

## On the use of prosody in automatic dialogue understanding

**Elmar Nöth, Anton Batliner, Volker Warnke, Jürgen Haas, Manuela Boros, Jan Buckow, Richard Huber, Florian Gallwitz, Matthias Nutt, Heinrich Niemann**

### Angaben zur Veröffentlichung / Publication details:

Nöth, Elmar, Anton Batliner, Volker Warnke, Jürgen Haas, Manuela Boros, Jan Buckow, Richard Huber, Florian Gallwitz, Matthias Nutt, and Heinrich Niemann. 2002. "On the use of prosody in automatic dialogue understanding." *Speech Communication* 36 (1-2): 45–62.  
[https://doi.org/10.1016/S0167-6393\(01\)00025-5](https://doi.org/10.1016/S0167-6393(01)00025-5).

# ON THE USE OF PROSODY IN AUTOMATIC DIALOGUE UNDERSTANDING

E. Nöth<sup>1</sup>    A. Batliner<sup>1</sup>    V. Warnke<sup>1</sup>    J. Haas<sup>1</sup>    M. Boros<sup>2</sup>    J. Buckow<sup>1</sup>    R. Huber<sup>1</sup>  
                     F. Gallwitz<sup>1</sup>    M. Nutt<sup>2</sup>    H. Niemann<sup>1</sup>

<sup>1</sup>University of Erlangen–Nuremberg, Chair for Pattern Recognition (Inf. 5), D-91058 Erlangen, Germany

<sup>2</sup>Bavarian Research Center for Knowledge Based Systems (FORWISS), D-91058 Erlangen, Germany

## ABSTRACT

In this paper, we show how prosodic information can be used in automatic dialogue systems and give some examples of promising new approaches. Most of these examples are taken from our own work in the VERBMÖBIL speech-to-speech translation system and the EVAR train timetable dialogue system. In a ‘prosodic orbit’, we first present units, phenomena, annotations and statistical methods from the signal (acoustics) to the dialogue understanding phase. We show then, how prosody can be used together with other knowledge sources for the task of resegmentation and how an integrated approach leads to better results than a sequential use of the different knowledge sources; then we present a hybrid approach which is used to perform a shallow parsing and which uses prosody to guide the parsing; finally, we show how a critical system evaluation can help to improve the overall performance of automatic dialogue systems.

## 1. INTRODUCTION

We describe the present state of the art of using prosody in automatic dialogue systems. By that, we give a rather personal view, exemplified with our own work in the VERBMÖBIL domain [12] and in a train timetable information system [15]. Older, well-known surveys on the use of prosody in automatic speech processing are [27, 39]; cf. [30] as well. Work on the use of prosody in automatic speech processing in general and in automatic dialogue understanding in particular has been, and is still quite often, ‘off-line’; this means that it cannot be used directly in fully automatic systems, because, e.g., manually corrected features are used, because it is based on the spoken word chain, because the correct segmentation is assumed, etc. On the other hand, there is an urgent need for ‘real life’ approaches that could be used in systems which really work and can be applied commercially. This means, in turn, that such ‘real life’ approaches have to be fully automatic and, e.g., have to work with word hypotheses graphs (WHG) which are the usual output of word recognition. Manual processing is only ‘allowed’ while testing the algorithms. In order to meet these requirements, all available knowledge should be used. In our presentation, we sketch those components that are necessary for such a use; this is done in Section 2. In Section 3, we focus on some promising trends.

\*This work was funded by the German Federal Ministry of Education, Science, Research and Technology (BMBF) in the framework of the VERBMÖBIL Project under Grant 01 IV 102 H/0 and by the DFG (German Research Foundation) under contract number 810 939-9. The responsibility for the contents lies with the authors.

## 2. THE PROSODIC ORBIT: FROM SIGNAL TO DIALOGUE

In Table 1, we try to sketch those units, phenomena, annotations and statistical modelling methods one normally has to deal with if one tries to use prosody in automatic dialogue systems. By that, we only want to *illustrate* different and possibly alternative procedures; we do not want to present an *exhaustive* overview; of course, a different terminology could be used. Some of the descriptive terms that are used here are intuitively clear, even if a precise description is practically impossible (what is a ‘word’?); some of them are rather vague and unclear (what precisely does ‘focus’ mean?). Still, we believe that all of these terms are well-known so that the reader can follow our argumentation. Some interesting topics where prosody can provide valuable information are not mentioned in Table 1, e.g., emotional state of the speaker or speaker identification/recognition. Such topics will be relevant for automatic dialogue systems in the near future.

We do not give every suitable **level of analysis** in Table 1, only those two which are the main topics of this workshop, i.e., prosody and dialogue, and one rather complex level in between, namely syntax/semantics which is traditionally — and in fact — the mediator between these two levels. We do believe, however, that these levels represent the core of most of the work that has been done in this area.

We usually presuppose that somehow the result of a **word recognition** is available. We can use the spoken word chain and by that assume one hundred percent correct word recognition (‘cheating’) if we want to concentrate on the other phenomena or if we want to determine an upper bound. For a ‘real life’ task, however, we have to deal with the output of a word recognizer, i.e., with a WHG with several alternative word chains. Sometimes, the WHG does not even contain the spoken word chain. Note that for prosodic processing, a representation of the spoken words is actually not necessary: ‘pure’ prosody can be used to recognize accentuation or prosodic boundaries, cf. [37]. Afterwards, however, this pure prosody approach has to be combined with word information.

Pitch, loudness, etc. are perceived prosodic properties. Actually, they are given in Table 1 only for ‘completeness’ because the methods used in automatic speech processing do, of course, not perceive; rather they measure the **acoustic** correlates of perception, i.e., F0, energy, duration, etc. These acoustic correlates have to be computed for a certain time dimension: either a fixed one, if they are measured in fixed time windows, or a flex-

levels of analysis			
	prosody	syntax/semantics	dialogue
acoustics	(segmental) units		
F0, energy, duration, .....	phones/phonemes syllables words phrases/sentences turns/utterances	morphemes words phrases/sentences	phrases/sentences turns/utterances
perception	phenomena		
pitch, loudness, duration, speaking rate, .....	boundaries/phrasing	constituents/phrases clauses/sentences	dialogue act boundaries
	accentuation	focus	saliency
	sentence mood	sentence mood	dialogue acts ( $\approx$ illocution)
extraction	annotations (exemplified with our own approach)		
automatically extracted/ manually corrected	boundaries: B3, B2, B0, B9	synt.–pros. M labels (M3, M0) $\rightarrow$ S labels	D3, D0
	accents: EC, PA, SA, NA	A3, A2, A0	—
	questions PQ	questions SQ	dialogue acts DA
	statistical modelling methods		
	NN, DT, LDA, HMM, ...	LM, DT, ...	LM, DT, ....

**Table 1:** units, phenomena, annotations and methods

ible one, if they are confined to certain segmental units, such as phones/phonemes, syllables, words, etc. A ‘pure’ prosody approach has to work with fixed time windows or, e.g., with independently extracted syllable boundaries. The **extraction** of prosodic features in automatic systems is — no wonder — automatic. For a training sample or a test sample that is used as reference, the extraction can be manual as well, or an automatic extraction can be corrected manually afterwards. (This does not happen too often because of the effort needed.)

Note that from an application point of view (i.e. for an automatic system), **perception** units are not ‘necessary’: if there is a mapping from acoustics onto perception, and again, a mapping from perception onto function, then statistical modelling should be able to directly map acoustics onto function. Of course, knowledge on perception can guide feature selection and feature transformation/normalization. It is, however, our experience that very often, raw feature values rather than transformed or combined feature values should be taken if the database is sufficiently large for the training of the statistical classifier; i.e., we leave it up to the classifier to learn the most appropriate transformation.

The same holds for the **phonological** level: to put it bluntly, phonological systems like the well-known ToBI-approach only introduce a ‘quantisation error’: the whole variety of F0 levels available in acoustics is reduced to a mere binary opposition, Low vs. High, and to some few additional, diacritic distinctions. In our opinion, this fact alone prevents tone levels (or any other ‘prosodic phonological’ concepts as, e.g., the one developed within the IPO-approach) from being a meaningful step that automatic processing should be based on; it seems better to leave it up to a large feature vector and to statistical classifiers to find the form to the function. Actually, to our knowledge, there is no existing approach which really uses such phonological units for the recognition of prosodic events. To prevent misunderstandings we want to stress that this caveat does not hold for *phonological knowledge*, which can be a valuable source, but only for

the direct use of *phonological theoretical concepts* in automatic speech recognition.

The **segmental units** in prosody can be very short — either a time window or a phone/phoneme — or they can constitute a whole turn/utterance. Larger units are normally only used for comparison/normalization. Dialogue units are higher level units and thus usually longer than those of syntax/semantics.

The **phenomena** we want to deal with are first **phrasing**, i.e., prosodic boundaries that mirror syntactic boundaries which, in turn, mirror dialogue act boundaries. ‘Mirror’ means here, that a rather high, albeit not perfect correlation is assumed — otherwise, the use of prosodic information in syntax and/or dialogue would not make much sense. Second comes **accentuation** and, by that, the most important information in a unit, e.g., in a sentence (focus) or in a dialogue act (saliency). Third, prosody can, for certain constellations, disambiguate between different **sentence moods/modalities** and, by that, different illocutionary/dialogue acts. For example, prosody can be used to decide whether an elliptic sentence (free phrase) is a statement or a question [5, 10].

In order to know what we are talking about, we have to have **labels** for our phenomena, and in order to know, whether we are on the right track or not, we have to **annotate** corpora with these labels which we then can use as training and test data. In Table 1, we give examples of our own work within the VERBMOBIL project which started in 1994 and will end in September 2000. The VERBMOBIL database contains spontaneous speech dialogues of German, English, and Japanese speakers. For each utterance, a basic transliteration is given containing the spoken words, the lexically correct word form, and several labels for (filled) pauses and non-verbal sounds. In addition to this basic transliteration, large parts of the corpus are further annotated with prosodic, syntactic, and dialogue act labels. All labels are word-based and normally introduced into the spoken word chain to the

right of the word they belong to, cf. Table 2. We started with a ToBI-like annotation scheme, cf. [34, 18]. Because of the caveats mentioned above, we only use the functional boundary tier comparable to the break index tier in ToBI, and the functional accent tier, comparable to the ‘starred’ tones in ToBI: strong boundary B3, medium boundary B2, no boundary B0, and irregular boundary B9, and primary (phrase) accent PA, emphatic/contrastive accent EC, secondary accent SA, and unaccentuated UA; as for details, cf. [8, 24, 26]. The boundary labels were used within the syntax modules of VERBMOBIL. Because prosodic boundaries do not always denote syntactic boundaries, we introduced another type of boundaries, the syntactic-prosodic, so-called M boundaries (‘M’ for language ‘M’odel). A total of 25 different subclasses were mapped onto three main classes: a main boundary class M3 (between clauses, free phrases, etc.), M0 (no boundary), and MU (ambiguous boundary). A detailed description of these M labels, including correlations with other label types and classification results, can be found in [8]. Alternatively, the M subclasses were mapped onto five syntactic ‘S’ boundary classes which can be described in an informal manner as follows: S0: no boundary, S1: at particles, S2: at phrases, S3: at clauses, S4: at main clauses and at free phrases. These S boundaries meet the special needs of some higher linguistic modules in the VERBMOBIL system. Based on the M boundaries and the prosodic-perceptual accent labels as a reference, we developed a rule-based system of accents with primary accent A3, secondary accent A2, and no accent A0 [9]. In addition, syntactic questions SQ are annotated in the basic transliteration. Sentence boundaries annotated with SQ and ending in a high boundary tone H% can be labelled as prosodic questions PQ. We thus have a complete set of boundary, accent, and question labels that is based on the **prosodic form** and an analogous set of labels that is based on **syntactic structure**, i.e. on the surface, on word ordering. Dialogue act (DA) classes were annotated independently; in this paper, we use the same 18 DA classes as in [22]; they are defined by their illocutionary force, such as “GREET, INIT, BYE, SUGGEST, REQUEST, ACCEPT, ...”. The criteria for the segmentation of turns into DAs are partly syntactic: for example, all ‘material’ that belongs to the verb frame of a finite verb belongs to the same DA. By that, we avoided to listen to the turns and could thus reduce the labelling effort. In [13], it is reported that DA segmentation changes only slightly when the annotators can listen to the speech data, but cf. [35]. DA boundaries D3 are, so to speak, a by-product of the DA annotation, as well as their complement, D0 (no DA boundary).

Of course, it is always desirable to have large-scaled annotations of exactly those units one has to deal with; this is not always realistic, however. We thus tried to aim at an **integrated** labelling approach: for example, prosodic, syntactic, and DA boundaries are highly correlated with each other; exact figures can be found in [8]. If enough material is available, we can use exactly those labels that model the units we are interested in; if not, we can use highly correlated labels. Generally, we try to use **overspecified** labels that are normally not classified as such but are mapped onto some few main classes. For example, we currently do not use D3 labels for the segmentation of DA units in the ‘official’ VERBMOBIL system, but S4 labels, which in more than 90% correspond to D3 labels. It is, however, no problem to use D3 labels directly in a later stage, if necessary. In analogy, we do not have to annotate saliency in DAs at all, because we can use our prosodic and/or rule-based accent labels instead. Table 2 shows a slightly sim-

plified example from the English VERBMOBIL database with all label types introduced above (The default classes B0, A0, etc. are not shown).

There is, of course, a wide variety of feature extraction algorithms which we do not want to deal with in this paper. Also, there is a wide variety of **statistical modelling methods** for (more or less unsupervised) clustering and subsequent classification of the phenomena. In Table 1, we only mention: Neural Networks (NN) — Multi-Layer-Perceptrons (MLP), a special kind of NN, are used by us to classify prosodic labels; Decision Trees (DT), Linear Discriminant Analysis (LDA), Hidden Markov Models (HMM), and Language Models (LM). Each of these general methods has a variety of sub-methods. Normally, NNs, LDA, and HMMs are used for acoustic data [17], although categorical labels can be incorporated as well. LMs are used for words (unigrams) and word sequences (bi-, trigrams etc.), and DTs are used for both.

Practically all studies on the use of prosody in speech processing, in general, and in automatic dialogue understanding, in particular, use one or more acoustic prosodic features F0, energy, duration, etc. (top left corner of Table 1) and try to recognize the kind of labels given under the heading ‘annotations’ that represent those events given under the heading ‘phenomena’ in Table 1. This is the common core, everything else differs: number and manner of features extracted, units, phenomena, and statistical methods. (Note that this fact makes it virtually impossible to compare classification results across studies in a strict sense!) Classification can be separated and sequential, e.g., first prosodic boundaries, then syntactic boundaries, then dialogue act boundaries, then dialogue acts, and independent from that, accent classification etc. Classification can be combined and integrated, e.g., one can combine boundary and accent classification, cf. [25], one can integrate DA boundary and DA classification etc., cf. below and [41], and one can even combine word recognition and boundary classification [16, 17].

Thus, out of each column and row in Table 1, we can choose one, more, or all items we want to use and/or recognize, and this can be done separately, or combined, or integrated. In [35, p. 446], e.g., it is reported that overall duration is the most important prosodic feature for the classification of DAs: “This is not surprising, as the task involves a seven-way classification including longer utterances (such as statements) and very brief ones (such as backchannels like “uh-huh”).” In [6], it is reported that three durational features alone (word based, syllable based and pause duration) yield an overall recognition rate of 86% for prosodic boundaries. So we *could* use only such duration features but, of course, the more (relevant) prosodic features we use, the better is the classification [6]. In our opinion, this result can reasonably be generalized to all other knowledge sources: the more knowledge sources we employ — and the better they are tuned to each other, the better the classification will be. This, of course, holds only if these knowledge sources are modelled adequately. This means, at least, that enough reliable training data are available, and that the statistical modelling is adequate as well.

### 3. SOME PRESENT AND FUTURE TRENDS

Based on the prosodic orbit put forth in the previous section, we now want to describe some promising trends exemplified with

turn with types of labels given in Table 1	Dialogue Acts
<i>two o'clock in the afternoon sounds fine</i> PA A3 B3 M3/S4 D3	ACCEPT
<i>where would you like</i> SA A2 M3/S3 <i>to meet</i> PA A3 B3 M3/S4 QBT D3	REQUEST

**Table 2:** Example turn with annotations

our own material and work. Again, we cannot give a complete overview; for that, we refer to the other ‘tutorials’ given by Julia Hirschberg, Herb Clark, and Stephen Pulman, and to the other papers at this workshop. We will concentrate on a **shallow** analysis; as for a **deep** (syntactic) analysis, we refer to [26].

From a phylogenetic as well as from an ontogenetic point of view, a dialogue with the parents or the peer group is the earliest and most natural way of communication for human beings. If we thus compare automatic dialogue systems with other automatic speech processing applications, we can say that they are ‘most natural’, i.e., rather close to the original function of natural, spontaneous language/speech, in contrast to other applications, e.g., automatic speaker or language identification. This means that we can find parallels between the behavior of humans in natural dialogues and features that should be incorporated in sophisticated automatic dialogue systems. In a human–human dialogue, the speakers

1. reanalyze, if they went the wrong track and notice that their analysis will not work
2. integrate different knowledge sources, and do most certainly not proceed in a strictly sequential manner
3. pay attention to salient parts of utterances and disregard non-salient ones
4. are content if and only if they ‘get what they wanted’, and tolerate non-fatal misunderstandings

User utterances, e.g. in VERBMOBIL, can be very long. Such utterances are always processed **incrementally** by the listener; that means that the hearer forgets more or less the exact wording of a sentence rather soon, and only stores its meaning and, if necessary, its illocution [19, p. 460ff]. In analogy, automatic systems should be able to process longer utterances in an incremental way as well. Otherwise, system responses would be delayed unduly and, by that, user acceptance would be rather low.

We will address all four topics and concentrate on the second and the third topic, where we present methodologies and experimental results for two domains. More details can be found in [41, 32]. The first and the last topic will be dealt with rather sketchy, since we can not yet provide substantial results.

### 3.1. Reanalysis within a Sequential Approach: the VERBMOBIL System

State of the art speech understanding systems use different knowledge sources to interpret a spoken utterance. In the field of human–human or human–machine dialogue processing, the most important tasks are the segmentation, classification and interpretation of automatically recognized user utterances using several different knowledge sources [35, 38, 40, 11]. Commonly, these different knowledge sources are applied sequentially. For example, in the VERBMOBIL speech-to-speech translation system

[40, 11], first a word recognizer generates a WHG using only acoustic and LM information. The word sequences are then segmented into syntactic–prosodic phrases using prosodic and LM information. Finally, these already segmented phrases are interpreted by a parser or a stochastic process with the use of several knowledge sources. Thus, it is impossible to incorporate the knowledge of the syntactic–prosodic process, the parser or any other later process to find the best word chain within the word recognition task. In a system like VERBMOBIL, which proceeds in a sequential manner with no back-tracking mechanism, the higher linguistic modules are thus sometimes faced with wrong segmentation. Here, we want to give examples for two different factors that can be responsible for a wrong segmentation.

Consider the following turn containing a repair; note that repairs are most of the time not marked by an edit term in the VERBMOBIL database [7].

*Treffen wir uns am Montag* B3/S4 — *am Dienstag*  
*Let's meet on Monday* B3/S4 — *on Tuesday*

Let us assume that both NN and LM classify this boundary as B3 and S2 or S4, resp. If the syntax module accepts this analysis, the phrase *on Tuesday* has to be interpreted as a free phrase or as a right dislocation [8], and by that, as a sort of — contradictory — addendum or specification of *on Monday*. In VERBMOBIL, a repair module is located between prosody and syntax [36]. It uses prosodic information, i.e. looks at boundary locations, to see if a repair occurs in their surrounding. If the repair module can mark the word boundary in question as a possible interruption point, it compares the part-of-speech labels of the constituent to the left and to the right of the end of the reparandum. If it can find the reparandum, these words can be cut out and the correct translation can be generated:

*Let's meet on Tuesday*

As a second example, consider the following turn with a word-by-word translation into English:

*Ja, ich habe am Montag* B3/S4 *oder am Donnerstag* Zeit  
*Well, I have on Monday* B3/S4 *or on Thursday* time

Prosody alone might not help because there is a pronounced pause after *Montag*. Here, the analysis window of the LM can be too small, and thus, a wrong segmentation within the verbal bracket can be generated; note that the verbal bracket (‘Verbracket’, i.e., a special bracketing for linguistic groupings) is a syntactic phenomenon that does not exist in English. In such cases, the syntax module will not simply rely on the output of the NN/LM but detect, that the right end of the verbal bracket has not been reached yet, and that a correct analysis can only be generated if this wrong segmentation is discarded [23].

### 3.2. An Integrated Approach: the $A^*$ Search

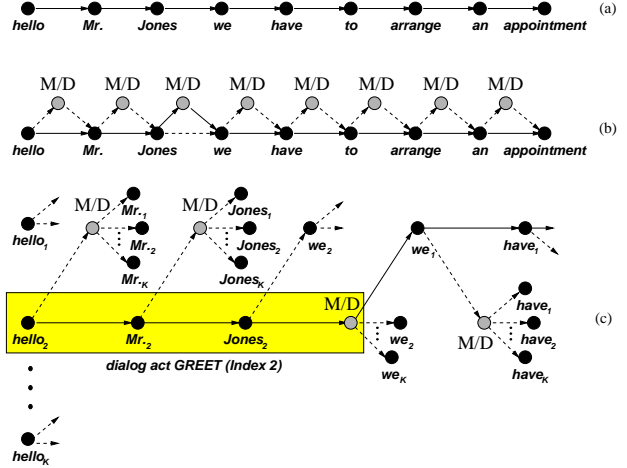
We have seen that in a sequential approach, we sometimes have to repair wrong analyses, e.g., a wrong segmentation, in a subsequent pass within the higher linguistic modules. Another way of combining higher with lower linguistic knowledge is an integrated approach. In such an approach, we integrate multiple knowledge sources into one  $A^*$  search to find, for example, the best word chain, the best syntactic-prosodic phrase or DA boundaries, and the best DA interpretation. The procedure is suitable for any type of WHG, e.g. a complex graph with a high number of word hypotheses, a flat graph containing only the best recognized word chain, or a manually transliterated spoken word chain. The phrase boundaries can be determined using a MLP with prosodic features and/or a LM using textual information. During the search, the possibility of a DA switch is taken into account at each hypothesized phrase boundary. For example, the LM score of the optimal path for the utterance “Good morning, my name is Jones” is determined using the DA specific LMs for GREETING and INTRODUCTION. This score is combined with the score of the DA transition from GREETING to INTRODUCTION, which is calculated using a DA sequence LM. During search, the individual cost functions are combined as a weighted sum. Thus, the search procedure implicitly determines not only the best word sequence, but also phrase boundaries and a rough semantic interpretation of the utterance, using all available knowledge sources.

A high correlation between different types of boundary labels can be found not only in the example given in Table 2, but also in the rest of the corpus (cf. [8] for a detailed analysis). On average, one of two M3 boundaries is also a D3 boundary, and practically all D3 boundaries are also M3 boundaries. This is the main reason why we started to combine the MLP of our prosodic classifier with a text-based LM classifier in previous work [28, 42]. For our experiments, we use the data from the German part of the VERBMobil database annotated in the manner described above. Because of different amounts of training data available for the different knowledge sources (790 turns for prosodic accents and boundaries, 12970 turns for M3, 5980 turns for D3) we have different training and validation sets for each classifier. Our experimental results, however, were always achieved on the same disjunctive test set with 1683 turns. In [41], the modelling of word and DA sequences and their boundaries is described in more detail. Here, we will concentrate on the  $A^*$  search procedure [31] and introduce it in an informal manner. The search proceeds left-to-right through a word graph.

#### The Expansion Procedure

The main difficulty with integrating several knowledge sources into one  $A^*$  search lies in the expansion procedure. In [42], the DA boundaries were modelled implicitly within the word nodes. In our new expansion procedure each phrase boundary is explicitly modelled as a node of its own. Thus, the costs of inserting a boundary can be computed directly, and a boundary node is now required at the end of each DA.

An example for the new expansion procedure is given in Figure 1. The best path is indicated with solid lines, dashed lines indicate alternative expansion rules. Figure 1 (a) shows an example utterance produced by a word recognizer (or the manually transliterated word chain) used as input to the search procedure.



**Figure 1:** (a) A flat word graph with the spoken or recognized word chain. (b) The expansion procedure for integrated boundary classification. (c) The expansion procedure for integrated boundary and DA classification.

In Figure 1 (b), the expansion step for the case of integrated word and boundary classification is depicted. After each word, a possible phrase boundary has to be modelled. If the boundary node has a better score than the following word node, the boundary is inserted into the graph, and the word node is expanded after the boundary node.

The complex expansion procedure for integrated boundary and DA classification is shown in Figure 1 (c). At the beginning of a turn, each DA is possible. Thus, we have to start the expansion with  $K$  alternative nodes (one for each DA). Now the costs for the different alternatives are computed, and the best scored node is expanded next. In our example, the node *hello*<sub>2</sub> (2 is the index for the DA GREET) achieves the best score. Because the current node is no boundary, there are only two alternatives to continue the search. Either there is a phrase boundary after *hello*, or the phrase continues with the word *Mr.*<sub>2</sub>. In this case, a change to another DA is not possible, because new DAs can only be started if a boundary node is expanded. In our example, this happens at the end of the first DA (GREET) at the boundary after the word *Jones*. Now, all  $K$  alternatives for the word *we* have to be generated, and the search again continues with the best scored node. The search is stopped as soon as an explicit goal node is scored best. As for the computation of the costs and the estimation of the remaining costs, we refer to [41].

#### Experiments and Results

All experiments were performed using the manually transliterated word chains as input. The aim of the experiments was to examine if the recognition rates for boundaries and DAs can be improved by adding further knowledge sources to the classification procedure. Analogously to ‘word accuracy’ and ‘word correct’ we evaluate the DA classification with ‘DA accuracy’ (DAA) and ‘DA correct’ (DAC); DAA takes insertions, deletion and substitutions into account while DAC gives the relative amount of correctly classified DAs. For the boundary (M3, D3) classification results we give precision (PR) and recall rate (RE).

First, we used word graphs annotated with D3 boundaries simulating 100% correct boundary classification to show how the

$\lambda_s$		DA class.		D3 class.	
<i>da</i>	<i>das</i>	DAA	DAC	PR	RE
1.00	0.00	68.3	70.0	100	100
0.50	0.50	59.9	62.0	100	100
0.80	0.20	69.9	71.5	100	100
<b>0.90</b>	<b>0.10</b>	<b>70.8</b>	<b>72.6</b>	<b>100</b>	<b>100</b>
0.98	0.02	69.6	71.4	100	100

**Table 3:** Recognition results in % using manually annotated word graphs.

recognition rates for DA classification improve, if only the 18 DA LMs (*da*) are used, and if the DA sequences LM (*das*) is added. The results are given in Table 3. The first line is the baseline system using only the 18 DA LMs and manually segmented word graphs. If we give an equal weight to both classifiers the results worsen, but a weight that compensates for the different value ranges yields improved recognition rates.

Second, we wanted to determine the best D3 segmentation and DA classification using the MLP (*mlp*) trained on B3 boundaries, the LM including M3 boundaries (*lm*), the boundary LM for D3 boundaries (*bound*), the 18 DA LM (*da*) and the DA sequences LM (*das*). This is done using an automatic optimization procedure to find the best weight configuration for  $\lambda_s$ . The optimization procedure minimizes the total costs of the best path for each utterance in a cross-validation set (here, the test set). Using the automatic optimization procedure, we achieved the results presented in Table 4.

Iteration	DAA	DAC	PR	RE
1	45.6	52.4	92	57
5	50.9	59.9	91	60
10	52.1	62.4	89	66
15	52.5	63.6	88	68
20	52.6	64.6	88	69

**Table 4:** Recognition results in % using an automatic optimization procedure for the weight configurations classifying DAs and boundaries.

One can see, that the recognition results for DA classification improve with each iteration. For the D3 segmentation the recall improves considerably with only a minor loss of precision. The results for DA classification are, of course, somewhat lower than the results shown in Table 3, because those experiments were performed based on manually DA-segmented utterances.

The best result was achieved using all knowledge sources with the following weight configuration:

<i>lm</i>	<i>da</i>	<i>das</i>	<i>mlp</i>	<i>bound</i>
0.25	0.27	0.06	0.22	0.20

In [28], we presented a sequential approach where a turn was first segmented and then the resulting segments were classified into DAs. If we proceed the same way on our new test set and use the same classifiers as for the integrated approach we achieve the results presented in Table 5.

DAA	DAC	PR	RE
47.3	62.0	71	73

**Table 5:** Recognition results in % achieved by performing segmentation and classification of DAs sequentially.

One can see that the integrated approach improves the DAA by over 5% points and the DAC by over 2% points. Even the segmentation accuracy improves a lot when both tasks are performed in an integrated procedure. These results show that the classification of boundaries and DAs based on the spoken word chain and the speech signal can be improved significantly by an integrated search procedure incorporating a number of knowledge sources.

The most important advantage of the  $A^*$  search is, however, not that it yields better results than a sequential approach but the possibility to work directly with the WHG. In a sequential approach, one has to work with the best word chain(s). This might do for basic research but it is, in the long run, not a feasible strategy for ‘real life’ systems.

### 3.3. A Hybrid Approach: Prosody, Statistics, and Partial Parsing

In this section, we want to focus not on the segmentation of turns with the help of boundary classification, but on accentuation, which is, in a way, an orthogonal complement to segmentation. Linguistic analysis in spoken dialogue systems has to cope with two main problems. First, spontaneous speech very often is fragmented, ungrammatical or exceeds the system’s boundaries (e.g. out-of-vocabulary words). Second, word recognition in spoken dialogue systems produces errors, thus rendering utterances ungrammatical on the syntactic as well as the semantic level. In order to cope with these problems, methods of robust parsing have been established. For example, partial parsing methods restrict syntactic analysis to sub-units of utterances only, therefore reducing the above mentioned problems to these sub-units. Different methods of partial parsing have been successfully employed in spoken dialogue systems, such as the systems described in [2] and [3].

Partial parsing in dialogue systems becomes even more efficient if more sophisticated sources of information, beyond acoustically scored word graphs and dialogue predictions, can be used to guide the linguistic processor. We concentrated on the integration of prosodic information, extracted from the speech signal, and statistically detected semantic concepts in utterances as additional support for the parser, thus resulting in a hybrid approach to language understanding. The units to be analyzed correspond to semantic concepts, e.g., time, date, source or target location for train timetable inquiries, or to DA classes, as, e.g., SUGGEST or ACCEPT, in the VERBMOBIL task. Such units are vital for the correct interpretation of the utterance in the application domain. The parser will identify and analyze these concepts, assigning a semantic representation to each.

For each concept and its possible surface realizations, grammar fragments are defined that may be used by the parser upon request. The parser is guided by prosodic information on phrase boundaries and phrase accents, telling it where to start the partial analysis. Statistical concept detection provides information on which semantic concepts are included by the current utterance, thus helping the parser to choose the appropriate grammar fragments. The use of grammar fragments has the following advantages: the danger of false alarms in parsing is drastically reduced, as well as the time consumed by the parser and the efforts for grammar development.

### Accent Information in Word Hypotheses Graphs

Here, we want to use prosodic information to determine the salient regions in a phrase. These regions are those parts of a sentence which hold the most important content words, e.g., time expressions and locations and which most of the time are ‘in focus’, i.e., are the carrier of the focal accent. To get information for those regions, we use a NN trained on a part of the VERBMobil database.

Using  $Score(A3 | w)$  and  $Score(\neg A3 | w)$  from the output nodes of the NN for each word  $w$  we can estimate the probability  $P(A3 | w)$  by using the following formula

$$P(A3 | w) = \frac{Score(A3 | w)}{Score(A3 | w) + Score(\neg A3 | w)}.$$

Now, we are able to estimate the probability  $P(A3 | w)$  for each word of an utterance. We decide for a focused region by using a threshold. In Figure 2, an example is given for a German utterance.

The estimation of stressed regions in a given utterance offers two possible methods of using this knowledge in combination with the parser:

1. The regions are ranked by their prosodic scores and the ranking list is given to the parser, which has to find the best expression for the given context.
2. A list of possible expressions from the parser is disambiguated using the prosodic score from the NN.

Both methods can efficiently be employed to find the best expression the parser is searching for in the context the concept predictor has estimated. The first way seems to be the better one if working on WHGs, because the parser only has to search in the best scored paths and, thus, search effort is smaller.

For the 18 DA classes sketched above, we estimated the most frequent stressed words of a subset of the VERBMobil database using the method described above. Only those words are considered, that exceed a threshold of 0.8 for the automatically calculated stress-probability in more than 80% of their occurrences. In Table 6 the ten most often observed words are shown, which fulfill this criterion when looking at all DAs. Table 7 shows the five most often observed stressed words for the most frequent DA classes SUGGEST and ACCEPT. In both tables the words are ranked by their frequency of occurrence in the observed data set.

The results from Tables 6 and 7 show that a successful classification of content words in an utterance is possible through determining the stressed words. Semantically important information can thus be obtained via the detection of the focused regions. This can be done by only using prosodic-acoustic features (Note that for the classification of DAs, function words that normally are not accentuated are important as well, cf. [33]).

### Statistical Concept Detection

As a second additional source of information for the hybrid partial parsing, we apply a statistical approach which uses  $n$ -gram LMs as semantic concept predictors. The model has to decide about the occurrence of special semantic concepts in word

$P(A3   w) > 0.8$		
Rank	% stressed	word (translation)
1	88.57	Freitag (Friday)
2	82.69	Wiederhören (bye)
3	84.31	Donnerstag (Thursday)
4	90.91	Samstag (Saturday)
5	95.35	neunzehnten (19th)
6	81.82	August (August)
7	96.15	vierundzwanzig. (24th)
8	87.50	achten (8th)
9	86.96	wunderbar (marvellous)
10	100.00	sechszwanzig. (26th)

**Table 6:** Automatically determined stressed words for all DAs.

ACCEPT		
$P(A3   w) > 0.8$		
Rank	% stressed	word (translation)
1	100.00	einverstanden (ok)
2	100.00	Ordnung (all right)
3	100.00	wunderbar (marvellous)
4	85.71	Freitag (Friday)
5	85.71	frei (free)

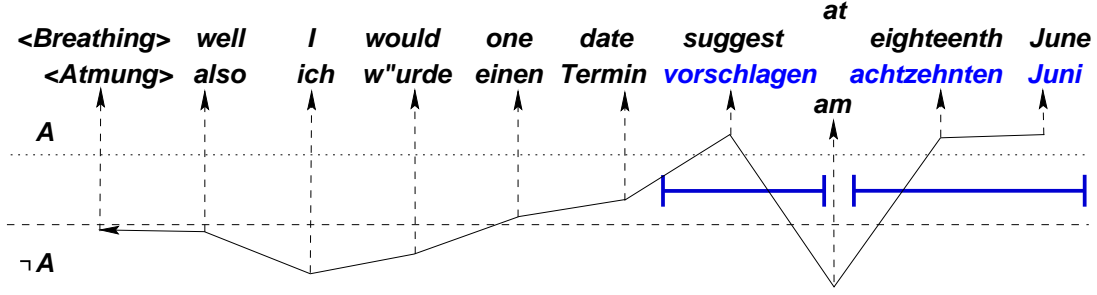
SUGGEST		
$P(A3   w) > 0.8$		
Rank	% stressed	word (translation)
1	82.22	Montag (Monday)
2	87.80	Freitag (Friday)
3	83.33	Donnerstag (Thursday)
4	82.76	Mittwoch (Wednesday)
5	93.10	Samstag (Saturday)

**Table 7:** Automatically determined stressed words for dialogue acts ACCEPT and SUGGEST.

chains. We show its usability on a corpus collected with the above mentioned information retrieval system containing the utterances used for the grammar development. In the following, we present two predictors, one for time expressions and one for date expressions. The predictor should be able to decide whether there appears such a time/date expression in an utterance or not.

If we use LMs as semantic concept predictors we have to claim for a word chain  $\mathbf{w}$  whether the concept we are looking for is expressed in  $\mathbf{w}$  or not. For this purpose we build two different LMs. The first one is trained with word chains expressing the semantic concept and the second one with the utterances not expressing it. During analysis the two scores for the incoming word chain are computed — for WHGs the best word chain in the graph is used — and the predictor with the higher probability is chosen. We apply category based LMs, rational interpolation for the LMs, and a context of three words. The ‘Semantic Concept Predictor’ results are shown in Table 8 as confusion matrices. One can see that our LM approach to the prediction task performs well enough and can therefore be used as a predictor for the semantic concept analysis, i.e. to select the correct grammar fragment. For the problem of detecting time expressions we obtain a recognition rate of 95.6%; if we measure the performance of the models





**Figure 2:** A German sentence from VERBMOBIL with probability  $P(A3 | w)$  for each word  $w$  and the two focused regions hypothesized (with word to word translation).

as being a time expression spotter we get a recall of 98.4% and a precision of 80.0%. For date expressions we have a recognition rate of 95.7%, a recall of 96.7%, and a precision of 82.4%.

	TIME	NOTIME
TIME	1127	18
NOTIME	282	5375
	DATE	NODATE
DATE	1191	41
NODATE	254	5316

**Table 8:** Confusion Matrices for time and date expressions

### Partial Parsing

The partial parser described here is an agenda driven chart parser, operating as an island parser (cf. [29]). As our approach restricts the linguistic analysis to the analysis of semantic concepts, lexicon and grammar of the parser only need to cover the relevant syntactic realizations for each concept, thus resulting in several grammar fragments, rather than one full grammar. Island parsing on the basis of these grammar fragments means that each of the maximal islands, found by the parser, corresponds to one relevant part of an utterance. We coded a grammar fragment for each of the semantic concepts in terms of a context-free phrase structure grammar. Thus, the predictions on the occurrence of concepts in user utterances can be used to guide the parsing process. This is done by using only those grammar fragments for parsing that correspond to semantic concepts predicted by the concept detection module. In order to further improve efficiency of the parsing process, prosodic information is included into the parsing process. Each word hypothesis contains a prosodic accent score, in addition to the usual acoustic score. This information is used for choosing the initial islands: only those hypotheses which are marked as accentuated are chosen as initial islands.

### The Parsing Algorithm

The chart is initialized with the lexical entries for the hypotheses in the WHG. As not every grammar fragment is used for each parse, many hypotheses are unknown, thus leaving gaps in the chart. In parallel to the chart, two agendas are initialized that guide the flow of the analysis. The first agenda (*seed agenda*) contains all hypotheses that serve as initial islands. The second agenda (*non-seed agenda*) contains the remaining hypotheses. Each hypothesis, whose accent score exceeds a given threshold, is inserted into the seed agenda, the remaining ones into the non-seed agenda. Within both agendas, entries are sorted according to their acoustic score. Agenda entries may not only be used as initial lexical entries (*seed entries*), but also pairs of chart edges

(*non-seed entries*) that comprise pointers to two adjacent chart edges and a list of grammar rules that might combine these two edges to a new one. The following steps are performed until no entries are left in the seed agenda.

1. Take best scored agenda entry  $E$  from seed agenda.
2. If  $E$  is a seed entry go to 3, else go to 4.
3. For each adjacent chart edge to  $E$  look for rules that can be applied to both and generate an agenda pair for both and sort it into seed agenda; go to 1.
4. For each grammar rule in  $E$ : apply this rule to both edges, insert new edge (if rule can be applied) into chart, generate new agenda pairs for this new edge and insert them into seed agenda; go to 1.

This is done for each of the predicted semantic concepts using the respective grammar fragments. Only if no valid semantic representation for a concept can be found in the chart after parsing, the process is restarted with the non-seed agenda.

First experiments were done for the semantic concepts *time* and *date*. Detailed results can be found in [32] and in Table 9 a limited version of the results is shown. The numbers in Table 9 denote the following: we examine the two concepts *time* and *date* and for our test database we count for each set of integrated knowledge how many sentences we have to parse i.e. how often the parser is applied. This number is given in the row *Parses*. Additionally we count how many words are on the initial seed agenda (row *Seeds*) as all these words must be considered when applying a grammar fragment. The used information set is composed from the following parts. As one information source we have the lexicons of the grammar fragments which is always applied except in the first column (NIL) where no knowledge is

	NIL	lex.	Pred.	Pros. 0.5	0.5+ Pred.
Parses time	871	346	136	219	124
date	871	292	169	244	133
Seeds time	2761	836	431	439	285
date	2761	605	416	447	308
D (time/date)	1/0	1/0	1/0	1/0	1/0
I (time/date)	2/2	2/2	3/1	2/2	3/1

**Table 9:** Number of necessary parses and possible island seeds with different levels of information sources and the number of deletions (D) and insertions (I) for *date* and *time*

used. As prediction (Pred.) we use our LM classifiers and for the prosodic scores we define a threshold and all words whose accentuation score is higher than that threshold are put on the seed agenda. In the last two rows we give the numbers of deletions (D) and insertions (I) the analysis results in for the two concepts. The first column with the results when no knowledge is used corresponds to an analysis using a full grammar. Through applying grammar fragments the numbers of column two (lex.) are obtained and so on. We see that the more knowledge we use for our hybrid approach the less parses we have to perform and the less island seeds have to be considered. As a consequence the processing time is reduced drastically without an increase of the error rate.

### 3.4. Looking back from the End: Adequate Evaluation and Adequate Design

In the VERBMOBIL system, each module evaluates and optimizes its analyses and classification results independently from the other modules, and there is an end-to-end-evaluation with the criterion: ‘Is the translation approximately correct?’. In addition, there is a more or less informal feed-back from the higher linguistic modules to the lower ones. If it is not (yet) possible to ‘formalize’ such a feed-back, it should at least be intensified. The criterion should not simply be the *correctness* of the translation, but the *success* of the communication. Sometimes, underspecification will do; this will be shown with the following example.

Let us assume the following utterance with correct word recognition, parse, and subsequent translation:

... *Um zwei Uhr nachmittags. Wollen wir uns am Berliner Hauptbahnhof treffen?*

... *At two p.m.. Should we meet at Berlin main station?*

If, however, the boundary between *nachmittags* and *wollen* is not recognized, and if there is no prosodic question at the end of the turn, the parse and the translation would result in:

... *Um zwei Uhr nachmittags wollen wir uns am Berliner Hauptbahnhof treffen.*

... *At two p.m., we want to meet at Berlin main station.*

Here, prosody and especially intonation are irrelevant for the classification of sentence mood because the speaker can produce a final rise or a final fall [5, 10]. If the turn is translated as a statement, then the segmentation is wrong, the proposition ( $\approx$  the salient words) is right; illocution and translation are wrong because the sentence mood is not reproduced correctly. The perlocution, however, is successful: no matter whether the correct or the wrong translation is generated, if the dialogue partner accepts, e.g., with an *ok.*, and if the dialogue partners meet at the given time and place, the communication is felicitous — even if the translation is wrong. This means that there are fatal and harmless errors which should be treated differently in the evaluation. In the design of the system, it might be better to leave such alternatives underspecified. We thus believe that a local optimization — e.g., recognition rates for DA classification — can only be an intermediate step towards the ‘ultimate’ evaluation within an existing dialogue system.

These arguments corroborate the findings presented in Section

3.3. In a deep linguistic analysis, we thus should leave underspecified certain distinctions, in a shallow analysis, we can concentrate on partial parses. A similar argumentation can be found in [4, p. 32] who compare the impact of insertions, substitutions, and deletions on user acceptance: deletions cause little problems, substitutions are more serious, and most serious are insertions, because “... the systems seems to ‘assume’ something the caller never mentioned.”

## 4. CONCLUDING REMARKS

In this paper, we focussed on the ‘What’ and on the ‘How’ concerning the use of prosody in automatic dialogue systems: what we are working with and what are we working on in Section 2, and how we are working, i.e., which methodology we should use, in Section 3. There, we gave some examples for the present state of the art and for promising trends out of our own work. We put a stronger emphasis on **shallow** analysis, automatic learnability and an easy adaptation to new applications. We aim at an integration of all available knowledge sources in a global search procedure; hard decisions should be taken as late as possible. A flexible use of knowledge sources means at least:

1. From a paradigmatic point of view, we should use those units we are interested in, if enough data are available. As a fall back, we can use substitutes that model these units indirectly, e.g., syntactic boundaries instead of dialogue act boundaries.
2. From a syntagmatic point of view, we should use the maximum context available for a given database. As a fall back, we can use less context if this can be modelled more adequately. Thus, for our language model for syntactic boundaries, we use trigrams, for dialogue act classification, we use 4-grams, and for dialogue act sequences, we use bigrams.
3. From a pragmatic point of view, we should concentrate on those parts of an utterance that contain the crucial information (e.g., partial parsing, accentuated words in focus).

We assume that such highly sophisticated methods correspond closely with the strategies of human beings in human-human communication — but this is, basically, yet another story.

## 5. REFERENCES

1. *Proc. ESCA Workshop on Dialogue and Prosody*, Eindhoven, Netherlands, 1999.
2. D. Albesano, P. Baggia, M. Danieli, R. Gemello, E. Gerbino, and C. Rullent. A Robust System for Human-Machine Dialogue in Telephony-Based Applications. *International Journal of Speech Technology*, 2:101–111, 1997.
3. H. Aust, M. Oerder, F. Seide, and V. Steinbiss. The Philips Automatic Train Timetable Information System. *Speech Communication*, 17:249–262, 1995.
4. H. Aust and O. Schröder. Application Development with the Philips Dialog System. In *Proceedings of 1998 International Symposium on Spoken Dialogue (ISSD 98)*, pages 27–34, Sydney, Australia, 1998.
5. A. Batliner. Ein einfaches Modell der Frageintonation und seine Folgen. In E. Klein, F. Pouradier Duteil, and K.H. Wagner, editors, *Betriebslinguistik und Linguistikbetrieb*, pages 147–160. Niemeyer, Tübingen, 1991.

6. A. Batliner, J. Buckow, R. Huber, V. Warnke, E. Nöth, and H. Niemann. Prosodic Feature Evaluation: Brute Force or Well Designed? In *Proc. of the 14th Int. Congress of Phonetic Sciences*, volume 3, pages 2315–2318, San Francisco, 1999.
7. A. Batliner, A. Kießling, S. Burger, and E. Nöth. Filled Pauses in Spontaneous Speech. In *Proc. of the 13th Intl. Congress of Phonetic Sciences*, volume 3, pages 472–475, Stockholm, August 1995.
8. A. Batliner, R. Kompe, A. Kießling, M. Mast, H. Niemann, and E. Nöth. M = Syntax + Prosody: A syntactic-prosodic labelling scheme for large spontaneous speech databases. *Speech Communication*, 25:193–222, 1998.
9. A. Batliner, M. Nutt, V. Warnke, E. Nöth, J. Buckow, R. Huber, and H. Niemann. Automatic Annotation and Classification of Phrase Accents in Spontaneous Speech. In *EUROSPEECH 99* [14], pages 519–522.
10. A. Batliner, C. Weiland, A. Kießling, and E. Nöth. Why sentence modality in spontaneous speech is more difficult to classify and why this fact is not too bad for prosody. In D. House and P. Touati, editors, *Working Papers, Prosody Workshop 1993*, pages 112–115, Lund, Sweden, 1993.
11. H.U. Block. The Language Components in Verbmobil. In *ICASSP 97* [20], pages 79–82.
12. T. Bub and J. Schwinn. Verbmobil: The Evolution of a Complex Large Speech-to-Speech Translation System. In *ICSLP 96* [21], pages 1026–1029.
13. J. Carletta, N. Dahlbäck, N. Reithinger, and M. Walker. Standards for Dialogue Coding in Natural Language Processing. Dagstuhl-Seminar-Report 167, 1997.
14. *Proc. European Conf. on Speech Communication and Technology*, Budapest, Hungary, 1999.
15. F. Gallwitz, M. Aretoulaki, M. Boros, J. Haas, S. Harbeck, R. Huber, H. Niemann, and E. Nöth. The Erlangen Spoken Dialogue System EVAR: A State-of-the-Art Information Retrieval System. In *Proceedings of 1998 International Symposium on Spoken Dialogue (ISSD 98)*, pages 19–26, Sydney, Australia, 1998.
16. F. Gallwitz, A. Batliner, J. Buckow, R. Huber, H. Niemann, and E. Nöth. Integrated Recognition of Words and Phrase Boundaries. In *Proc. Int. Conf. on Spoken Language Processing*, volume 7, pages 2883–2886, Sydney, Australia, December 1998.
17. F. Gallwitz, H. Niemann, E. Nöth, and V. Warnke. Prosodic Information for Integrated Word-and-Boundary Recognition. In *Proc. ESCA Workshop on Dialogue and Prosody* [1], pages 163–168.
18. M. Grice, M. Reyelt, R. Benz Müller, J. Mayer, and A. Batliner. Consistency in Transcription and Labelling of German Intonation with GToBI. In *ICSLP 96* [21], pages 1716–1719.
19. H. Hörmann. *Meinen und Verstehen*. Suhrkamp Taschenbuch Wissenschaft. Suhrkamp-Verlag, Frankfurt, 1978.
20. *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Munich, April 1997. IEEE Computer Society Press.
21. *Proc. Int. Conf. on Spoken Language Processing*, Philadelphia, September 1996.
22. S. Jekat, A. Klein, E. Maier, I. Maleck, M. Mast, and J. Quantz. Dialogue Acts in VERBMOBIL verbmobil-report-65–95, April 1995.
23. W. Kasper, B. Kiefer, H.U. Krieger, C.J. Rupp, and K. Worm. Charting the Depths of Robust Speech Parsing. In *Proc. of the 37th meeting of the ACL*, page (to appear), 1999.
24. A. Kießling. *Extraktion und Klassifikation prosodischer Merkmale in der automatischen Sprachverarbeitung*. Berichte aus der Informatik. Shaker Verlag, Aachen, 1997.
25. A. Kießling, R. Kompe, A. Batliner, H. Niemann, and E. Nöth. Automatic Labeling of Phrase Accents in German. In *Proc. Int. Conf. on Spoken Language Processing*, volume 1, pages 115–118, Yokohama, 1994.
26. R. Kompe. *Prosody in Speech Understanding Systems*. Lecture Notes for Artificial Intelligence. Springer-Verlag, Berlin, 1997.
27. W. Lea. Prosodic Aids to Speech Recognition. In W. Lea, editor, *Trends in Speech Recognition*, pages 166–205. Prentice-Hall Inc., Englewood Cliffs, New Jersey, 1980.
28. M. Mast, R. Kompe, St. Harbeck, A. Kießling, H. Niemann, and E. Nöth. Dialog Act Classification with the Help of Prosody. In *ICSLP 96* [21], pages 1728–1731.
29. K. Mecklenburg, P. Heisterkamp, and G. Hanrieder. A Robust Parser for Continuous Spoken Language using Prolog. In *Proceedings of the Fifth International Workshop on Natural Language Understanding and Logic Programming (NLULP 95)*, pages 127–141, Lisbon, 1995.
30. H. Niemann, E. Nöth, A. Batliner, J. Buckow, F. Gallwitz, R. Huber, and V. Warnke. Using Prosodic Cues in Spoken Dialog Systems. In *Proceedings of the International Workshop on Speech and Computer*, pages 17–28, St. Petersburg, Russia, October 1998.
31. N.J. Nilsson. *Principles of Artificial Intelligence*. Springer-Verlag, Berlin, 1982.
32. E. Nöth, M. Boros, J. Haas, V. Warnke, and F. Gallwitz. A Hybrid Approach to Spoken Dialogue Understanding: Prosody, Statistics and Partial Parsing. In *EUROSPEECH 99* [14], pages 2019–2022.
33. M. Nutt, A. Batliner, V. Warnke, and E. Nöth. Using Phrase Accent Information for Dialogue Act Recognition in Spontaneous German Speech. In *Proc. ESCA Workshop on Dialogue and Prosody* [1], pages 151–155.
34. M. Reyelt and A. Batliner. Ein Inventar prosodischer Etiketten für Verbmobil. Verbmobil Memo 33, 1994.
35. E. Shriberg, R. Bates, P. Taylor, A. Stolcke, D. Jurafsky, K. Ries, N. Cocarro, R. Martin, M. Meteer, and C. Van Ess-Dykema. Can Prosody Aid the Automatic Classification of Dialog Acts in Conversational Speech? *Language and Speech* 41, pages 439–487, 1998.
36. J. Spilker, H. Weber, and G. Görz. Detection and Correction of Speech Repairs in Word Lattices. In *EUROSPEECH 99* [14]. (to appear).
37. V. Strom and C. Widera. What's in the "Pure" Prosody? In *ICSLP 96* [21], pages 1497–1500.
38. P. Taylor, S. King, S. Isard, and H. Wright. Intonation and Dialogue Context as Constraints for Speech Recognition. *Language and Speech* 41, pages 489–508, 1999.
39. J. Vaissière. The Use of Prosodic Parameters in Automatic Speech Recognition. In H. Niemann, M. Lang, and G. Sagerer, editors, *Recent Advances in Speech Understanding and Dialog Systems*, volume 46 of *NATO ASI Series F*, pages 71–99. Springer-Verlag, Berlin, 1988.
40. W. Wahlster, T. Bub, and A. Waibel. Verbmobil: The Combination of Deep and Shallow Processing for Spontaneous Speech Translation. In *ICASSP 97* [20], pages 71–74.
41. V. Warnke, F. Gallwitz, A. Batliner, J. Buckow, R. Huber, E. Nöth, and A. Höthker. Integrating Multiple Knowledge Sources for Word Hypotheses Graph Interpretation. In *EUROSPEECH 99* [14], pages 235–239.
42. V. Warnke, R. Kompe, H. Niemann, and E. Nöth. Integrated Dialog Act Segmentation and Classification using Prosodic Features and Language Models. In *Proc. European Conf. on Speech Communication and Technology*, volume 1, pages 207–210, 1997.