

The Use of Lexical Semantics in Interlingual Machine Translation

Bonnie J. Dorr
Institute for Advanced Computer Studies
University of Maryland
A.V. Williams Building
College Park, Maryland 20742
(301) 405-6768

Abstract

This paper describes the lexical-semantic basis for UNITRAN, an implemented scheme for translating Spanish, English, and German bidirectionally. Two claims made here are that the current representation handles many distinctions (or *divergences*) across languages without recourse to language-specific rules and that the lexical-semantic framework provides the basis for a systematic mapping between the interlingua and the syntactic structure. The representation adopted is an extended version of *lexical conceptual structure* which is suitable to the task of translating between divergent structures for two reasons: (1) it provides an *abstraction* of language-independent properties from structural idiosyncrasies; and (2) it is *compositional* in nature. The lexical-semantic approach addresses the divergence problem by using a linguistically grounded mapping that has access to parameter settings in the lexicon. We will examine a number of relevant issues including the problem of defining primitives, the issue of interlinguality, the cross-linguistic coverage of the system, and the mapping between the syntactic structure and the interlingua. A detailed example of lexical-semantic composition will be presented.

1 Introduction

This paper describes the lexical-semantic basis for UNITRAN, an implemented scheme for translating Spanish, English, and German bidirectionally. Lexical semantics provides a useful foundation for representing meaning in an interlingual representation that includes, among other things, the participants in the activities or states described by verbs. Two claims made here are that the current representation handles many distinctions (or *divergences*) across languages without recourse to language-specific rules and that the lexical-semantic framework provides the basis for a systematic mapping between the interlingua and the syntactic structure. We will focus both on the representation itself as well as how the representation is used during translation. Thus, we will discuss the theoretical basis and cross-linguistic applicability of the interlingual representation and we will also describe

how the lexical-semantic composition process maps the source language into this representation.

Many machine translation systems operate on the basis of non-compositional representations that are specifically tailored to each of the source and target languages. Frequently, such systems map between lexical entries without accounting for cases in which the arguments themselves have a special compositional status that needs to be considered during the translation mapping. Using a compositional approach allows one to define a recursive translation mapping that treats arguments of verbs as compositional units in their own right. Thus, the properties of the verb coupled with the properties of the verb's arguments are considered during each step of the translation process.

The approach taken here is to use a single underlying representation for all three languages, Spanish, English, or German, and to factor out language-specific information by means of parametric markers in the lexicon. The representation that has been adopted is an extended version of *lexical conceptual structure* (henceforth LCS) based on work by Jackendoff (1983, 1990) and further studied by Hale and Keyser (1986a, 1986b, 1989), Hale and Laughren (1983), Levin and Rappaport (1986), and Zubizarreta (1982, 1987).¹ This representation is suitable to the task of translating between divergent structures for two reasons: (1) it provides an *abstraction* of language-independent properties from structural idiosyncrasies; and (2) it is *compositional* in nature. The advantage to using this representation is that source-to-target transfer rules are not required. Instead, the system maps between the LCS and the surface syntactic form by means of a generalized linking routine that is grounded in linguistic theory.

The class of problems that this approach addresses are translation *divergences* such as those shown in figure 1.² (Literal translations are included for the Spanish and German cases.) A translation divergence arises when the natural translation of one language into another results in a very different form than that of the original. (See Dorr (1990) for a more formal discussion about divergences.) Consider the first divergence example of figure 1, *i.e.*, *conflational* divergence:

- (1) E: I stabbed John \Leftrightarrow S: Yo le di puñaladas a Juan
'I gave knife-wounds to John'

Conflation is the incorporation of necessary components of meaning (or arguments) of a given action. Here, English uses the single word *stab* for the two Spanish words *dar* (*give*) and *puñaladas* (*knife-wounds*). The *knife-wounds* component of meaning is not overtly realized in English, but is considered to be *conflated* into the main verb. By contrast, this component of meaning is not conflated in Spanish, but is overtly realized on the surface. Such an example is translated naturally by the compositional approach, which readily lends itself to the specification of arguments that may or may not be realized on the surface. We will return to this example in section 5.

¹For alternative (lexical-)semantic and case representations, see, for example, Fillmore (1968), Gruber (1965), Schank (1972), and Wilks (1973).

²Many sentences fit into these divergence classes, not just the ones listed here. Also, a single sentence may exhibit any or all of these divergences. Throughout this paper, the abbreviations E, G, and S will be used to stand for English, German, and Spanish, respectively.

Conflational	E: I stabbed John ⇕ S: Yo le di puñaladas a Juan 'I gave knife-wounds to John'
Structural	E: John entered the house ⇕ S: Juan entró en la casa 'I saw to John'
Thematic	E: I like Mary ⇕ S: Me gusta María 'Mary pleases me'
Categorial	E: I am hungry ⇕ G: Ich habe Hunger 'I have hunger'
Demotional	E: I like to eat ⇕ G: Ich esse gern 'I eat likingly'
Promotional	E: John usually goes home ⇕ G: Juan suele ir a casa 'John tends to go (to) home'
Lexical	E: John broke into the room ⇕ S: Juan forzó la entrada al cuarto 'John forced entry to the room'

Figure 1: The LCS approach is suited to the task of translating between divergent structures.

To further clarify the class of problems handled by the current approach, we will describe the rest of the divergence types shown in figure 1. The second divergence type is *structural*: the verbal object is realized as a noun phrase (*the house*) in English and as a prepositional phrase (*en la casa*) in Spanish. The third divergence type is *thematic*: the theme is realized as the verbal object (*Mary*) in English but as the subject (*María*) of the main verb in Spanish. The fourth divergence type is *categorial*: the predicate is adjectival (*hungry*) in English but nominal (*Hunger*) in German. The fifth divergence type, *demotional*, is one of two *head swapping* divergence types: the word *like* is realized as a main verb in English but as an adverbial modifier (*gern*) in German. The sixth divergence type, *promotional*, is the second *head swapping* divergence type: the modifier (*usually*) is realized as an adverbial phrase in English but as the main verb *soler* in Spanish.³ Finally, the seventh divergence type is a *lexical* divergence: the main verb is *break* in English but a different verb *forzar* (literally *force*) in Spanish.

The key to being able to handle these divergence types is modularity: the UNTRAN system consists of two processing components that allow for a decoupling of syntactic decisions from lexical-semantic decisions. The syntactic component parses

³The distinction between demotional and promotional divergences is not obvious at first glance. In both examples in figure 1, the translation mapping associates a main verb with an adverbial satellite, or *vice versa*. The justification for distinguishing between these two *head swapping* cases is given in Dorr (in press).

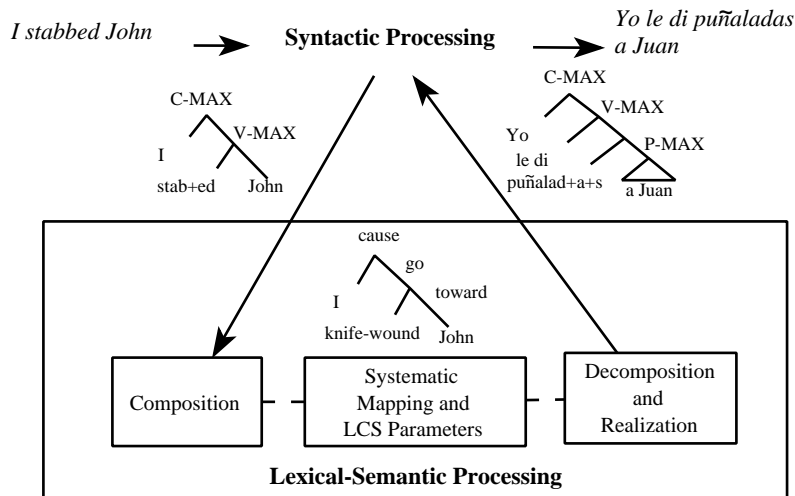


Figure 2: The lexical-semantic processing component is designed so that the composition and decomposition/realization processes both rely on the same systematic mapping and LCS parameters.

and generates sentences on the basis of syntactic principles that are parameterized to handle language-specific idiosyncrasies. The lexical-semantic component uses the output of the parser to build a syntax-free interlingua that is then used as input to the generator. We will focus primarily on the second component of the system, which is illustrated in figure 2 for example (1).

Because the syntactic variation between the source and target languages has been factored out at the level of lexical-semantic processing, the lexical-semantic component is concerned solely with the tasks of composing and decomposing the interlingual representation without regard to the syntactic form of the source- and target-language sentences. Two top-level tasks are performed during lexical-semantic processing: (1) composition of a single interlingual representation from the lexical forms associated with the syntactic tree nodes; and (2) decomposition of the interlingual representation into its surface-level realization.⁴ This paper demonstrates that, just as the syntactic component relies on parameters of variation, so too does the lexical-semantic component. In particular, the composition and decomposition processes make use of parameter settings in the lexicon in order to accommodate divergences such as those shown in figure 1.

It should be noted that the representation relies solely on lexical-semantic information. “Deeper” notions of meaning, *e.g.*, aspectual, contextual, discourse, domain, and world knowledge, are not included in the representation. While such notions are arguably necessary for the general solution to the problem of machine translation, they are not necessary for the solution to the class of problems that are addressed here, *i.e.*, translation divergences. Nevertheless, the use of a lexical-semantic representation does not preclude the possibility of superimposing “deeper” knowledge onto the current framework. On the contrary, a knowledge-based meaning representation such as that of Nirenburg and Levin (1989) could significantly enhance the translation mapping, particularly during the processes of lexical selection and generation, which will not be discussed here.

We will examine a number of issues concerning the lexical-semantic framework

⁴This paper will only discuss the first task of lexical-semantic processing. The decomposition and realization processes are described in more detail in Dorr (in press).

adopted in UNITRAN. Section 2 provides the motivation for the interlingual approach and discusses issues relevant to translation divergences. Section 3 addresses the problem of defining primitives for the interlingua, the issue of interlinguality, and the coverage of the scheme. Section 4 describes the structure of the interlingual representation, the use of the representation in lexical entries, the organization of the lexicon, and the extensions that have been made to the original formulation of the LCS. Section 5 defines the parameterized mapping between the syntactic structure and the interlingua and provides an example of lexical-semantic composition in detail. Finally, section 6 addresses issues relevant to the approach and discusses the limitations of, and future work on, the UNITRAN system.

2 Motivation for the Interlingual Approach

Machine translation has been a particularly difficult problem in the area of natural language processing for over four decades. Early approaches to translation failed in part because interaction effects of complex phenomena made translation appear to be unmanageable. Later approaches to the problem have achieved varying degrees of success. (See Hutchins (1986), King (1987), Maxwell *et al.* (1988), Nirenburg (1987), Nirenburg *et al.* (1992), Slocum (1988), and Hutchins and Somers (1992) for cogent reviews of this area.) Only certain issues of relevance to the current approach will be discussed in this section. In particular, the motivation for the interlingual approach will be presented with specific reference to the divergence classification shown in figure 1.

One might argue that several of the examples from figure 1 could be translated using the direct translation approach with little loss of information. For example, if we choose English to be the target language, we could directly translate the three of the examples from figure 1 as follows:

- (2) S: Yo le di puñaladas a Juan \Rightarrow E: I inflicted knife wounds on John
- (3) S: Me gusta María \Rightarrow E: Mary pleases me
- (4) G: Ich esse gern \Rightarrow E: I eat enjoyably

The problem with taking such an approach is that it is not general enough to handle a wide range of cases. For example, if we translate the word *gern* directly to the word *enjoyably* as in case (4), we will run into problems when we try to translate *gern* in other contexts. As it turns out, the adverb *gern* may be used in conjunction with *haben* to mean *like*: *Ich habe Marie gern* ('I like Mary'). The literal translation, *I have Mary enjoyably*, is not only stylistically unattractive, but it is not a valid translation for this sentence. Another problem with the direct-mapping approach is that it is not bidirectional in the general case. Thus, even if we did take (2), (3), and (4) as the desired translations, we would not be able to apply the same direct mapping in the reverse direction (*i.e.*, translating from the more standard English version to Spanish and German). For example, the same mapping could not be used to translate *stab* and *like* into Spanish and German. It is clear that a uniform method for bidirectional translation is required.

The direct approach to translation has been largely discounted as an alternative to the interlingual approach. However, there have been a number of arguments for the more commonly used transfer approach as an alternative to the interlingual approach. (See, *e.g.*, Arnold and Sadler (1990), Boitet (1988), and Vauquois and Boitet (1985).) Paradoxically, these anti-interlingual arguments are based precisely on the same types of examples that have motivated the current research (*e.g.*, those sentences that exhibit the types of divergences shown above). The assumption of previous approaches is that it would be too difficult to design an interlingual representation that includes enough information to accommodate complex divergences. This study argues to the contrary, adopting the view that complex divergences are precisely what necessitates the use of an interlingual representation because the interlingual approach allows surface syntactic distinctions to be represented at a level that is independent from that of the underlying “meaning” of the source and target sentences. Factoring out these surface-level distinctions allows cross-linguistic generalizations to be captured at the level of lexical-semantic structure.

Some examples of the types of translations that are used to justify the use of a transfer approach are:⁵

- (5) **Demotional divergence:**
G: Johann küßt Marie gern⁶ ⇔ E: John likes to kiss Mary
‘John kisses Mary likingly’
- (6) **Promotional divergence:**
F: Jean a failli finir le livre⁷ ⇔ E: John has almost finished the book
‘John has missed to finish the book’
- (7) **Conflational divergence:**
F: Ils entrent dans la salle en courant⁸ ⇔ E: They run into the room
‘They enter the room in running’

The implicit assumption rejected by the current approach is that it would be impossible to design an interlingual system that supports a systematic mapping between the source and target languages in cases such as (5), (6), and (7). As we will see, these types of divergences are precisely the ones that are handled by the UNITRAN system.

⁵These examples are taken from Arnold *et al.* (1988) and Arnold and Sadler (1990). The divergence types are specified using the terminology of the present author, not that of the authors from whom these (and later) examples are taken.

⁶The example given by Arnold *et al.* (1988) is actually a translation from the Dutch equivalent of the German sentence (*i.e.*, *Jan kust Marie graag*), but the construction is entirely analogous.

⁷This example is an adapted version of one taken from Arnold and Sadler (1990): *Jean a manqué de finir le livre*. Judgments as to the naturalness and acceptability of this sentence differ among native French speakers; the verb *manquer* does not convey the “almost” meaning for all speakers. Thus, the more acceptable version that uses *faillir* is given here. In any case, the concept of promotion is valid, regardless of which main verb is used.

⁸Again, this is an adapted version of a sentence taken from Arnold and Sadler (1990); the original form of the sentence did not include the preposition *dans* which forced the sentence to have a questionable status among native French speakers. Thus, the more acceptable version that uses *dans* is given here. In any case, the concept of conflation is valid with, or without, the preposition since the lexical items of interest in this phenomenon are *entrer* and *en courant*.

For the purposes of comparison, we will briefly discuss the approach taken by the MiMo system (Arnold *et al.* (1988), Arnold and Sadler (1990), van Noord *et al.* (1989), van Noord *et al.* (1990), and Sadler *et al.* (1990)) for these examples. Other systems that have attempted to deal with these (and other) divergence types are TAUM (Colmerauer *et al.* (1971)), GETA/ARIANE (Vauquois and Boitet (1985) and Boitet (1987)), LMT (McCord (1989)), LFG-MT (Kaplan *et al.* (1989)), METAL (Alonso (1990) and Thurmair (1990), and LTAG (Abeillé *et al.* (1990)). Specific examples of these systems are discussed in Dorr (in press).

The MiMo system developed from an effort that had its beginnings within the Eurotra framework. (For cogent descriptions of the Eurotra project, see, *e.g.*, Arnold and des Tombe (1987), Copeland *et al.* (1991), and Johnson *et al.* (1985).) Whereas the ‘mainstream’ Eurotra work stops short of describing how translation divergences are handled, the MiMo project has extended the Eurotra formalism to handle a wide range of divergence classes. MiMo is a structural transfer approach that addresses the solution to divergence cases such as (5), (6), and (7) above. To achieve the translation mapping in these examples, the transfer approach requires the existence of transfer entries. In particular, the verb *like* must be related to the adverb *gern* by means of a transfer entry of the form:

$$(8) \quad r!((\text{cat} = \text{S}).[\text{mod} = \text{GERN}]) \Leftrightarrow \text{LIKE}((\text{cat} = \text{S}).[r!\text{arg1}])$$

This rule indicates that the word *gern* would be realized as an adverbial modifier in German, whereas the English counterpart would be realized as a main verb that takes a sentential argument.

As for (6) and (7), the required transfer rules are (9) and (10), respectively:

$$(9) \quad \text{FAILLIR}([1!\text{arg1}], [!\text{arg2}=r![2!\text{arg2}]] \Leftrightarrow r!([1!\text{arg1}], [2!\text{arg2}][\text{mod} = \text{ALMOST}])$$

$$(10) \quad \text{ENTRER}([1!\text{arg2}], [\text{mod} = \text{EN}[!\text{arg1}=\text{COURIR}]] \Leftrightarrow \text{RUN}([\text{mod}=\text{INTO}[1!\text{arg2}]]$$

Rule (9) ensures that the French main verb *faillir* (in (6)) is translated as an English adverbial *almost* in English. (The *r!* variable is bound to *finir* in French and *finish* in English.) Rule (10) maps the compound French construction *entrer en courant* (in (7)) into the simple English construction *run into*.

We need not dwell on the well-known problem that such transfer entries are required for each source-language/target-language pair of the system, including pairs that exhibit less complicated divergences such as in the *like-gustar* case shown in figure 1 (as well as pairs for which there is no divergence at all). Suffice it to say that specifying transfer mappings for all of the lexical items of each source-language/target-language pair is very tedious work. (See, *e.g.*, Bennett *et al.* (1986) for additional discussion.) There are, however, other problems with the transfer approach that are worth mentioning here. One problem is that much of the information that is, or could be, lexically stored (*e.g.*, the argument structure of the word *entrer*) is included in the transfer rules as well. The use of transfer rules results in a severe proliferation of redundancy on a per-language basis. For example, *faillir* is not the only verb that gives rise to the type of construction shown in (6); in fact, Arnold and Sadler (1990) present a similar verb, *venir* (which translates to the adverbial *just*) that is precisely analogous to the verb *faillir*. Thus, rules analogous to (9) must be constructed for *venir*, differing only in that *venir* is used

in place of *faillir* and *just* in place of *almost*.⁹ The multiplicative effect of all of these combinations wreaks havoc with the number of transfer rules that are required on a per-language basis.

Another problem with the transfer approach is that it misses a number of cross-linguistic generalizations. For example, not only do constructions such as (5), (6) and (7) arise *within* a single language, but such constructions exist *across* languages as well. The root of the problem with the MiMo approach is that it does not use a canonical representation that would factor out this redundancy. As we will see shortly, the UNITRAN system takes advantage of such a representation to map between languages in a more systematic and uniform fashion.

3 Interlinguality and Linguistic Coverage

In order to adopt an interlingual approach to machine translation, one must construct a language-independent representation that lends itself readily to the specification of a systematic mapping that operates uniformly across all languages. To meet this objective, one needs a clear characterization of the entire range of divergences that could possibly arise in machine translation. This characterization can emerge only from a serious cross-linguistic investigation into the adequacy of different lexical representations for natural language. Such a representation has been constructed for use in the UNITRAN system *i.e.*, the *lexical conceptual structure* (LCS). This representation has been adapted to the UNITRAN machine translation model in that it is associated with an algorithm for recursive composition and decomposition of the interlingual form, and it is linked systematically to the syntactic structure, both during parsing as well as during generation. The LCS abstracts away from syntax just far enough to enable language independent encoding, while retaining enough structure to be sensitive to the requirements for language translation and, in particular, the resolution of divergences such as those shown in figure 1.

The approach described here is an alternative to interlingual solutions that have been criticized in the past. For example, Bennett *et al.* (1986, p. 86) criticizes two possible approaches to an interlingual system: (1) using a canonical syntactic structure (because of the degree of language specificity); and (2) defining an interlingua that is specifically geared toward the languages in the system (because of the difficulty of adding a new language). The solution described here does not fall into either of these two categories. Rather, it appeals to a different notion of interlingua that excludes language-specific syntactic properties and includes language-independent lexical-semantic properties.

The remainder of this section examines the primitive building blocks of the interlingua and the issues of interlinguality and linguistic coverage of the system.

3.1 Primitive Building Blocks of the Interlingua

The field of machine translation has (almost from the beginning) been concerned with the use of a “deep semantic representation” and with looking for “universals”

⁹In addition, unlike *faillir*, the verb *venir* requires the particle *de* to be inserted: *Jean vient de tomber* (John just fell).

for translation. One of the biggest objections to the use of an interlingual representation is that it requires the system designer to define a set of primitives (to represent the information to be translated) which allows the mapping to and from the languages in question. Because it is generally difficult to define such a set, many researchers have abandoned this model. (See, *e.g.*, Vauquois and Boitet (1985).) However, recently, there has been a resurgence of interest in the area of lexical representation and organization (with special reference to verbs) that has initiated an ongoing effort to delimit the classes of lexical knowledge required to process natural language. (See, *e.g.*, Grimshaw (1990), Hale and Laughren (1983), Hale and Keyser (1986a, 1986b, 1989), Jackendoff (1983, 1990), Levin and Rappaport (1986), Levin (1985, in press), Pustejovsky (1988, 1989, 1990, 1991), Rappaport *et al.* (1987), Rappaport and Levin (1988), Olsen (1991), and Zubizarreta (1982, 1987).) As a result of this effort, it has become increasingly more feasible to isolate the components of meaning common to verbs participating in particular classes. These components of meaning can then be used to determine the lexical representation of verbs across languages.

The LCS approach views semantic representation as a subset of conceptual structure, *i.e.*, the language of mental representation. Jackendoff's approach includes *types* such as Event and State, which are specialized into *primitives* such as GO, STAY, BE, GO-EXT, and ORIENT. As an example of how the primitive GO is used to represent sentence semantics, consider the following sentence:

- (11) (i) The ball rolled toward Beth.
(ii) [_{Event} GO ([_{Thing} BALL],
[_{Path} TOWARD ([_{Position} AT ([_{Thing} BALL], [_{Thing} BETH])))]])

This representation illustrates one dimension (*i.e.*, the *spatial* dimension) of Jackendoff's representation. Another dimension is the *causal* dimension, which includes the primitives CAUSE and LET. These primitives take a Thing and an Event as arguments. Thus, we could embed the structure shown in (11)(ii) within a causative construction:

- (12) (i) John rolled the ball toward Beth.
(ii) [_{Event} CAUSE
([_{Thing} JOHN],
[_{Event} GO ([_{Thing} BALL],
[_{Path} TOWARD [_{Position} AT ([_{Thing} BALL], [_{Thing} BETH]))]])]

Jackendoff includes a third dimension by introducing the notion of *field*. This dimension extends the semantic coverage of spatially oriented primitives to other domains such as Possessional, Temporal, Identificational, Circumstantial, and Existential.¹⁰ For example, the primitive GO_{POSS} refers to a GO event in the Possessional field as in the following sentence:

¹⁰The label *Loc* has been adopted to distinguish the spatial field from the non-spatial fields. Note that the spatial field is used to denote the primitives that fall in the spatial dimension. Jackendoff argues that spatial primitives are more fundamental to those of other domains (*e.g.*, Possessional); in particular, all primitives from other domains pattern after those in the spatial domain with respect to argument-structure constraints (to be discussed shortly). Thus, spatial primitives have their own special status as an independent dimension.

- (13) (i) Beth received the doll.
(ii) $[\text{Event GO}_{\text{Poss}}$
 $([\text{Thing DOLL}],$
 $[\text{Path TO}_{\text{Poss}}([\text{Position AT}_{\text{Poss}}([\text{Thing DOLL}], [\text{Thing BETH}]])])]$

To further illustrate the notion of field, the GO primitive can be used in the Temporal and Identificational fields:

- (14) (i) The meeting went from 2:00 to 4:00.
(ii) $[\text{Event GO}_{\text{Temp}}$
 $([\text{Thing MEETING}],$
 $[\text{Path FROM}_{\text{Temp}}([\text{Position AT}_{\text{Temp}}([\text{Thing MEETING}], [\text{Time 2:00}])])]$
 $[\text{Path TO}_{\text{Temp}}([\text{Position AT}_{\text{Temp}}([\text{Thing MEETING}], [\text{Time 4:00}])])]$
- (15) (i) The frog turned into a prince.
(ii) $[\text{Event GO}_{\text{Ident}}$
 $([\text{Thing FROG}],$
 $[\text{Path TO}_{\text{Ident}}([\text{Position AT}_{\text{Ident}}([\text{Thing FROG}], [\text{Thing PRINCE}])])]$

As these examples illustrate, there are also other primitives that are included in the LCS framework. In particular, the Position and Path types are used to include primitives such as AT and TO.¹¹ Furthermore, the Thing, Location, Time, Manner, and Property types are used. Figure 3 shows a subset of the types and primitives that are currently used in the LCS scheme.¹²

A crucial point concerning the LCS representation is that, although it appears to be somewhat “English-like,” it is only superficially so by virtue of the labels of the primitives (*e.g.*, GO, TO, *etc.*) that are used in the representation. These primitives were chosen on the basis of an extensive cross-linguistic investigation, though they may be labeled and used in a fashion that appears to be modeled after a particular language.

An additional point about the representation is that a large percentage of primitives fall under the Manner category, which may appear to be peculiar at first glance. Because Jackendoff says very little about the function of the Manner type in his description of conceptual structure, the approach taken here is to use the

¹¹The Position type corresponds to the Place type used by Jackendoff (1983). An extension that has been made to the Position type is that it is a two-place predicate rather than a one-place predicate. For example, in (14)(ii), the MEETING argument appears both internally and externally to the AT_{Temp} node. This is due to the observation that primitives such as AT, IN, ON, *etc.* are actually relations between two arguments (*e.g.*, the representation for *the book is on the table* incorporates the relation ON(BOOK, TABLE) as part of its meaning). The use of the two-place predicate also allows for additional type-checking when the LCS representation for a word (*e.g.*, *on*) is composed with the LCS representation for another word (*e.g.*, *put*) in order to derive the underlying representation for the entire concept (*e.g.*, *put the book on the table*). Both argument positions must be checked for a match before the composition can take place. In the case of *put the book on the table*, the two arguments associated with the predicate ON are: a movable Thing and a locative Thing, respectively. It cannot be the case that one of these arguments is, for example, a State or an Event. (The LCS composition process will be discussed in section 5.)

¹²Note that primitives have not been included for *John*, *Beth*, *I*, *me*, *it*, *etc.* In actuality, proper names are represented by the PERSON primitive and referring expressions (*e.g.*, pronouns) are represented by the REFERENT primitive. For notational convenience, the examples and figures will continue to use the more informative labels in place of these primitives.

Type	Primitives
Event	CAUSE, LET, GO, STAY
State	BE, GO-EXT, ORIENT
Position	AT, IN, ON
Path	TO, FROM, TOWARD, AWAY-FROM, VIA
Thing	BOOK, PERSON, REFERENT, KNIFE-WOUND, KNIFE, SHARP-OBJECT, WOUND, FOOT, CURRENCY, PAINT, FLUID, ROOM, SURFACE, WALL, HOUSE, BALL, DOLL, MEETING, FROG
Property	TIRE, HUNGRY, PLEASED, BROKEN, ASLEEP, DEAD, STRETCHED, HAPPY, RED, HOT, FAR, BIG, EASY, CERTAIN
Location	HERE, THERE, LEFT, RIGHT, UP, DOWN
Time	TODAY, SATURDAY, 2:00, 4:00
Manner	FORCEFULLY, LIKINGLY, WELL, QUICKLY, DANCINGLY, SEEMINGLY, HAPPILY, LOVINGLY, PLEASINGLY, GIFTINGLY, UPWARD, DOWNWARD, WITHIN, HABITUALLY

Figure 3: LCS primitives are divided into a handful of types. A subset of the types and primitives used by UNITRAN is shown here.

Manner component to distinguish between verbs that fall within the same linguistic class when no other distinguishing features are available. Additional discussion about this point is given in section 4.4.

3.2 Cross-Linguistic Adequacy and Coverage

This section addresses the issue of the overall adequacy of the LCS scheme for interlingual machine translation. In particular, we will discuss the status of the primitive building blocks as well as the degree of linguistic coverage provided by the representation.

In defining a set of primitives for an interlingua, one must abide by a number of general restrictions such as those proposed by Wilks (1987): *finitude*, *comprehensiveness*, *independence*, *noncircularity*, and *primitiveness*. While the LCS scheme has been designed to fulfill these requirements, it has also been designed to fulfill an additional requirement which is considered to be central to the current approach:

- (16) **Linguistic Generalization:** The primitives should be defined so that their combination captures both conceptual and syntactic generalities of actions or entities that might otherwise be represented differently.

That is, each action (*e.g.*, GO) and entity (*e.g.*, PERSON) must be associated with a representation that is both conceptually plausible and systematically related to a syntactic structure. The choice of one representation over another, then, depends on the degree to which the representation captures conceptual and syntactic generalities. The LCS approach fulfills this requirement by imposing certain constraints on the way the primitives may be combined in the conceptual structure and realized in the syntactic structure. These constraints are available along the three dimensions described in the last section. The *spatial* dimension, which provides the basic set of

primitive building blocks, GO, STAY, BE, GO-EXT, and ORIENT, must adhere to the following constraints on argument structure:

Events		
Primitive	Argument 1	Argument 2
GO	Thing	Path
STAY	Thing	Position

(17) (i)

States		
Primitive	Argument 1	Argument 2
BE	Thing	Position
ORIENT	Thing	Path
GO-EXT	Thing	Path

(ii)

The *causal* dimension must adhere to the following constraints on argument structure:¹³

Primitive	Argument 1	Argument 2
CAUSE	{ Thing Event }	{ Event State }
LET	{ Thing Event }	{ Event State }

(18)

Finally, the *field* dimension imposes the following constraints:

Field	Argument 1	Argument 2 ¹⁴
Locational	{ Thing Event }	Location
Possessional	Thing	Thing
Temporal	{ Event State }	Time
Identificational	Thing	{ Thing Property }
Circumstantial	Thing	{ Event State }
Existential	Thing	EXT

(19)

A crucial point concerning the linguistic generalization proposed in (16) is that it requires the rules for combining primitives to have a direct correspondence to the syntactic structure. This is precisely the nature of the lexical-semantic constraints shown in (17)–(19). One might think of these constraints as a means of specifying argument-selection restrictions that are imposed on the lexical-semantic representation and that have their corresponding syntactic reflexes in the surface sentence. For example, the ORIENT primitive selects a Thing and a Path as in the following lexical conceptual structure:

(20) [_{State} ORIENT_{Loc}
 ([_{Thing} SIGN],
 [_{Path} TOWARD_{Loc} ([_{Position} AT_{Loc} ([_{Thing} SIGN], [_{Location} BOSTON])])])]

¹³The {} notation denotes choice. For example, the first argument of the CAUSE primitive may either be a Thing or an Event.

¹⁴Technically, the second argument for each of these fields is a Path or a Position. For the purposes of the current description, the column under “Argument 2” refers to the lowest leaf node embedded inside of the second argument.

These conceptual constituents are then realized in the syntactic structure, respectively, as a noun phrase subject and a prepositional phrase object:

(21) [C-MAX [I-MAX [N-MAX The sign] [V-MAX points [P-MAX toward [N-MAX Boston]]]]]

Thus, the constraints support the linguistic generalization requirement in (16) by allowing the primitives to be combined in such a way as to capture both conceptual and syntactic generalities. In the next section we will see that the constraints are consistent with a systematic correspondence between syntactic structure and lexical-semantic structure.

Clearly, such an approach could not hope to distinguish among all possible meanings of every lexical item. However, the intent of the current scheme is to provide a linguistically motivated approach to handling translation divergences, not to describe the deep semantic content of lexical items. While the primitives proposed in the current approach have been defined specifically for the task of resolving translation divergences, there is clearly additional information that would be required for solving other pieces of the translation problem. In the spirit of Wilks (1987), the current approach adopts the view that there may be other representational *levels*, with their own primitive definitions, that might be superimposed on top of the current framework to handle other types of problems (*e.g.*, distinguishing between verbs, like *donate* and *give*, which occur in the same lexical-semantic class).

One could argue against this approach on the grounds that the LCS representation has been constructed in such a way that is biased toward the particular translation problem to be solved. That is, the choice of lexical-semantic primitives and their allowable combinations appear to rely heavily on linguistic knowledge about the nature of translation divergences. However, it is not clear that developing the interlingua on the basis of this knowledge is a drawback to the approach. It is much more worthwhile to construct a representation on the basis of a carefully studied, and finitely specified, classification system than on the basis of vague intuitions about the nature of what a word means in a particular language. Furthermore, since the divergence types addressed by this approach (*i.e.*, the classification given in figure 1) are expected to cover all potential source-language/target-language distinctions based on properties of lexical items, it is considered an advantage to use this classification as the basis for the representation. Since the range of divergence possibilities has been carefully delimited, it is feasible to use the classification to isolate components of verb meaning for the construction of an appropriate interlingua.

An additional argument for constructing the lexical-semantic representation based on a study of cross-linguistic divergences is that it is consistent with the philosophy behind well-grounded lexical-semantic theories such as *meaning-text theory* (MTT) by Mel'čuk and Polguère (1987, p. 266); that is, rather than postulating a set of semantic primitives *a priori*, the intention is to discover them by means of a "painstaking process of semantic decomposition applied to thousands of actual lexical items." Given that recent research has made it increasingly more feasible to isolate components of verb meaning, this approach to the construction of an interlingual representation is no longer impractical. Note that this is in direct contrast to other approaches that employ primitives that are chosen *a priori* such as the conceptual dependency approach by Schank (1972, 1973, 1975) and Schank and Abelson

(1977), for which there has never been a cross-linguistic investigation into the applicability of the primitives. (We will return to further discussion of the conceptual dependency approach shortly.)

A final point to be made about the current interlingual representation is that it is intended to include those aspects of lexical knowledge related to argument structure, not “deeper” notions of meaning such as aspectual, contextual, domain, and world knowledge. Previous approaches that have employed deep semantic representations say very little about how to construct a single systematic mapping that operates cross-linguistically between the syntactic structure and the lexical-semantic representation. The types of divergences that are addressed here do not require knowledge about the world, but rather about general lexical-semantic properties that hold across languages. The current approach strives to fill an obvious gap corresponding to the linking rules between the lexicon and syntactic structure. Clearly, the techniques used in deeper knowledge approaches (*e.g.*, KBMT by Carbonell and Tomita (1987) and Nirenburg *et al.* (1992)) are necessary for filling other gaps in the translation process, most notably the process of lexical selection, which generally requires more than just argument-structure information for disambiguation of nouns and attachment of modifiers.¹⁵ The current approach is intended to be a supplement, not a substitute, for knowledge-based techniques. It is expected that a fully interlingual translation system would require knowledge-based techniques to operate in tandem with the techniques described here.

A number of machine translation approaches that have used an interlingua based on “deeper” knowledge representations are presented in Wilks (1973), Vauquois (1975) (CETA),¹⁶ and more recently in Carbonell and Tomita (1987) and Nirenburg *et al.* (1992) (KBMT), Nirenburg *et al.* (1987) (TRANSLATOR), Muraki (1987) (PIVOT), Uchida (1989) (ATLAS), among others. We will not survey all of the different ways that one can construct an interlingua here. However, given that the LCS representation has been commonly compared to the *conceptual dependency* (CD) representation, a brief comparison is presented here.

Early translation approaches used the CD representation as the basis for interlingual machine translation. (See, for example, Schank (1975), Schank and Abelson (1977), and Lytinen and Schank (1982).) These approaches were similar to that of UNITRAN in that they relied on a compositional representation based on a small set of primitives. However, a well-known problem with the traditional CD-based approach is that it provides a target-language paraphrase of the source-language sentence rather than a target-language translation of the source-language sentence; it is now widely accepted that this approach is not adequate for machine translation. The paraphrased output is a symptom of a more serious problem, *i.e.*, that CD-based systems lack a canonical mapping from the interlingual representation to the syntactic structure. For example, there is no uniform mechanism for handling even the simplest divergences such as the thematic reversal in the translation of *I like Mary* to Spanish (see figure 1) because there is no systematic relationship between

¹⁵The issue of disambiguation will be discussed in section 6.

¹⁶Although the CETA system has been classified as interlingual here, it should be pointed out that there have been a number of persuasive arguments against this classification due to the fact that the lexicon used in the CETA system had a bilingual transfer-like mechanism (see, *e.g.*, Hutchins (1986) and Nirenburg *et al.* (1992)).

Class of Verb	Examples
position	be, remain, ...
change of position	fall, throw, drop, change, move, slide, float, roll, fly, bounce, move, drop, turn, rotate, shift, ...
directed motion	enter, break into, bring, carry, remove, come, go, leave, arrive, descend, ascend, put, raise, lower, ...
motion with manner	sail, walk, stroll, jog, march, gallop, jump, float, dance, run, skip, ...
exchange	buy, sell, trade, ...
physical state	be, remain, keep, leave, ...
change of physical state	open, close, melt, redden, soften, break, crack, freeze, harden, dry, whiten, grow, change, become, ...
orientation	point, aim, face, ...
existence	exist, build, grow, shape, make, whittle, spin, carve, weave, bake, fashion, create, appear, disappear, reappear, persist, ...
circumstance	be, start, stop, continue, keep, exempt, ...
range	go, last, extend, intend, aim, range, ...
change of ownership	give, take, receive, relinquish, borrow, lend, loan, steal, ...
ownership	belong, remain, keep, own, ...
ingestion	eat, drink, smoke, gobble, munch, sip, ...
psychological state	like, fear, admire, detest, despise, enjoy, esteem, hate, honor, love please, scare, amuse, astonish, bore, surprise, stun, terrify, thrill, ...
perception and communication	see, hear, smell, feel, look, watch, listen, learn, sniff, show, tell, talk, speak, shout, whisper, scream, ...
mental process	know, learn, ...
cost	cost, charge, ...
load/spray	smear, load, cram, spray, stuff, pile, stack, splash, ...
contact/effect	cut, stab, crush, smash, pierce, bite, shoot, spear, ...

Figure 4: A subset of the linguistic classes from Levin (1985, in press) has been implemented using the LCS framework.

the conceptual arguments of the CD representation and the corresponding syntactic positions in which these arguments are realized. The LCS approach differs in that it provides a mapping between the underlying concept and the syntactic structure that is systematically defined in all languages. Section 5 describes this LCS-syntax mapping in more detail.

Traditional interlingual approaches pay a high price for incorporating too much “deeper” knowledge without preserving structurally defined lexical-semantic information. This is not to say that research should head in the direction of “shallow” interlingual representations (*e.g.*, see Sharp (1985)). The current approach attempts to achieve a middle ground between representations that encode too much information and those that encode too little. Specifically, the current representation captures parametric information that accommodates cross-linguistic divergences without losing the systematic relation between the interlingual representation and the syntactic structure. Moreover, the representation lends itself readily to the specification of a systematic mapping that operates uniformly across all languages.

In attempting to cover a wide range of phenomena, a subset of the linguistic classes from Levin (1985, in press) has been implemented using the LCS framework. The implemented classes are summarized in figure 4 with examples of verbs given in English. It is expected that this verb classification scheme will need further refinement as more properties of verbs are identified. Some examples of the types of sentences that are represented using the LCS framework are shown in figure 5 with

Type	Primitive-Field	Example
Event	GO _{Poss}	Beth received the doll.
	GO _{Ident}	Elise became a mother.
	GO _{Temp}	The meeting went from 2:00 to 4:00.
	GO _{Loc}	We moved the statue from the park to the zoo.
	GO _{Circ}	John started shipping goods to California.
	GO _{Exist}	John built a house.
	STAY _{Poss}	Amy kept the doll.
	STAY _{Ident}	The coach remained a jerk.
	STAY _{Temp}	We kept the meeting at noon.
	STAY _{Loc}	We kept the statue in the park.
	STAY _{Circ}	John kept shipping goods to California.
	STAY _{Exist}	The situation persisted.
State	BE _{Poss}	The doll belongs to Beth.
	BE _{Ident}	Elise is a pianist.
	BE _{Temp}	The meeting is at noon.
	BE _{Loc}	The statue is in the park.
	BE _{Circ}	John is shipping goods to California.
	BE _{Exist}	Descartes exists.
	GO-EXT _{Ident}	Our clients range from psychiatrists to psychopaths.
	GO-EXT _{Temp}	The meeting lasted from noon to night.
	GO-EXT _{Loc}	The road went from Boston to Albany.
	ORIENT _{Loc}	The sign points to Philadelphia.
	ORIENT _{Circ}	John intended to ship goods to California.

Figure 5: A number of different types of sentences are represented using the Event and State primitives of the LCS framework.

their corresponding primitive-field combinations.¹⁷

4 Lexical-Semantic Representation

We have already seen an informal description of the interlingua in section 3. The first two parts of this section present a more formal description of the interlingual structure and demonstrate how this structure is used in lexical entries. The second two parts describe the organization of the lexicon and the extensions that have been made to the original formulation of the LCS.

4.1 Formal Structure of the Interlingua

The interlingua is defined primarily on the basis of *domination* relations in the LCS; unlike the syntactic structure, *linear ordering* is not as critical as hierarchical structure in preserving the meaning of the LCS. However, certain positioning conventions are retained in order to preserve uniformity during the mapping between the interlingua and the syntactic structure.

An LCS is potentially comprised of four types of lexical-semantic tokens. The first is the *logical head*, *i.e.*, the root node of an LCS structure (or substructure). The second token is the *logical subject*; there is only one such argument, and it is always the highest/left-most constituent under the logical head. The third type of lexical-semantic token is the *logical argument*; there may be any number of these,

¹⁷The example sentences were taken from Siskind (1989), with minor modifications.

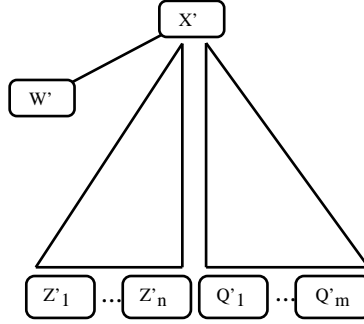


Figure 6: The formal structure of the interlingua defines canonical positions for the logical head, logical subject, logical arguments, and logical modifiers.

all occurring to the right of the logical subject. Finally, there are any number of (optional) *logical modifiers*, which are ordered, by convention, to the right of logical arguments. Thus, the overall structure is defined as follows:

$$(22) \quad [_{T(X')} X' \\ ([_{T(W')} W'], \\ [_{T(Z'_1)} Z'_1], \dots, [_{T(Z'_n)} Z'_n] \\ [_{T(Q'_1)} Q'_1], \dots, [_{T(Q'_m)} Q'_m])]$$

where X' is the logical head, W' is the logical subject, Z'_1, \dots, Z'_n are the logical arguments, Q'_1, \dots, Q'_m are the logical modifiers, and $T(\phi)$ is the LCS type corresponding to the primitive ϕ . The tree-like representation corresponding to this structure is shown in figure 6.¹⁸

Consider the the following translation example:

$$(23) \quad \text{E: John happily entered the room} \Leftrightarrow \text{S: Juan felizmente entr3 al cuarto}$$

The LCS representation corresponding to this example is the following:

$$(24) \quad [_{\text{Event}} \text{GO}_{\text{Loc}} \\ ([_{\text{Thing}} \text{John}], \\ [_{\text{Path}} \text{TO}_{\text{Loc}} ([_{\text{Position}} \text{IN}_{\text{Loc}} ([_{\text{Thing}} \text{JOHN}], [_{\text{Location}} \text{ROOM}])])]) \\ [_{\text{Manner}} \text{HAPPILY}]]]$$

where X' corresponds to GO_{Loc} , $T(X')$ corresponds to Event , W' corresponds to JOHN , $T(W')$ corresponds to Thing , Z'_1 corresponds to TO_{Loc} , $T(Z'_1)$ corresponds to Path , Q'_1 corresponds to HAPPILY , and $T(Q'_1)$ corresponds to Manner .

There are two properties of the LCS representation that enable the system to derive an appropriate interlingual representation even in cases where the source and target sentences are not structurally equivalent. The first is that the LCS representation provides an *abstraction* of language-independent properties from structural idiosyncrasies. For example, the location in example (24) (*i.e.*, where John

¹⁸For ease of illustration, this diagram shows the primitive-field without the type. We will retain this convention throughout the paper. In addition, the ordering given in the syntactic tree is the one used for English and Spanish. In actuality, the syntactic component determines the appropriate syntactic ordering from the setting of a parameter that is not described here. (See Dorr (in press) for more details about the syntactic component.) For the purposes of illustration, we will use this ordering throughout the rest of the paper.

is entering) is associated with a single, canonical representation in the interlingua regardless of how it is syntactically realized on the surface. In the case of English, the location (*room*) is realized as an object of the main verb, whereas in the case of Spanish, the location (*cuarto*) is realized as an object of the preposition *a*. The second property is that the LCS representation is *compositional* in nature. Through a process called lexical-semantic composition (to be described in section 5.1), the implicit $[_{\text{Path}} \text{TO}_{\text{Loc}} ([_{\text{Position}} \text{IN}_{\text{Loc}} \dots])]$ argument associated with the word *enter* is made available for realization as the word *a* in Spanish, even though the argument does not appear overtly in the English sentence.

4.2 Lexical Entries

We have just seen how the LCS representation is used as the basis of the interlingua. We now turn to the use of the LCS in lexical entries. We shall distinguish between these two cases by using the term CLCS (*composed* LCS) for the interlingua and the term RLCS (*root word* LCS) for the lexical entry.

Each language processed by the system requires a dictionary of RLCS entries. An RLCS has two levels of description: the first is the language-independent LCS representation of the lexical word (which conforms to (22) above) and the second is the language-specific parametric specification that guides the syntactic realization of the word and its arguments. Consider the RLCS entry for the English verb *stab*:

- (25) $[_{\text{Event}} \text{CAUSE}$
 $([_{\text{Thing}} * \text{W}],$
 $[_{\text{Event}} \text{GO}_{\text{Poss}}$
 $([_{\text{Thing}} \text{KNIFE-WOUND}],$
 $[_{\text{Path}} \text{TOWARD}_{\text{Poss}} ([_{\text{Position}} \text{AT}_{\text{Poss}} ([_{\text{Thing}} \text{KNIFE-WOUND}], [_{\text{Thing}} * \text{Z}]])])]),$
 $[_{\text{WITH}}_{\text{Instr}} * ([_{\text{Event}} * \text{HEAD}*], [_{\text{Thing}} \text{U SHARP-OBJECT}])])]$

The first level of description specifies the language-independent meaning which, in this case, roughly corresponds to “thing W causes thing Z to possess a knife-wound by means of a sharp object U.”¹⁹ Note that the SHARP-OBJECT modifier is included in the definition even though this constituent does not show up in the CLCS for *I stabbed John*:

- (26) $[_{\text{Event}} \text{CAUSE}$
 $([_{\text{Thing}} \text{I}],$
 $[_{\text{Event}} \text{GO}_{\text{Poss}}$
 $([_{\text{Thing}} \text{KNIFE-WOUND}],$
 $[_{\text{Path}} \text{TOWARD}_{\text{Poss}} ([_{\text{Position}} \text{AT}_{\text{Poss}} ([_{\text{Thing}} \text{KNIFE-WOUND}], [_{\text{Thing}} \text{John}])])])])]$

Because modifiers are considered to be optional constituents, they may be omitted from the interlingua. Thus, the same RLCS may be used to build different CLCS representations. In particular, the RLCS provides flexibility in representing the source and target language sentences, which may or may not include certain

¹⁹The variable U refers to a locally introduced modifier. The *HEAD* symbol is a place-holder that points to the root (CAUSE) of the overall *stab* event (*i.e.*, the *stab* event is performed with a sharp object).

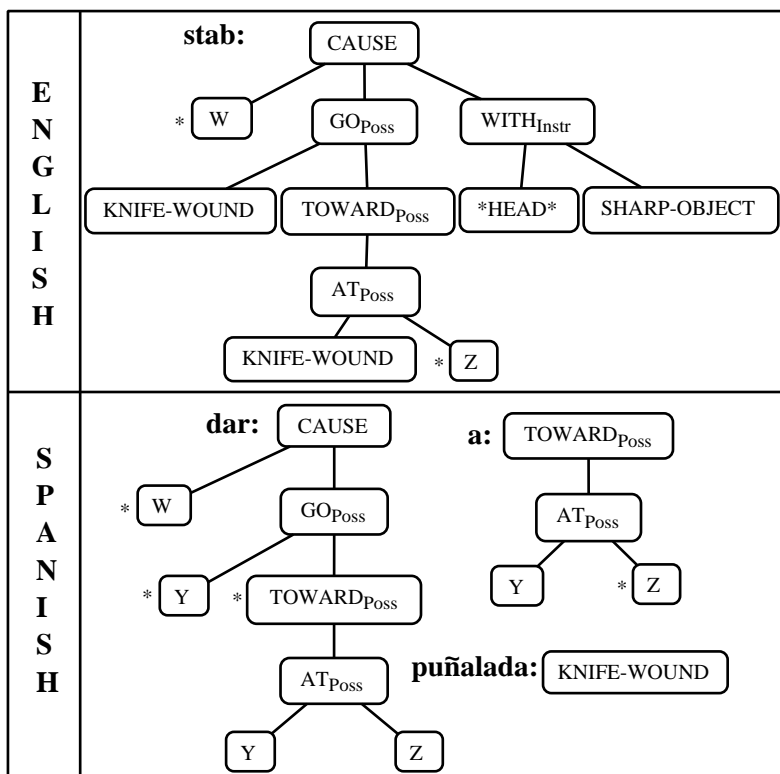


Figure 7: The English and Spanish lexical entries for *stab* event demonstrate the utility of the * marker for pinpointing language-specific distinctions.

modifiers. This means it would be possible to generate either *he stabbed the robber*, or *he stabbed the robber with a knife (scissors, poker, etc.)* from this single lexical entry.

The second level of description in the RLCS is imposed by means of the “*” (pronounced “star”) notation, which is used to specify the language-specific correspondence between LCS arguments and the syntactic structure. Figure 7 illustrates the use of the * notation for the English and Spanish entries corresponding to the *stab* event in example (1) given earlier. In section 5, we will see how the * marker acts as a parameter of variation that allows the system to accommodate this divergence example during the lexical-semantic composition process.

4.3 Organization of the Lexicon

Lexical entries are organized into *LCS classes*. These classes are not identical to the linguistic classes that are typically studied by researchers of the lexicon (*e.g.*, the classes shown in figure 4). Each LCS class is based, not on syntactic distribution and alternation constraints, but on a template that conforms to the well-formedness constraints shown in (17)–(19). Thus, each LCS class may include verbs from more than one linguistic class, and conversely, each linguistic class may include verbs from more than one LCS class. For example, the verb *give* (from the linguistic class of change of ownership) and the verb *stab* (from the linguistic class of contact/effect) are in the same LCS class associated with the GO_{Poss} primitive. Conversely, the linguistic class *directed motion* includes verbs associated with the GO_{Loc} primitive

L = GO_{Poss} Template: [Event GO _{Poss} ([Thing Y], [Path P])]	
C = Causative Template: [Event CAUSE ([Thing W], [Event L])]	
T = Permissive Template: [Event LET ([Thing W], [Event L])]	

P = Path Template: [Path TOWARD _{Poss} ([Position AT _{Poss} ([Thing Y], [Thing Z])])]	
Root Word:	<i>go</i> (* = [Thing Z]; * = [Path P])
	<i>receive</i> (* = [Thing :INT Y]; * = [Thing :EXT Z])
Causative:	<i>stab</i> (W ≠ U, Y, Z; Y = KNIFE-WOUND; * = [Thing W]; * = [Thing Z]; Modifier: [WITH _{Instr} * ([Event *HEAD*], [Thing U SHARP-OBJECT])])
	<i>cut</i> (W = U, or W ≠ U, Y, Z; Y = KNIFE-WOUND; * = [Thing W]; * = [Thing Z]; Modifier: [WITH _{Instr} * ([Event *HEAD*], [Thing U SHARP-OBJECT])])
	<i>give</i> (W ≠ Y, Z; * = [Thing W]; * = [Thing Z]; * = [Path P] or [Thing Z]) ²⁰
	<i>repossess</i> (W = Z; * = [Thing W]; * = [Thing Y])
	<i>obtain</i> (W = Z; * = [Thing W]; * = [Thing Z])
Permissive:	<i>accept</i> (W = Z; * = [Thing W]; * = [Thing Z])

P = Path Template:	
[Path AWAY-FROM _{Poss} ([Position AT _{Poss} ([Thing Y], [Thing Z])])]	
Root Word:	<i>lose</i> (* = [Thing Y]; * = [Thing Z])
Permissive:	<i>relinquish</i> (W = Z; * = [Thing W]; * = [Thing Y])
	<i>surrender</i> (W = Z; * = [Thing W]; * = [Thing Y])

Figure 8: The lexical entries for English verbs in the GO_{Poss} LCS class are characterized by an LCS template L and a path pointer P. The two path types shown here are TOWARD_{Poss} (e.g., *go*, *receive*, *stab*, *cut*, *give*, *repossess*, *obtain*, and *accept*) and AWAY-FROM_{Poss} (e.g., *lose*, *relinquish*, and *surrender*).

(e.g., *enter*) and verbs associated with the GO_{Poss} primitive (e.g., *receive*).

All lexical items belonging to an LCS class are grouped together in the dictionary according to a particular LCS template. For example, the lexical entries defined within the GO_{Poss} class are shown in figure 8. (For brevity, a small subset of the verbs from this class is shown.) In this example, L is the LCS template and the words defined in this class are grouped into different path types by means of the path pointer P. The two path types shown here are TOWARD_{Poss} (e.g., *go*, *receive*, *stab*, *cut*, *give*, *repossess*, *obtain*, and *accept*) and AWAY-FROM_{Poss} (e.g., *lose*, *relinquish*, and *surrender*). Note that each lexical entry includes the surface realization information (specified by means of the “* =” notation).²¹

In addition to specifying the LCS and path templates, L and P, the LCS class specifies a causative template, C, and a permissive template, T. These templates may be used compositionally with the LCS template L to form more complex LCS constructions. For example, the representation for the words *go* and *receive* has a causative form that corresponds to the words *stab*, *cut*, *give*, *repossess*, and *obtain*

²⁰The disjoint * specification used in the entry for *give* allows the verb to be realized both in the dative (e.g., *I gave him the book*) and in the non-dative (e.g., *I gave the book to him*). The special use of the * notation will be discussed in section 6.

²¹In addition to the * marker, the lexicon also makes use of the :INT and :EXT markers (e.g., in the lexical entry for *receive*). The :INT and :EXT markers are not relevant to the focus of this paper, but see Dorr (in press) for more details.

and a permissive form that corresponds to the word *accept*. Causative and permissive entries specify coreference information about the logical subject (*i.e.*, the variable W in the current example). For example, in the case of *give*, the logical subject is not coreferential with any other arguments (as in *I gave John the gift*). In contrast, *repossess* requires the logical subject to be coreferential with the recipient Z since the subject of the event is also the recipient of the possessional transfer (as in *I repossessed John's car*).²² Note that, in the case of *stab* and *cut*, the subject is non-coreferential (*I stabbed/cut John*), but in the case of *cut*, the subject may also be coreferential with the instrument (*the knife cut John*).²³ For illustrative purposes, figure 9 shows some examples of causative trees implicitly represented by the specification given in the GO_{POSS} LCS class. For reasons of readability, we will continue to use the tree-like representations rather than the LCS template format used for lexical entries.

Optional modifiers are not included in the lexical entries unless they provide immutable information that is idiosyncratic to a particular verb (such as the SHARP-OBJECT modifier included in *cut* and *stab* entries). In general, modifiers are available through an inheritance mechanism based on the following field specifications:

- (27) (i) **Events and States:** Inherit modifiers from the Intentional, Instrumental, Temporal, and Locational Fields.
- (ii) **Things:** Inherit modifiers from the Identificational, Possessional, Temporal, and Locational Fields.

For example, the verb *stab* is an Event; thus, it automatically inherits optional modifiers in the Intentional field (*e.g.*, *I stabbed John for Harry*), Instrumental field (*e.g.*, *I stabbed John with a knife*), Temporal field (*e.g.*, *I stabbed John at 9:00*), and Locational field (*e.g.*, *I stabbed John in the street*). Note that the use of LCS fields rather than LCS types provides a more flexible means of specifying LCS modifiers. For example, the noun *book* may be modified by a number of different types in the possessional field including Position (*e.g.*, *the book of John's*), State (*e.g.*, *John's book*), Event (*e.g.*, *the book that John owns*), *etc.*

The system of LCS classes and modifier inheritance offers a number of benefits. First, arguments and modifiers need not be stated for every word in the dictionary because lexical entries automatically inherit arguments from the LCS template and modifiers from the field specifications shown in (27). This eliminates proliferation of redundancy across lexical entries. Second, the LCS classes allow causative and permissive forms to specify their idiosyncratic argument information by means of coreference and other types of argument indexation. Thus, words with related meanings may be grouped together despite different argument specifications; this avoids having to define these words independently in unrelated lexical entries. Finally, the LCS classes allow syntactically distinct lexical items to be defined in the same class

²²The LCS representation used for the word *repossess* might also be used for the word *steal*. In order to distinguish between these, the entry for *steal* would require an additional modifier, STEALINGLY.

²³Note that the coreference here refers to the values of the arguments as defined in the lexical entry. It does not refer to cases of Binding coreference in the surface structure (*i.e.*, coindexation of instantiated arguments). In particular, the lexical entries do not rule out cases such as *I stabbed myself* where the subject and object become coindexed during syntactic processing.

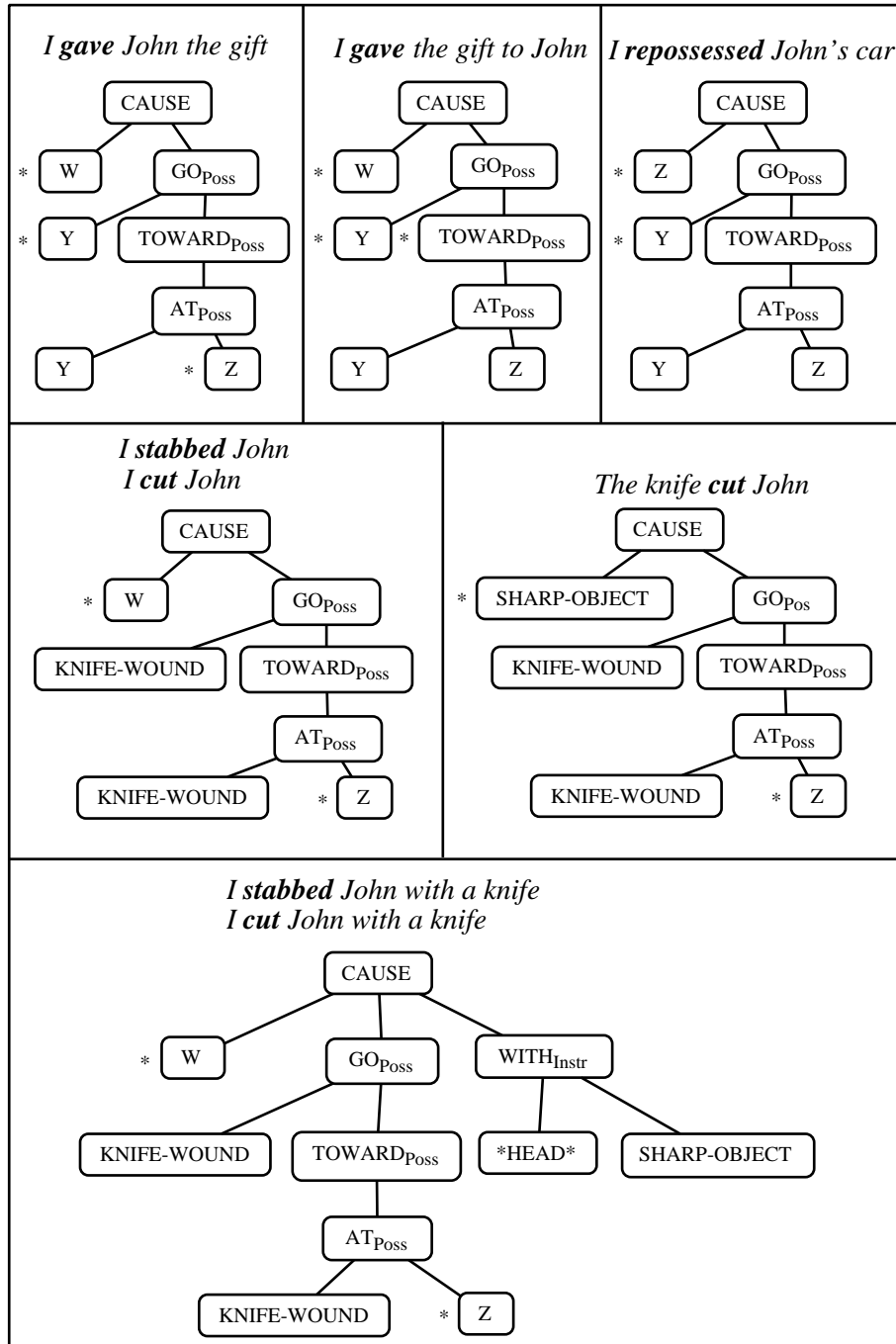


Figure 9: A number of different causative trees are implicitly represented by GO_{Poss} class.

if they have the same word sense. For example, the deverbalized noun *entrance* is defined in the same LCS class (GO_{Loc}) as its verbal counterpart *enter* even though these two words are not of the same syntactic category.

4.4 Extensions to Jackendoff's LCS Representation

The use of the LCS as an interlingua for more than one language is generally considered to be nonstandard, and in fact, is a possibility that is not even discussed by Jackendoff himself. Because the lexicon has been parameterized to account for lexical-semantic divergences, this possibility has been realized in the UNITRAN system with a minimal amount of engineering. One of the reasons that the latest version of the LCS framework by Jackendoff (1990) is not well-suited to an interlingual specification of more than one language is that it retains the language-specific syntactic and semantic subcategorization information in the frame of each lexical entry rather than deriving syntactic properties via linking rules. Other proposals have been made for using a lexical-semantic representation that does not retain syntactic information in lexical entries (see, *e.g.*, Pinker (1989), Rappaport and Levin (1988), and White (1992)), but so far, most approaches have concentrated on representing English and have not entertained the possibility of parameterizing the representation for use as an interlingua.

The primary extension to Jackendoff's original framework is the association of the representation with an algorithm for recursive composition and decomposition of the interlingual form. The algorithm is defined on the basis of a systematic linking of the LCS to the syntactic structure, both during parsing as well as during generation. The most critical mechanism that brings about this linking is the * marker described in section 4.2. We will temporarily defer the discussion about the linking process until section 5.

We turn now to a discussion about the primitives that have been added to the framework beyond those shown in figure 3. In particular, the primitives CAUSE-EXCHANGE, DO, EAT, SEE, HEAR, FEEL, SEARCH, UP, DOWN, ACROSS, and ALONG have been added. In addition to new primitives, the fields Perceptual, Intentional, and Instrumental have been added to the original framework. We will see that the addition of these new primitives and fields has induced further extensions to the Manner type. In addition, the ORIENT primitive has been extended to include the Identificational and Temporal fields. Also, a new type called Intensifier has been added. Finally, a number of primitives corresponding to question words have been added. All of these extensions came about as a result of an in-depth investigation of translation divergences. In order to demonstrate the effectiveness of the current approach, a number of verbs needed to be added that were not adequately characterized by Jackendoff's primitives. Thus, the original framework needed to be augmented. We will briefly sketch the general intuition behind these extensions.

The CAUSE-EXCHANGE primitive has been added to include verbs such as *buy* and *sell*. This primitive is a shorthand notation for the CAUSE primitive coupled with a subordinating function called EXCH as presented by Jackendoff (1990). The representation distinguishes between *buy* and *sell* by using a different causative agent (*i.e.*, logical subject) in each case:

- (28) (i) John sold the book to Mary
(ii) $[_{\text{Event}} \text{ CAUSE-EXCHANGE}$
 $([_{\text{Thing}} \text{ JOHN}],$
 $[_{\text{Event}} \text{ GO}_{\text{POSS}}$
 $([_{\text{Thing}} \text{ BOOK}],$
 $[_{\text{Path}} \text{ FROM}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ BOOK}], [_{\text{Thing}} \text{ JOHN}]])]),$
 $[_{\text{Path}} \text{ TO}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ BOOK}], [_{\text{Thing}} \text{ MARY}]])]),$
 $[_{\text{Event}} \text{ GO}_{\text{POSS}}$
 $([_{\text{Thing}} \text{ MONEY}],$
 $[_{\text{Path}} \text{ FROM}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ MONEY}], [_{\text{Thing}} \text{ MARY}]])]),$
 $[_{\text{Path}} \text{ TO}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ MONEY}], [_{\text{Thing}} \text{ JOHN}]])])])])])$
- (29) (i) Mary bought the book from John
(ii) $[_{\text{Event}} \text{ CAUSE-EXCHANGE}$
 $([_{\text{Thing}} \text{ MARY}],$
 $[_{\text{Event}} \text{ GO}_{\text{POSS}}$
 $([_{\text{Thing}} \text{ BOOK}],$
 $[_{\text{Path}} \text{ FROM}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ BOOK}], [_{\text{Thing}} \text{ JOHN}]])]),$
 $[_{\text{Path}} \text{ TO}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ BOOK}], [_{\text{Thing}} \text{ MARY}]])]),$
 $[_{\text{Event}} \text{ GO}_{\text{POSS}}$
 $([_{\text{Thing}} \text{ MONEY}],$
 $[_{\text{Path}} \text{ FROM}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ MONEY}], [_{\text{Thing}} \text{ MARY}]])]),$
 $[_{\text{Path}} \text{ TO}_{\text{POSS}} ([_{\text{Position}} \text{ AT}_{\text{POSS}} ([_{\text{Thing}} \text{ MONEY}], [_{\text{Thing}} \text{ JOHN}]])])])])$

The primitive DO has been added to cover the English, Spanish, and German verbs *do*, *hacer*, and *tun*, respectively. Thus, the representation for *John did the wash* is:

- (30) $[_{\text{Event}} \text{ DO} ([_{\text{Thing}} \text{ John}], [_{\text{Event}} \text{ WASH} ([_{\text{Thing}} \text{ John}], [_{\text{Thing}} \text{ Y}])])]$ ²⁴

It might be argued that the verb *do* should not be represented at the level of conceptual structure. That is, instead of using the representation in (30) for the sentence *John did the wash*, we could simply represent this sentence in the same way that we would represent *John washed*:

- (31) $[_{\text{Event}} \text{ WASH} ([_{\text{Thing}} \text{ John}], [_{\text{Thing}} \text{ Y}])]$

However, without the DO primitive, these two sentences would be translated interchangeably, which would not be appropriate in all contexts. Furthermore, without such a primitive there would be no way of representing a sentence in which the main action is not stated, such as *John did it*. Finally, the DO primitive is completely analogous to the CAUSE and LET primitives which are independently motivated by Jackendoff (1983). To the extent that these primitives can be justified in an interlingual system, so too can the DO primitive. This primitive is analogous to CAUSE and LET in three ways: (1) it does not have a field dimension; (2) it takes the same number and type of arguments (*i.e.*, a Thing and an Event or State); and (3) it is optionally realizable on the surface depending on the requirements of the

²⁴Note that the logical Thing argument is left uninstantiated since the object of the washing action is not stated.

language. Regarding this final point, note that the same types of cross-linguistic realization alternatives are available for CAUSE and DO:

- (32) (i) S: John forzó la entrada al cuarto (CAUSE realized as *forzar*)
 E: John broke into the room (CAUSE not realized)
- (ii) E: John did the wash yesterday (DO realized as *do*)
 S: Juan lavaba ayer (DO not realized)

Another primitive that has been added to the system is EAT because no such action was included in the original Jackendoff framework.²⁵ The EAT primitive does not have a field dimension given that no field seems to apply (with the possible exception of the Locational field). Thus, it is assumed that this primitive already has its field incorporated into it. This incorporated field, whatever it is called, is likely to be the same one that DRINK and SMOKE use. (Note that these primitives are more specific than those that can be moved around from field to field, such as GO and BE.) An example of the use of the EAT primitive is the following:

- (33) (i) John ate beans
 (ii) [_{Event} EAT ([_{Thing} JOHN], [_{Thing} BEANS])]

The Perceptual field has been created in order to accommodate verbs such as *see*, *hear*, and *feel*.²⁶ However, unlike other Perceptual verbs (to be described next), these verbs are characterized by the new primitives SEE, HEAR, and FEEL. The reason these new primitives are adopted instead of using the BE primitive (*e.g.*, BE_{Perc} SEEINGLY) is that these verbs are considered to be events, not states. A syntactic test, attributed to Dowty (1979), that teases apart events and states is the progressive construction:

- (34) (i) I was hearing the sirens as we danced. (event)
 (ii) *I was wanting him to talk softly as we danced. (state)

This Perceptual field is also used with the primitive BE to represent verbs such as *believe*, *know*, *want*, *etc.*, *i.e.*, it is used to characterize verbs of *mental* perception as well as *visual*, *aural*, and *tactile* perception. When the BE primitive is used in the Perceptual field, it is generally used in conjunction with a manner component such as BELIEVINGLY, KNOWINGLY, WANTINGLY, *etc.* in order to distinguish between the different types of perceived notions. (The manner component will be discussed shortly.) The Perceptual field has also been shown to be useful for paths: TOWARD_{Perc} (*e.g.*, *look toward*), AWAY-FROM_{Perc} (*e.g.*, *look away from*), ABOUT_{Perc} (*e.g.*, *know about*), and IN_{Perc} (*e.g.*, *believe in*). An example of the use of BE in the Perceptual field is the following:

- (35) (i) John believes in unicorns

²⁵Note that this primitive might be better formulated as the more general INGEST primitive (as in that of Schank (1972, 1973, 1975)) where the ingested object is FOOD (for *eat*), FLUID (for *drink*), SMOKE (for *smoke*), *etc.*, but it has been left as the more specific formulation in the current implementation.

²⁶There are surely others (*e.g.*, *smell*, *taste*, *etc.*) that can be added to this set.

- (ii)
$$\begin{aligned} & [\text{Event BE}_{\text{Perc}} \\ & \quad ([\text{Thing JOHN}], \\ & \quad [\text{Path IN}_{\text{Perc}} \\ & \quad \quad ([\text{Position AT}_{\text{Perc}} ([\text{Thing JOHN}], [\text{Thing UNICORNS}]])]) \\ & \quad [\text{Manner BELIEVINGLY}]] \end{aligned}$$

Instead of using a Perceptual field, Jackendoff represents mental notions by means of the Rep operator which provides a referential reading of something that might occur in someone’s mind. Thus, sentence (35)(i) might be represented as:

- (36)
$$[\text{State BE} ([\text{Rep UNICORN}], [\text{Position IN} ([\text{Location JOHN'S MIND}]])])]$$

There are two problems with using this representation. The first is that it relies on an operator that exists at a meta-level above the conceptual representation (*i.e.*, the Rep operator). If this meta-level were used in UNITRAN, another tier of lexical-semantic processing would be required for translation since the conceptual representation currently adopted is intended to be mapped uniformly to the syntactic structure. (We will see this in section 5.) The additional tier would greatly complicate the mapping between the interlingua and the syntactic structure. The second problem, as acknowledged by Jackendoff, is that this representation does not differentiate the thematic analyses of *believe*, *imagine*, *remember*, and so forth. In the simplified representation used by UNITRAN, these notions are differentiated by means of a manner component.²⁷ This representation has the advantage that it is mapped to the syntactic structure uniformly by the same mapping that is used for non-mental verbs.

Another new primitive, SEARCH, has been added to represent the notion of *searching* or *looking for* something. This primitive is used only in the Possessional field and is generally used in conjunction with the new position primitive, FOR:²⁸

- (37) (i) John searched / looked for the book
(ii)
$$\begin{aligned} & [\text{Event SEARCH}_{\text{Poss}} \\ & \quad ([\text{Thing JOHN}], \\ & \quad [\text{Path FOR}_{\text{Poss}} ([\text{Thing JOHN}], [\text{Thing BOOK}])]) \end{aligned}$$

The new paths UP, DOWN, and ALONG have been added to the system in order to handle sentences such as the following:

- (38) (i) John went up/down the stairs (UP_{Loc}/DOWN_{Loc})

²⁷The differentiation of mental states is also a problem in Schank’s system (see Schank (1972, 1973, 1975)) which relies on the same primitive (MTRANS) for a number of different verbs such as *want*, *know*, *think*, *believe*, *etc.* Later versions of Schank’s model (*e.g.*, Rieger (1975)) extended the framework to include primitives such as WANT. In effect, this is a close approximation to the solution adopted here.

²⁸An alternative to defining the new SEARCH primitive is to define *search* and *look for* as a causative form of the GO_{Poss} primitive:

$$\begin{aligned} & [\text{Event CAUSE} \\ & \quad ([\text{Thing W}], \\ & \quad [\text{Event GO}_{\text{Poss}} \\ & \quad \quad ([\text{Thing W}], [\text{Path TOWARD}_{\text{Poss}} ([\text{Position AT}_{\text{Poss}} ([\text{Thing W}], [\text{Thing Z}])])])]) \\ & \quad [\text{Manner SEARCHINGLY}]] \end{aligned}$$

The SEARCH primitive can be thought of as an abbreviation for this longer expression.

- (ii) John walked along the river (ALONG_{Loc})

Note that these paths are used only in the Locational field.²⁹

Another modification that has been made to the primitives is the use of the ORIENT primitive in the identificational and temporal fields to account for the following cases:

- (39) (i) The book costs \$10.00 (ORIENT_{Ident})
- (ii) John aims to start at 2:00 (ORIENT_{Temp})

In addition to new primitives, two more fields have been added to the system: Instrumental and Intentional. The Instrumental field is used to represent various types of instrumental modifiers (as we will see below in the lexical entry for *stab*). The primitives that are used in this field, FOR, WITH, and CO, are also new:

- (40) (i) John bought the book for \$5.00 (FOR_{Instr})
- (ii) John stabbed Mary with a knife (WITH_{Instr})
- (iii) John walked with Mary (CO_{Instr})

The Intentional field is used for cases where an action is done for the purpose of something or for the benefit of someone. The two primitives used in this field are FOR and AT:

- (41) (i) John signed the book for Mary (FOR_{Intent})
- (ii) John ate because he was hungry (FOR_{Intent})
- (iii) John turned the light on so that he could see (FOR_{Intent})
- (iv) The book cost me \$10.00 (AT_{Intent})

The justification for adding these fields is that there are a number of conjunctions and particles (such as *so that*, *because*, *for*, *with*, *etc.*) that must be made distinguishable in order for the system to make the appropriate lexical selection during generation. For example, the word *for* (FOR_{Poss}), not *because* (FOR_{Intent}), is used to represent the modifier phrase *for John* in the sentence *John bought the book for Mary*.

Another extension to the Jackendoff framework is the augmentation of the Manner type to include a wide range of primitives. While the LCS framework currently provides a means for distinguishing between verbs *across* LCS classes, it does not yet provide a principled account of constraints *within* LCS classes. Verbs within a particular class are frequently distinguishable by some feature corresponding to Manner (*e.g.*, *walk vs. run*), yet Jackendoff says very little about the function of the Manner type in his description of conceptual structure. In fact, he claims that “it is not the business of conceptual structure” to encode different Manner types since conceptual structure is designed “to encode primarily an appropriate argument structure . . .”; the lexicon, then, must be linked to “a more detailed spatial structure encoding” in order to distinguish between verbs in a particular verb class (*e.g.*, *wriggle* and *wiggle*) (see Jackendoff (1990, p. 88)). Although I whole-heartedly

²⁹It might also be possible to use these paths in other fields such as the Perceptual field *e.g.*, *John looked up* (UP_{Perc}).

agree with this view, we are left with the question of what is meant by a “detailed spatial structure encoding,” *i.e.*, how such a representation would be constructed and how it interacts with the conceptual structure. Moreover, this approach does not address the problem of representing the Manner type in other categories outside of the spatial domain (*e.g.*, *like vs. love* or *repossess vs. steal*).

The current solution has been to use the Manner component to distinguish between two verbs that fall within the same linguistic class when no other distinguishing features are available. A side effect of this solution is that the representation appears to be too specific in certain cases. While Manner primitives are no more specific than many other open-ended primitives (*e.g.*, HOUSE, ASLEEP, 9:00, HERE, *etc.*), one could conceive of a worst case scenario in which the number of Manner primitives grows linearly with respect to the number of verbs that are added to each verb class. It could even be argued that the work of categorizing predicates such as *believe* and *want* into the BE_{Perc} class is wasted since the use of BELIEVINGLY and WANTINGLY is equivalent to categorizing the predicates into two distinct classes (*i.e.*, we could just use primitives called BELIEVE and WANT and get rid of the BE_{Perc} primitive).

These concerns have recently been addressed by Siskind (1992) who argues that the Manner component of certain types of verbs can be broken down into primitives of a different type. For example, Siskind is able to differentiate among such spatial verbs as *turn*, *spin*, *revolve*, and *rotate* simply by using the GO_{Loc} primitive with a Manner component that is defined in terms of physical notions such as “orientation at the end of the action with respect to the starting point.” (The formal mechanism for capturing this physical notion is not discussed here.) It is expected that the same principle carries over to non-spatial classes of predicates such as *believe* and *want*. Clearly this is an area that requires further investigation, but while there are still a number of open questions, the current solution has proven to be sufficient for addressing the problem of translation divergences. That is, Manner components such as BELIEVINGLY and WANTINGLY can be taken to be “macros” that may be expanded into some form of primitive notion analogous to the physical descriptions used in Siskind’s work.

Note that there are a number of alternatives for representing the Manner component. For example, the Manner component might be represented by means of feature values (*i.e.*, \pm believe, \pm want, *etc.*) that are independent from the lexical representation. A number of researchers have taken a feature-setting approach to distinguish among verbs that belong to the same lexical class. (See, *e.g.*, Bennett *et al.* (1990).) The reason for adopting the approach described here is that the Manner component often shows up in the surface syntactic structures of different languages (*e.g.*, German *gern* corresponds to LIKINGLY, French *a la nage* corresponds to SWIMMINGLY, *etc.*). In order to provide a uniform mapping between LCS representations and their surface realizations, we need a lexical place holder for Manner components as well as argument components. Although these lexical tokens are not “reusable” in the sense of the more general tokens (GO, BE, TO, AT, *etc.*), it is expected that these tokens may ultimately be decomposed into primitives that *are* reusable at a different level of representation (as discussed above). In any case, these tokens are necessary to support the general scheme of LCS composition

and decomposition of surface structures that include a corresponding Manner component of meaning. The mapping between the LCS representation and the surface syntactic structure is set up so that it is easy to determine whether the Manner component is suppressed or overtly realized in the surface sentence.

One final comment in defense of the “macro” version of the Manner component is that it does avoid some of the well-known problems of extreme decomposition. (For a detailed discussion, see *e.g.*, Sproat (1985).) One such problem is the potential for deep recursion that is inherent in the CD framework of Schank. As noted by Schank himself (1973, p. 201), this is particularly a problem with instrumentality:

- (42) “If every ACT requires an instrumental case which itself contains an ACT, it should be obvious that we can never finish diagramming a given conceptualization. For [the] sentence [John ate the ice cream with a spoon], for example, we might have ‘John ingested the ice cream by transing the ice cream on a spoon to his mouth, by transing the spoon to the ice cream, by grasping the spoon, by moving his hand to the spoon, by moving his hand muscles, by thinking about moving his hand muscles,’ and so on . . . These instrumental actions are not really needed and are rarely actively thought about . . . but we shall retain the ability to retrieve these instruments should we find this necessary.”

While the current approach does not pretend to solve the problem of representing the deep meaning of the Manner component, it does at least avoid the infinite recursion problem by leaving out the detailed mechanics underlying the modifying action (*e.g.*, that SWIMMINGLY involves certain leg and arm motions that rely on moving certain muscles, *etc.*).

In addition to new fields and primitives, the system also has a new type, Intensifier, that currently allows Properties, Manners, and Intensifiers themselves to be intensified. The intensifier is placed in the Identificational field for modifying properties and the Instrumental field for modifying manners:

- (43) (i) John was very happy
(ii) [State BE_{Ident}
([Thing JOHN],
[Position AT_{Ident}
([Thing JOHN],
[Property HAPPY ([Intensifier VERY_{Ident}])])])]
- (44) (i) John ate very happily
(ii) [Event EAT
([Thing JOHN],
[Thing FOOD],
[Manner HAPPILY ([Intensifier VERY_{Instr}])])]
- (45) (i) John was so very happy
(ii) [State BE_{Ident}
([Thing JOHN],
[Position AT_{Ident}
([ctype:Thing JOHN],
[Property HAPPY ([Intensifier SO_{Ident} ([Intensifier VERY_{Ident}])])])])]

- (46) (i) John ate so very happily
(ii) [Event EAT
([Thing JOHN],
[Thing FOOD],
[Manner HAPPILY ([Intensifier SO_{Instr} ([Intensifier VERY_{Instr}]))]])]

A final augmentation to the set of primitives is the addition of the capability to handle question words referring to Intensifiers, Manners, Things, Locations, Times, Properties, and Purposes.³⁰ The relevant primitives are: WH-INTENSIFIER (*e.g.*, *how*), WH-MANNER (*e.g.*, *how*), WH-THING (*e.g.*, *what*), WH-LOCATION (*e.g.*, *where*), WH-TIME (*e.g.*, *when*), WH-PROPERTY (*e.g.*, *how*), and WH-PURPOSE (*e.g.*, *why*). For example, the WH-THING primitive would be used as follows:

- (47) (i) What did John eat
(ii) [Event EAT ([Thing JOHN], [Thing WH-THING])]

To sum up, the primitive/field combinations that have been added to the system are shown in figure 10 along with relevant examples.³¹ The full extended set of primitives is shown in figure 11. There are 33 primitives multiplied out into 9 fields, not including the more open-ended types (*i.e.*, Manners, Things, Locations, Times, Properties, and Purposes). Although the system currently uses a small set of lexical-semantic primitives, the set is quite adequate for defining a wide variety of words due to the compositional nature of LCS's. The types of sentences handled by the extended set of primitives are given in appendix A.

Now that we have examined the lexical-semantic representation in more detail, we turn to the mapping between this representation and the syntactic structure.

5 Mapping Between the Interlingua and the Syntactic Structure

In order to relate source- and target-language structures to an interlingual form, lexical entries must specify certain language-specific syntactic information. This is the nature of the second level of lexical description alluded to in section 4.2. An example of a mechanism that is used by this level is the * marker, which is required for every explicitly realized argument and modifier in the lexical entries for each language.³² This marker is used as a means of parameterization in the lexicon in that it specifies the LCS positions that have corresponding syntactic realizations in a particular language.

³⁰The Purpose type is not included in Jackendoff's description of the model; in the current model, the Purpose is intended to be the reference object of the Intentional field, just as the Location is considered to be the reference object of the Locational field.

³¹This list is by no means exhaustive. In particular, the primitives corresponding to Manners, Things, Locations, Times, Properties, and Purposes are too numerous to list here given that these categories are intended to be open-ended. Note that open-ended categories correspond precisely to those primitives that do not have field specifications, *i.e.*, the taking primitives that do not take any arguments.

³²There are also other lexical markers that will not be discussed here. These are presented in detail in Dorr (in press).

Type	Primitive-Field	Example	
Event	CAUSE-EXCHANGE	John sold the book to Mary	
	DO	John did the wash	
	EAT	John ate breakfast	
	SEE _{Perc}	John saw the rain	
	HEAR _{Perc}	John heard the rain	
	FEEL _{Perc}	John felt the rain	
	SEARCH _{Poss}	John looked for the book	
	ORIENT _{Ident}	The book costs \$10.00	
State	ORIENT _{Temp}	John aims to start at 2:00	
	BE _{Perc}	John believed Mary	
Path	TOWARD _{Perc}	John looked toward the sun	
	AWAY-FROM _{Perc}	John looked away from the sun	
	ABOUT _{Perc}	John knew about Mary	
	IN _{Perc}	John believed in Mary	
	UP _{Loc}	John went up the stairs	
	DOWN _{Loc}	John went down the stairs	
Position	ALONG _{Loc}	John walked along the river	
	FOR _{Instr}	John bought the book for \$5.00	
	WITH _{Instr}	John stabbed Mary with a knife	
	CO _{Instr}	John walked with Mary	
	FOR _{Intent}	John signed the book for Mary	
Intensifier	AT _{Intent}	The book cost me \$10.00	
	FOR _{Poss}	John bought the book for Mary	
	VERY _{Ident}	John was very happy	
	VERY _{Instr}	John ate very happily	
	SO _{Ident}	John was so happy	
	SO _{Instr}	John ate so happily	
Manner	WH-INTENSIFIER _{Ident}	How happy was John	
	WH-INTENSIFIER _{Instr}	How happily did John eat	
	BELIEVINGLY	John believed Mary	
	WANTINGLY	John wanted Mary	
	KNOWINGLY	John knew Mary	
	READINGLY	John read the book	
Thing	WRITINGLY	John wrote the book	
	WH-MANNER	How did John write the book	
	WH-THING	What did John eat	
	Location	WH-LOCATION	Where did John go
	Time	WH-TIME	When did John write the book
	Property	WH-PROPERTY	How did John feel
	Purpose	WH-PURPOSE	Why did John write the book

Figure 10: The extended set of fields and primitives includes types from the following categories: Events, States, Paths, Positions, Intensifiers, Manners, Things, Times, Locations, Properties, and Purposes.

Type	Primitive
Event	CAUSE, LET, GO, STAY, CAUSE-EXCHANGE, DO, EAT, SEE, HEAR, FEEL, SEARCH
State	BE, GO-EXT, ORIENT
Path	TO, FROM, TOWARD, AWAY-FROM, VIA, ABOUT, IN, UP, DOWN, ALONG
Position	AT, IN, ON, WITH, CO, FOR
Intensifier	VERY, SO, WH-INTENSIFIER
Manner	BELIEVINGLY, WANTINGLY, KNOWINGLY, READINGLY, WRITINGLY, HAPPILY, WH-MANNER
Thing	BOOK, PERSON, REFERENT, KNIFE-WOUND, WH-THING
Location	HERE, THERE, LEFT, RIGHT, UP, DOWN, WH-LOCATION
Time	TODAY, SATURDAY, 2:00, 4:00, WH-TIME
Property	TIRED, HUNGRY, PLEASED, BROKEN, ASLEEP, DEAD, HAPPY, WH-PROPERTY
Purpose	WH-PURPOSE

Figure 11: The extended set of primitives includes 33 primitives multiplied out into 9 fields, not including the more open-ended types (*i.e.*, Manners, Things, Locations, Times, Properties, and Purposes).

Earlier in figure 7, we saw that the * notation is used in the English and Spanish lexical entries for the *stab* event. Note that the * marker makes two crucial language-specific distinctions: (1) the TO_{Poss} portion of the RLCS is *-marked for *dar* but not for *stab*; and (2) the KNIFE-WOUND argument is not *-marked for *stab* whereas the corresponding Y position is *-marked for *dar*. These parametric distinctions force the representation for the word *dar* to be combined with the representations of other words such that the composite representation is equivalent to the single lexical structure specified in the entry for *stab*. In particular, the TO_{Poss} portion of the lexical entry for *stab* is forced to unify with the lexical entry for the word *a* and the Y position must be filled in with the LCS for *puñalada*. Note that the lexical entry for *a* also has its own * marker, which must be filled in recursively by means of the same mechanism.

The * marker is used in conjunction with a generalized linking routine that specifies the mapping between the syntactic structure and the lexical-semantic representation. The routine assumes the following structural correspondence:³³

$$\begin{aligned}
 (48) \quad & [_{Y-MAX} \\
 & \quad Q-MAX_{j+1} \dots Q-MAX_k \\
 & \quad [_{Y-MAX} \\
 & \quad \quad W-MAX \\
 & \quad \quad [_{X-MAX} [x \quad Q-MAX_1 \dots Q-MAX_i \quad X \quad Q-MAX_{i+1} \dots Q-MAX_j] \\
 & \quad \quad \quad Z-MAX_1 \dots Z-MAX_n]] \\
 & \quad Q-MAX_{k+1} \dots Q-MAX_m] \\
 & \text{corresponds to:}
 \end{aligned}$$

³³The reader is referred to Dorr (in press) for details concerning the linguistic theory underlying the syntactic structure (*i.e.*, \bar{X} Theory).

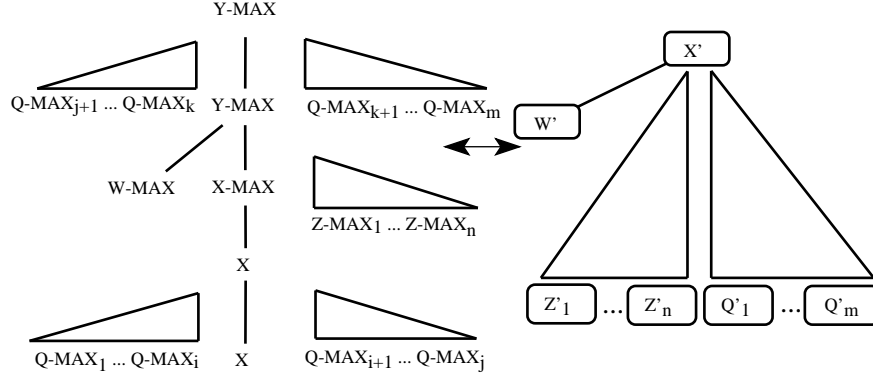


Figure 12: The structure of the LCS lends itself to a systematic specification of the correspondence between syntactic structure and lexical-semantic structure.

$$\begin{aligned}
 & [_{T(X')} X' \\
 & \quad ([_{T(W')} W'], \\
 & \quad [_{T(Z'_1)} Z'_1], \dots, [_{T(Z'_n)} Z'_n], \\
 & \quad [_{T(Q'_1)} Q'_1], \dots, [_{T(Q'_m)} Q'_m])]
 \end{aligned}$$

where:

1. X is the *head* of the X-MAX phrase.
2. W-MAX is the *external argument* (or subject) of X.
3. Z-MAX₁ ... Z-MAX_n are the *internal arguments* (or objects) of X.
4. Q-MAX₁ ... Q-MAX_j are the *minimal adjuncts* of X and Q-MAX_{j+1} ... Q-MAX_m are the *maximal adjuncts* of X.
5. X' and T(X') are, respectively, the primitive-field and type corresponding to the syntactic constituent X. (X' is the *logical head*.)
6. W' and T(W') are, respectively, the primitive-field and type corresponding to the syntactic constituent W-MAX. (W' is the *logical subject*.)
7. Z'₁, ..., Z'_n and T(Z'₁), ..., T(Z'_n) are, respectively, the primitive-fields and types corresponding to the syntactic constituents Z-MAX₁ ... Z-MAX_n. (Z'₁, ..., Z'_n are the *logical arguments*.)
8. Q'₁, ..., Q'_m and T(Q'₁), ..., T(Q'_m) are, respectively, the primitive-fields and types corresponding to the syntactic constituents Q-MAX₁ ... Q-MAX_m. (Q'₁, ..., Q'_m are the *logical modifiers*.)

Figure 12 shows the tree-like structures corresponding to (48).

Specifically, the generalized linking routine is defined as follows:

(49) **Generalized Linking Routine:**

- a. Relate the syntactic head X to the logical head position X'.
- b. Relate the syntactically external position W to the logical subject position W'.
- c. Relate the syntactically internal positions Z₁, ..., Z_n to the logical argument positions Z'₁, ..., Z'_n.
- d. Relate the syntactic adjunct positions Q₁, ..., Q_m to the logical modifier positions Q'₁, ..., Q'_m.

Compose_LCS (X-MAX)

- 1 Let X = head of X-MAX.
- 2 Let X' = RLCS corresponding to head X .
- 3 Let W = external argument of X-MAX (*i.e.*, phrase outside of X-1).
- 4 Let $Z_1 \dots Z_n$ = internal arguments of X-MAX (*i.e.*, phrases inside of X-1).
- 5 Let $Q_1 \dots Q_m$ = adjuncts of X-MAX (*i.e.*, phrases adjoined to X-MAX or X).
- 6 For $i \in \{W, Z_1 \dots Z_n, Q_1 \dots Q_m\}$
 - 6.a Determine highest *-marked RLCS position i' in X' that corresponds to i using the generalized linking routine (49).
 - 6.b Let $L = \mathbf{Compose_LCS}(i)$.
 - 6.c Unify L with position i' in X' .
- 7 Return X' .

Figure 13: The algorithm for LCS composition takes a syntactic tree and composes the interlingua on the basis of unification with *-marked positions.

This routine is used in both directions, *i.e.* composition of the interlingua and decomposition of the interlingua. Only the first of these two directions will be discussed for the remainder of this section; this should be sufficient to illustrate the general nature of the process.

5.1 Lexical-Semantic Composition

Lexical-semantic composition (henceforth LCS composition) is performed by the second module of the lexical-semantic component shown in figure 2. We will introduce the algorithm for LCS composition and illustrate how this procedure works for our translation example repeated here for convenience:

(50) E: I stabbed John \Leftrightarrow S: Yo le di puñaladas a Juan
'I gave knife-wounds to John'

5.1.1 Algorithm for LCS Composition

Figure 13 shows the top-level LCS-composition procedure, **Compose_LCS**. This procedure takes a syntactic tree whose top phrasal node is X-MAX and finds the RLCS associated with the lexical item X (steps 1 and 2). It then extracts the arguments and adjuncts of X-MAX, if there are any (steps 3–5), and recursively passes these to the same procedure (step 6). Finally, it returns the CLCS X' (*i.e.*, the instantiated RLCS) corresponding to the head of the top phrasal node (step 7).

Step 6 requires further elaboration. First, step 6.a uses the linking routine to determine the mapping between the positions in the syntactic structure and the positions in the RLCS. This step ensures that only those RLCS positions that are marked with a * will be instantiated. Then step 6.b calls the procedure, recursively, in order to construct the CLCS representations for the logical subject, arguments, and modifiers. This step ensures that the procedure will be called exactly once for each syntactic head appearing in the source-language structure. Finally, step

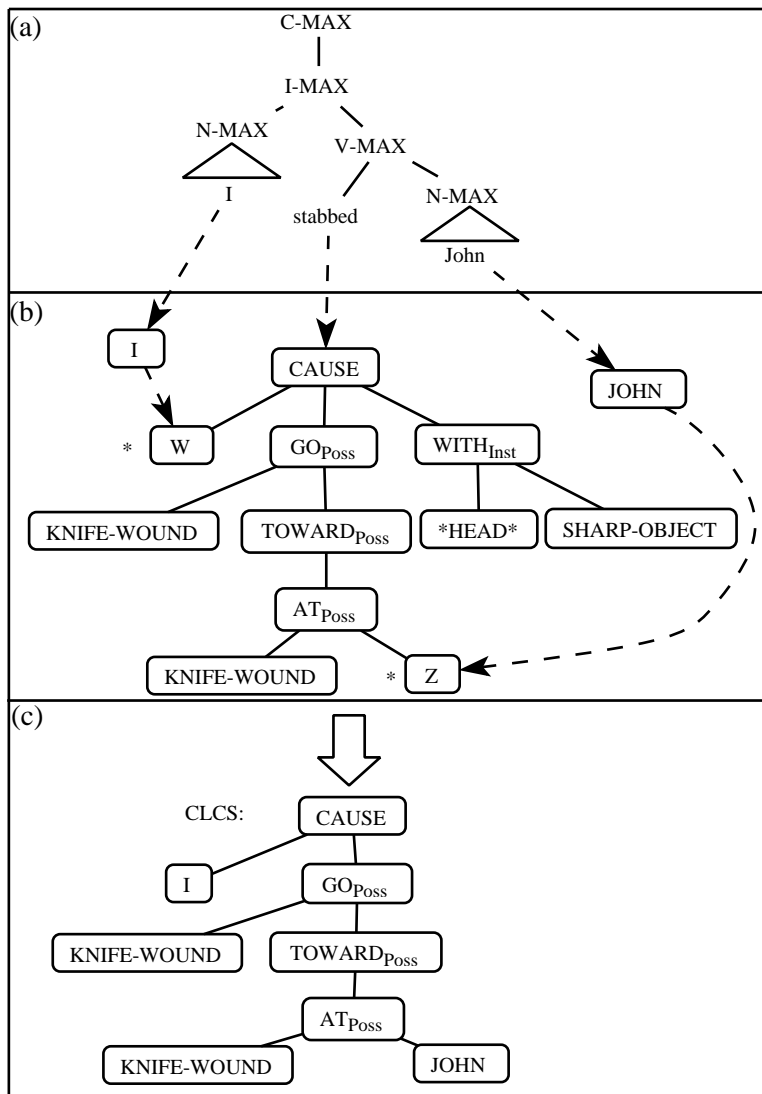


Figure 14: The LCS composition for *I stabbed John* requires recursive unification of *-marked arguments.

6.c fills the *-marked positions with the CLCS representations constructed by step 6.b, thus producing a new CLCS. Note that this final step is a “pseudo-unification” step in the sense that the *-marked position may be instantiated directly, if the position constitutes a leaf node of the RLCS, or indirectly (through unification), if the position constitutes a non-leaf node of the RLCS.

5.1.2 LCS Composition for Stab-Dar Example

Now that the algorithm for LCS Composition has been presented, we will see how it applies to the translation example (50). We will not discuss the generation of the target-language sentence for this example. Rather, we will show that the algorithm produces the same CLCS (*i.e.*, the interlingua) for both of these sentences. The intent here is to demonstrate the adequacy of the representation and the accompanying procedure for both sentences of this example.

In the case of the English sentence, the parser in the syntactic component of

UNITRAN supplies the source-language syntactic tree shown in figure 14(a). When this tree is passed to the LCS component, the RLCS's corresponding to the words *I*, *stab*, and *John* are selected as shown in figure 14(b). Finally, a single CLCS is composed from these RLCS's as shown in figure 14(c). The trace of this process is shown in appendix B.³⁴

The key point to note about this trace is that the two *-marked positions in the RLCS for *stab* are instantiated on the fourth and fifth entries to the **Compose_LCS** routine. This difference is due to the position of the * specification in the RLCS for *stab*. Recall that step 6.a of the algorithm selects an RLCS position by appealing to the generalized linking routine in conjunction with the position of the *-marker. Since there are only two *-marked positions identified by the linking routine ($[_{\text{Thing}} * W]$ and $[_{\text{Thing}} * Z]$), these are the only ones that are instantiated during the unification step 6.c; the first is instantiated as $[_{\text{Thing}} I]$ and the second is instantiated as $[_{\text{Thing}} JOHN]$.

The instantiation step differs crucially from the analogous step in the composition of the interlingua for the Spanish sentence. Figure 15(a) shows the syntactic tree for the Spanish case. This tree is passed to the LCS component which selects the RLCS's corresponding to the words *Yo*, *dar*, *puñaladas*, *a*, and *Juan* as shown in figure 15(b). A single CLCS is then composed from these RLCS's as shown in figure 15(c).

The trace of this process is shown in appendix C. There are two differences between the English and Spanish composition processes. The first is that, because there is an additional *-marked token in Spanish (*i.e.*, the token corresponding to the word *puñaladas*), there is an additional invocation of the procedure (*i.e.*, the fifth entry). This invocation results in the instantiation of the $[_{\text{Thing}} \text{KNIFE-WOUND}]$ constituent which fills the $[_{\text{Thing}} * Y]$ position of the RLCS corresponding to *dar*. Thus, the different positioning of the * accounts for the fact that the Spanish verb *dar* must explicitly realize the $[_{\text{Thing}} \text{KNIFE-WOUND}]$ argument as *puñaladas*, whereas no corresponding syntactic constituent is required for the English verb *stab*.

The second difference is that, in the Spanish case, the sixth entry to the procedure (which is analogous to the fifth entry in the English case) calls the **Compose_LCS** procedure recursively (*i.e.*, the seventh entry). This is because the * marker is positioned at the level of $[_{\text{Path}} * \text{TOWARD}_{\text{POSS}} \dots]$, not at the level of $[_{\text{Thing}} Z]$, in the lexical entry for *dar*. The recursive entry allows the RLCS corresponding to the word *a* to be unified with the matching portion of the *stab* RLCS. Note that the *a* RLCS has its own *-marked $[_{\text{Thing}} * Z]$ that is unified with $[_{\text{Thing}} JOHN]$ after exiting the sixth invocation. Thus, the positioning of the * marker accounts for the fact that the Spanish verb *dar* must realize the recipient of the action inside of a prepositional phrase, whereas the recipient of the English verb *dar* is realized directly as a noun phrase.

We have just seen how two divergences, conflational (*i.e.*, the suppression/realization of an argument) and structural (*i.e.*, the noun-phrase/prepositional-phrase distinc-

³⁴The first two calls to **Compose_LCS** do little more than “peel off” a level of maximal projection. This is because both the C-MAX and the I-MAX phrases contain an *empty* head (*i.e.*, a head with no lexical content). Such phrases are considered to be projections of their internal argument, which means they inherit their result from the structure produced by the recursive call on their internal argument.

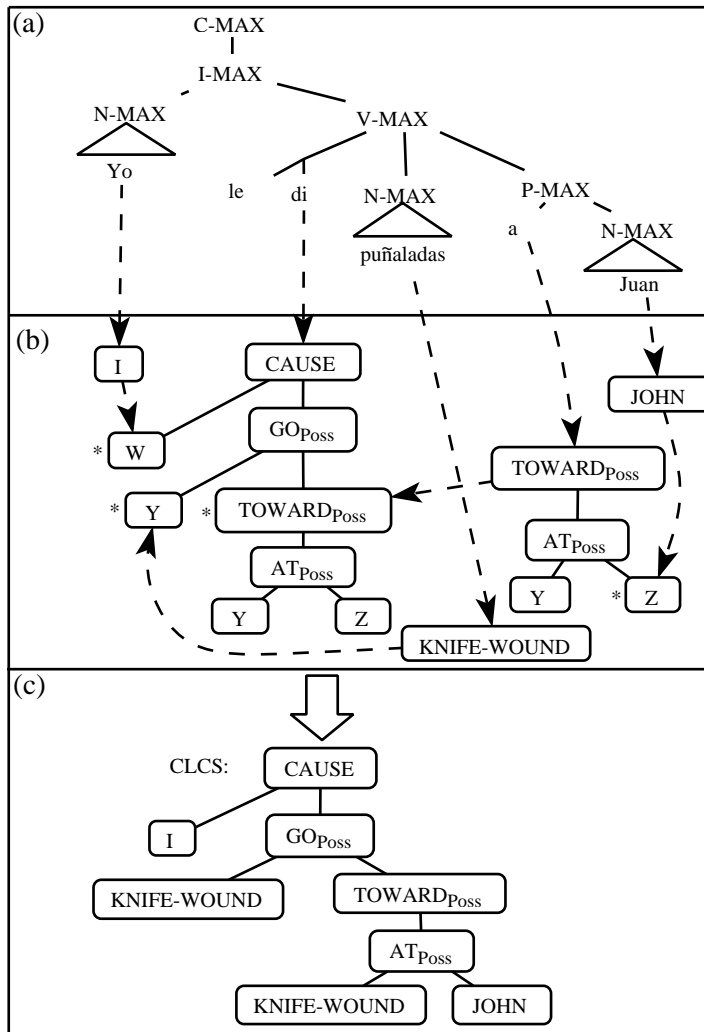


Figure 15: The LCS composition for *Yo le di puñaladas a Juan* requires recursive unification of *-marked arguments.

tion) are handled by means of an *abstraction* of language-independent properties (*e.g.*, the use of general representational constituents such as $[_{\text{Path}} \text{TOWARD}_{\text{Poss}} \dots]$) from language-specific idiosyncrasies (*e.g.*, the position of the * marker with respect to such constituents). Furthermore, we have seen that the *compositional* nature of the LCS representation allows these divergences to be handled by means of a recursive process that operates on each subcomponent of the representation first before combining them together by means of unification. In the next section we will examine issues and future work concerning the use of the LCS framework for machine translation.

6 Issues and Future Work

The LCS is suited to the description of information concerning verbs and their arguments as well as possible syntactic realizations of such constituents. However, there are clearly components of meaning that are not definable in terms of the LCS

representation. Consider the two verbs *fight* and *follow*. One might consider both of these verbs to be in the Locational class, which up until this point, have been the easiest to describe in the current framework. However, there are components of meaning for these two verbs that may exist at a level that is distinct from that of the LCS representation. A number of researchers, notably Mieizitis (1988), have argued that the verb *fight* could be inferred from other forms of knowledge representation, such as knowledge about punching and kicking:

(51) (punch Mary John) AND (kick John Mary) \rightarrow (fight Mary John)

Clearly, the LCS representation for *fight* could not be composed from the events corresponding to punching and kicking since these are not considered to be primitive activities, nor are they considered to be inherent in the lexical-semantic structure of the verb. Mieizitis' system achieves the mapping in (51) by means of abstraction to the next highest level in a concept network: *punch* and *kick* are instances of *violent-action*, which is an instance of *fight* (as long as there are two participants).

Similarly, the word *follow* has components of meaning that exist at a level that is independent of the LCS representation. For example, the sentence *John follows Mary* may be derived from two ideas:

(52) (walk John) AND (walk Mary (behind John Mary))

These implicational relations are not part of the LCS representation, which serves only to capture the relation between a predicate and its arguments. While it might be possible to extend the system to handle these implications, such an extension would probably be made at a representational level that is distinct from that of the LCS.

Another verbal construction that is not addressed in the current framework is that of light verbs (see *e.g.*, Grimshaw and Mester (1988), and Grimshaw (1990)). For example, we have seen that the verb *give* currently exists in the system as a verb of possessional transfer, but this verb may also be viewed as a light verb in constructions such as “give X a kiss” or “give X a kick.” In these cases, it is the direct object noun phrase, not the main verb, that supplies the information pertaining to the affected entity (*i.e.*, X). The current implementation does not handle light verbs because the notion of “affected entity” — the traditional role of Patient — is omitted from the theory of LCS. However, Jackendoff (1990) presents a considerably richer version of conceptual structure in which this notion is now included. He proposes that conceptual roles fall into two tiers, a *thematic tier* dealing with motion and location (*i.e.*, the LCS), and an *action tier* dealing with Actor-Patient relations. Light verbs are then handled by expressing the possessional transfer portion of the verb *give* in the thematic tier (which he claims must remain intact) while simultaneously expressing the role of the affected entity of the nominal constituent in the action tier (*i.e.*, Beneficiary for the nominal *kiss*, Patient for the nominal *kick*, *etc.*). The design of the current model does not preclude the possibility of superimposing this additional level of representation on the pre-existing LCS representation. Whether this enriched representation is general enough to serve as the foundation for interlingual machine translation is an area that has been left open for future investigation.

An additional extension that might be made to the current framework is one that concerns the handling of verbal alternations such as the locative and conative, illustrated here in (53) and (54), respectively:

- (53) John smeared paint on the wall
 John smeared the wall with paint

- (54) John cut the rope
 John cut at the rope

While the syntactic realization of alternating types is readily accommodated within the current framework, the semantic ramifications of these alternations are not fully addressed. The alternating surface forms are handled by using the * marker in such a way as to allow for more than one surface realization of the relevant arguments.³⁵ For example, the lexical entry for *smeare* specifies that the argument corresponding to *paint* may be realized in the “normal” object position or in an optional modifier slot headed by WITH_{Instr}. In addition, this entry specifies that the argument corresponding to *wall* may be realized as the argument of the directional path *on* or it may be realized without a path. The full entry is specified as follows:³⁶

- (55) [_{Event} CAUSE
 ([_{Thing} * W],
 [_{Event} GO_{Loc} ([_{Thing} * Y FLUID],
 [_{Path} * TOWARD_{Loc}
 ([_{Position} ON_{Loc} ([_{Thing} Y], [_{Thing} * Z])))]),
 [WITH_{Instr} * ([_{Event} *HEAD*], [_{Thing} Y])]]]

The advantage to representing alternating verbs this way is that they do not require a different lexical entry for each alternation. Instead, the syntactic information that differentiates the altered and unaltered forms is specified by the * mechanism within the same entry. However, using the same representation for both forms does not allow the system to distinguish between meaning-preserving alternations and non-meaning preserving alternations. For example, meaning is preserved in cases such as the following:

- (56) I entrusted him with my keys
 I entrusted my keys to him

³⁵We have already seen a case where the * marker was used in this way to accommodate the dative alternation for the verb *give* (see figure 8 and fn. 20).

³⁶Note that the representation in (55) would allow an unconstrained syntactic processor to incorrectly accept or generate the following sentences:

- * John smeared paint the wall
- * John smeared the wall paint
- * John smeared with paint on the wall
- * John smeared on the wall with paint

The current approach avoids such cases by applying a set of syntactic constraints; these constraints are outside of the scope of this paper. (See Dorr (1992b) for more details.)

This is not true in cases such as (54) above: in the first sentence the rope was actually cut, whereas in the second sentence the rope was not necessarily cut. Clearly, some mechanism beyond lexical semantics is required in order to capture this information. Brent (1988) proposes a theory that accounts for the difference in meaning between locative and conative alternations based on the degree of guaranteed completion for given word senses. A future area of investigation would be to incorporate such a notion of “completion” into the LCS framework.

Another issue concerning the adequacy of the representation is the question of whether the current framework is applicable to machine translation of non-European languages. Unlike the representation used in the Eurotra project, the LCS is intended to be more than just a “euroversal” representation (see Copeland *et al.* (1991)). An interesting example that demonstrates the potential applicability of the LCS representation to a non-European language is the translation of *I stabbed John* in Japanese:

- (57) E: I stabbed John
 J: watashi ga John ni (naifu de) kizu tsukemashita
 ‘I to John (with knife) cut attached’

The *kizu tsukemashita* portion of the Japanese is a compound construction formed from the phrase *kizu o tsukemashita* where *o* is the direct object marker. Thus, the word *kizu* is a noun that corresponds to the KNIFE-WOUND concept and the verb *tsukemashita* corresponds the possessional transfer portion of the LCS representation given earlier for the verbs *dar* and *stab*.³⁷ Several researchers have investigated the use of an LCS for languages such as Warlpiri (Hale and Laughren (1983)), Urdu (Husain (1989)), Greek (Olsen (1991)), Arabic (Shaban (1991)), and Chinese, Indo-European, and North American Indian (Talmy (1983, 1985)), among others. These investigations will be useful for future extensions of UNITRAN.

Another issue that must be addressed with respect to the LCS translation model is that of semantic ambiguity resolution. Whereas disambiguation strategies have been the focus of many previous semantics-based approaches (see *e.g.*, the Preference Semantics approach (Wilks (1973), Wilks and Fass (1992))), such strategies have not been the focus of the current investigation. However, there is no reason to assume that such strategies would be incompatible with the current framework. Most types of disambiguation require the use of a deeper type of knowledge that currently is not available in the LCS representation. Clearly, a fully interlingual system would require knowledge-based techniques to operate in tandem with the techniques described here in order to handle semantic ambiguity.

Two types of ambiguity that are particularly problematic for the LCS framework are lexical ambiguity and multiple modifier attachment. The first type of ambiguity arises most frequently in the context of selecting appropriate nominal and prepositional constituents for a given verb. It is clear that, at best, this type of ambiguity requires a richer knowledge-representation scheme.³⁸ The semantic

³⁷This example was given to me by Noyuri Soderland of the University of Massachusetts (personal correspondence).

³⁸Most likely, other types of knowledge (*e.g.*, discourse knowledge) would be needed as well. For example, in sentence (58), it would be useful to have more knowledge about the context of the

network of nominal primitives used in the preference semantics approach by Fass (1988) is one example of such a scheme.

Two examples of lexical ambiguity are the following:

(58) John borrowed money from the bank

(59) The man in the picture was tall

In (58), a strong semantic link (or *preference*, in the terminology of Wilks) exists between the act of *borrowing money* and the nominal *bank*. If such a link were made accessible to the LCS representation, the system would be able to determine that the word *bank* corresponds to a financial institution, and not to the border of a river. Such information is critical to machine translation. For example, the Spanish translation would be *banco* in the former case and *orilla* in the latter case. The use of a rich knowledge representation coupled with a preference semantics scheme would facilitate the process of lexical selection in such cases.

In (59), the preposition *in* does not have the same *contained-in* sense that is found in sentences such as *the man in the room was tall*. Such distinctions are not made in the LCS representation, though they are clearly necessary for translation. Consider the following English-French translations:

- | | | | | |
|------|-------|-------------------|---|---------------------------|
| | (i) | in a photo | ⇒ | sur la photo |
| | (ii) | in the paper | ⇒ | dans le journal |
| (60) | (iii) | in Canada | ⇒ | au Canada |
| | (iv) | in Spain | ⇒ | en Espagne |
| | (v) | in a tobacco shop | ⇒ | chez un marchand de tabac |

In order to select the appropriate preposition in French, the underlying meaning of the preposition is a critical component of the translation process. Research is currently underway (see Dorr and Voss (to appear)) to investigate a possible augmentation to the LCS scheme that allows for finer distinctions among spatial prepositions.

Despite these problematic cases, the LCS representation succeeds in filling the gap between the lexicon and syntactic structure, particularly in the context of the translation divergences presented earlier in figure 1. Clearly, the techniques used in deeper knowledge approaches are necessary for filling other gaps, including disambiguation.

An area that is currently under investigation is the development of a large-scale lexicon for processing within the LCS framework. Automatic lexical acquisition is becoming more critical to the success of machine translation because it is a tedious undertaking to construct dictionary representations by hand for each language.³⁹ Research is currently underway to investigate the possibility for scaling up the system through automatic means so that a wider range of phenomena and languages may be handled. A program has been developed (see Dorr (1992a), Dorr and Lee (1992)) that automatically acquires aspectual representations from corpora (currently the tagged version of the Lancaster/Oslo-Bergen (LOB) corpus)⁴⁰ by examining the context in which all verbs occur and then dividing them into four groups:

borrowing event such as the fact that John had just entered a building and was in front of the teller's window. For this, a fully developed context theory would be necessary.

³⁹The lexicon has been shown to be a major bottleneck for the development of UNITRAN: it took more than one person-month to define just 150 words per language.

⁴⁰The LOB corpus was obtained from the Norwegian Computing Center for the Humanities.

state, activity, accomplishment, and achievement. The division of verbs into these four groups is based on several syntactic tests that are well-defined in the linguistic literature such as those by Dowty (1979). This research has led to the discovery of a fundamental link between a subset of Jackendoff's primitives (those in the Circumstantial field) and the features of Dowty's aspectual scheme. If the link turns out to be generalizable to other fields, then the LCS framework could prove to be well-suited to the task of automatic construction of conceptual structures from corpora.

7 Summary

This paper has demonstrated that the LCS framework provides the appropriate level of abstraction for the application of a language-independent mapping routine. We have addressed the problem of translation divergences and have shown that the compositional nature of the LCS readily lends itself to the specification of lexical parameter settings that ultimately determine the shape of the surface syntactic structure.

The definition of a potentially large set of words is supported by the ability to combine the same lexical-semantic primitives in an indefinite number of ways. Although the number of primitives is small, the multiplicative effect provided by the fields and the recursive well-formedness constraints allows a number of linguistic classes of verbs to be represented.

The approach presented here tries to incorporate some of the more promising syntactic and semantic aspects of existing translation systems. Although several other approaches have attempted to find a systematic relation between syntactic structure and conceptual structure, UNITRAN is unique in that it attempts to account for a number of different language-specific syntactic phenomena by a simple parametric mechanism while remaining general and uniform enough to apply a single linking routine cross-linguistically.

8 Acknowledgements

This paper describes research done at the University of Maryland Institute for Advanced Computer Studies and the Artificial Intelligence Laboratory of the Massachusetts Institute of Technology. Support for this research has been provided in part by NSF Grant DCR-85552543 at MIT and by NSF Grant IRI-9120788 at the University of Maryland. Useful guidance and commentary during the research and preparation of this document were provided by Bob Berwick, Gary Coen, Bruce Dawson, Ken Hale, Joel Hoffman, Paola Merlo, Patrick Saint-Dizier, Noyuri Soderland, Stephen Soderland, Clare Voss and Amy Weinberg.

9 References

- Abeillé Anne, Yves Schabes, and Aravind K. Joshi (1990) "Using Lexicalized Tags for Machine Translation," *Proceedings of COLING-90, volume 3*, Helsinki, Finland, 1–6.
- Alonso, Juan Alberto (1990) "Transfer InterStructure: Designing an 'Interlingua' for Transfer-based MT Systems," *Proceedings of the Third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages*, Linguistics Research Center, The University of Texas, Austin, TX, 189–201.
- Arnold, Doug and Louisa Sadler (1990) "Theoretical Basis of MiMo," *Machine Translation* 5:3, 195–222.
- Arnold, Doug and Louis des Tombe (1987) "Basic Theory and Methodology in Eurotra," in *Machine Translation: Theoretical and Methodological Issues*, Sergei Nirenburg (ed.), Cambridge University Press, Cambridge, England.
- Arnold, Doug, Steven Krauwer, Louis des Tombe, and Louisa Sadler (1988) "Relaxed Compositionality in Machine Translation," *Proceedings of the Second International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages*, Carnegie Mellon University, Pittsburgh, PA.
- Bennett, P. A., R. L. Johnson, J. McNaught, J. M. Pugh, J. C. Sager, H. L. Somers (1986) *Multilingual Aspects of Information Technology*, Gower, Brookfield, VT.
- Bennett, Winfield S., Tanya Herlick, Katherine Hoyt, Joseph Liro and Ana Santisteban (1990) "A Computational Model of Aspect and Verb Semantics," *Machine Translation* 4:4, 247–280.
- Boitet, Christian (1987) "Research and Development on MT and Related Techniques at Grenoble University (GETA)," in *Machine Translation: The State of the Art*, Margaret King (ed.), Edinburgh University Press, Edinburgh.
- Boitet, Christian (1988) "Pros and Cons of the Pivot and Transfer Approaches in Multilingual Machine Translation," in *Recent Developments in Machine Translation*, Dan Maxwell, Klaus Schubert, and Toon Witkam (eds.), Foris, Dordrecht.
- Brent, Michael R. (1988) "Decompositional Semantics and Argument Expression in Natural Language," Master of Science thesis, Massachusetts Institute of Technology.
- Carbonell, Jaime G. and Masaru Tomita (1987) "Knowledge-based Machine Translation, the CMU Approach," in *Machine Translation: Theoretical and Methodological Issues*, Sergei Nirenburg (ed.), Cambridge University Press, Cambridge, England, 68–89.
- Colmerauer, A. (1971) "Les systèmes-Q ou un formalisme pour analyser et synthétiser des phrases sur ordinateur," *TAUM*, 1–45.
- Copeland, C., J. Durand, S. Krauwer, B. Maegaard (1991) "The Eurotra linguistic Specifications," in *Studies in machine Translation and Natural Language Processing, Volume 1*, Erwin Valentini (ed.), Commission of the European Communities, Brussels.
- Dorr, Bonnie J. (1990) "Solving Thematic Divergences in Machine Translation," *Proceedings of the 28th Annual Conference of the Association for Computational Linguistics*, University of Pittsburgh, Pittsburgh, PA, 127–134.
- Dorr, Bonnie J. (1992a) "A Parameterized Approach to Integrating Aspect with Lexical-Semantics for Machine Translation," *Proceedings of 30th Annual Conference of the Association of Computational Linguistics*, University of Delaware, Newark DE, 257–264.
- Dorr, Bonnie J. (1992b) "Parameterization of the Interlingua in Machine Translation," *Proceedings of Fourteenth International Conference on Computational Linguistics*, Nantes, France, 624–630.
- Dorr, Bonnie J. (in press) "Machine Translation: A View from the Lexicon," MIT Press.
- Dorr, Bonnie J. and Ki Lee (1992) "Building a Lexicon for Machine Translation: Use of Corpora for Aspectual Classification of Verbs," Institute for Advanced Computer Studies, University of Maryland, UMIACS TR 92-41, CS TR 2876.

- Dorr, Bonnie J. and Clare Voss (to appear) "Constraints on the Space of MT Divergences," *Working Notes for the AAAI Spring Symposium on Building Lexicons for Machine Translation*, Stanford University, CA.
- Dowty, David (1979) *Word Meaning and Montague Grammar*, Reidel, Dordrecht, Netherlands.
- Fass, D. C. (1988) "Collative Semantics: A Semantics for Natural Language Processing," Computing Research Laboratory, New Mexico State University, Memorandum in Computer and Cognitive Science, MCCS-88-118.
- Fillmore, Charles J. (1968) "The Case for Case," in *Universals in Linguistic Theory*, Bach, E., and R. T. Harms (eds.), Holt, Rinehart, and Winston, 1-88.
- Grimshaw, Jane (1990) *Argument Structure*, MIT Press, Cambridge, MA.
- Grimshaw, Jane and Armin Mester (1988) "Light Verbs and θ -Marking," *Linguistic Inquiry* 19:2, 205-232.
- Gruber, J. S. (1965) "Studies in Lexical Relations," Ph.D. thesis, Department of Information Science Department, Massachusetts Institute of Technology.
- Hale, Kenneth and Jay Keyser (1986a) "Some Transitivity Alternations in English," Center for Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA, Lexicon Project Working Papers #7.
- Hale, Kenneth and Jay Keyser (1986b) "A View from the Middle," Center for Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA, Lexicon Project Working Papers #10.
- Hale, Kenneth and S. Jay Keyser (1989) "On Some Syntactic Rules in the Lexicon," Center for Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA, manuscript.
- Hale, Kenneth and Mary Laughren (1983) "The Structure of Verbal Entries: Preface to Dictionary Entries of Verbs," Massachusetts Institute of Technology, Cambridge, MA, Warlpiri Lexicon Project.
- Husain, Saadia (1989) "A Lexical Conceptual Structure Editor," Bachelor of Science thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology.
- Hutchins, W. J. (1986) *Machine Translation: Past, Present, Future*, Ellis Horwood Limited, Chichester, England.
- Hutchins, J. W. and H. L. Somers (1992) *An Introduction to Machine Translation*, Academic Press, London, England.
- Jackendoff, Ray S. (1983) *Semantics and Cognition*, MIT Press, Cambridge, MA.
- Jackendoff, Ray S. (1990) *Semantic Structures*, MIT Press, Cambridge, MA.
- Johnson, Rod, Maghi King, and Louis des Tombe (1985) "EUROTRA: A Multilingual System under Development," *Computational Linguistics* 11:2-3, 155-169.
- Kaplan, Ronald M., Klaus Netter, Jürgen Wedekind, Annie Zaenen (1989) "Translation By Structural Correspondences," *Proceedings of Fourth Conference of the European Chapter of the Association for Computational Linguistics*, Manchester, 272-281.
- King, Margaret, (ed.) (1987) *Machine Translation: The State of the Art*, Edinburgh University Press, Edinburgh.
- Levin, Beth (1985) "Lexical Semantics in Review," Center for Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA, Lexicon Project Working Papers #1.
- Levin, Beth (in press) *English Verb Classes and Alternations: A Preliminary Investigation*, University of Chicago Press, Chicago, IL.
- Levin, Beth and Malka Rappaport (1986) "The Formation of Adjectival Passives," *Linguistic Inquiry* 17, 623-662.

- Lytinen, Steven and Roger Schank (1982) "Representation and Translation," Department of Computer Science, Yale University, New Haven, CT, Technical Report 234.
- Maxwell, Dan, Klaus Schubert, and Toon Witkam (1988) *Recent Developments in Machine Translation*, Foris, Dordrecht.
- McCord, Michael C. (1989) "Design of LMT: A Prolog-Based Machine Translation System," *Computational Linguistics* 15:1, 33–52.
- Mel'čuk, Igor and Alain Poguère (1987) "A Formal Lexicon in Meaning-Text Theory (Or How to Do Lexica with Words)," *Computational Linguistics* 13:3–4, 261–275.
- Miezitis, Mara (1988) "Generating Lexical Options by Matching in a Knowledge Base," Technical Report CSRI-217, Master of Science thesis, Department of Department of Computer Science, University of Toronto.
- Muraki, K. (1987) "PIVOT: A Two-Phase Machine Translation System," *Machine Translation Summit – Manuscripts and Program*, Japan, 81–83.
- Nirenburg, Sergei (ed.) (1987) "Machine Translation: Theoretical and Methodological Issues," Cambridge University Press.
- Nirenburg, Sergei and Lori Levin (1989) "Knowledge Representation Support," *Machine Translation* 4:1, 25–52.
- Nirenburg, Sergei, Victor Raskin, and Allen B. Tucker (1987) "The Structure of Interlingua in TRANSLATOR," in *Machine Translation: Theoretical and Methodological Issues*, Sergei Nirenburg (ed.), Cambridge University Press, Cambridge, England, 90–113.
- Nirenburg, Sergei, Jaime Carbonell, Masaru Tomita, and Kenneth Goodman (1992) *Machine Translation: A Knowledge-Based Approach*, Morgan Kaufmann, San Mateo, CA.
- Noord, Gertjan van, Joke Dorrepaal, Doug Arnold, Steven Krauwer, Lousia Sadler, and Louis des Tombe (1989) "An Approach to Sentence-Level Anaphora in Machine Translation," *Proceedings of Fourth Conference of the European Chapter of the Association for Computational Linguistics*, Manchester.
- Noord, Gertjan van, Joke Dorrepaal, Pim van der Eijk, Maria Florenza, and Louis des Tombe (1990) "The MiMo2 Research System," *Proceedings of the Third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages*, Linguistics Research Center, The University of Texas, Austin, TX, 213–233.
- Olsen, Mari Broman (1991) "Lexical Semantics, Machine Translation, and Talmy's Model of Motion Verbs," Northwestern University, Evanston, IL, Linguistics Working Paper, Volume 3.
- Pinker, Steven (1989) *Learnability and Cognition: The Acquisition of Argument Structure*, MIT Press, Cambridge, MA.
- Pustejovsky, James (1988) "The Geometry of Events," Center for Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA, Lexicon Project Working Paper 24.
- Pustejovsky, James (1989) "The Semantic Representation of Lexical Knowledge," *Proceedings of the First International Lexical Acquisition Workshop*, IJCAI-89, Detroit, MI.
- Pustejovsky, James (1990) "The Generative Lexicon," *Computational Linguistics* 17:4, 409–441.
- Pustejovsky, James (1991) "The Syntax of Event Structure," *Cognition* 41.
- Rappaport, Malka and Beth Levin (1988) "What to Do with Theta-Roles," in *Thematic Relations*, Wendy Wilkins (ed.), Academic Press.
- Rappaport, Malka, Mary Laughren, and Beth Levin (1987) "Levels of Lexical Representation," Center for Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA, Lexicon Project Working Papers #20.
- Rieger, Charles J. III (1975) "Conceptual Memory and Inference," in *Conceptual Information Processing*, Schank, Roger (ed.), Elsevier Science Publishers, Amsterdam, The Netherlands.

- Sadler, Louisa, Ian Crookston, Doug Arnold, and Andy Way (1990) "LFG and Translation," *Proceedings of the Third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages*, Linguistics Research Center, The University of Texas, Austin, TX, 121–130.
- Schank, Roger (1972) "Conceptual Dependency: A Theory of Natural Language Understanding," *Cognitive Psychology* 3, 552–631.
- Schank, Roger C. (1973) "Identification of Conceptualizations Underlying Natural Language," in *Computer Models of Thought and Language*, Roger C. Schank and K. M. Colby (eds.), Freeman, San Francisco, CA, 187–247.
- Schank, Roger C. (ed.) (1975) *Conceptual Information Processing*, Elsevier Science Publishers, Amsterdam, Holland.
- Schank, Roger C. and Robert Abelson (1977) *Scripts, Plans, Goals, and Understanding*, Lawrence Erlbaum Associates, Inc., Hillsdale, NJ.
- Shaban, Marwan (1991) "GB Parsing of Arabic," Master of Science thesis, Department of Computer Science Department, Boston University.
- Sharp, Randall M. (1985) "A Model of Grammar Based on Principles of Government and Binding," Master of Science thesis, Department of Computer Science, University of British Columbia.
- Siskind, Jeffrey Mark (1989) "Decomposition," Massachusetts Institute of Technology, Cambridge, MA, area exam paper.
- Siskind, Jeffrey Mark (1992) "Naive Physics, Event Perception, Lexical Semantics, and Language Acquisition," Ph.D. thesis, Department of Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA.
- Slocum, Jonathan (1988) *Machine Translation Systems*, Cambridge University Press, Cambridge.
- Sproat, Richard (1985) "Identification of Conceptualizations Underlying Natural Language," in *Lexical Semantics in Review*, Beth Levin, Lexicon Project Working Paper 1, Center for Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA, 115–124.
- Talmy, Leonard (1983) "How Language Structures Space," in *Spatial Orientation: Theory, Research, and Application*, Pick, Herbert L., Jr., and Linda P. Acredolo (eds.), Plenum Press, New York.
- Talmy, Leonard (1985) "Lexicalization Patterns: Semantic Structure in Lexical Forms," in *Grammatical Categories and the Lexicon*, Timothy Shopen (ed.), University Press, Cambridge, England.
- Thurmair, Gregor (1990) "Complex Lexical Transfer in METAL," *Proceedings of the Third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages*, Linguistics Research Center, The University of Texas, Austin, TX, 91–107.
- Uchida, H. (1989) "ATLAS: Fujitsu Machine Translation System," *Machine Translation Summit II – Manuscripts and Program*, Japan, 129–134.
- Vauquois Bernard (1975) *La Traduction Automatique à Grenoble*, Dunod, Paris.
- Vauquois Bernard, and Christian Boitet (1985) "Automated Translation at Grenoble University," *Computational Linguistics* 11:1, 28–36.
- White, Michael (1992) "Conceptual Structures and CCG: Linking Theory and Incorporated Argument Adjuncts," *Proceedings of Fourteenth International Conference on Computational Linguistics*, Nantes, France, 246–252.
- Wilks, Yorick (1973) "An Artificial Intelligence Approach to Machine Translation," in *Computer Models of Thought and Language*, R. C. Schank and K. M. Colby (eds.), Freeman, San Francisco, CA, 114–151.

- Wilks, Yorick (1987) "Primitives," in *Encyclopedia of Artificial Intelligence*, S.C. Shapiro (ed.), John Wiley and Sons, New York, NY, 759–761.
- Wilks, Yorick and Dan Fass (1992) "Preference Semantics," in *Encyclopedia of Artificial Intelligence*, S. C. Shapiro (ed.), John Wiley and Sons, New York, NY, 1182–1194.
- Zubizarreta, Maria Luisa (1982) "On the Relationship of the Lexicon to Syntax," Ph.D. thesis, Massachusetts Institute of Technology.
- Zubizarreta, Maria Luisa (1987) *Levels of Representation in the Lexicon and in the Syntax*, Foris Publications, Dordrecht, Holland/Cinnaminson, USA.

A Linguistic Coverage of Extended Set of Primitives

Class of Verb	Primitive	Examples
position	STAY _{Temp}	Basic: The meeting remained at noon. Causative: We kept the meeting at noon.
	STAY _{Loc}	Basic: The statue remained in the park. Causative: We kept the statue in the park.
	BE _{Temp}	Basic: The meeting is at noon.
	BE _{Loc}	Basic: The statue is in the park.
change of position	GO _{Loc}	Basic: The rock fell from the roof to the ground. Causative: John threw the rock from the roof to the ground. Permissive: John dropped the rock from the roof to the ground.
	GO _{Temp}	Basic: The meeting changed from 2:00 to 4:00. Causative: We changed the meeting from 2:00 to 4:00. Permissive: We allowed the meeting to change from 2:00 to 4:00.
directed motion	GO _{Loc}	Basic: John entered the room. Causative: John broke into the room. Permissive: John let Beth enter the room.
	GO _{Poss}	Basic: Beth received the doll. Causative: John gave the doll to Beth. Permissive: Beth relinquished the doll.
motion with manner	GO _{Loc}	Basic: The boat sailed to Cuba. Causative: John sailed the boat to Cuba. Permissive: John let the boat sail to Cuba.
exchange	CAUSE-EXCHANGE	Basic: Beth bought the doll for Mary.
physical state	BE _{Ident}	Basic: The door is open.
	STAY _{Ident}	Basic: The door remained open. Causative: John kept the door open. Permissive: John left the door open.
change of physical state	GO _{Ident}	Basic: The door opened. Causative: John opened the door. Permissive: John let the door open.
orientation	ORIENT _{Loc}	Basic: The sign points to Philadelphia. Causative: John pointed the sign to Philadelphia. Permissive: John let the sign point to Philadelphia.
existence	BE _{Exist}	Basic: The house exists on my property. Causative: John built the house on my property. Permissive: John allowed the house to exist on my property.
	GO _{Exist}	Basic: The house appeared. Causative: The magician made the house disappear. Permissive: The magician allowed the house to reappear.
	STAY _{Exist}	Basic: The situation persisted. Causative: Bill caused the situation to persist. Permissive: Bill allowed the situation to persist.
circumstance	BE _{Circ}	Basic: John is shipping goods to California.
	GO _{Circ}	Basic: John started shipping goods to California. Causative: John forced Beth to ship goods. Permissive: Beth allowed John to start shipping goods.
	STAY _{Circ}	Basic: John continued shipping goods. Causative: John kept Beth from shipping goods. Permissive: Beth exempted John from shipping goods.
range	GO-EXT _{Ident}	Basic: Our clients range from psychiatrists to psychopaths. Causative: The sun caused the leaves to range from green to brown. Permissive: Bill allowed the situation to range from bad to worse.
	GO-EXT _{Temp}	Basic: The meeting lasted from noon to night. Causative: John made the meeting last from noon to night. Permissive: John allowed the meeting to last from noon to night.
	GO-EXT _{Loc}	Basic: The road went from Boston to Albany. Causative: The workers made the road extend from Boston to Albany.
intention	ORIENT _{Circ}	Basic: John intended to ship goods to California.
	ORIENT _{Temp}	Basic: John aims to start at 2:00.
ownership	BE _{Poss}	Basic: The doll belongs to Beth.
	STAY _{Poss}	Basic: The doll remained in her hands. Causative: Amy kept the doll.
ingestion	EAT	Basic: John ate breakfast.
psychological state	BE _{Ident}	Basic: Beth liked John. Causative: John pleased Beth.
perception and communication	HEAR _{Perc}	Basic: John heard Mary. Causative: John listened to Mary. Permissive: Mary told the story to John.
	SEE _{Perc}	Basic: John saw Mary. Causative: John watched Mary. Permissive: Mary showed the dress to John.
mental process	BE _{Perc} GO _{Perc}	Basic: Beth knew the lesson. Basic: Beth learned the lesson.
cost	ORIENT _{Ident}	Basic: The book costs \$10.00. Causative: Beth charged me \$10.00 for the book.
load/spray	GO _{Loc}	Causative: Bill smeared the wall with paint.
contact/effect	GO _{Poss}	Basic: The knife cut Bill.
		Causative: Mary stabbed Bill with a knife.

B Trace of LCS Composition: I stabbed John

1. Entering **Compose_LCS**:

$$\begin{aligned} & [C\text{-MAX} [I\text{-MAX} [N\text{-MAX} I] [V\text{-MAX} stabbed [N\text{-MAX} John]]]] \\ & X = \text{empty}(\text{projection of } I\text{-MAX}) \\ & Z_1 = [I\text{-MAX} [N\text{-MAX} I] [V\text{-MAX} stabbed [N\text{-MAX} John]]] \end{aligned}$$
2. Entering **Compose_LCS**:

$$\begin{aligned} & [I\text{-MAX} [N\text{-MAX} I] [V\text{-MAX} stabbed [N\text{-MAX} John]]] \\ & X = \text{empty}(\text{projection of } V\text{-MAX}) \\ & Z_1 = [V\text{-MAX} stabbed [N\text{-MAX} John]] \end{aligned}$$
3. Entering **Compose_LCS**: $[V\text{-MAX} stabbed [N\text{-MAX} John]]$

$$\begin{aligned} & X = \text{stab} \\ & X' = [Event \text{ CAUSE} ([Thing * W], \\ & \quad [Event \text{ GOPoss} ([Thing \text{ KNIFE-WOUND}], \\ & \quad \quad [Path \text{ TOWARDPoss} \\ & \quad \quad \quad ([Position \text{ ATPoss} ([Thing \text{ KNIFE-WOUND}], [Thing * Z])))]), \\ & \quad [WITH_{Instr} ([Event *HEAD*], [Thing \text{ U SHARP-OBJECT}])]])] \end{aligned}$$

$$\begin{aligned} W &= [N\text{-MAX} I] \\ Z_1 &= [N\text{-MAX} John] \\ i' &= [Thing * W] \end{aligned}$$
4. Entering **Compose_LCS**: $[N\text{-MAX} I]$

$$\begin{aligned} & X = I \\ & X' = [Thing I] \end{aligned}$$
4. Exiting **Compose_LCS**: $L = [Thing I]$

$$\begin{aligned} & X' = [Event \text{ CAUSE} ([Thing I], \\ & \quad [Event \text{ GOPoss} ([Thing \text{ KNIFE-WOUND}], \\ & \quad \quad [Path \text{ TOWARDPoss} \\ & \quad \quad \quad ([Position \text{ ATPoss} ([Thing \text{ KNIFE-WOUND}], [Thing * Z])))]), \\ & \quad [WITH_{Instr} ([Event *HEAD*], [Thing \text{ U SHARP-OBJECT}])]])] \end{aligned}$$

$$i' = [Thing * Z]$$
5. Entering **Compose_LCS**: $[N\text{-MAX} John]$

$$\begin{aligned} & X = John \\ & X' = [Thing JOHN] \end{aligned}$$
5. Exiting **Compose_LCS**: $L = [Thing JOHN]$

$$\begin{aligned} & X' = [Event \text{ CAUSE} ([Thing I], \\ & \quad [Event \text{ GOPoss} ([Thing \text{ KNIFE-WOUND}], \\ & \quad \quad [Path \text{ TOWARDPoss} \\ & \quad \quad \quad ([Position \text{ ATPoss} ([Thing \text{ KNIFE-WOUND}], [Thing JOHN])))]), \\ & \quad [WITH_{Instr} ([Event *HEAD*], [Thing \text{ U SHARP-OBJECT}])]])] \end{aligned}$$
3. Exiting **Compose_LCS**:

$$\begin{aligned} L = X' &= [Event \text{ CAUSE} ([Thing I], \\ & \quad [Event \text{ GOPoss} ([Thing \text{ KNIFE-WOUND}], \\ & \quad \quad [Path \text{ TOWARDPoss} \\ & \quad \quad \quad ([Position \text{ ATPoss} ([Thing \text{ KNIFE-WOUND}], [Thing JOHN]))]]))]^{41} \end{aligned}$$
2. Exiting **Compose_LCS**:

$$\begin{aligned} L = X' &= [Event \text{ CAUSE} ([Thing I], \\ & \quad [Event \text{ GOPoss} ([Thing \text{ KNIFE-WOUND}], \\ & \quad \quad [Path \text{ TOWARDPoss} \\ & \quad \quad \quad ([Position \text{ ATPoss} ([Thing \text{ KNIFE-WOUND}], [Thing JOHN]))]]))] \end{aligned}$$
1. Exiting **Compose_LCS**:

$$\begin{aligned} L = X' &= [Event \text{ CAUSE} ([Thing I], \\ & \quad [Event \text{ GOPoss} ([Thing \text{ KNIFE-WOUND}], \\ & \quad \quad [Path \text{ TOWARDPoss} \\ & \quad \quad \quad ([Position \text{ ATPoss} ([Thing \text{ KNIFE-WOUND}], [Thing JOHN]))]]))] \end{aligned}$$

⁴¹ Note that the $WITH_{Instr}$ modifier that appears in the RLCS is not included in the final CLCS since it is an uninstantiated optional modifier.

C Trace of LCS Composition: Yo le di puñaladas a Juan

1. Entering **Compose_LCS**:
 [C-MAX [I-MAX [N-MAX Yo] [V-MAX le di [N-MAX puñaladas] [P-MAX a [N-MAX Juan]]]]
 X = empty (projection of I-MAX)
 Z₁ = [I-MAX [N-MAX Yo] [V-MAX le di [N-MAX puñaladas] [P-MAX a [N-MAX Juan]]]]
2. Entering **Compose_LCS**:
 [I-MAX [N-MAX Yo] [V-MAX le di [N-MAX puñaladas] [P-MAX a [N-MAX Juan]]]
 X = empty (projection of V-MAX)
 Z₁ = [V-MAX le di [N-MAX puñaladas] [P-MAX a [N-MAX Juan]]]
3. Entering **Compose_LCS**:
 [V-MAX le di [N-MAX puñaladas] [P-MAX a [N-MAX Juan]]]⁴²
 X = dar
 X' = [Event CAUSE ([Thing * W],
 [Event GOPoss ([Thing * Y],
 [Path * TOWARDPoss ([Position ATPoss ([Thing Y], [Thing Z])]])])]
- W = [N-MAX Yo]
 Z₁ = [N-MAX puñaladas]
 Z₂ = [P-MAX a [N-MAX Juan]]
 i' = [Thing * W]
4. Entering **Compose_LCS**: [N-MAX Yo]
 X = Yo
 X' = [Thing I]
4. Exiting **Compose_LCS**: L = [Thing I]
 X' = [Event CAUSE ([Thing I],
 [Event GOPoss ([Thing * Y],
 [Path * TOWARDPoss
 ([Position ATPoss ([Thing Y], [Thing Z])]])])]
- i' = [Thing * Y]
5. Entering **Compose_LCS**: [N-MAX puñaladas]
 X = puñaladas
 X' = [Thing KNIFE-WOUND]
5. Exiting **Compose_LCS**: L = [Thing KNIFE-WOUND]
 X' = [Event CAUSE ([Thing I],
 [Event GOPoss ([Thing KNIFE-WOUND],
 [Path * TOWARDPoss
 ([Position ATPoss ([Thing KNIFE-WOUND], [Thing Z])]])])]
- i' = [Path * TOWARDPoss ([Position ATPoss ([Thing KNIFE-WOUND], [Thing Z])]])]
6. Entering **Compose_LCS**: [P-MAX a [N-MAX Juan]]
 X = a
 X' = [Path TOWARDPoss ([Position ATPoss ([Thing Y], [Thing * Z])]])
 Z₁ = [N-MAX Juan]
- i' = [Thing * Z]
7. Entering **Compose_LCS**: [N-MAX Juan]
 X = Juan
 X' = [Thing JOHN]
7. Exiting **Compose_LCS**: L = [Thing JOHN]
 X' = [Path TOWARDPoss ([Position ATPoss ([Thing Y], [Thing JOHN])]])]
6. Exiting **Compose_LCS**:
 L = [Path TOWARDPoss ([Position ATPoss ([Thing Y], [Thing JOHN])]])
 X' = [Event CAUSE ([Thing I],
 [Event GOPoss ([Thing KNIFE-WOUND],
 [Path TOWARDPoss
 ([Position ATPoss ([Thing KNIFE-WOUND], [Thing JOHN])]])])]
3. Exiting **Compose_LCS**:
 L = X' = [Event CAUSE ([Thing I],
 [Event GOPoss ([Thing KNIFE-WOUND],
 [Path TOWARDPoss
 ([Position ATPoss ([Thing KNIFE-WOUND], [Thing JOHN])]])])]
2. Exiting **Compose_LCS**:
 L = X' = [Event CAUSE ([Thing I],
 [Event GOPoss ([Thing KNIFE-WOUND],
 [Path TOWARDPoss
 ([Position ATPoss ([Thing KNIFE-WOUND], [Thing JOHN])]])])]
1. Exiting **Compose_LCS**:
 L = X' = [Event CAUSE ([Thing I],
 [Event GOPoss ([Thing KNIFE-WOUND],
 [Path TOWARDPoss ([Position ATPoss ([Thing KNIFE-WOUND], [Thing JOHN])]])])]

⁴²The handling of the clitic pronoun *le* is not described here. For all intents and purposes, we may ignore this pronoun since it is coreferential with *Juan*. Clitics are handled by an independent process that is outside of the scope of this paper.