

Anke Lüdeling / Seanna Doolittle / Hagen Hirschmann /
Karin Schmidt / Maik Walter¹

Das Lerner korpus Falko

0 Lerner korpora in der Erforschung des Fremdspracherwerbs

In diesem Artikel wird das frei zugängliche fehlerannotierte **Lerner korpus Falko** (<http://www.lingustik.hu-berlin.de/institut/professuren/korpuslinguistik/forschung/falko>) vorgestellt, das Texte von fortgeschrittenen Lernern des Deutschen als Fremdsprache enthält. Auch wenn es bisher noch kaum elektronisch verfügbare Lerner korpora für DaF gibt, steht natürlich die empirische Beschäftigung mit Sammlungen von Lerner daten in langer Tradition.

Viele Aspekte des Erst- oder Zweit-/ Fremdspracherwerbs können nur auf der Basis von tatsächlich vorkommenden Lerner äußerungen untersucht werden. Dies betrifft einerseits die genaue Be trachtung von Fehlern bzw. Abweichungen (in einem Kontext) und andererseits vor allem auch die Ermittlung quantitativer Eigenschaften von Lernersprache (im Vergleich zur L1 oder im Vergleich mit anderen Lerner varietäten). Lerner äußerungen können entweder in (mehr oder weniger) authentischen Situationen entstehen oder in gezielten Fragebogenstudien oder psycholinguistischen Studien eliziiert werden.² Seit sich die Wissenschaft systematisch mit der Erforschung des Spracherwerbs beschäftigt, werden Lerner äußerungen gesammelt und ausgewertet (eine frühe Sammlung von Erstspracherwerbsdaten ist Stern/Stern 1907; zur Rolle von Daten im Zweitspracherwerb vgl. z.B. Dietrich 2006).

Viele Sammlungen von Lerner äußerungen sind aber nicht systematisch – die Erkenntnisse, die sich daraus ableiten lassen, sind daher oft nicht ohne weiteres generalisierbar. Außerdem sind die Daten oft nicht allgemein zugänglich, sodass Ergebnisse nicht reproduzierbar sind. In den letzten Jahren werden deshalb Lerner korpora – also Sammlungen von Lerner daten, nach genauen Kriterien und mit ausführlichen Metadaten – zunehmend kontrolliert erhoben.³

Im Folgenden werden wir uns nur auf Lerner korpora mit Daten von Fremdsprach lernern beziehen (wir kürzen diese Lerner daten als L2-Daten ab, unabhängig von der Anzahl der vorher gelernten Sprachen).

Während es für das Englische als Fremdsprache bereits mehrere Lerner korpora und viele Studien gibt,⁴ liegen für das Deutsche amtierte LEAP-Korpus (vgl. Milde/Gut 2002) und die ESF-Korpora des MPI Nijmegen (http://corpus.uni.nl/ds/midi_browser). Um diese

¹ Unser Dank geht an Emil Kroymann, Marc Reznick und alle anderen Falkos für ihre Hilfe.

² Ergänzend dazu gibt es weitere psycholinguistische Studien, wie zum Beispiel die Studien zur Konzeptualisierung in der L2 (vgl. von Stutterheim/Carroll 2006).

³ Neben Lerner korpora spielen natürlich auch Mutter sprachler korpora (L1-korpora) im Fremdsprachenbe reich eine Rolle, allerdings eher im Bereich Sprachver mittlung. Wir können hier auf L1-korpora oder L2-korpora in der Sprachvermittlung eingehen; vgl. dazu zum Beispiel Nesselhauf (2004); Fandrych/Tschimer (2007); Römer (i.Dr.).

⁴ Das bekannteste Lerner korpus des Englischen als Fremdsprache ist sicher das International Corpus of Learner English (ICLE), das an der Universität Lou vain-la-Neuve erstellt wird (vgl. Granger/Dagneau/Meanier 2002). Das Korpus besteht aus Essays von Englis chlern mit unterschiedlichen Muttersprachen und ist Grundlage für sehr viele Untersuchungen (eine Sammlung von Literaturhinweisen findet sich unter <http://cecl.fltr.ucl.ac.be/Ced/icle/icle.htm>). Für weitere Korpora siehe Granger (2003); Granger (i.Dr.) und die Proceedings der TALC(Teaching and Language Corpora)-Konferenzen.

⁵ Neben privaten, nicht öffentlich zugänglichen Sammlungen (zum Beispiel in Wegener 1995; Birkner/Dinroth-Dittmar 1995; Weinberger 2002; Belz 2004) gibt es so weit wir wissen, nur die folgenden: das phonologisch annotierte LEAP-Korpus (vgl. Milde/Gut 2002) und die ESF-Korpora des MPI Nijmegen (http://corpus.uni.nl/ds/midi_browser).



Fun with languages!
Иностранные языки – с удовольствием!
Che spasso con le lingue!
Šála sa jecíma! Le plaisir des langues!
Diversão com idiomas!
Spaß am Sprachen!
Keori ya үйләсөг!
Dil ögemenin keyfi!

angues
TUM
Langenscheidt

ISSN: 0011-9741
Artikelnummer: 1622
Erscheinungsweise: jährlich 4 Hefte
Preis: für die Bundesrepublik Deutschland inkl. Mehrwertsteuer Jahresabonnement € 36,- zuzügl. Porto
Bestellung: in jeder Buchhandlung; für Privatpersonen auch im Verlag (schriftlich; telefonisch: +49-(0)341-973752; Fax: +49-(0)341-9737548;
Redaktion: Prof. Dr. Dr. h. c. Gerhard Helbig (Chefredakteur), Dr. Bernd Skibitzki; Herder-Institut der Universität Leipzig, Beethovenstraße 15, 04107 Leipzig;
Telefon: +49-(0)341-973752; Fax: +49-(0)341-9737548;
E-Mail: daRed@server1.rz.uni-leipzig.de

Redaktionsbeiträte: Prof. Dr. Lutz Götz (Saarbrücken), Prof. Dr. Ursula Hirschfeld (Halle), Prof. Dr. Werner Reinecke (Leipzig), Prof. Dr. Dietmar Rosler (Gießen), Prof. Dr. Dr. h. c. Horst Sitta (Zürich), Prof. Dr. Peter Suchstand (Jena). Prof. Dr. Erwin Ischirner (Leipzig), Prof. Dr. Barbara Wojak (Leipzig)



SSG
Deutsche Gesellschaft für Linguistik
und Sprachwissenschaft

Deutsch als Fremdsprache – Zeitschrift zur Theorie und Praxis des Deutschunterrichts für Ausländer
© 2008 by Langenscheidt KG Berlin und München
Druck: CS Druck, Berlin
Printed in Germany
Vervielfältigung (auch fotomechanisch) und Verarbeitung des Inhalts dieser Zeitschrift nur mit Genehmigung des Verlages.
ISSN: 0011-9741
Artikelnummer: 1622
Erscheinungsweise: jährlich 4 Hefte
Preis: für die Bundesrepublik Deutschland inkl. Mehrwertsteuer Jahresabonnement € 36,- zuzügl. Porto
Bestellung: in jeder Buchhandlung; für Privatpersonen auch im Verlag (schriftlich; telefonisch: +49-(0)89-56096333)

Mit Beginn des Jahres 2007 werden eingesandte Manuskripte **doppelblind** begutachtet.
Die Redaktion behält sich das Recht auf Kürzung und Bearbeitung eingesandter Manuskripte vor.
Autonome Beiträge geben nicht unbedingt die Meinung der Redaktion wieder.

Redaktionsbeiträte: Prof. Dr. Lutz Götz (Saarbrücken), Prof. Dr. Ursula Hirschfeld (Halle), Prof. Dr. Werner Reinecke (Leipzig), Prof. Dr. Dietmar Rosler (Gießen), Prof. Dr. Dr. h. c. Horst Sitta (Zürich), Prof. Dr. Peter Suchstand (Jena). Prof. Dr. Erwin Ischirner (Leipzig), Prof. Dr. Barbara Wojak (Leipzig)

Lücke zu schließen, wurde das Lernerkorpus Falko in einer Zusammenarbeit der Freien Universität (FU) Berlin und der Humboldt-Universität (HU) zu Berlin konzipiert und erhoben.

Wir möchten zunächst in Abschn. 1 einige Aspekte der Erstellung von Lerner korpora ansprechen und erläutern, wie wir uns jeweils für das Falko-Korpus entschieden haben. Fehlerannotation ist für viele Forschungsfragen hilfreich, aber konzeptuell sicherlich problematisch. Daher werden wir in Abschn. 2 auf die Fehlerannotation in Falko genauer eingehen. In Abschn. 3 werden dann exemplarisch zwei kurze Studien mit Falkodaten vorge stellt.

1 Design und Architektur von Falko

1.1 Korpusdesign

Es gibt viele Definitionen des Begriffs „Korpus“. In der Korpuslinguistik hat sich aber in den letzten Jahren weitgehend durchgesetzt, dass eine Textmenge dann als Korpus bezeichnet wird, wenn sie für einen bestimmten Forschungszweck nach vorher definierten externen oder internen Kriterien zusammenge stellt wurde (vgl. zum Beispiel die EAGLES-Definiti on unter <http://nljissi/et/talks/korpus/eagles-corpus.html>; Hunston i.Dr.). Für Lerner korporpora einer gegebenen L2 sind hier – jeden nach Forschungsfrage – externe Kriterien (mehr als eine L1, gesprochene oder geschriebene Texte, Erhebungssituation, Lernstand) genauso wie interne Kriterien (Texte, die bestimmt sprachliche Strukturen aufweisen)

¹ Die Fortgeschrittenheit wird unterschiedlich festgestellt, s. unten; zur Diskussion von Fortgeschrittenheit vgl. auch Walter/Gronnes (2008a).

² Einige Texte wurden handschriftlich erhoben, andere am Computer. Bei handschriftlicher Erhebung fallen Tippfehler weg. Allerdings ist es oft kaum möglich, die Handschrift eindeutig zu entziffern. In der Digitalisierung müssen daher Entscheidungen getroffen werden. Außerdem können hier Übertragungsfehler entstehen. Direkt am Computer entstandene Texte können Tippfehler ent halten. In unseren Erhebungen wurde immer darauf geachtet, dass kein Textverarbeitungsprogramm mit eingebauter Rechtschreibkorrektur verwendet wurde.

³ Das frei zugängliche Korpus Akademisches Deutsch umfasst 1,8 Millionen Tokens. Es ist automatisch mit Wortarten annotiert und lemmatisiert, ein nach Fachgebieten ausgewogenes Subkorpora Akademisches Deutsch 2006 umfasst 841.483 Tokens.
⁴ Bisher haben wir Daten aus Dänemark, Kenia, Usbekistan, der Türkei und Südafrika sowie Daten aus Ferienstunden an der HU und der FU. Weitere Erhebungen sind geplant.

Die Erhebungen für beide Subkorpora waren stark kontrolliert: Es durften keine Hilfsmittel verwendet werden, und die Zeit war vorgegeben.²

Das *Zusammenfassungskorpus* ent hält Zusammenfassungen von linguistischen oder literaturwissenschaftlichen Fachtexten, die im Rahmen einer obligatorischen Sprachprüfung von nichtmuttersprachlichen Studierenden der Germanistik an der FU erstellt wurden. Die Fortgeschrittenheit wird hier definiert, dass die Studierenden die DSH-Prüfung bestanden und ihr Grundstudium in Deutschland absolviert haben. Zug ordnet ist ein Kontrollkorpus mit Daten von Muttersprachlern, das unter vergleichbaren Bedingungen (dieselben Vorlagentexte, gleiche Bearbeitungszeit etc.) erhoben wurde. Außerdem sind die Vorlagentexte abfragbar. Das Zusammenfassungskorpus ist abgeschlossen.

Als externes Vergleichskorpus haben wir ein Korpus mit deutschsprachigen Dissertations abstractis zusammengestellt.³

Das *Essaykorpus* enthält Essays zu vier kontrovers zu diskutierenden Themen, die sich an Themen des ICLE-Korpus anlehnen. Die Daten wurden in unterschiedlichen Ländern erhoben,⁴ das Korpus wächst noch. Die Fortgeschrittenheit wird mit Hilfe eines C-Tests überprüft; nur Texte von Lernern, die einen bestimmten Wert überschreiten, werden aufgenommen. Zugordnet ist wieder ein L1-Korpus, das an Brandenburger und Berliner Gymnasien (12./13. Klasse) erhoben wurde.

fügt werden.² Dies ist insbesondere für die Fehlerannotation wichtig, auf die wir in Abschn. 2 eingehen. Alle Annotationsebenen können bei der Suche kombiniert werden.

1.3 Annotation

In nichtannotierten Korpora kann man nur nach Zeichenketten suchen. Um Strukturen oder Muster suchbar zu machen, werden Korpora annotiert, d. h. mit Metainformationen versehen. Neben Header-Informationen, die sich auf ganze Dokumente beziehen, können auch kleinere Einheiten annotiert werden. Dazu muss ein Korpus zunächst in kleinste annotierbare Einheiten (Tokens) zerlegt werden.

Wir haben uns in Falko für eine automatische Tokenisierung in graphemische Wörter (Ab rücken von Satzzeichen, stattdessen Zeichenketten zwischen jeweils zwei Leerzeichen) entschieden, weil das in den meisten Korpora europäischer Sprachen so gemacht wird. Die bekannten Probleme einer solchen Tokenisierung (Mehrwortlexeme und mehrdeutige Zeichen; vgl. Schmid i.Dr.) nehmen wir in Kauf.

Falko wird automatisch mit Wortarten annotiert. Dazu verwenden wir den TreeTagger (vgl. Schmid 1994) mit dem Stuttgart-Tübinger-Tagger (<http://www.ims.uni-stuttgart.de/projekte/complex/TagSets/sts-table.html>). Die Fehlerrate im Wortarten-Tagging ist in einem Lernerkorpus wegen der orthographischen Fehler und der Wortstellungsfehler höher als bei Zeitungskorpora. Aus diesem Grund wurde in Teilen des Korpus die automatische Annotation manuell nachkorrigiert (und diese zusätzliche Annotationsebene kann bei Suchanfragen berücksichtigt werden).

Ebenfalls automatisch erfolgt die Annotation der Lemmata. D. h., jeder Wortform wird eine Grundform zugewiesen, sodass nach allen Formen, die zu denselben Grundform gehören, gesucht werden kann.

In Falko wurden darüber hinaus bisher ZIELhypthesen (s. Abschn. 2) und topologische Felder (vgl. Höhle 1986) annotiert. Mit Letzteren ist es beispielsweise möglich, Wortstellungsfehler, die Besetzung von Vorfeldern oder von Satzklammern effizient zu unterscheiden. Durch eine Kombination der einzelnen Ebenen können auch Lemmata in den einzelnen Feldern (z. B. alle Flexionsformen des Nomens *Problem* im Vorfeld eines Matrixsatzes) oder aber spezifische Wortartbesetzungen in den Feldern (z. B. alle attributiven Adjektive,

Subkorpora	Texte	Tokens	Types
Zusammenfassungen L1, Version 1.1	107	41075	5181
Essays L2, Version 1.0	57	21211	4099
Essays L1, Version 0.5	132	65387	7720
Essays L1, Version 0.5	39	33806	5831

Tab. 1: Korpusgrößen für die einzelnen Subkorpora in Falko
(Das Zusammenfassungskorpus ist abgeschlossen; das Essaykorpus wird weiter wachsen.)

1.2 Korpusarchitektur

Die Forschungsfragen beeinflussen nicht nur das Korpusdesign, sondern auch die Korpusarchitektur, d. h. das formale Modell, in dem das Korpus und seine Annotation gespeichert sind. Unterschiedliche Architekturen ermöglichen bzw. verhindern bestimmte Annotationen. Die meisten Korpora sind flach annotiert, d. h., die Annotation erfolgt mit besonderer Kennzeichnung in derselben Datei wie die Daten, sei es in einem Tabellenmodell oder in einem Baummodell (XML).¹ Das ist für große, einfach annotierte Korpora eine gute Option. Bestimmte Möglichkeiten gibt es aber in derart annotierten Korpora nicht (vgl. Lüdeling/Walter/Kroymann/Adolphs 2005; Lüdeling 2007). Ein großes Problem ist, dass alternative Analysen für die gleichen Daten auf diese Weise nicht annotiert werden können. Für ein kleines, spezialisiertes Korpus wie Falko mit Daten, deren Interpretation oft umstritten sein kann (wie zum Beispiel die Ziellhypothese, s. Abschn. 2), bietet sich daher eine Korpusarchitektur an, in der die Annotation getrennt von den eigentlichen Daten gespeichert werden. Mit einer solchen Mehrebenenarchitektur können zu den Daten beliebig viele voneinander unabhängige Annotationsebenen einge

¹ Dies gilt für alle uns bekannten Lerner korpora außer dem LEAP Korpus (vgl. Milde/Gut 2002).

² Wir verwenden zur Annotation das an der Universität Hamburg entwickelte Programm EXMARALDA (vgl. Schmidt/Wörner 2005).

hinaus auch einer Qualitätssicherung in der deskriptiven Analyse.

3 Die Analyse von Lernersprachen – zwei Beispiele

Tokens	Und	immer	noch	kann	man	eine	unzufriedenheit	spüren
Lemma	und	immer	noch	können	man	ein	unknown	spüren
Wortart	KON	ADV	ADV	VMFIN	PIS	ART	NN	VVINF
Zielhypothese							Umzufriedenheit	
Topologische Felder	Vorfeld		linker Satzklammer		Mittelfeld		rechte Satzklammer	

Tab. 2: Schematische Darstellung eines Ausschnitts aus dem annotierten Lernertext Falko Essays L2 1.0, fk024, 247

die im Mittelfeld eines Konstituentensatzes erscheinen) herausgefiltert werden.

In Tab. 2 wird eine Lerneräußerung auf den verschiedenen Ebenen dargestellt. Die Notation wird in der Dokumentation, die sich auf der Falko-Homepage befindet, ausführlich beschrieben. In der ersten Zeile befinden sich die einzelnen Tokens, in der zweiten und dritten Zeile das Ergebnis der automatischen Lemmatisierung bzw. der Wortartenerkennung. Vom Lemmatisierer nicht erkannte Einheiten (z.B. *unzufriedenheit*) werden als „unknown“ markiert. Die Ziellhypothese und die topologischen Felder werden manuell in den beiden letzten Zeilen hinzugefügt.

2 Ziellhypothese und Fehlerannotation

Wie in der Einleitung bereits angesprochen, kann man die Lernendaten verwenden, um Abweichungen (Fehler) zu untersuchen und um quantitative Eigenschaften der Lerner- sprache zu ermitteln (vgl. z.B. Granger 2002: 11ff.; Walter/Grommes 2008a: 15ff.). Für die Fehlerforschung ist eine Fehlerannotation hilfreich, mithilfe deren Fehler eines bestimmten Typs zuverlässig gefunden werden können.

An dieser Stelle wollen wir die Diskussion um den Begriff „Fehler“ nicht aufgreifen (vgl. aber z.B. Cherubim 1980; Ellis 1994: 50ff.; Lennon 1991). Wir definieren „Fehler“ als Abweichung von einer Ziellhypothese, welche von uns auf der Basis der Lerneräußerung rekonstruiert wird (s. unten). Die Entwicklung von Fehler-Tagsets kann nach unterschiedlichen Kriterien (linguistische Ebene, formaler FehlerTyp, Prinzipienverletzung etc.) und in unterschiedlicher Granularität erfolgen. Die meisten Fehler-Tagsets, wie z.B. das Tagset für ICLÉ (vgl. Dagneaux/Denness/Granger/Meunier 1996) oder das Tagset in Weinberger (vgl. 2002), sind Mischungen aus mehreren Kriterien und daher nicht immer leicht nachvoll-

Genre bezogenen ausgewertet, mit Daten aus dem L2-Essay-Korpus kontrastiert und bezüglich ausgewählter Kontextmerkmale analysiert. Besonders häufig formulieren die Lerner mit Hilfe der Konstruktion in den Essays rhetorische Fragen, wie das folgende Beispiel illustriert:

(2) Und wo haben diese berühmte Menschen ihre Ideen entwickelt, wenn nicht in der Universität? (Falko Essays L2 1.0, fk012, 630)

Weitere Faktoren, z.B. das Auftreten der Konstruktion in einer elliptischen Struktur oder aber die unterschiedliche Verwendung in Konstituent- bzw. Matrixsätzen, werden in der Studie ausgewertet und diskutiert; für eine vertiefte Darstellung der Ergebnisse vgl. Walter/Schmidt (i.Dr.). Studien wie diese zeigen, dass Lernerkorpora verwendet werden können, um Präferenz (Overuse), aber auch Vermeidung (Underuse) von ausgewählten Konstruktionen bei den Lernern festzustellen. Die Ergebnisse dieser Analysen können beispielweise in der Lernerlexikographie genutzt werden, indem Lerner in Lernerwörterbüchern explizit auf diese lernersprachlichen Besonderheiten aufmerksam gemacht werden. Für das Englische wurde dies auch schon richtungweisend im Macmillan English Dictionary für das satzinitiale *and* in der Auswertung der ICLE-Daten vollzogen (vgl. Rundell/Fox 2007: IW4).

3.2 Orthographie

Während die Orthographie lange als „idiosynkratisch“ oder als für strukturelle Unterschiede uninteressant galt, wurden in den letzten Jahren mehr und mehr Regelmäßigkeiten gefunden und Verbindungen zu den anderen grammatischen Gebieten gezeigt. Da-

her kann die Orthographie auch in der Sprachentwicklungsorschung interessante Hinweise auf phonologische, morphologische und syntaktische Hypothesen der Lerner liefern. In Hirschmann (i.Vorb.) werden die Orthographiefehler im L2-Essay-Subkorporus analysiert und mit L1-Vergleichsdaten¹ kontrastiert. Derzeit wird darauf aufbauend ein Fehler-Tagset bzw. ein Annotationsschema entwickelt, um eine effiziente automatische Abfrage der Orthographiefehler in Falko zu ermöglichen.

Die Orthographiefehler werden zunächst in kontextunabhängige (Wortschreibungs-)Fehler und kontextabhängige (syntaktische) Fehler kategorisiert. So beinhaltet die Schreibunterschiede zwischen dem norwegischen ASK-Korpus (<http://decentius.saksis.uib.no/corpus/askdeu/home.xml>). Da diese Korpora keine Mehrebenenarchitektur haben, kann jeweils nur eine Ziellhypothese angegeben werden.

¹ Hierzu wurde das auf S. 68 in Fußnote 3 bereits angeführte Korpus Akademisches Deutsch 2006 verwendet.
² Als L1-Vergleichskorpus dienen die Daten aus Falko Essay L1 0,5, kann die Leistungskurschülern des Faches Deutsch produziert wurden.

gen <nammen> und <kontrollieren> in (1) orthographische Fehler, die ohne den syntaktischen Kontext diagnostiziert werden können. Die Wortfolge <in den> in (3) hingegen ist nur unter der Berücksichtigung eines Kontextes (wie in Beispiel (3)) als (Getrenntschreibungs-)Fehler zu identifizieren:

- (3) ... aber auf dieser mitleidlosen Welt sollen sie ihre Zierlichkeit nicht verlieren, in dem sie alles versuchen zu machen. (Falko Essays 1.0, trk001, 369)

Die grundlegende Unterscheidung zwischen diesen beiden Fehlertypen wird häufig übergangen (vgl. z.B. Eisenberg 2004: 332), ist jedoch gerade hinsichtlich des Schriftspracherfahrlieferfehler auf (das Verhältnis brägt 92 zu 19), wohingegen im L1-Korpus die Relation umgekehrt (67 zu 101) ist.

In Hirschmann (i.Vorb.) werden verschiedene Fehlerkategorien qualitativ und quantitativ analysiert und lerntheoretische Erklärungsvorschläge unterbreitet. Studien wie diese zeigen, dass korpusbasierte (Fehler-)Untersuchungen zu neuen, differenzierten Ergebnissen führen können: Wortschreibungsfehler sind oft eigentlich phonologische oder morphologische Fehler (sehr oft Transferfehler), während kontextabhängige Fehler auf syntaktische Probleme hinweisen können. In der Vermittlung können solche Probleme dann gezielt besprochen werden.

4 Zusammenfassung

Kontrolliert erhobene und transparent ammontierte Lerner korpora werden in der Untersuchung der Lernersprache immer wichtiger (vgl. Fandrych/Tschirner 2007). Bestimmt komplexe Erwerbsphe nomen, die quantitative und qualitative Merkmale von Lernern betreffen, können wahrscheinlich am besten anhand von Lernerkorpus-Studien erforscht werden. Die aus den Korpusstudien gewonnenen Hypothesen müssen dann oft in gezielten Erhebungen überprüft werden. Wir denken zum Beispiel an integrierte Registeruntersuchungen, die mehrere Eigenschaften der Lernersprache gleichzeitig betreffen (das oben beschriebene „mündliche“ satzinitiale *und* kann dann mit weiteren Beobachtungen verbunden werden).

Das in diesem Artikel vorgestellte Lerner korpus Falko kann aufgrund seiner Designkriterien und seiner vielschichtigen Annotation mit einer expliziten Zielhypothese für viele Zwecke genutzt werden. Das Korpus ist zwar noch klein, aber es wächst!¹

¹ Über Unterstützung bei der Datenerhebung würden wir uns sehr freuen.

Literatur

- Milde, Jan-Thorsten/Gut, Ulrike (2002): A prosodic corpus of non-native speech. In: B. Bel/ I. Marien (Hg.): Proceedings of the Speech Prosody 2002 Conference. Aix-en-Provence, 503–506.
- Nesselhauf, Nadja (2004): Learner Corpora and their Potential in Language Teaching. In: J. Sinclair (Hg.): How to Use Corpora in Language Teaching. Amsterdam, 125–152.
- Römer, Ute (i.Dr.): Corpora and Language Teaching. In: A. Lüdeling (M. Kyö (Hg.); Rundell, Michael/Fox, Gwyneth (Hg.): (2007): Macmillan English Dictionary For Advanced Learners. Second edition. Oxford.
- Schmid, Helmut (1994): Probabilistic Part-of-Speech Tagging Using Decision Trees. In: Proceedings of the International Conference on New Methods in Language Processing. Manchester, 44–49 (erweiterte Fassung verfügbar unter: <http://www.ims.uni-stuttgart.de/fip/pubs/corpora/tree-tagger1.pdf>).
- Schmid, Helmut (i.Dr.): Tokenizing and part-of-speech tagging. In: A. Lüdeling (M. Kyö (Hg.); Schmidt, Thomas/Wörner, Kai (2005): Erstellen und Analysieren von Gesprächskorpora mit EX-MARaLDA. In: Gesprächsforschung. Online-Zeitschrift zur verbalen Interaktion 6, 171–195 (online verfügbar unter: <http://www.ge sprachsforschung-ozs.de/heft2005/jpx-woerner.pdf>).
- Stifter, Clara/Stern, William (1907): Die Kindersprache. Eine psychologische und sprachtheoretische Untersuchung. Leipzig.
- von Stutterheim, Christiane/Carroll, Mary (2006): The Impact of Grammatical Temporal Categories on Ultimate Attainment in L2 Learning. In: H. Byrnes et al. (Hg.): Educating for Advanced Foreign Language Capacities. Washington, D.C., 40–53.
- Walter, Maike/Grommes, Patrick (2008a): Die Entwicklung des fortgeschrittenen Lernens in der Varietätenlinguistik. In: M. Walter/P. Grommes (Hg.) (2008b), 3–27.
- Walter, Maike/Grommes, Patrick (Hg.): (2008b): Fortgeschrittenen Lernvarietäten. Zweitspracherwerbsforschung und Korpuslinguistik. Tübingen.
- Walter, Maike/Schmidt, Karin (i.Dr.): *Und das ist auch gut so!* Der Gebräuch des satzinitialen und bei fortgeschrittenen Lernern des Deutschen als Fremdsprache. In: B. Ahrenholz et al. (Hg.), Wegener, Heide (1995): Die Nominalflexion des Deutschen – verstanden als Lerngegenstand. Tübingen.
- Lüdeling, Anke (2007): Das Zusammenspiel von qualitativen und quantitativen Methoden in der Korpuslinguistik. In: W. Kallmeyer/G. Zifonun (Hg.): Sprachkorpora – Datenmengen und Erkenntnisfortschritt. Berlin/New York, 28–48 (Institut für Deutsche Sprache, Jahrbuch 2006).
- Lüdeling, Anke (2008): Mehrdeutigkeiten und Kategorisierung. Probleme bei der Annotation von Lerner korpora. In: M. Walter/P. Grommes (Hg.) (2008b), 119–140.
- Lüdeling, Anke/Kyö, Merija (Hg.): Corpus Linguistics. An International Handbook. Berlin/New York.
- Dietrich, Rainer (2006): Second Language Acquisition in the 20th Century. In: S. Autroux et al. (Hg.), History of the Language Sciences. Vol. 3. Berlin, 2705–2728.
- Eisenberg, Peter (2004): Grundriß der deutschen Grammatik. Das Wort. 2. Aufl. Stuttgart/Weimar.
- Ellis, Rod (1994): The Study of Second Language Acquisition. Oxford.
- Fandrych, Christian/Tschirner, Erwin (2007): Korpuslinguistik und Deutsch als Fremdsprache. Ein Perspektivenwechsel. In: DaF 4, 195–204.
- Granger, Sylviane (2002): A bird's-eye view of learner corpus research. In: S. Granger/J. Hung/S. Peitch-Tyson (Hg.): Computer Learner Corpora, Second Language Acquisition and Foreign Language Teaching. Amsterdam/Philadelphia, 3–33.
- Granger, Sylviane (2003): The International Corpus of Learner English. A New Resource for Foreign Language Learning and Teaching and Second Language Acquisition Research. In: TESOL Quarterly 37/3, 538–546.
- Granger, Sylviane (i.Dr.): Learner Corpora. In: A. Lüdeling/M. Kyö (Hg.).
- Granger, Sylviane/Dagneau, Estelle/Meunier, Fanny (Hg.): (2002): The International Corpus of Learner English. Louvain-la-Neuve.
- Hirschmann, Hagen (i.Vorb.): Orthographische Kompetenz zwischen Deutsch als L1 und Deutsch als L2. Hähle, Tilman (1986): Der Begriff „Mittelfeld“. Anmerkungen über die Theorie der topologischen Felder. In: Akten des Siebten Internationalen Germanistikkongresses. Göttingen, 329–340.
- Hunston, Susan (i.Dr.): Collection strategies and design decisions. In: A. Lüdeling (M. Kyö (Hg.); Izumi, Emi et al. (2005): Error Annotation for a Corpus of Japanese Learner English. In: Proceedings of the Sixth International Workshop on Linguistically Interpreted Corpora (LINC 2005). Jeju Island, Südkorea. ACL Anthology IOC 6009, 72–80 (online verfügbar unter: <http://acl.ldc.upenn.edu/I/IOC105-6009.pdf>).
- Lemmon, Paul (1991): Error and the very advanced learner. In: International Review of Applied Linguistics 29/1, 31–44.
- Lüdeling, Anke (2007): Das Zusammenspiel von qualitativen und quantitativen Methoden in der Korpuslinguistik. In: W. Kallmeyer/G. Zifonun (Hg.): Sprachkorpora – Datenmengen und Erkenntnisfortschritt. Berlin/New York, 28–48 (Institut für Deutsche Sprache, Jahrbuch 2006).
- Cherubini, Dieter (Hg.) (1980): Fehlerlinguistik. Beiträge zum Problem der sprachlichen Abweichung. Tübingen.
- Dagneau, Estelle et al. (1996): Error Tagging Manual Version 1.1. Centre for English Corpus Linguistics, Université Catholique de Louvain. Louvain-la-Neuve.
- Baumgarten, Nicole (2006): Konventionen der Kohäsionsbildung in deutschen und englischen Texten. AND und UND als makrosyntaktische Verknüpfung in populärwissenschaftlichen Zeitschriftenartikeln. In: D. Wolff (Hg.): Mehrsprachige Individuen – vielsprachige Gesellschaften. Frankfurt a.M., 133–153.
- Bell, David (2007): Sentence-initial *and* and *but* in academic writing. In: Pragmatics 17/2, 183–201.
- Belz, Judy A. (2004): Learner Corpus Analysis and the Development of Foreign Language Proficiency. In: System. An International Journal of Educational Technology and Applied Linguistics 32/4, 577–591.

Grammatik. Das Wort. 2. Aufl. Stuttgart/Weimar.

Ellis, Rod (1994): The Study of Second Language Acquisition. Oxford.

Fandrych, Christian/Tschirner, Erwin (2007): Korpuslinguistik und Deutsch als Fremdsprache. Ein Perspektivenwechsel. In: DaF 4, 195–204.

Granger, Sylviane (2002): A bird's-eye view of learner corpus research. In: S. Granger/J. Hung/S. Peitch-Tyson (Hg.): Computer Learner Corpora, Second Language Acquisition and Foreign Language Teaching. Amsterdam/Philadelphia, 3–33.

Granger, Sylviane (2003): The International Corpus of Learner English. A New Resource for Foreign Language Learning and Teaching and Second Language Acquisition Research. In: TESOL Quarterly 37/3, 538–546.

Granger, Sylviane (i.Dr.): Learner Corpora. In: A. Lüdeling/M. Kyö (Hg.).

Granger, Sylviane/Dagneau, Estelle/Meunier, Fanny (Hg.): (2002): The International Corpus of Learner English. Louvain-la-Neuve.

Hirschmann, Hagen (i.Vorb.): Orthographische Kompetenz zwischen Deutsch als L1 und Deutsch als L2.

Hähle, Tilman (1986): Der Begriff „Mittelfeld“. Anmerkungen über die Theorie der topologischen Felder. In: Akten des Siebten Internationalen Germanistikkongresses. Göttingen, 329–340.

Hunston, Susan (i.Dr.): Collection strategies and design decisions. In: A. Lüdeling (M. Kyö (Hg.); Izumi, Emi et al. (2005): Error Annotation for a Corpus of Japanese Learner English. In: Proceedings of the Sixth International Workshop on Linguistically Interpreted Corpora (LINC 2005). Jeju Island, Südkorea. ACL Anthology IOC 6009, 72–80 (online verfügbar unter: <http://acl.ldc.upenn.edu/I/IOC105-6009.pdf>).

Lemmon, Paul (1991): Error and the very advanced learner. In: International Review of Applied Linguistics 29/1, 31–44.

Lüdeling, Anke (2007): Das Zusammenspiel von qualitativen und quantitativen Methoden in der Korpuslinguistik. In: W. Kallmeyer/G. Zifonun (Hg.): Sprachkorpora – Datenmengen und Erkenntnisfortschritt. Berlin/New York, 28–48 (Institut für Deutsche Sprache, Jahrbuch 2006).

Cherubini, Dieter (Hg.) (1980): Fehlerlinguistik. Beiträge zum Problem der sprachlichen Abweichung. Tübingen.

Dagneau, Estelle et al. (1996): Error Tagging Manual Version 1.1. Centre for English Corpus Linguistics, Université Catholique de Louvain. Louvain-la-Neuve.

Baumgarten, Nicole (2006): Konventionen der Kohäsionsbildung in deutschen und englischen Texten. AND und UND als makrosyntaktische Verknüpfung in populärwissenschaftlichen Zeitschriftenartikeln. In: D. Wolff (Hg.): Mehrsprachige Individuen – vielsprachige Gesellschaften. Frankfurt a.M., 133–153.

Bell, David (2007): Sentence-initial *and* and *but* in academic writing. In: Pragmatics 17/2, 183–201.

Belz, Judy A. (2004): Learner Corpus Analysis and the Development of Foreign Language Proficiency. In: System. An International Journal of Educational Technology and Applied Linguistics 32/4, 577–591.

Birkner, Karin et al. (1995): Der adversative Konjunkt *aber* eines italienischen und zweier polnischer Lerner des Deutschen. In: B. Handwerker (Hg.): Fremde Sprache. Tübingen, 65–118.

Cherubini, Dieter (Hg.) (1980): Fehlerlinguistik. Beiträge zum Problem der sprachlichen Abweichung. Tübingen.

Dagneau, Estelle et al. (1996): Error Tagging Manual Version 1.1. Centre for English Corpus Linguistics, Université Catholique de Louvain. Louvain-la-Neuve.

Dietrich, Rainer (2006): Second Language Acquisition in the 20th Century. In: S. Autroux et al. (Hg.), History of the Language Sciences. Vol. 3. Berlin, 2705–2728.

Eisenberg, Peter (2004): Grundriß der deutschen Grammatik. Das Wort. 2. Aufl. Stuttgart/Weimar.

Ellis, Rod (1994): The Study of Second Language Acquisition. Oxford.

Fandrych, Christian/Tschirner, Erwin (2007): Korpuslinguistik und Deutsch als Fremdsprache. Ein Perspektivenwechsel. In: DaF 4, 195–204.

Granger, Sylviane (2002): A bird's-eye view of learner corpus research. In: S. Granger/J. Hung/S. Peitch-Tyson (Hg.): Computer Learner Corpora, Second Language Acquisition and Foreign Language Teaching. Amsterdam/Philadelphia, 3–33.

Granger, Sylviane (2003): The International Corpus of Learner English. A New Resource for Foreign Language Learning and Teaching and Second Language Acquisition Research. In: TESOL Quarterly 37/3, 538–546.

Granger, Sylviane (i.Dr.): Learner Corpora. In: A. Lüdeling/M. Kyö (Hg.).

Granger, Sylviane/Dagneau, Estelle/Meunier, Fanny (Hg.): (2002): The International Corpus of Learner English. Louvain-la-Neuve.

Hirschmann, Hagen (i.Vorb.): Orthographische Kompetenz zwischen Deutsch als L1 und Deutsch als L2.

Hähle, Tilman (1986): Der Begriff „Mittelfeld“. Anmerkungen über die Theorie der topologischen Felder. In: Akten des Siebten Internationalen Germanistikkongresses. Göttingen, 329–340.

Hunston, Susan (i.Dr.): Collection strategies and design decisions. In: A. Lüdeling (M. Kyö (Hg.); Izumi, Emi et al. (2005): Error Annotation for a Corpus of Japanese Learner English. In: Proceedings of the Sixth International Workshop on Linguistically Interpreted Corpora (LINC 2005). Jeju Island, Südkorea. ACL Anthology IOC 6009, 72–80 (online verfügbar unter: <http://acl.ldc.upenn.edu/I/IOC105-6009.pdf>).

Lemmon, Paul (1991): Error and the very advanced learner. In: International Review of Applied Linguistics 29/1, 31–44.

Lüdeling, Anke (2007): Das Zusammenspiel von qualitativen und quantitativen Methoden in der Korpuslinguistik. In: W. Kallmeyer/G. Zifonun (Hg.): Sprachkorpora – Datenmengen und Erkenntnisfortschritt. Berlin/New York, 28–48 (Institut für Deutsche Sprache, Jahrbuch 2006).

Cherubini, Dieter (Hg.) (1980): Fehlerlinguistik. Beiträge zum Problem der sprachlichen Abweichung. Tübingen.

Dagneau, Estelle et al. (1996): Error Tagging Manual Version 1.1. Centre for English Corpus Linguistics, Université Catholique de Louvain. Louvain-la-Neuve.

Baumgarten, Nicole (2006): Konventionen der Kohäsionsbildung in deutschen und englischen Texten. AND und UND als makrosyntaktische Verknüpfung in populärwissenschaftlichen Zeitschriftenartikeln. In: D. Wolff (Hg.): Mehrsprachige Individuen – vielsprachige Gesellschaften. Frankfurt a.M., 133–153.

Bell, David (2007): Sentence-initial *and* and *but* in academic writing. In: Pragmatics 17/2, 183–201.

Belz, Judy A. (2004): Learner Corpus Analysis and the Development of Foreign Language Proficiency. In: System. An International Journal of Educational Technology and Applied Linguistics 32/4, 577–591.