### From Sentences to Scope Relations and Backward

Gábor Alberti, Márton Károly and Judit Kleiber

Department of Linguistics, University of Pécs, 6 Ifjúság Street, 7624 Pécs, Hungary

**Abstract.** As we strive for sophisticated machine translation and reliable information extraction, we have launched a subproject pertaining to the revelation of *reference* and *information structure* in (Hungarian) declarative sentences. The crucial part of information extraction is a procedure whose input is a sentence, and whose output is an information structure, which is practically a set of possible operator scope orders (*acceptance*). A similar procedure forms the first half of *machine translation*, too: we need the information structure of the source-language sentence. Then an opposite procedure should come (*generation*), whose input is an information structure, and whose output is an intoned word sequence, that is, a sentence in the target language. We can base the procedure of acceptance (in the above sense) upon that of generation, due to the reversibility of Prolog mechanisms. And as our approach to grammar is "totally lexicalist", the lexical description of verbs is responsible for the order and intonation of words in the generated sentence.

## **1** Generating and Accepting Hungarian Sentences

As we strive for a sophisticated level of machine translation and reliable information extraction, we have launched a subproject pertaining to the revelation of reference and information structure in declarative sentences.

We are primarily working with data from Hungarian, which is known to be a language with a very rich and explicit information structure (consisting of different types of topics, quantifiers and foci) [1], [2], [3], [4] and an also quite explicit system of four degrees of referentiality [5], [6], [7], [8], including the indefinite specific degree [9] [10]. The kind of input we consider is an ordered set of (Hungarian) words furnished with four stress marks ("unstressed" / "STRESSED" / "FOCUS-STRESSED" / "CONTRASTIVELY STRESSED  $\downarrow$ ") – and our program decides if they constitute a wellformed sentence at all, with arguments of appropriate degrees of referentiality and a possible information structure, and delivers these semantic data, including the possible scope orders of topics, quantifiers and foci. We call this direction acceptance. We also try to "accept" sequences of words without stress marks: in this case the first step is furnishing them with all possible intonation patterns. A further kind of input is the opposite direction, which can be called generation, whose output is an intoned sentence. Generation is based upon the rich lexical description of tensed verbs pertaining to the sentence-internal arrangement and checking of their arguments; and – in harmony with our "totally lexicalist" approach to grammar [11], which can be regarded as a successor of Hudson's [12] Word Grammar or

Karttunen's [13] *Radical Lexicalism*, and a formal execution of cognitive ideas similar to those of Croft's [14] *Radical Construction Grammar* – special intra-lexical *generator rules* are responsible for the development of the intricate pre-verbal operator zone of sentences.

In what follows, Sec2 provides a review of the relevant linguistic phenomena, then Sec3 elucidates what we mean by "accepting" potential sentences with or without stress marks and "generating" sentences; and finally we speak about implementation, our theoretical and practical work in progress driven by computational aims.

#### 2 Referentiality Requirements and Information Structure

Hungarian, similar to English in this respect, has an indefinite article (egy `a(n)`) and a definite article (a(z) `the`) to distinguish different *degrees of referentiality*. This fact seems to suggest two degrees of referentiality, but a closer look to complex facts (in English, in Hungarian and even in Finnish, which lacks articles) proves that there are (at least) three degrees of positive referentiality in the semantic background of Universal Grammar (see example series (1-5) below), besides the lack of referentiality as a fourth degree, which occurs in Hungarian even in the case of countable nouns, as will be shown in (7) below [8]:

|  | Table 1. The four degrees | s of referentiality (and the | eir expression in Hungarian). |
|--|---------------------------|------------------------------|-------------------------------|
|--|---------------------------|------------------------------|-------------------------------|

| non-referential             | referential  |            |            |
|-----------------------------|--------------|------------|------------|
| non-s                       | pecific      | spe        | cific      |
|                             | non-definite |            | definite   |
|                             |              |            |            |
| $\emptyset$ (bare singular) | egy 'a(n)'   | egy 'a(n)' | a(z) 'the' |

The indefinite article is claimed to refer to a specific referentiality: "its referent is a subset of a set of referents already in the domain of discourse" [10] in the English sentence (1e) below, in opposition to the one in the *there* construction (1b):

Example 1. Degrees of referentiality – in English: three (positive) degrees.

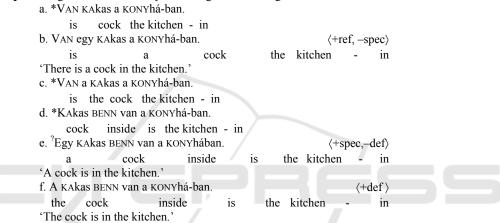
| a. *There is cock in the kitchen.           |   |
|---|---|
| b. There is a cock in the kitchen.          | $\langle +ref, -spec \rangle$               |
| c. *There is the cock in the kitchen.       |   |
| d. *Cock is in the kitchen.                 |   |
| e. <sup>(?)</sup> A cock is in the kitchen. | $\langle +\text{spec}, -\text{def} \rangle$ |
| f. The cock is in the kitchen.              | $\langle +def \rangle$                      |
|   |   |

Even without articles, Finnish can differentiate the three degrees of positive referentiality, by means of word order ( $\langle -spec \rangle$ : (2a) vs.  $\langle +spec \rangle$ : (2b-c)) and number agreement ( $\langle -def \rangle$ : (2b) vs.  $\langle +def \rangle$ : (2c)):

| Example 2. Degrees of referentiality - in Finnish: three (positive) degr                    | rees, too.                                  |
|---|---|
| a. Tul - i kaksi suomalais-ta tyttö-ä.  | $\langle +ref, -spec \rangle$               |
| come-Past3Sg two Finnish - Part girl-Part came.'  | 'Two Finnish girls                          |
| b. Kaksi suomalais-ta tyttö-ä tul - i.  | $\langle +\text{spec}, -\text{def} \rangle$ |
| two Finnish - Part girl-Part come-Past3Sg<br>'Two Finnish girls came (out of the four, say, | , that we expect to come).'                 |
| c. Kaksi suomalais-ta tyttö-ä tul - i - vat.  | $\langle +def \rangle$                      |
| two Finnish - Part girl-Part come-Past-3Pl came.'   | 'The two Finnish girls                      |

In Hungarian, too, there are similar constructions triggering the Non-Specificity Effect (3b), as well as the Specificity Effect (3e):

Example 3. Degrees of referentiality - in Hungarian: I. Being.



The system is that Patients of verbs expressing being (3a-c), coming into being (4a) and bringing into being (5a) show the Non-Specificity Effect, whilst these verbs regularly have counterparts with Patients showing an opposite, Specificity, effect: see (3e) above, and (4b), (5b) below.

**Example 4-5.** Degrees of referentiality – in Hungarian: II-III.. Coming / Bringing into being 4a. ÉRkez-ett [\*Ø/egy/\*a] MExikói a KONferenciá-ra.

arrive-Past [ $\emptyset$  / a / the ] Mexican the conference - onto

'A Mexican arrived at the conference.'

b. MEG-érkez-ett [\*Ø / <sup>(?)</sup>egy / a] MExikói a KONferenciá-ra. Perf - arrive - Past [Ø / a / the] Mexican the conference - onto

'A/The Mexican has/had arrived at the conference.'

5a. A GYErek-ek Alakít-ott-ak [ \*Ø / egy / \*az ] Énekkar-t a MŰsor-ra.

the child - ren form - Past - Pl [ $\varnothing$  / a / the ] choir - Acc the show-onto

'The children formed a choir for the show.'

b. A GYErek-ek MEG-alakít-ott-ak [\*Ø / <sup>(?)</sup>egy / az ] Énekkar-t a MŰsor-ra. the child - ren Perf - form - Past-3Pl [Ø / a / the ] choir - Acc the show-onto

'The children have formed a/the choir for the show.'

102

Certain argument positions, thus, undergo positive or negative specificity requirements. Alike, referentiality itself can be studied. As a default, arguments seem to be required to be referential, at least in the post-verbal zone of neutral Hungarian sentences [8]. See the summary in (6):

Example 6. Pos. and neg. ref.-degree requirements in the post-verbal zone of Hung. Sentences.

+ref: (3a), (3d), (4), (5) +spec: (3d-f), (4b), (5b) -spec: (3a-c), (4a), (5a)

Things are even more complicated, however: the above listed requirements can all be neutralized in the pre-verbal operator zone of Hungarian sentences; see (7-9) below. The neutralization of the  $\langle +\text{ref} \rangle$  requirement may result in well-formed nominal expressions without any kind of article, as is shown in (7). Such non-referential nominal expressions can occur even in neutral sentences, due to the special pre-verbal position (drawing the stress to itself from the verb stem), occupied by the Patient in (7a), for instance. The neutral meaning content in (8a) below, thus, can be realized in the three word orders listed in the example. Whilst in the case of an adjectival argument, see (8b), the pre-verbal position is the only "shelter" from the Referentiality Requirement, illustrated in (5) above.

**Example 7.** A few ways of neutralizing positive referentiality-degree requirements in the preverbal zone of Hungarian sentences

| a. V-modi     | fier  |   |
|---------------|---|---|
| (M)           | A GYErek-ek <i>énekkar-t</i> alakít-ott-ak a MŰsor-                 | ra.   |
|               | the child - ren choir - Acc form-Past-3Pl the show-on               | to  |
|               | 'The children form  | ed a choir for the show.'                                     |
| b. Focus (1   | F) A GYErek-ek <u>Énekkar-t</u> alakít-ott - ak a műső              | or - ra.  |
|               | the child - ren choir - Acc form-Past-3Pl the show-on               |   |
|               | 'It was a choir that the childre                                    | en formed for the show.'                                      |
| c. Quantif    |   |   |
|               | A GYErek-ek énekkar-t is Alakít-ott-ak a Műse                       | or-ra. BLICATIONS   |
|               | the child - ren choir - Acc also form-Past-3Pl the show             |   |
|               | 'The children formed al   | so a choir for the show.'                                     |
| d. Contras    | tive <i>↑Énekkar-t</i> ↓ Alakít-hat-tok a MŰsor-ra!                 |   |
| topic (K      |   |   |
| 1 `           | 'As for a choir, you are allowed                                    | to form it for the show.'                                     |
|               |   |   |
| Example 8.    | Consequence of the neutralization of the Referentiality I           | Requirement.  |
| a. A GYErek-e | ek [Ø/egy] Énekkar-t alakít-ott-ak a MŰsor-ra.                      | $\langle + ref \rangle$ , $\langle - spec \rangle$            |
| the child -   | - ren $[ \emptyset / a ]$ choir - Acc form - Past-3Pl the show-onto |   |
| A GYErek-e    | ek Alakít-ott-ak <i>egy Énekkar-t</i> a MŰsor-ra.                   | $\langle +ref \rangle, \langle -spec \rangle$                 |
| the child -   | ren form - Past - 3Pl a choir - Acc the show-onto                   |   |
|               | 'The childre  | n formed a choir for the show.'                               |
| b. *A GYErek  | -ek FESt-ett-ék ZÖLD-re a KErítés-t.                                | $\langle +ref \rangle, \langle -ref \rangle$                  |
|               | - ren paint-Past-3Pl green-onto the fence-Acc                       |   |
| A GYErek-     | ek ZÖLD-re fest-ett-ék a KErítés-t.                                 | $\langle +\text{ref} \rangle$ , $\langle -\text{ref} \rangle$ |
| the child     | - ren green-onto paint-Past-3Pl the fence-Acc 'The children pa      | inted the fence green.'                                       |

**Example 9.** The neutralization of negative referentiality requirements in the pre-verbal zone of Hungarian sentences (due to another argument's coming into F or K).

| a. A <u>NAGY</u> szoba-ban van <i>a kakas</i> .  |                               | ct. (3c)                  |
|--|-------------------------------|---------------------------|
| the sitting-room - in is the cock                | 'It is in the sitting-room    | that <i>the cock</i> is.' |
| b. TEGnap érkez - ett a mexikói a konfer         | enciá-ra.                     | cf. (4a)                  |
| yesterday arrive-Past3Sg the Mexican the confere | ence-onto                     |                           |
| 'It was yes                                      | sterday that the Mexican arri | ived at the conf.'        |
| c. A GYErek-ek alakít-ott-ák az énekkar-t        | a műsor-ra.                   | cf. (5a)                  |

the child - ren form - Past-3Pl the choir-Acc the show-onto

'The choir was formed for the show BY THE CHILDREN.'

Table 2 below serves as an illustration of the overall hypothesis concerning the distribution of  $\pm$ -referentiality restrictions. In a prototypical neutral sentence (type A. below), the sentence-initial topic zone and the postverbal complement zone are devoted to the task of *anchoring* referents, which requires  $\pm e^+$  arguments, and the tensed verb is the *assertive* center, which contains the new piece of information about the anchored referents. Patients of being, however, straightforwardly belong to the assertive center (see type B). The position immediately preceding the verb stem also belongs to the assertive center, and hence provides "shelter" for genuinely non-referential arguments (type C). What is common in types D and E, is that some assertive operator appears in the sentence (e.g. a focus), which draws the assertive center to itself from other zones of the sentence. As a consequence, positive referentiality-degree requirements are neutralized in the new assertive zone (D), whilst negative ones are neutralized elsewhere in the sentence structure (E).

| A GYErek-ek <sub>+ref</sub>  | MEG-alakít-ott-ák                           | uments and an assertive ver                               |                              |
|--|---|---|------------------------------|
| the child-ren  | Perf-form-Past-Pl3                          | az Énekkar-t <sub>+ref</sub> a MŰsor-ra <sub>+ref</sub> . |                              |
|  |   | the choir-Acc the show-onto                               |                              |
| B. Neutral sentence with a   |   | , hence belonging to the asse                             | rtive zone (5a):             |
| A GYErek-ek+ref  | Alakít-ottak egy                            | Énekkar-t+ref-spec  | a MŰsor-ra <sub>+ref</sub> . |
| the child-ren  | form-Past-3Pl                               | a choir - Acc   | the show-onto                |
| C. Neutral sentence with   | an argument expressing be                   | ing in the preverbal modifie                              | er position, hence           |
| belonging to the assertive   | zone (8a):                                  |   |                              |
| A GYErek-ek+ref  | Ø/egy Énekkart                              | -spec alakít-ott-ak                                       | a MŰsor-ra+ref.              |
| the child-ren  | Ø/a choir-Acc                               | form-Past-3Pl   | the show-onto                |
| D. Focused sentence I: the assertive zone is occupied by a Non-Specificity Effect argument, due to its |   |   |                              |
| focus status (while the ve   | rb gets out of the assertive                | zone also due to the focus of                             | construction) (7b):          |
| A GYErek-ek+ref  | <u>É</u> nekkar-t <sub>+ref</sub>           | alakít-ott-ak a műsor-ra.                                 |                              |
| the child-ren  | choir-Acc                                   | form-Past-3Pl the show-onto                               |                              |
| E. Focused s. II: the asser  | tive zone is occupied by a                  | focused constituent, while t                              | he verb and a Specificity-   |
| Effect argument get out o  | f the assertive zone (due to                | the focus construction) (90                               | l):                          |
| A GYErek-ek  | alakít-ott-ák az énekkar-t.spec a műsor-ra. |   |                              |
| the child-ren  | form-Past-3Pl the choir-Acc the show-onto   |   |                              |

 Table 2. Anchoring (grey) and assertive pieces of information in different sentence types.

# **3** Generating Sentences, Accepting (Intoned) Word Sequences

The crucial part of *information extraction* is a procedure whose input is a sentence, and whose output is an information structure, which is practically a set of possible

operator scope orders (*acceptance*). A similar procedure forms the first half of *machine translation*, too: we need the information structure of the source-language sentence. Then an opposite procedure should come (*generation*), whose input is an information structure, and whose output is an intoned word sequence, that is, a sentence in the target language. Now let us consider this latter procedure, because the former procedure will be based upon it.

As our approach to grammar is "totally lexicalist" [11], similar to our earlier attempts [15], the lexical description of a verb is responsible for the order and intonation of words in the *generated* sentence. What is demonstrated in (10a) below is the requirement, registered in the *core lexicon*, that the subject of *alakít* 'form' should be the (stressed) topic in the initial part of the sentence ( $(\langle 1,T \rangle)$ ), the object should occupy the (stressed) verbal modifier position ( $(\langle 2,M \rangle)$ ) (with an unstressed verb stem following it), and the *-rA* expression should remain in an (also stressed) post-verbal argument position. The numbers provide a scope order, which is still, in a neutral sentence, practically irrelevant. Let us call this lexical rule the *generator*. In (10b), the default generator in the core lexicon requires the *-rA* argument to occupy the verbal modifier position. In (10c) the Hungarian neutral word order requires a generator that ensures that the prefix will occupy the position next to the verb stem and the (non-agentive) subject will remain in a post-verbal A position.

**Example 10.** A few Hungarian lexical items with default scope order in the core lexicon.

 $\langle \langle 1, M \rangle, \langle 2, A \rangle, \langle 3, A \rangle \rangle$ 

MEGérkezett a MExikói a KONferenciára. 'The Mexican arrived at the conference.' (4b)

Two further types of generators produce verb variants, located in an extended lexicon. What is shown in the first row of (11a) is that an *extending* generator inserts certain adjuncts in the argument structure as "fake arguments". Then an *inducing* generator is shown in the second row, responsible for transporting the two fake arguments into the sentence-initial topic zone, the subject into a quantifier position, and the object into the focus. This generator is responsible for the distribution of appropriate stresses. (11b-c) also show an extending generator and two inducing ones. Formula ' $\langle 1, K \rangle$ ' refers to a contrastive topic position in the initial part of the preverbal Hungarian operator zone.

Example 11. Lexical rules extending argument structures and inducing non-neutral ones.

a. FORM(Arg $_{\oslash}$ , Arg<sub>-t</sub>, Arg<sub>-ra</sub>) $\land$ (Arg<sub>Time</sub>, Arg<sub>Place</sub>) (extending generator)  $\langle\langle 3,Q \rangle, \langle 4,F \rangle, \langle 5,A \rangle\rangle, \langle \langle 1,T \rangle, \langle 2,T \rangle\rangle$  (inducing generator) TEGnap a KLUB-ban a GYErek-ek is <u>É</u>nekkar-t alakít-ott-ak a műsor-ra.

yesterday the club - in the child - ren also choir - Acc form - Past-3Pl the show-onto 'It was a choir that yesterday in the club the children, too, formed for the show.'

b. FORM( $\operatorname{Arg}_{\varnothing}, \operatorname{Arg}_{-t}, \operatorname{Arg}_{-ra}$ ) $\land \langle \operatorname{Arg}_{Place} \rangle$ 

 $\langle \langle 2,F \rangle, \langle 3,M \rangle, \langle 4,A \rangle \rangle, \langle \langle 1,K \rangle \rangle$ 

A ↑ĸ⊔uB-ban↓ a <u>GYE</u>rek-ek alakít-ott-ak énekkar-t a műsor-ra.

the club - in the child - ren form-Past-3Pl choir - Acc the show-onto

'As for the club, a choir was formed there for the show BY THE CHILDREN.'

c. PAINT(Arg<sub>Ø</sub>, Arg<sub>-t</sub>, Arg<sub>-ra</sub>)

 $\langle \langle 1,T \rangle, \langle 2,Q \rangle, \langle 3,F \rangle \rangle$  'The children painted each fence GREEN.' A GYErek-ek MINdegyik KErítés-t <u>ZÖLD</u>-re fest-ett-ék. the child - ren every fence - Acc green-onto paint-Past-3Pl

What a generator produces is generally a *set* of sentences, typically ones with different word orders, arranged in a preference order. In Hungarian, it is the *quantifier* that is responsible for this phenomenon, because a quantifier can choose between occupying its preverbal operator position according to the scope order ( $\sigma_1$ ,  $\sigma_{10}$ ,  $\sigma_{11}$ ,  $\sigma_{20}$ ,  $\sigma_{30}$ ) or remaining in the post-verbal zone ( $\sigma_2$ ,  $\sigma_3$ ,  $\sigma_{20}$ ,  $\sigma_{30}$ ,  $\sigma_{40}$ ,  $\sigma_{50}$ ).

**Example 12.** Generating intoned sentences:  $v \rightarrow \langle \sigma_1, \sigma_2 \dots \sigma_k \rangle$ .

- a. PAINT(Arg<sub>Ø</sub>, Arg<sub>-t</sub>, Arg<sub>-ra</sub>)
  - $\nu: \quad \langle \langle 3, Q \rangle, \langle 1, F \rangle, \langle 2, A \rangle \rangle \to \langle \sigma_1, \sigma_2, \sigma_3 \rangle$
  - $\sigma_1: \quad \text{MINdegyik KErités-t ZÖLD-re fest-ett-ék a GYErek-ek}.$ 
    - each fence-Acc green-onto paint-Past-3Pl the child-ren
  - σ<sub>2</sub>: Zöldre festették a GyErekek MINdegyik KErítést.
  - σ<sub>3</sub>: Zöldre festették MINdegyik KErítést a GYErekek.

'Each fence has been painted green by the children.'

b. PAINT(Arg<sub>Ø</sub>, Arg<sub>-ra</sub>)

- $\nu: \quad \langle \langle 1, Q \rangle, \langle 2, Q \rangle, \langle 3, M \rangle \rangle \rightarrow \langle \sigma_{10}, \sigma_{20}, \sigma_{30}, \sigma_{40}, \sigma_{50} \rangle$
- $\sigma_{10}$ : A GYErekek is ('also') MINdegyik KErítést ZÖLDre festették.
- σ<sub>20</sub>: A GYErekek is ZÖLDre festették MINdegyik KErítést.
- σ<sub>30</sub>: MINdegyik KErítést ZÖLDre festették a GYErekek is.
- $\sigma_{40}$ : ZÖLDre festették a GYErekek is MINdegyik KErítést.
- σ<sub>50</sub>: ZÖLDre festették MINdegyik KErítést a GYErekek is.

'The children, too, have painted each fence green.'

c.  $PAINT(Arg_{\emptyset}, Arg_{t}, Arg_{ra})$  'Each fence has been painted green also by the children.'

 $\nu : \quad \langle \langle 2, Q \rangle, \langle 1, Q \rangle, \langle 3, M \rangle \rangle \rightarrow \langle \sigma_{11}, \sigma_{30}, \sigma_{20}, \sigma_{50}, \sigma_{40} \rangle$ 

 $\sigma_{11}$ : MINdegyik KErítést a GYErekek is ZÖLDre festették.

The generated set can also be empty  $(\lambda = \langle \rangle)$ . The problem in (13a) below is that an adjectival (hence,  $\langle -\text{ref} \rangle$ ) argument cannot accept the argument position suggested by the inducing generator (5b). In (13b) the inducing generator doubly violates the Referentiality Requirement (5), which is not neutralized in a (non-contrastive) topic position (7). (13c) illustrates the violation of the possible operator order (in Hungarian), which is as follows (7): {T, K}\*  $\wedge$  {Q, F}\*  $\wedge$  (M)A\*.

**Example 13.** Generating an empty set of intoned sentences:  $v \to \langle \sigma_1, \sigma_2 \dots \sigma_k \rangle = \lambda$ . PAINT(Arg<sub>Ø</sub>, Arg<sub>-t</sub>, Arg<sub>-ra</sub>) the child - ren paint-Past-3Pl green-onto the fence-Acc a.  $v_1$ :  $\langle \langle 1, T \rangle, \langle 3, A \rangle, \langle 2, A \rangle \rangle \to \lambda$ : \*A GYErek-ek FESt-ett-ék *ZÖLD-re* a KErítés-t. b.  $v_2$ :  $\langle \langle 1, T \rangle, \langle 3, A \rangle, \langle 2, M \rangle \rangle \to \lambda$ : \*GYErek-ek ZÖLD-re fest-ett-ek KErítés-t. child - ren green-onto paint-Past-3Pl fence-Acc c.  $v_3$ :  $\langle \langle 1, Q \rangle, \langle 2, T \rangle, \langle 3, M \rangle \rangle \to \lambda$ : \*A GYErek-ek is a KErítés-t ZÖLD - re fest-ett-ék.

the child - ren also the fence-Acc green-onto paint-Past-3Pl

*Acceptance* of an intoned word sequence is the reverse of generation, and can be based upon the latter: we should collect the potential core-lexical verbs, and then collect potential generators, and finally test all combinations in some effective way: see (14) below. If the input contains no reference to intonation, the first step of acceptance is to furnish the word sequence with all potential stress patterns – in some effective way, of course. (15) will serve with some illustration.

**Example 14.** Accepting scope orders (input: intoned sentence):  $\sigma \rightarrow \langle v_1, v_2 \dots v_k \rangle$ . a.  $\sigma_{10}$ : A GYErekek is MINdegyik KErítést ZÖLDre festették.  $\sigma_{10} \rightarrow \langle v_1 \rangle$ ; see (12) above

- PAINT(Arg<sub>Ø</sub>, Arg<sub>-t</sub>, Arg<sub>-ra</sub>)
- $v_1$ :  $\langle \langle 1, Q \rangle, \langle 2, Q \rangle, \langle 3, M \rangle \rangle$
- b.  $\sigma_{11}$ : MINdegyik KErítést a GYErekek is ZÖLDre festették.  $\sigma_{11} \rightarrow \langle \nu_2 \rangle$  $\nu_2$ :  $\langle \langle 2, Q \rangle, \langle 1, Q \rangle, \langle 3, M \rangle \rangle$

c.  $\sigma_{20}$ : A GYErekek is ZÖLDre festették MINdegyik KErítést.  $\sigma_{20} \rightarrow \langle v_1, v_2 \rangle$ 

d.  $\sigma_{30}$ : MINdegyik KErítést ZÖLDre festették a GYErekek is.  $\sigma_{30} \rightarrow \langle v_2, v_1 \rangle$ 

e.  $\sigma$ ': \*GYErekek <u>ZÖLD</u>re MINdegyik KErítést festették.  $\sigma$ '  $\rightarrow \emptyset$ 

**Example 15.** Accepting scope orders (input: a word order with no intonation):  $\kappa \rightarrow \langle \sigma_1, \sigma_2 \dots \sigma_n \rangle$ , where  $\sigma_1 \rightarrow \langle v_{1,1}, v_{1,2} \dots v_{1,k_1} \rangle \dots \sigma_n \rightarrow \langle v_{n,1}, v_{n,2} \dots v_{n,k_n} \rangle$ .

 $\kappa$ : a gyerek-ek is zöld-re fest-ett-ék mindegyik kerítés-t

the child - ren also green-onto paint-Past-3Pl each fence-Acc  $\kappa \rightarrow \langle \sigma_{20}, \sigma_{22}, \sigma_{23}, \ldots \rangle$ ; PAINT(Arg<sub> $\varnothing$ </sub>, Arg<sub>-t</sub>, Arg<sub>-t</sub>)

a.  $\sigma_{20}$ : A GYErekek is ZÖLDre festették MINdegyik KErítést.

 $\text{`The children, too, painted each fence green.'} \\ \sigma_{20} \rightarrow \langle v_1, v_2 \rangle; v_1: \langle \langle Q, 1 \rangle, \langle Q, 2 \rangle, \langle M, 3 \rangle \rangle; v_2: \langle \langle Q, 2 \rangle, \langle Q, 1 \rangle, \langle M, 3 \rangle \rangle$ 

b.  $\sigma_{22}$ : A GYErekek is <u>ZÖLD</u>re festették MINdegyik KErítést.

'It is green that the children, too, painted each fence.'  $\sigma_{22} \rightarrow \langle v_{12}, v_{22} \rangle; v_{12}: \langle \langle Q, 1 \rangle, \langle Q, 2 \rangle, \langle F, 3 \rangle \rangle; v_{22}: \langle \langle Q, 2 \rangle, \langle Q, 1 \rangle, \langle F, 3 \rangle \rangle;$ ' $v_{23}: \langle \langle Q, 1 \rangle, \langle Q, 3 \rangle, \langle F, 2 \rangle \rangle$ 

c.  $\sigma_{23}$ : <sup>??</sup> A GYErekek is  $\downarrow$  ZÖLDre festették MINdegyik KErítést.

'What is true for a set containing children and other people, is nothing else but that it is green that they painted each fence.'

$$\begin{split} \sigma_{23} & \rightarrow \langle \nu_{13}, \nu_{23} \rangle; \ ^{?}\nu_{13} : \langle \langle K, 1 \rangle, \langle Q, 2 \rangle, \langle F, 3 \rangle \rangle; \ ^{??}\nu_{23} : \langle \langle K, 1 \rangle, \langle Q, 3 \rangle, \langle F, 2 \rangle \rangle \\ d. \ \sigma_{44} : *A \ \underline{GYE} \text{rekek is ZOLDre festették MINdegyik KErítést.} \ \sigma_{44} \rightarrow \emptyset \end{split}$$

The last example in this section concerns translation between Hungarian and English. As English word order is very strict, what corresponds to an inducing generator in Hungarian is not simply the same generator but its combination with what we also regard as a lexical generator: one producing new argument structure versions by passivization or dative shift. This approach is also supported by Croft's ([14], Section 8) typological analyses: the best way of understanding the numerous intermediate forms across languages of the world between the active version and the English-type standard passive form ('voice continuum') is in terms of demands for topicalization pertaining to different argument positions.

Example 16. Hungarian operators ~ English argument structure versions.

a. GIVE(Arg<sub>Ø</sub>, Arg<sub>-nAk</sub>, Arg<sub>-t</sub>) GIVE(Arg<sub>Ø</sub>, Arg<sub>Obj1</sub>, Arg<sub>Obj2</sub>)  $\nu_2$ :  $\langle \langle 1,T \rangle, \langle 3,A \rangle, \langle 2,F \rangle \rangle$  $v_2$ :  $\langle \langle 1,T \rangle, \langle 3,A \rangle, \langle 2,F \rangle \rangle$ Péter egy KÖNYv-et ad-ott Mari-nak. Peter gave Mary a BOOK. Peter a book - Acc give-Past3sg Mary-Dat  $GIVE(Arg_{\emptyset}, Arg_{Obj_1}, Arg_{Obj_2})$ b. GIVE(Arg<sub>Ø</sub>, Arg<sub>-nAk</sub>, Arg<sub>-t</sub>)  $\nu_3$ :  $\langle \langle 1,T \rangle, \langle 2,F \rangle, \langle 3,A \rangle \rangle$  $\langle Arg_{\varnothing}, Arg_{to}, Arg_{Obi} \rangle + v_3$ Péter MAri-nak ad-ott egy könyv-et. Peter gave a book to MAry. Peter Mary-Dat give-Past3sg a book - Acc c. GIVE(Arg<sub>Ø</sub>, Arg<sub>-nAk</sub>, Arg<sub>-t</sub>) GIVE(Arg<sub>Ø</sub>, Arg<sub>Obj1</sub>, Arg<sub>Obj2</sub>)  $v_4$ :  $\langle \langle 2, F \rangle, \langle 1, T \rangle, \langle 3, A \rangle \rangle$  $\langle Arg_{by}, Arg_{\varnothing}, Arg_{Obj} \rangle + v_4$ MAri-nak Péter ad-ott egy könyv-et. Mary was given a book by PEter. Mary-Dat Peter give-Past3sg a book-Acc d. GIVE(Arg<sub>Ø</sub>, Arg<sub>-nAk</sub>, Arg<sub>-t</sub>)  $GIVE(Arg_{\emptyset}, Arg_{Obj_1}, Arg_{Obj_2})$  $\langle Arg_{by}, Arg_{to}, Arg_{\varnothing} \rangle + v_5$  $\nu_5$ :  $\langle \langle 2, F \rangle, \langle 3, A \rangle, \langle 1, T \rangle \rangle$ The book was given to Mary by PEter. A KÖNYv-et PÉter ad - t - a Mari-nak.

We note at this point, following our anonymous reviewer's advice: what (16) illustrates is a simplified map of the relevant phenomena; as 'intonation' includes not only stress/prominence but also, as a separate set of choices, tone (e.g. falling/rising). We say that an 'intoned' word sequence is generated, but it would be more accurate to say that it is a 'word sequence with prominence'. An exception is the Hungarian contrastive topic with its rising-falling tone, which we consider; but the special tonal pattern of yes/no questions and ironic performances has not been considered yet.

In (spoken) English prominence alone can signal information structure, in a way that is not possible in Hungarian, e.g. (16c) could be expressed in English as <u>**PE**</u>ter gave Mary a book. We assume in our simplified approach that tonic focus is always on the last lexical word in the sentence, being aware of the fact that it is only the unmarked case. Considering the marked cases is a task postponed to future research.

#### 4 Implementation

The book - Acc Peter give-Past-3sg Mary-Dat

There are only a few works pertaining to deep parsing of scope order and referentiality in connection with word order and intonation. An (excellent) example is shown by Traat and Bos [16], but we are the first to attempt to build a similar system for Hungarian (whose relevant and advantageous properties are discussed in Sec. 2).

First of all we analyze sentences phonologically and morphologically. Our approach is "totally lexicalist" – grammars based on "total lexicalism" need *not* build phrase structures [11], [12], [13], [14]. In our system, word order is handled by *rank parameters*, instead. In general, we use whole numbers from 1 to 7 for ranks, 1 being

the strongest; by default, it means direct adjacency. Our lexicon contains *morphemes* instead of words and they can search for any other morpheme inside or outside the word which they are part of. In our approach, the only difference between morphology and syntax is that in the syntactic subsystem, the program searches for morphemes being in a certain grammatical relation in two *different* words. But because of this, morphological and syntactic rank parameters are handled separately.

The lexicon is extendable in multiple ways. Apart from adding more words and morphemes, new features can be added to the data structure, thus improving the precision of the analysis. Third, the *core lexicon* can be extended via generating new lexical elements by rules of generation, thus forming the *extended lexicon*. The core lexicon contains all the basic properties of a morpheme including its default behavior (e.g. argument structure of a verb). This applies primarily for intonation: the core lexicon contains the property of "stress" (some words cannot be stressed at all while others are stressed in a neutral sentence) – but it can be overwritten if the sentence is not neutral. For instance, the entity with the property value "focus-stressed" is generated automatically and inserted into the extended lexicon.

Obviously, it is not very effective to try out all of the possible intonation schemes on a long sentence. However, there are two things which must be taken into account: apart from morphemes (such as articles) which can never be stressed or focusstressed, intonation of arguments before and after the verb is constrained.

By default, the normal Hungarian argument position is *after* the verb. There are, however, arguments which prefer being in the verbal modifier or the topic position. These preferences should be stored in the core lexicon along with the arguments, as was illustrated in (10) in Section 3. If the sentence does not fit into the preferred schema, it is best to use heuristics to create the appropriate generator (11-12).

Prolog (a language for logical programming), including Visual Prolog 7 (which is probably the most elaborate version of Prolog and which we use), usually allows writing predicates which can be invoked "backwards" (generation (12) vs. acceptance (14)). We suppose that the evaluation of predicates can be reversed. Based on this principle, the future machine translator can be symmetrical. By using the keyword anyflow, all arguments of the predicate can be used for both input and output, allowing even entire programs to be executed "backwards". The keywords procedure, determ, multi, nondeterm etc. describe whether the predicate can fail (a procedure must always succeed), have multiple backtrack points (nondeterm) or not (determ).

Let us turn to the particular phases of parsing (acceptance).

**Phase 0.** Before taking intonation etc. into account, all words are segmented and analyzed phonologically and morphologically, allowing the class of each word to be determined. Practically, the last class-changing morpheme (derivative affix) has the relevant "class" output feature. The input and output classes are stored in the core lexicon of every class-changing morpheme.

**Phase 1.** During the actual syntactic analysis, arguments of a verb are searched with rank 7 (weak). This is a bi-directional search: all non-predicative nominal expressions look for their predicate, too, and if appropriate, the result is stored in the memory of the computer. The adverbial adjuncts search for the verb with rank 7, too, but if it is found, the extending generator must be invoked to modify the default argument structure of the verb and insert the form into the extended lexicon (11a-b).

**Phase 2.** The inducing generator tries to determine the discourse function -T, Q, F, M, A, see (7) – of the arguments of a verb by creating all possible patterns and trying to apply them onto the sentence. If intonation is present in the input, it can be considered here, making the analysis faster. But shortly, if the Hungarian operator order is violated, the analysis is apparently wrong and the backtracking mechanism of Prolog should be activated before we reach the end of the sentence (13c).

**Phase 3.** As for the referentiality degrees (6), two of them (non-referential and definite) can be handled during Phase 1. First of all we suppose that all nouns are non-referential by default. Since the article searches the noun, it changes the features 'positive ref. degree' and 'negative ref. degree' of the noun and the new form should be inserted into the extended lexicon.

The only remaining problem is (non-)specificity. Let us take (3). If a non-definite determinant is present, this can be determined only (partly) during or after Phase 2. If no other constraints exist and a verbal modifier is present in a neutral sentence, the argument is specific. So if we generate two instances of *egy mexikói* 'a Mexican', one with features 'spec' and 'non-def' and one with 'non-spec' and 'ref', the analysis of (3b) must fail with the latter. Of course, the generated instances must be inserted into the extended lexicon. If the verb has a modifier, its argument structure *differs* from that of the *same* verb without a modifier because of the specificity criteria. An easier but slower way is to insert two *egy*'s 'a(n)' into the core lexicon, one being specific and one not. If the result of Phase 1 or 2 is wrong, much of the analysis, including parts of Phase 0 will start over because even the most basic unifications (belonging to the same unit in the *core lexicon*) will have to be broken up.

#### 5 Conclusions

Very few computational systems exist which aim at deep parsing of scope order and referentiality in connection with word order and intonation, although this is a crucial step in reliable information extraction and sophisticated machine translation. We strive for building a system based – at first – primarily upon Hungarian, which is famous for expressing scope relations (in its pre-verbal operator zone) in an explicit way [2], [3] and referentiality degrees also in a quite straightforward manner [5], [8]. We also began to extend our approach to other languages (e.g. English), pointing to similarities, differences and correspondences, for instance, between Hungarian word-order changing operations and English argument-structure changing operations.

To achieve this goal in our totally lexicalist approach [11], the implementation requires a double-layered lexicon. This consists of a core lexicon containing the default behaviour of morphemes (e.g. verb stems) and an extended lexicon, whose elements are generated by a generic lexical-rule module. The elements of the core lexicon are responsible for accounting for neutral sentences, whilst those of the extended lexicon are to handle sentences with different topic, quantifier and focus constructions, where referentiality requirements are often different from those in neutral sentences. We thus primarily generate neutral and non-neutral sentences from

the elements of the double-layered lexicon, and, based upon this generation in our Prolog environment, we extract the information structure of +/- intoned sentences.

### Acknowledgements

We are grateful to OTKA60595 for their contribution to all our costs at NLPCS 2010.

#### References

- Kiefer, F., É. Kiss, K. (eds.): The Syntactic Structure of Hungarian. Syntax and Semantics, Vol. 27. Academic Press, New York (1994)
- 2. É. Kiss, K.: Hungarian Syntax. Cambridge University Press, Cambridge (2001)
- Szabolcsi, A.: Strategies for Scope Taking. In: Szabolcsi, A. (ed.): Ways of Scope Taking, SLAP 65. Kluwer, Dordrecht (1997) 109–154
- 4. Alberti, G., Medve, A.: Focus Constructions and the "Scope-Inversion Puzzle" in Hungarian. In: Approaches to Hungarian, Vol.7. Szeged (2000) 93–118
- Szabolcsi, A.: From the Definitness Effect to Lexical Integrity. In: Abraham, W., de Meij, S. (eds.): Topic, Focus, and Configurationality. John Benjamins, Amsterdam (1986) 321– 348
- 6. É. Kiss, K.: Definiteness Effect Revisited. In: Appr. to Hung., Vol. 5. (1995) 63-88
- Kálmán, L.: Definiteness Effect Verbs in Hungarian. In: Appr. to Hung., Vol. 5. (1995) 221–242
   Alberti, G.: Restrictions on the Degree of Referentiality of Arguments in Hungarian Sentences. In: Acta Linguistica Hungarica, Vol. 44/3-4. (1997) 341–362
- 9. de Jong, F., Verkuyl, H.: Generalized Quantifiers: the Properness of Their Strength. In: van Benthem, J., ter Meulen, A. (eds.): GRASS 4. Foris, Dordrecht (1984)
- 10. Enç, M.: The semantics of specificity. In: Linguistic Inquiry, 22. (1991) 1–25
- 11. Alberti, G., Kleiber, J. Viszket, A.: GeLexi project: Sentence Parsing Based on a
- GEnerative LEXIcon. In: Acta Cybernetica, Vol 16. (2004) 587–600
- 12. Hudson, R.: Word Grammar. Blackwell, Oxford (1984)
- 13. Karttunen, L.: Radical Lexicalism. Report No. CSLI 86-68. Stanford (1986)
- Croft, W.: Radical Construction Grammar. Syntactic Theory in Typological Perspective. Oxford University Press, Oxford (2001)
- 15. Alberti, G. Kleiber, J.: The GeLexi MT Project. In: Hutchins, J. (ed.): Proceedings of EAMT 2004 Workshop (Malta), Univ. of Malta, Valletta (2004) 1–10
- Traat, M. Bos, J.: Unificational Combinatory Categorial Grammar. In: Proceedings of the 20<sup>th</sup> International Conference on Computational Linguistics. Geneva (2004)