# Unificational Combinatory Categorial Grammar:
## Combining Information Structure and Discourse Representations

**Maarika Traat**
The University of Edinburgh
2 Buccleuch Place
Edinburgh EH8 9LW,
United Kingdom,
M.Traat@ed.ac.uk

**Johan Bos**
The University of Edinburgh
2 Buccleuch Place
Edinburgh EH8 9LW,
United Kingdom,
jbos@inf.ed.ac.uk

## Abstract

In this paper we present a grammar formalism that combines the insights from Combinatory Categorial Grammar with feature structure unification. We show how information structure can be incorporated with syntactic and semantic representations in a principled way. We focus on the way theme, rheme, and focus are integrated in the compositional semantics, using Discourse Representation Theory as first-order semantic theory. UCCG can be used for parsing and generating prosodically annotated text, and therefore has the potential to advance spoken dialogue systems.

## 1 Introduction

The integration of information structure (the way information is "packaged" in a sentence) in practical formalisms in computational linguistics has long been ignored. There are two main reasons for this: (1) formalisations of information structure often use variants of higher-order logic to characterise its semantic impact (Krifka, 1993; Kruijff-Korbayova, 1998; Steedman, 2000), which limits the use of inference in practice (Blackburn and Bos, 2003); and (2) the effect of information structure on the compositional semantics of an utterance is rarely worked out in enough detail useful for computational implementation. On the other hand, exploring information structure in spoken dialogue systems is becoming realistic now because of the recent advances made in text-to-speech synthesisers and automated speech recognisers — hence there is a growing need for computational implementations of information structure in grammar formalisms.

In this paper we present Unificational Combinatory Categorial Grammar (UCCG), which integrates aspects of Combinatory Categorial Grammar (Steedman, 2000), Unification Categorial Grammar (Zeevat, 1988; Calder et al., 1988), and Discourse Representation Theory (Kamp and Reyle, 1993). It offers a compositional analysis of information structure, a semantics compatible with first-order logic, and a computational implementation for a fragment of English, using unification in combining grammatical categories. As we will show, this makes UCCG easy to implement, and allows us to integrate prosodic information in the semantics in a transparent and systematic way. Although based on first-order logic, we claim that UCCG has enough expressive power to model information structure such that it has the potential to improve speech generation with context appropriate intonation in spoken dialogue systems.

## 2 Background

Categorial Grammars (CG) (Wood, 2000) are lexicalised theories of grammar. The notion of "category" refers to the functional type that is associated with each entry in the lexicon which determines the ability of a lexical item to combine with other lexical items. CGs also have a set of rules defining the syntactico-semantic operations that can be performed on the categories.

Combinatory Categorial Grammar (CCG) is a generalisation of CG (Steedman, 2000). While the pure CG only involved functional application rules for combining categories, CCG introduces several additional combinatory rules for both syntactic and semantic composition — forward and backward composition, and crossed composition, as well as substitution rules. As a result, CCG covers a wide range of linguistic phenomena, including various kinds of coordination. For building semantic representation CCG uses the lambda calculus, although unification has been proposed as well (Steedman, 1990). Moreover, CCG has a built-in theory of intonation and information structure (Steedman, 2000), that we will use as the basis for our computational treatment of theme, rheme and focus.

Unification Categorial Grammar (UCG) uses Head-Driven Phrase Structure Grammar type of feature structures, called *signs*, to represent the categories of lexical items (Zeevat, 1988; Calder et al., 1988). The directionality of the attributes of a functor category is marked by the features *pre* and *post* on its attributes rather than by the directionality of the slashes as it is done in CCG. In contrast to CCG, UCG only uses forward and backward application as means for combining categories. The use of signs makes it straightforward to define the syntax-semantic interface.

The formalism that we introduce in this paper, UCCG, aims to marry the best parts of CCG and UCG. Following UCG, we use signs to represent the linguistic data, and both semantics and syntax are built up simultaneously via unification. From CCG we inherit the directional slash notation, the additional combinatory rules, and the analysis of intonation. UCCG employs DRT (Kamp and Reyle, 1993) with neo-davidsonian style event semantics as semantic formalism, but extends the basic DRS language to allow integration of prosodic information in syntactic and semantic analysis.

## 3  Unificational CCG

### 3.1  Signs

UCCG makes use of feature structures called *signs* in its linguistic description. There are two types of signs: basic and complex signs. A basic sign is a list of attributes or features describing the syntactic and semantic characteristics of a lexical expression, in the spirit of UCG. We deviate from UCG in the way we define complex signs, which is done recursively:

> - If X and Y are signs then X/Y is a complex sign.
> - If X and Y are signs X\Y is a complex sign.
> - All basic and complex signs are signs.

A basic sign can have a varied number of features, depending on the syntactic category of the lexical expression the sign is characterising. There are three obligatory features any sign must have, namely PHO, CAT and DRS. PHO stands for the phonological form, CAT for the syntactic category of the lexical expression, and DRS for its semantic representation. Besides the above three a sign can also have the following features:[1]
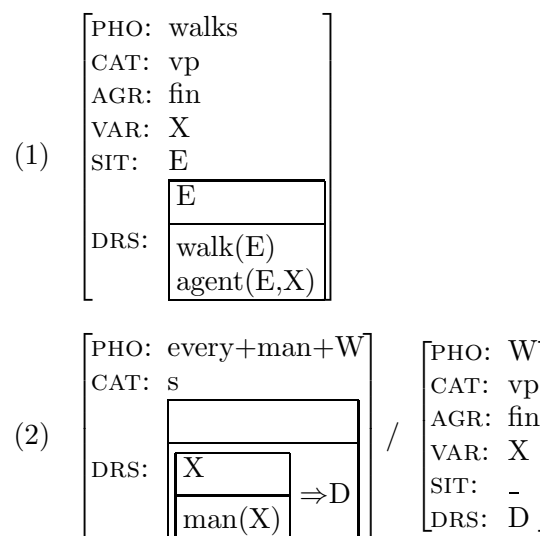
- AGR to mark the inflectional characteristics of categories;
- VAR for discourse referents ranging over individuals;
- SIT for discourse referents ranging over eventualities (events or states).

In our notation inside the feature structures we use the following convention: constants start with a lower case letter, and variables start with an upper case letter. The feature names are written using small capitals. To make the feature structures more easily readable we narrow the choice of possible variable names for each type of variables:

- (PHO) variables: W, W1, W2, etc.
- (AGR) variables: A, A1, A2, etc.
- (DRS) variables: D, D1, D2, etc.
- (SIT) variables: E, E1, E2, etc.
- Discourse referents (VAR) use any other capital letter with the preference for the characters towards the end of the alphabet.

There are three kinds of basic signs in UCCG, corresponding to the basic categories — those with CAT feature sentence (s), those with CAT feature noun (n), and those with CAT feature verb phrase (vp). A basic sign for verb phrases is shown in (1), and a complex sign for noun phrases is shown in (2).

(1)
$$\begin{bmatrix} \text{PHO:} & \text{walks} \\ \text{CAT:} & \text{vp} \\ \text{AGR:} & \text{fin} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{DRS:} & \begin{array}{|l|} \hline \text{E} \\ \hline \text{walk(E)} \\ \text{agent(E,X)} \\ \hline \end{array} \end{bmatrix}$$

(2)
$$\begin{bmatrix} \text{PHO:} & \text{every+man+W} \\ \text{CAT:} & \text{s} \\ \text{DRS:} & \boxed{\boxed{\begin{array}{|l|} \hline \text{X} \\ \hline \text{man(X)} \\ \hline \end{array}} \Rightarrow \text{D}} \end{bmatrix} \Big/ \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{vp} \\ \text{AGR:} & \text{fin} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \_ \\ \text{DRS:} & \text{D} \end{bmatrix}$$

---

[1]Depending of the needs of a specific application and language for which a UCCG grammar is constructed many more features could be introduced in basic signs.

The above examples illustrate the role of unification by creating a link between syntax and semantics. UCCG explores the fact that the same variables can be used at several different levels. For example, the variables standing for discourse referents serve as a link between syntax and semantics — the variable in the VAR feature in the feature structure fits into its corresponding slot in the DRS in the DRS feature. We use this technique to integrate information structure as well.

## 3.2 Categories

Each sign corresponds to a related CCG category. The category of a basic sign is the value of its CAT feature. The category of a complex sign it is made up of the CAT feature values of all the component parts of the complex sign, separated by the slashes and brackets used in the complex sign, resulting in a complex category. For instance, the the syntactic category of the sign in (1) is *vp*, and in (2) the category is *s/vp*. The three basic categories used in UCCG are thus *s*, *n* and *vp*, while all other categories are formed by combining the above three, using backward and forward slashes.

Note that noun phrase is not among the basic categories. In UCCG We use its 'type-raised' variant *s/vp* (corresponding to the CCG category *s/(s\np)*). This choice is motivated by the need to determine quantifier scope in the semantics of quantified noun phrases. The somewhat unconventional basic category *vp* is a byproduct of the above.

## 3.3 Feature Values

In order to make it easier to refer to parts of complex signs later, we introduce the following terminology:

- X is the **result** of a sign X/Y or X\Y.
- Y is the **argument** of a sign X/Y or X\Y.

The value of the VAR and the SIT features is always a variable, while other features can have a number of constant values. The PHO feature holds the string value of the linguistic expression represented by the given feature structure. Presently, we use the orthographic form of words. In basic signs the PHO feature is filled by lexical items, in complex signs it also contains variables, which get constant values when the complex sign is combined with its argument signs. The PHO feature in result parts of complex signs is of the form:

$$\ldots + W1 + \textit{word} + W2 + \ldots$$

where *word* is a lexical item, and W1 and W2 are variables that get values through unification in the categorial combination process. The item unifying with W1 precedes and the one unifying with W2 follows the lexical item *word*. The exact number and order of the variables the PHO feature contains depends on the category of the given sign.

In the present implementation the AGR feature is only used in connection with verb phrases and can take constant values *fin* (finite) or *non-fin* (non finite).

The DRS feature, if it is not a variable itself, holds a DRS corresponding to the semantics of the lexical item(s) characterised by the given sign. DRSs are constructed in a compositional way using the VAR and SIT features of the sign to take care of predicate argument structure, and the merge operator (;) to construct larger DRSs from smaller ones. Merge-reduction is used to eliminate merge operators introduced in the composition process. This is also the stage where discourse referents are renamed to avoid accidental clashes of variables introduced by unification (Blackburn and Bos, 2003).

## 3.4 The Combinatory Rules

Presently we have introduced the following four CCG combinatory rules in UCCG: forward application, backward application, forward composition, and backward composition. Other CCG combinatory rules could be introduced equally easily should the need arise.

| $X/Y \; Y \Longrightarrow X$ |
|---|
| *Forward application* ——————> |
| $Y \; X\backslash Y \Longrightarrow X$ |
| *Backward application* <—————— |
| $X/Y \; Y/Z \Longrightarrow X/Z$ |
| *Forward composition* ———*Comp>* |
| $Y\backslash Z \; X\backslash Y \Longrightarrow X\backslash Z$ |
| *Backward composition* *<Comp*——— |

The rule boxes above are to be interpreted in the following way: in the first row there is the rule, on the left in the second row there is the name of the rule and on the right the marking for it as used in the derivations. The variables $X$, $Y$ and $Z$ in the rules above stand for (basic or complex) signs.

Some of the combinatory rules can be seen in action on UCCG signs in Figures 1 to 3 below.

# 4 Adding Information Structure

By *information structure* we mean the way information is packaged in a sentence. We use the terms *theme* and *rheme* as introduced by the Prague circle of linguists. Theme is the central question or topic the sentence is about, while rheme is the novel contribution of the sentence.

In many languages, including English, prosody is the main means of indicating the information structure of the sentence. In other languages additional or alternative means may be available, such as word order, and the use of specific lexical items. Example (3) illustrates the connection between information structure and prosody in English.

(3)  Who taught Alexander the Great?
     [ARISTOTLE]$_{rh}$ [taught Alexander the Great.]$_{th}$
     *[Aristotle taught]$_{th}$[ALEXANDER the GREAT.]$_{rh}$

The lexical items in capital letters in (3) carry the main rhematic accent of the sentence. As illustrated by this example, the placement of this accent determines whether the answer given to the question is appropriate or not.
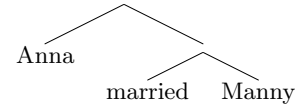
## 4.1 Information Structure in CCG

Steedman introduces information structure as an integral part of the CCG formalism (Steedman, 2000). He argues that there is a specific set of pitch accents in English that can accompany theme, and another set that accompany rheme, the most common theme pitch accent being L+H* and the most common rheme pitch accent being H*.[2] The main pitch accent of the intonational phrase combined with a boundary tone gives us a complete intonational phrase.

There are various boundary tones, the most frequently occurring ones being a low boundary LL% and a rising boundary LH%. There is a tendency for LH% to occur at the end of an intonational phrase containing the theme pitch accent L+H*, and for LL% to occur after the rheme pitch accent H*.
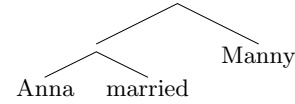
According to the prosodical phrasing, CCG provides different parses for the same string of words, giving rise to different interpretation with respect to information structure:

---

[2]The intonational notation used is due to Pierrehumbert (Pierrehumbert, 1980). According to her intonational phrases are made up of the following components: pitch accent(s), phrasal tone and boundary tone. In Steedman's (Steedman, 2000) representation the last two have been joined together under the name '*boundary tone*'. L stands for low pitch, and H for high pitch.

(4)  *Anna    married   Manny.*
     H* LL%             L+H* LH%

     Anna
          married   Manny

(5)  *Anna   married   Manny.*
     L+H*       LH%   H* LL%

                    Manny
          Anna    married

Parsing according to intonational phrasing in CCG is achieved in the following way: the categories of lexical items can be either theme marked by a theme accent, rheme marked by a rheme accent, or unmarked (i.e., unaccented). Theme and rheme marked categories can freely combine with adjacent categories with the same marking or adjacent categories with no intonational marking. If a theme or rheme marked category combines with an intonationally unmarked category, the result category inherits the themeness or rhemeness from the marked category that participated in the combination process.

While pitch accents are seen as properties of words that carry them, boundary tones are seen as individual lexical entries, and have their own category of the form S\$$_\phi$\S\$$_{\theta/\rho}$, where S\$ is a variable that stands for any category that is a function whose result is $S$ (i.e., sentence), $\phi$ stands for phrase, $\theta$ for theme and $\rho$ for rheme (Steedman, 2000). The effect this category achieves is copying the category to its left it combines with, and replacing its intonational marking by *phrase*. Phrase marked categories can only combine with other phrase marked categories, and hence avoid combination over intonational phrase boundaries. In other words, boundary tones function like "stoppers" of theme and rheme categories, preventing theme and rheme to be further spread along sub-phrases of the sentence.

## 4.2 Information Structure in UCCG

When introducing prosodical and information-structural features to UCCG we follow the theory of CCG, with a few exeptions. As we also aim to derive a computational implementation of UCCG in the form of a parser we need to be concrete about how sign unification in UCCG interacts with CCG's theory of information structure.

Adding intonation to UCCG raises several problems, as combination of signs only via straightforward unification is not possible any more. We have to give prosodical signs the ability to alter prosodical feature values in the result signs they produce when combining with a lexical sign. We will do this using recursive unification—the details of this process will be discussed in Sections 4.3 and 4.4.

Integrating information structure with the UCCG sign representation brought along some additions. Firstly, we introduce two new features in the sign. The first of them is called INF and expresses information structure. It can either be a variable — in the case of unmarked expressions, or it can take the following values $\theta$ (theme), $\rho$ (rheme), or $\phi$ (phrase). The second newly introduced feature is FOC. This feature indicates focus, i.e. whether the particular word carries a pitch accent or not. This feature is only present on lexical signs.

The second change involves introducing information structural labels on DRS conditions (except on those expressing the semantic roles of verbs). The labels are of the form *Cond:Inf Foc*, where *Cond* is a DRS condition, *Inf* stands for the information-structure value ($\theta$, $\rho$, or $\phi$), and *Foc* for the value of the focus (+ or −). The information-structure label in the DRS is tied to the INF feature through the use of the same variable, and gets its constant value from the feature by unification.

### 4.3 Pitch Accents

CCG views pitch accents as properties of words and introduces multiple entries for each lexical item in the lexicon, whether it is theme marked, rheme marked, or unmarked. We do not oppose CCG's view of pitch accents, but we chose a slightly different approach in UCCG: pitch accents get similar treatment as boundary tones — they are independent entries in the lexicon. This way we avoid having to expand the lexicon. For instance, the lexical sign for the proper name *Manny* is shown in (6).

$$(6) \quad \begin{bmatrix} \text{PHO:} & \text{Manny+W} \\ \text{CAT:} & \text{s} \\ \text{INF:} & \text{I} \\ \text{FOC:} & \text{F} \\ \text{DRS:} & \boxed{\begin{array}{l} \text{X} \\ \hline \text{manny(X):I F} \end{array}}\text{;D} \end{bmatrix} \Big/ \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \text{I} \\ \text{FOC:} & \text{F} \\ \text{DRS:} & \text{D} \end{bmatrix}$$

Like all lexical signs, the sign in (6) shows that the values for FOC and INF are still uninstantiated. Once it combines with the sign for a pitch accent, both of these features will get instantiated. For example, (7) shows the result of combining the above lexical sign with a the sign for L+H*:

$$(7) \quad \begin{bmatrix} \text{PHO:} & \text{Manny+W} \\ \text{CAT:} & \text{s} \\ \text{INF:} & \theta \\ \text{DRS:} & \boxed{\begin{array}{l} \text{X} \\ \hline \text{manny(X):}\theta\text{+} \end{array}}\text{;D} \end{bmatrix} \Big/ \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \text{I} \\ \text{DRS:} & \text{D} \end{bmatrix}$$

Note that signs for pitch accents need to be combined first with the signs of the lexical items the accents appear on. Otherwise it would be impossible to tell which item actually carries the accent for larger phrases such as *married Manny H* LL%* , where without the above mentioned constraint we could combine *married* and *Manny* first to form the unit *married Manny*, and only then combine this two word unit with the pitch accent. However, this is not what we want, because this way we cannot determine any more which of the two words was accented. Note also, that the FOC feature only appears in lexical signs.

So what does a sign for pitch accents look like? Borrowing from Steedman's notation, the sign for L+H* has the following format:

$$(8) \quad \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{C} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \theta \\ \text{DRS:} & \text{D} \end{bmatrix}\text{\$} \Big\backslash \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{C} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \theta \\ \text{FOC:} & \text{+} \\ \text{DRS:} & \text{D} \end{bmatrix}\text{\$}$$

The idea behind the sign in (8) is the following: the sign $\mathcal{X}\$$ stands for unification of $\mathcal{X}$ with a (basic or complex) sign. In the case of basic signs, ordinary unification on the level of signs applies, in the case of complex signs, unification of S also applies to sub-signs. Through unification of variables the information structural marking also finds its way to the DRS in the form of labels on the appropriate DRS conditions.

Combining the sign for *'Manny'* (6) with the sign for the theme accent *'L+H*'* (8) results in the unit *'Manny L+H*'*, shown in (7). Notice how through unification also the information
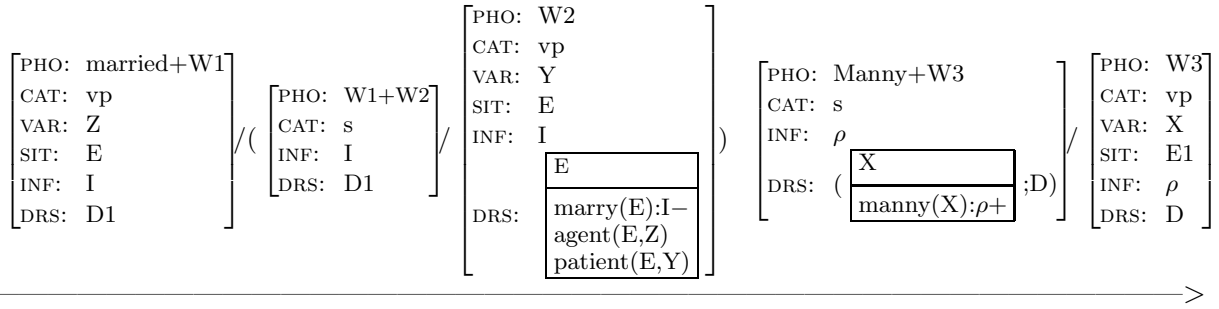
$$\begin{bmatrix} \text{PHO:} & \text{married+W1} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{Z} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \text{I} \\ \text{DRS:} & \text{D1} \end{bmatrix} / (\begin{bmatrix} \text{PHO:} & \text{W1+W2} \\ \text{CAT:} & \text{s} \\ \text{INF:} & \text{I} \\ \text{DRS:} & \text{D1} \end{bmatrix} / \begin{bmatrix} \text{PHO:} & \text{W2} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{Y} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \text{I} \\ \text{DRS:} & \boxed{\begin{array}{l} \text{E} \\ \hline \text{marry(E):I--} \\ \text{agent(E,Z)} \\ \text{patient(E,Y)} \end{array}} \end{bmatrix}) \qquad \begin{bmatrix} \text{PHO:} & \text{Manny+W3} \\ \text{CAT:} & \text{s} \\ \text{INF:} & \rho \\ \text{DRS:} & (\boxed{\begin{array}{l} \text{X} \\ \hline \text{manny(X):}\rho+ \end{array}}\text{;D}) \end{bmatrix} / \begin{bmatrix} \text{PHO:} & \text{W3} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E1} \\ \text{INF:} & \rho \\ \text{DRS:} & \text{D} \end{bmatrix}$$

---$>$

$$\begin{bmatrix} \text{PHO:} & \text{married+Manny} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{Z} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \rho \\ \text{DRS:} & (\boxed{\begin{array}{l} \text{X} \\ \hline \text{manny(X):}\rho+ \end{array}}\ ;\ \boxed{\begin{array}{l} \text{E} \\ \hline \text{marry(E):}\rho- \\ \text{agent(E,Z)} \\ \text{patient(E,X)} \end{array}}) \end{bmatrix}$$

Figure 1: Derivation for *married Manny H\** using Forward Application

structural label of the DRS condition manny(X) gets the value $\theta+$ (theme and focus).

### 4.4 Boundary Tones

In essence the signs for boundary tones are similar to the pitch accent signs, except that they do not contain a FOC feature in the argument part. They take the following form:

$$(9) \quad \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{C} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \phi \\ \text{DRS:} & \text{D} \end{bmatrix} \$ \ \backslash \ \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{C} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \rho \\ \text{DRS:} & \text{D} \end{bmatrix} \$$$

As with pitch accents, when combining signs of boundary tones the argument sign will unify recursively with all sub-signs of the lexical sign, effectively replacing the value of the INF feature by $\phi$ (phrase).

Hence, the constant value $\phi$ for the INF feature only serves the purpose of keeping the full intonational phrase from combining with any other signs than similarly phrase marked signs, and it has no impact on the semantics. There are two signs for each boundary tone: one that deals with boundary tones occurring at the end of a rheme marked intonational phrase (as shown above), and another one that deals with boundary tones after themes.

We have restricted the variable in the argument part of the boundary signs to only be able to combine with themes and rhemes, assum-

ing that in the case of unmarked themes (as here is no pitch accent, there is no theme marking on the sign) we do not encounter boundary tones after the theme part, and therefore we are dealing with genuinely ambiguous information structure. An unmarked theme will in our approach be automatically marked as part of the rheme. For illustration of combining a lexical sign with a boundary tone sign see Figure 2.
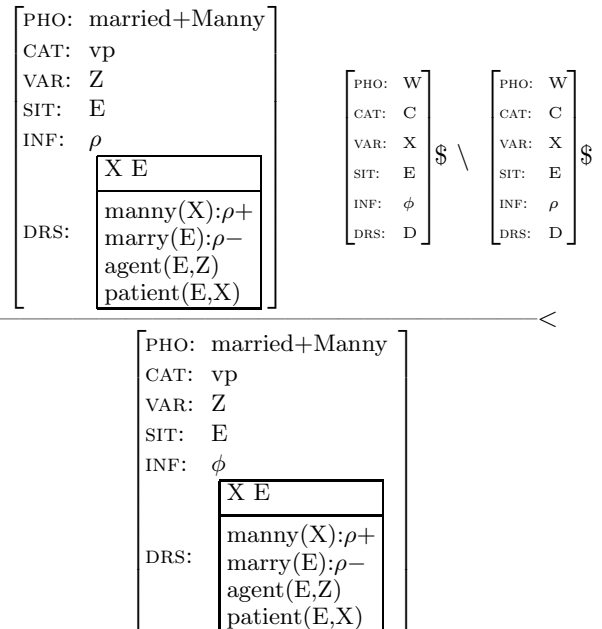
$$\begin{bmatrix} \text{PHO:} & \text{married+Manny} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{Z} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \rho \\ \text{DRS:} & \boxed{\begin{array}{l} \text{X E} \\ \hline \text{manny(X):}\rho+ \\ \text{marry(E):}\rho- \\ \text{agent(E,Z)} \\ \text{patient(E,X)} \end{array}} \end{bmatrix} \quad \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{C} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \phi \\ \text{DRS:} & \text{D} \end{bmatrix} \$ \ \backslash \ \begin{bmatrix} \text{PHO:} & \text{W} \\ \text{CAT:} & \text{C} \\ \text{VAR:} & \text{X} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \rho \\ \text{DRS:} & \text{D} \end{bmatrix} \$$$

---$<$

$$\begin{bmatrix} \text{PHO:} & \text{married+Manny} \\ \text{CAT:} & \text{vp} \\ \text{VAR:} & \text{Z} \\ \text{SIT:} & \text{E} \\ \text{INF:} & \phi \\ \text{DRS:} & \boxed{\begin{array}{l} \text{X E} \\ \hline \text{manny(X):}\rho+ \\ \text{marry(E):}\rho- \\ \text{agent(E,Z)} \\ \text{patient(E,X)} \end{array}} \end{bmatrix}$$

Figure 2: Derivation of *married Manny H\* LL%* using Backward Application

Finally, Figures 1 to 3 show a complete parse of the prosodically marked sentence *'Anna L+H\* LH% married Manny H\* LL%'*. Due to space considerations we omitted the two initial steps that involve combining the sign of *'Anna'* with the sign of the theme accent *'L+H\*'* to form a new theme unit *'Anna L+H\*'*, and then combining this unit with the sign of the boundary tone *'LH%'* to form the full intonational phrase *'Anna L+H\* LH%'*. (These steps are similar to the ones illustrated in Figure 2.) Due to variable unification in the features VAR and SIT, while performing the syntactic combination of the lexical signs, we simultaneously construct the semantic representation in the DRS.
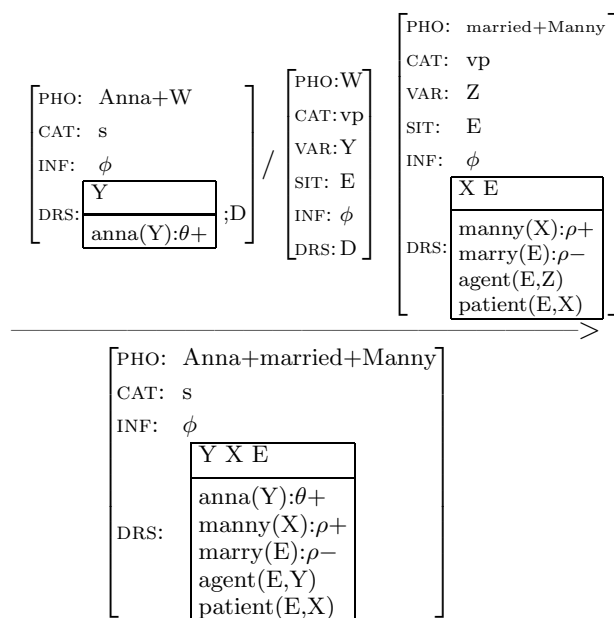


Figure 3: Derivation for *'Anna L+H\* LH% married Manny H\* LL%'*, using Forward Application and Merge-Reduction

## 5 Conclusions and Future work

The present paper described the Unificational Combinatory Categorial Grammar (UCCG) formalism, which was developed bearing in mind its future application in parsing and generating prosodically annotated text. One of the key features of UCCG is the novel use of Discourse Representation Theory combined with a theory of information structure. We believe that UCCG has the potential to advance spoken language dialogue systems, both in natural language analysis and generation. Although current automatic speech recognisers do not output prosodic information, some of the state-of-the-art speech synthesisers handle prosodically annotated input strings.

We have implemented a UCCG parser for a fragment of English that takes prosodically annotated strings as input and generates DRSs enriched with information structure. Future work involves implementing a generation component based on UCCG, evalating the expressive power of UCCG with respect to information structure on a selected corpus, and using the formalism in existing spoken dialogue systems.

### References

Patrick Blackburn and Johan Bos. 2003. Computational semantics. *Theoria*, 18(46):27–45.

Jonathan Calder, Ewan Klein, and Henk Zeevat. 1988. Unification categorial grammar: A concise, extendable grammar for natural language processing. In *Proceedings of the 12th International Conerence on Computational Linguistics*, Budapest, August.

Hans Kamp and Uwe Reyle. 1993. *From Discourse to Logic*. Kluwer Academic Publishers, London.

Manfred Krifka. 1993. Focus and Presupposition in Dynamic Interpretation. *Journal of Semantics*, 10(4):269–300.

Ivana Kruijff-Korbayova. 1998. *The Dynamic Potential of Topic and Focus: A Praguian Approach to Discourse Representation Theory*. Ph.D. thesis, Faculty of Mathematics and Physics, Charles University, Prague.

Janet Pierrehumbert. 1980. *The Phonology and Phonetics of English Intonation*. Ph.D. thesis, Massachusetts Institute of Technology, Bloomington, IN. Published 1988 by Indiana University Linguistics Club.

Mark Steedman. 1990. Gapping as constituent coordination. *Linguistics and Philosophy*, 13.

Mark Steedman. 2000. *The Syntactic Process*. The MIT Press, Cambridge, Massachusetts.

Mary McGee Wood. 2000. Syntax in categorial grammar: An introduction for linguists. ESSLLI 2000, Birmingham, England. ESSLLI coursebook.

Henk Zeevat. 1988. Combining categorial grammar and unification. In U.Reyle and C.Rohrer, editors, *Natural Language Parsing and Linguistic Theories*. D.Reidel Publishing Company.