

MANUELA SCHÖNENBERGER UND ERIC HAEBERLI

Ein geparstes und grammatisch annotiertes Korpus schweizerdeutscher Spontansprachdaten

Das Projekt ‚Studying variation in syntax: a parsed corpus of Swiss German‘ wird vorgestellt, dessen Ziel es ist, ein geparstes und grammatisch annotiertes Korpus von ca. 1 Million Wörtern zu erstellen. Spontansprachdaten von Gewährspersonen unterschiedlichen Alters, die alle Muttersprachler der lokalen Varietät des St. Galler Deutschen sind, die in Wil gesprochen wird, werden durch informelle Interviews erhoben und aufgezeichnet. Nach Transkription dieser Interviews mit EXMARaLDA (s. SCHMIDT/WÖRNER 2009) werden die einzelnen Wörter grammatisch annotiert, damit sie später von einem Parser verarbeitet werden können. Aufgrund dieses Korpus hoffen wir, Variation innerhalb eines Sprechers (*intra-speaker variation*) sowie zwischen Sprechern (*inter-speaker variation*) aufzudecken und eventuell einen Einblick in den Zusammenhang zwischen syntaktischer Variation und Sprachwandel zu gewinnen. In diesem Beitrag steht das Erstellen der Datenbank im Vordergrund.

1. Einleitung

In diesem Beitrag stellen wir das Projekt ‚Studying variation in syntax: a parsed corpus of Swiss German‘ vor, das vom Schweizerischen Nationalfonds (SNF Projekt 146450) während einer Periode von 3 Jahren gefördert wird. Ziel des Projekts ist es, ein geparstes und grammatisch annotiertes Korpus von ca. 1 Million Wörtern zu erstellen. Aufgrund dieses Korpus hoffen wir, Variation innerhalb eines Sprechers (*intra-speaker variation*) sowie zwischen Sprechern (*inter-speaker variation*) aufzudecken und eventuell einen Einblick in den Zusammenhang zwischen syntaktischer Variation und Sprachwandel zu gewinnen. Wir konzentrieren uns in diesem Beitrag auf das Erstellen der Datenbank und stellen erste Ergebnisse vor.

Dieser Beitrag ist wie folgt gegliedert. Abschnitt 2 enthält eine Beschreibung der Gewährspersonen, des Interviewablaufs sowie der Tran-

skription der im Interview erhobenen Daten. Abschnitt 3 befasst sich mit der grammatischen Annotation einzelner Wörter und liefert einen kurzen Einblick ins Parsen. In Abschnitt 4 werden erste Ergebnisse, die aufgrund dieser Datenerhebung gewonnen wurden, zusammengefasst. Abschnitt 5 enthält unsere Schlussfolgerungen.

2. Erhebung und Transkription von Spontansprachdaten

In diesem Abschnitt beschreiben wir, wie wir Gewährspersonen für unsere Studie gesucht haben und welche an der Studie teilgenommen haben (2.1) sowie die Interviewsituation (2.2). Wie diese gesprochenen Daten verschriftlicht werden und welches Programm dazu verwendet wird, wird in 2.3 thematisiert.

2.1 Gewährspersonen

Alle unsere Gewährspersonen stammen aus Wil (Kanton St. Gallen), einer Kleinstadt mit 24 000 Einwohnern, und scheinen, so nach erster Einschätzung der Interviewerinnen, die selbst in Wil aufgewachsen sind, den lokalen Dialekt zu sprechen. Von Anfang an war klar, dass es nicht möglich sein würde, absolut ‚reine‘ DialektsprecherInnen zu finden, d. h. Gewährspersonen, die ihr ganzes Leben in Wil verbracht haben und die ausschließlich Kontakt zu Personen, auf die diese Charakterisierung auch zutreffen würde, gepflegt haben. Deshalb war der ursprüngliche Plan, Gewährspersonen, die in Wil aufgewachsen sind und die mindestens einen Elternteil haben, der auch in Wil aufgewachsen ist, zu interviewen. Offensichtlich spielen jedoch mehr oder weniger intensive Kontakte zu Dialektsprechern aus anderen Dialektregionen eine gewisse Rolle und manche Personen scheinen empfänglicher für solche Einflüsse zu sein als andere, eventuell auch abhängig von der Lebensphase (z. B. Kindheit, Pubertät), zu welcher diese Kontakte stattgefunden haben. In einem Interview in der Anlaufzeit der Projektphase mit einem Ehepaar, das selbst und auch deren Eltern aus Wil stammen, wurde bereits deutlich, dass sich die beiden Ehepartner hin-

sichtlich der Aussprache gewisser Wörter unterscheiden und dass der Ehemann gewisse Wörter auf eine nicht-Wilerische Art ausspricht. Da unser Augenmerk jedoch auf die Syntax des Wilerischen gerichtet ist, die weniger von der Syntax anderer Dialekte beeinflusst zu werden scheint, haben wir unser Auswahlkriterium gelockert. Um an der Studie teilzunehmen muss eine Gewährsperson in Wil aufgewachsen sein. Wichtig war auch, dass die Gewährspersonen nicht alle einen ähnlichen Beruf ausüben (z. B. im kaufmännischen Bereich) und sich im schulischen Bildungsstand kaum unterscheiden (z. B. alle mit einem Fachhochschul- oder Universitätsabschluss). Nicht nur syntaktische Variation sondern auch das Phänomen, ob der Dialekt sich in den letzten Jahrzehnten eventuell verändert hat, soll untersucht werden. Deshalb wurden Gewährspersonen gesucht, die sich grob in 3 Altersgruppen einteilen lassen: junge Gewährspersonen (20–30), Gewährspersonen mittleren Alters (45–55) und ältere (70+). Ziel ist es, Gruppen von ca. 20 Personen mit einem gleichen Anteil an Männern und Frauen zu bilden. Einige wenige Personen wurden interviewt, obwohl sie in keine dieser Altersgruppen passen. Es gibt verschiedene Gründe dafür, z. B. weil wir Daten von einem Geschwisterpaar erheben wollten und nur eine der beiden Personen in eine der Altersgruppen gepasst hat oder weil eine Person uns andere Kontakte vermitteln konnte und selbst unbedingt auch an der Studie teilnehmen wollte. Die Gewährspersonen haben wir durch persönliche Kontakte (Bekannte aus der Schulzeit oder ehemalige Nachbarn), durch die Zeitung und eben auch durch Gewährspersonen, die selbst andere für die Teilnahme an dieser Dialektstudie begeistern konnten, gefunden.

Normalerweise nimmt nur eine Gewährsperson an einem Interview teil, aber an 3 haben zwei Gewährspersonen gleichzeitig teilgenommen und zwar in zwei mit jeweils einem Ehepaar und in einem mit zwei Personen, die seit der Kindheit befreundet sind. Bislang wurden 58 Interviews durchgeführt, 3 von diesen werden eventuell nicht transkribiert werden, weil sich die Aussprache der Gewährsperson in z. T. häufig benutzten Wörtern zu sehr von der lokalen Varietät unterscheidet. Tabelle 1 gibt eine Übersicht über die Gewährspersonen, deren Inter-

views entweder bereits transkribiert wurden oder die noch zu transkribieren sind. Aufgrund des Fragebogens haben wir Informationen unter anderem zu folgenden Punkten erhalten: Herkunft (Stammt ein Elternteil oder stammen beide aus Wil oder nicht aus Wil?), Aufenthalt von mindestens einem Jahr in einem anderen Dialektgebiet der Schweiz (dt. CH) oder in einem Gebiet der Schweiz, in dem nicht schweizerdeutsch gesprochen wird oder im Ausland (\neq dt. CH), Schulbildung und berufliche Aus- oder Weiterbildung (Lehre oder Lehre und Weiterbildung an einer Fachhochschule, Lehrerseminar/Matura oder gleichwertiger Schulabschluss, Studium an einer Fachhochschule/Universität).

	20–30		45–55		70+	
	m	w	m	w	m	w
	8	10	8	8	9	7
Elternteil Wil	2	7	1	2	6	2
Eltern Wil	3	0	5	3	1	3
Eltern \neq Wil	3	3	2	3	2	2
dt. CH (1J+)	0	1	4	3	5	3
\neq dt. CH (1J+)	0	1	5	2	0	5
Lehre	3	4	5	3	4	6
Lehre+FH/Uni	3	0	0	0	0	0
Semi/Matura	0	1	0	4	1	1
FH/Uni	2	5	3	1	4	0

Tab. 1: Übersicht Gewährspersonen (nach Altersgruppen und Geschlecht)

Folgende 8 Gewährspersonen passen in keine dieser Altersgruppen: 1m (40–45), 2m (56–59), 1w (56–59), 3w (60–65), 1w (66–69). Die 3 Interviews, die eventuell nicht transkribiert werden, betreffen 1w (45–55), 1m (56–59) – passt zusätzlich auch nicht in eine der Altersgruppen – und 1w (70+).

2.2 Interviewablauf

Bei der ersten Kontaktaufnahme wurde die Gewährsperson über das Ziel des Projektes informiert sowie über den Ablauf des Gesprächs, z. B. dass es mit einem Aufnahmegerät aufgezeichnet würde, dass aber der Name sowie genaue Angaben zur Person nirgendwo auftreten würden. Interessierte Personen wurden dann gebeten, einen Informantenfragebogen auszufüllen entweder in elektronischer Form (als Anhang an eine e-mail) oder in Papierformat (vorher als Brief zugestellt) oder vor Interviewbeginn. In diesem Informantenfragebogen wurden Informationen zur Person (Name, Geburtsdatum), schulischer und beruflicher Werdegang (Lehre, Beruf), Aufenthalt in einem anderen Dialektgebiet oder im Ausland und Kenntnisse anderer Sprachen erfragt, sowie ob die Eltern aus Wil stammen. Je nach Wunsch fand das Gespräch entweder bei der Interviewerin oder bei der Gewährsperson oder deren Arbeitsplatz statt. Vor Interviewbeginn wurde noch einmal betont, dass das Gespräch anonymisiert würde.

Die Interviews wurden von Interviewerinnen durchgeführt, die selbst den Wiler Dialekt sprechen, die meisten davon von zwei Interviewerinnen mittleren Alters und ein paar wenige von einer jungen Interviewerin. Im Allgemeinen dauerten die Gespräche um die 90 Minuten und Themen wie Kindheit, Beruf, Hobbys, Reisen und Leben in Wil wurden diskutiert. Alle Gespräche wurden mit einem Tascamgerät aufgezeichnet, dessen eingebautes Mikrofon von sehr guter Qualität ist. Wir haben uns gegen die Verwendung eines Knopflochmikrofons entschieden, da diese Art von Mikrofon nicht-Sprachsignale teils amplifizieren (z. B. Atmen, Schlucken, Lachen) und das Gefühl einer Interviewsituation noch verstärken kann. Auch das eingebaute Mikrofon des Tascamgeräts nimmt Nebengeräusche auf, z. B. Schnalzlauten, in die Hände klatschen oder das mit-den-Fingern-auf-den-Tisch-Trommeln. Obwohl solche Nebengeräusche als störend erscheinen mögen, können sie etwas über das Befinden der Gesprächsteilnehmer aussagen (z. B. Begeisterung, Nervosität). Bei zwei Gesprächen gab es Interferenzerscheinungen (ein irritierendes elektronisches Nebengeräusch), welche eventuell durch ein eingeschaltetes Smartphone hervorgerufen wurden. Die Teilnehmer

wurden deshalb jeweils gebeten, ihr Handy vor Gesprächsbeginn in einem anderen Raum zu platzieren.

Nach Gesprächsende wurde die Gewährsperson noch einmal auf die wesentlichen Punkte der Dialektstudie hingewiesen und danach gebeten, eine schriftliche Einwilligung zu unterschreiben mit der Wahl zwischen „Einwilligung zur Verschriftlichung der Aufnahme“ und „Einwilligung zur Verschriftlichung der Aufnahme und zur Freistellung der Tonaufnahme“ in eine von der Universität Genf verwalteten Datenbank. Für den Gesamtzeitaufwand von jeweils 2–3 Stunden erhielten die Gewährspersonen eine kleine Vergütung.

2.3 Transkription

Im Gegensatz zum Deutschen gibt es keine Rechtschreiberegeln für das Schweizerdeutsche. Die Transkription des Schweizerdeutschen orientiert sich deshalb an den allgemeinen Richtlinien der Dialektschreibung in Dieth (1986). Es wurde eine sehr ‚weite‘ Dialektschrift gewählt, da Details in der Aussprache für unsere Untersuchung nicht relevant sind und eine weite Dialektschrift im Gegensatz zu einer engen einerseits einfacher anzuwenden und andererseits auch weniger zeitaufwändig ist. Die Variation in der Aussprache gewisser, vor allem hochfrequenter Wörter, kann groß sein, sowohl zwischen Sprechern als auch innerhalb desselben Sprechers. Diese Variation haben wir versucht in der Transkription ein wenig zu widerspiegeln: das Pendant zu deutschem *ein wenig* wird beispielsweise als *achli*, *ächli*, *echli*, *chli* und manchmal auch als *òchli* ausgesprochen. Die Entscheidung, ob in einem konkreten Fall eher „echli“ oder „ächli“ vorliegt, ist nicht einfach aber auch nicht relevant für unsere Untersuchungen. Zugegebenermaßen ist die Transkription gewisser Wörter unscharf aber sie ist trotzdem völlig ausreichend für syntaktische Studien. Hätten wir von Anfang an gewusst, in welchem Ausmaß die Aussprache variieren kann, hätten wir vermutlich ein Wort immer auf dieselbe Weise transkribiert und zwar so, wie es am häufigsten ausgesprochen wird, was wir aber zu dem Zeitpunkt noch nicht wissen konnten. Das übergeordnete Ziel der Verschriftli-

chung der Gespräche besteht darin, leicht leserlich zu sein. Es wurden deshalb fast keine Diakritika verwendet und die Satzzeichensetzung orientiert sich an derjenigen der deutschen Schriftsprache. SprachwissenschaftlerInnen, die sich z. B. für die Aussprache einzelner Wörter oder Intonation interessieren, können diese selbst überprüfen, da die Gespräche nicht nur in verschriftlichter Form vorliegen sondern auch als auditive Daten (.wav Dateien) zur Verfügung stehen werden.¹

Für die Verschriftlichung wurde EXMARaLDA, ein Programm zur Verschriftlichung gesprochener Sprache, das von Thomas Schmidt und Kai Wörner (2009) an der Universität Hamburg entwickelt wurde, verwendet. EXMARaLDA ist unter dem Link www.exmaralda.org frei verfügbar, läuft auf PCs, Macs und Linux und ist in seiner Anwendung einfach. Das Programm läuft einwandfrei auf Macs mit älteren Versionen des Betriebssystems OS X (z. B. Maverick), es können aber Probleme mit neueren Versionen (z. B. Yosemite, El Capitan oder Sierra) auftreten.

Im Partitur-Editor von EXMARaLDA können beliebig viele Spuren pro Sprecher angelegt werden. Wir verwenden jeweils 4 Spuren pro Sprecher: eine verbale Spur [v], eine suprasegmentale Spur [sup], eine nicht-verbale Spur [nv] und eine Kommentarspur [comment]. In der [v] Spur wird das Gesagte transkribiert. Die [sup] Spur enthält Kommentare zur Aussprache, z. B. genuschelt, lachend, nicht-Wilerisch. In der [nv] Spur werden nicht-verbale Handlungen kommentiert, z. B. klatscht in die Hände, lacht. In der [comment] Spur notieren wir, was uns auffällt und eventuell von Interesse sein könnte, z. B. V2 in w-Komplement, Apokoinu. Abb. 1 zeigt einen Ausschnitt einer transkribierten Passage im Partitur-Editor.

¹ Das Korpus soll nach Fertigstellung anderen SprachwissenschaftlerInnen auf Anfrage zur Verfügung gestellt werden.

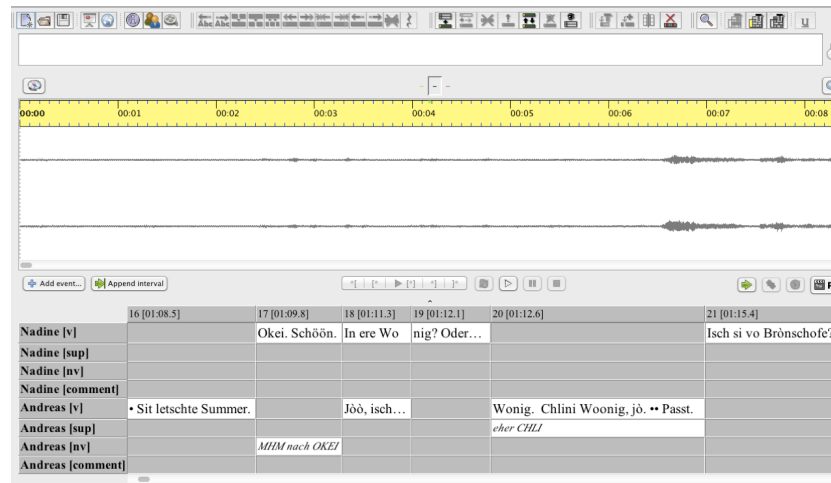


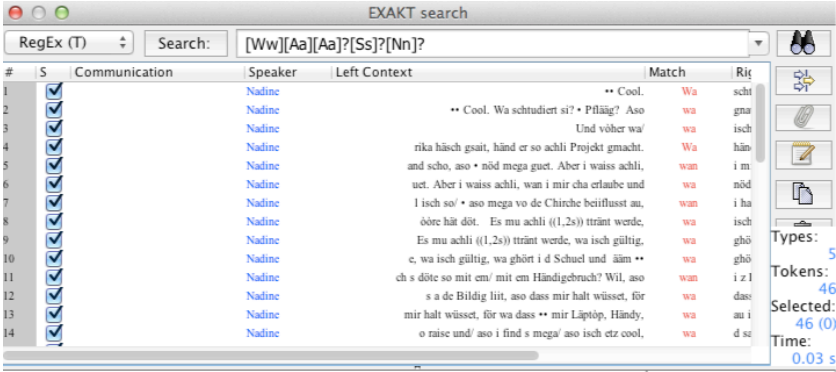
Abb. 1: Screenshot einer transkribierten Passage im Partitur-Editor

Wenn nur eine Person am Sprechen ist, kann ein Intervall gewählt werden, das alles umfasst, was diese – ohne unterbrochen zu werden – gesagt hat. Auch wenn dieser Redefluss mehr als 10 Sekunden übersteigt, haben wir normalerweise keine Intervalle gesetzt, die mehr als 10 Sekunden umfassen (was ungefähr der Bildschirmbreite eines Laptops entspricht) und die Intervallgrenze so gesetzt, dass sie entweder mit dem Satzende oder einer Pause zusammenfällt. Da der Redefluss oft durch kurzes Atmen unterbrochen wird, ist es einfach, bei einer solchen Atempause eine Intervallgrenze zu setzen. Oft fällt diese nicht mit dem Satzende zusammen, was vermutlich damit zu tun hat, dass der Sprecher signalisieren möchte, dass er noch nicht zu Ende gesprochen hat. Oft folgt nämlich am Ende einer Äußerung ein *und* oder manchmal auch ein nicht-kausales *wil* ‘weil’ und danach eine kurze oder längere Pause. Bei simultanem Sprechen muss dann ein neues Intervall begonnen werden. Gesprächsteilnehmer tendieren dazu, das vom Gesprächspartner Gesagte in einer gewissen Weise kurz zu kommentieren, z. B. Zustimmung (*jò*, *mhm*), Begeisterung (*waa*, *wau*, *super*), Skepsis (*eerlech?*, *scho?*). Es wäre nicht sinnvoll, diese Kurzkommentare jeweils

als gleichzeitiges Sprechen in separaten Intervallen festzuhalten. Solche stereotypen Kurzkommentare werden deshalb jeweils in der [nv] Spur kommentiert.

Basierend auf den Konventionen für EXMARaLDA werden stark betonte Silben mit Großbuchstaben geschrieben. Abbrüche einer Äußerung werden durch „...“ und *false starts* durch „/“ gekennzeichnet, wobei diese Unterscheidung in der Praxis nicht immer einfach ist. Kurze Pausen von 0.3–0.9 s werden durch einen bis drei Punkte (•) und längere durch eine Zeitangabe (z. B. 1,3s) signalisiert. Das Markieren einer Pause dient vor allem dazu, anzuzeigen, dass ein kurzer Unterbruch stattgefunden hat, der entweder von kurzer Dauer (meistens um Atem zu holen) oder länger (der Sprecher überlegt sich etwas) war.

Die Benutzung von EXMARaLDA bringt viele Vorteile mit sich. Folgende waren für uns besonders wichtig. In EXMARaLDA ist es möglich, nach bestimmten Mustern mit EXAKT search zu suchen. Mit dem in Abb. 2 gezeigten regulären Ausdruck kann gleichzeitig nach *WA*, *WAA*, *wa*, *waa*, *was*, *waas* und *wan*, die alle für deutsches ‘was’ stehen, gesucht werden.² Die Suchergebnisse können dann beispielsweise in eine Excel Tabelle zur weiteren Bearbeitung eingefügt werden.



The screenshot shows the EXAKT search window with the following data:

#	S	Communication	Speaker	Left Context	Match	Riq
1	<input checked="" type="checkbox"/>		Nadine		** Cool.	Wa scht
2	<input checked="" type="checkbox"/>		Nadine	** Cool. Wa schrudiert si? • Pfläg? Aso	wa gna	
3	<input checked="" type="checkbox"/>		Nadine	Und vöher wa/	wa isch	
4	<input checked="" type="checkbox"/>		Nadine	rika häisch gsait, händ er so achli Projekt gmacht.	Wa hän	
5	<input checked="" type="checkbox"/>		Nadine	and scho, aso • nöd mega guet. Aber i weiss achli,	wan i m	
6	<input checked="" type="checkbox"/>		Nadine	uet. Aber i weiss achli, wan i mir cha erlaube und	wa nöd	
7	<input checked="" type="checkbox"/>		Nadine	i isch so/ • aso mega vo de Chirche beiflusst au,	wan i ha	
8	<input checked="" type="checkbox"/>		Nadine	ööre hät döt. Es mu achli ((1,2s)) trünt werde,	wa isch	
9	<input checked="" type="checkbox"/>		Nadine	Es mu achli ((1,2s)) trünt werde, wa isch gültig.	wa ghö	
10	<input checked="" type="checkbox"/>		Nadine	e, wa isch gültig, wa ghört i d Schued und ääm **	wa ghö	
11	<input checked="" type="checkbox"/>		Nadine	ch s döte so mit em/ mit em Händgebruch? Wil, aso	wan i z l	
12	<input checked="" type="checkbox"/>		Nadine	s a de Bildig lit, aso dass mir halt wüsstet, für	wa dass	
13	<input checked="" type="checkbox"/>		Nadine	mir halt wüsstet, für wa dass ** mir Läptöp, Händy,	wa au i	
14	<input checked="" type="checkbox"/>		Nadine	o raise und/ aso i find s mega/ aso isch etz cool,	wa d se	

Summary statistics on the right side of the window:

- Types: 5
- Tokens: 46
- Selected: 46 (0)
- Time: 0.03 s

² Durch diesen regulären Ausdruck werden auch Strings, die *waan*, *wAn*, *wAan* usw. enthalten, gefunden. Falls kein Schreibfehler vorliegt, sollten jedoch keine Suchergebnisse mit solchen strings gefunden werden.

Abb. 2: Screenshot eines Suchmusters in EXAKT search

Auch lässt sich ein Transkript, das zugrundeliegend eine .txt Datei ist, im Partitur-Editor unter **File: print** als PDF Datei speichern, wobei einzelne Spuren vorher ausgeblendet werden können, z. B. die [comment] Spur, falls man diese in der PDF Datei nicht mitabbilden möchte (siehe Abb. 3).

		12 [01:01.9]			
Nadine [nv]		<i>AA SCHO mit SIT</i>			
Andreas [v]		• Immer achli dò obe, so Näugruebe und so. • Denn etz äbe z Brönschofe.			
[12]					
		13 [01:04.7]		14 [01:06.1]	
Nadine [v]		Und etz woonsch älai?		Cool.	
Nadine [nv]				<i>OK</i>	
Andreas [v]		Etz bin i jò uszzòge sit...		Mit de Fründin.	
[13]					
		16 [01:08.5]		17 [01:09.8]	
Nadine [v]		Wie lang scho?		Okei. Schöön. In ere Wo nig?	
Andreas [v]		• Sit letschte Summer.		Jòò, isch...	
Andreas [nv]		<i>MHM nach OKEI</i>			
[14]					
		20 [01:12.6]		21 [01:15.4]	
Nadine [v]		Oder...		Isch si vo Brönschofe?	
Andreas [v]		Wonig. Chlini Woonig, jò. •• Passt.			
Andreas [sup]		<i>eher CHLI</i>			
[15]					

Abb. 3: Screenshot einer Transkriptpassage als PDF Datei

Der Befehl *word count*, der im Partitur-Editor unter **Transcription: Word list...** aufgerufen werden kann, ist auch sehr nützlich, denn mit diesem lässt sich einfach feststellen, wie viele Wörter und was für welche ein Sprecher produziert hat.

3. Grammatische Annotation

Nach der Verschriftlichung eines Gesprächs kann das Transkript als .txt Datei exportiert und grammatisch annotiert werden (3.1), bevor dieser grammatisch annotierten .txt Datei eine hierarchische Struktur verliehen wird (3.2). Ziel ist es, ein Korpus nach dem Muster jener an der Universität von Philadelphia kreierten zu erstellen.

3.1 Grammatische Annotation (*tagging*) in .txt Dateien

Wir haben zwei der mit Hilfe von EXMARaLDA verschriftlichten Gespräche als .txt Dateien exportiert und danach die einzelnen Wörter und Satzzeichen von Hand grammatisch annotiert (*getaggt*). Größtenteils haben wir die POS Tags übernommen, die für andere an der Universität von Philadelphia kreierten Korpora verwendet wurden³, aber auch einige für das Schweizerdeutsche spezifische Tags hinzugefügt, z. B. „DV“ für Verdoppelungsverben wie *go* in Beispiel (1).

- (1) *Si isch gange go schpile.*
sie ist gegangen gehen spielen
'Sie ist spielen gegangen.'

Die verwendeten POS Tags sind nicht sehr differenziert. So wird beispielsweise zwischen einem Personalpronomen und einem Possessivpronomen unterschieden: *du*/PRO vs. *mini*/PRO\$, aber das POS Tag enthält keine Angaben über Person, Numerus und Kasus des Pronomens.

Diese beiden manuell getaggtten Dateien enthielten ca. 40 000 Tags, die dann als Trainingskorpus für einen Tagger dienten, um 4 weitere .txt Dateien automatisch zu taggen. In einem Vergleich von drei verschiedenen Taggern erwies sich der B Tagger (<http://clcl.unige.ch/SOFTWARE.html>) als für unsere Zwecke am effizientesten.⁴ Die vom B Tagger annotierten Dateien wurden manuell

³ www.ling.upenn.edu/hist-corpora/annotation/index.html

⁴ Wir danken Yves Scherrer für die Hilfe bei der Auswahl der Tagger sowie die Outputs dieser Tagger in Tabellenformat.

korrigiert. Danach wurde in allen getaggen Dateien die *disfluencies*, worunter nicht nur Abbrüche sondern auch wortgetreue Wiederholungen, Einschübe und Elaborationen fallen, manuell gekennzeichnet und mit einem entsprechenden Tag versehen, z. B. ELAB/CODE und ELAB\$\$/CODE um den Anfang, respektive das Ende, einer Elaboration anzuzeigen. Das Markieren von *disfluencies* ist sehr zeitaufwendig, denn in vielen Fällen erfordert es, dass man noch einmal in die Audio-datei hineinhört, um eine Entscheidung zu fällen. Manchmal hilft auch dieses wiederholte Hineinhören nicht weiter, vor allem, wenn man eine eher verwirrte Äußerung wie in (2) nicht einfach als Abbruch einstufen möchte. Nebst Abbrüchen und Fragmenten, produziert der Sprecher *klaare Sune* ‘klarer Sonne’ und nicht *klaari Sune* ‘klare Sonne’, was vermutlich kein Genusfehler sondern eher ein Abbruch von *klaare Sunetag* ‘klarer Sonnentag’ ist, eventuell intendiert als ‘klarer Himmel’.

(2) • *Und die händ • a däm/ ((1,3s)) a däm Taag, wo meer um die Insle umegfaare sind, isch en Taag •• braandhaiss • klaare Sune, also isch würkli en/ en waansinig schööne • Sunetaag gsi.* (Leonard, 45–55)

und die haben an dem an dem Tag wo wir um die Insel herumgefahren sind ist ein Tag brandheiss klarer Sonne also ist wirklich ein ein wahnsinnig schöner Sonnentag gewesen

Gleichzeitig wurden in diesen Dateien Äußerungen, die nicht aufgrund der Satzzeichensetzung automatisch in satzähnliche Tokens aufgegliedert wurden, manuell nachgebessert. Die Aufbereitung in satzähnliche Tokens und Fragmente vereinfacht die Arbeit beim Parsen.

3.2 Erste geparste Daten

Nach dem Taggen und Markieren von Disfluencies wurden diese 6 individuellen Dateien zu einer Datei zusammengefügt, die dann insgesamt aus ca. 110 000 Tags bestand. Beatrice Santorini, die unter anderem bereits das Parsen eines sehr ähnlichen Korpus unternommen hat,

nämlich jenes von Appalachian English von Christina Tortora (<http://csivc.csi.cuny.edu/aapcapp/>), das auch aus Spontansprachdaten besteht, hat dann unsere schweizerdeutschen Daten geparst. Es wird bei dieser Art von Parsen nicht das Ziel verfolgt, eine möglichst detaillierte syntaktische Analyse einzelner Äußerungen zu liefern, sondern sie so zu parsen, dass danach Recherchen nach gewissen Strukturen und Häufigkeit ihres Vorkommens sehr schnell durchgeführt werden können, z. B. mit Hilfe von CorpusSearch (<http://corpussearch.sourceforge.net/>). Abb. 4 zeigt eine geparste Äußerung.

```
Mini Töchter isch letschthin cho: "Du, wörsch die nöd mit der nee?"

(IP-MAT (NP-SBJ (PRO$ Mini) (N Töchter))
  (BEP isch)
  (ADVP (ADV letschthin))
  (VBN cho)
  (. :)
  (" "))
(CP-QUE-SPE
  (NP-VOC (PRO Du))
  (.....)
  (IP-QUE-SPE (RDKII wörsch)
    (NP-SBJ *pro*)
    (NP-OB (D die))
    (NEG nöd)
    (PP (P mit)
      (NP (PRO der)))
    (VB nee)
    (.....?)))
  (" "))
```

Abb. 4: Screenshot einer geparsten Äußerung

Die restlichen Daten werden nun nach derselben Methode getaggt und geparst.

4. Erste Ergebnisse

Basierend auf verfügbaren transkribierten Daten wurden erste Hypothesen zur Variation im Verbalkomplex (4.1), Verbstellung in Sätzen, die durch *wil* 'weil' eingeleitet werden (4.2), Gebrauch von *doubly-filled*

Comps in w-Komplementen (4.3) sowie Art und Häufigkeit von Apokoinus (4.4) aufgestellt. Weil die meisten Daten noch nicht in gearter Form vorliegen, wurden Abfragen mit Hilfe von EXAKT search durchgeführt.

4.1 Variation im Verbalkomplex

SCHÖNENBERGER/HAEBERLI (2015) untersucht die Variation im Verbalkomplex in den verschriftlichten Daten von 9 Gewährspersonen (7 mittleren Alters und 2 ältere) und zwar, was die Variation in der Abfolge der Verben als auch des nicht-verbalen Materials betrifft. Welche Abfolgen mit zwei beziehungsweise drei Verben im Nebensatz erlaubt sind, wird für den Wiler Dialekt in (3) und (4) gezeigt. Diese Akzeptabilitätseinstufung beruht auf einer informellen Umfrage mit nur wenigen Dialektsprechern aus Wil.

- (3)a. *dass er wött jòdle (V1 V2)*
 b. *dass er jòdle wött (V2 V1)*
 ‘dass er jodeln will’
- (4)a. *dass er hät wöle jòdle (V1 V2 V3)*
 b. *dass er hät jòdle wöle (V1 V3 V2)*
 c.? *dass er jòdle hät wöle (V3 V1 V2)*
 d.?? *dass er jòdle wöle hät (V3 V2 V1)*
 e.* *dass er wöle hät jòdle (V2 V1 V3)*
 f.* *dass er wöle jòdle hät (V2 V3 V1)*
 ‘dass er jodeln hat wollen’

Obwohl mehrere Abfolgen als akzeptabel betrachtet werden, haben wir in den Spontansprachdaten fast ausschließlich die aufsteigende Abfolge gefunden: V1 V2 wie in (3a) und V1 V2 V3 wie in (4a). Diese wird in 531 von 532 Beispielen produziert.

Die Abfolge von Verben im Verbalkomplex kann durch nicht-verbales Material unterbrochen werden. Im Allgemeinen können Kom-

plemente und Adjunkte des lexikalischen Verbs in einem Verbalkomplex auftreten, was auch als Verbprojektionsanhebung bezeichnet wird (HAEGEMAN/VAN RIEMSDIJK 1986). In (5) taucht ein Adverbial (*am Wuchenend*) und ein direktes Objekt (*e Rüebliort*) zwischen dem flektierten Modalverb und dem lexikalischen Verb auf.

- (5) *dass er ebe immer mue am Wuchenend e Rüebliort mache*
 dass er eben immer muss am Wochenende eine Karottentorte machen

In den Daten der 9 Gewährspersonen zeichnen sich gewisse Tendenzen ab. Verbpartikeln, Prädikate von *si* 'sein' und PP-Komplemente werden fast immer und nicht-pronominale Objekte und PP-Adjunkte werden meistens angehoben, d. h. sie treten im Verbalkomplex auf. Im Gegensatz dazu werden Adverbien und DP-Adjunkte eher selten angehoben. Negation und Diskurspartikeln werden praktisch nie und Klitika gar nie angehoben.

4.2 Verbstellung in adverbialen Kausalsätzen mit *wil*

Wie im gesprochenen Deutschen lässt ein Satz, der durch *wil* eingeleitet wird, prinzipiell die Verbzweitstellung (V2) als auch die Verbendstellung (VE) zu. Die Wahl der Verbstellung ist nicht frei, denn sie geht mit einem Bedeutungsunterschied einher. Beispiel (6a) mit *wil+VE* bedeutet, dass der Grund für Rochus Kopfschmerzen darin zu finden ist, dass er zu viel getrunken hat. Die Bedeutung von (6b) mit *wil+V2* kann wie folgt umschrieben werden: die Tatsache, dass Rochus Fieber hat, ist der Grund dafür, dass der Sprecher sagt, dass Rochus krank ist.

- (6)a. *De Rochus hät Chöpfwee, wil er z tüüf is Glaas glueget hät.*
 der Rochus hat Kopfschmerzen weil er zu tief ins Glas geschaut hat
- b. *De Rochus isch chränk, wil er hät Fieber.*
 der Rochus ist krank weil er hat Fieber

Anders ausgedrückt liefert *wil*+VE den Grund für einen Sachverhalt, der im Hauptsatz beschrieben wird, wohingegen *wil*+V2 normalerweise eine Begründung beinhaltet, weshalb eine Äußerung gemacht wird. Ersteres wird oft als faktisches *weil* und Letzteres als epistemisches *weil* bezeichnet (siehe u. a. WEGENER 1993, GÜNTNER 1996, UHMANN 1998).

Eine Auswertung der verschriftlichten Daten von 23 Gewährspersonen (2 junge, 14 mittleren Alters und 7 ältere) in SCHÖNENBERGER (i. E. a) zeigt, dass die meisten sowohl die VE als auch die V2 Stellung in *wil*-Sätzen verwenden. Zudem scheint die Verbstellung mit dem oben erläuterten Unterschied in der Bedeutung Hand in Hand zu gehen. Es kamen keine Beispiele mit epistemischem *weil* aber mit VE-Stellung wie in (7) vor.⁵

- (7)a. *Er hat zu viel getrunken, weil er gar so daherdefert* [dumm daherredet]. (Schierling/Landkreis Mallersdorf)
 b. *Hot dei Frau a Stoffwechselkrankheit, weils jedn Tog a anders Gwand ohot?* (Altbayerische Heimatpost 2007, Nr. 15, S.24)

Ältere Sprecher benutzten *wil*+V2 nur in etwa 27 % (22/83) der Fälle und deutlich weniger oft als jüngere Sprecher. Vor allem wurde *wil*+V2 von den Sprechern mittleren Alters gebraucht (73 %, 539/740). Bei den jungen Sprechern liegt der Prozentsatz mit 48 % (39/81) deutlich tiefer. Die Anzahl der Sprecher pro Altersgruppe ist jedoch sehr unausgewogen. Interessanterweise gab es eine ältere Sprecherin, die beide Verbstellungen nicht nur mit *wil* sondern auch mit *wäge* 'wegen' gebrauchte wie in (8a) und (8b). Wie im Deutschen wird *wäge* im Wiler Dialekt normalerweise ausschließlich mit nominalem Objekt wie in (8c) verwendet.

⁵ Wir danken Michael Schnabel (p.c. 3.11.2016) für den Hinweis, dass es sehr wohl solche Beispiele in der gesprochenen Sprache gibt sowie für einige Belege dafür (wie auch jene in [7]).

- (8)a. *Ich waiss es nu, wäg die Zäddel nò dine sind.*
ich weiss es nur wegen diese Zettel noch drinnen sind
- b. *Wo min Vatter gschtòrben isch, sind s mit zwee Wäge cho, wäg er hät so vil Chränz gha.*
als mein Vater gestorben ist sind sie mit zwei Wagen gekommen wegen er hat so viele Kränze gehabt
- c. *Aber er hät denn nöd ghüròòte scho wäg em Rutli nöd.*
aber er hat dann nicht geheiratet schon wegen dem Rutli nicht
(alle Beispiele von Hedda, 70+)

In den meisten Beispielen mit *wil+V2* – in 63 % (381/600) – wird die Vorfeldposition von einem Subjekt besetzt. Die Gewährspersonen nutzen aber auch die Möglichkeit, diese Position anders zu füllen. Die meisten Nicht-Subjekte, welche die Vorfeldposition einnehmen, sind Adjunkte (27 %, 161/600), und unter den Adjunkten sind dies meistens kurze, eher farblose Adverbien wie *susch/süsch/süs* ‘sonst’, *denn* ‘denn’, *etz/etzt/jetz* ‘jetzt’, *dòò* ‘dort’.

4.3 Doubly-filled Comps

In verschiedenen süddeutschen Dialekten aber nicht nur dort kann eine W-Konstituente in einem eingebetteten Satz zusammen mit *dass* auftreten (siehe u. a. BAYER/BRANDNER 2008 für Bodenseealemannisch und Bayerisch). Dieses Phänomen wurde in SCHÖNENBERGER (i. E. b) in den Daten von 35 Gewährspersonen (7 junge, 17 mittleren Alters und 11 ältere) untersucht. Wie im Bodenseealemannischen und Bayerischen kamen einsilbige W-Konstituenten im Gegensatz zu nicht-einsilbigen W-Konstituenten fast nie mit *dass* vor, wie die Beispiele in (9) zeigen. Weniger als 3 % (24/825) aller W-Komplemente mit einsilbiger W-Konstituente treten mit *dass* auf, wohingegen über 90 % (218/241) aller W-Komplemente mit nicht-einsilbiger W-Konstituente dies tun.

- (9) a. *I waiss nüme, wie si hässt mit Nòòchname.*
ich weiss nicht wie sie heisst mit Nachnamen

(Anna, 45–55)

- b. *I waiiss nöd, wie alt dass si isch.*
 ich weiss nicht wie alt dass sie ist
 (Thea, 45–55)

Es wird die Hypothese vertreten, dass das Einfügen von *dass* die Prosodie des W-Komplementes ‚positiv‘ beeinflussen kann, in dem es in einem Kontext entweder betont oder unbetont in Bezug auf die es umgebenden Silben sein kann, und so zu einem trochäischen Muster beitragen kann.

4.4 Apokoinus

SCHÖNENBERGER/HAEBERLI (eingereicht) befasst sich mit Apokoinus, einer Konstruktion, die den Eindruck erweckt, dass ein Sprecher zwei Sätze miteinander vermischt hat. Was Apokoinus auszeichnet ist, dass eine Konstituente, das Koinon, als zwei Sätzen zugehörig interpretiert wird wie *wacher* in Beispiel (10).

- (10) *Denn bisch es Zitli lang wie „wacher“ isch vilicht en Usdruck.*
 dann bist(du) eine Zeit lang wie wacher ist vielleicht ein Ausdruck
 (Leo, 45–55)

Die Analyse von MEINUNGER (2011) Apokoinus im Deutschen orientiert sich an einem Vorschlag von VAN RIEMSDIJK (2006), die dieser für ein anderes Phänomen entwickelt hat. Nach diesem Ansatz lassen sich Apokoinus als veredelte Strukturbäume (*grafted trees*) abbilden (siehe Abb. 5).

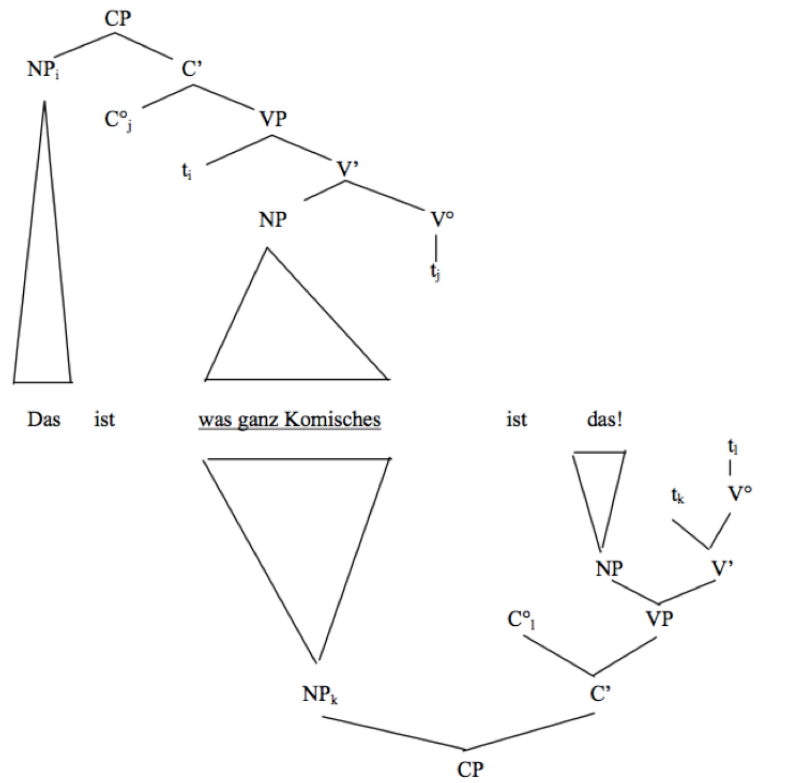


Abb. 5: Veredelter Strukturbaum (aus MEINUNGER 2011)

In den 6 Transkripten, ca. 100 000 Wörter, die bereits in grammatisch annotierter und gearparter Form vorliegen, gab es 187 Belege von Apokoinus, die sich grundsätzlich in zwei Typen aufgliedern lassen. In Beispiel (11a) – Typ A – ist das Koinon (*übrigens übermòrn*) die letzte Konstituente des ersten Satzes und zugleich die erste Konstituente des zweiten Satzes. In Beispiel (11b) – Typ B – werden Konstituenten vor dem Koinon (*Handòrgle*) danach wiederholt, (in (11b) fett gedruckt). Es liegen wenige Beispiele von Typ A aber viele von Typ B vor.

- (11) a. Typ A (6 Beispiele)
Dä gseen i übrìgens übermòrn gang i mit em go ässe.
 den sehe ich übrìgens übermorgen gehe ich mit ihm (gehen)
 essen
 (Otto, 70+)
- b. Typ B (181 Beispiele)
Er hüt irgèndwie Handòrgle hüt er organisiert.
 er hat irgèndwie Handorgel hat er organisiert
 (Leo, 45–55)

Diese Belege stammen von 8 Gewährspersonen (7 mittleren Alters und eine ältere). Alle bisher interviewten Gewährspersonen scheinen Apokoinus zu produzieren. Die Beliebtheit dieser Konstruktion könnte daran liegen, dass die Verwendung eines Apokoinus dem Sprecher die Möglichkeit eröffnet, Einschübe und Erklärungen, die nicht im voraus geplant waren, einzubauen, ohne die Sprachverarbeitung des Zuhörers zu erschweren (siehe z. B. AUER 2009).

Obwohl die meisten Apokoinus in den schweizerdeutschen Daten nach Meiningers Ansatz analysiert werden können, erscheint dieser Ansatz als zu permissiv, denn es ist nicht klar, an welchen Stellen und wie oft ein ursprünglicher Strukturbaum veredelt werden darf.

5. Schlussfolgerungen

Diese Art der Datenerhebung ist zeitaufwändig, vor allem was die Transkription und die Aufbereitung der grammatisch annotierten Dateien mit Markierungen von *disfluencies* betrifft. Gewisse Phänomene wie Artikelverdoppelung (*da isch de vil de besser Chòch* ‘das ist (der) viel der bessere Koch’) können basierend auf diesen Spontansprachdaten nicht im Detail untersucht werden, da Kontexte, in denen Artikelverdoppelung überhaupt möglich ist, nur sehr selten vorkommen. Andere Phänomene wie jene, die in Abschnitt 4 diskutiert wurden, eignen sich sehr wohl für eine genauere Untersuchung, da sie relativ häufig auftreten und von allen Gewährspersonen verwendet werden.

Auch sollte es möglich sein, potentiellen Sprachwandel in solchen Spontansprachdaten, die von Gewährspersonen verschiedenen Alters stammen – die jüngste 21 Jahre und die älteste 90 Jahre alt – zu untersuchen. Obwohl wir noch keine Studie zur Morphologie gemacht haben, scheint sich in diesem Bereich ein gewisser Wandel abzuzeichnen, was Singular und Plural von Nomen betrifft. Die Singular- und Pluralform von vielen Nomen ist homophon, z. B. *Ross* (SG/PL) von ‘Pferd’ und *Fründ* (SG/PL) von ‘Freund’. Die Tendenz, Plural overt zu markieren, ist vor allem bei jungen Sprechern auffällig (*Rösser*; *Fründe*). Einige wenige Nomen im Wiler Dialekt weisen eine Singularform auf, die wie eine Pluralform aussieht wie *Aier* SG/PL von ‘Ei’ und *Töchter* SG/PL von ‘Tochter’. Diese Formen werden zwar zum Teil noch von älteren Gewährspersonen und solchen mittleren Alters verwendet, sind aber bei jungen Gewährspersonen nicht mehr anzutreffen. Letztere verwenden im Singular *Ai* und *Tochter*.

Wir sind überzeugt, dass sich diese Art der Datenerhebung und Aufbereitung lohnt, obwohl sie zeitaufwändig ist, und hoffen, dass dieses Korpus nicht nur für uns sondern auch für andere SprachwissenschaftlerInnen von Interesse sein wird, auch für solche, die sich mit anderen linguistischen Bereichen wie Phonologie oder Soziolinguistik befassen.

Literaturangaben

- AUER, PETER. (2009): On-line syntax: Thoughts on the temporality of spoken language. In: *Language Sciences* 31, 1–13.
- BAYER, JOSEF/BRANDNER, ELLEN (2008): Wie oberflächlich ist die syntaktische Variation zwischen Dialekten? Doubly-filled COMP revisited. In: PATOCKA, FRANZ/SEILER, GUIDO (Hrsg.): *Dialektale Morphologie, dialektale Syntax*. Vienna: Praesens, 9–26.
- DIETH, EUGEN (1986): Schwyzertütschi Dialäktschrift. Dieth-Schreibung. In: SCHMID-CADALBERT, CHRISTIAN (Hg.): *Lebendige Mundart*. Band 1. Aarau/Frankfurt am Main: Verlag Sauerländer.
- GÜNTNER, SUSANNE (1996): From subordination to coordination? Verb-Second position in German causal and concessive constructions. In: *Pragmatics* 6:3, 323–356.
- HAEGEMAN, LILIANE/VAN RIEMSDIJK, HENK (1986): Verb Projection Raising, scope, and the typology of rules affecting verbs. In: *Linguistic Inquiry* 17, 417–466.
- MEINUNGER, ANDRÉ (2011): Das ist was ziemlich Komisches ist das! The syntax of apokoinu-constructions in colloquial German and other languages. In: BREINDL, EVA/FERRARESI, GISELLA/VOLODINA ANNA (Hrsg.): *Satzverknüpfung – Zur Interaktion von Form, Bedeutung und Diskursfunktion*. Berlin: Mouton de Gruyter, 351–378.
- SCHMIDT, THOMAS/WÖRNER, KAI (2009): EXMARaLDA – Creating, analysing and sharing spoken language corpora for pragmatic research. In: *Pragmatics* 19.4, 565–582.
- SCHÖNENBERGER, MANUELA (i. E. a): Verbstellung in *weil*-Sätzen des Schweizerdeutschen. In: NEFEDOV, SERGEJ/GRIGORIEVA, LJUBOV/BOCK, BETTINA (Hrsg.): *Deutsch als Bindeglied zwischen Inlands- und Auslandsgermanistik*. Beiträge zu den 23. GeSuS-Linguistik-Tagen in Sankt Petersburg, 22.-24. Juni 2015. Hamburg: Verlag Dr. Kovač, 395–404.
- SCHÖNENBERGER, MANUELA (i. E. b): Are doubly-filled COMPs governed by prosody in Swiss German? The chameleonic nature of *dass* ‘that’. In: ENOCH, ABOH/HAEBERLI, ERIC/PUSKAS, GENOVEVA/SCHÖNENBERGER, MANUELA (Hrsg.): *Elements of Comparative Syntax: Theory and Description*. Berlin: Mouton de Gruyter.
- SCHÖNENBERGER, MANUELA /HAEBERLI, ERIC (2015): Studie zur Sprachvariation im Schweizerdeutschen: Erste Ergebnisse. In: KLAUSMANN, HUBERT (Hg.): *Alemannentagung 2014: Sprache und Öffentlichkeit*. Universität Tübingen, Tübingen, 1–13.
- SCHÖNENBERGER, MANUELA /HAEBERLI, ERIC (eingereicht): Aso du chasch nõcher chasch du überaal mitrede – und andere Apokoinus in Spontansprachdaten des Schweizerdeutschen. Tagungsband der GeSuS Tagung in Brno.

- UHMANN, SUSANNE (1998): Verbstellungsvarianten in *weil*-Sätzen: Lexikalische Differenzierung mit grammatischen Folgen. In: Zeitschrift für Sprachwissenschaft 17.1, 92–139.
- VAN RIEMSDIJK, HENK (2006): Grafts from Merge. In: FRASCARELLI, MARA (Hg.): Phases of interpretation. Berlin: Mouton de Gruyter, 17–44.
- WEGENER, HEIDE (1993): Weil – das hat schon seinen Grund. Zur Verbstellung in Kausalsätzen mit *weil* im gegenwärtigen Deutsch. In: Deutsche Sprache 4, 289–305.