

Prosodische Etikettierung des Deutschen mit ToBI

Matthias Reyelt, Martine Grice, Ralf Benzmüller, Jörg Mayer, Anton Batliner

Angaben zur Veröffentlichung / Publication details:

Reyelt, Matthias, Martine Grice, Ralf Benzmüller, Jörg Mayer, and Anton Batliner. 1996. "Prosodische Etikettierung des Deutschen mit ToBI." In *Natural language processing and speech technology: results of the 3rd KONVENTS Conference, Bielefeld, October 1996*, edited by Dafydd Gibbon, 144–55. Berlin: de Gruyter. <https://doi.org/10.1515/9783110821895-016>.

Nutzungsbedingungen / Terms of use:

licgercopyright



15 Prosodische Etikettierung des Deutschen mit ToBI

*Matthias Reyelt, Martine Grice, Ralf Benzmüller,
Jörg Mayer, Anton Batliner*

For the automatic analysis of prosody, large hand labelled speech databases are needed. The production of such databases for the German language requires the development of a multi-site labelling system for German prosody which, for practical reasons, can be learned relatively quickly. To accomplish this task, a consensus prosodic transcription system for German, based on the English ToBI-System, has been developed by members of the Universities of Stuttgart, Saarbrücken, Braunschweig and Munich. This system, GToBI, was then evaluated in a transcription experiment.

Die automatische Analyse prosodischer Information erfordert umfangreiche Datenbasen etikettierter Sprache. Um diese zu erzeugen, muß ein Standard-Transkriptionssystem entwickelt werden, mit dem Transkribenten nach kurzer Einarbeitung Etikettierungen von guter Qualität erstellen können. Ein gemeinsames System der Universitäten Stuttgart, Saarbrücken, Braunschweig und München wird hier vorgestellt, das sich am englischen ToBI-System orientiert. Zur Evaluierung des Systems wurde ein Transkriptionsexperiment durchgeführt, bei dem teils Experten, teils wenig erfahrene Testpersonen das gleiche Material parallel etikettierten.

15.1 Einleitung

In den letzten Jahren hat die automatische Spracherkennung große Fortschritte gemacht und ist dabei, den Übergang vom Einzelworterkennner zum sprachverstehenden System zu vollziehen.¹ Mittlerweile erreichen statistische Worterkennner auch für Spontansprache Erkennungsraten, die eine weitgehende syntaktisch–semantische Analyse erlauben. Problematisch ist dabei, daß die zu analysierende Wortkette ohne zusätzliche Information nur schwer in syntaktische Einheiten segmentierbar ist. Kompe et al. (1995) geben als Beispiel die folgende Wortkette:

ja | zur | Not | geht's | auch | am Samstag |

Die Balken geben dabei mögliche Satzgrenzen an, die in geschriebener Sprache durch Komma, Punkt oder Fragezeichen gekennzeichnet würden.

¹Ein Teil dieser Arbeiten wurde mit Mitteln des Bundesministers für Bildung und Forschung unter den Förderkennzeichen 01 IV 101 N/0 und 01 IV 102 F/4 gefördert. Die Verantwortung für den Inhalt liegt bei den Autoren. Die Sprachdaten für dieses Experiment sowie GToBI–Trainingsmaterial ist via ftp verfügbar. Zur Visualisierung der Daten kann das Programm *fish* ebenfalls via ftp bezogen werden.

Insgesamt geben Kompe et. al. 36 (!) syntaktisch korrekte Alternativen an, zwei davon sind:

*Ja? Zur Not geht's? Auch am Samstag?
Ja. Zur Not. Geht's auch am Samstag?*

Da Satzgrenzen häufig prosodisch – sei es durch Intonation, Längung oder Pausen – markiert sind, kann eine automatische Prosodieanalyse hier wertvolle zusätzliche Information liefern. Batliner et al. (1996) zeigen zum Beispiel, daß sich die Dauer des syntaktischen Parsing teilweise stark verkürzt, wenn in das Wortgitter zusätzlich prosodische Information über Phrasengrenzen eingefügt wird. Batliner et. al. weisen auf die Bedeutung der Intonation für die Disambiguierung des Satzmodus hin (Batliner et al. 1992). Die Integration eines Prosodiemoduls in ein automatisches Übersetzungssystem wird z.B. im Projekt VERBMOBIL (Wahlster 1993) durchgeführt. Das Prosodiemodul reichert dabei den Worthypothesengraphen durch Informationen über prosodische Grenzen und Akzente an, bevor dieser syntaktisch-semantisch analysiert wird (vgl. auch Hess et al. 1996).

Allerdings wird zum Training und Test statistischer Erkennner Datenmaterial in großen Mengen benötigt. Dieses Material muß zuvor von Hand prosodisch etikettiert werden, eine Arbeit, die notwendigerweise von mehreren trainierten Hilfskräften – evtl. sogar an unterschiedlichen Orten – durchgeführt wird. Trotzdem müssen die dabei entstehenden Transkriptionen „einheitlich“ sein, d.h. die Etiketten müssen von verschiedenen Personen gleich angewandt werden. Eine völlige Übereinstimmung der Transkriptionen, ganz gleich welcher Art, ist allerdings illusorisch. Daher sind Untersuchungen zur Konsistenz wichtig, um den Grad der Übereinstimmung abzuschätzen. Derartige Untersuchungen wurden von Tillmann et al. (1992) für phonetisch enge Transkriptionen und von Reyelt (1994) für prosodische Etiketten durchgeführt und sind ein wichtiges Maß für die Qualität der Verfahren.

Ein Inventar prosodischer Etiketten, das für die Bearbeitung großer Korpora geeignet sein soll, muß daher neben der prinzipiellen phonologisch adäquaten Beschreibung noch weitere Anforderungen erfüllen (vgl. hierzu auch die in Silverman et al. (1992a) geforderten Kriterien):

Das Inventar muß relativ schnell erlernbar sein. Es muß einheitliches Trainingsmaterial geben, das die Anwendung der Etiketten möglichst umfassend an repräsentativen Sprachbeispielen beschreibt, und zwar so, daß das Training der Transkribenten an unterschiedlichen Orten durchgeführt werden kann. Um die Verarbeitung der Daten mit statistischen Methoden zu ermöglichen, müssen die Datenformate maschinell verarbeitbar sein. Weiterhin sollten Evaluierungen der Transkriptionen durchgeführt werden, um Anhaltspunkte über deren Qualität zu erhalten und ggf. Schwachpunkte zu verbessern.

15.2 Das ToBI-System

Eine Beschreibung der Intonation durch eine Folge hoher und tiefer Töne, wie sie im von Pierrehumbert (1980) entwickelten Tonsequenzansatz verwendet wurde, bildet den Ausgangspunkt für eine Reihe prosodischer Beschreibungen. Anfang der neunziger Jahre wurde in den USA unter dem Namen ToBI (Tone and Break Indices) ein einheitlicher Standard für ein Inventar prosodischer Etiketten entwickelt, um die dort existierenden prosodisch etikettierten Datenbasen zusammenzufassen (Silverman et al. 1992a,b). Neben Tonakzenten, Phrasenakzenten (intermediärer Grenzton) und Grenztönen werden prosodische Grenzen verschiedener Stärke (*intermediate* und *international phrase boundaries*) markiert. Das ToBI-System für das Englische wurde in umfangreichen Transkriptionsexperimenten evaluiert; Pitrelli et al. (1994) geben für Tonakzente eine Übereinstimmung von 68% an, für Phrasenakzente 85% und für Grenztöne 91%.

Für das Deutsche wurden an mehreren Instituten unabhängig voneinander Adaptionen des ToBI-Systems durchgeführt, u.a. in Saarbrücken, Stuttgart und Braunschweig/München (VERBMOBIL). Auf einem Workshop in Stuttgart im Februar 1995 wurden anhand von Sprachdaten, die an den drei Instituten etikettiert wurden, die Beschreibungssysteme verglichen und Unterschiede bzw. Gemeinsamkeiten herausgearbeitet. Auf dem *One Day Workshop on Prosodic Labelling* im August 1995 in Stockholm zeigte sich, daß auch in Hamburg und Nijmegen prosodische Etikettierungen des Deutschen nach dem ToBI-System durchgeführt werden.² Auf diesem Workshop wurde auch der Vergleich mit einem konturbasierten Ansatz (PROLAB, vgl. Kohler 1995) vorgenommen. Auf einem weiteren Workshop in Stuttgart wurde ein gemeinsames ToBI-System für das Deutsche (GToBI) entworfen und ein Trainingskorpus zusammengestellt. Weiterhin sollte das System in einem ersten Transkriptionsexperiment überprüft werden, bei dem das Sprachmaterial aus einer Mischung unterschiedlicher Korpora bestand.

15.3 Teilnehmer und Korpora

Die Teilnehmer entwickelten ihre prosodischen Etikettensysteme anhand unterschiedlicher Sprachdaten mit unterschiedlichen Zielsetzungen:

Das Saarbrücker System (Grice und Benzmüller 1995) wurde zur Transkription spontansprachlicher Dialoge entworfen. Es handelt sich um task-orientierte Dialoge (vgl. HCRC Map Task (Anderson et al. 1991)), in denen die Versuchspersonen ohne Sichtkontakt (d.h. nur über den akustischen Kanal) Informationen über die Route auf einer Landkarte austauschen. Dadurch, daß sich die Versuchspersonen kennen und die Aufgabe spielerisch

²In der Vorbereitungsphase dieses Workshops wurde unter den deutschen Teilnehmern eine Umfrage über die verwendeten Etikettensysteme durchgeführt. Die Ergebnisse sind in Mayer (1995a) zusammengefaßt.

ausgeführt wird, sind relativ viele stilisierte Intonationskonturen (zu Stilisierung (*stylized contours*) vgl. Gibbon 1976) enthalten. Parallel dazu werden vergleichbare Dialoge in verschiedenen Varietäten des Deutschen, Italienischen und Bulgarischen aufgenommen.

In Stuttgart wird hauptsächlich gelesene Sprache transkribiert. Das Transkriptionssystem (Mayer 1995b) wurde anhand konstruierter, vorgelesener Sätze entwickelt und wird nun vor allem zur Etikettierung von Radionachrichten verwendet. Die Leistungsfähigkeit des Stuttgarter Systems bei der Transkription von Spontansprache wird derzeit getestet.

Im Projekt VERBMOBIL werden Terminabsprachen zwischen zwei Versuchspersonen aufgenommen. Die Dialoge sind zwar spontan, die Dialogpartner müssen jedoch einen Knopf drücken, bevor sie sprechen können. Die Daten werden in Kiel, Bonn, München und Karlsruhe aufgenommen und in Braunschweig prosodisch etikettiert; das dazu verwendete Transkriptionssystem (Reyelt und Batliner 1994) wurde in Zusammenarbeit mit Projektpartnern in München und Erlangen entwickelt.

Es zeigte sich, daß sich die unterschiedlichen Schwerpunkte der einzelnen Transkriptionssysteme gut ergänzen. Während das VERBMOBIL-System auf vielen nord- und süddeutschen Varietäten basiert, sind im Saarbrücker Transkriptionssystem stilisierte Konturen sehr genau beschrieben. Das Stuttgarter Transkriptionssystem wiederum ist stärker phonologisch orientiert.

15.4 Beschreibung des GToBI-Systems

In GToBI werden die Tonhöhenverläufe mit zwei verschiedenen Tönen beschrieben: H und L. H steht für einen hohen und L für einen tiefen Zielpunkt (*target*) eines Tonhöhenverlaufs. Die Rekonstruktion der beschriebenen Kontur erfolgt durch regelhafte Interpolation zwischen den Zielpunkten. Die Verknüpfung der tonalen mit der textuellen Ebene geschieht über Diakritika. *Akzenttöne* werden durch einen „*“ (*Sternchen*) mit der akzentuierten Silbe verknüpft, der *Phrasenton* bei intermediären Phrasen (B3) wird durch „–“ (*minus*) an das letzte Wort der Phrase gebunden. Grenzen von Intonationsphrasen (B4) werden immer aus zwei Tönen gebildet, dem *Phrasenton* und einem *Grenzton*, der mit „%“ gekennzeichnet wird. Der Grenzton beschreibt dabei den Tonverlauf am Ende der Intonationsphrase, der Phrasenton den Verlauf zwischen letztem Akzentton und Grenzton.

15.4.1 Tonakzente

Die folgenden Tonakzente können auftreten:

H* : Dies ist der normale *Gipfelakzent*. Er zeichnet sich durch eine Abweichung der Tonhöhe nach oben aus. Folgen zwei oder mehr H*-Akzente aufeinander, so kann die Tonhöhe zwischen diesen leicht abfallen (*sagging*).

- L+H* : Bei diesem Akzent findet auf der akzentuierten Silbe ein steiler Anstieg statt, so daß der Höhepunkt der Bewegung erst recht spät in der akzentuierten Silbe stattfindet. Dieser Akzent tritt häufig als emphatischer oder kontrastiver Akzent auf.
- L*+H : Hier gibt es einen tiefen Zielpunkt früh in der akzentuierten Silbe. Der Gipfel ist hinter die akzentuierte Silbe verschoben (vgl. *scooped accent* bei Ladd et al. (1983)).
- L* : Dieser Akzenttyp zeichnet sich durch eine Abweichung der Tonhöhe auf der akzentuierten Silbe nach unten im unteren Bereich des Stimmumfangs aus, *Talakzent*.
- H+L* : Bei diesem Akzent fällt die Tonhöhe von einem hohen Zielpunkt vor der Akzentsilbe in den unteren Bereich des Stimmumfangs ab (*early peak* bei Féry (1993)).
- H+!H* : Dieser Akzent ähnelt dem H+L*, fällt aber nur bis in den mittleren Bereich des Stimmumfangs ab.

15.4.2 Tonale Grenzmarkierungen

Während Grenzen zwischen intermediären Phrasen nur durch einen *Phrasenton* (also H- bzw. L-) gekennzeichnet werden, sind Grenzen an Intonationsphrasen immer bitonal durch *Phrasenton* und *Grenzton* markiert.³ Zusammen mit H- und L- an intermediären Phrasengrenzen gibt es sechs mögliche tonale Grenzmarkierungen:

- L-L% : terminaler Fall, Grenze im unteren Bereich des Stimmumfangs.
- L-H% : Nach Gipfelakzent fallend-steigende Grenze, nach L* leichter Anstieg in der letzten Silbe vor der Grenze.
- H-L% : ebene Grenze im mittleren Bereich des Stimmumfangs, „progredient“.
- H-H% : Anstieg nach dem letzten Akzent bis in die obersten Bereiche des Stimmumfangs, „interrogativ“.

15.4.3 Downstepping

In einer Folge von Gipfelakzenten können diese jeweils „treppenartig“ ein Stück nach unten versetzt sein. Dieser Effekt wird als *Downstepping* bezeichnet und diakritisch durch ein dem hohen Ton des Akzents vorgestelltes „i“ etikettiert. Weiterhin kann Downstepping auch an hohen Phrasentönen markiert werden.

15.4.4 Upstepping

Ein treppenartiger Anstieg der Akzente wird durch ein vorgestelltes „^“ etikettiert.

³Im Englischen *phrase accent* und *boundary tone*. Ladd (1995) führt allgemein für tonale Markierung prosodischer Grenzen den Begriff *edge tone* ein.

15.4.5 Grenztypen (break indices)

Die Grenztypen stellen ein (subjektives) Maß für die „Stärke“ der Trennung dar.

- B1 : normale Wortgrenze. Sie trägt *keine* prosodische Markierung und wird auch normalerweise nicht etikettiert. Der Vollständigkeit halber ist sie hier trotzdem aufgeführt.
- B2 : irreguläre Grenze, Abbruch, Häsitation.
- B3 : intermediäre Phrasengrenze. Leichte tonale Markierung innerhalb einer Intonationsphrase. B3 werden normalerweise durch einen Phrasenton (H- bzw. L-) markiert.
- B4 : Grenze einer Intonationsphrase. Starke tonale Markierung, häufig verbunden mit Dehnung oder Pause. B4 werden normalerweise bitonal durch einen Phrasenton (H-/L-) und einen Grenzton (H%/L%) markiert.

Grenztypen und Grenztöne sind nur teilweise voneinander unabhängig. Dadurch, daß B4 immer bitonal und B3 immer durch einen Phrasenton markiert sind, brauchen sie nicht unbedingt explizit etikettiert werden. Sie lassen sich aus der tonalen Etikettierung rekonstruieren. Es wird dann nur eine explizite Etikettierung irregulärer Grenzen (B2) benötigt. Das heißt jedoch nicht, daß die in GToBI etikettierten prosodischen Grenzen rein auf der intonatorischen Markierung beruhen. Zum Beispiel liegt der Unterschied zwischen H- (B3) und H-L% (B4) hauptsächlich in der unterschiedlich wahrgenommenen „Tiefe“ der Grenze (die auch durch Dehnung bzw. Pause verursacht sein kann).⁴

15.5 Transkriptionsexperiment

Um die Verwendbarkeit des Inventars zu untersuchen, wurde ein Transkriptionsexperiment durchgeführt. Dazu wurde zunächst ein Trainingskorpus festgelegt. Dieses enthielt zu den einzelnen Etiketten passende Sprachbeispiele. Weiterhin wurde eine kurze Beschreibung ausgearbeitet, in der die Anwendung der Etiketten in den Sprachbeispielen erläutert wurde.

Das Testmaterial enthielt zum Teil gelesene Sprache:

1. Nachrichten, gesprochen von einem trainierten Sprecher (DLF),
2. Buchausschnitt, gelesen von einem trainierten Schauspieler,
3. Abschnitte aus einem Reiseführer, gelesen von einem untrainierten Sprecher.

⁴Im VERBMOBIL GToBI werden die Grenztypen etwas anders etikettiert: normale Wortgrenze als B1, intermediäre PG als B2, Intonationsphrasengrenze als B3. Die irreguläre Phrasengrenze wird mit B9 etikettiert. Der Unterschied ist aber rein formal, die Übersetzung nach GToBI ist eindeutig.

Der andere Teil bestand aus Spontandialogen: 1) eine Terminabsprache aus dem VERBMOBIL Korpus, 2) ein Dialog aus dem Map Task Korpus.

Insgesamt bestand das Material aus 35 Äußerungen, 733 Wörtern und 304 Sekunden. Das Material wurde parallel von 13 Transkribenten etikettiert. Zu den Transkribenten gehörten drei der AutorInnen, die restlichen waren Studenten mit nur wenig Transkriptionserfahrung, die sich anhand des Trainingsmaterials eingearbeitet hatten.

Wegen des beschränkten Materialumfangs und um die Studenten nicht zu überfordern, wurde ein eingeschränktes Inventar verwendet: Upstepping wurde nicht etikettiert, ein diakritisches Fragezeichen konnte verwendet werden, um Unsicherheit der Transkription anzuzeigen. Auch wurden für das Sprachmaterial nur norddeutsche Sprecher verwendet.

15.6 Arbeitsumgebung

Die Transkriptionen wurden in Stuttgart und Saarbrücken mit *ESPS xuvaves(tm)* durchgeführt. In Braunschweig wurde das Programm *fish* verwendet, ein freies Softwarepaket zur Etikettierung von Sprachsignalen. Die Transkriptionen wurden akustisch mit visueller Unterstützung (Sprachsignal und Sprachgrundfrequenz) erstellt.

15.7 Resultate

Die resultierenden 13 Transkriptionen wurden miteinander verglichen und die Korrespondenzen zwischen den einzelnen Transkribenten ermittelt. In Tabelle 15.1 sind zunächst für die 13 Transkribenten die Anzahl vergebener Akzente und Grenzen angegeben. Von den 733 Wörtern markierten die Transkribenten im Mittel 301 als akzentuiert und etikettierten 60 B3 und 90 B4. Die bei Tabelle 15.1 angegebenen Standardabweichungen (σ) zeigen, daß die Anzahl vergebener Etiketten unterschiedlich streut. Während Tonakzente und B4 relativ zum Mittelwert wenig streuen, liegt die Standardabweichung bei B3 mit 25 bei im Mittel 60 etikettierten B3 relativ hoch. Dies zeigt, daß die Etikettierung der intermediären Phrasengrenze weniger verlässlich ist als die der anderen Kategorien.

Die prozentualen Korrespondenzen zwischen den Transkribenten wurden analog dem in Silverman et al. (1992a) für das englische ToBI-System beschriebenen Verfahren ermittelt. Bei diesem Verfahren werden für jedes der 733 Wörter alle Transkribentenpaare (78 bei 13 Transkribenten) miteinander verglichen. Insgesamt wurden also 57174 Transkriptionspaare ausgewertet; die Korrespondenz ergibt sich dabei aus der Anzahl gleich transkribierter Wörter bezogen auf die Gesamtheit. Die Korrespondenzen wurden für die unterschiedlichen Etiketten einzeln ermittelt. Sie sind in Tabelle 15.2 angegeben. Dabei wird zunächst verglichen, ob jeweils zwei Transkribenten überhaupt einen Akzent auf einem Wort etikettiert hatten.

Tabelle 15.1: Anzahl vergebener Etiketten für die 13 Transkribenten. Der Mittelwert für die Anzahl der Akzente liegt bei 301 ($\sigma = 17$), für B3 liegt der Mittelwert bei 60 ($\sigma = 25$) und für B4 bei 90 ($\sigma = 18$).

	tr1	tr2	tr3	tr4	tr5	tr6	tr7
Anz. Akzente	295	304	295	301	319	312	288
Anz. B3	74	50	31	62	38	73	34
Anz. B4	104	88	56	92	100	110	87
	tr8	tr9	tr10	tr11	tr12	tr13	
Anz. Akzente	279	290	328	271	319	312	
Anz. B3	50	65	70	131	53	54	
Anz. B4	104	93	91	62	113	63	

Tabelle 15.2: Korrespondenzen insgesamt und getrennt für Experten und Studenten. Die Spalten geben an: (1) Kategorie, (2) Anzahl der unterschiedlichen Etiketten (jeweils +1 für keine Markierung), (3) Gesamtkorrespondenz, (4) Expertengruppe, (5) Studentengruppe

	Auswahl	Korr. ges.	Exp.	Stud.
akz./unakz.	2	87%	91%	86%
Tonakzente	7	74%	81%	73%
Tonakz. + Downst.	10	71%	78%	70%
tonale Grenzmarkierung	7	86%	90%	85%
Grenztypen	4	87%	89%	86%

Der zweite Vergleich bezog sich auf die Art des Tonakzentes (dabei wurden alle Wörter, also auch die unakzentuierten mit einbezogen). Der dritte Vergleich berücksichtigte zusätzlich, ob beide übereinstimmend Downstepping etikettierten. Für die Grenzmarkierungen wurde einmal die tonale Markierung verglichen und dann der Grenztyp, also die „Tiefe“ der Grenze. Zunächst wurde dieser Vergleich für alle 13 Transkribenten gemeinsam durchgeführt, dann wurden zwei Gruppen gebildet, einmal die 3 Experten (Autoren) und die 10 Studenten.

Mit Ausnahme der Tonakzente liegen die Korrespondenzen durchgehend deutlich über 80%. Die Werte sind denen in Pitrelli et al. (1994) angegebenen für englisches ToBI vergleichbar, ein detaillierter Vergleich beider Experimente wird von Grice et al. (1996) durchgeführt. Allein die Tonakzente liegen nur bei etwa 70%. Dies liegt zum Teil an dem reichhaltigen Inventar, zum Teil sind die Kategoriengrenzen zwischen den Tonakzenten aber auch unscharf, sie haben Randbereiche, in denen die Unterscheidung schwerfällt. Zum Beispiel verfügen H* und L+H* beide über einen hohen Zielpunkt in der akzentuierten Silbe und sind, gerade am Beginn einer Intonationsphrase, beide tendenziell steigende Akzente. Der durch den tiefen

Tabelle 15.3: Korrespondenzen für die einzelnen Tonakzente, Häufigkeit der Tonakzente (Gesamtsumme ist 733 Wörter x 13 Transkriptionen = 9529).

unakz.	ohne Downstepping 90%	Anzahl 5757	mit Downstepping 90%	Anzahl 5757
H*	62%	2272	56%	1914
!H*			28%	358
L+H*	38%	610	38%	576
(L+!H*)			6%	34
L*+H	35%	408	34%	382
L*+!H			10%	26
H+!H*	14%	171	14%	171
H+L*	14%	72	14%	72
L*	18%	239	18%	239

Führungston hervorgerufene stärkere Anstieg bei L+H* ist nicht immer klar auszumachen.

Es zeigt sich weiterhin, daß es zwar, wie zu erwarten war, einen Unterschied zwischen der Studenten- und der Expertengruppe gibt, daß dieser Unterschied aber nicht sehr groß ist. Die Studenten scheinen trotz des kurzen Trainings keine prinzipiellen Probleme mit der Anwendung des GToBI-Etikettensystems zu haben. Durch ein erweitertes und verbessertes Training könnte der Unterschied zur Expertengruppe noch verringert werden. Eine Korrespondenz von 90% wäre hier ein erreichbares Ziel (außer wohl für die Etikettierung der Tonakzente).

In Tabelle 15.3 wurden die Korrespondenzen für die Tonakzente einzeln ermittelt. Für die Berechnung konnte das oben verwendete Verfahren nicht benutzt werden. Es wurde daher ein Verfahren angewandt, das an die Erkennungsrate bei der automatischen Spracherkennung angelehnt ist. Die Berechnung erfolgt gemäß Gleichung 15.1.⁵

$$(15.1) \quad corr_{1,2,label} = \frac{n_{corr(1,2),label}}{(n_{1,label} + n_{2,label})/2}$$

Werden die Übereinstimmungen für die Tonakzente wie in Tabelle 15.3 einzeln durchgeführt, so zeigen sich ganz erhebliche Unterschiede zur Gesamtkorrespondenz. In letztere gehen die einzelnen Akzente mit ihrer Häufigkeit gewichtet ein, so daß selten vorkommende Akzente wie z.B. H+L* trotz sehr geringer Übereinstimmung den Wert kaum verschlechtern. Gerade die selten vorkommenden Akzente werden aber offensichtlich leicht verwechselt.

⁵ Das Prinzip ist folgendermaßen: zunächst wird die Etikettierung des einen Transkribenten als Referenz angenommen und die „Erkennungsrate“ des anderen ermittelt. Danach wird getauscht und die „Erkennungsrate“ des ersten berechnet. Zuletzt werden beide Werte gemittelt.

Interessant ist auch die Frage, ob sich Unterschiede zwischen den einzelnen Transkribenten zeigen. Frühere Untersuchungen (Reyelt 1993) zeigten, daß die Streuung sehr uneinheitlich verteilt sein kann. Dazu wurde jeweils die Korrespondenz eines Transkribenten zu allen anderen berechnet und gemittelt. Die Ergebnisse zeigt Tabelle 15.4.

Tabelle 15.4: Korrespondenz jedes der 13 Transkribenten zu den restlichen. Alle Werte in Prozent. Die Experten sind mit „*“ markiert.

tr1	tr2	tr3	tr4	tr5	tr6	tr7	tr8	tr9	tr10	tr11	tr12	tr13
86	87*	82	87	84	86*	86	87	85	86*	84	82	83

Offensichtlich gibt es unter den Transkribenten keine „Ausreißer“, die deutlich anders etikettieren als der Rest. Auch die Experten zeigen weder bessere noch schlechtere Werte als die Studenten, was ein Indiz dafür wäre, daß beide Gruppen unterschiedlich etikettieren. Das Training scheint also zumindest so beschaffen zu sein, daß die Transkribenten die Etiketten zwar etwas unsicherer, aber nicht falsch anwandten.

15.8 Zusammenfassung

Mit GToBI ist ein System zur prosodischen Transkription des Deutschen entwickelt worden, das bereits an mehreren Instituten in Deutschland angewandt wurde. Das hier beschriebene Experiment zeigte, daß Transkribenten schon nach einer kurzen Einarbeitungszeit mit dem GToBI-System Etikettierungen von guter Konsistenz durchführen können. Damit ist eine Vorbedingung für die Erstellung umfangreicher prosodisch etikettierter Datenbasen erfüllt. Um die Konsistenz weiter zu verbessern, muß an Hand der Ergebnisse des hier beschriebenen Transkriptionsexperiments das Trainingsprogramm überarbeitet und ergänzt werden. Das System soll auch noch in einigen weiteren Aspekten erweitert werden (Markierung von Kern- vs. Randtypen sowie der Sicherheit der Entscheidung).

Literaturverzeichnis

- A. Anderson, M. Bader, E. Bard, E. Boyle, G. Doherty, S. Garrod, S. Isard, J. Kowtko, J. McAllister, J. Miller, C. Sotillo, H. Thompson und R. Weinert (1991). The HCRC map task corpus. *Language and Speech* 34(4): 351–366.
- A. Batliner, A. Feldhaus, S. Geißler, T. Kiss, R. Kompe und E. Nöth (1996). Prosody, empty categories and parsing – A success story. In: *Proc. ICSLP*, Philadelphia. In Druck.
- A. Batliner, A. Kießling, R. Kompe, E. Nöth und B. Raithel (1992). Wann geht der Sonderzug nach Pankow? (Uhrzeitangaben und ihre prosodische Markierung in der Mensch–Mensch– und in der Mensch–Maschine–Kommunikation). In: *Fortschritte der Akustik – DAGA92*, S. 541–544, Bad Honnef. DPG–GmbH.

- C. Féry (1993). *German Intonational Patterns*. Linguistische Arbeiten 285. Niemeier, Tübingen.
- D. Gibbon (1976). *Perspectives of Intonation Analysis*. Lang, Hamburg.
- M. Grice und R. Benzmüller (1995). Transcription of German using ToBI-tones: The Saarbrücken System. *PHONUS* Institut für Phonetik, Universität Saarbrücken.
- M. Grice, M. Reyelt, R. Benzmüller, J. Mayer und A. Batliner (1996). Consistency in transcription and labelling of German intonation with GToBI. In: *Proc. ICSLP*, Philadelphia. In Druck.
- W. Hess, A. Batliner, A. Kießling, R. Kompe, E. Nöth, A. Petzold, M. Reyelt und V. Strom (1996). Prosodic modules for speech recognition and understanding in Verbmobil. In: Y. Sagisaka, N. Campbell und N. Higuchi, Hg., *Computing Prosody*, S. 363–383. Springer, New York.
- K. Kohler (1995). PROLAB – The Kiel system of prosodic labelling. In: *Proc. International Congress of Phonetic Sciences*, Band 3, S. 162–165, Stockholm.
- R. Kompe, A. Kießling, H. Niemann, E. Nöth, E. Schukat-Talamazzini, A. Zottmann und A. Batliner (1995). Prosodic scoring of word hypotheses graphs. In: *Proc. EUROSPEECH*, Band 2, S. 1333–1336.
- D. R. Ladd (1995). “linear” and “overlay” descriptions: An autosegmental-metrical middle way. In: *Proc. ICPHS*, Band 2, S. 116–123, Stockholm.
- D. R. Ladd, K. E. Silverman und K. R. Scherer (1983). Parametrische und kategoriale Ansätze bei der Erforschung intonatorischer Funktion. *Zeitschrift für Literaturwissenschaft und Linguistik* 49: 124–133.
- J. Mayer (1995a). Towards the workshop on prosodic labelling. Results of a questionnaire concerning current works on the prosodic labelling of German speech data bases. Technischer Bericht, Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart, Stuttgart.
- J. Mayer (1995b). Transcription of German intonation: The Stuttgart System. Technischer Bericht, Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart, Stuttgart.
- J. Pierrehumbert (1980). *The Phonology and Phonetics of English Intonation*. Dissertation, M.I.T.
- J. Pitrelli, M. Beckman und J. Hirschberg (1994). Evaluation of prosodic transcription labeling. In: *Proc. ICSLP*, S. 123–126, Yokohama.
- M. Reyelt (1993). Experimental investigation on the perceptual consistency and the automatic recognition of prosodic units in spoken German. Working papers 41, Lund University, Dept. of Linguistics.
- M. Reyelt (1994). Ein flexibles Programm Paket zur Visualisierung von Sprachdaten. In: K. Fellbaum, Hg., *Tagungsband Elektronische Sprachsignalverarbeitung*. S. 358–365, Berlin.
- M. Reyelt und A. Batliner (1994). Ein Inventar prosodischer Etiketten für Verbmobil. Verbmobil-Memo 34, Techn. Univ. Braunschweig, LM-Universität München.
- K. Silverman, M. Beckman, J. Pitrelli, M. Ostendorf, C. Wightman, P. Price, J. Pierrehumbert und J. Hirschberg (1992a). TOBI: A standard for labeling

- English prosody. In: *Proceedings of the 1992 International Conference on Spoken Language Processing*, S. 867–870.
- K. Silverman, E. Blaauw, J. Spitz und J. Pitrelli (1992b). Prosodic comparison of spontaneous speech and read speech. In: *Proceedings of the 1992 International Conference on Spoken Language Processing*, S. 1299–1302.
- H. Tillmann, B. Eisen und C. Draxler (1992). Consistency of judgements in manual labelling of phonetic segments: The distinction between clear and unclear cases. In: *Proc. ICSLP 92*, S. 871–874.
- W. Wahlster (1993). Verbmobil: Translation of face-to-face dialogs. In: *Proc. EUROSPEECH*.