

Zielhypothesen und Lernersprache

—

Das Falko-Lernerkorpus

Marc Reznicek

22.06.12

offene Berliner DaF-Reihe

Mit Folien des gesamten
Korpuslinguistikteams der HU-Berlin

Übersicht

- Fragestellungen der Lernerkorpusforschung
- Falko - Design
- Annotation
 - Automatische Vorverarbeitung
 - Fehler und Zielhypothesen
- Fragestellungen und Analysen
 - Lernaltersyntax
 - Produktivität

Fragestellungen

- Welche sprachlichen Strukturen sind für Deutschlerner schwierig?
- Sind bestimmte Fehler L1-abhängig?
- Liegt das an der Form oder an der Funktion dieser Strukturen?
- Wie produktiv sind Fremdsprachenlerner im Sprachgebrauch?

Methoden

Contrastive Interlanguage Analysis (CIA: Granger 2008)

- Muster in abstrakten Repräsentationen finden
- Quantitative Unterschiede zwischen Lernern und Muttersprachlern aufdecken
- L2-Texte unterschiedlicher Muttersprachler kontrastieren

Fehleranalyse (EA) (Corder 1981, Izumi et al. 2005 u.v.m.)

- Welche Fehler sind lernertypisch
- Abhängigkeit von L1 der Lerner?

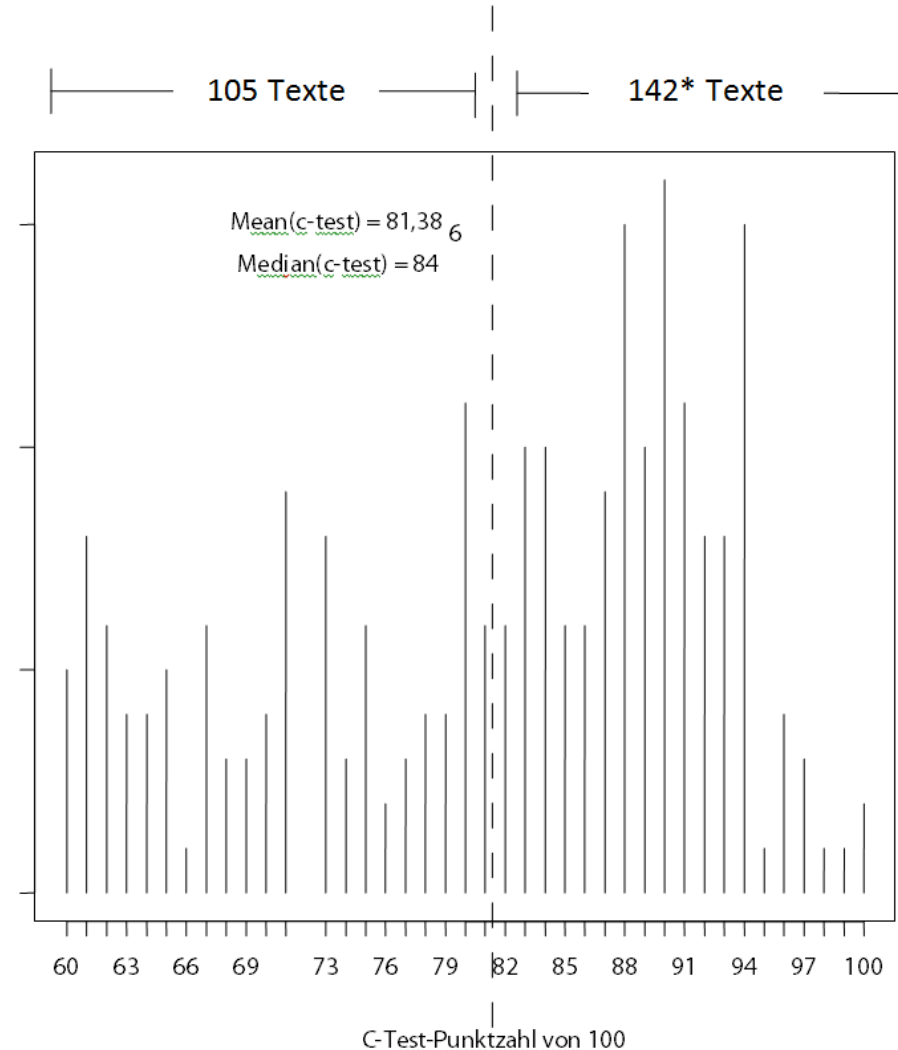
Lernerkorpora



- Kontrollierte und meist digitalisierte Sammlungen von Lernertexten
- Das Design hängt von der Fragestellung ab
 - gesprochen vs. geschrieben / Textsorte / Aufgabe
 - Fortgeschrittenheit
 - L1 der Lerner
 - ...
- die meisten Korpora auf Englisch
 - immer mehr Lernerkorpora auch für andere Sprachen

Granger/Hung/Petch-Tyson (2002), Cobb (2003), Tono (2003), Myles/Mitchell (2004), Nesselhauf (2004), Tenfjord/Meurer/Hofland (2004), Granger (2008), Lüdeling/Walter (2009) etc.

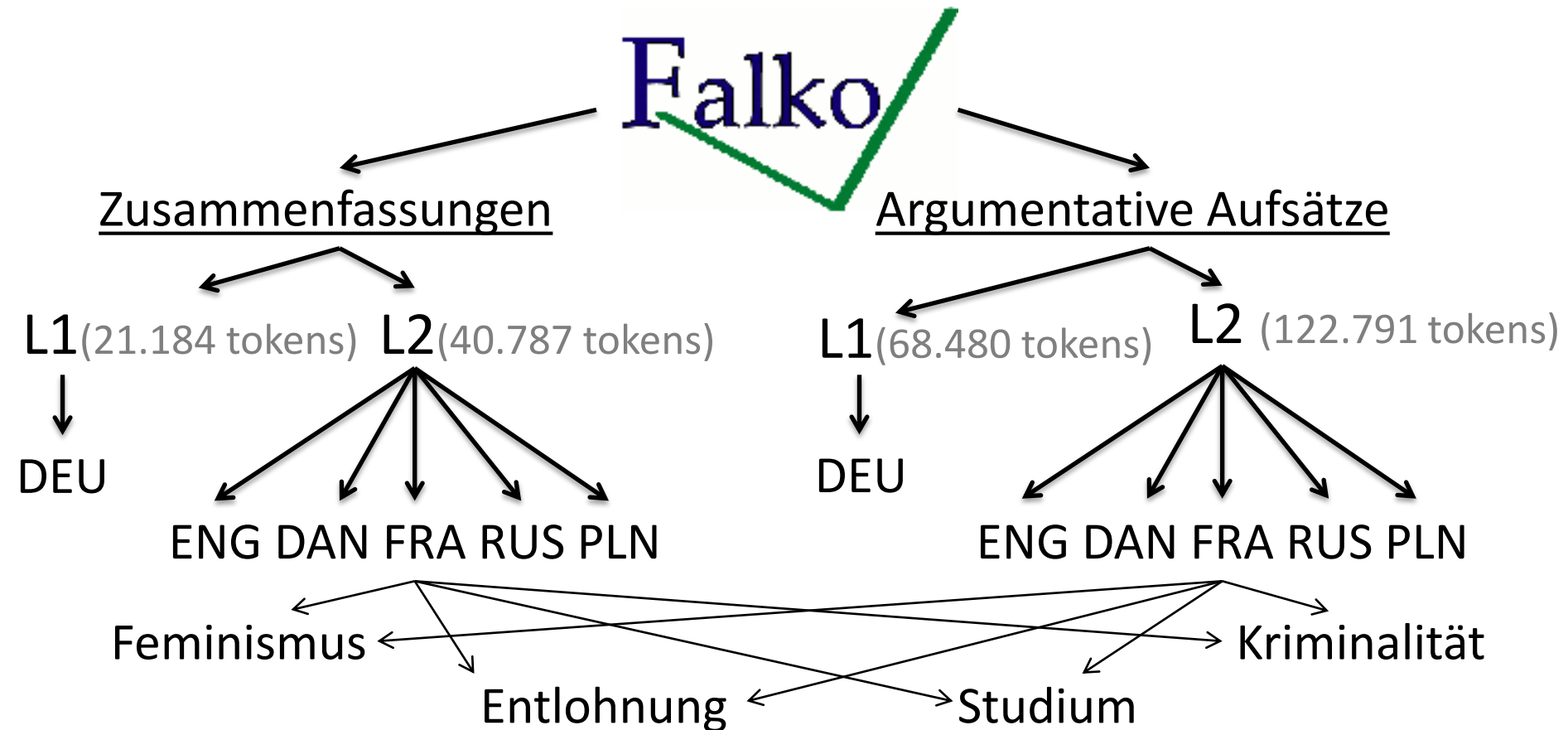
- frei verfügbares Lerner-
korpus des Deutschen als
Fremdsprache
- Fortgeschrittene Lerner
(DHS oder C-Test: HU-
Sprachenzentrum)
~B1-C1 (CEFR)



Falko



- Subkorpora 2 x 2 x 4 (16)



Metadaten



- Erhebungen
 - 90 Minuten
 - keine Hilfe (Internet, Wörterbücher, Spell-checker etc.)
 - handschriftlich (nur Zusammenfassungen) & computergeschrieben
- Dokumentation (Reznicek et al. 2010)
http://www.linguistik.hu-berlin.de/institut/professuren/korpuslinguistik/forschung/falko/Falko-Handbuch_Korpusaufbau%20und%20Annotationen_v1.0.1
- Projektseite:
<http://www.linguistik.hu-berlin.de/institut/professuren/korpuslinguistik/forschung/falko>

Spracherwerbsbiographie



- Lerner mit 49 Muttersprachen
- Für jede Sprache
 - Dauer des Erwerbs
 - Unterricht
 - Dauer
 - Auslandsaufenthalte
 - Beherrschungsgrad

Falko Welche Sprachen können sie?
moderne Sprachen und Muttersprachen

REIHENFOLGE = Grad der Beherrschung 1. beste 2. zweitbeste etc. <small>Bitte geben Sie den Namen der Sprache an.</small>	Ab welchem Alter haben Sie diese Sprache gebraucht? „seit der Geburt“ = „0“	Ist diese Sprache für Sie eine Muttersprache? <input type="checkbox"/> Ja <input type="checkbox"/> Nein	Wurde Ihnen die Sprache jemals unterrichtet? <input type="checkbox"/> Ja <input type="checkbox"/> Nein	Wenn ja, wie lange haben Sie darin Unterricht erhalten? (Jahre: Monate) z.B. 3 Jahre und 3 Monate = 3:3	Wenn ja, wo fand der Unterricht statt: Schule= SH Universität = UV Sprachschule=SP mehrere Kreuze möglich! <small>siehe Anmerkung(!)</small>
* z.B.: Englisch	ab 0 Jahren	<input checked="" type="checkbox"/>	<input type="checkbox"/>	13 Ja:Mo	SH <input checked="" type="checkbox"/> UV <input type="checkbox"/> SP <input type="checkbox"/>
* z.B.: Deutsch	ab 15 Jahren	<input type="checkbox"/>	<input checked="" type="checkbox"/>	5:5 Ja:Mo	SH <input checked="" type="checkbox"/> UV <input checked="" type="checkbox"/> SP <input type="checkbox"/>
1	ab _____ Jahren	<input type="checkbox"/>	<input type="checkbox"/>	Ja:Mo	SH <input type="checkbox"/> UV <input type="checkbox"/> SP <input type="checkbox"/>
2	ab _____ Jahren	<input type="checkbox"/>	<input type="checkbox"/>	Ja:Mo	SH <input type="checkbox"/> UV <input type="checkbox"/> SP <input type="checkbox"/>
3	ab _____ Jahren	<input type="checkbox"/>	<input type="checkbox"/>	Ja:Mo	SH <input type="checkbox"/> UV <input type="checkbox"/> SP <input type="checkbox"/>
4	ab _____ Jahren	<input type="checkbox"/>	<input type="checkbox"/>	Ja:Mo	SH <input type="checkbox"/> UV <input type="checkbox"/> SP <input type="checkbox"/>
5	ab _____ Jahren	<input type="checkbox"/>	<input type="checkbox"/>	Ja:Mo	SH <input type="checkbox"/> UV <input type="checkbox"/> SP <input type="checkbox"/>
6	ab _____ Jahren	<input type="checkbox"/>	<input type="checkbox"/>	Ja:Mo	SH <input type="checkbox"/> UV <input type="checkbox"/> SP <input type="checkbox"/>

WHIG – What is hard in German?

- Weitere Erhebungen englischer Deutschlernern
- ~120.000 tokens
- Datenerhebung und Aufbereitung nach Falko - Standard



Arts & Humanities
Research Council

HU-Berlin

Anke Lüdeling
Marc Reznicek

Bangor

University

Astrid Ensslin
Cedric Krummes

Annotation von Lernersprache

- Viele Lernerkorpora sind nicht annotiert
 - Manche haben Fehlerannotationen
 - Wenige (vor allem junge) Korpusprojekte integrieren weitere Annotationsebenen
(ALESKO, KOBALT, BEMATAC, DALEKO, KanDel)
 - Falko: Standoff-Format (jede Annotation wird unabhängig gespeichert)
- **Neue Annotationsebenen können unbeschränkt hinzugefügt werden.** (Lüdeling et al. 2005)

automatische Annotation – Wortarten/Lemma

- Meiste Lernerkorpusforschung konzentrierte sich auf Oberflächenformen (Möllering 2004, Vyatkina 2007 etc.)
- Aber interessanter z.B.:
 - Welche Wortarten und Wortartenketten werden vermieden? (Oberflächensyntax: Borin/Prütz 2004)
- Automatische Wortartentagger werden immer besser: Zeitungssprache ~98% (Kübler et al. 2010)

→ automatische Wortarten(POS)- und Lemma-Annotation

Fehlerannotation in Falko

An der anderen Seite, wenn da kein Feminismus wäre, stünden wir noch nur in der Küche und köchten wir. (fkb034_2008_07)

Fehlerannotationen beziehen sich immer auf eine (zumindest implizite) korrekte Entsprechung der Lerneräußerung. → **Zielhypothese**

Wie würden Sie korrigieren?

Zielhypothesen

An der anderen Seite, wenn da kein Feminismus wäre, stünden wir noch nur in der Küche und köchten wir.
(fkb034_2008_07)

Falko: explizite Zielhypothese

- oft konkurrierende Korrekturen möglich

ZH1: *An der anderen Seite, wenn da kein Feminismus wäre, stünden wir **nur noch** in der Küche und **köchten**.*

ZH2: ***Andererseits** stünden wir, wenn es keinen **Feminismus** gäbe, **nur noch** in der Küche und **köchten**.*

Zielhypothesen in Falko



- ZH1:** satzbasiert, nah an Lernerstruktur:
Orthographie, Morphosyntax
- ZH2:** text-basiert, nah an Lernerintention:
Semantik, Pragmatik, Stilistik
- Die Unterschiede zwischen ZH und Lernertext wird automatisch mit **edit tags** (CHAnge, INSert, DELeTe etc.) markiert
 - Alle ZHs werden mit POS-und Lemma annotiert (TreeTagger, Schmid 1994)
 - weitere **manuelle Fehlertags** für einige Phänomene

LT	TH1	TH1Diff	TH2	TH2Diff
An	Auf	CHA		
der	der		Andererseits	MERGE
anderen	anderen			
Seite	Seite			
,	,			
			stunden	MOVT
			wir	MOVT
			,	INS
wenn	wenn		wenn	
da	da			DEL
			es	INS
kein	kein		inen	CHA
Feminismus	Feminismus		Feminismus	
wäre	wäre		gäbe	CHA
,	,		,	
stunden	stunden			MOVS
wir	wir			MOVS
	nur	MOVT	nur	MOVT
noch	noch		noch	
nur		MOVS		MOVS
in	in		in	

LT	pos	lemma	TH1	TH1Diff	TH1pos	TH1posDiff
An	APPR	an	Auf	CHA	APPR	
der	ART	d			ART	
anderen	ADJA	andere	anderen		ADJA	
Seite	NN	Seite	Seite		NN	
,	\$,	,	,		\$,	
wenn	KOUS	wenn	wenn		KOUS	
da	PAV	da	da		PAV	
kein	PIAT	kein	kein		PIAT	
Feminismus	NN	Feminismus	Feminismus		NN	
wäre	VAFIN	sein	wäre		VAFIN	
,	\$,	,	,		\$,	
stunden	VVFIN	stehen	stunden		VVFIN	
wir	PPER	wir	wir		PPER	
			nur	MOVT	ADV	MOVT
noch	ADV	noch	noch		ADV	
nur	ADV	nur		MOVS		MOVS
in	APPR	in	in		APPR	

Exkurs: Suche in ANNIS

ANNIS? Tutorial logged in as "marc"

Search Form

AnnisQL: `pos="ADV" &
pos="ADV" &
ZH1Diff="MOV" &
#1.#2 &
#2=_#3`

Suchfenster

Query Builder:

Result: 23

Trefferzahl

More Corpora

Name	Texts
<input checked="" type="checkbox"/> FalkoEssayL2V2_0	248
<input type="checkbox"/> Register	22

Korpuswahl

Search

Context Left: 5

Context Right: 5

Results per page: 10

Search Result - pos="ADV" & pos="ADV" & ZH1Diff="MOV" & #1.#2 & #2=_#3

Token: `APPART NN VVINF $ ADV ADV A`

Text	Token
Außerdem	jetzt
zur Schule schicken	außerdem
jetzt	jetzt

Meta Data for id 303

Name	Value
SPK0:birth-year	1983
SPK0:cstest	64
SPK0:degree	N/A
SPK0:first-name	aa3ec21ddc1cdf8a9ebc617a19c
SPK0:I1_1	tur
SPK0:I1_1_away	N/A
SPK0:I1_1_away	N/A
SPK0:I1_1_durati	N/A
SPK0:I1_1_langs	N/A
SPK0:I1_1_schoc	N/A
SPK0:I1_1_since	N/A
SPK0:I1_1_univei	N/A
SPK0:L1index	tur:N/A:N/A:N/A:N/A:N/A:N/A
SPK0:I2_1	eng
SPK0:I2_1_away	N/A
SPK0:I2_1_away	N/A
SPK0:I2_1_durati	84

Meta Data for id 272

Name	Value
projectName	FALKO Essay Corpus L2 2.0
projectURL	Falko project site

Korpusmetadaten

Universitätsabschlüsse
NN
MOV
Universitätsabschlüsse
NN
Universitätsabschlüsse
MOV

Displaying Results 1 - 10 of 23



Exkurs: Suche in ANNIS

Search Form

AnnisQL: `pos="ADV" & pos="ADV" & ZH1Diff="MOVS"& #1.#2 & #2=_#3`

Query Builder: Show>>

Result: 23

Name	Texts	Tokens
FalkoEssayL2V2_0	248	131599
Register	22	19342

Search Export

Context Left: 5

Context Right: 5

Results per page: 10

ShowResult

Search Result - pos="ADV" & pos="ADV" & ZH1Diff="MOVS"& #1.#2 & #2=_#3 (5, 5)

Page 2 of 3

Original text & token

wäre , stünden wir noch nur in der Küche und köchten
 sein , stehen wir noch nur in d Küche und

VF	VF	PP	ADV	ADV	APPR	ART	NN	KON	VFIN	
						in	der	Küche	und	köchten
						in	d	Küche	und	köchen
			MOVS	MOVS	MOVT					CHA

Partituransicht

ZH2 & Fehlerannotationen

ZH2pos	V	V	PP	ADV	ADV	APPR	ART	NN	KON	VFIN
ZH2	g	g				in	der	Küche	und	köchten
ZH2lem	geben			nur	noch					
ZH2Diff	CHA			MOVS	MOVS	MOVT				MOVS

Partituransicht

ZH1 & Fehlerannotationen

ZH1lem	sein	,	st					Küche	und	köchen
ZH1Diff										CHA

linker & rechter Kontext

Partituransicht

ZH1 & Fehlerannotationen

ZH1lem	sein	,	st					Küche	und	köchen
ZH1Diff										CHA

Resultate anzeigen oder exportieren

Original text & token basierte Annotationen

Partituransicht ZH2 & Fehlerannotationen

Partituransicht ZH1 & Fehlerannotationen

Exkurs: Suche in ANNIS



väre , stünden v noch nur
sein , stehen v noch nur
VAFIN \$, VFIN PPE ADV ADV
+ ZHverb (grid)
+ ZH2 (grid)
+ falko (grid)
- ZH1 (grid)

noch nur

nur noch

Select Displayed Annotation Levels ▾

ZH1lemma	sein	,	stehen	wir	nur	noch
ZH1Diff					MOV	M
ZH1pos	VAFIN	\$,	VFIN	PPE	ADV	ADV
ZH1	väre	,	stünden	wir	nur	noch
tok	väre	,	stünden	wir		noch

MOVT = MOVEDtarget
Token sollte

source
wegt

+ text (grid)
- Volltext

Der Feminismus hat den Interessen der Frauen mehr geschadet als genützt Was heißt eigentlich Feminismus ? Ich meine es gibt unterschiedliche Stufen von diesem Fennomen . An einer Seite muss ich mit der Anzeige zustimmen . Der Feminismus hat uns - den Frauen - um einige Rechte geräubert . Oder Vorteile besser zu sagen . Wir können , sogar müssen , die männliche Arbeiten beherrschen , wir müssen schwere Sachen tragen und selbst die immer bereit sind , uns mit den Kofern und mit den Türen zu helfen . Die Frage ist eine gleichgerechte Gesellschaft schaffen ? An der anderen Seite , wenn da ke und köchten wir . Kein Studium , kein Selbstbewusstsein und die einzigen Gipfel , die den wir aber sogar selbst nicht gewählt könnten) und die Kinder zu gebären . Mein Frauen . Die Männer haben sich auch " feminisiert " . So dass heutige Generation der männer mit den Frauen in der Haushalt sicher mehr als die ältere . Mein Vater war anderer Meinung . Ich weiß , dass er selbst die Haushalt beherrschen konnte , z. B. wenn er unterwegs ohne Mutti war

Treffer im Volltext

noch nur in der Küche

Mögliche Studie

- Machen **romanische** Deutschlerner signifikant weniger Artikelfehler als **germanische**?
- Finde alle Artikel in der ZH2, die mit "**INS**" markiert sind
- Vergleich die Ergebnisse in Texten von **spanischen** und **italienischen** Lernern mit denen in Texten von **dänischen** und **afrikaans** Lernern.

Contrastive Interlanguage Analysis

- Auffinden struktureller Lernerschwierigkeiten im L2-Deutsch
- Strukturelle Schwierigkeiten sind solche, die
 - **unabhängig** sind von der **Lerner-L1**
 - daher **abhängig** sind von der **Grammatik der Zielsprache**



Arts & Humanities
Research Council

HU-Berlin

Anke Lüdeling
Marc Reznicek

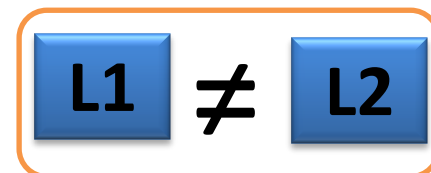
Bangor

University

Astrid Ensslin
Cedric Krummes

Wie finden wir komplizierte Strukturen? ^{Falko ✓}

- L1-korrelierte Strukturen (Transfer)



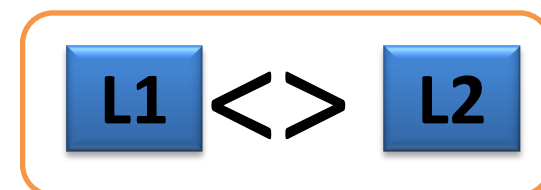
- Von den Lernern als kompliziert eingestufte Strukturen



- Fehlerhafte Strukturen

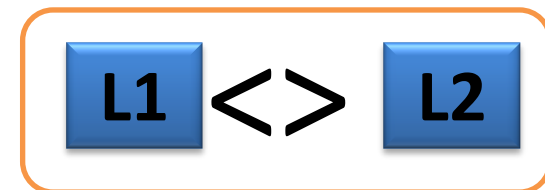
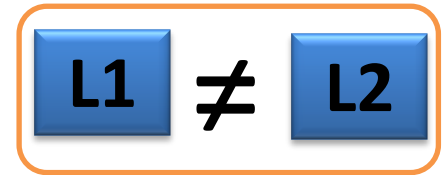


- Unterrepräsentierte Strukturen



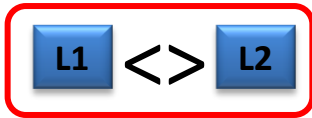
Wie finden wir komplizierte Strukturen?

- L1-korrelierte Strukturen (Transfer)
grammatische Analysen
→ nur sehr beschränkte Vorhersagekraft
- von den Lernern als kompliziert eingestufte Strukturen
Lernerintuition
→ unsystematisch, lehrformabhängig
- fehlerhafte Strukturen
Lehrerintuition
→ unsystematisch, normenabhängig
- unterrepräsentierte Strukturen
Korpusanalyse: Frequenzvergleiche
→ Overuse/Underuse



Mindergebrauch (Underuse)

Falko ✓

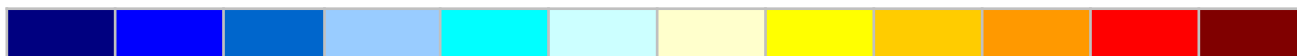


- Frequenzen von Strukturen im L2-Subkorpus werden verglichen mit Strukturen im L1-Subkorpus
 - Overuse & underuse sind definiert als (statistisch signifikante) Unterschiede zwischen zwei Varietäten
 - Eine Kategorie kann mindergebraucht sein, weil:
 - die Lerner sie nicht kennen
 - die Lerner sie zwar kennen, aber (unbewusst?!?) vermeiden
- Diagnostik für das Auffinden schwieriger Strukturen

Overuse/Underuse-Visualisierung

- underuse: kalte Farben
- overuse: warme Farben
- Intensität der Farbe signalisiert Stärke des overuse/underuse

Underuse



Overuse

Excel –AddIn (Amir Zeldes) verfügbar unter:

<http://korpling.german.hu-berlin.de/~amir/uoadin.htm>

Overuse/Underuse-Visualisierung lexikalische Einheiten

lemma	tot_norm	deu	dan	eng	fra	pln	rus
in	0.013188	0.012261	0.014041	0.014247	0.015272	0.012135	0.009534
es	0.010897	0.011945	0.010900	0.011379	0.013347	0.008163	0.012385
sie	0.010618	0.008193	0.010643	0.008835	0.010909	0.006067	0.005613
man	0.010164	0.007900	0.012438	0.008742	0.009754	0.006950	0.007306
dass	0.009522	0.007404	0.012823	0.008789	0.009625	0.008880	0.009890
von	0.007982	0.007122	0.007309	0.006846	0.007315	0.010259	0.007930
auch	0.007028	0.008362	0.008527	0.005828	0.005775	0.005461	0.004455
für	0.006683	0.007201	0.006091	0.007216	0.006802	0.005736	0.004188
sind	0.006465	0.004271	0.008976	0.007308	0.006930	0.004964	0.005346
sich	0.006309	0.011697	0.006283	0.006291	0.006930	0.007170	0.005435
ich	0.006262	0.003877	0.013272	0.005366	0.003465	0.001434	0.001426
aber	0.006048	0.003347	0.007309	0.006245	0.007315	0.003365	0.003831

sich ist in allen L1-Subkorpora mindergebraucht

Stuttgart-Tübingen-Tagset (STTS)

(Schiller et al. 1995)

ADJektiv	Nomen	Pronomen	Verb	Partikel	Konjunktion
ADJA	NN	PDS	VVFIN	PTKZU	KOUI
ADJD	NE	PDAT	VVIMP	PTKNEG	KOUS
		PIS	VVINFL	PTKVZ	KON
		PIAT	VVIZU	PTKANT	KOKOM
		PIDAT	VVPP	PTKA	
		PPER	VAFIN		
		PPOSS	VAIMP		
		PPOSAT	VAINFL		
		PRELS	VAPP		
		PRELAT	VMFIN		
		PRF	VMINFL		
		PWS	VMPP		
		PWAT			
		PWAV			

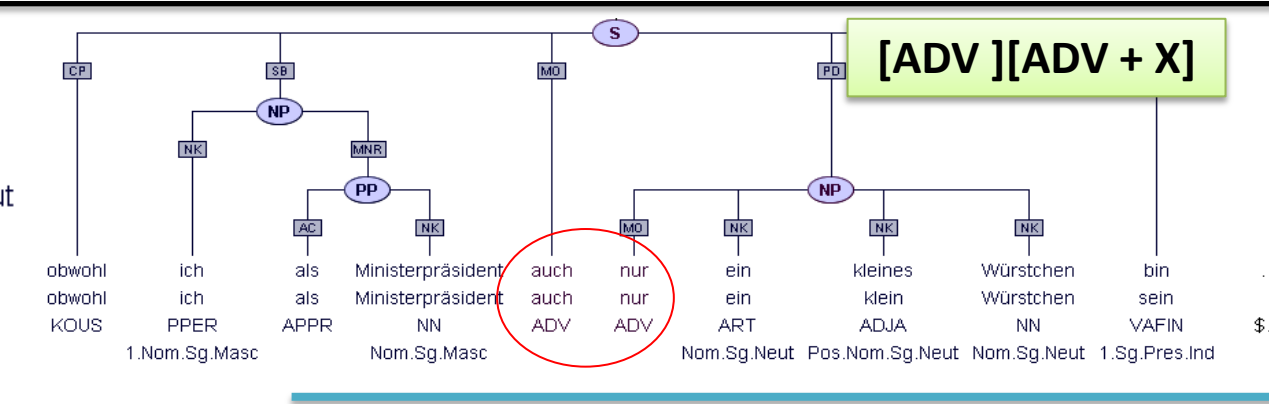
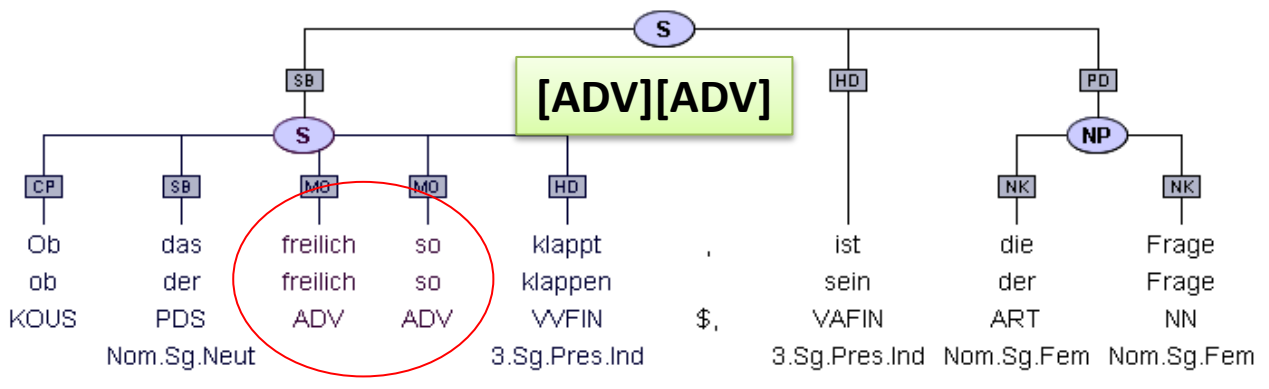
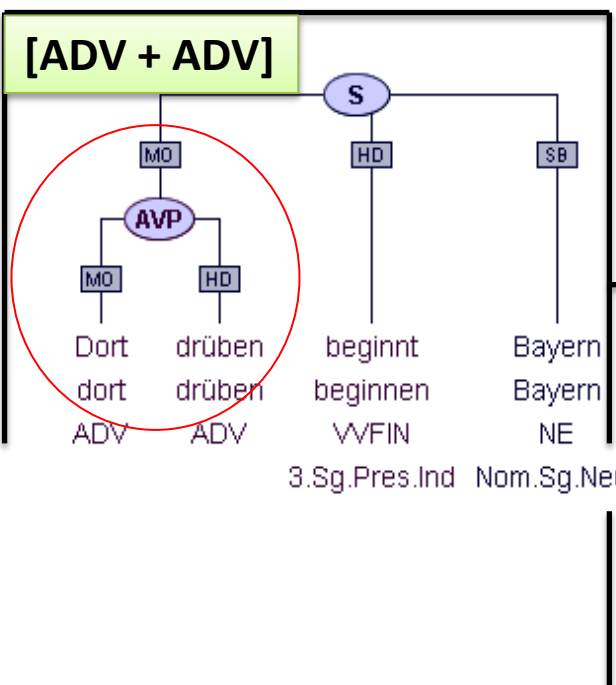
Overuse/Underuse-Visualisierung POS-Bigramme

bigram	tot_norm	de	da	en	fr	pl	ru
\$.PPER	0.042384	0.005297	0.009748	0.007963	0.006166	0.005801	0.007409
VVFIN-\$,	0.042131	0.006457	0.00776	0.006343	0.006937	0.006243	0.008391
PPOSAT-NN	0.041739	0.008058	0.007247	0.007269	0.007066	0.006298	0.005802
ADV-ADV	0.041604	0.012858	0.010518	0.006111	0.006166	0.003094	0.002856
ADV-APPR	0.039742	0.009117	0.008016	0.005324	0.007837	0.004807	0.004642
PDAT-NN	0.03956	0.005409	0.004233	0.005509	0.007837	0.007735	0.008837
ADV-ART	0.037125	0.007629	0.006349	0.006898	0.005653	0.006133	0.004463

Adverbketten sind in allen L1-Subkorpora mindergebraucht

Fazit

- **Adverbketten** werden von allen Lernern unabhängig von der Muttersprache vermieden.



Fazit

- Adverbketten werden von allen Lernern unabhängig von der Muttersprache vermieden.

[ADV + ADV]

[ADV][ADV]

[ADV][ADV + X]

→ Strukturen mit sehr variabler Tiefenstruktur werden besonders häufig vermieden.

→ Werden die **Formen** vermieden oder die **Funktionen**?

Hypothese: Modifikation wird generell vermieden.

CIA: syntaktische Funktion

Untersuchung zu Underuse feinerer Adverbklassen (Hirschmann erscheint) zeigt einen Unterschied zwischen adverbialen Funktionen.

→ Wie sieht der Vergleich rein syntaktischer Funktionen unabhängig von der Füllerkategorie aus?

→ Methode: tiefe syntaktische Annotation der Lernerdaten

Falko – syntaktische Vorverarbeitung

- Direkte Verarbeitung von Lernaltersprache ist problematisch

- Äußerungen gehorchen nicht der beschreibenden Grammatik

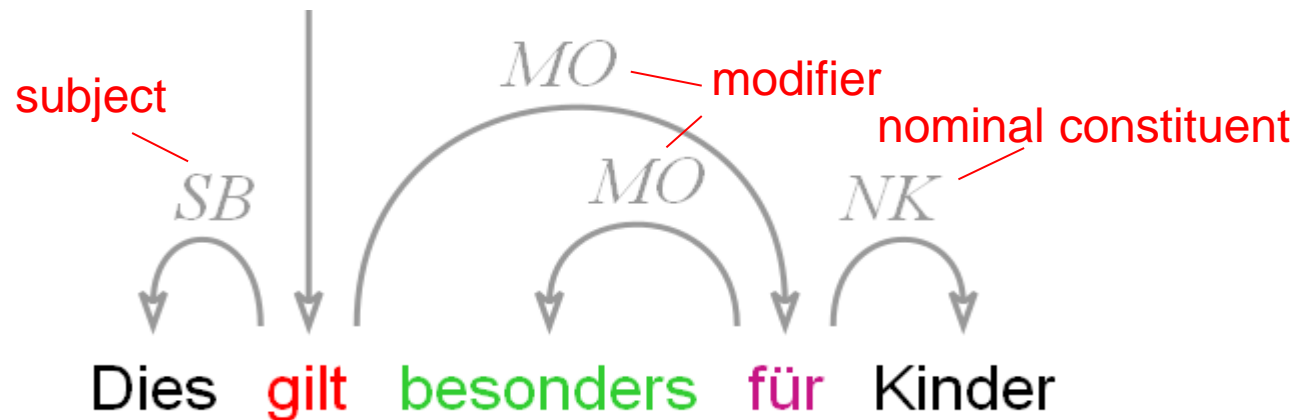
(siehe Hirschmann et al. 2009, Díaz-Negrillo et al. 2010)

→ ZH1 als Grundlage für das Parsing

→ Rückübertragung der Annotationen auf den Lernertext (ZH0)

Syntaxschema (Kurzformat)

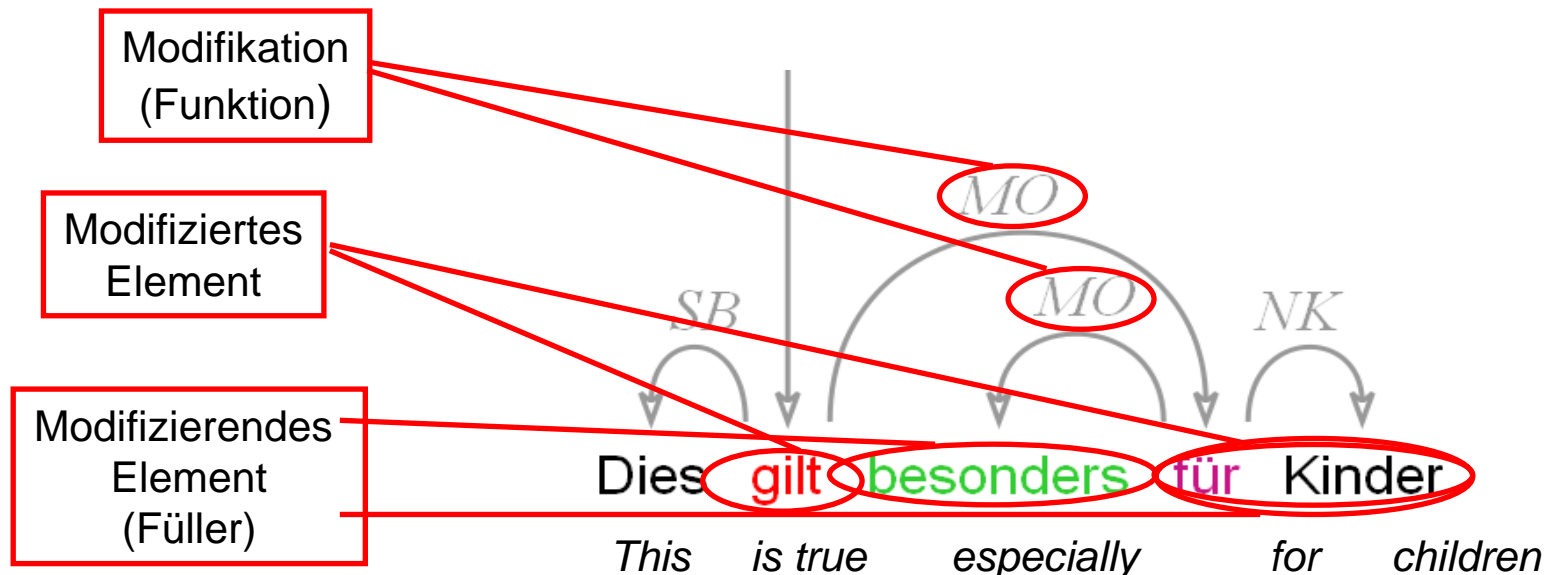
- Jedes Wort ist mit den von ihm abhängigen Worten verbunden.
- Pfeile zeigen zu hierarchisch niedrigeren Dependents.
- Jeder Pfeil (Dependenz) trägt eine Funktionsbezeichnung.



Suche nach Modifikationen

Verschiedene Aspekte des Problems

- Wird die **syntaktische Funktion** vermieden?
- Werden bestimmte **Ziele der Modifikation** vermieden?
- Werden bestimmte **modifizierende Kategorien** vermieden?



Underuse/Overuse von Funktionen

label	de	da	en	fr	ru	usb
NK	0,264067	0,278546	0,284881	0,303271	0,29552	0,295136
HD	0,156192					
MO	0,141968	0,12789	0,113704	0,110112	0,112513	0,108707
SB	0,07398	0,078506			0,078852	0,085512
CJ	0,059604	0,053397	0,056411	0,050632		0,072183
AC	0,057051					0,04916
OC	0,050335					0,040679
OA	0,044213					
CD	0,026549			0,022156		
CP	0,017653	0,021732	0,020325			
PD	0,014435		0,015943		0,016947	0,018002
NG	0,011065					
MNR	0,010995	0,013707	0,013429	0,013383		
RC	0,010051				0,006268	0,005366

MO (modification) ist signifikant mindergebraucht
unabhängig von L1

Modifikationen in Falko

- Alle Kategorien werden oft modifiziert in L1- und L2-Daten
- aber **alle syntaktischen Modifikationsrelationen** zeigen einen Underuse
- Adverbmodifizierer zeigen den stärksten Underuse
- Das gilt **unabhängig** von der **L1 der Lerner**

Produktivität bei Lernern

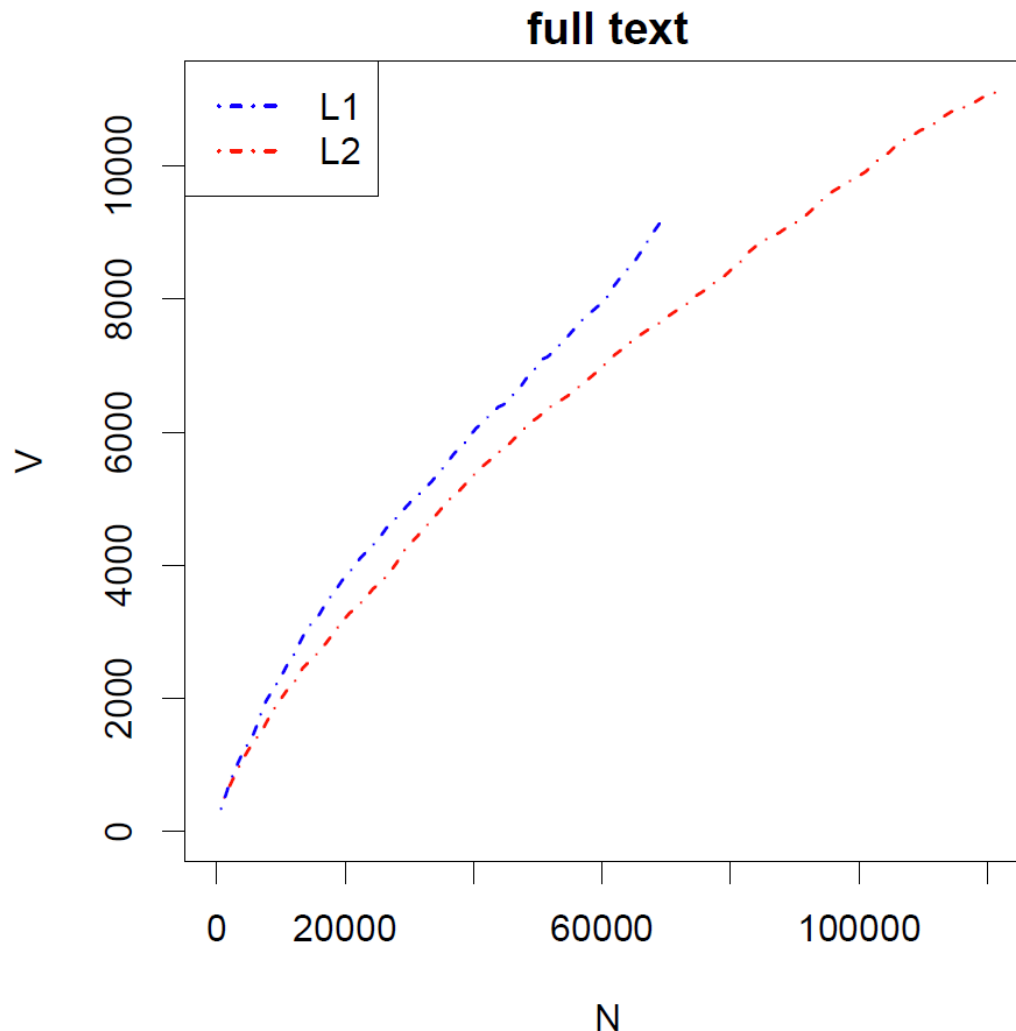
- Viele Studien zeigen, dass Muttersprachler oft Worte neubilden
 - sind sich (implizit) Wortbildungsregeln bewusst
 - wahrscheinlich aus Input abgeleitet

(Baayen 1992, 2001 etc., Plag 1998, Bauer 2001, Lüdeling & Evert 2005, Kiss 2007, Zeldes to appear etc.)
- Gilt das auch für Lerner?
- Wie können Lernerkorpora uns bei der Lösung helfen?

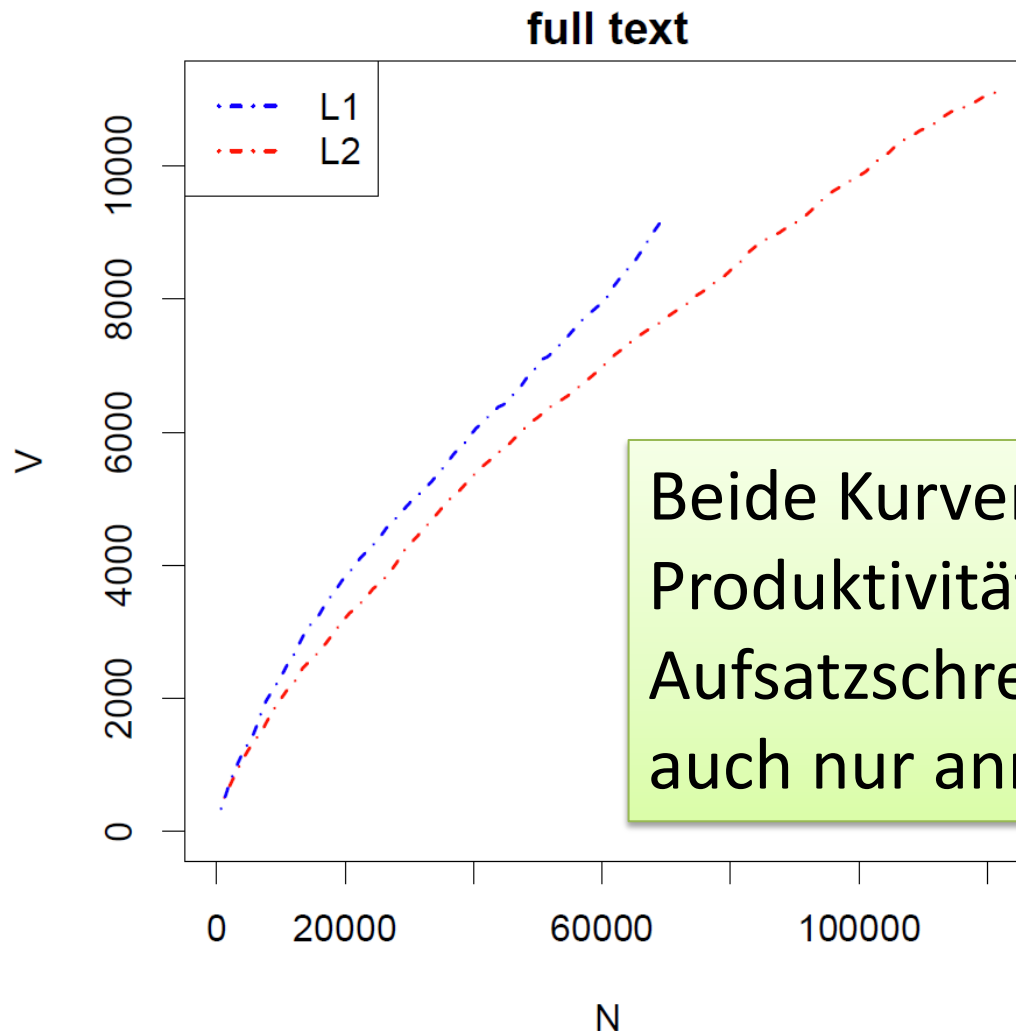
Vokabelwachstumskurven

- Viele korpusbasierte Maße für morphologische Produktivität (Baayen 2001, 2009 etc.)
- Kurz: Vokabelwachstumskurven
 - **N** = Token vs. **V** = Typen
- Zeigen ob die Stichprobe (Korpus) das gesamte Vokabular ausreicht, oder noch viele (mögliche) Worte in der Gesamtpopulation gibt.

Vokabelwachstumskurve für Falko L1 und L2



Vokabelwachstumskurve für Falko L1 und L2



Beide Kurven zeigen eine hohe Produktivität – weder L1 noch L2 Aufsatzschreiber kommen erreichen auch nur annähernd die Vokabelgrenze

Komplexe Verben im Deutschen

Präfixverben

- *ver*•*kaufen*
- [...] dass Peter Schokolade *verkauft*.
- Peter *verkauft* Schokolade.
- Infinitive: zu *verkaufen*
- Partizip: *verkauft*

Partikelverben

- *auf*•*essen*
- [...] dass Peter die Schokolade *aufisst*.
- Peter *isst* die Schokolade *auf*.
- Infinitive: *aufzuessen*
- Partizip: *aufgegessen*

ZHverb: komplexe Verben in Falko

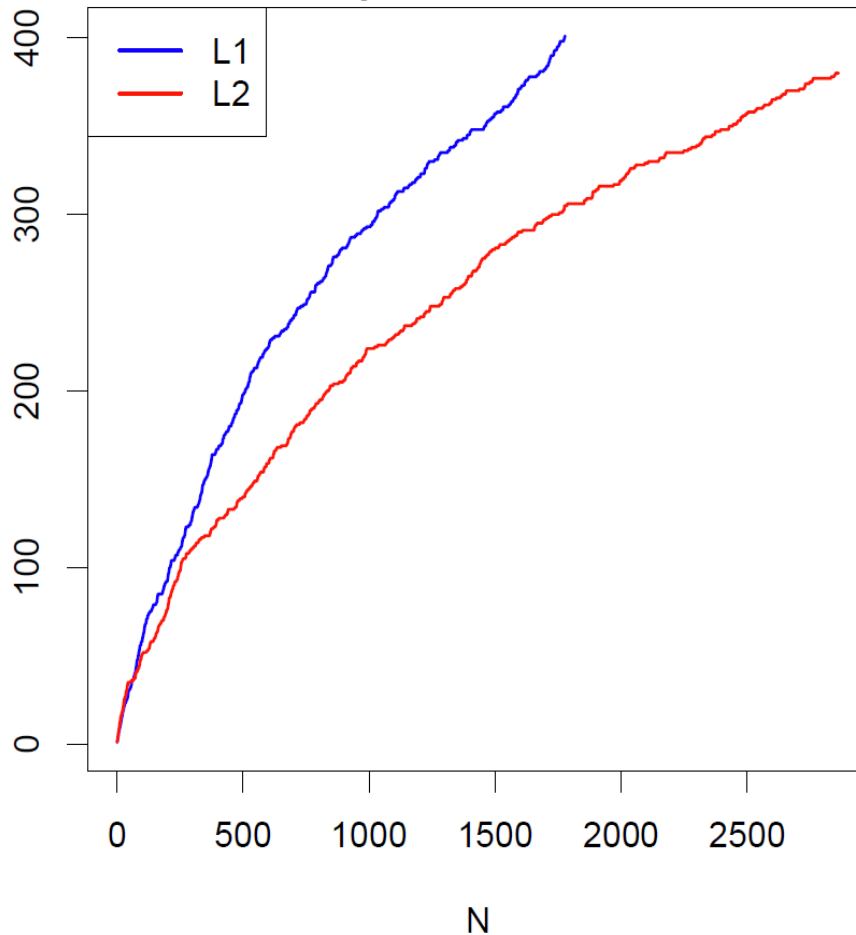
- Grundlage ist ZHverb (konzentriert sich auf komplexe Verben)
- manuelle Annotation
 - Typ
 - Partikel vs. Präfix
 - Lemma
 - Fehlerklasse
 - Orthographie, Flexion, Argumentstruktur, Trennbarkeit etc.
 - Form
 - finit, infinit, Partizip, getrennt etc.

Exkurs: Fehler

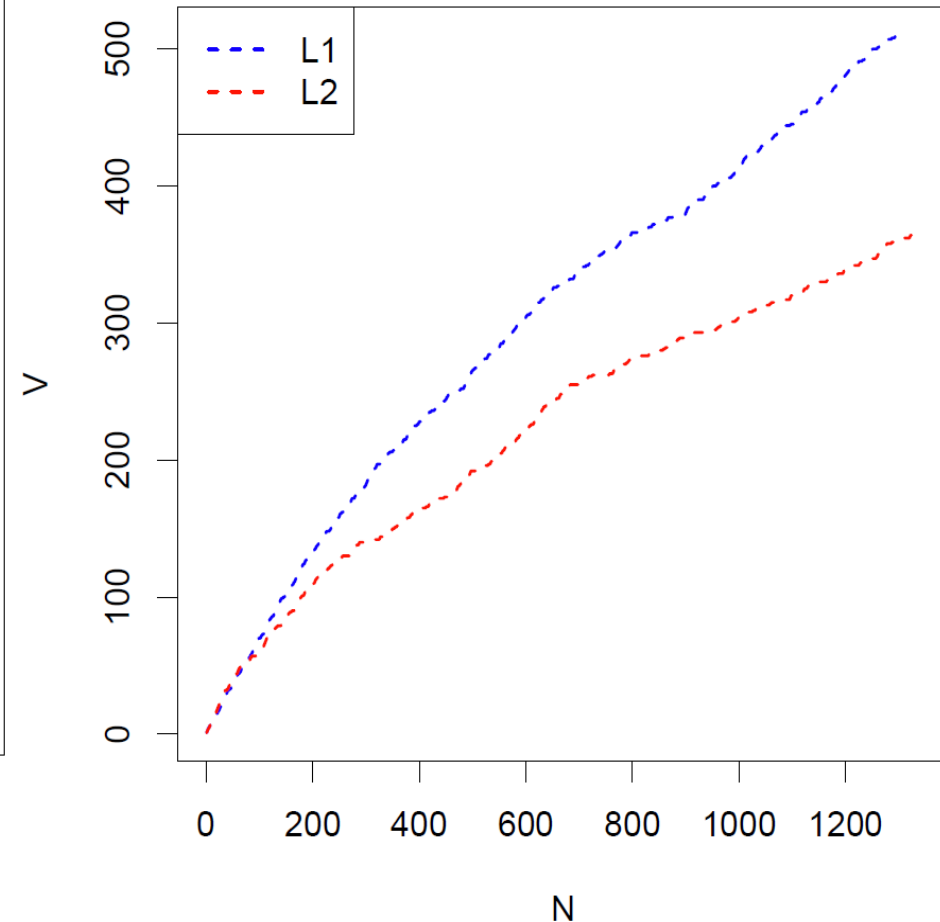
- Alle Fehlerdifferenzen sind signifikant
 - Lerner machen signifikant **mehr strukturelle** und **semantische Fehler**
 - Muttersprachler machen signifikant **mehr Orthographiefehler**

Vokabelwachstum

prefix verbs



particle verbs



Komplexe Verben

- **Komplexe Verben** sind für L1 und L2 **produktiv**
 - L1 benutzen mehr produktive Verben
 - Kennen L1-Sprecher mehr komplexe Verben?
 - Wissen L1-Sprecher besser wie man komplexe Verben bildet?
- a) Analyse produktiver Muster
- b) Analyse klarer Fälle von Wortneubildung

Komplexe Verben

a) Produktives Muster : hinein + Verb

hineinbewegen – move into

hineingehen – go into

hineintreten – step into

hineinbringen – bring into

hineinfallen – fall into

hineingeraten – get into

hineinpassen – fit into

hineinversetzen – put yourself in
someones position

hineinwachsen – grow into

hineinziehen – pull into

- für viele mögliche Muster zeigt L1 mehr Typen
- L1 verwendet mehr Muster

!!! Nicht klar, welche davon produktiv gebildet!!!

Komplexe Verben

b) 'unbekannte' Worte

L2

L1

34 'neue' Form (uns unbekannt, Token)

... was eine sozialere oder gerechtere Entlohnung benachträchtigt. (hu_005_2006_09)

274 neue Lesarten für bekannte Formen (Token)

Viele Leute sich bewundern, ob ... (hu_006_2006_10)

23 'neue' Formen (uns unbekannt, Token)

Sie vollrichten in ihrer Arbeitszeit fast doppelt so viel ... (dhw_031_2007-06)

34 neue Lesarten für bekannte Formen (Token)

Die Regierung kann dieses große Netz der Machenschaften meist nicht entspinnen. (dew_03_2007_09)

Komplexe Verben

- Lerner bilden 'neue' Wörter signifikant häufiger als benutzen 'neue' Lesarten häufiger als Muttersprachler
- machen Fehler beim Anwenden der Muster

Zusammenfassung

Tief annotierte Lernerkorpora erlauben die Beantwortung neuer Forschungsfragen.

- Abweichungen in der Frequenz bestimmter Muster auf unterschiedlichen Abstraktionsebenen
 - Lemmata, Wortarten, Wortartenketten, syntaktische Funktionen
- Unterschiede in der graduellen Beschreibung von Produktivität Produktionsmustern
- Auffinden von in den Lernertexten NICHT vorhandenen Elementen (unterlassene Artikel)

Zusammenfassung

Sowohl für Fehleranalysen als auch für rein kontrastive Analysen ist eine explizite Zielhypothese nötig

Fehler: Grundlage für Abweichungen

Kontrast: Beschreibungsgrundlage für zielsprachliche Annotationsschemata

Einladung

There is no data like more data.

Wir suchen fleißige Kollaborateure bei der
Vergrößerung des Korpus.

Jeder ist eingeladen, mitzumachen

Danke

Literatur:

- Borin, Lars; Prütz, Klas (2004): New wine in old skins? A corpus investigation of L1 syntactic transfer in learner language. In: Guy Aston, Silvia Bernardini, Dominic Stewart (Eds.): Corpora and language learners. Amsterdam, Philadelphia: John Benjamins (Studies in corpus linguistics, 17), 67–87.
- Granger, Sylviane (2008): Learner Corpora. In: Anke Lüdeling, Merja Kytö (Eds.): Corpus linguistics. Berlin, New York: Mouton de Gruyter (Handbücher zur Sprach- und Kommunikationswissenschaft / Handbooks of Linguistics and Communication Science, 29,1), 259–275.
- Kübler, Sandra; Scheutz, Matthias; Baucom, Eric; Israel, Ross (2010): Adding Context Information to Part of Speech Tagging for Dialogues. In: Proceedings of the 9th International Workshop on Treebanks and Linguistic Theories (TLT 10). Tartu, Estonia.
- Schiller, Anne; Teufel, Simone; Thielen, Christine (1995): The Stuttgart-Tübingen Tagset (STTS).