hat.“ Art: „Beratung“ - Paraphrase: „Anwalt berät einen Klienten über rechtliche Möglichkeiten im Konflikt mit dessen Nachbarn.“

Informationen über den Verlauf des Sprechereignisses kann man im gleichnamigen Feld notieren. Das können einfache Hinweise wie z.B. „Sehr turbulent“ oder komplexere Angaben über die Entwicklung sein.

„Elizitierung“ ist eine Technik zur Erhebung sprachlicher Daten, bei der die Informanten systematisch zu Äußerungen veranlasst werden. Wir haben die Werte „Elizitiert“ und „Nicht elizitiert“ vorgesehen. „Institutionalisierung“ verstehen wir als Zugehörigkeit zu bzw. Erwartbarkeit eines Sprechereignisses im Rahmen einer Institution (im Sinne von „Organisation“). So fanden z.B. für das Korpus „Deutsch heute“ aufgezeichnete Sprechereignisse in Institutionen wie Schulen und Volkshochschulen statt, gelten in diesem Zusammenhang jedoch als nicht institutionell. Für Sprechereignisse, die im IDS-Projekt „Schlichtung“ aufgezeichnet wurden, wäre hier „Institutionell“ einzutragen.

„Mediale_Realisierung“ steht für den jeweiligen Kommunikationskanal (wie z.B. „Face to Face“, „Telefon“, „Hörfunk“). Für das Feld „Öffentlichkeitsgrad“ werden die Werte „Öffentlich“ und „Nicht öffentlich“ bereitgestellt. Im Feld „Sprachen“ sind die im Sprechereignis verwendeten Sprachen zu verzeichnen. Über Instruktionen von Sprechern durch Aufnahmeleiter und ggf. auch über Materialien, die den Sprechern zur Lösung bestimmter Aufgaben vorgelegt wurden, kann man im Feld „Vorgaben“ informieren.

5.1.6.3. Sprecher

Im Rahmen des Ereignisschemas werden hauptsächlich sprechereignisspezifische Sprecherdaten erfasst. Für allgemeine Informationen über Sprecher steht ein separates (XML-) Schema zur Verfügung (vgl. Abschnitt 6.). Der Sprecherkomplex ist fakultativ, für den Fall, dass keine

Sprecherdaten erhoben wurden. Da an einem Sprechereignis mehrere Sprecher beteiligt sein können, wurde der Komplex als iterativ gekennzeichnet.

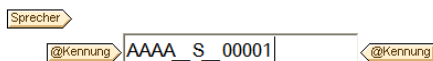


Abb. 28, Sprecher - Kennung

Die an erster Stelle der Sprecherdaten erwartete Kennung enthält die vierstellige Korpuskennung, den Kennbuchstaben S (für „Sprecher“) und eine fünfstelligen laufende Nummer: AAAA_S_00001, AAAA_S_00002 usw. Da mittels dieser Kennung eine Verbindung zu den allgemeinen Sprecherdaten (vgl. Abschnitt 6.) hergestellt wird, muss sie mit der dort für den Sprecher vergebenen Kennung übereinstimmen.

5.1.6.3.1. Basisdaten

In zwei Fällen weichen wir von unserem Prinzip, sprechereignisunabhängige Sprecherdaten in einem separaten Schema zu speichern, ab. Wir übernahmen „Alter“ und „Geschlecht“ in den Sprecherblock des Ereignisschemas aus praktischen Gründen, um Nutzern, die an diesen elementaren Daten interessiert sind, ein Umschalten auf die Dokumentation allgemeiner Sprecherdaten zu ersparen.

Das Element „Rolle“ ist für Angaben über die Beteiligungsrolle des jeweiligen Sprechers in dem zu dokumentierenden Sprechereignis vorgesehen. In den Basisdaten gibt es die Möglichkeit, Informationen über Besonderheiten eines Sprechers im Sprechereignis, wie z.B. „War anfangs sehr nervös“ oder „Stellte sich in hohem Maße auf seinen Gesprächspartner ein“, „Vermied Blickkontakte“ etc., zu erfassen. Ein Feld für sprachliche Besonderheiten ist in den „Sprachdaten“ (vgl. 5.1.6.3.2.) enthalten.

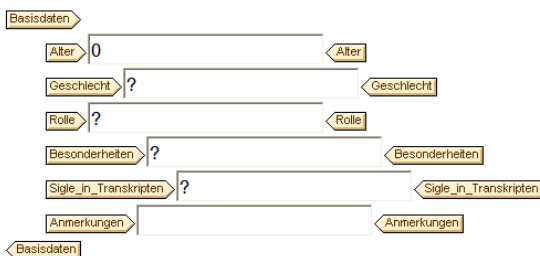


Abb. 29, Sprecher - Basisdaten

Am Ende der Basisdaten soll die in Transkripten verwendete Sprechersigle notiert werden. An dieser Stelle möchten wir darauf hinweisen, dass für einen Sprecher nicht mehrere Siglen vergeben werden sollten.

5.1.6.3.2. Sprachdaten

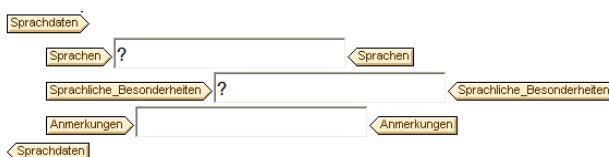


Abb.30, Sprechereignis - Sprecher - Sprachdaten

„Sprachdaten“ ist ein Bereich zur Erfassung von Angaben über die von einem Sprecher im jeweiligen Sprechereignis gesprochene(n) Sprache(n) und seine sprechereignisspezifischen sprachlichen Besonderheiten.

Zunächst sind die Namen der verwendeten Sprachen (z.B. „Deutsch ; Englisch ; Türkisch“) einzutragen. Unter dem Stichwort „Sprachliche Besonderheiten“ können weitere für das Sprechereignis charakteristische sprachliche Merkmale notiert werden. Die folgenden Beispiele stammen aus der 2004 veröffentlichten Dokumentation des Korpus „Emigrantendeutsch in Israel“ von Anne Betten: „Wiener Verkehrsmundart mit einzelnen österreichischen Dialekteigenheiten. Lebhafter Erzählstil.“ (IS010) „Sehr gewandtes, nuancenreiches Sprechen; überwiegend standardsprachlich korrekt formulierend, aber themenabhängig variierend von sehr verhaltenen bis zu drastisch-lebendigen Ausdrucksweisen. Stimme stark modulierend.“ (IS015) [12]

5.1.6.4. Sprechereignisspezifische Aufnahme

Wir nehmen an, dass es zu jedem dokumentierten Sprechereignis (mindestens) eine Aufnahme (SE-Aufnahme) gibt, die in einem bestimmten Verhältnis zu einer Quellaufnahme steht. Meist haben wir es mit einem sprechereignisspezifischen Teil in einer Quellaufnahme bzw. einer Kopie dieses spezifischen Teils zu tun. Es kommt allerdings auch vor, dass SE-Aufnahmen mit Quellaufnahmen vollkommen übereinstimmen. In solchen Fällen handelt es sich bei der angenommenen SE-Aufnahme um ein Konstrukt, das es uns erlaubt, an dieser Stelle einen Bezug zu einer Quellaufnahme herzustellen. Da mehrere SE-Aufnahmen vorliegen können, haben wir diesen Komplex im Schema als iterativ gekennzeichnet.

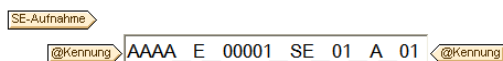


Abb. 31, SE-Aufnahme - Kennung

Die Kennung für eine SE-Aufnahme umfasst die jeweils aktuelle Sprechereigniskennung, den Kennbuchstaben A (für Aufnahme) und eine zweistellige laufende Nummer. Ein Beispiel dafür finden Sie in Abb. 31. Die Kennung für eine zweite SE-Aufnahme lautet in unserem Beispiel: AAAA_E_00001_SE_01_A_02.

5.1.6.4.1. Basisdaten

An erster Stelle der Basisdaten steht das Feld „Sonstige_Bezeichnung“. Damit sind eventuell im Projekt vergebene Kurzbezeichnungen gemeint. SE-Aufnahmen können unterschiedlichen Typs sein: Audioaufnahme, Videoaufnahme und ggf. auch Tonkopie von Videoaufnahme.

Für Angaben über die Dauer der jeweiligen SE-Aufnahme wurde ein fakultatives Zeitfeld mit dem Format hh:mm:ss vorbereitet. SE-Aufnahmen können Daten enthalten, die nach dem Willen der Urheber und aus datenschutzrechtlichen Gründen Außenstehenden nicht kenntlich werden dürfen, wie z.B. persönliche Sprecherdaten. Für entsprechende Informationen wurde das fakultative Feld „Schutzbedürftige_Daten“ bereitgestellt. Wenn die dokumentierten SE-Aufnahmen mit den Quellaufnahmen vollständig übereinstimmen, können diese beiden Felder bei der Erstellung projektspezifischer Schemata übergangen werden.

Unter der Überschrift „Relation_zu_Quellaufnahme“ kann man Informationen über das Verhältnis der jeweiligen SE-Aufnahme zu einer oder mehreren Quellaufnahmen erfassen. Dieses Modul ist iterativ. An erster Stelle ist die Kennung der Quellaufnahme (z.B. AAAA_E_00001_A_01) einzutragen. Eine SE-Aufnahme kann mit einer Quellaufnahme vollständig übereinstimmen oder ein Segment in einer Quellaufnahme sein. Diese Angaben sind im Feld „Vollständigkeit“ mit den Werten „Vollständig“ bzw. „Segment“ zu notieren. Die Angabe „Segment“ sollte im Feld

„Zeitabschnitt“ präzisiert werden. Bei vollständigen Aufnahmen wird hier der Wert „Vollständig“ erwartet.

Abb. 32, SE-Aufnahme - Basisdaten (1)

Unter der Überschrift „Relation_zu_SE“ werden Informationen über das Verhältnis der jeweiligen SE-Aufnahme zum Sprechereignis erfasst. Im Feld „Vollständigkeit“ wird eingetragen, ob eine vollständige oder eine unvollständige Aufnahme eines Sprechereignisses dokumentiert wird. Die Angabe „Unvollständig“ sollte man im Feld „Zeitabschnitt“ präzisieren. „Zeitabschnitt“ meint den in der jeweiligen SE-Aufnahme aufgezeichneten Zeitabschnitt des Sprechereignisses. Wenn die genaue Zeit nicht zu ermitteln ist, kann hier darüber informiert werden, welcher Abschnitt des Sprechereignisses in der jeweiligen Aufnahme festgehalten ist, z.B. „1. Abschnitt“, „2. Abschnitt“ usw. Bei vollständigen Aufnahmen wird der Wert „Vollständig“ erwartet.

5.1.6.4.2. Technische Fassungen

SE-spezifische Kopien von Quellaufnahmen liegen in ganz bestimmten technischen Fassungen vor. Die Komponenten „Analoge_Fassung“ und „Digitale_Fassung“ sind fakultativ und iterativ. Sollten alle SE-Aufnahmen eines Korpus mit den Quellaufnahmen identisch oder keine SE-spezifischen Kopien vorhanden sein, können beide Komponenten bei der Erstellung korpuspezifischer Schemata übergangen werden.

Die Strukturen dieser Komponenten stimmen mit den entsprechenden Strukturen im Komplex Quellaufnahmen überein. Erläuterungen dazu finden Sie in Abschnitt 5.1.4.3. Die Kennungen sind allerdings ebenenspezifisch und setzen sich hier zusammen aus der Kennung der SE-Aufnahme, dem Kürzel AF (für „Analoge_Fassung“) bzw. DF (für „Digitale_Fassung“) und einer zweistelligen Nummer (vgl. Abb. 33 und 34).

Analoge_Fassung

@Kennung AAAA_E_00001_SE_01_A_01_AF_01 @Kennung

Typ AFT_01 Typ

Datum

YYYY-MM-DD 9999-01-01 YYYY-MM-DD

Anmerkungen Anmerkungen

Datum

Datenschutz ? Datenschutz

Kanäle ? Kanäle

Qualität

Bewertung

Allgemein ? Allgemein

Akustisch ? Akustisch

Optisch ? Optisch

Technisch ? Technisch

Bewertung

Probleme ? Probleme

Anmerkungen Anmerkungen

Qualität

Datenträger

@Inventarnummer ? @Inventarnummer

Sonstige_Bezeichnung ? Sonstige_Bezeichnung

Typ ? Typ

Anmerkungen Anmerkungen

Datenträger

Anmerkungen Anmerkungen

Analoge_Fassung

Abb. 33, SE-Aufnahme - Analoge Fassung

Digitale_Fassung

@Kennung AAAA_E_00001_SE_01_A_01_DF_01 @Kennung

Basisdaten

Typ DFT_01 Typ

Datum

YYYY-MM-DD 9999-01-01 YYYY-MM-DD

Anmerkungen Anmerkungen

Datum

Digitalisierungssoftware ? Digitalisierungssoftware

Datenschutz ? Datenschutz

Kanäle ? Kanäle

Qualität

Bewertung

Allgemein ? Allgemein

Akustisch ? Akustisch

Optisch ? Optisch

Technisch ? Technisch

Bewertung

Probleme ? Probleme

Anmerkungen Anmerkungen

Qualität

Datenträger

@Inventarnummer ? @Inventarnummer

Sonstige_Bezeichnung ? Sonstige_Bezeichnung

Typ ? Typ

Anmerkungen Anmerkungen

Datenträger

Abb. 34, SE-Aufnahme - Digitale Fassung (1)

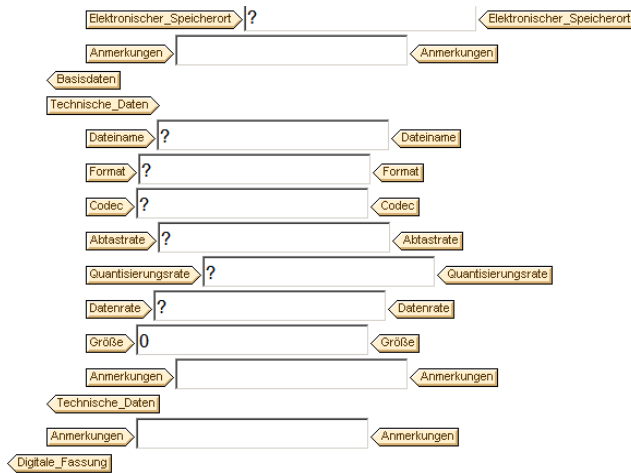


Abb. 35, SE-Aufnahme - Digitale Fassung (2)

5.1.6.4.3. Archivierung und Distribution

Auch mit den Modulen „Archivierung“ und „Distribution“ haben wir Sie schon bekannt gemacht. Sie werden in der Beschreibung aller Korpusbestandteile verwendet und in Abschnitt 5.1.4.4. eingeführt. Für Fälle, in denen alle SE-Aufnahmen eines Korpus mit den Quellaufnahmen identisch sind, wurden diese Verwaltungsdaten für SE-Aufnahmen als fakultativ gekennzeichnet.

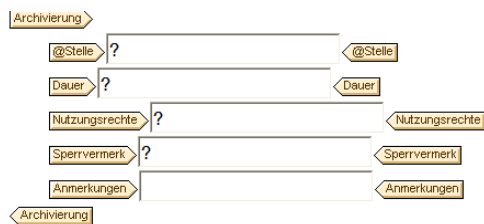


Abb. 36, SE-Aufnahme - Archivierung

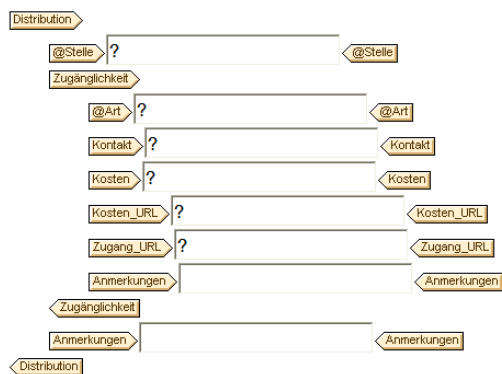


Abb. 37, SE-Aufnahme - Distribution

5.1.6.5. Transkript

Der gesamte Transkriptkomplex ist fakultativ. Da zu einer SE-Aufnahme mehrere Transkripte vorliegen können, haben wir ihn im Schema als iterativ gekennzeichnet.

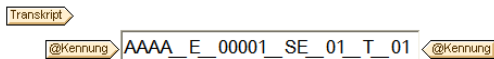


Abb. 38, Transkript - Kennung

An erster Stelle der Transkriptdaten wird eine Kennung vergeben, die die jeweilige Sprecher-eigniskennung, das Kürzel T (für „Transkript“) und eine zweistellige laufende Nummer umfasst. Ein Beispiel finden Sie in Abb. 38. Die Kennung eines zweiten Transkripts lautet im o.g. Bei-spiel: AAAA_E_00001_SE_01_T_02.

5.1.6.5.1. Basisdaten

Die Felder „Sonstige_Bezeichnung“ und „Titel“ in den Transkript-Basisdaten stehen für even-tuell im Projekt vergebene Kurzbezeichnungen und Transkripttitel.

Transkripte können typisiert werden, wobei die Extension sowie Art und Anzahl der Annotati-onen relevant werden können. Für die Bezeichnung von Transkripttypen wird das Kürzel TT in Verbindung mit einer zweistelligen Nummer verwendet.

Transkripte können Daten enthalten, die nach dem Willen der Urheber und aus datenschutz-rechtlichen Gründen Außenstehenden nicht kenntlich werden dürfen, wie z.B. persönliche Sprecherdaten, über die man im Feld „Schutzbedürftige_Daten“ informieren kann.

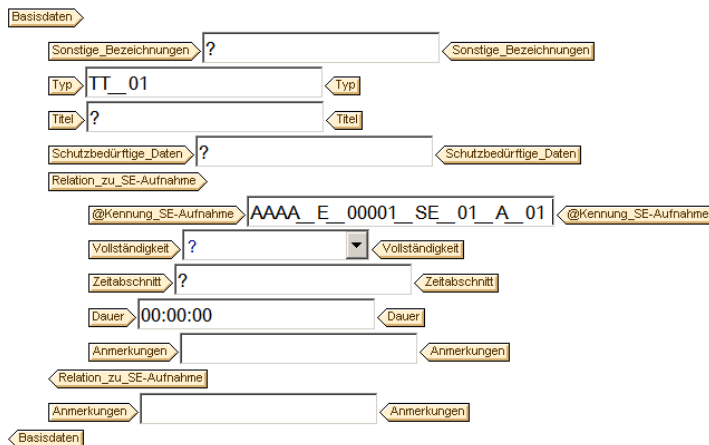


Abb. 39, Transkript - Basisdaten

Im Abschnitt „Relation_zu_SE-Aufnahme“ wird zunächst die Kennung der transkribierten SE-Aufnahme abgefragt. Im Feld „Vollständigkeit“ soll notiert werden, ob es sich um ein vollständi-ges oder ein Teiltranskript handelt. „Zeitabschnitt“ meint den Zeitabschnitt des transkribierten Teils der SE-Aufnahme. Wenn sich ein Transkript auf unterschiedliche Teile einer SE-Aufnahme bezieht (auf Auslassungen wird dann i.d.R. im Transkripttext hingewiesen), sind hier mehrere Werte zu erwarten. Für Angaben über die Dauer einer transkribierten SE-Aufnahme bzw. eines Aufnahmeausschnittes wurde ein Zeitfeld mit dem Format hh:mm:ss vorbereitet. Bei vollständig transkribierten Aufnahmen sollte die Angabe zur Dauer der SE-Aufnahme über-nommen werden.

5.1.6.5.2. Annotation

Wir verwenden die Bezeichnung „Annotation“ für inhaltlich und formal charakterisierte Ebenen eines Transkripts, wie z.B. Aufzeichnungen des Wortlauts in orthographischer, literarischer oder phonetischer Umschrift, syntaktische Angaben, Notationen suprasegmentaler oder nonverbaler

Phänomene, Übersetzung des Wortlautes etc. [13] Da u.U. mehrere Annotationen pro Transkript zu dokumentieren sind, wurde der Komplex im Schema als iterativ gekennzeichnet. [14]



Abb. 40, Annotation - Typ

Zur einfachen Benennung unterschiedlicher Annotationen haben wir eine Typenbezeichnung eingeführt. Diese Bezeichnung besteht aus dem Kürzel „ANT“ (für „Annotation_Typ“) und einer zweistelligen laufenden Nummer. Ein Beispiel finden Sie in Abb. 40.

5.1.6.5.2.1. Basisdaten

An erster Stelle der Basisdaten kann eine im Transkript enthaltene Bezeichnung für die jeweilige Annotation eingetragen werden. Im Feld „Spezifikation“ wird eine Charakterisierung der Annotation erwartet. Hier sollte man über den Gegenstand (z.B. „Wortlaut“), die Umschrift (z.B. „Literarisch“) und die Reichweite (z.B. „Ohne Interviewerbeiträge“, „Nur für Sprecher XY“) informieren.



Abb. 41, Annotation - Basisdaten

Auf die Konventionen für die jeweilige Annotation kann man im gleichnamigen Feld hinweisen. Beispiele für solche Hinweise wären: „Projektspezifisch“, „DIDA, Version vom Januar 2001“, „GAT“ etc. Über eine URL kann man ggf. direkt auf diese Konventionen zugreifen. Unter „Zeicheninventar“ ist das Inventar an Schriftzeichen zu verstehen, das bei der Wiedergabe des Wortlautes verwendet wurde. Das sind i.d.R. standardisierte Inventare wie z.B. der IPA-Zeichensatz oder ein spezifisches Alphabet.

5.1.6.5.2.2. Erstellung

Die Erstellung einer Annotation umfasst nach unserem Verständnis neben der Ersterstellung auch Ergänzungen und Korrekturen. Da man eine Annotation mehrmals überarbeiten kann, wurde die Komponente „Erstellung“ im Schema als iterativ gekennzeichnet.

Um verschiedene Erstellungsprozesse einfach benennen zu können, wurde auch an dieser Stelle eine Typenbezeichnung eingefügt. Sie setzt sich zusammen aus dem Kürzel ERT (für „Erstellung_Typ“) und einer zweistelligen Nummer.

Das Feld „Spezifikation“ wurde für Informationen über die Art der Erstellung (z.B. „Ersterfassung“, „1. Korrektur“, „Endkorrektur“, „Überarbeitung für Publikation xy“) und mögliche besondere Umstände (z.B. „halbautomatisch“) bereitgestellt. Die Werte dieses Feldes sind grundlegend für die Typisierung. Für den Namen eines Erstellungsprojekts wurde das Feld „Projekt“ eingefügt. Die Abschnitte „Bearbeiter“ und „Datum“ müssen vermutlich nicht erläutert werden. Im Feld

„Instrumente“ kann man Angaben über den genutzten Editor und evtl. weitere Hinweise auf die Systemumgebung notieren.

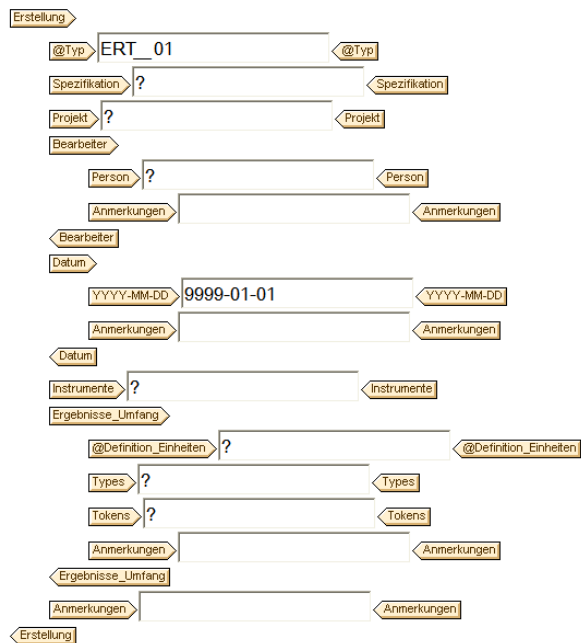


Abb. 42, Transkript - Annotation - Erstellung

Informationen über den Umfang der Ergebnisse einer Erstellung sollten eine Definition der gezählten Einheiten, Angaben über die Anzahl unterschiedlicher Einheiten (Types) und die Anzahl aller gezählten Einheiten (Tokens) umfassen.

5.1.6.5.2.3. Alignment

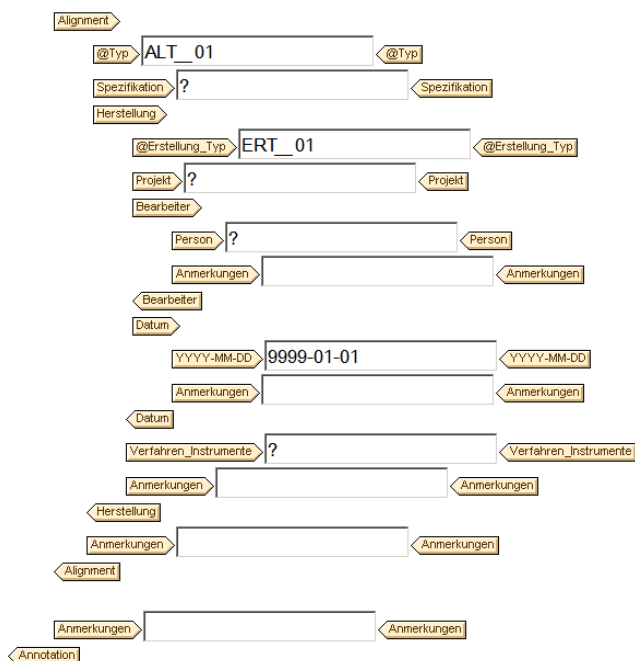


Abb. 43, Transkript - Annotation - Alignment

Wir verwenden die Bezeichnung „Alignment“ in der Dokumentation für die Text-Ton-Synchronisation, also die Koppelung von Aufnahmen und Transkripten auf Phon-, Phonem-, Wort- oder Phrasenbasis, wobei Transkriptsegmenten Zeitmarken zugeordnet werden. Die Komponente ist fakultativ und iterativ.

Alignmentprozessen werden Typenbezeichnungen zugeordnet, die aus dem Kürzel ALT (für „Alignment_Typ“) und einer zweistelligen Nummer bestehen. Grundlage für die Typisierung sind die Angaben im Feld „Spezifikation“, wo über die alignierten Segmente (z.B. „Phonweise“, „Wortweise“) informiert wird.

Die Komponente „Herstellung“ ist iterativ. Zunächst wird der Typ der Erstellung verzeichnet, dessen Ergebnisse aligniert wurden. Im Feld „Projekt“ wird der Namen des Projekts, in dem das Alignment vorgenommen wurde, erwartet. Die Abschnitte „Bearbeiter“ und „Datum“ müssen vermutlich nicht erläutert werden. Im Feld „Verfahren_Instrumente“ kann man Angaben darüber machen, ob manuell oder automatisch aligniert wurde, auf die genutzte Software hinweisen und ggf. weitere Informationen über die Systemumgebung erfassen.

5.1.6.5.3. Technische Fassungen

Transkripte können in verschiedenen technischen Fassungen vorliegen. Um diese dokumentieren zu können, haben wir auch an dieser Stelle die Komponenten „Analoge_Fassung“ und „Digitale_Fassung“ eingerichtet. Beide Komponenten sind fakultativ und iterativ, können aber nicht beide bei der Erstellung eines korpuspezifischen Schemas übergangen werden.

Zunächst wird eine Kennung für die zu dokumentierende technische Fassung generiert. Diese Kennung besteht aus der Kennung des Transkripts, dem Kürzel AF (für „Analoge Fassung“) bzw. DF (für „Digitale Fassung“) und einer zweistelligen laufenden Nummer. Beispiele finden Sie in den Abb. 44 und 45.

Es folgt ein Feld für Informationen über den jeweiligen Typ der technischen Fassungen. Für die Typisierung analoger Fassungen sind die Werte für Datenschutz und Inhalt relevant, bei der Typisierung digitaler Fassungen auch technischen Daten (wie z.B. das Dateiformat).

The image shows a form for 'Analoge_Fassung' with the following fields and labels:

- @Kennung**: AAAA_E_00001_SE_01_T_01_AF_01
- Typ**: AFT_01
- Datum**: 9999-01-01
- Anmerkungen**: (empty)
- Inhalt**: ANT_01, ERT_01
- Datenschutz**: ?
- Seitenzahl**: ?
- Datenträger**: (empty)
- @Inventarnummer**: ?
- Sonstige_Bezeichnungungen**: ?
- Anmerkungen**: (empty)
- Anmerkungen**: (empty)

Abb. 44, Transkript - Analoge Fassung

digitale_Fassung

@Kennung AAAA_E_00001_SE_01_T_01_DF_01 @Kennung

Basisdaten

Typ DFT_01 Typ

Datum

YYYY-MM-DD 9999-01-01 YYYY-MM-DD

Anmerkungen Anmerkungen

Datum

Digitalisierungssoftware ? Digitalisierungssoftware

Inhalt ANT_01, ERT_01, ALT_01 Inhalt

Datenschutz ? Datenschutz

Seitenzahl ? Seitenzahl

Datenträger

@Inventarnummer ? @Inventarnummer

Sonstige_Bezeichnung ? Sonstige_Bezeichnung

Typ ? Typ

Anmerkungen Anmerkungen

Datenträger

Abb. 45, Transkript - Digitale Fassung (1)

Elektronischer_Speicherort ? Elektronischer_Speicherort

Anmerkungen Anmerkungen

Basisdaten

Technische_Daten

Dateiname ? Dateiname

Format ? Format

Character_Encoding ? Character_Encoding

Größe 0 Größe

Anmerkungen Anmerkungen

Technische_Daten

Anmerkungen Anmerkungen

Digitale_Fassung

Abb. 46, Transkript - Digitale Fassung (2)

Im iterativen Feld „Inhalt“ sollte man die Annotationstypen sowie die Typen der Erstellungs- und Alignmentprozesse notieren, deren Ergebnisse in der jeweiligen Fassung gespeichert sind. Im Anschluss daran wird eine Information über Maßnahmen zum Datenschutz (wie z.B. die Maskierung von Personennamen) erwartet. Für den Fall, dass seitenformatierte Transkripte vorliegen und der Seitenumfang zu dokumentieren ist, haben wir das Feld „Seitenzahl“ eingefügt.

Die anderen Komponenten dieser Module stimmen mit den entsprechenden Komponenten für Zusatzmaterial überein. Erläuterungen dazu finden Sie in Abschnitt 5.1.5.2.

5.1.6.5.4. Archivierung und Distribution

Archivierung

@Stelle ? @Stelle

Dauer ? Dauer

Nutzungsrechte ? Nutzungsrechte

Sperrvermerk ? Sperrvermerk

Anmerkungen Anmerkungen

Archivierung

Abb. 47, Transkript - Archivierung

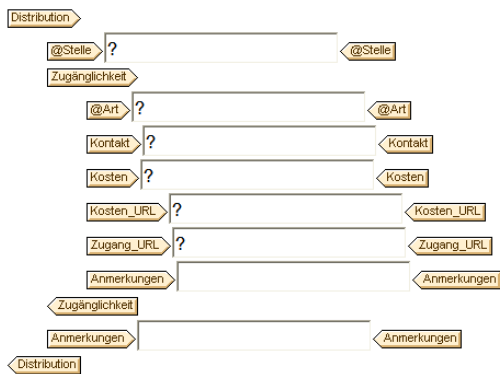


Abb. 48, Transkript - Distribution

Die Module „Archivierung“ und „Distribution“ sind in der Dokumentationsstruktur für alle Korpusbestandteile enthalten ist. Die Erläuterungen in Abschnitt 5.1.4.4. gelten auch für Transkripte.

5.1.6.6. Zusatzmaterial

Unter „Zusatzmaterial“ auf der Sprechereignisebene verstehen wir solche Dokumente, die zusätzlich zu Aufnahmen und Transkripten vorhanden sein können. Für ein Sprechereignis spezifische Unterlagen, wie z.B. schriftliche Vorbereitungen eines Interviews oder längere Aufzeichnungen über den Verlauf eines Sprechereignisses, sollten als Zusatzmaterial dokumentiert werden. Der gesamte Komplex „Zusatzmaterial“ ist fakultativ. Da zu einem Sprechereignis mehrere Dokumente vorliegen können, wurde er im Schema als iterativ gekennzeichnet.

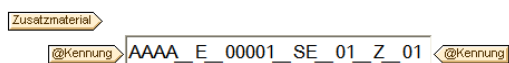


Abb. 49, Sprechereignis - Zusatzmaterial -Kennung

Die hier geforderte Kennung setzt sich zusammen aus der Sprechereigniskennung, dem Kennbuchstaben Z (für „Zusatzmaterial“) und einer zweistelligen laufenden Nummer. Ein Beispiel finden Sie in Abb. 49. Die Kennung eines zweiten sprechereignisspezifischen Dokuments müsste in unserem Beispiel AAAA_E_00001_SE_01_Z_02 lauten.

5.1.6.6.1. Basisdaten

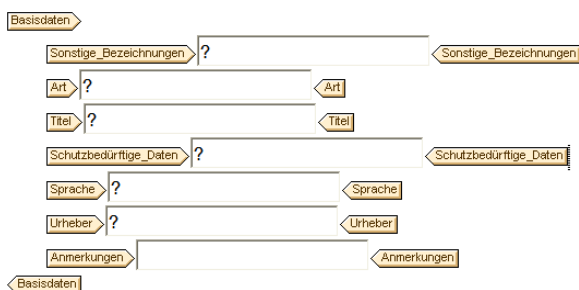


Abb. 50, Sprechereignis - Zusatzmaterial - Basisdaten

Das Modul „Basisdaten“ ist für Zusatzmaterial auf der Ereignisebene und der Sprechereignisebene gleich strukturiert. Eine Beschreibung der Struktur finden Sie in Abschnitt 5.1.5.1.

5.1.6.6.2. Technische Fassungen

Auch die Strukturen der Dokumentation technischer Fassungen von Zusatzmaterial auf Ereignis- und auf Sprechereignisebene stimmen überein. In beiden Fällen haben wir die Komponenten „Analoge_Fassung“ und „Digitale_Fassung“ im Schema als fakultativ und iterativ gekennzeichnet. Bei der Erstellung korpuspezifischer Schemata muss wenigstens eine Komponente gewählt werden.

Die Kennungen der technischen Fassungen sind allerdings ebenenspezifisch. An dieser Stelle besteht die Kennung aus der Kennung des Zusatzmaterials auf Sprechereignisebene, dem Kürzel AF (für „Analoge_Fassung“) bzw. DF (für „Digitale_Fassung“) und einer zweistelligen Nummer. Beispiele sind in den Abb. 51 und 52 enthalten. Erläuterungen der anderen Felder in diesen Modulen finden Sie in Abschnitt 5.1.5.2.

The screenshot shows a data entry form for 'Analoge_Fassung'. The form is structured as follows:

- Header:** 'Analoge_Fassung' (left), '@@Kennung' (left), 'AAAA_E_00001_SE_01_Z_01_AF_01' (center), '@@Kennung' (right).
- Typ:** 'AFT_01'.
- Datum:** '9999-01-01'.
- Anmerkungen:** Two empty text boxes.
- Datenschutz:** '?'.
- Datenträger:** '@@Inventarnummer' (left), '?' (center), '@@Inventarnummer' (right).
- Sonstige_Bezeichnungungen:** '?'.
- Typ:** '?'.
- Anmerkungen:** One empty text box.
- Datenträger:** One empty text box.
- Anmerkungen:** One empty text box.
- Footer:** 'Analoge_Fassung' (left).

Abb. 51, Sprechereignis - Zusatzmaterial - Analoge Fassung

The screenshot shows a data entry form for 'Digitale_Fassung'. The form is structured as follows:

- Header:** 'Digitale_Fassung' (left), '@@Kennung' (left), 'AAAA_E_00001_SE_01_Z_01_DF_01' (center), '@@Kennung' (right).
- Basisdaten:** 'DFT_01'.
- Datum:** '9999-01-01'.
- Anmerkungen:** Two empty text boxes.
- Datum:** 'Digitalisierungssoftware' (left), '?' (center), 'Digitalisierungssoftware' (right).
- Datenschutz:** '?'.
- Datenträger:** '@@Inventarnummer' (left), '?' (center), '@@Inventarnummer' (right).
- Sonstige_Bezeichnungungen:** '?'.
- Typ:** '?'.
- Anmerkungen:** One empty text box.
- Datenträger:** One empty text box.

Abb. 52, Sprechereignis - Zusatzmaterial - Digitale Fassung (1)

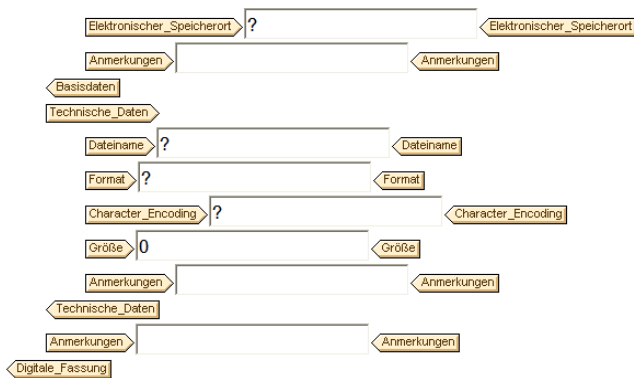


Abb. 53, Sprechereignis - Zusatzmaterial - Digitale Fassung (2)

5.1.6.6.3. Archivierung und Distribution

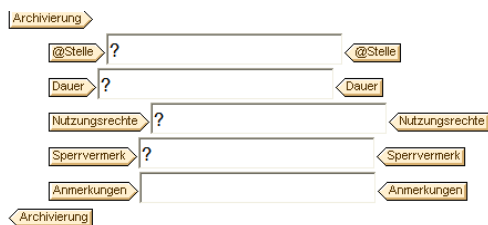


Abb. 54, Sprechereignis - Zusatzmaterial - Archivierung

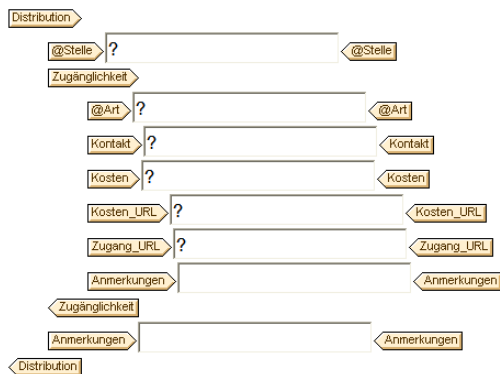


Abb. 55, Sprechereignis - Zusatzmaterial - Distribution

„Archivierung“ und „Distribution“ sind Bausteine, die in der Dokumentation aller Korpusbestandteile gebraucht werden. Die Erläuterungen in Abschnitt 5.1.4.4. gelten auch für Zusatzmaterial auf Sprechereignisebene.

5.1.7. Dokumentationsgeschichte

Informationen über Arbeitsstand und Bearbeiter der Dokumentation werden bei der manuellen Dateneingabe automatisch in einer (Oracle-)Datenbank gespeichert, sollten nach unserer Vorstellung jedoch auch in den Dokumenten sichtbar sein. Daher haben wir am Ende des Ereignisschemas den Komplex „Dokumentationsgeschichte“ eingebaut. Die Komponente „Update“ und beide Teile der Korrekturkomponente wurden im Schema als iterativ gekennzeichnet.

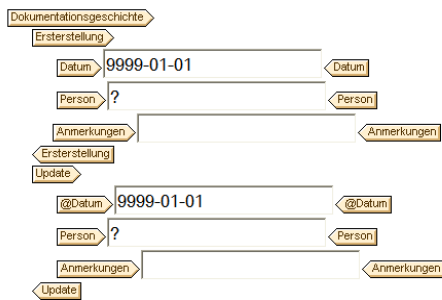


Abb. 56, Dokumentationsgeschichte (1)

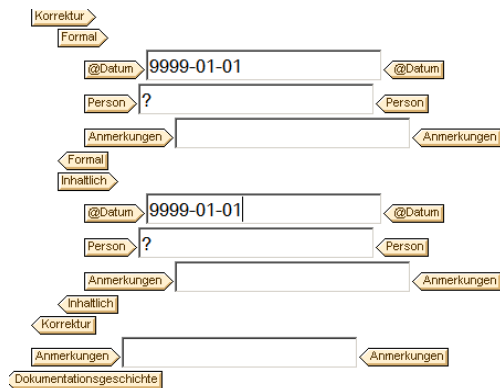


Abb. 57, Dokumentationsgeschichte (2)

6. Generisches Schema für die Dokumentation allgemeiner Sprecherdaten

Neben einem generischen Schema für die Dokumentation von Korpusbestandteilen auf der Ereignisebene (vgl. Abschnitt 5.) gibt es ein Schema für die Dokumentation ereignis- und sprechereignisübergreifender Sprecherdaten. Auch für dessen Darstellung haben wir ein projektneutrales Erfassungsformular erstellt, an dem wir uns im Folgenden orientieren. Einen Überblick über die Gesamtstruktur finden Sie in Anhang 2.

Die Kategorie „Sprecher“ dient als Startknoten eines (XML-)Schemas, das folgende Informationen vorsieht: Angaben über den jeweiligen Sprecher (Basisdaten, Ortsdaten, Sprachdaten), Beziehungen dieses Sprechers zu anderen Sprechern, sonstige Bezugspersonen des Sprechers, Vereinbarungen über Datenschutz und Nutzungsrechte, Zusatzmaterial auf Sprecher-ebene sowie eine Dokumentationsgeschichte.

Auch dieses Schema enthält obligatorische und fakultative Komponenten. Obligatorische Komponenten sind in allen projektspezifischen Subschemata zu berücksichtigen, fakultative Komponenten stehen zur Wahl und müssen in den Subschemata nicht verwendet werden. Wenn man sie verwendet, sind alle Kennungsfelder und die mit ? gekennzeichneten Felder zu bearbeiten.

Eingaben für fehlende Daten in diesen Feldern sind standardisiert: „Nicht dokumentiert“ bedeutet: Es kann ein Datum geben, das bei der Datenerfassung nicht bekannt ist. Ein Beispiel dafür wäre: „Ethnische Zugehörigkeit: Nicht dokumentiert“ - zu lesen als: „Die ethnische Zugehörigkeit des Sprechers ist nicht bekannt.“ „Nicht vorhanden“ bedeutet: Es gibt kein Datum. Ein Beispiel dafür wäre: „Aktuell ausgeübter Beruf: Nicht vorhanden“ - zu lesen als: „Der Sprecher ist nicht berufstätig.“

Das an vielen Stellen vorgesehene Feld „Anmerkungen“ ist für Anmerkungen zu Angaben in anderen Feldern und für nicht kategorisierte Angaben vorgesehen. Das Feld kann leer bleiben.

Einzelne Komponenten des Schemas wurden als iterativ gekennzeichnet, d.h. dass sie bei der korpuspezifischen Datenerfassung vervielfältigt werden können.

6.1. Sprecher

IDS-Schema für ereignisübergreifende Sprecherdaten

The screenshot shows a form with a label 'Sprecher' and a field for '@Kennung' containing the value 'AAAA_S_00001'. The field is flanked by arrows pointing to the label and the value.

Abb. 58, Sprecher - Kennung

An erster Stelle der Sprecherdaten steht eine Sprecher-Kennung, die eine vierstellige Korpuskennung, den Kennbuchstaben S (für „Sprecher“) und eine fünfstelligen laufende Nummer umfasst. Ein Beispiel finden Sie in Abb. 58. Diese Kennung ist auch im Sprecherblock des Ereignisschemas (vgl. 5.1.6.3.) zu verzeichnen. Die Kennung eines zweiten Sprechers müsste im o.g. Beispiel AAAA_S_00002 lauten.

6.1.1. Basisdaten

The screenshot shows a form for 'Basisdaten' with several fields: 'Sonstige_Bezeichnungungen' (optional), 'Name' (value: 'Anonym'), 'Früherer_Name' (value: 'Anonym'), 'Pseudonym' (optional), 'Geschlecht' (optional), 'Geburtsdatum' (value: '9999-01-01'), 'Anmerkungen', 'Aufällige_Merkmale' (optional), 'Bildungsabschluss' (optional), and 'Berufe' (optional).

Abb. 59, Sprecher - Basisdaten (1)

The screenshot shows a form for 'Basisdaten' with several fields: 'Ethnische_Zugehörigkeit' (optional), 'Gruppenzugehörigkeit' (optional), 'Staatsangehörigkeit' (optional), 'Weitere_biographische_Daten' (optional), 'Zuschreibungen' (optional), 'Sigle_in_Transkripten' (optional), and 'Anmerkungen'.

Abb. 60, Sprecher - Basisdaten (2)

An erster Stelle der Sprecher-Basisdaten steht das Feld „Sonstige_Bezeichnungungen“. Damit sind projektinterne Kurzbezeichnungen des Sprechers gemeint. Die Felder „Name“ und „Früherer_Name“ wurden mit dem Wert „Anonym“ vorbelegt und werden nur dann anders genutzt, wenn Sprechernamen Außenstehenden kenntlich werden dürfen. Für den Fall, dass maskiert wurde, haben wir das Feld „Pseudonym“ eingerichtet.

Im Anschluss daran wird das Geschlecht erfasst. Es folgt das Geburtsdatum, wobei das zugehörige Feld „Anmerkungen“ die Möglichkeit bietet, auf ungenaue Angaben hinzuweisen. Unter der Bezeichnung „Auffällige Merkmale“ kann man z.B. über körperliche Behinderung oder auffällige Kleidung informieren. Es folgen Angaben über Bildungsabschluss und Berufe.

In das Feld „Ethnische_Zugehörigkeit“ sollten Selbsteinschätzungen der Sprecher eingetragen werden. Im Feld „Gruppenzugehörigkeit“ kann man Informationen über die Zugehörigkeit eines Sprechers zu sozialen Gruppen und seine Positionen in diesen Gruppen, wie z.B. „Mitglied des Gemeinderats“ oder „Vorsitzende des örtlichen Tierschutzvereins“, ablegen. Für das iterative Feld „Staatsangehörigkeit“ gibt es eine ISO-Länderliste.

Weitere biographische Daten, für die keine Kategorien bereitgestellt wurden, können im gleichnamigen Feld notiert werden. Mit „Zuschreibungen“ sind Attribute gemeint, die sich der Sprecher selbst zuschreibt und / oder die dem Sprecher von anderen zugeschrieben werden, wie z.B. „Manager des Jahres“ oder „Bürgerschreck“.

Das Feld „Sigle_in_Transkripten“ muss vermutlich nicht erläutert werden. Da in älteren Korpora z.T. mehrere Siglen pro Sprecher verwendet wurden, haben wir das Feld auch im Sprecherblock des Ereignisschemas (vgl. 5.1.6.3.) bereitgestellt. Wir möchten an dieser Stelle noch einmal darauf hinweisen, dass für einen Sprecher nicht mehrere Siglen vergeben werden sollten.

6.1.2. Sprecher - Ortsdaten

Der Komplex „Ortsdaten“ ist fakultativ. Damit Informationen über mehrere sprachlich relevante Orte erfasst werden können, wurde er im Schema als iterativ gekennzeichnet.

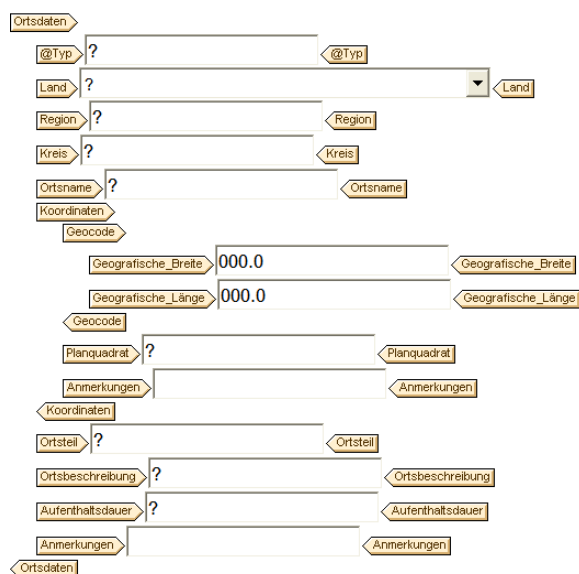


Abb. 61, Sprecher - Ortsdaten

An erster Stelle der Ortsdaten wird der Ortstyp abgefragt. Für dieses Feld sind Werte wie „Geburtsort“, „Wohnort“, „Arbeitsort“ relevant. Für das Feld „Land“ gibt es eine ISO-Liste. Das Feld „Region“ ist für die amtliche Bezeichnung eines Bundeslandes, eines Kantons oder einer Provinz vorgesehen. In die Felder „Kreis“ und „Ortsname“ sollen ebenfalls amtliche Bezeichnungen eingetragen werden.

Die Komponente „Koordinaten“ ist fakultativ. Wenn sie von einem Projekt gewählt wird, ist entweder der Geocode oder das Planquadrat (Kategorie im DSAv-Katalog [9]) des Ortes zu verzeichnen. Der „Geocode“ umfasst die Felder „Geographische_Breite“ und „Geographi-

sche_Länge“. Die Felder „Ortsteil“ und „Aufenthaltsdauer“ müssen vermutlich nicht erläutert werden. Weitere Informationen über den sprachlich relevanten Ort kann man im Feld „Ortsbeschreibung“ erfassen.

6.1.3. Sprecher - Sprachdaten

Der fakultative Bereich „Sprachdaten“ umfasst die drei fakultativen Komplexe „Sprachkenntnisse“, „Sprachproduktion“ und „Sprachgebrauch“.

6.1.3.1. Sprachkenntnisse

Abb. 62, Sprecher - Sprachdaten - Sprachkenntnisse

Um mehrsprachigen Personen gerecht werden zu können, wurde der Komplex „Sprachkenntnisse“ als iterativ gekennzeichnet.

Nach dem Sprachnamen (z.B. „Deutsch“, „Englisch“, „Türkisch“) ist der Sprachstatus anzugeben (z.B. „Muttersprache“, „Erstsprache“, „Zweitsprache“, „1. Fremdsprache“). Die Komponente „Kenntnisgrade“ ist fakultativ, alle darin enthaltenen Felder außer den Anmerkungen werden verbindlich, wenn sie gewählt wird. In Abb. 63 sehen Sie die für die Einschätzung der Kenntnisgrade vorgegebenen Werte.

Abb. 63, Werte für die Einschätzung von Kenntnisgraden

Im Anschluss an die Kenntnisgrade kann man Informationen über sprachliche Besonderheiten, wie z.B. dialektale Aspekte, erfassen.

6.1.3.2. Sprachproduktion

Die in diesem Abschnitt vorgestellte Systematik rekuriert u.a. auf einen Sprecherfragebogen für die Erhebung „Deutsch heute“ (DH) [15]. Wir haben allerdings die in die Schemaentwicklung

eingespeisten DH-Kategorien daraufhin geprüft, ob sie auch für andere Anwendungen brauchbar sind, und im Hinblick darauf überarbeitet.

Der Abschnitt Sprachproduktion ist fakultativ. Er besteht aus dem ebenfalls fakultativen Komplex „Einflussfaktoren“ und dem Element „Sprachliche_Besonderheiten“.

Der Komplex „Einflussfaktoren“ umfasst die fünf fakultativen Komponenten „Körpermaße“, „Beeinträchtigung“, „Drogen_Medikamente“, „Gebrauch_von_Hilfsmitteln“ und „Unterricht_Korrekturen“.

In Abb. 64 sehen Sie zunächst die Struktur des Moduls „Körpermaße“. Die Körpergröße wird in cm, das Körpergewicht in kg angegeben.

Das Modul „Beeinträchtigung“ wurde im Schema als iterativ gekennzeichnet. Hier können Informationen über körperliche und psychische Beeinträchtigungen erfasst werden. Wir denken dabei an Werte wie z.B. „Zahnlücke im vorderen Unterkiefer“, „Schwerhörigkeit“, „Asthma“ oder „Depression“. Im Feld „Häufigkeit_Umfang“ kann notiert werden, wie häufig bzw. in welchem Umfang eine Beeinträchtigung auftritt. Eine Information darüber, wie lange die Beeinträchtigung schon vorhanden ist, kann man im Feld „Dauer“ erfassen.

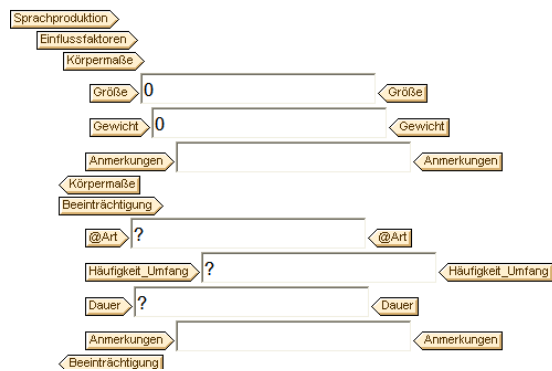


Abb. 64, Sprecher - Sprachdaten - Sprachproduktion (1)

Abb. 65 zeigt die Struktur der Komponente „Drogen_Medikamente“. Auch diese Komponente wurde im Schema als iterativ gekennzeichnet. Im Feld „Art“ kann z.B. „Nikotin“ eingetragen werden. Im Feld „Häufigkeit_Umfang“ kann man dann festhalten, welche und wie viele Rauchwaren am Tag konsumiert werden. Im Feld „Dauer“ wird vermerkt, wie lange der Sprecher schon raucht.

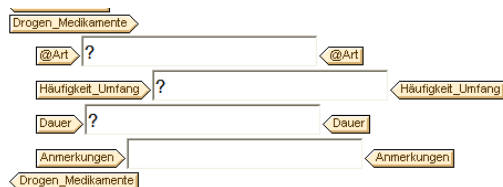


Abb. 65, Sprecher - Sprachdaten - Sprachproduktion (2)

Die Struktur der iterativen Komponente „Gebrauch_von_Hilfsmitteln“ ist in Abb. 66 dargestellt. Bei sprachlich relevanten Hilfsmitteln ist v.a. an Wörterbücher aller Art, Grammatiken, Stillehren und Kommunikationsratgeber zu denken, hier könnten aber auch andere Hilfsmittel, wie z.B. ein Hörgerät, verzeichnet werden. Für entsprechende Werte steht das Feld „Art“ bereit. Unter „Häufigkeit_Umfang“ kann vermerkt werden, wie oft ein Hilfsmittel genutzt wird. Die übrigen Felder stimmen mit den in Abb. 65 gezeigten überein.

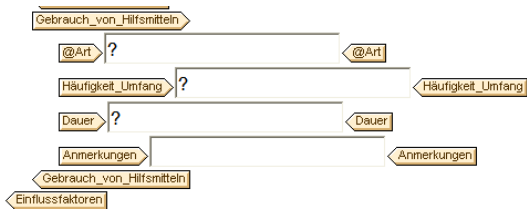


Abb. 66, Sprecher - Sprachdaten - Sprachproduktion (4)

In Abb. 67 ist die Struktur der iterativen Komponente „Unterricht_Korrekturen“ zu sehen.

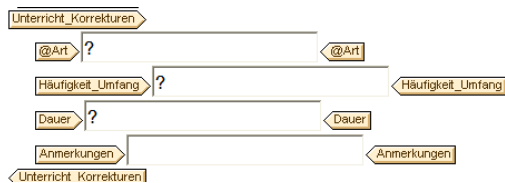


Abb. 67, Sprecher - Sprachdaten - Sprachproduktion (3)

Abb. 68 zeigt die im Erfassungsformular des Projekts „Deutsch heute“ für das Feld „Art“ im Komplex „Unterricht_Korrekturen“ bereitgestellten Werte. „Fremdkorrektur“ meint Korrektur des Sprechers durch andere Personen, „Selbstkorrektur“ heißt, dass der Sprecher seine sprachlichen Äußerungen gewöhnlich selbst korrigiert.

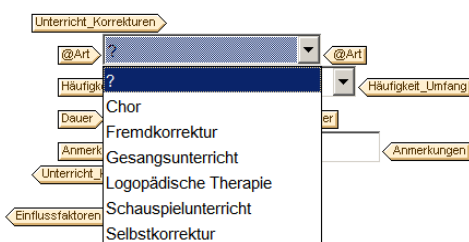


Abb. 68, Beispielwerte für Unterricht_Korrekturen

In das Feld „Sprachliche_Besonderheiten“ können produktionsrelevante Angaben wie z.B. „Stottern“, „Lispeln“, „Nuscheln“ etc. eingetragen werden.



Abb. 69, Sprecher - Sprachdaten - Sprachproduktion (5)

6.1.3.3. Sprachgebrauch

Abb. 70 zeigt die Struktur des fakultativen Abschnitts „Sprachgebrauch“. Eine wichtige Kategorie in diesem Abschnitt ist „Domäne“, worunter wir einen relevanten Kommunikationsbereich verstehen. Da i.d.R. Angaben über mehrere Kommunikationsbereiche zu erfassen sind, wurde der Abschnitt im Schema als iterativ gekennzeichnet. Als „Domäne“ gelten z.B. „Familie“, „Nachbarschaft“ und „Arbeitsplatz“.

Im Abschnitt „Sprachen“ werden zunächst die Namen der verwendeten Sprachen (z.B. „Deutsch ; Türkisch“) erwartet. Unter dem Titel „Sprachliche_Besonderheiten“ kann man Hinweise auf weitere Aspekte, wie z.B. Dialektgebrauch, verzeichnen.

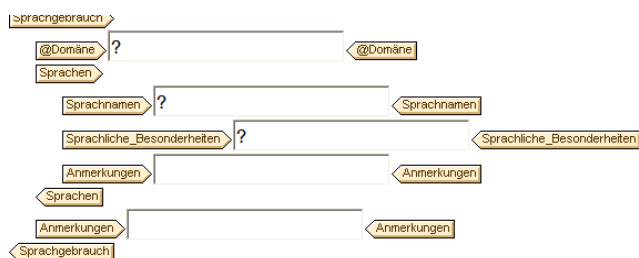


Abb. 70, Sprecher - Sprachdaten - Sprachgebrauch

6.1.4. Beziehungen zu anderen Sprechern

In Abb. 71 sehen Sie die Struktur des fakultativen Bereichs, in dem über Beziehungen zwischen Sprechern informiert werden kann. Damit mehrere Beziehungen dokumentiert werden können, wurde der Abschnitt im Schema als iterativ gekennzeichnet.

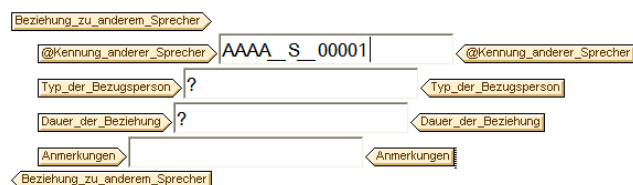


Abb. 71, Beziehung zu anderem Sprecher

An erster Stelle wird die Kennung des anderen Sprechers registriert, dann ist über den „Typ_der_Bezugsperson“ zu informieren. In diesem Feld könnte z.B. „Freund“, „Nachbar“, „Tochter“ etc. stehen. Für Angaben über die Dauer der Beziehung gibt es ein gleichlautendes Feld.

6.1.5. Sonstige Bezugspersonen

Neben dem Bereich der eigentlichen Sprecherdaten haben wir einen fakultativen Bereich für Daten über sonstige Bezugspersonen vorgesehen, der die fakultativen Teile „Bezugspersonen kompakt“ und „Einzelne Bezugsperson“ umfasst. Die Unterscheidung wurde angesichts der Notwendigkeit getroffen, verschiedenen Datenbeständen gerecht zu werden. Wir müssen zum einen kompakte Informationen über Gruppen von Bezugspersonen, zum anderen Daten für einzelne Bezugspersonen erfassen können.

6.1.5.1. Bezugspersonen kompakt

Wir gehen davon aus, dass sich auch bei kompakten Angaben über Gruppen von Bezugspersonen (z.B. „Eltern“, „Kinder“, „Freunde“, „Kollegen“) Personen- und Ortsdaten einerseits sowie Sprachdaten andererseits unterscheiden lassen. Falls der Gesamtkomplex gewählt wird, werden Personen- und Ortsdaten obligatorisch, Sprachdaten sind fakultativ, können also bei der Erstellung projektspezifischer Schemata übergangen werden.

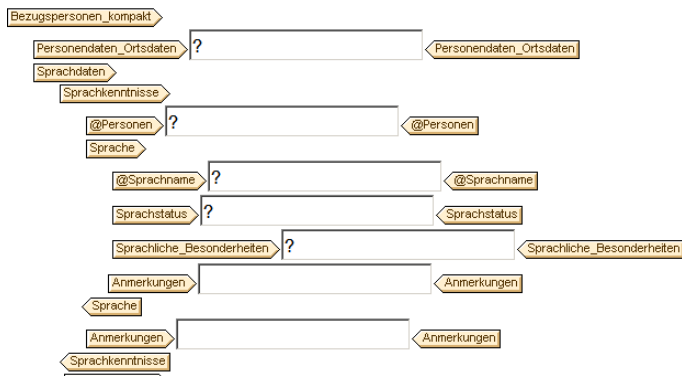


Abb. 72, Bezugspersonen kompakt (1)

An erster Stelle in Abb. 72 ist die Sammelkategorie „Personendaten_Ortsdaten“ zu sehen, die persönliche Angaben über Personengruppen und Informationen über sprachlich relevante Orte aufnehmen soll. Ein fiktives Beispiel dafür wäre: „Die Eltern sind seit 60 Jahren verheiratet und stammen beide aus Freiburg. Die drei Kinder leben und arbeiten in Frankfurt, Hamburg und München und besuchen die Herkunftsfamilie nur noch zu besonderen Anlässen.“

Sprachdaten gliedern wir hier in Sprachkenntnisse und Sprachgebrauch. Auf die Komponente „Sprachproduktion“, die für den Sprecher angelegt wurde, haben wir in diesem Bereich verzichtet, da wir davon ausgehen, dass diese Daten für Bezugspersonen nicht erhoben werden.

Die Struktur der Komponente „Sprachkenntnisse“ im Bereich „Bezugspersonen_kompakt“ ist in Abb. 72 dargestellt. Diese Komponente ist fakultativ und - für den Fall, dass Sprachkenntnisse mehrerer Personengruppen zu dokumentieren sind - iterativ. Im Unterschied zu Abschnitt 6.1.3.1., wo diese Komponente eingeführt wurde, haben wir hier auf den Teil „Kenntnisgrade“ verzichtet, da für Personengruppen vermutlich keine (einheitlichen) Kenntnisgrade zu ermitteln sind.

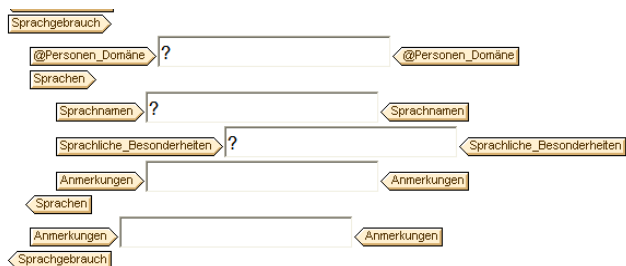


Abb. 73, Bezugspersonen kompakt (2)

An erster Stelle der Komponente „Sprachgebrauch“ steht das Feld „Personen_Domäne“, das Werte wie z.B. „Eltern untereinander“, „Eltern mit Kindern“ oder „Eltern mit Großeltern“ aufnehmen soll. Die übrigen Elemente wurden schon im Abschnitt 6.1.3.3. vorgestellt.

6.1.5.2. Einzelne Bezugsperson

Der fakultative Bereich „Einzelne_Bezugsperson“ umfasst das Feld „Typ“ (der Bezugsperson) für Werte wie z.B. „Mutter“ oder „Tante“ sowie Basisdaten, Ortsdaten und Sprachdaten. Das Feld „Typ“ ist obligatorisch, d.h. dass es berücksichtigt werden muss, wenn der Bereich „Einzelne_Bezugsperson“ gewählt wird, Basisdaten, Ortsdaten und Sprachdaten sind fakultativ, können aber nicht alle bei der Erstellung eines projektspezifischen Schemas ausgeschlossen werden.

Abb. 74, Einzelne Bezugsperson - Typ

6.1.5.2.1. Basisdaten

Besonderheiten der fakultativen Basisdaten für einzelne Bezugspersonen sind am Anfang von Abb. 75 dargestellt. Dort finden Sie die Felder „Status_der_Bezugsperson“ (für den Sprecher), wofür die Werte „Relevant“ und „Peripher“ vorgesehen sind, sowie „Dauer der Beziehung“.

Zu allen anderen Feldern in diesem Komplex gibt es Entsprechungen in den Sprecher-Basisdaten, die in Abschnitt 6.1.1. beschrieben wurden.

Abb. 75, Einzelne Bezugsperson - Basisdaten (1)

Abb. 76, Einzelne Bezugsperson - Basisdaten (2)

6.1.5.2.2. Ortsdaten

In Abb. 77 sehen Sie die Struktur der fakultativen Ortsdaten für einzelne Bezugspersonen. Sie stimmt mit der der Ortsdaten für Sprecher überein, über die Sie sich in Abschnitt 6.1.2. informieren können.

Ortsdaten

@Typ ? @Typ

Land ? Land

Region ? Region

Kreis ? Kreis

Ortsname ? Ortsname

Koordinaten

Geocode

Geografische_Breite 000.0 Geografische_Breite

Geografische_Länge 000.0 Geografische_Länge

Geocode

Planquadrat ? Planquadrat

Anmerkungen Anmerkungen

Koordinaten

Ortsteil ? Ortsteil

Ortsbeschreibung ? Ortsbeschreibung

Aufenthaltsdauer ? Aufenthaltsdauer

Anmerkungen Anmerkungen

Ortsdaten

Abb. 77, Einzelne Bezugsperson - Ortsdaten

6.1.5.2.3. Sprachdaten

Der fakultative Komplex „Sprachdaten“ im Bereich „Einzelne Bezugspersonen“ umfasst die fakultativen Komponenten „Sprachkenntnisse“ und „Sprachgebrauch“. Auf die Komponente „Sprachproduktion“, die für den Sprecher angelegt wurde, haben wir in diesem Bereich verzichtet, da wir davon ausgehen, dass diese Daten für Bezugspersonen nicht erhoben werden.

6.1.5.2.3.1. Sprachkenntnisse

Für Informationen über Sprachkenntnisse von Sprechern und einzelnen Bezugspersonen haben wir eine einheitliche Struktur gewählt. Erläuterungen dazu finden Sie in Abschnitt 6.1.3.1.

Sprachkenntnisse

@Sprachname ? @Sprachname

Sprachstatus ? Sprachstatus

Kenntnisgrade

Allgemeine_Einschätzung ? Allgemeine_Einschätzung

Hörverstehen ? Hörverstehen

Lesen ? Lesen

Schreiben ? Schreiben

Sprechen ? Sprechen

Anmerkungen Anmerkungen

Kenntnisgrade

Sprachliche_Besonderheiten ? Sprachliche_Besonderheiten

Anmerkungen Anmerkungen

Sprachkenntnisse

Abb. 78, Einzelne Bezugsperson - Sprachdaten - Sprachkenntnisse

6.1.5.2.3.2. Sprachgebrauch

Auch Informationen über den Sprachgebrauch von Sprechern und einzelnen Bezugspersonen werden in einer einheitlichen Struktur erfasst. Erläuterungen dazu gibt es in Abschnitt 6.1.3.3.

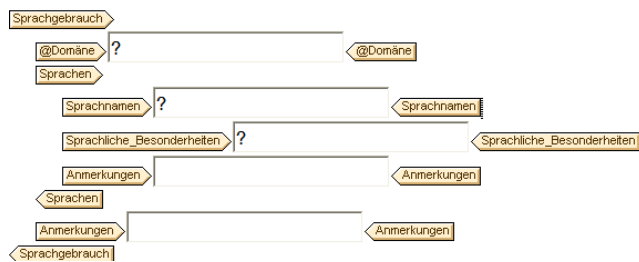


Abb. 79, Einzelne Bezugsperson - Sprachdaten - Sprachgebrauch

6.1.6. Rechteverwaltung

Im Komplex „Rechteverwaltung“ sollen rechtliche Aspekte der Datenerhebung sowie rechtsrelevante Vereinbarungen mit Sprechern und ggf. auch Bezugspersonen über Schutz und Verwendung ihrer Daten dokumentiert werden. Im Folgenden unterscheiden wir zwischen personenbezogenen Daten und Korpusbestandteilen, wobei Korpusbestandteile personenbezogene Daten enthalten können.

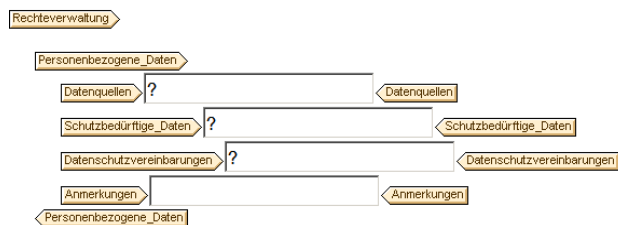


Abb. 80, Sprecher - Rechteverwaltung - Personenbezogene Daten

Zu den rechtsrelevanten Aspekten zählt die Frage, aus welchen Quellen die erhobenen personenbezogenen Daten stammen. In der Regel stammen sie von den Sprechern, es können aber auch Bezugspersonen befragt oder schriftliche Quellen ausgewertet worden sein. Im Feld „Schutzbedürftige_Daten“ kann man festhalten, ob nur besondere Arten oder alle personenbezogener Daten zu schützen sind. Mit „Datenschutzvereinbarungen“ meinen wir Vereinbarungen mit Sprechern und Bezugspersonen über den Schutz der im vorigen Feld genannten Daten. Solche Vereinbarungen können vorsehen, dass diese Daten nur im Rahmen des erhebenden Projekts verwendet und danach gelöscht werden oder auf bestimmten Wegen für bestimmte Zwecke an Dritte weitergegeben werden dürfen.

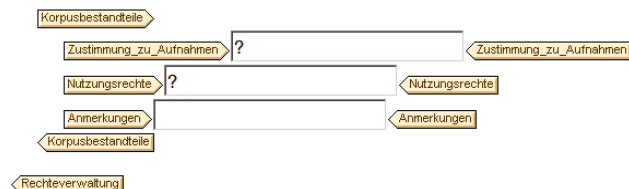


Abb. 81, Sprecher - Rechteverwaltung - Korpusbestandteile

Die Zustimmung der Sprecher zu den Aufnahmen ist eine wesentliche Voraussetzung für die Verwendung von Aufnahmen und Transkripten. Im entsprechenden Feld sollte man daher festhalten, ob und wann die Sprecher über den Zweck der Aufnahmen informiert wurden und in welcher Form - schriftlich oder (aus welchen Gründen?) mündlich - sie den Aufnahmen zugestimmt haben. Schließlich werden Nutzungsrechte an Korpusbestandteilen dokumentiert, die von der wissenschaftlichen Auswertung im Daten erhebenden Projekt bis zur Veröffentlichung im Internet reichen können.

6.1.7. Zusatzmaterial

Unter Zusatzmaterial auf Sprecherebene verstehen wir Dokumente wie z.B. Sprecherfotos oder schriftliche Vereinbarungen mit den Sprechern, die als Korpusbestandteile gelten können. Die Struktur für die Dokumentation von Zusatzmaterial auf Sprecherebene stimmt mit der für Zusatzmaterial auf Ereignis- und auf Sprechereignisebene (vgl. Abschnitte 5.1.5. und 5.1.6.6.) überein. Der Komplex ist fakultativ. Da pro Sprecher mehrere zusätzliche Dokumente vorliegen können, wurde er im Schema als iterativ gekennzeichnet.

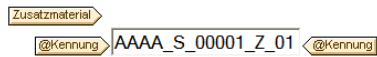


Abb. 82, Sprecher - Zusatzmaterial - Kennung

Die Kennung des Zusatzmaterials ist ebenenspezifisch und umfasst hier die Sprecherkennung, den Kennbuchstaben Z (für „Zusatzmaterial“) und eine zweistellige laufende Nummer. Ein Beispiel finden Sie in Abb. 82. Bei der Beschreibung eines zweiten sprecherspezifischen Dokuments müsste die Kennung im AAAA_S_00001_Z_02 lauten.

6.1.7.1. Basisdaten

An erster Stelle der Komponente „Basisdaten“ können Bezeichnungen eingetragen werden, die vor der Kennung vergeben wurden. Zusatzmaterialien können Daten enthalten, die nach dem Willen der Urheber und aus datenschutzrechtlichen Gründen Außenstehenden nicht kenntlich werden dürfen, wie z.B. persönliche Sprecherdaten. Für entsprechende Informationen wurde das Feld „Schutzbedürftige_Daten“ bereitgestellt. „Urheber“ steht für Autoren, Grafiker, Fotografen etc. Die übrigen Felder in diesem Komplex müssen vermutlich nicht erläutert werden.

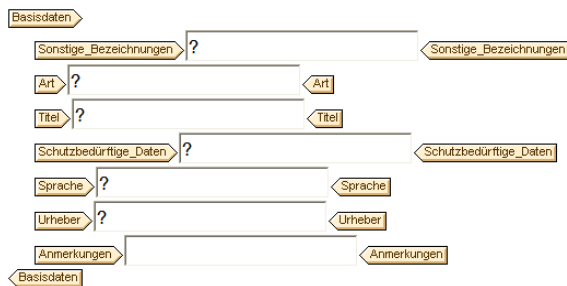


Abb. 83, Sprecher - Zusatzmaterial - Basisdaten

6.1.7.2. Technische Fassungen

Da zusätzliche Dokumente in verschiedenen technischen Fassungen vorliegen können, haben wir die fakultativen und iterativen Komponenten „Analoge_Fassung“ und „Digitale_Fassung“ vorgesehen. Wenigstens eine Komponente muss bei der Erstellung eines projektspezifischen Schemas gewählt werden.

Die Kennungen der technischen Fassungen bestehen aus der Kennung des Zusatzmaterials auf Sprecherebene, dem Kürzel AF (für „Analoge_Fassung“) bzw. DF (für „Digitale_Fassung“) und einer zweistelligen laufenden Nummer.

Im Anschluss an die Kennungen werden Typen analoger und digitaler Fassungen benannt. Grundlage für eine Typisierung analoger Fassungen ist der Datenschutz, für die Typisierung digitaler Fassungen sind außerdem noch technische Daten (z.B. das Dateiformat) relevant. Als Typenbezeichnungen dienen die Kürzel AFT (analoge Fassung) und DFT (digitale Fassung) in Verbindung mit einer zweistelligen Nummer.

Beispiele für Kennungen und Typenbezeichnungen finden Sie in den Abb. 84 und 85.

Analoge_Fassung

@Kennung AAAA_S_00001_Z_01_AF_01 @Kennung

Typ AFT_01 Typ

Datum

YYYY-MM-DD 9999-01-01 YYYY-MM-DD

Anmerkungen Anmerkungen

Datum

Datenschutz ? Datenschutz

Datenträger

@inventarnummer ? @inventarnummer

Sonstige_Bezeichnung ? Sonstige_Bezeichnung

Typ ? Typ

Anmerkungen Anmerkungen

Datenträger

Anmerkungen Anmerkungen

Analoge_Fassung

Abb. 84, Sprecher - Zusatzmaterial - Analoge Fassung

Digitale_Fassung

@Kennung AAAA_S_00001_Z_01_DF_01 @Kennung

Basisdaten

Typ DFT_01 Typ

Datum

YYYY-MM-DD 9999-01-01 YYYY-MM-DD

Anmerkungen Anmerkungen

Datum

Digitalisierungssoftware ? Digitalisierungssoftware

Datenschutz ? Datenschutz

Datenträger

@inventarnummer ? @inventarnummer

Sonstige_Bezeichnung ? Sonstige_Bezeichnung

Typ ? Typ

Anmerkungen Anmerkungen

Datenträger

Elektronischer_Speicherort ? Elektronischer_Speicherort

Anmerkungen Anmerkungen

Basisdaten

Abb. 85, Sprecher - Zusatzmaterial - Digitale Fassung (1)

Technische_Daten

Dateiname ? Dateiname

Format ? Format

Character_Encoding ? Character_Encoding

Größe 0 Größe

Anmerkungen Anmerkungen

Technische_Daten

Anmerkungen Anmerkungen

Digitale_Fassung

Abb. 86, Sprecher - Zusatzmaterial - Digitale Fassung (2)

Das Modul „Datum“ ist für Angaben über das Erstellungsdatum der technischen Fassung vorgesehen. „Datenschutz“ meint technische Maßnahmen zum Datenschutz, wie z.B. die Maskierung von Personennamen in Texten.

Da eine technische Fassung auf mehreren Datenträgern gespeichert sein kann, wurde dieser Abschnitt im Schema als iterativ gekennzeichnet. An erster Stelle dieses Abschnitts wird eine eindeutige Inventarnummer des zu dokumentierenden Datenträgers erwartet. Im Feld „Sonstige_Bezeichnung“ können weitere in einem Korpusprojekt möglicherweise generierte Ordnungskennzeichen (Name, Ordnungsnummer, etc.) erfasst werden. Im nächsten Schritt ist über den Typ des Datenträgers (z.B. Papier, Mikrofilm, Daten-CD) zu informieren.

Die Komponenten „Digitalisierungssoftware“, „Elektronischer Speicherort“ und „Technische_Daten“ sind nur für digitale Fassungen relevant. „Digitalisierungssoftware“ meint das Programm, mit dem eine analoge Fassung digitalisiert wurde. In das Feld „Elektronischer Speicherort“ kann man eine URL oder einen Pfadnamen eintragen.

Das erste Feld im Abschnitt „Technische_Daten“ ist für den Dateinamen vorgesehen, an den sich eine Information über das Dateiformat anschließen sollte. „Character_Encoding“ steht für die Zeichencodierung in einer Textdatei (z.B. ASCII oder UTF-16BE). Im Feld „Größe“ ist die Dateigröße (Anzahl von Bytes) anzugeben.

6.1.7.3. Archivierung und Distribution

In den Abb. 87 und 88 werden die Bausteine „Archivierung“ und „Distribution“ vorgestellt, die Sie möglicherweise schon aus der Beschreibung des Ereignisschemas kennen. Sie sollen Informationen über rechtliche und organisatorische Aspekte der Korpusbestandteile aufnehmen.

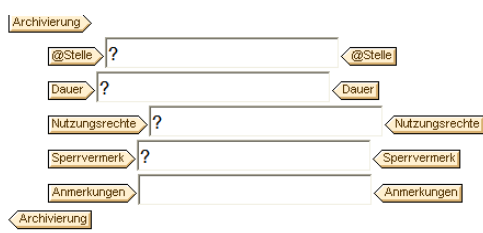


Abb. 87, Sprecher - Zusatzmaterial - Archivierung

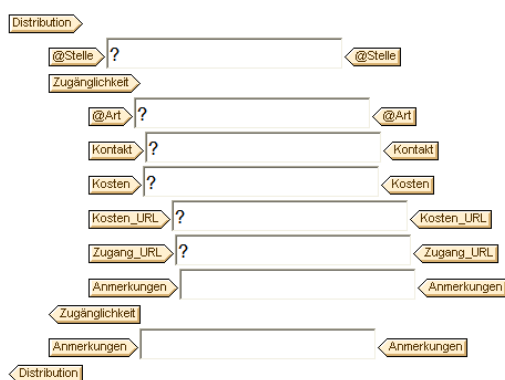


Abb. 88, Sprecher - Zusatzmaterial - Distribution

Das Modul „Archivierung“ wurde im Schema als iterativ gekennzeichnet. Zunächst soll der Name der archivierenden Stelle vermerkt werden. Es folgt ein Feld für Informationen über die vorgesehene Archivierungsdauer. Hier könnte z.B. „Bis 2018“ oder „Langfristig“ stehen. Die Nutzungsrechte der archivierenden Stelle können von der ausschließlichen Nutzung durch den Projektleiter bis hin zur Veröffentlichung eines Dokuments im Internet reichen. Sperrvermerke, wie z.B. „Bis 2010 für Externe gesperrt“, können die Nutzungsmöglichkeiten einschränken.

Das Modul „Distribution“ ist ebenfalls iterativ und umfasst neben einem Feld für den Namen der für die Distribution zuständigen Stelle die iterative Komponente „Zugänglichkeit“. In dieser

Komponente sollen folgende Angaben verzeichnet werden: Art der Zugänglichkeit, E-Mail-Kontaktadresse, Angaben über die Kosten, ggf. eine URL dieser Angaben sowie ggf. eine URL, die den direkten Zugang zum jeweiligen Korpusbestandteil ermöglicht.

6.1.8. Dokumentationsgeschichte

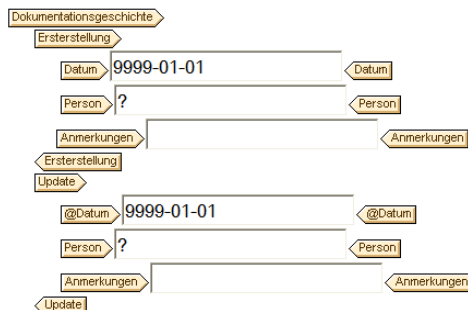


Abb. 89, Sprecher - Dokumentationsgeschichte (1)

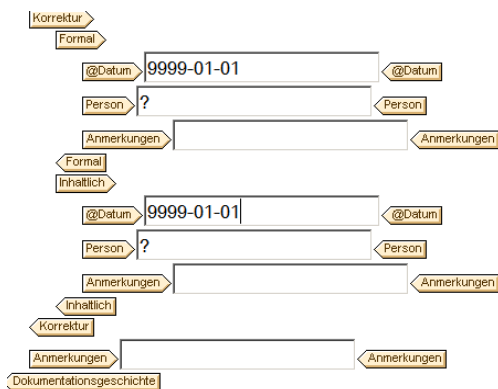


Abb. 90, Sprecher - Dokumentationsgeschichte (2)

Das Schema für die Erfassung allgemeiner Sprecherdaten enthält wie das Schema für die Dokumentation von Korpusbestandteilen auf Ereignisebene den Komplex „Dokumentationsgeschichte“. Die Komponente „Update“ und beide Teile der Korrekturkomponente wurden im Schema als iterativ gekennzeichnet.

7. Generisches Schema für die Dokumentation von Zusatzmaterial auf der Korpus ebene

Unter „Zusatzmaterial auf der Korpusebene“ verstehen wir Dokumente, die zusätzlich zu Aufnahmen, Transkripten und Zusatzmaterialien auf Ereignis-, Sprechereignis- und Sprecherebene vorhanden sein können. Das sind z.B. Transkriptionskonventionen, ein Interviewleitfaden, Wortlisten, verschiedene Varianten der Wenckersätze, ggf. auch Spezifikationen für die Validierung von Korpusdaten sowie Dokumente, die die Ergebnisse solcher Qualitätsprüfungen enthalten [16].

Die Struktur für die Dokumentation von Zusatzmaterial auf Korpusebene stimmt mit der für Zusatzmaterial auf Ereignis-, Sprechereignis- und Sprecherebene überein. Lediglich die Kennungen sind ebenenspezifisch. Wenn Sie den Komplex „Zusatzmaterial“ des Ereignisschemas oder des Schemas für allgemeine Sprecherdaten kennengelernt haben, können Sie weite Teile der folgenden Darstellung übergehen.

Das Schema enthält obligatorische und fakultative Komponenten. Obligatorische Komponenten sind in allen projektspezifischen Subschemata zu berücksichtigen, fakultative Komponenten stehen zur Wahl und müssen in den Subschemata nicht verwendet werden. Wenn man sie verwendet, sind alle Kennungsfelder und die mit ? gekennzeichneten Felder zu bearbeiten.

Eingaben für fehlende Daten in diesen Feldern sind standardisiert: „Nicht dokumentiert“ bedeutet: Es kann ein Datum geben, das bei der Datenerfassung jedoch nicht bekannt ist. Ein Beispiel dafür wäre: „Urheber: Nicht dokumentiert“ - zu lesen als: „Der Name des Urhebers ist nicht dokumentiert.“ „Nicht vorhanden“ bedeutet: Es gibt kein Datum. Ein Beispiel dafür wäre: „Schutzbedürftige Daten: Nicht vorhanden“ - zu lesen als: „Es gibt in diesem Dokument keine schutzbedürftigen Daten.“

Das an vielen Stellen vorgesehene Feld „Anmerkungen“ ist für Anmerkungen zu Angaben in anderen Feldern und für nicht kategorisierte Angaben vorgesehen. Das Feld kann leer bleiben.

Einzelne Komponenten des Schemas wurden als iterativ gekennzeichnet, d.h. dass sie bei der Datenerfassung vervielfältigt werden können.

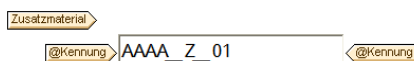


Abb. 91, Zusatzmaterial - Kennung

Zuerst wird eine Kennung für Zusatzmaterial auf Korpusebene generiert. Diese Kennung setzt sich zusammen aus der Korpuskennung, dem Kennbuchstaben Z (für „Zusatzmaterial“) und einer zweistelligen laufenden Nummer. Ein Beispiel finden Sie in Abb. 91. Bei der Beschreibung eines zweiten korpuspezifischen Dokuments lautet die Kennung im o.g. Beispiel AAAA_Z_02.

7.1. Basisdaten

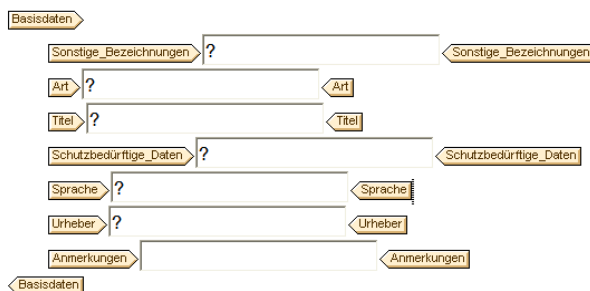


Abb. 92, Zusatzmaterial - Basisdaten

An erster Stelle der Komponente „Basisdaten“ können Bezeichnungen eingetragen werden, die vor der Kennung vergeben wurden. Zusatzmaterialien können Daten enthalten, die nach dem Willen der Urheber und aus datenschutzrechtlichen Gründen Außenstehenden nicht kenntlich werden dürfen, wie z.B. persönliche Sprecherdaten. Für entsprechende Informationen wurde das Feld „Schutzbedürftige_Daten“ bereitgestellt. „Urheber“ steht für Autoren, Grafiker, Fotografen etc. Die übrigen Felder in diesem Komplex müssen vermutlich nicht erläutert werden.

7.2. Technische Fassungen

Da zusätzliche Dokumente in verschiedenen technischen Fassungen vorliegen können, werden die fakultativen und iterativen Komponenten „Analoge_Fassung“ und „Digitale_Fassung“ bereitgestellt. Wenigstens eine Komponente muss bei der Erstellung eines projektspezifischen Schemas gewählt werden.

Analoge_Fassung

@Kennung AAAA_Z_01_AF_01 @Kennung

Typ AFT_01 Typ

Datum

YYYY-MM-DD 9999-01-01 YYYY-MM-DD

Anmerkungen

Datum

Datenschutz ? Datenschutz

Datenträger

@inventarnummer ? @inventarnummer

Sonstige_Bezeichnung ? Sonstige_Bezeichnung

Typ ? Typ

Datenträger

Anmerkungen

Analoge_Fassung

Abb. 93, Zusatzmaterial - Analoge Fassung

Digitale_Fassung

@Kennung AAAA_Z_01_DF_01 @Kennung

Basisdaten

Typ DFT_01 Typ

Datum

YYYY-MM-DD 9999-01-01 YYYY-MM-DD

Anmerkungen

Datum

Digitalisierungssoftware ? Digitalisierungssoftware

Datenschutz ? Datenschutz

Datenträger

@inventarnummer ? @inventarnummer

Sonstige_Bezeichnung ? Sonstige_Bezeichnung

Typ ? Typ

Anmerkungen

Datenträger

Elektronischer_Speicherort ? Elektronischer_Speicherort

Anmerkungen

Basisdaten

Abb. 94, Zusatzmaterial - Digitale Fassung (1)

Technische_Daten

Dateiname ? Dateiname

Format ? Format

Character_Encoding ? Character_Encoding

Größe 0 Größe

Anmerkungen

Technische_Daten

Anmerkungen

Digitale_Fassung

Abb. 95, Zusatzmaterial - Digitale Fassung (2)

In den Abbildungen 93 und 94 sehen Sie Kennungen für analoge und digitale technische Fassungen von Zusatzmaterial auf Korpusebene. Diese Kennungen enthalten die Kennung des Zusatzmaterials auf Korpusebene, das Kürzel AF (für „Analoge_Fassung“) bzw. DF (für „Digitale_Fassung“) und eine zweistellige laufende Nummer.

Im Anschluss an die Kennungen werden Typen analoger und digitaler Fassungen benannt. Grundlage für eine Typisierung analoger Fassungen ist Datenschutz, für die Typisierung digitaler Fassungen sind außerdem noch technische Daten (z.B. das Dateiformat) relevant. Als Ty-

penbezeichnungen dienen die Kürzel AFT (analoge Fassung) und DFT (digitale Fassung) in Verbindung mit einer zweistelligen Nummer.

Das Modul „Datum“ ist für Angaben über das Erstellungsdatum der technischen Fassung vorgesehen. „Datenschutz“ meint technische Maßnahmen zum Datenschutz, wie z.B. die Maskierung von Personennamen in Texten.

Da eine technische Fassung auf mehreren Datenträgern gespeichert sein kann, wurde dieser Abschnitt im Schema als iterativ gekennzeichnet. An erster Stelle dieses Abschnitts wird eine eindeutige Inventarnummer des zu dokumentierenden Datenträgers erwartet. Im Feld „Sonstige_Bezeichnung“ können weitere in einem Korpusprojekt möglicherweise generierte Ordnungskennzeichen (Name, Ordnungsnummer, etc.) erfasst werden. Im nächsten Schritt kann man über den Typ des Datenträgers (z.B. Papier, Mikrofilm, Daten-CD) informieren.

Die Komponenten „Digitalisierungssoftware“, „Elektronischer_Speicherort“ und „Technische_Daten“ sind nur für digitale Fassungen relevant. „Digitalisierungssoftware“ meint das Programm, mit dem eine analoge Fassung digitalisiert wurde. Im Feld „Elektronischer_Speicherort“ kann man eine URL oder einen Pfadnamen erfassen.

Das erste Feld ist im Abschnitt „Technische_Daten“ für den Dateinamen vorgesehen, an den sich eine Information über das Dateiformat anschließen sollte. „Character_Encoding“ steht für die Zeichencodierung in einer Textdatei (z.B. ASCII oder UTF-16BE). Im Feld „Größe“ ist die Dateigröße (Anzahl von Bytes) anzugeben.

7.3. Archivierung und Distribution

Auch dieses Schema enthält die Module „Archivierung“ und „Dokumentation“.

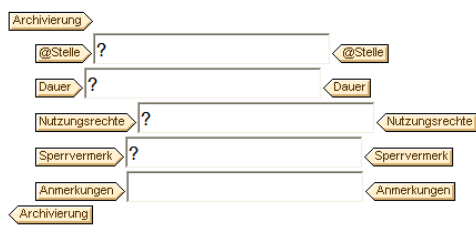


Abb. 96, Zusatzmaterial - Archivierung

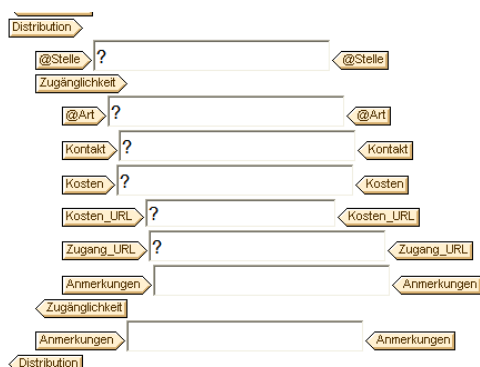


Abb. 97, Zusatzmaterial - Distribution

„Archivierung“ ist iterativ, kann also bei der korpuspezifischen Datenerfassung vervielfältigt werden. Zunächst soll der Name der archivierenden Stelle vermerkt werden. Es folgt ein Feld für Informationen über die vorgesehene Archivierungsdauer. Hier könnte z.B. „Bis 2018“ oder „Langfristig“ stehen. Die Nutzungsrechte der archivierenden Stelle können von der ausschließlichen Nutzung durch den Projektleiter bis hin zur Veröffentlichung eines Dokuments im Internet

reichen. Sperrvermerke, wie z.B. „Bis 2010 für Externe gesperrt“, können die Nutzungsmöglichkeiten einschränken.

Das Modul „Distribution“ ist ebenfalls iterativ und umfasst neben einem Feld für den Namen der für die Distribution zuständigen Stelle die iterative Komponente „Zugänglichkeit“, wo folgende Angaben verzeichnet werden können: Art der Zugänglichkeit, E-Mail-Kontaktadresse, Angaben über die Kosten, ggf. eine URL dieser Angaben sowie ggf. eine URL, die den direkten Zugang zum jeweiligen Korpusbestandteil ermöglicht.

7.4. Dokumentationsgeschichte

Das Schema für die Dokumentation von Zusatzmaterial auf der Korpusebene umfasst den in allen Schemata enthaltenen Komplex „Dokumentationsgeschichte“. Die Komponente „Update“ und beide Teile der Korrekturkomponente sind iterativ, d.h. dass sie bei der korpuspezifischen Datenerfassung vervielfältigt werden können.

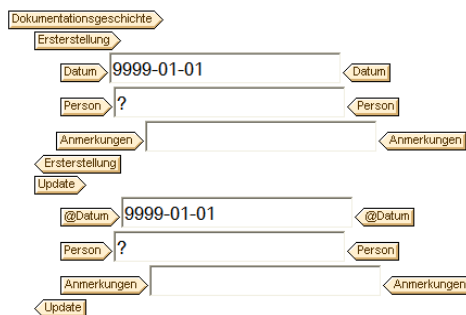


Abb. 98, Zusatzmaterial - Dokumentationsgeschichte (1)

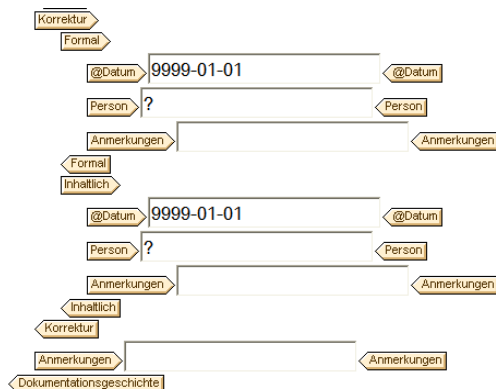


Abb. 99, Zusatzmaterial - Dokumentationsgeschichte (2)

8. Generisches Schema für die Korpusbeschreibung

Der vierte Bereich in unserem Datenmodell, der mithilfe eines (XML-)Schemas gestaltet wurde, ist die Korpusbeschreibung, die einen systematischen Überblick über Erstellung, Zusammensetzung, Bearbeitungsstand und Verwaltung eines Korpus ermöglichen soll. Bei der Gestaltung dieses Schemas haben wir auf Konzepte und Strukturelemente zurückgegriffen, die v.a. im 5. Abschnitt des vorliegenden Textes beschrieben sind

Auch dieses Schema enthält obligatorische und fakultative Komponenten. Obligatorische Komponenten sind in allen projektspezifischen Subschemata zu berücksichtigen, fakultative Komponenten stehen zur Wahl. Wenn man sie nutzt, sind alle Kennungsfelder und die mit ? gekennzeichneten Felder zu bearbeiten.

Eingaben für fehlende Daten in diesen Feldern sind standardisiert: „Nicht dokumentiert“ bedeutet: Es kann ein Datum geben, das bei der Datenerfassung jedoch nicht bekannt ist. Ein Beispiel dafür wäre: „Laufzeit: Nicht dokumentiert“ - zu lesen als: „Die Laufzeit des Erstellungsprojekts ist nicht dokumentiert.“ „Nicht vorhanden“ bedeutet: Es gibt kein Datum. Ein Beispiel dafür wäre: „Beschreibung: Nicht vorhanden“ - zu lesen als: „Es gibt keine Beschreibung des Erstellungsprojekts bzw. der Korpusarbeiten.“ In ein Feld („Datenrate“) kann auch der Wert „Nicht relevant“ eingegeben werden.

Das an vielen Stellen vorgesehene Feld „Anmerkungen“ ist für Anmerkungen zu Angaben in anderen Feldern und für nicht kategorisierte Angaben vorgesehen. Das Feld kann leer bleiben.

Einzelne Komponenten des Schemas wurden als iterativ gekennzeichnet, d.h. dass sie bei der Datenerfassung vervielfältigt werden können.

IDS-SCHEMA FÜR DIE KORPUSBESCHREIBUNG

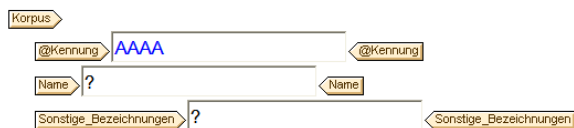


Abb. 100, Korpusbeschreibung (1)

Jede Korpusbeschreibung beginnt mit einer vierstelligen Korpuskennung, dem Korpusnamen und möglichen früheren Bezeichnungen.

8.1. Erstellungsprojekt

Unter „Erstellungsprojekt“ verstehen wir das Projekt, das ein Korpus aufgebaut hat. Bei der Korpuserstellung können Materialien aus anderen Projekten verwendet worden sein, was in bestimmten Komponenten des Komplexes „Korpusbestandteile“ (vgl. 8.4.) vermerkt werden kann.

Der in Abb. 101 gezeigte Abschnitt „Erstellungsprojekt“ wurde im Schema als iterativ gekennzeichnet, um Projektkooperationen bei der Korpuserstellung dokumentieren zu können. Er ist an der „IDS-Dokumentation zur Germanistischen Sprachwissenschaft - Sprachwissenschaftliche Forschungsvorhaben“ [17] orientiert, deren Struktur wir für unsere Zwecke erweitert haben.

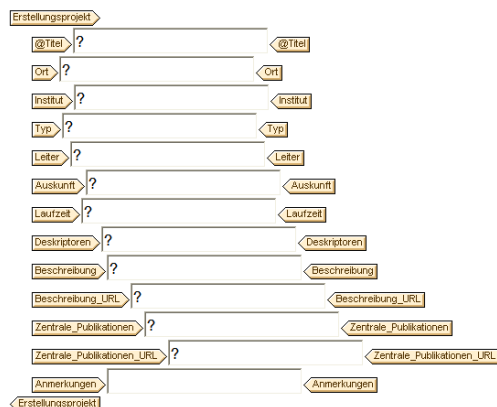


Abb. 101, Korpusbeschreibung - Erstellungsprojekt

Korpora entstehen zum einen in Projekten, in denen Datenbanken für eigene Forschungen benötigt werden (a). Erstellungsprojekte können allerdings auch dem alleinigen Zweck dienen, ein Korpus aufzubauen (b). Informationen über die Korpusarbeiten sollten in jedem Fall in dem gleichnamigen Abschnitt (vgl. 8.2.) notiert werden. Im Fall (b), kann man in den Feldern „Beschreibung“ und „Beschreibung_URL“ des Abschnitts „Erstellungsprojekt“ auf Informationen im Abschnitt „Korpusarbeiten“ verweisen, sofern solche Informationen vorliegen.

8.2. Korpusarbeiten

Wir nehmen an, dass wir die in Abb. 102 gezeigte Struktur nicht erläutern müssen, und notieren dazu lediglich einen u.E. wichtigen Hinweis: Bei der Beschreibung von Korpusarbeiten sollte man berücksichtigen, dass z.B. die Digitalisierung von Aufnahmen, die Transkription, die Text-Ton-Synchronisation sowie eine Validierung von Korpusdaten [16], d.h. eine Prüfung der Übereinstimmung der Daten mit den für das Korpus geltenden Spezifikationen, auch außerhalb von Erstellungsprojekten geleistet werden können.

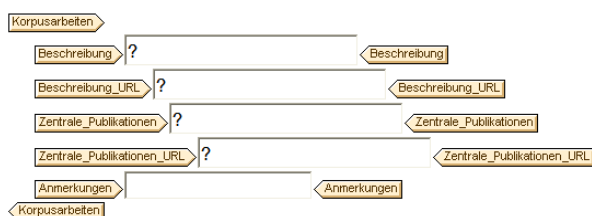


Abb. 102, Korpusarbeiten

8.3. Aufzeichnungsobjekte

In diesem Komplex greifen wir die Konzepte „Ereignis“ und „Sprechereignis“ wieder auf, die in den Abschnitten 5.1. und 5.1.6. eingeführt wurden. Zur Erinnerung daran wiederholen wir im Folgenden die o.g. Definitionen.

Unter „Ereignis“ (E) verstehen wir eine Phase des sozialen Geschehens, die von Beteiligten bzw. Korpusproduzenten als abgrenzbare Einheit wahrgenommen und aufgezeichnet wird. Unter „Sprechereignis“ (SE) verstehen wir den aufgezeichneten sprachlichen / kommunikativen Gehalt eines Ereignisses bzw. Segmente dieses Gehalts. Hier wie in den Abschnitten 5.1. und 5.1.6. gilt: Die Definitionen sind bewusst sehr allgemein gehalten. Wir stellen lediglich für die Dokumentation von Korpusbestandteilen relevante Konzepte bereit, keine linguistischen Segmentierungskriterien.

Im Feld Anzahl der Ereignis_Basisdaten kann man die Anzahl der aufgezeichneten Ereignisse notieren. Im Feld „Beschreibung“ wird eine kurze inhaltliche Charakterisierung der Ereignisse erwartet, die auf den Werten der Kategorie „Beschreibung“ der korpuspezifischen Ereignisdokumentationen (vgl. 5.1.3.) basieren sollte. „Institution“ verstehen wir im Sinne von „Organisation“. Im Feld Räumlichkeiten kann man über das räumliche Umfeld der Ereignisse berichten. „Zeit“ steht für den Zeitraum, in den die dokumentierten Ereignisse stattgefunden haben. Das Feld „Rundfunksendungen“ wurde für Korpora mit Mitschnitten solcher Sendungen eingerichtet. Hier sollten ggf. Anzahl und Typen der Sendungen (Hörfunksendung, Fernsehsendung) eingetragen werden.

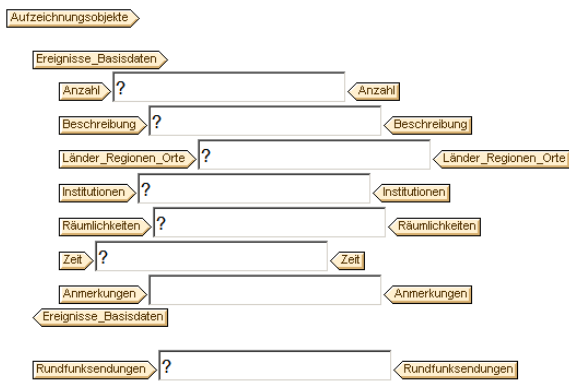


Abb. 103, Korpusbeschreibung - Aufzeichnungsobjekte (1)

Zusammenfassende Informationen über die aufgezeichneten Sprechereignisse sollten auf den Daten in der Komponente „Beschreibung“ der korpuspezifischen Sprechereignisdokumentationen beruhen. Zur Erinnerung wiederholen wir im Folgenden Erläuterungen von Kategorien aus dem Abschnitt vgl. 5.1.6.2.

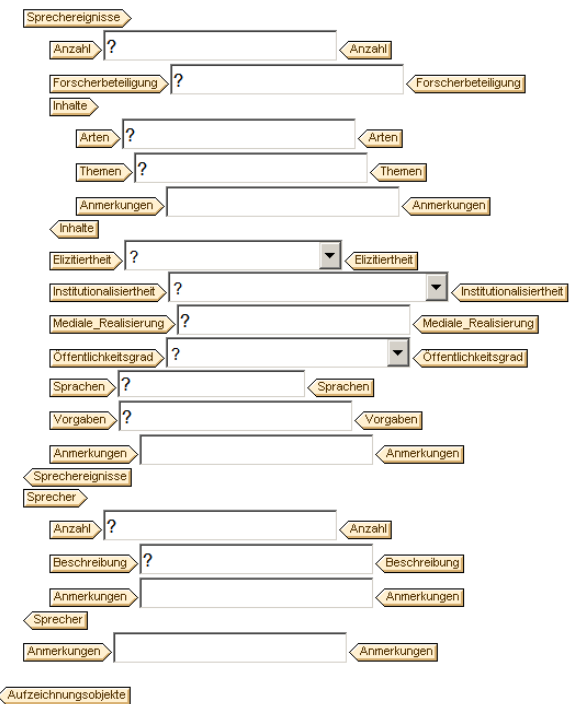


Abb. 104, Korpusbeschreibung - Aufzeichnungsobjekte (2)

Für „Forscherbeteiligung“ haben wir die Werte „Verbal beteiligt“, „Nicht verbal beteiligt“ und „Nicht vorhanden“ (für „Forscher nicht anwesend“) vorgesehen.

In der Korpusbeschreibung umfasst das Modul „Inhalt“ die Elemente „Art“ und „Themen“. Wir verwenden „Art“ anstelle von Kategorien wie „Textsorte“, „Texttyp“, „Interaktionstyp“, „Gesprächstyp“, „Diskurstyp“, „Genre“, „Gattung“, die aus verschiedenen Forschungsansätzen stammen, um für Daten aus allen Bereichen offen bleiben zu können. Bei „Art“ denken wir an Bezeichnungen wie „Erzählung“, „Rede“, „Anleitung“, „Beschreibung“, „Benennung“, „Übersetzung“, „Interview“, „Beratung“, „Diskussion“, „Begrüßung“ etc. Mit diesen Beispielwerten wollen wir keine Vorentscheidung über eine im Einzelfall anzuwendende Systematik treffen. Themenangaben können stichwortartig sein (z.B. „Politik“, „Recht“, „Studium“, „Lebenslauf“).

„Elizitierung“ ist eine Technik zur Erhebung sprachlicher Daten, bei der die Informanten systematisch zu Äußerungen veranlasst werden. Wir haben die Werte „Elizitiert“ und „Nicht elizitiert“ vorgesehen. „Institutionalisiertheit“ verstehen wir als Zugehörigkeit zu bzw. Erwartbarkeit eines Sprechereignisses im Rahmen einer Institution (im Sinne von „Organisation“). So fanden z.B. für das Korpus „Deutsch heute“ aufgezeichnete Sprechereignisse in Institutionen wie Schulen und Volkshochschulen statt, gelten in diesem Zusammenhang jedoch als nicht institutionell. Im Korpus „Schlichtungs- und Gerichtsverhandlungen“ dagegen sind Aufzeichnungen institutioneller Sprechereignisse in Vergleichsbehörden, Schlichtungsstellen und Gerichten enthalten.

„Mediale_Realisierung“ steht für den jeweiligen Kommunikationskanal (wie z.B. „Face to Face“, „Telefon“, „Hörfunk“). Für das Feld „Öffentlichkeitsgrad“ werden die Werte „Öffentlich“ und „Nicht öffentlich“ bereitgestellt. Im Feld „Sprachen“ sind die im Sprechereignis verwendeten Sprachen zu verzeichnen. Unter „Vorgaben“ sind Instruktionen der Sprecher durch die Aufnahmeleiter und ggf. auch Materialien, die den Sprechern zur Lösung bestimmter Aufgaben vorgelegt wurden, zu verstehen.

Im Unterschied zum generischen Schema für die Dokumentation von Korpusbestandteilen auf Ereignisebene und zum generischen Schema für allgemeine Sprecherdaten wurden hier nur zwei Felder für Sprecherdaten vorgesehen. Im ersten Feld soll die Anzahl der Sprecher verzeichnet werden, im zweiten Feld kann die Sprecherauswahl beschrieben sowie zusammenfassend über Sprechermerkmale und sprachliche Besonderheiten informiert werden.

8.4. Korpusbestandteile

Zu den Korpusbestandteilen zählen wir Quellaufnahmen von Ereignissen, sprechereignisspezifische Aufnahmen, Transkripte und Zusatzmaterial auf Ereignis-, Sprechereignis-, Sprecher- und Korpusebene. Die Strukturen der entsprechenden Komplexe in der Korpusbeschreibung unterscheiden sich kaum von den entsprechenden Strukturen in den oben beschriebenen Schemata. Das ermöglicht in der Korpusbeschreibung eine Bilanzierung des jeweiligen Bestandes aufgrund der für die einzelnen Bezugsobjekte gesammelten Daten.

8.4.1. Quellaufnahmen

Rohdaten, Originalaufnahmen von Ereignissen oder Aufnahmen, die für die dokumentierende Stelle Originalcharakter haben und die Quellen für sprechereignisspezifische Aufnahmen sein können, bezeichnen wir als Quellaufnahmen. Da nicht jedes Korpus Quellaufnahmen umfasst und da Quellaufnahmen unterschiedlichen Typs vorliegen können, wurde der Komplex im Schema als fakultativ und iterativ gekennzeichnet.



Abb. 105, Korpusbestandteile - Quellaufnahmen (1)

Der Typ der Aufnahmen (Audio, Video und ggf. Tonkopie von Video) sollte im gleichnamigen Feld notiert werden.

8.4.1.1. Basisdaten

Die Anzahl der Quellaufnahmen des genannten Typs kann man im gleichnamigen Feld verzeichnen. Im Anschluss daran wird dokumentiert, ob das Korpus vollständige und / oder unvollständige Aufnahmen der Ereignisse umfasst. Im nächsten Schritt sollte man Angaben über die

Dauer der einzelnen Aufnahmen (z.B. „zwischen 4 und 50 Minuten“) und die Gesamtdauer der Quellaufnahmen des genannten Typs notieren.

Abb. 106, Korpusbestandteile - Quellaufnahmen - Basisdaten

Quellaufnahmen müssen nicht unbedingt aus dem dokumentierten Erstellungsprojekt stammen. Sie können aus anderen Projekten bzw. Korpora übernommen worden sein. Um solchen Fällen gerecht werden zu können, wurde das Feld „Herkunft“ in die Basisdaten für Quellaufnahmen eingefügt. Quellaufnahmen können Daten enthalten, die nach dem Willen der Urheber und aus datenschutzrechtlichen Gründen Außenstehenden nicht kenntlich werden dürfen, wie z.B. persönliche Sprecherdaten. Für entsprechende Informationen wurde das Feld „Schutzbedürftige Daten“ bereitgestellt.

8.4.1.2. Aufnahmetechnik

Abb. 107, Korpusbestandteile - Quellaufnahmen - Aufnahmetechnik

Unter dem Stichwort „Aufnahmetechnik“ werden Informationen über die Aufnahmeapparatur (Aufnahmegerat, Mikrofone), eine ggf. eingesetzte Aufzeichnungssoftware, die Aufnahmege-
schwindigkeit (bei Spulentonbandaufnahmen relevante Angabe in cm/s) und Rauschunterdrückungsverfahren (z.B. Dolby B) zusammengefasst.

8.4.1.3. Technische Fassungen

Quellaufnahmen liegen in bestimmten technischen Fassungen vor. Das können analoge und / oder digitale Fassungen sein. Für jeden Typ gibt es einen eigenen Abschnitt im Schema. Beide Abschnitte sind fakultativ und iterativ, wenigsten einen Abschnitt muss bei der Erstellung projektspezifischer Schemata übernommen werden.

In allen Abschnitten werden Typen definiert. Bei analogen Fassungen sind die Werte in den Feldern „Datenschutz“ und „Kanäle“ für eine Typisierung relevant, bei digitalen Fassungen auch Informationen über die technischen Daten (z.B. das Dateiformat). Die Typenbezeichnungen setzen sich zusammen aus „AFT“ (für „Analoge_Fassungen_Typ“) bzw. „DFT“ (für „Digitale_Fassungen_Typ“) und einer zweistelligen Nummer. Beispiele sehen Sie in den Abb. 108 und

109. Diese Typenbezeichnungen werden auch bei der Dokumentationen einzelner Aufnahmen verwendet.

Im Feld „Anzahl“ soll über die Anzahl der technischen Fassungen des jeweiligen Typs informiert werden. „Datenschutz“ meint technische Maßnahmen zum Datenschutz, wie z.B. die Anonymisierung von Personennamen durch Verzerrung. Das Feld „Kanäle“ steht für Angaben wie „Mono“ oder „Stereo“ bereit. Für Angaben über die Qualität der Fassungen wurden zwei Felder bereitgestellt: „Bewertung“ und „Probleme“. Über die Datenträger kann im gleichnamigen Feld informiert werden.

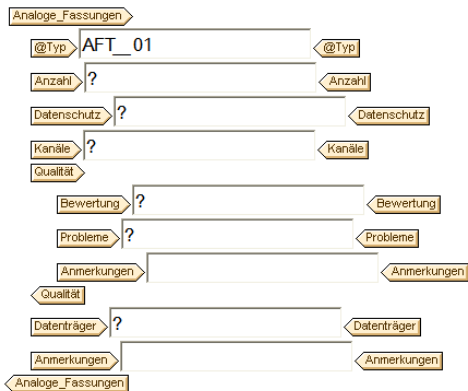


Abb. 108, Korpusbestandteile - Quellaufnahmen - Analoge Fassungen

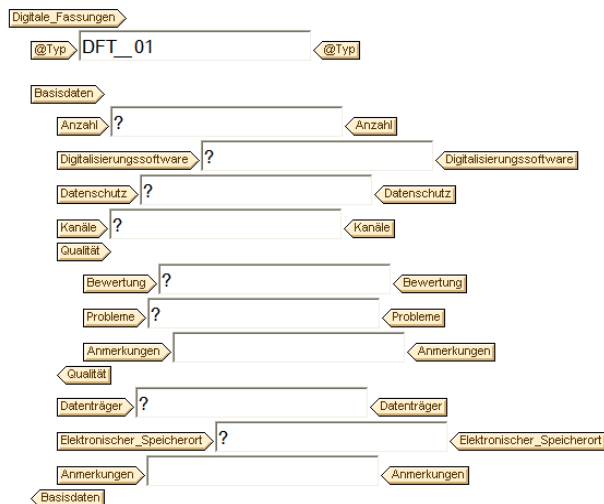


Abb. 109, Korpusbestandteile - Quellaufnahmen - Digitale Fassungen (1)

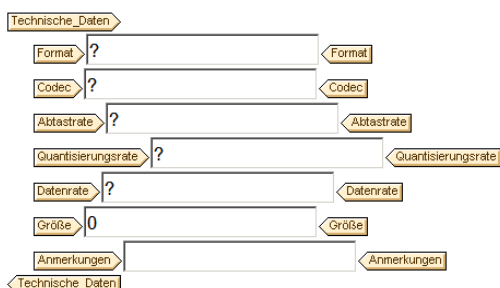


Abb. 110, Korpusbestandteile - Quellaufnahmen - Digitale Fassungen (2)

Zur Erinnerung an die für digitale Fassungen von Quellaufnahmen relevanten Daten wiederholen wir im Folgenden die in Abschnitt 5.1.4.3. notierten Erläuterungen:

„Digitalisierungssoftware“ steht für die Software, die ggf. bei der Digitalisierung einer analogen Fassung verwendet wurde. [11] In das Feld „Elektronischer Speicherort“ kann man eine URL oder einen Pfadnamen eintragen.

Im Feld „Format“ werden Angaben über das Dateiformat und das Digitalisierungsverfahren bzw. Audioformat erfasst. Ein möglicher Feldwert wäre: „WAVE (Linear PCM)“. Informationen über relevante Audioformate finden Sie unter der Adresse <http://de.wikipedia.org/wiki/Audioformat>

Als „Codec“ bezeichnet man ein Verfahren bzw. Programm, das Daten oder Signale digital kodiert und dekodiert. Unter der Adresse <http://de.wikipedia.org/wiki/Codec> finden Sie eine Liste mit Namen gängiger Codecs.

„Abtastung“ (engl. sampling) bezeichnet die Registrierung von Messwerten zu diskreten, meist äquidistanten Zeitpunkten. Aus einem zeitkontinuierlichen Signal wird so ein zeitdiskretes Signal gewonnen. Die Anzahl der Abtastungen pro Zeiteinheit wird Abtastrate genannt und meist in Hertz (Hz = Anzahl pro Sekunde) angegeben. Mögliche Werte sind z.B. „44100“ oder „48000“. Nach der Abtastung erfolgt die Quantisierung des zeitdiskreten, aber noch wertkontinuierlichen Signals. Dadurch entsteht ein zeit- und wertdiskretes Signal. Die Quantisierungsrate (auch Samplingtiefe oder Bittiefe) gibt die Anzahl der Bits an, die bei der Quantisierung pro Abtastwert verwendet werden. Typische Quantisierungsraten sind 8, 16 und 24 Bit.

Bei komprimierten Daten wird die Datenrate relevant - die Anzahl der Informationseinheiten, die pro Zeiteinheit gespeichert werden. Sie wird in kBit/s angegeben. Der Gesamtumfang der digitalen Fassungen eines Typs soll in MegaBytes erfasst werden.

8.4.1.4. Archivierung und Distribution

Auch die Module „Archivierung“ und „Distribution“, die Sie möglicherweise schon gesehen haben, tauchen hier wieder auf.

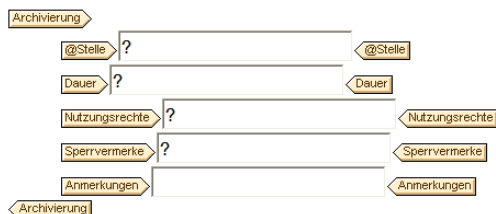


Abb. 111, Korpusbestandteile - Quellaufnahmen - Archivierung

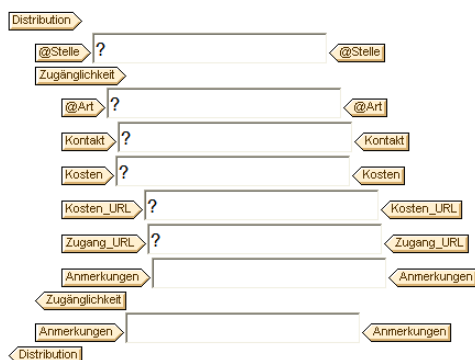


Abb. 112, Korpusbestandteile - Quellaufnahmen - Distribution

„Archivierung“ wurde als iterativ gekennzeichnet. Zunächst soll der Name der archivierenden Stelle vermerkt werden. Es folgt ein Feld für Informationen über die vorgesehene Archivierungsdauer. Hier könnte z.B. „Bis 2018“ oder „Langfristig“ stehen. Die Nutzungsrechte der archivierenden Stelle können von der ausschließlichen wissenschaftlichen Auswertung durch den Aufnahmeleiter bis hin zur Veröffentlichung einer Aufnahme im Internet reichen. Sperrvermerke, wie z.B. „Bis 2010 für Externe gesperrt“, können die Nutzungsmöglichkeiten einschränken.

„Distribution“ ist ebenfalls iterativ und umfasst neben einem Feld für den Namen der für die Distribution zuständigen Stelle die iterative Komponente „Zugänglichkeit“. In dieser Komponente sollen folgende Angaben verzeichnet werden: Art der Zugänglichkeit, E-Mail-Kontaktadresse, Angaben über die Kosten, ggf. eine URL dieser Angaben sowie ggf. eine URL, die einen direkten Zugang zu den Aufnahmen ermöglicht.

8.4.2. Sprechereignisspezifische Aufnahmen (SE-Aufnahmen)

Neben Quellaufnahmen dokumentieren wir sprechereignisspezifische Aufnahmen (SE-Aufnahmen), die in einem bestimmten Verhältnis zu den Quellaufnahmen stehen. Das können Segmente in den Quellaufnahmen bzw. Kopien dieser Segmente sein. Es kommt allerdings auch vor, dass alle SE-Aufnahmen mit den Quellaufnahmen übereinstimmen. Für solche Fälle haben wir die Struktur des iterativen Komplexes SE-Aufnahmen besonders flexibel gestaltet.



Abb. 113, Korpusbestandteile - SE-Aufnahmen (1)

Der Typ der Aufnahmen (Audio, Video und ggf. Tonkopie von Video) sollte im gleichnamigen Feld notiert werden.

8.4.2.1. Basisdaten

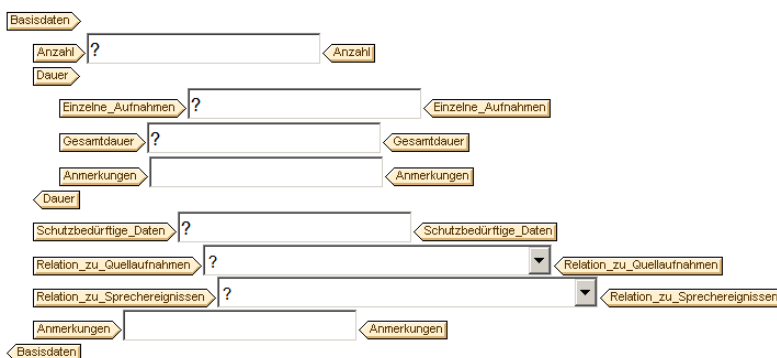


Abb. 114, Korpusbestandteile - SE-Aufnahmen - Basisdaten

Die Anzahl der SE-Aufnahmen des genannten Typs kann man im gleichnamigen Feld notieren. Im nächsten Schritt werden Informationen über die Dauer einzelner Aufnahmen (z.B. „zwischen 4 und 50 Minuten“) und die Gesamtdauer der SE-Aufnahmen des genannten Typs erwartet.

Die Felder „Dauer“ und „Schutzbedürftige_Daten“ sind in diesem Komplex fakultativ, d.h. dass sie bei der Erstellung korpuspezifischer Schemata übergangen werden können, wenn Quellaufnahmen dokumentiert wurden und sich die SE-Aufnahmen eines Korpus von den Quellaufnahmen nicht unterscheiden.

Im Feld „Relation_zu_Quellaufnahmen“ sollte vermerkt werden, ob die SE-Aufnahmen mit den Quellaufnahmen übereinstimmen oder ob es sich um Segmente in den Quellaufnahmen handelt. Sprechereignisse können in SE-Aufnahmen vollständig oder unvollständig aufgezeichnet sein. Für solche Angaben gibt es das Feld „Relation_zu_Sprechereignissen“.

8.4.2.2. Transkribierte SE-Aufnahmen

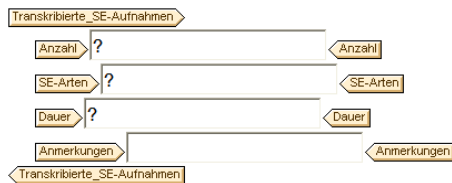


Abb. 115, Korpusbestandteile - SE-Aufnahmen - Transkribierte SE-Aufnahmen

Damit dokumentiert werden kann, wie viele SE-Aufnahmen und welche Arten von Sprechereignissen transkribiert sind, und für eine Information über die Dauer der transkribierten Aufnahmen wurde das Modul „Transkribierte_SE-Aufnahmen“ eingefügt.

8.4.2.3. Technische Fassungen

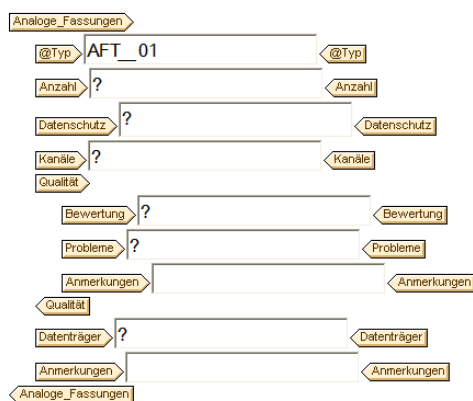


Abb. 116, Korpusbestandteile - SE-Aufnahmen - Analoge Fassungen

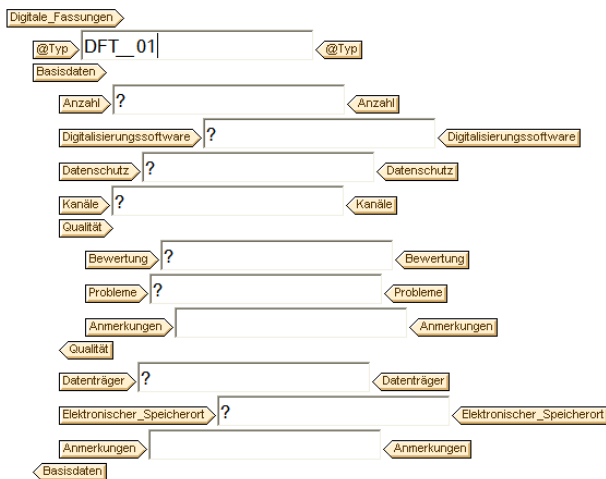


Abb. 117, Korpusbestandteile - SE-Aufnahmen - Digitale Fassungen (1)

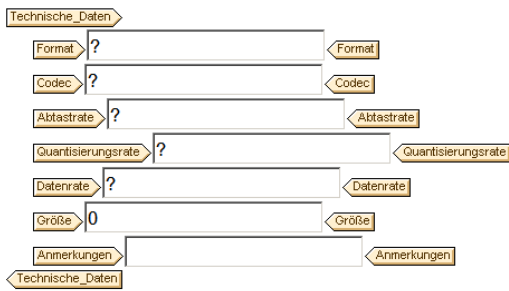


Abb. 118, Korpusbestandteile - SE-Aufnahmen - Digitale Fassungen (2)

Der Komplex „SE-Aufnahmen“ enthält wie der Komplex „Quellaufnahmen“ die fakultativen und iterativen Module „Analoge_Fassungen“ und „Digitale_Fassungen“. Auch diese Teile können übergangen werden, wenn die Quellaufnahmen dokumentiert wurden und sich die SE-Aufnahmen eines Korpus von den Quellaufnahmen nicht unterscheiden.

Die in den Abb. 116 bis 118 gezeigten Strukturen stimmen mit den in Abschnitt 8.4.1.3. erläuterten überein.

8.4.2.4. Archivierung und Distribution

Die Komponenten „Archivierung“ und „Distribution“ wurden zuletzt in Abschnitt 8.4.1.4. vorgestellt. Wenn die Quellaufnahmen dokumentiert wurden und sich die SE-Aufnahmen eines Korpus von den Quellaufnahmen nicht unterscheiden, können diese Module bei der Beschreibung von SE-Aufnahmen übergangen werden.

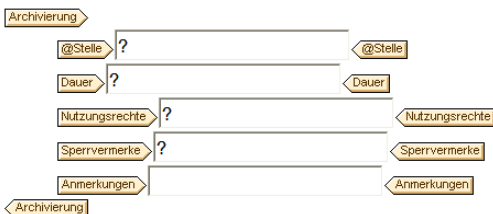


Abb. 119, Korpusbestandteile - SE-Aufnahmen - Archivierung

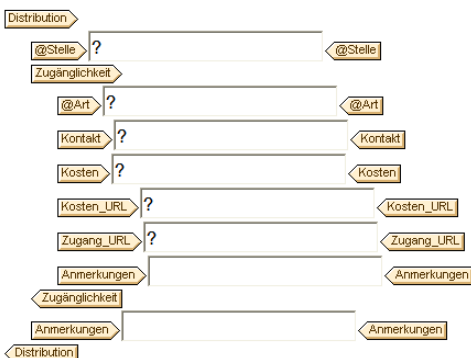


Abb. 120, Korpusbestandteile - SE-Aufnahmen - Distribution

8.4.3. Transkripte

Zu den Korpusbestandteilen zählen auch Transkripte, die allerdings nicht in jedem Korpus enthalten sind, weshalb der Transkriptkomplex im Schema als fakultativ gekennzeichnet wurde.



Abb. 121, Korpusbestandteile - Transkripte (1)

An erster Stelle des Transkriptkomplexes wird die Bezeichnung eines Transkripttypes erwartet. Die Typisierung kann zum einen über die Extensionen (vollständige Transkripte vs. Teiltranskripte) erfolgen, zum anderen über Art und Anzahl der Annotationen. Typenbezeichnungen wie die in Abb. 121 werden auch bei der Dokumentation einzelner Transkripte verwendet.

8.4.3.1. Basisdaten

Im Modul „Basisdaten“ der Transkriptdokumentation kann man die Anzahl der Transkripte des genannten Typs, einen Hinweis auf möglicherweise in den Transkripten enthaltene schutzbedürftige Daten sowie eine Information darüber notieren, ob es sich um vollständige Transkripte oder Teiltranskripte von SE-Aufnahmen handelt.

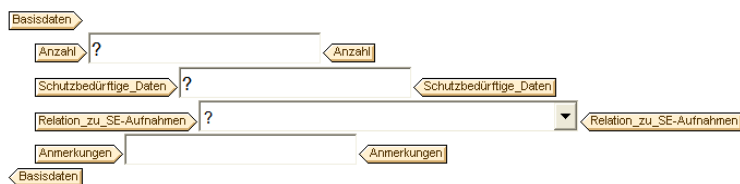


Abb. 122, Korpusbestandteile - Transkripte - Basisdaten

8.4.3.2. Annotationen

In der Beschreibung des Moduls „Annotationen“ greifen wir Erläuterungen aus dem Abschnitt 5.1.6.5.2. wieder auf.

Wir verwenden die Bezeichnung „Annotation“ für inhaltlich und formal charakterisierte Ebenen eines Transkripts, wie z.B. Aufzeichnungen des Wortlauts in orthographischer, literarischer oder phonetischer Umschrift, syntaktische Angaben, Notationen suprasegmentaler oder nonverbaler Phänomene, Übersetzung des Wortlautes etc. [13] Da u.U. mehrere Annotationen zu dokumentieren sind, wurde der Komplex im Schema als iterativ gekennzeichnet. [14]

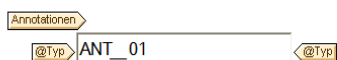


Abb. 123, Korpusbestandteile - Transkripte - Annotationen - Typ

Für jeden ermittelten Annotationstyp wird eine Bezeichnung generiert, die sich zusammensetzt aus dem Kürzel ANT (für „Annotation_Typ“) und einer zweistelligen Nummer. Ein Beispiel finden Sie in Abb. 123. Diese Typenbezeichnungen werden auch in die Dokumentation einzelner Transkripte (vgl. 5.1.6.5.2.) eingesetzt.

8.4.3.2.1. Basisdaten

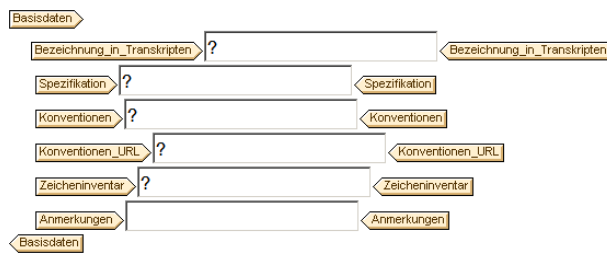


Abb. 124, Korpusbestandteile - Transkripte - Annotationen - Basisdaten

Im ersten Feld der Basisdaten kann man die in Transkripten verwendete Bezeichnung des Annotationstyps notieren. Eine „Spezifikation“ der Annotationen des genannten Typs sollte Angaben über den Gegenstand (z.B. „Wortlaut“), die Umschrift (z.B. „Literarisch“) und die Reichweite (z.B. „Ohne Interviewerbeiträge“) enthalten.

Auf die dem jeweiligen Annotationstyp zugrunde liegenden Konventionen ist im gleichnamigen Feld hinzuweisen. Beispiele für solche Hinweise wären: „Projektspezifisch“, „DIDA, Version vom Januar 2001“, „GAT“ etc. Unter „Zeicheninventar“ ist das Inventar an Schriftzeichen zu verstehen, das bei der Wiedergabe des Wortlautes verwendet wurde. Das sind i.d.R. standardisierte Inventare wie z.B. der IPA-Zeichensatz oder ein spezifisches Alphabet.

8.4.3.2.2. Erstellung

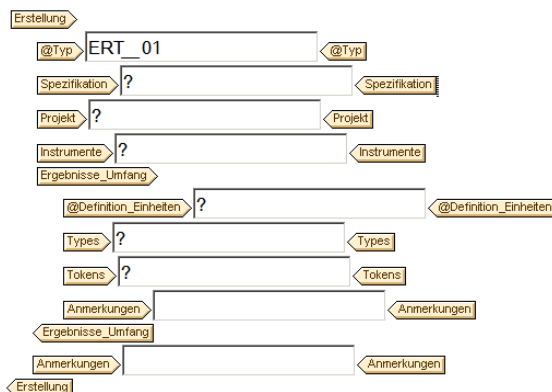


Abb. 125, Korpusbestandteile - Transkripte - Annotationen - Erstellung

Auch im iterativen Modul „Erstellung“ der Korpusbeschreibung werden Typen benannt. Für die Typisierung relevant sind die Spezifikationen der Erstellungsprozesse. Abb. 125 enthält ein Beispiel für eine Typenbezeichnung, die auch bei der Dokumentation einzelner Transkripte eingesetzt .

Im Feld „Spezifikation“ sollten v.a. der Arbeitsstand (z.B. „Ersterfassung“, „1. Korrektur“, „Endkorrektur“, „Überarbeitung für Publikation xy“) und mögliche besondere Umstände (z.B. „halbautomatisch“) der Erstellung dokumentiert werden. Dann folgen Fragen nach dem für die Erstellung des genannten Typs zuständigen Projekt sowie die bei der Erstellung genutzten Instrumente (Editoren und ggf. Systemumgebung).

Die Information über den Umfang der Ergebnisse eines Erstellungstyps umfasst eine Definition der gezählten Einheiten sowie Felder für Angaben über die Anzahl unterschiedlicher Einheiten (Types) und die Anzahl aller gezählten Einheiten (Tokens).

8.4.3.2.3. Alignment

Wir verwenden die Bezeichnung „Alignment“ in der Dokumentation für die Text-Ton-Synchronisation, also die Koppelung von Aufnahmen und Transkripten auf Phon-, Phonem-, Wort- oder Phrasenbasis, wobei Transkriptsegmenten Zeitmarken zugeordnet werden. Die entsprechende Komponente ist fakultativ und iterativ.

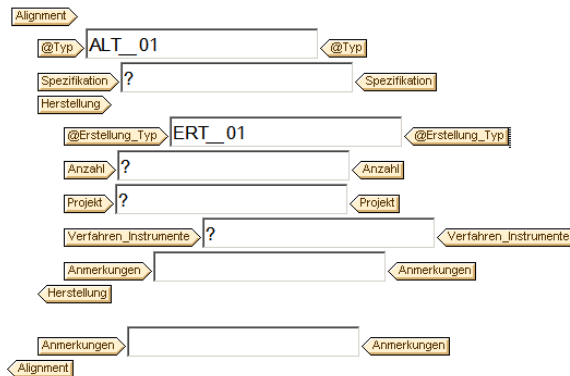


Abb. 126, Korpusbestandteile - Transkripte - Annotationen - Alignment

Auch im Modul „Alignment“ ist eine Typisierung vorgesehen, die sich auf die Spezifikation von Alignmentprozessen stützt. Die hier definierten Typen sollen auch bei der Dokumentation einzelner Transkripte berücksichtigt werden. Ein Beispiel für eine Typenbezeichnung finden Sie in Abb. 126. Im Feld „Spezifikation“ werden Angaben über die für den genannten Typ relevanten Segmente (z.B. „Phonweise“, „Wortweise“) erwartet.

Die Komponente „Herstellung“ ist iterativ. Zunächst wird der Typ der Erstellungen verzeichnet, deren Ergebnisse aligniert wurden. Im Feld „Projekt“ wird der Namen des Projekts, in dem das Alignment vorgenommen wurde, erwartet. Im Feld „Verfahren_Instrumente“ kann man Angaben darüber machen, ob manuell oder automatisch aligniert wurde, auf die genutzte Software hinweisen und ggf. weitere Informationen über die Systemumgebung erfassen.

8.4.3.3. Technische Fassungen

Um die technischen Fassungen von Transkripten dokumentieren zu können, wurden auch in den Transkriptkomplex die Module „Analoge_Fassungen“ und „Digitale_Fassungen“ eingefügt. Beide Abschnitte sind fakultativ und iterativ, wenigstens ein Abschnitt muss bei der Erstellung projektspezifischer Schemata übernommen werden.

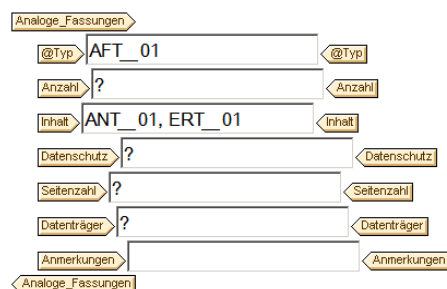


Abb. 127, Korpusbestandteile - Transkripte - Analoge Fassungen

Abb. 128, Korpusbestandteile - Transkripte - Digitale Fassungen (1)

Abb. 129, Korpusbestandteile - Transkripte - Digitale Fassungen (2)

In jedem Abschnitt wird ein Typ abgefragt. Bei analogen Fassungen sind die Werte in den Feldern „Inhalt“ und „Datenschutz“ für eine Typisierung relevant, bei digitalen Fassungen auch Informationen über die technischen Daten (z.B. das Dateiformat). Die Typenbezeichnungen werden auch bei der Dokumentation einzelner Transkripte verwendet. Die Anzahl der Fassungen des genannten Typs kann man im gleichnamigen Feld notieren.

Im Feld „Digitalisierungssoftware“ kann man über das Programm, mit dem eine analoge Fassung digitalisiert wurde, informieren. Im iterativen Feld „Inhalt“ sollte man alle Annotationen sowie die Typen der Erstellungen und der Alignments verzeichnen, deren Ergebnisse in den technischen Fassungen des genannten Typs gespeichert sind. „Datenschutz“ meint technische Maßnahmen zum Datenschutz, wie z.B. die Maskierung von Personennamen in Texten.

Bei seitenformatierten Texten gibt eine Information über die Seitenzahl einen groben Überblick über den Umfang des Materials. Nach einer Information über den oder die Datenträger wird der elektronische Speicherort der dokumentierten digitalen Fassungen erfasst. An dieser Stelle kann man eine URL oder einen Pfadnamen eintragen. Das erste Feld im Abschnitt „Technische_Daten“ ist für den Namen des Datenformats vorgesehen. „Character_Encoding“ steht für die Zeichencodierung (z.B. ASCII oder UTF-16BE). Im Feld „Größe“ kann man den Gesamtumfang der jeweiligen digitalen Fassungen (in MegaBytes) angeben.

8.4.3.4. Archivierung und Distribution

In den Abb. 130 und 131 sehen Sie die Strukturen der schon bekannten Bausteine „Archivierung“ und „Distribution“, die zuletzt in Abschnitt 8.4.1.4. erläutert wurden.

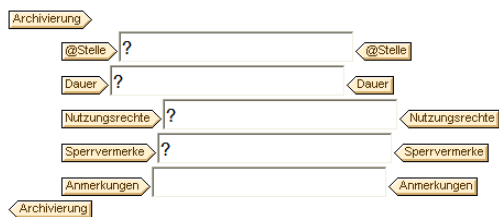


Abb. 130, Korpusbestandteile - Transkripte - Archivierung

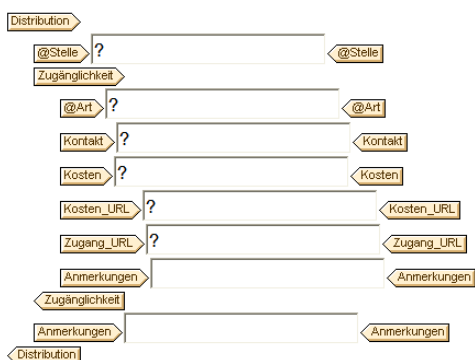


Abb. 131, Korpusbestandteile - Transkripte - Distribution

8.4.4. Zusatzmaterial

Unter „Zusatzmaterial“ verstehen wir Dokumente, die zusätzlich zu Aufnahmen und Transkripten vorhanden sein können. Zusatzmaterialien können auf verschiedenen dokumentarischen Ebenen angesiedelt sein: auf der Ebene der aufgezeichneten Ereignisse (z.B. Skizzen von Sitzordnungen), der Sprechereignisse (z.B. Ablaufprotokolle), der Sprecher (z.B. Sprecherfotos) und auf der Korpusebene (z.B. Transkriptionskonventionen). Der Komplex „Zusatzmaterial“ wurde im Schema als fakultativ und iterativ gekennzeichnet, d.h. dass er bei der Erstellung korpuspezifischer Schemata übergangen und, wenn er gewählt wird, bei der Dateneingabe vielfältig werden kann.



Abb. 132, Korpusbestandteile - Zusatzmaterial (1)

Im ersten Feld dieses Komplexes sollte die Art des Zusatzmaterials (vgl. o.g. Beispiele) beschrieben werden.

8.4.4.1. Basisdaten

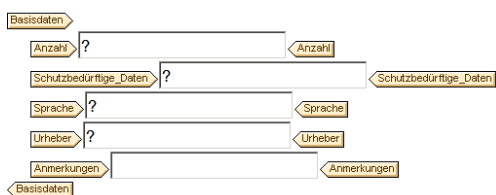


Abb. 133, Korpusbestandteile - Zusatzmaterial - Basisdaten

In der Komponente „Basisdaten“ kann man festhalten, wie viele Zusatzmaterialien der genannten Art vorhanden sind, ob und wenn ja welche schutzbedürftigen Daten diese Dokumente ent-

halten und welcher Sprache Textdokumente abgefasst wurden. Informationen über den oder die Urheber können im gleichnamigen Feld erfasst werden.

8.4.4.2. Technische Fassungen

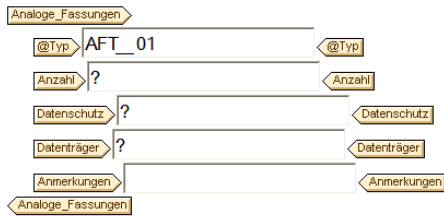


Abb. 134 , Korpusbestandteile - Zusatzmaterial - Analoge Fassungen

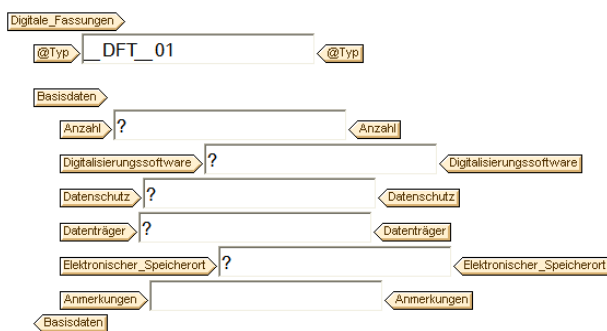


Abb. 135 , Korpusbestandteile - Zusatzmaterial - Digitale Fassungen (1)

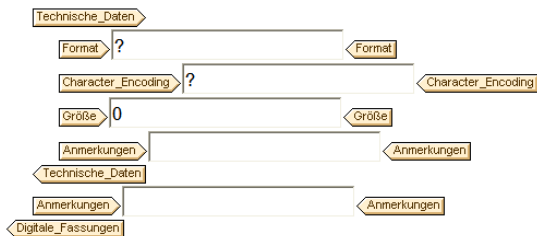


Abb. 136 , Korpusbestandteile - Zusatzmaterial - Digitale Fassungen (2)

Die für die Beschreibung technischer Fassungen von Zusatzmaterial bereitgestellten Module „Analoge_Fassungen“ und „Digitale_Fassungen“ stimmen mit den entsprechenden Modulen im Komplex „Transkripte“ (vgl. 8.4.3.) weitgehend überein.

8.4.4.3. Archivierung und Distribution

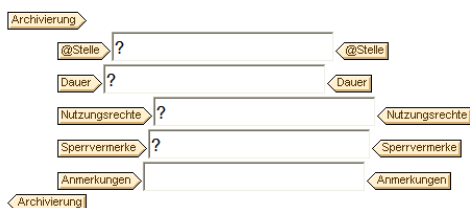


Abb. 137, Korpusbestandteile - Zusatzmaterial - Archivierung

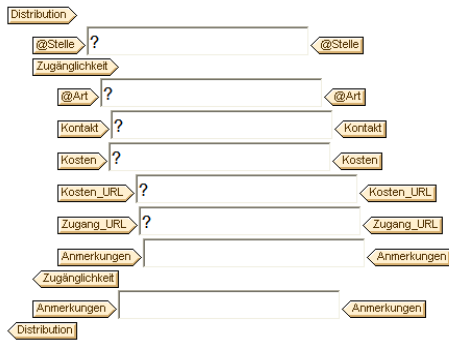


Abb. 138, Korpusbestandteile - Zusatzmaterial - Distribution

Auch mit den Bausteinen „Archivierung“ und „Distribution“ haben wir Sie schon bekannt gemacht, Erläuterungen finden Sie in 8.4.1.4.

8.7. Dokumentationsgeschichte

Informationen über Arbeitsstand und Bearbeiter der Dokumente werden bei der manuellen Dateneingabe automatisch in einer (Oracle-)Datenbank gespeichert, sollten nach unserer Vorstellung jedoch auch in den Dokumenten sichtbar sein. Daher haben wir am Ende aller Schemata den Baustein „Dokumentationsgeschichte“ eingebaut. Die Komponente „Update“ und beide Teile der Korrekturkomponente wurden als iterativ gekennzeichnet.

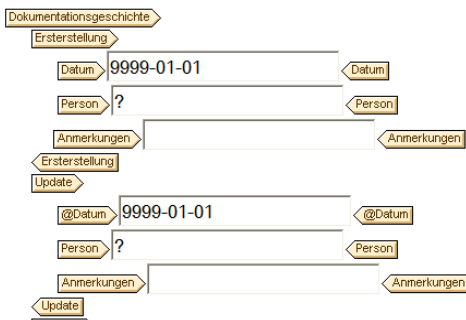


Abb. 139, Dokumentationsgeschichte - Ersterstellung, Update

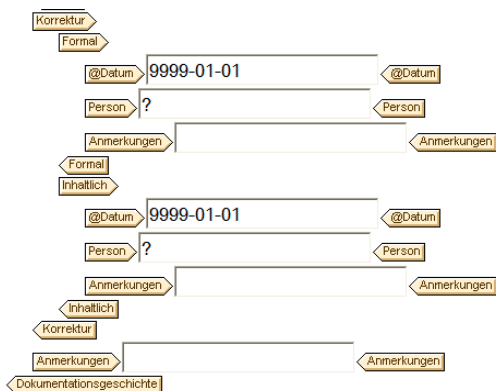


Abb. 140, Dokumentationsgeschichte - Korrektur

Mit dieser letzten Ansicht auf die Struktur des Moduls „Dokumentationsgeschichte“ beenden wir unsere Metadatenbeschreibung und gehen über zu einigen abschließenden Bemerkungen.

9. Abschließende Bemerkungen

Wir haben im vorliegenden Text Grundlagen der Metadatenkomponente der Datenbank für Gesprochenes Deutsch (DGD 2.0) vorgestellt. Im Mittelpunkt der Beschreibung standen vier generische (XML-)Metadaten-Schemata - das Schema für die Dokumentation von Korpusbestandteilen auf Ereignisebene, kurz „Ereignisschema“, das Schema für die Dokumentation ereignis- und sprechereignisübergreifender Sprecherdaten, kurz „Sprecherschema“, das Schema für die Dokumentation von Zusatzmaterial auf Korpusebene sowie ein Schema für überblicksartige Korpusbeschreibungen. In allen Fällen handelt es sich um weitreichende Kategoriensammlungen in flexiblen Strukturen, von denen korpuspezifische Subschemata abgeleitet werden können.

Die auf den oben vorgestellten Schemata basierenden korpuspezifischen Dokumente dienen in erster Linie der projekt- und archivinternen Dokumentation. Für die Außendarstellung von Korpora, Korpusbestandteilen und Sprechern, über die Externe informiert werden können, werden reduzierte Ansichten entwickelt, d.h. dass nicht alle Daten sichtbar werden.

Auf eine Besonderheit der Entwicklung möchten wir an dieser Stelle noch einmal aufmerksam machen: Wir orientierten uns zunächst u.a. an der ISLE MetaData Initiative (IMDI) [5], die Konventionen für die Veröffentlichung von Metadaten linguistischer Ressourcen vorgelegt hat, sind aber angesichts anderer Aufgaben zu anderen Ergebnissen gekommen als IMDI.

Wir setzten auf vergleichsweise strikte dokumentarische Vorgaben, die wir aus mehreren Gründen für nötig erachten. Die vielfältigen technologischen Möglichkeiten, große Korpora mit einzelnen Bestandteilen in unterschiedlichen Fassungen aufzubauen, spezielle Erfordernisse der Langzeitarchivierung digitaler Daten sowie das Bestreben, solche Bestandteile projektübergreifend nutzbar zu machen, führen zu hohen Anforderungen an die Qualität von Metadaten, mit denen einzelne Projekte nach unserer Erfahrung i.d.R. überfordert sind. Hier besteht ein Regelungsbedarf, dem wir gerecht werden wollen.

Im Rahmen der Begutachtung eines Vortragskonzepts wurden wir gefragt: „Are you suggesting this annotation as a standard?“ Die Antwort lautet: Ja, wir möchten Standards für die Dokumentation von IDS-Korpora der gesprochenen Sprache bereitstellen und hoffen, korpuserstellende Projekte für unser Konzept gewinnen zu können.

Zu guter Letzt bedanken wir uns bei den KollegInnen im Archiv für Gesprochenes Deutsch (AGD) für ihre Unterstützung dieser Arbeit.

10. Anmerkungen

[1] Dublin Core Metadata Initiative. <http://dublincore.org/>

[2] Simons, Gary & Steven Bird (Hg.) (2006): OLAC Metadata. <http://www.language-archives.org/>

[3] Text Encoding Initiative. <http://www.tei-c.org/>

[4] Martínez, José M. (Hg.) (2004): MPEG-7 Overview (version 10). <http://www.chiariglione.org/mpeg/standards/mpeg-7/mpeg-7.htm>

[5] ISLE MetaData Initiative. <http://www.mpi.nl/IMDI/>

[6] Trippel, Thorsten & Tanja Baumann (2003): Metadaten für Multimodale Korpora: Verwendung im Modellex-Projekt. Technisches Dokument 4, Universität Bielefeld. http://www.spectrum.uni-bielefeld.de/modellex/publication/techdoc/modellex_techrep4/

Die Originalformatierung der zitierten Textpassage wurde aus Platzgründen nicht übernommen.

[7] IMDI, Part 1, Metadata Elements for Session Description, Version 3.0.4, October 2003, S. 10. <http://www.mpi.nl/IMDI/>

[8] Schiel, Florian & Christoph Draxler (2004): The production of speech corpora. Version 2.5. <http://www.phonetik.uni-muenchen.de/forschung/BITS/TP1/Cookbook/>

[9] Gesamtkatalog der Tonaufnahmen des Deutschen Spracharchivs. Erarbeitet von Mitarbeiterinnen und Mitarbeitern des Instituts für Deutsche Sprache. Phonai Bde. 38 u. 39 - Tübingen: Niemeyer, 1992.

[10] Schreibkonventionen findet man im „Regelwerk Mediendokumentation“ <http://rmd.dra.de/arc/php/main.php>

[11] Wir nennen hier lediglich zur Veranschaulichung einige Produkte: Adobe Audition, Audacity, Audiograbber, Digidesign ProTools, No23 Recorder, Steinberg Wavelab.

[12] <http://agd.ids-mannheim.de/html/korpora/pdf/isdok.pdf>

[13] Das bedeutet, dass hier unter „Annotation“ nicht etwa die einzelne, segmentbezogene und auf einer semantisch definierten und einem Sprecher zugeordneten Annotationsspur eingetragene „Annotation“ verstanden wird, wie der Sprachgebrauch in ELAN wäre. (Mitteilung von Wilfried Schütte)

[14] Verschiedenen Annotationen können vermischt sein (z.B. Wortlaut in literarischer Umschrift plus Intonationsnotation). In solchen Fällen empfiehlt es sich, in der Dokumentation mit einem Annotationskomplex zu arbeiten, in dem alle Angaben zusammengefasst werden. Wenn verschiedenen Annotationen getrennt gehalten sind - z.B. a) Wortlaut in literarischer Umschrift plus Intonationsnotation, b) morphosyntaktische Annotation, c) Notation nonverbaler Phänomene - sollten für a), b) und c) verschiedene Annotationskomplexe angelegt werden.

[15] Caren Brinckmann, Stefan Kleiner, Ralf Knöbl and Nina Berend (2008): German Today: an areally extensive corpus of spoken Standard German. In: Proceedings 6th International Conference on Language Resources and Evaluation (LREC 2008), Marrakesch, Marokko. <http://www.lrec-conf.org/proceedings/lrec2008/summaries/806.html>

[16] Florian Schiel, Angela Baumann, Christoph Draxler, Tania Ellbogen, Phil Hoole, Alexander Steffen: The Validation of Speech Corpora. Version 1.11 : June 3, 2004.
<http://www.phonetik.uni-muenchen.de/forschung/BITS/TP2/Cookbook/Tp2.html>

[17] <http://www.ids-mannheim.de/oea/forsch/>