



LUND UNIVERSITY

Computational extraction of lexico-grammatical information for generation of Swedish intonation

Horne, Merle; Filipsson, Marcus

Published in:
Proceedings of ESCA/IEEE Workshop on Speech Synthesis

1994

[Link to publication](#)

Citation for published version (APA):
Horne, M., & Filipsson, M. (1994). Computational extraction of lexico-grammatical information for generation of Swedish intonation. In *Proceedings of ESCA/IEEE Workshop on Speech Synthesis* (pp. 220-223). ESCA.

Total number of authors:
2

General rights

Unless other specific re-use rights are stated the following general rights apply:
Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: <https://creativecommons.org/licenses/>

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

LUND UNIVERSITY

PO Box 117
221 00 Lund
+46 46-222 00 00

COMPUTATIONAL EXTRACTION OF LEXICO-GRAMMATICAL INFORMATION FOR GENERATION OF SWEDISH INTONATION

Merle Horne and Marcus Filipsson

Department of Linguistics and Phonetics, Helgonabacken 12, S-223 62 Lund, Sweden

ABSTRACT

This article presents a discussion of a number of algorithms being developed which will enable the generation of prosodic structure for Swedish restricted texts. These algorithms, including a word-class tagger, a complex-word identifier and a prosodic parser form part of a linguistic preprocessor to a text-to-speech system for generation of intonation.

PROSODIC STRUCTURE

One of the goals of current research in text-to-speech systems is to improve the quality of intonation by developing algorithms for preprocessing texts in order to extract grammatical and discourse information necessary for the generation of appropriate prosodic patterns. The present article will describe the work we are currently carrying out aimed at using the information on coreferentiality obtained from the referent tracking algorithm previously developed (see Horne et al. 1993) together with further information on lexico-syntactic category designation to group words together into a hierarchy of prosodic constituents. Whereas the referent-tracking process is important to the Fo-generating component in order to be able to predict the distribution of focal and non-focal accents, information on prosodic structure is needed in order to better predict the location as well as the particular form of tone accents associated with utterance-internal prosodic boundaries.

Following an approach similar to Bachenko & Fitzpatrick (1990), Quené & Kager (1993) and inspired by concepts within prosodic phonology (e.g. Nespor & Vogel 1986), we are attempting to determine how one, using a minimal amount of parsing, can obtain enough information to construct a hierarchical prosodic structure for each sentence in a text. Unlike other researchers, however, we are also using contextual information such as coreference in our approach to generating prosodic structure.

At least three levels of prosodic structure are required for Swedish in order to model all the prosodic information observed in our data (Horne 1994). The smallest of these is the Prosodic Word which we will define as corresponding to a content word and any following function words up to the next content word within a given clause. At the beginning of a clause, the Prosodic Word can also begin with one or more function words. The Prosodic Word is characterized by a word accent and potentially a focal accent (Accent 1 = HL*(HL⁻), Accent 2 = H*L(H⁻L⁻) (We use H⁻ and L⁻ to represent respectively a focal high and the low tone accent following a focal high in order to distinguish them from the H and L associated with the word accents.). It is also marked by a boundary tone which is realized by a final rise in the case where the content word is not focussed (i.e. contextually given) (H#) or a fall when the content word is focussed (L#). These boundary tones, we claim, play an important role in creating the transitions between consecutive Prosodic Words in a larger Prosodic Phrase. They are also points for potential pauses, e.g. before focussed content words (see Gårding 1967, Strangert 1993). The unit does not necessarily correspond to a syntactic constituent as the example in (1) illustrates ('-' represents the boundary between Prosodic Words). This type of 'nonsyntactic' grouping is perhaps more characteristic of well-planned read texts or spontaneous speech than of non well-planned texts read e.g. by a non-expert/non-professional.

- (1) Kurserna på – Stockholmsbörsen – fortsätter att – falla.
Rates(det) on – Stockholm Stock Exchange(det) – continue to – fall
'Rates on Stockholm's Stock Exchange continue to fall'

One or more Prosodic Words make up a Prosodic Phrase which is marked by a final L% or H% boundary tone accent. Factors which determine the location of Prosodic Phrase boundaries include the following: a) sentence boundary: A sentence boundary corresponds to the end of a Prosodic Phrase, b) new/given distinction: A Prosodic Phrase must contain at least one focussed Prosodic Word, c) length: A Prosodic Phrase will not exceed x syllables at a given rate of speech y . Finally, one or more Prosodic Phrases make up a Prosodic Utterance, which is bounded by pauses. It is further generally assumed that each prosodic constituent is characterized by a certain amount of preboundary lengthening (Gussenhoven & Rietveld 1992, Wightman et al. 1992), and although we have not as yet made any detailed investigations of the phenomenon in our data which would allow us to quantify a lengthening index, we are assuming that, all other things being equal, the higher up in the hierarchy a prosodic constituent is placed, the greater the relative duration associated with its final syllable(s) will be.

Figure 1 presents in schematic form the prosodic constituents assumed for Swedish and their phonetic correlates. The tone accents (H and L) are assumed to be associated with syllables (S) according to principles outlined in Bruce (1977). It is also assumed that the realization of the tone accents is dependent to some extent on the number of syllables present in a particular word, i.e. the number of syllables in a given word dictates to a great extent how many tones will be realized phonetically.

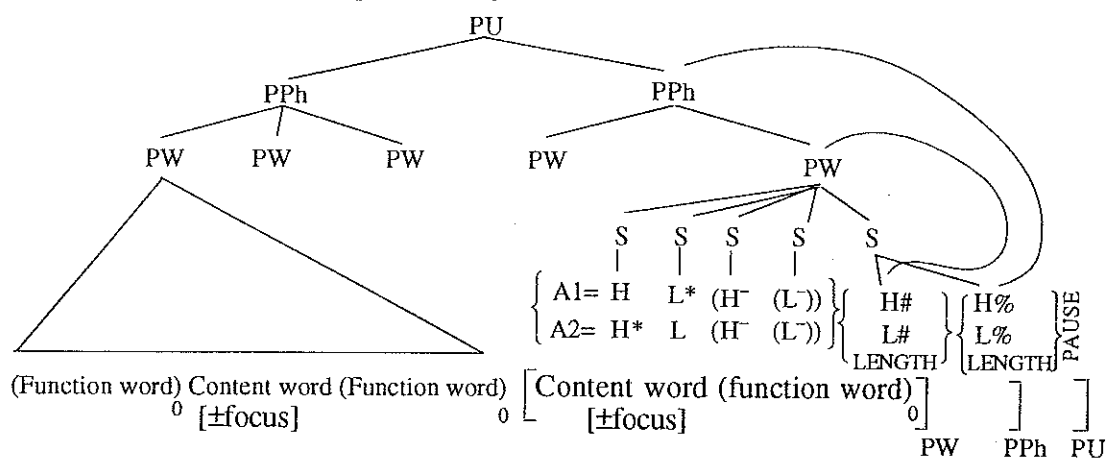


Figure 1. Schematic presentation of the prosodic hierarchy assumed for Swedish and the associated phonetic correlates. Accent 1 is represented as $HL^*(H^-L^-)$ and Accent 2 as $H^*L(H^-L^-)$, where (H^-L^-) represents the focal High (H^-) and potential Low (L^-) associated with the focal accent. $H\#$ and $L\#$ represent the Prosodic Word boundaries and $H\%$ and $L\%$ designate the Prosodic Phrase boundaries. PW stands for Prosodic Word, PPh for Prosodic Phrase and PU for Prosodic Utterance. $(\text{Function word})_0$ stands for zero or more function words.

DESIGN OF THE PROSODIC STRUCTURE COMPONENT

In order to construct these prosodic constituents automatically, a number of different analyses are required. The present system is based on a strictly modular approach, with each module having well-defined input/output formats. This will enable us to easily replace a module with a new one if a more efficient algorithm is developed at a later stage.

The first task is to tokenize the text into a list of words. At the same time, punctuation marks and paragraph boundaries are recognized. The next step is to look up the words in our domain-specific lexicon, which is an expanded subset of a larger computerized lexicon (Hedelin et al. 1987). This process will generate multiple tags for some words. The next step is therefore the disambiguation of these. In this endeavour, we are currently testing the performance of a stochastic parser based on lexical and sequential occurrence probabilities as well as overall tag probability (Eeg-Olofsson 1991). The algorithm implements a first-order Markov chain and uses dynamic programming to estimate the best hypothesis for the whole sentence. A set of approximately 30 lexico-syntactic tags based on Ejerhed et al.'s tag set (1992) have been

chosen to train the system. These have been further assigned to the tagged words' lemma representations in the computerized lexicon, thus allowing recognition of all morphologically derived forms of a given head-word. Preliminary results indicate that the algorithm works quite well, but we intend to compare it with other approaches. One involves a Hidden Markov model such as in the Xerox Part-of-Speech Tagger (Cutting et al. 1987). Another approach is a rule-based one. Finally, we are considering combinations of these, e.g. using a rule-based system as the default method, and a probabilistic algorithm for cases where the rules fail.

After word classes are determined, the next stage is to recognize complex words, i.e. strings of content words that function as a single prosodic unit. In the stock-market domain, these correspond to proper names (i.e. company/bank names and stock designations, e.g. 'Avesta Sheffield', 'S-E Banken', 'Hennes & Mauritz', 'Hasselförs Förvaltnings AB' (AB 'CO.')). These strings are assigned a specific tag ('CX'-complex word) which, although it is not a lexical tag, is a member of the class of content word tags together with those associated with nouns, adjectives, verbs, adverbs, etc.

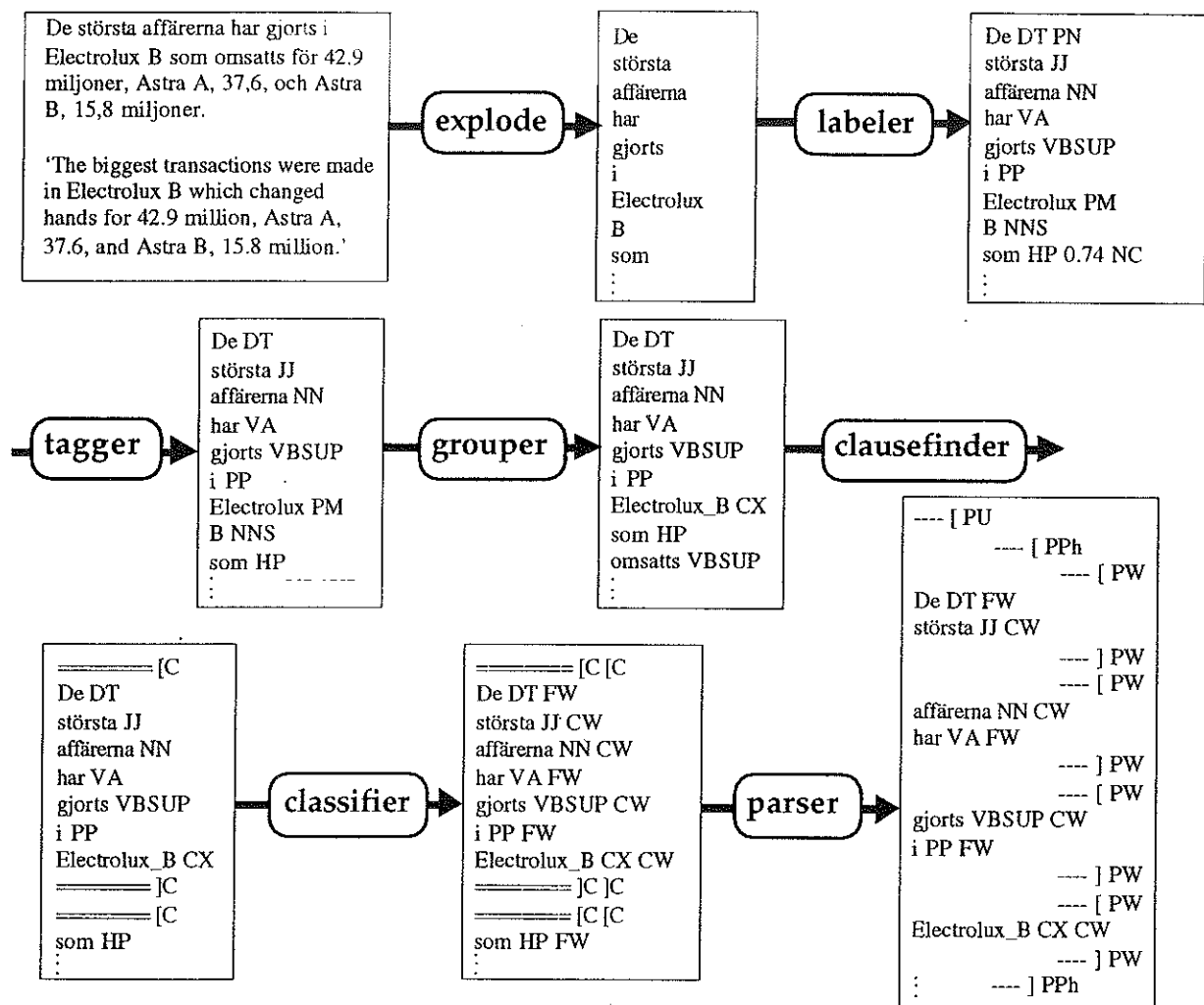


Figure 2. Schematic presentation of the present computer system for prosodic parsing. The modules in the system are represented by a rounded corner rectangle. An excerpt of the output of each module (and consequently the input to the next module) is shown between each pair of modules. The first input is the stock market report newspaper text; the final output is the prosodically parsed text. (DT=Determiner, PN=Pronoun, JJ=Adjective, NN=Noun, VA=Auxiliary Verb, VBSUP=Supine form of Verb, PP=Preposition, PM=Proper Noun, NNS=Noun specifier, HP=Relative Pronoun, NC='non-clausal' conjunction)

The next step is to recognize clause boundaries since the clause is the domain over which Prosodic Words are defined. Clause boundaries occur at certain punctuation marks, e.g. full stop, colon, semicolon, some commas (those not occurring in lists of words having the same word class), as well as before coordinate and subordinate conjunctions (*och* 'and', *men* 'but', *fast* 'although', *att* 'that'), and relative pronouns (e.g. *som* 'that', 'who').

The following stage involves classifying each word as either a content word ('CW') or a function word ('FW'). The assignment of words to one of these classes is not always straightforward, but one can say that in general, content words include the traditional categories of nouns, verbs, adjectives, adverbs, numerals, whereas function words consist of prepositions, pronouns, determiners, auxiliary verbs, interrogative/relative adverbs, deictic adverbs, quantifiers, etc. Domain-specific considerations lead to the introduction of a number of unconventional tags, for example 'specifier' nouns and adjectives that occur after the head noun in complex proper names like *B fria* in the name *Electrolux B fria* 'Electrolux B free (shares)'.

The final stage of the system is the actual prosodic parser, which parses the list of words into a hierarchical structure with three levels: Prosodic Word, Prosodic Phrase and Prosodic Utterance. First, content words and function words are grouped together to form Prosodic Words (see Fig. 1). Second, clause boundaries currently generate Prosodic Phrase boundaries, although other factors such as length must also be taken into consideration when determining the location of these boundaries. These are currently being incorporated into the parser. Finally, a Prosodic Utterance boundary is generated at each sentence boundary in the present algorithm.

Figure 2 presents all the modules in the system and the output from each module.

ACKNOWLEDGEMENTS

This research has been supported by a grant from the HSFR/NUTEK Language Technology Programme.

REFERENCES

- Bachenko, J. & E. Fitzpatrick. 1990. 'A computational grammar of discourse-neutral prosodic phrasing in English'. *Computational Linguistics* 16, 155-170.
- Bruce, G. 1977. *Swedish accents in sentence perspective*. Lund: Gleerups.
- Cutting, D., J. Kupiec, J. Pedersen, & P. Sibun. 1992. 'A practical part-of-speech tagger'. *Proceedings of the Third Conference on Applied Natural Language Processing, Trento, Italy, April 1992*. Also available as Xerox PARC technical report SSL-92-01.
- Eeg-Olofsson, M. 1991. *Word-class tagging. Some computational tools*. Univ. of Göteborg: Dept. of Linguistics.
- Ejerhed, E., G. Källgren, O. Wennstedt, & M. Åström. 1992. *The linguistic annotation system of the Stockholm-Umeå corpus project*. Umeå: Dept. of Linguistics Report No. 33.
- Gussenhoven, C. & A.C.M. Rietveld. 1992. 'Intonation contours, prosodic structure and preboundary lengthening'. *Journal of Phonetics* 20, 283-303.
- Gårding, E. 1967. 'Prosodiska drag i spontant och uppläst tal'. In G. Holm (ed.) *Svenskt talspråk*, 40-85. Stockholm: Almqvist & Wiksell.
- Hedelin, P., A. Jonsson, & P. Lindblad. 1987. *Svenskt uttalslexikon: 3 ed.* Tech. Report, Chalmers' Univ. of Technology.
- Horne, M. 1994. Generating prosodic structure for synthesis of Swedish intonation. *Working Papers* (Dept. Ling., Univ. of Lund) 43, 72-75.
- Horne, M., M. Filipsson, M. Ljungqvist, & A. Lindström. 1993. 'Referent tracking in restricted texts using a lemmatized lexicon: implications for generation of prosody'. *Proceedings Eurospeech '93* (Berlin) Vol. 3, 2011-2014.
- Nespor, M. & I. Vogel. 1986. *Prosodic phonology*. Dordrecht: Foris.
- Quené, H. & R. Kager. 1993. 'Prosodic sentence analysis without parsing'. In Vincent van Heuven & Louis Pols (eds.), *Analysis and synthesis of speech*, 115-130. Berlin: Mouton de Gruyter.
- Strangert, E. 1993. 'Speaking style and pausing'. *Phonum* 2, 121-137.
- Wightman, C.W., S. Shattuck-Hufnagel, M. Ostendorf & P. Price. 1992. 'Segmental durations in the vicinity of prosodic phrase boundaries'. *J. Acoust. Soc. Am.* 91, 1707-1717.