

---

# Interlingua based English-Hindi Machine Translation and Language Divergence

Shachi Dave

Jignashu Parikh

Pushpak Bhattacharyya<sup>1</sup>

Department of Computer Science and Engineering,  
Indian Institute of Technology,  
Bombay.

**Keywords:** Interlingua, Language Divergence, Language Analysis, Language Generation, Universal Networking Language

## Abstract

Interlingua and transfer based approaches to machine translation have long been in use in competing and complimentary ways. The former proves economical in situations where translation among multiple languages is involved, while the latter is used for pair specific translation tasks. The additional attraction of an interlingua is that it can be used as a knowledge representation scheme. But given a particular interlingua, its adoption depends on its ability to (a) capture the knowledge in texts precisely and accurately and (b) handle cross language divergences. This paper studies the language divergence between English and Hindi and its implication to machine translation between these languages using the Universal Networking Language (UNL). UNL has been introduced by the United Nations University (UNU), Tokyo, to facilitate the transfer and exchange of information over the internet in the natural languages of the world. The representation works at the level of single sentences and defines a semantic net like structure in which nodes are word concepts and arcs are semantic relations between these concepts. Hindi belongs to the Indo European family of languages. The language divergences between Hindi and English can be considered as representing the divergences between SOV and SVO class of languages. The work presented here is the only one to our knowledge that describes language divergence phenomena in the framework of computational linguistics through a South Asian language.

## 1 Introduction

The *digital divide* among people arises not only from the infrastructural factors like personal computers and high speed networks, but also from the *Language Barrier*. This barrier appears whenever the language in which information is presented is not known to the receiver of that information. The Web contents are mostly in English and cannot be accessed without some proficiency in this language. This is true for other languages too. The Universal Networking Language (UNL) has been proposed by the United Nations University (UNU) for overcoming the language barrier. However, a particular interlingua can be adopted only if it can capture the knowledge present in natural language documents precisely and accurately. Also it should have the ability to handle cross language divergences. Our work investigates the efficacy of the UNL as an Interlingua in the context of the language divergences between Hindi and English. The language divergence between these two languages can be considered representative of the divergences between the SOV and SVO class of languages.

Researchers have long been investigating the Interlingua approach to MT and some of them have considered the widely used transfer approach as the better alternative (Arnold and Sadler 1990; Boitet 1988; Vauquois and Boitet 1985). In the transfer approach, **some** amount of text analysis is done in the context of the source language and then some processing is carried out on the translated text in the context of the target language. But the bulk of the work is done on the comparative information on the specific pair of languages. The arguments in favour of the transfer approach to MT are (a) the sheer difficulty of designing a single interlingua that can be all things to all languages and (b) the fact that translation is, by its very nature, an exercise in comparative linguistics. The Eurotra system (Schutz, Thurmair, *et. al.*,

---

<sup>1</sup> Author for Correspondence, pb@cse.iitb.ac.in

---

1991; Arnold and des Tombes, 1987; King and Perschke, 1987; Perschke, 1989) in which groups from all the countries of the European Union participated is based on the transfer approach. So is the Verbmobil system (Wahlster 1997) sponsored by the German Federal Ministry for Research and Technology.

However, since the late eighties, the interlingua approach has gained momentum with commercial interlingua based machine translation systems being implemented. PIVOT of NEC (Okumura, Muraki, *et. al.*, 1991; Muraki, 1989), ATLAS II of Fujitsu (Uchida, 1989), ROSETTA of Phillips (Landsbergen, 1987) and BSO (Witkam, 1988; Schubert, 1988) in the Netherlands are the examples in point. In the last mentioned, the interlingua is not a specially designed language, but Esperanto. It is more economical to use an interlingua if translation among multiple languages is required. Only  $2N$  converters will have to be written, as opposed to  $N \times (N - 1)$  converters in the transfer approach, where  $N$  is the number of languages involved.

The interlingua approach can be broadly classified into (a) primitive based and (b) *deeper* knowledge representation based. (Schank 1972, 1973, 1975; Schank and Abelson 1977; Lytinen and Schank 1982) using Conceptual Dependency, the UNITRAN system (Dorr 1992, 1993) using the LCS and Wilk's system (Wilks 1972) are the examples of the former, while CETA (Vauquois 1975), (Carbonell and Tomita 1987), KBMT (Nirenburg *et. al.* 1992), TRANSLATOR (Nirenburg *et al.* 1987), PIVOT (Muraki 1987) and ATLAS (Uchida 1989) are the examples of the latter. The UNL falls into the latter category.

(Dorr 1993) describes how language divergences can be handled using the Lexical Conceptual Structure (LCS) as the interlingua in the UNITRAN system. The argument is that it is the complex divergences that necessitate the use of an interlingua representation. This is because of the fact that such a representation allows surface syntactic distinctions to be represented at a level that is independent of the underlying *meanings* of the source and target sentences. Factoring out these distinctions allows cross linguistic generalizations to be captured at the level of the lexical semantic structure.

The work presented here is the only one to our knowledge that describes language divergences between Hindi and English in a formal way from the point of view of computational linguistics. However, several studies by the linguistic community bring out the differences between the western and Indian languages (Bholanath 1987, Gopinathan 1993). These are presented in section 5.

Many systems have been developed in India for translation to and from Indian languages. The Anusaaraka system- based on the Paninian Grammar (Akshar Bharati *et. al.*, 1996)- renders text from one Indian language into another. It analyses the source language text and presents the information in the target language retaining a flavour of the source language. The grammaticality constraint is relaxed and a special purpose notation is devised. The aim of this system is to allow language access and not machine translation. IIT Kanpur is involved in designing translation support systems called *Anglabharati* and *Anubharati*. These are for MT between English and Indian languages and also among Indian languages (Sinha 1994). The approach is based on the word expert model utilizing the *Karaka* theory, a pattern directed rule base and a hybrid example base. In MaTra (Rao *et. al.* 2000)- a human-aided translation system for English to Hindi- the focus is on the innovative use of man machine synergy. The system breaks an English sentence into chunks and displays it using an intuitive browser like representation which the user can verify and correct. The Hindi sentence is generated after the system has resolved the ambiguities and the lexical absence of words with the help of the user.

We now give a brief introduction to the **Universal Networking Language**. It is an interlingua that has been proposed by the United Nations University to access, transfer and process information on the internet in the natural languages of the world. UNL represents information sentence by sentence. Each sentence is converted into a hyper graph having

---

concepts as nodes and relations as directed arcs. Concepts are called *Universal Words (UWs)*. The knowledge within a document is expressed in three dimensions:

- a. Word Knowledge is represented by **Universal Words (UWs)** which are language independent. These UWs have *restrictions* which describe the sense of the word. For example, *drink(icl>liquor)* denotes the noun *liquor*. *icl* stands for inclusion and forms an *is-a* structure as in semantic nets (woods 1985). The UWs are picked up from the lexicon during the analysis into or generation from the UNL expressions. The entries in the lexicon have syntactic and semantic attributes. The former depends on the language word while the latter is obtained from the language independent ontology.
- b. Conceptual Knowledge is captured by relating UWs through the standard set of **Relations Labels (RLs)** (UNL 1998). For example, *Humans affect the environment* is described in UNL as

**agt(affect(icl>do)).@present.@entry:01, human(icl>animal).@pl:I3)**  
**obj(affect(icl>do)).@present.@entry:01, environment(icl>abstract thing).@pl:I3)**

**agt** means the *agent* and **obj** the *object*. *affect(icl>do)*, *human(icl>animal)* and *environment(icl>abstract thing)* are the UWs denoting concepts.

Speaker's view, aspect, time of the event, *etc.* are captured by **Attribute Labels (ALs)**. For instance, in the above example, the attribute *@entry* denotes the main predicate of the sentence, *@present* the *present tense* and *@pl* the *plural number*.

The total number of relations in the UNL is currently 41. All these relations are binary and are expressed as *rel(UW<sub>1</sub>, UW<sub>2</sub>)*, where *UW<sub>1</sub>* and *UW<sub>2</sub>* are universal words or compound UW labels. A compound UW is a set of binary relations grouped together and regarded as one Universal Word. UWs are made up of a character string (usually an English-language word) followed by a list of restrictions. When used in UNL expressions, a list of attributes and often an instance ID follow these UWs.

**<UW>::=<Head Word>[<Constraint List>][":" <UW ID>][":" <Attribute List>]**

We explain the entities in the above BNF rule. The Head Word is an English word or a phrase or a sentence that is interpreted as a label for a set of concepts. This is also called *A Basic UW* (which is without restrictions). For example, the Basic UW *drink*, with no Constraint List, denotes the concepts of *putting liquids in the mouth*, *liquids that are put in the mouth*, *liquids with alcohol*, *absorb* and so on.

The constraint list restricts the interpretation of a UW to a specific concept. For example, the restricted UW *drink(icl>do, obj>liquid)* denotes the concept of *putting liquids into the mouth*. Words from different languages are linked to these disambiguated UWs and are assigned syntactic and semantic attributes. This forms the core of the lexicon building activity.

The UW ID is an integer, preceded by a “:”, which indicates the occurrence of two different instances of the same concept. The Constraint List can be followed by a list of attributes, which provides information about how the concept is being used in a particular sentence. A UNL Expression can also be expressed as a UNL graph. For example,

*John, who is the chairman of the company, has arranged a meeting at his residence.*

The UNL expressions for this sentence are as follows:

```

;===== UNL =====
;John who is the chairman of the company has arranged a meeting at his residence.
[S]
mod(chairman(icl>post):01.@present.@def,company(icl>institution):02.@def)
aoj(chairman(icl>post):01.@present.@def, John(icl>person):00)
agt(arrange(icl>do):03.@entry.@present.@complete.@pred,John(icl>person):00)
pos(residence(icl>shelter):04, John(icl>person):00)
obj(arrange(icl>do):03.@entry.@present.@complete.@pred,meeting(icl>conference):05.@indef)
plc(arrange(icl>do):03.@entry.@present.@complete.@pred,residence(icl>shelter):04)
[/S]
;=====

```

The UNL graph for the sentence is given in figure 1.

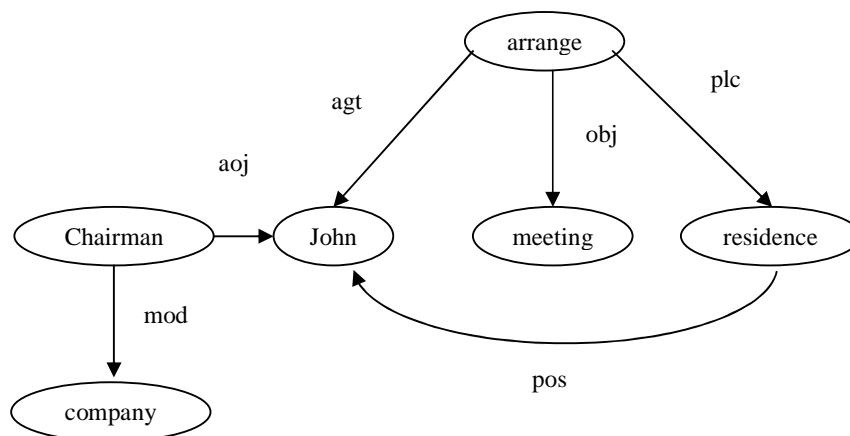


Figure 1: UNL graph

In the figure above, *agt* denotes the *agent* relation, *obj* the *object* relation, *plc* the *place* relation, *pos* is the *possessor* relation, *mod* is the *modifier* relation and *aoj* is the *attribute-of-the-object* (used to express constructs like *A is B*) relation.

The **international project on the Universal Networking Language** involves researchers from 14 countries of the world and includes 12 languages. For almost all the languages, the generator from the UNL expressions is quite mature. For the process of analysis into the UNL form, classical and difficult problems like ambiguity and anaphora are being addressed. All the research groups have to use the same repository of the universal words which is maintained by the UNDL foundation at Geneva and the UNU at Tokyo. When a new UW is coined by a research team it is placed in the UW repository at the UNU site. The restrictions are drawn from the *knowledge base* which again is maintained by the UNU. Individual teams have the responsibility of creating their *local language servers* which provide the services with respect to the analysis into and generation from UNL expressions.

The paper is organized as follows. The conceptual foundations, dealing with the formalisation of the UNL system and the universality of the lexicon, are given in section 2. Section 3 describes the use of lexical resources in semi-automatically constructing a semantically rich dictionary. Section 4 explains the working of the language independent analyser and generator tools as well as the actual Hindi and English Analysers and the Hindi generator. An overview of the major differences between Hindi and English is given in section 5. This is followed by a detailed description of the syntactic and lexical-semantic divergences between Hindi and English from a computational linguistics perspective in section 6. Section 7 describes our experiences in developing an MT system using the UNL. Section 8 deals with issues of disambiguation in the system. The paper ends with conclusions and future directions in section 8.

---

## 2 Conceptual Foundations

The strongest criticism against the interlingua based approach is that it requires the system designer to define a set of primitives which allow cross language mappings. This task is looked upon as a very hard one (Vauquois and Boitet 1985). (Wilks 1987) says,

*The notion of primitives in AI NL systems might be that they constitute not some special language, or another realm of objects, but are no more than a specialised sublanguage consisting of words of some larger standard language which plays a special organizing role in a language system.*

Since UNL is an interlingua we need to address this criticism. Rather than being based on primitives, the UNL system depends on a *large repository* of word concepts that occur in different languages. Such concepts are termed *Universal Words*. Thus words like *Ikebana* and *Kuchipudi* get included in this repository as *ikebana(icl>art form)* and *kuchipudi(icl>dance form)*. These word concepts are unambiguous, since every UW has a restriction which defines the sense of the basic UW used. For example, *spring* is a basic UW, which is disambiguated when it is restricted as *spring(icl>season)* meaning *spring* included in the class of *seasons*. The word concepts *spring* and *season* are ambiguous individually, but the combination *spring(icl>season)* is unambiguous. This can be further disambiguated as *spring(icl>(season(icl>time)))*.

No attempt is made in the UNL system to decompose *concepts (acts, objects, states and manner)* into primitives. A particular action, say *stab*, is represented using a single UW *stab(icl>do)*. This results in a representation that is more elegant and economical than some primitive based systems like Conceptual Dependency (Schank 1972, 1973, 1975).

### 2.1 Theoretical Background

UNL expressions are made of binary relations. The relation labels are designed to capture syntactic and semantic relations between Universal Words consistent with our knowledge of concepts (UWs) and gathered from the corpus of languages. The relations are chosen keeping in mind the following principles:

#### Principle 1) Necessary Condition

The necessary condition is something that characterizes separate relations: a relation is necessary, if one cannot do without it.

#### Principle 2) Sufficient Condition

The sufficient condition characterizes the whole set of relations: the set meets this condition if one need not add anything to it.

#### Explanation:

Let,

$U = \{UW_1, UW_2, \dots, UW_n\}$  be the UW Lexicon and  
 $C = \{C_1, C_2, C_3, \dots, C_m\}$  be the set of all possible contexts.

The set of relation labels  $\{RL_i\}$  in an interlingua IL defines functions of the following form:

$$RL_i : U \times U \rightarrow C$$

Let there be  $p$  such relation labels. We can call this set  $R$  where,

$$R = \{RL_1, RL_2, \dots, RL_p\}$$

---

Relating this to the UNL,  $RL_1$  could be *agt*,  $RL_2$  could be *obj*,  $RL_3$  could be *ins* and so on. Also concretely, contexts could be subsets of all possible sentences in all languages at all times. Each  $C_i$  is the set of *all sentences* in which each  $RL_i$  consists of tuples of the form,

$$\{((UW_{a_1}, UW_{a_2}), C_a), ((UW_{b_1}, UW_{b_2}), C_b), \dots\}$$

where, every  $((UW_{x_1}, UW_{x_2}), C_x)$  is unique across the members of the set R. Each  $C_x$  is the set of all possible sentences in which  $UW_{x_1}$  and  $UW_{x_2}$  appear. In this theoretical framework, contexts are language independent. Thus *John is driving a car* and *ज्होन गाड़ी चला रहा है* [John gaadi chalaah rahaa hai] belong to the same context  $C_q$ , say. From this definition it is clear what the necessity and sufficiency conditions mean.

The necessity condition implies that if a relation label  $RL_x$  is removed from the inventory the corresponding set,

$$\{((UW_{a_1}, UW_{a_2}), C_a), ((UW_{b_1}, UW_{b_2}), C_b), \dots\}$$

cannot be expressed in the IL.

Similarly sufficiency condition implies that if we add another relation  $RL_y$  then every element in the set  $RL_y$  will be present in some existing set  $RL_x$ .

The UNL expressions are binary and do not include the context information that has been referred to in the above discussion. Actually, the UNL reflects the context information through the **semantic types** of the UWs and the relation labels. For example, when we say *agt( $UW_1$ ,  $UW_2$ )*, it is clear that  $UW_1$  is an event of which the volitional entity  $UW_2$  is the agent. Thus, while encoding natural language sentences in the UNL, word and world knowledge will be used for implicitly capturing the context which has been described above in a hypothetical setting.

## 2.2 How Universal is the UW Lexicon?

An obvious question that arises for the UWs is *Why call these universal, since they are based on English?*. However, (Katz 1966) says:

*Although the semantic markers are given in the orthography of a natural language, they cannot be identified with the words or expressions of the language used to provide them with suggestive labels.*

This means that the primitives exist independently of the words used to describe, locate or interpret them. The UWs- though represented using Roman characters and English lexemes- are actually language independent concepts.

However, a problem arises when a group of words has to be used in a language whose lexical equivalent is a single word in another language. For example, for the Hindi word देवर [devar] the English meaning is *husband's younger brother*. Now, if we keep the universal word *husband's younger brother(icl>relative)* in the Hindi-UW dictionary and link it to देवर [devar], the analysis of the Hindi sentence H1 shown below will produce a set of UNL expressions in which the UW *husband's younger brother(icl>relative)* appears. From this set, an English language generator generates the sentence E1:

---

H1. लक्ष्मण सीता का देवर है <sup>2</sup>

laxman sita kaa devar hai<sup>3</sup>

Laxman Sita-of husband's-younger-brother-is

E1. Laxman is Sita's husband's younger brother.

Now, the English analyser, while analysing E1, will have the option of generating:

**aoj(young(icl>state).@comparative, brother(icl>relative))  
mod(brother(icl>relative), husband(icl>relative))**

**OR**

**husband's younger brother(icl>relative)**

*devar* was an example of conflation in noun for Hindi. As for verb, we can take औसाना [ausaanaa] which translates to English as *to ripen by covering in straw*. Thus *ausaanaa* has a conflational meaning. The UW for this could be

[औसाना] "**ripen(met>cover(ins>straw))**"

Now if the UNL expressions contain the words *ripen*, *cover* and *straw* separately, then it is a non-trivial problem for the generator to produce the conflated verb "ausaanaa". But if the above UW is used, then this can be done very easily.

One of the key assumptions about the UNL lexicon system is that the L-UW dictionaries should be usable *without change* in both analysis and generation. However, as is apparent from the discussion above, achieving this kind of universality is an idealisation.

A general decision taken in the present work is to introduce the language specific word as such in the UW dictionary, *if the corresponding English description is long-winded and cumbersome*. For example, we keep *kuchipudi(icl>dance)* in the dictionary instead of *an Indian dance form originating in the state of Andhra*. But, we do not keep *billi(icl>animal)*, where *billi* means a *cat* in Hindi, because *cat(icl>animal)* is available.

It should be noted that, the headwords are not always English words. English alphabets are USED to represent ALL the concepts which are found in ALL the languages at ALL times. Thus, *ikebana* and *kuchipudi* which are not English words are also stored in the dictionary. The disambiguation is done by a construct called the *restriction*. Restrictions are made of English alphabets. But they DO NOT DEPEND on English. The senses are not the ones which are peculiar to the English language. For example, one of the senses found in India of the word *back bencher* is *students who are not serious in their studies and while away their time sitting at the back of the class*. This additional sense is included in the UW dictionary as *back-bencher(icl>student)*. Thus if a particular word *w* in English has acquired an additional sense in another language, this sense is introduced into the UW dictionary by tagging the appropriate restriction. The words in specific languages get mapped to specific word senses and not to the basic UWs. The basic UWs are ambiguous and the linking process is carried out only after disambiguating.

We have given the example of *devar* (*husband's younger brother*) in Hindi. This illustrates the case where there is no direct mapping from Hindi to an English word. We have to discuss the reverse case where for an English word there is no direct mapping in another language. This is important since the UWs are primarily constructed from English lexemes. We have decided that if an English word is commonly used in Hindi, we keep the Hindi

---

<sup>2</sup> H[No.] indicates the Hindi sentence number and E[No.] the English sentence number. This is followed consistently through the paper.

<sup>3</sup> pronounce t as in *Taiwan* and T as in *Tokyo*

---

transliterated word in the dictionary. For example, for the word *mouse* used in the sense of an input device for the computer- we keep in the lexicon

[माउस] "mouse(icl>device)";

The same strategy is adopted if a word is very specific to a language and culture. For example, for the English word "blunderbuss" (an old type of gun with a wide mouth that could fire many small bullets at short range), there is no simple Hindi equivalent and so we keep in the lexicon the transliteration

[ब्लान्डरबस] "blunderbuss(icl>gun)";

The topic of multiple words for *snow* in Eskimo languages is very popular in NLP, MT and Lexical Semantics literature. We have discussed how to link these words with the appropriately formed UWs. In the Eskimo language *Inuit*, following are a few examples for the word *snow*:

'snow (in general)' *aput*, 'snow (like salt)' *pukak*, 'soft deep snow' *mauja*, 'soft snow' *massak*, 'watery snow' *mangokpok*.

The rich set of relation labels of UNL are exploited to form the UWs which in this case respectively are:

[aput] "snow(icl>thing)";  
[pukak] "snow(aj<salt like)";  
[mauja] "snow(aj<soft, aj<deep)";  
[massak] "snow(aj<soft)";  
[mangokpok] "snow(aj<watery)";

Note the disambiguating constructs for expressing the UWs. The relation labels of the UNL are used liberally. *aj* is the label for adjective-noun relation.

The issue of shades of meaning is a very important one, and the main idea again is that the RELATION LABELS OF UNL CAN BE USED IN THE LEXICON TOO. Here are some examples which have been added in the paper (the gloss sentences are attached for clarifying the meaning, which anyway gets communicated through the restrictions)

The verb *get off*:

[प्रस्थान करना] "get off(icl>leave)"; We got off after breakfast  
[बचना] "get off(icl>be saved)"; lucky to get off with a scar only  
[भेजना] "get off(icl>send)"; Get these parcels off by the first post  
[बन्ध करना] "get off(icl>stop)"; get off the subject of alcoholism  
[काम रोकना] "get off(icl>stop,obj>work)"; get off (the work) early tomorrow.

The noun *shadow*:

[अन्धेरा] "shadow(icl>darkness)"; the place was now in shadow  
[काले धब्बा] "shadow(icl>patch)"; shadows under the eyes.  
[परछाइ] "shadow(icl>atmosphere)"; country in the shadow of war  
[रन्धमाल] "shadow(icl>iota)"; not a shadow of doubt about his guilt  
[संकेत] "shadow(icl>hint)"; the shadow of the things to come  
[साया] "shadow(icl>close company)"; the child was a shadow of her mother  
[छाया] "shadow(icl>deterrant)"; a shadow over his happiness  
[शरण] "shadow(icl>refuge)"; he felt secure in the shadow of his father  
[आभास] "shadow(icl>semblance)"; shadow of power  
[भूत] "shadow(icl>ghost)"; seeing shadows at night

Again, note should be made of how the restrictions disambiguate and address the meaning shade.

### 2.3 Possibility of Representational Variations

Another important consideration while accepting UNL as an interlingua is the way it represents a particular sentence. UNL gives an unambiguous semantic representation of a sentence, but it does not claim uniqueness of the representation. Justifying the need for



---

primitives in an Interlingua, Hardt (Hardt 1987) says, *The requirement that sentences that have the same meaning be represented in the same way cannot be satisfied without some set of primitive ACTs.* This requirement may be a necessary condition for a knowledge representation scheme, but surely not for an Interlingua. For example, consider the following sentences:

- a. John gave a book to Mary.
- b. The book was given by John to Mary.
- c. Mary received a book from John.
- d. Mary took a book from John.

All these sentences have similar meanings, but are different from the point of view of the stylistics, focus and aspect. This is reflected in the UNL representation:

*John gave a book to Mary.*

```
[S]
agt(give(icl>do).@entry.@past, John(icl>person))
obj(give(icl>do) .@entry.@past, book(icl>text) .@def)
ben(give(icl>do) .@entry.@past, Mary(icl>person))
[/S]
```

*The book was given by John to Mary*

```
[S]
agt(give(icl>do) .@entry.@past, John(icl>person))
obj(give(icl>do) .@entry.@past, book(icl>text).@def.@topic)
ben(give(icl>do) .@entry.@past, Mary(icl>person))
[/S]
```

@*topic* is used for sentences in passive form to give more importance to the object than to the subject.

*Mary received a book from John.*

```
[S]
agt(receive(icl>do) .@entry.@past, Mary(icl>person))
obj(receive (icl>do) .@entry.@past, book(icl>text).@def)
src(receive (icl>do) .@entry.@past, John(icl>person))
[/S]
```

*Mary took a book from John.*

```
[S]
agt(take(icl>do) .@entry.@past, Mary(icl>person))
obj(take (icl>do) .@entry.@past, book(icl>text).@def)
src(take(icl>do) .@entry.@past, John(icl>person))
[/S]
```

Using these UNLs, a generator can generate an exact translation of the respective sentences and not its paraphrase, as it happens with CD based generators.

Although UNL represents similar information in different ways as above, its utility as a knowledge representation scheme does not get affected. Seniappan *et. al.* (Seniappan 2000) have investigated the use of UNL for automatic intra-document hypertext linking and have claimed that their system has an ability to extract anchors which are relevant but do not surface when frequency based methods are used.

As a summary of this section on conceptual foundations we mention the following points:

1. The UNL system strives to achieve language independence through its vast and rich repository of universal words.

- 
2. The basic UWs, *i.e.*, the unrestricted headwords, are mostly English words. But this does not make the UW dictionary an English language lexicon, since the concepts denoted by these UWs are valid for all languages.
  3. Whenever a language-specific word is cumbersome to express in English, the word is introduced into the UW repository after placing the proper restriction which clarifies the meaning of the particular UW and classifies it in a particular domain.
  4. The relation labels have stabilised to 41 and seem adequate to capture semantic relations between concepts across all languages. This is, however, only an empirical statement keeping in mind the necessity and the sufficiency conditions.
  5. A large portion of the burden of expressiveness in the UNL is carried by the attribute labels that indicate how the word is used in the sentence.
  6. The UW repository is the UNION of ALL concepts existing in ALL languages at ALL times.

### 3 L-UW Dictionary and The Universal Lexicon

In this section, we discuss the structure of a Language-UW (L-UW) Dictionary, its language dependent and independent parts and the associated attributes. The restriction attached with every word not only disambiguates it, but also puts it under a predefined hierarchy of concepts, called the *knowledge Base* in the UNL parlance. To construct the L-UW dictionary, the UWs are linked with the language words. Morphological, syntactic and semantic attributes are then added. For example, for the UW *dog(icl>mammal)*, the Hindi word कुत्ता [kutta](dog) is the language word, the morphological attribute is *NA* (indicating word ending with आ), the syntactic attribute is *NOUN* and the semantic attribute is *ANIMATE*. A part of the entry is

[कुत्ता] “dog(icl>mammal)” (NOUN, NA, ANIMATE);

The language independent part of this entry are *dog(icl>mammal)* and *ANIMATE*, while the language dependent parts are कुत्ता [kutta](dog) and *NA*. The same language-UW dictionary is used for the analysis and the generation of sentences for a particular language.

#### 3.1 The Architecture of the L-UW Development System

Figure 1 shows the architecture of the L-UW development system with both language dependent and language independent components. The Language independent parts are:

1. The Ontology Space.
2. The Set of UWs

The language dependent parts are:

1. The Language Specific Dictionary
2. The Syntactic and Morphological attributes

The process of L-UW dictionary construction can be partially automated. This achieves accuracy and exhaustiveness. Lexicon developers find it difficult to manually, consistently and exhaustively insert hundreds of semantic attributes required for the accurate analysis of the sentences. Also it is difficult to achieve uniformity in putting the restrictions. For example, for the noun *book*, a lexicon developer may restrict the meaning of *book* as *book(icl>concrete thing)*, *book(icl>textbook)*, *book(icl>register)*, *etc.*. This leads to a non-uniformity in the UWs which can be avoided by standardizing the knowledge base, *i.e.*, the UW repository. A brief description of the various components of the dictionary construction system now follows:

### 3.1.1 Language Independent Components

- **The Ontology Space**

The Ontology Space refers to a hierarchical classification of the word concepts. This Ontology is in the form of a Directed Acyclic Graph (DAG). Our system uses the upper CYC Ontology (Guha *et. al.* 1990) which has around 3000 concepts. This ontology is language independent and provides the semantic attributes.

- **The Set of UWs or the Knowledge base**

The Set of Basic UWs, *i.e.*, the unrestricted UWs are mostly the root words of English Language. Also, there are words from other languages, which do not have simple English equivalents, *e.g.*, *ikebana* from Japanese and *Kuchipudi* from Telugu. Basic UWs generally have more than one meaning. They are disambiguated by adding restrictions. These restricted UWs are language independent. A new knowledge base is in the process of being introduced and the UWs will be drawn from this resource.

### 3.1.2 Language Dependent Components

- **Language Specific Word Dictionary**

After selecting the UW, the corresponding language specific string is found by consulting the dictionary of the particular language and by translating the gloss attached.

- **Syntactic and Morphological Attributes**

This set includes attributes like *part of speech, tense, number, person, gender, etc.* and *morphological attributes* which describe *paradigms* of morphological transformations. These attributes are language specific and are inserted by the lexicon developer.

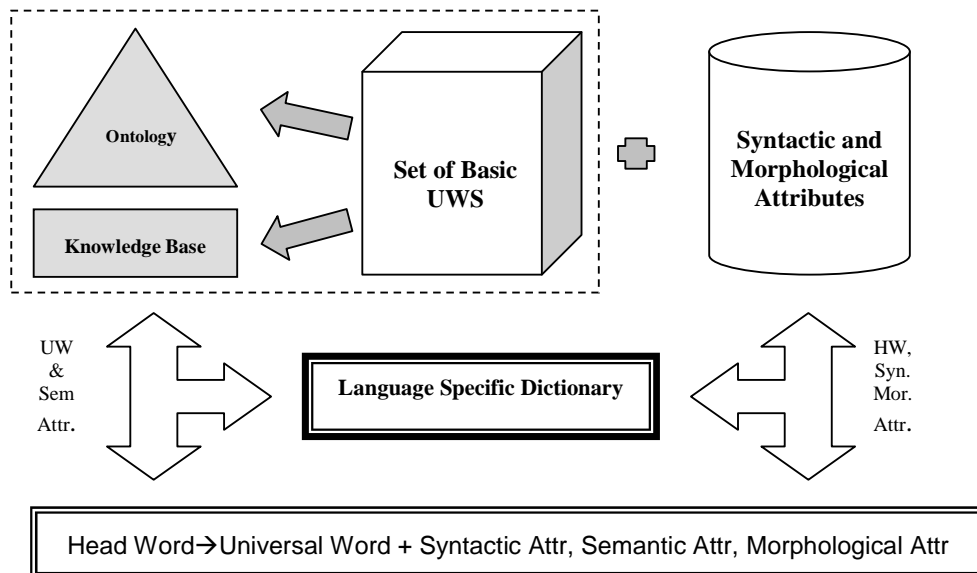


Figure 1: Integrated system for Language-UW Lexicon building

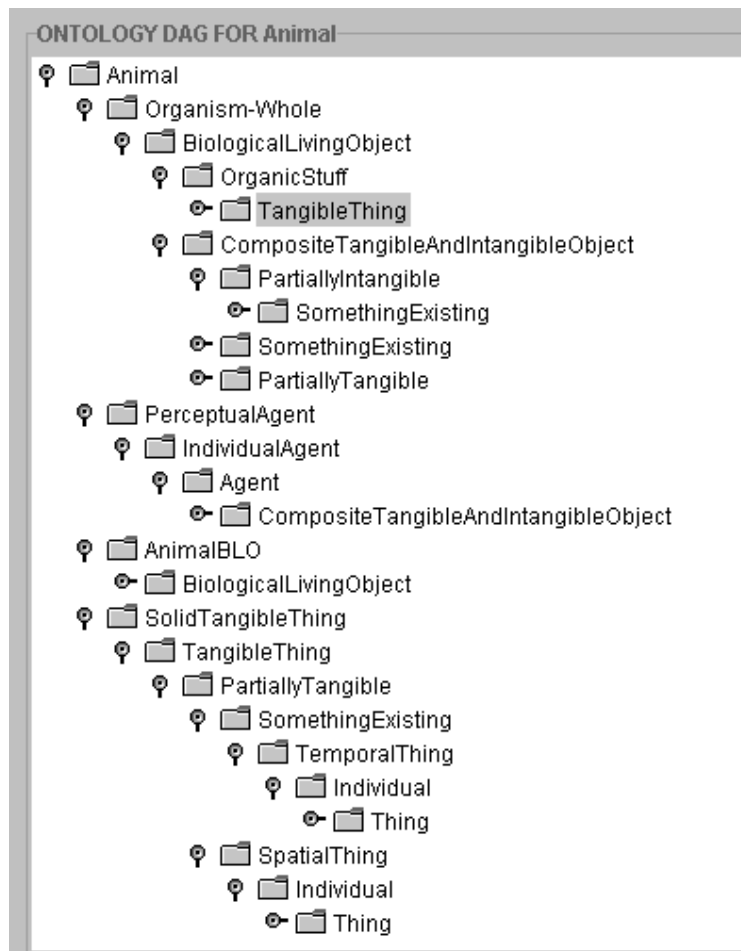
---

## 3.2 Constructing Dictionary Entries

The procedure of constructing dictionary entries is partially automated as follows:

1. The Human expert selects a UW from the knowledge base and finds for this sense the position of the basic UW (the portion left after stripping the restriction) as a leaf in the Ontology. Consider a snapshot of the CYC Ontology DAG given in Figure 2. Suppose we want to make a dictionary entry for the word *animal*. The word is found as a leaf in the Ontology. The UW is *animal (icl>living thing)*.
2. The semantic attributes of this UW are the nodes traversed while following all paths from the leaf to the root (*thing* in this case). For example, the following attributes are generated for the word *Animal*:

**SolidTangibleThing, TangibleThing, PartiallyTangible, PartiallyIntangible, CompositeTangibleAndIntangibleObject, AnimalBLO, BiologicalLivingObject, PerceptualAgent, IndividualAgent, Agent, Organism-Whole, OrganicStuff, SomethingExisting, TemporalThing, SpatialThing, Individual, Thing**



**Figure 2: A Snapshot of Cyc Upper Level Ontology**

3. The work of the human expert is now limited to adding the syntactic and morphological attributes. These attributes are far less in number than semantic attributes. Thus, the labour of making semantically rich dictionary entries is reduced.

---

An example of a dictionary entry generated by the above process is:

प्राणी{ }”animal(icl>organism whole)”(Noun, NI, SolidTangibleThing, TangibleThing, PartiallyTangible, PartiallyIntangible, CompositeTangibleAndIntangibleObject, AnimalBLO, BiologicalLivingObject, PerceptualAgent, IndividualAgent, Agent, Organism-Whole, OrganicStuff, SomethingExisting, TemporalThing, SpatialThing, Individual, Thing)

प्राणी [praanee](animal) is the Hindi equivalent for *animal*. *Noun* and *NI*<sup>4</sup> are the syntactic and morphological attributes added by the human lexicon developer.

## 4 The System

We describe here the systems we built, *viz.*, the Hindi Analyser (*HA*) which converts Hindi sentences into UNL expressions, the English Analyser (*EA*) which produces UNL expressions from English sentences and the Hindi Generator (*HG*) which generates Hindi sentences from UNL expressions. The Analysers use a software called the *EnConverter* while the Generator uses the *DeConverter*<sup>5</sup>. These tools are language independent systems which are driven by the language dependent rule-base and the L-UW dictionaries. We first give an overview of the working of the *EnConverter* and *DeConverter* engines. Then we explain in brief the three systems. Space restriction does not permit detailed description of all the three systems.

### 4.1 The Analyser Machine

The *EnConverter* is a language independent analyser which provides a framework for morphological, syntactic and semantic analysis synchronously. It analyses sentences by accessing a knowledge rich **L-UW lexicon** and interpreting the **Analysis Rules**. The process of formulating the rules is in fact programming a sophisticated symbol-processing machine.

The *EnConverter* can be likened to a multi-head Turing machine. Being a Turing Machine, it is equipped to handle phrase structured (*type 0*) grammars (Martin 1991) and consequently the natural languages. The *EnConverter* delineates a sentence into a tree- called the nodenet tree- whose traversal produces the UNL expressions for the sentence. **During the analysis, whenever a UNL relation is produced between two nodes, one of these nodes is deleted from the tape and is added as a child of the other node to the tree.** It is important to remember this basic fact to be able to understand the UNL generation process in myriad situations.

The *EnConverter* engine has two kinds of heads- *processing heads* and *context heads*. There are two processing heads- called *Analysis Windows*. The nodes under these windows are processed for linking by a UNL relation label and/or for attaching UNL attributes to. A node consists of the language specific word, the universal word and the attributes appearing in the dictionary as well as in the UNL expressions. The context heads are located on either sides of the processing heads and are used for look ahead and look back. The machine has functions like *shifting the windows right or left by one node*, *adding a node to the node-list* (tape of the Turing machine), *deleting a node*, *exchange of nodes under processing heads*, *copying a node* and *changing the attributes of the nodes*. The complete description of the structure and working of the *EnConverter* can be found in (*EnConverter 2000*).

---

<sup>4</sup> *NI* indicates that the noun ends with a *I* (Romanised hindi). This information helps in morphological analysis.

<sup>5</sup> *EnConverter* and *DeConverter* are tools provided by the UNL Project, Institute for Advanced Studies, United Nations University, Tokyo (*EnConverter 2000*).

---

## 4.2 The English Analyser

The English Analyser makes use of the English-UW dictionary and the rule base for English Analysis, which contains rules for morphological, syntactic and semantic processing. At every step of the analysis, the rule base drives the EnConverter to perform tasks like completing the morphological analysis (e.g., combine *Boy and 's*), combining two morphemes (e.g., *is* and *working*) and generating a UNL expression (e.g., *agt* relation between *he* and *is working*). Many rules are formed using Context Free (CFG)-like grammar segments, the productions of which help in clause delimitation, prepositional phrase attachment, part of speech (pos) disambiguation and so on. This is illustrated with the example of noun clause handling:

*The boy who works here went to school.*

Example Grammar:

CL → V	; e.g., <i>The boy who works ...</i>
ADV V N	; e.g., <i>The boy who <b>fluently speaks English</b></i>
V ADV	; e.g., <i>The boy who <b>works here</b></i>
V ADV ADV	; e.g., <i>The boy who <b>ran very quickly</b></i>

The processing goes as follows.

1. The clause *who works here* starts with a relative pronoun and its end is decided by the system using the grammar. There is no rule like

$CL \rightarrow V ADV V$

and so the system does not include *went* in the subordinate clause.

2. The system detects *here* as an adverb of place from the lexical attributes and generates *plc* (*place* relation) with the main verb *work* of the subordinate clause. After that, *work* is related with *boy* through the *agt* relation. At this point the analysis of the clause finishes.
3. *boy* is now linked with the main verb *went* of the main clause. Here too the *agt* relation is generated.
4. The main verb is then related with the preposition phrase to generate *plt* (indicating *place to*), taking into consideration the preposition *to* and the noun *school* (which has PLACE as a semantic attribute in the lexicon). The analysis process thus ends.

A typical example of the ability of the system for *part of speech* (*pos*) disambiguation is shown below:

```
===== UNL =====
The soldier went away to the totally deserted desert to desert the house in the desert
[S]
mod(deserted(icl>vacant):11,total(icl>complete):0T)
aoj(deserted(icl>vacant):11,desert(icl>landscape):1A.@def)
plc(go(icl>event):0C.@entry.@past.@pred,away(icl>logical place):0H)
obj(desert(icl>do):1K.@present.@pred,house(icl>place):1V.@def)
plc(desert(icl>do):1K.@present.@pred,desert(icl>landscape):28.@def)
plt(go(icl>event):0C.@entry.@past.@pred,desert(icl>landscape):1A.@def)
pur(go(icl>event):0C.@entry.@past.@pred,desert(icl>do):1K.@present.@pred)
agt(go(icl>event):0C.@entry.@past.@pred,soldier(icl>human):04.@def)
[/S]
;=====
```

The adjectival form of *desert* is represented as *deserted(icl>vacant)*. The noun form is *desert(icl>landscape)*, while the verb form is *desert(icl>do)*. The analysis rules make use of the linguistic clues present in the sentence. Thus, the adverb *totally* preceded by the article *the* makes the *desert+ed* an adjective, which in turn makes the following *desert* a noun.

---

The system can also convert sentences in which relative pronouns do not occur in the sentence explicitly. For example,

1. The study (which was) published in May issue was exhaustive.
2. He lives at a place (where) I would love to be at.
3. He gave me everything (that) I asked for.
4. The cabbage (which was) fresh from the garden was tasty.

Various heuristics are used to decide the start of clause and the relative pronoun that is implicit. Some of these are:

- Presence of two verbs with a single subject as in 1.
- A noun followed by a pronoun as in 2.
- Quantifiers like *all*, *everything* and *everyone* followed by another pronoun or noun as in 3.
- An adjective following a noun as in 4.

Semantic attributes stored in the dictionary are exploited to solve ambiguities of prepositional phrase and clausal attachment. For example,

*He went to my home **when I was away**.*  
*He met me at a **time when I was very busy**.*

The structures of the two sentences are similar, but semantic attributes indicate that *when* qualifies temporal nouns like *time*, *hour*, *second*, etc. Thus, in the first sentence the system attaches the clause *when I was away* to the verb considering it an adverb clause of time, while in the second it attaches the clause *when I was very busy* to the noun considering it an adjective clause.

Anaphora resolution is dealt with in a limited way at the sentence level. This can be seen from the UNL expressions produced by the system for the sentence given below:

```
===== UNL =====  
;He built his house in a very short span of time.  
[S]  
mod(house(icl>place):0D, he(icl>person):09)  
agt(built(icl>event):03.@entry.@past.@pred, he(icl>person):09)  
mod(short(icl>less):0T,very:0O)  
aoj(short(icl>less):0T,span(icl>duration):0Z.@indef)  
obj(built(icl>event):03.@entry.@past.@pred, house(icl>place):0D)  
dur(built(icl>event):03.@entry.@past.@pred,span(icl>duration):0Z.@indef)  
mod(span(icl>duration):0Z.@indef,time(icl>abstract thing):AB)  
[S]  
=====
```

The UW-IDs (a form of identifier) of both the instances of *he(icl>person)* in the above sentence are the same, viz., :09. The system does not do the same for the sentence *John built his house*, since it is not certain whether *John* and *he* refer to the same person.

Ellipses handling is done for various kinds of sentences. Few examples are given below:

1. I reached there before he could (reach).
2. (I am) Sorry, I did it.
3. I went to Bombay and then (I went) to Delhi.

For the first sentence, the implicit *reach* is produced explicitly in the UNL expressions. The second sentence obviously does not generate an extra *I*, but adds the attribute *@apology* to the verb *do*. Since there are two events of *going* in the third sentence, an explicit *go* is produced but not the extra *I* as the agent is the same for both the instances of *go*.

---

Thus, the EA is capable of handling many complex phenomena of the English language. The system also can guess a UW for a word not present in the lexicon. Currently, it has around 5800 rules. A detailed explanation of the system can be found in (Parikh 2000, 2001).

### 4.3 The Hindi Analyser

The rule base that drives the Hindi Analyser (HA) uses strategies different from its English counterpart. This is due to the numerous structural differences between Hindi and English (*vide* section 5). But the fundamental mechanism of the system is the same, *i.e.*, it performs morphological, syntactic and semantic analysis synchronously.

The rule base of the HA can be broadly divided into three categories – *morphological rules*, *composition rules* and *relation resolving rules*. Morphology rules have the highest priority. This is because unless we have the morphed word, we cannot decide upon the part of speech of the word and its relation with the adjacent words. Hindi has a rich morphological structure. Information regarding person, number, tense and gender can be extracted from the morphology of nouns, adjectives and verbs. An exhaustive study of the morphology is done for this purpose and appropriate rules are incorporated into the system (Monju *et. al.* 2000). To illustrate the process of Hindi analysis, we consider the following example of a Hindi sentence with an explicit pronoun.

H2. मैंने देखा कि सीता सब्जी खरीद रही है।  
mai ne dekhaa ki seetaa sabjee khareed rahee hai  
I saw that sita vegetable buying-is

E2. I saw that Sita is buying vegetables.

The processing of this sentence is carried out as follows:

1. The beginning of the clause is marked by the presence of the relative pronoun *ki* (that).
2. The analysis windows right shift till the predicate *dekhaa* is reached.
3. All the relations of the previous nodes with this predicate are resolved. In this case, *mai* (I) being *first person singular* and *animate* pronoun, *agt* relation is produced between *mai ne* and *dekhaa*.
4. The relative pronoun *ki* is now detected and the analysis heads right shift. It combines *ki* with *dekhaa* and adds a dynamic attribute *kiADD* to *dekhaa*.
5. The clause following *ki* is now resolved. The analysis windows right shift till the main predicate of the sentence- *khareed rahee hai*- is reached.
6. It combines the nodes *sabjee* and *khareed rahee hai* with the *obj* relation seeing the *inanimate* attribute of *sabjee*.
7. It then resolves the *agt* relation between *seetaa* and *khareed rahee hai* seeing the *animate* attribute of *seetaa*.
8. At the end of its analysis, its main predicate is retained which in this case is *khareed rahee hai*. Finally the *obj* relation is generated between this verb and *dekhaa*.

Composition rules are used to combine a noun or a pronoun in a sentence with a postposition or case-marker following it. During combination, the case marker is deleted from the Node-list and appropriate attributes are added to the noun or pronoun to retain the information that the particular noun or pronoun had a postposition marker following it. For example, consider the following sentences:



- H3. राम ने रावण को तीर से मारा ।  
raam ne raavan ko teer se maaraa  
Ram Ravan-to arrow-with killed  
 E3. Ram killed Ravan with an arrow.
- H4. पेड़ से पत्ते बाग में गीरे ।  
ped se patte baag mein geeere  
Tree-from leaves garden-in fell  
 E4. Leaves fell in the garden from the trees.
- H5. पीटर सुबह से काम कर रहा है ।  
peeTar subah se kaam kar rahaa hai  
Peter morning-since working-is  
 E5. Peter is working since morning.
- H6. बच्चे से ताला खुला ।  
bachche se taalaa khulaa  
child-by lock opened-was  
 E6. The lock was opened by the child.

In the above sentences, तीर [*teer*](arrow), पेड़ [*ped*](tree), सुबह [*subah*](morning) and बच्चा [*bachchaa*](child) are nouns and are followed by the same postposition marker से [*se*](with, from, since, by). However, as it is evident from the English translation, the meaning of से [*se*] is different in each sentence viz. with, from, since and by respectively. Hence, the noun preceding it forms a different relation with the main verb in each case as follows.

1. **ins(kill(icl>do)).@past, arrow(icl>thing)**
2. **plf(fall(icl>occur)).@past, tree(icl>place)**
3. **tmf(work(icl>do)).@present,@progress, morning(icl>time)**
4. **agt(open(icl>do)).@past, child(icl>person)**

These nouns have the semantic attributes *INSTRU* (can be used as an instrument), *PLACE*, *TIME* and *ANI* (animate entities) respectively in the lexicon. They help deciding upon the sense of the case marker and thus the role of the noun in the particular sentence. When the case marker से [*se*] is combined with the noun preceding it, attributes *INS*-instrument, *PLF* - place from which an event occurs, *TMF*- time from which an event has started and *AGT*-agent of the event, are added to the respective nouns. These attributes then lead to the production of the above UNL relations for the respective sentences.

Now we describe the various Hindi language phenomena handled by the system. Hindi is a null subject language [*vide* section 6.1.4]. This means that it allows the syntactic subject to be absent. For example, the following sentence is valid in Hindi.

- H7. जा रहा हूँ ।  
jaa rahaa hun  
going-am  
 E7. \*am going<sup>6</sup>

<sup>6</sup> \* indicates incorrect grammatical construct

---

The system makes the implicit subject explicit in the UNL expressions. The procedure to do this is discussed in section 6.1.4. The UNL expression produced by the system in this case is:

```
[S]
agt(go(icl>do).@entry.@present.@progress, I(icl>person))
[/S]
```

The system can also handle limited amount of Anaphora resolution. For example, consider the following sentence:

H8. मेरी ने अपनी किताब जीम को दी है।  
meree ne apanee kitaab jeem ko dee hai  
Mary her book Jim-to given-has

E8. Mary has given her book to Jim.

The corresponding UNL relations generated are:

```
[S]
pos(book(icl>publication):0C,Mary(icl>person):00)
ben(give(icl>do):0R.@entry.@present.@pred,Jim(icl>person):0J)
obj(give(icl>do):0R.@entry.@present.@pred,book(icl>publication):C)
agt(give(icl>do):0R.@entry.@present.@pred,Mary(icl>person):00)
[/S]
```

That resolution of the anaphora is apparent from the fact that the UW *she(icl>person)* for *her* is replaced by *Mary(icl>person)* in the *pos* relation.

One of the major differences between Hindi and English is that a single pronoun वह [vah](*he* or *she*) in Hindi is mapped to two pronouns *he* and *she* of English. The gender of the pronoun in Hindi can be known only from the verb morphology. So the system defers the generation of the UW for वह [vah](*he* or *she*) until the verb morphology is resolved. At the end of the analysis, the correct *he(icl>person)* or *she(icl>person)* is produced. For example,

H9. वह शाम को आएगी।  
vah shaam ko aaegee  
She evening-in will come

E9. She will come in the evening.

The UNL expressions are:

```
[S]
tim(come(icl>do):0D.@entry.@future,evening(icl>time):05.@def)
agt(come(icl>do):0D.@entry.@future,she(icl>person):00)
[/S]
```

Hindi uses the word-forms आएगा [aaegaa] and आएगी [aaegee](both meaning *will come*) for the verb आ [aa] (come) for a male subject and female subject respectively. Thus, in the above sentence, the verb आएगी [aaegee] causes the UW *she(icl>person)* to be generated for वह [vah](*he* or *she*).

Hindi being a relatively free word-ordered language, the same sentence can be written in more than one way by changing the order of words. For example,

H10. (A) तुम कहाँ जा रहे हो?  
tum kahaan jaa rahe ho?  
You where going are

- (B) कहाँ तुम जा रहे हो?  
kahaan tum jaa rahe ho?  
where you going-are
- (C) कहाँ जा रहे हो तुम?  
kahaan jaa rahe ho tum?  
where going-are you

E10. Where are you going?

The output in all cases is:

[S]  
**plc(go(icl>do):07.@entry.@interrogative.@pred.@present.@progress,**  
**where(icl>place):00)**  
**agt(go(icl>do):07.@entry.@interrogative.@pred.@present.@progress, you(icl>male):0I)**  
 [/S]

This is achieved as follows. Additional rules are added for each combination of the word types. Also the rules are prioritised such that the right rules are picked up for specific situations. For the sentence H10(A), first the rule for generating *plc* relation between *kahaan* and *jaa rahe ho* is fired, followed by the rule for generating *agt* relation between *tum* and *jaa rahe ho*. In H10(B), first *agt* and then *plc* are resolved. In H10(C), a rule first exchanges the positions of *jaa rahe ho* and *tum*. After that the rules fire as before for setting up the relations. Use is made of the question mark at the end of the sentence.

Hindi allows two types of constructions for adjective clauses– one with explicit clause markers like जो [jo](*who*), जिसकी [jisakee](*whose*), जिसे [jise](*whom*), etc. and the other with the वाला [vaalaa](*ing*) construction. Our analyser can handle both. For example,

- H11. पीटर जो लंडन में रहता है वह यहाँ काम करता है।  
peeTar jo london mein rahataa hai vah yahaan kaam karataa hai  
Peter who London-in stays he here work-do-is

E11. Peter who stays in London works here.

- H12. लंडन में रहनेवाला पीटर यहाँ काम करता है।  
london mein rahanevaalaa peeTar yahaan kaam karataa hai  
London-in staying Peter here work-do-is

E12. Peter who stays in London works here.

The system produces the following UNL relations for both these:

[S]  
**agt(work(icl>do).@entry.@present, Peter(icl>person))**  
**plc(work(icl>do) .@entry.@present, here)**  
**agt(stay(icl>do) .@present, Peter(icl>person))**  
**plc(stay(icl>do) .@present, London(icl>place))**  
 [/S]

The two incoming arrows into *Peter(icl>person)* provides the clue to the system to correctly identify the adjective clause in each sentence.

---

Unlike English, Hindi has a way of showing respect to a person (*vide* Section 5). This is conveyed through the verb morphology. For example,

H13. मेरे चाचा पढ़ रहे हैं।

mere chaachaa padh rahe hai

my uncle reading-are

E13. My uncle is reading.

The verb form here is for the subject in plural form. But since *uncle* is singular, the system infers that the speaker is showing respect and generates @respect attribute for *uncle(icl>person)*.

The HA can deal with simple, complex, compound, interrogative as well as imperative sentences. Currently the number of rules in HA is about 3500 and the lexicon size is around 70,000.

#### 4.4 The Generator Machine

The DeConverter is a language independent generator which provides a framework for morphology generation and syntax planning synchronously. It generates sentences by accessing a knowledge rich L-UW dictionary and interpreting the Generation Rules.

The working and the structure of the DeConverter are very similar to that of the EnConverter. It processes the UNL expressions on the input tape. It traverses the input UNL graph and generates the corresponding target language sentence. Thus, during the course of the generation, whenever a UNL relation is resolved between two nodes, one of the nodes is inserted into the tape.

Like the EnConverter, the DeConverter also has two types of heads- *processing heads* and *context heads*. There are two processing heads- called *generation windows*-, and only the nodes under these take part in any generation tasks like the left or right placement of the words and the resolution of attributes into morphological strings. The context heads- called the *condition windows*- are located on either sides of the processing heads and are used for look ahead and look back. The machine has functions of *shifting right or left by one node*, *adding a node to the node-list (tape of the Turing machine)*, *deleting a node*, *exchange of nodes under processing heads*, *copying a node* and *changing attributes of the nodes*. The complete description of the structure and working of the DeConverter can be found in (DeConverter 2000).

#### 4.5 Hindi Generator

The HG attempts to generate the most natural Hindi sentence from a given set of UNL expressions. The generation process is based on the predicate-centric nature of the UNL. It starts from the UW of the main predicate and the entire UNL graph is traversed in stages producing the complete sentence. The rule base contains the syntax planning rules and the morphology rules. Syntax planning is in general achieved with a very high degree of accuracy using two fundamental concepts called *parent-child relationships* and *Matrix based priority of relations* (Rayner 2001).

In a UNL relation  $rel(UW_1, UW_2)$ , the  $UW_1$  is always the parent node and  $UW_2$  the child. The syntax planning task is to decide upon the right or left insertion of the of the child with respect to its parent. The UNL specification puts constraints on the possible types of UWs that can occur as  $UW_1$  and  $UW_2$  of a particular relation. Using this information and the relation between the two UWs, the position of the child relative to the parent is arrived at.

Another important consideration is the traversal of the UNL graph. The path is decided based on the relative priority of UNL relations which is in turn decided by the *priority matrix*. An example matrix is given in Table 1.

	Agt	obj	ins
agt	-	L	L
obj	R	-	R
ins	R	L	-

**Table 1: An example priority matrix**  
where, L means *placed-left* and R means *placed-right*

This matrix is read as:

**agt placed-left-of obj OR obj placed-right-of agt**  
**agt placed-left-of ins OR ins placed-right-of agt**  
**ins placed-left-of obj OR obj placed-right-of ins**

Such an exhaustive matrix is produced for all the 41 relations.

According to the above matrix,

*child(agt)* is the leftmost element,  
*child(ins)* is the middle element and  
*child(obj)* is the rightmost element of the three

For example, consider the following UNL expressions

[S]  
agt(eat(icl>do).@entry.@past, Mary(icl>person))  
ins(eat(icl>do).@entry.@past, spoon(icl>thing).@indef)  
obj(eat(icl>do).@entry.@past, rice(icl>food))  
[/S]

The sentence generated according to the above matrix is,

H14. मेरी ने चम्मच से चावल खाया ।  
meree ne chammach se chaaval khaayaa  
Mary spoon-with rice ate

E14. Mary ate the rice with a spoon.

The rule writer uses the above matrix to decide upon the priorities of the rules. The relation for which the child is placed leftmost in the sentence has the highest priority and is resolved first, while the relation for which the child is placed rightmost, *i.e.*, nearest to the verb, has the lowest priority.

Morphology generation not only transforms the target language words for each UW, but also introduces case markers, conjunctions and other morphemes according to the relation labels- a procedure reified as *relation label morphology*. Table 2 gives an idea of this process. UNL attributes reflecting the aspect, tense, number, *etc.* also play a major role in the morphology processing.

The HG can produce both complex and compound sentences. The presence of a clause in the sentence is detected in two different ways: (i) presence of a *scope*, *i.e.*, a *compound universal word which is a label for more than one UNL expressions* or (ii) presence of two incoming arrows from two different predicates. For example, *He scolded the boy who had hit John* can be represented in the UNL in two different ways:

Relation M	Position of the word wrt child(M)	Word to be introduced
Agt	L	ने [ne]
And	R	और [aur](and)
Bas	L	से [se](as compared to)
Cag	L	के साथ [ke saath](with)
Cob	L	के साथ [ke saath](with)
Con	L	यदि [yadi]UW2 तो [to] UW <sub>1</sub> (if UW2 then UW <sub>1</sub> )
Coo	R	और [aur](and)/ null
fnt	R	से [se](to)
Gol	L	में [mein](into)
Ins	L	से [se](using)
Mod	L	का [kaa](of) / के [ke](of) / की [kee](of) / null (depends on gender and number)

Table 2: Relation Label Morphology

[S]  
**agt(scold(icl>do).@past.@entry, he(icl>person))**  
**obj(scold(icl>do).@past.@entry, boy(icl>person))**  
**agt(hit(icl>do).@pred.@complete.@past, boy(icl>person))**  
**obj(hit(icl>do).@pred.@complete.@past,John(icl>person))**

[/S]

**OR**

[S]  
**agt(scold(icl>do).@past.@entry, he(icl>person))**  
**obj(scold(icl>do).@past.@entry, :01)**  
**agt:01(hit(icl>do).@pred.@complete.@past.@entry, boy(icl>person))**  
**obj:01(hit(icl>do).@pred.@complete.@past.@entry,John(icl>person))**

[/S]

In the first representation, *boy(icl>person)* has two incoming arrows from *scold(icl>do)* and *hit(icl>do)*. The second representation explicitly marks the presence of the clause using the scope :01. The system generates the same sentence for both representations.

The HG is also capable of handling imperative, passive and interrogative sentences. The current system has around 5000 rules and uses the same Hindi-UW dictionary used by the Hindi Analyser.

## 5 Major Differences between Hindi and English

The basic difference between Hindi and English is the sentence structure. Hindi has a Subject-Object-Verb (SOV) structure for sentences, while English follows the Subject-Verb-Object (SVO) order. (Rao *et. al.* 2000) gives the following structure for English sentences:

**S S<sub>m</sub> V V<sub>m</sub> O O<sub>m</sub> C<sub>m</sub>**

where,

S: Subject

O: Object

V: Verb

S<sub>m</sub>: subject post-modifiers

O<sub>m</sub>: object post-modifiers

V<sub>m</sub>: the expected verb post-modifiers

C<sub>m</sub>: the optional verb post-modifiers

---

For example,

E15. The President of America will visit the capital of Rajasthan in the month of December.

(S) (S<sub>m</sub>) (V) (O) (O<sub>m</sub>) (C<sub>m</sub>)

On the other hand, Hindi has the following structure:

**C<sub>m</sub> S<sub>m</sub> S O<sub>m</sub> O V<sub>m</sub> V**

H15. दिसम्बर के महिने में अमरिका के राष्ट्रपति राजस्थान की राजधानी की सैर करेंगे।

disambar ke mahine mein amarikaa ke raashtrapati raajasthaan kee raajadhaani kee sair karenge

(C<sub>m</sub>) (S<sub>m</sub>) (S) (O<sub>m</sub>) (O) (V)

December-of month-in America-of President Rajasthan-of capital-of tour will-do

The morphological variations are richer in Hindi than in English. The case markers में [mein], से [se], को [ko], का [kaa] *etc.* are placed *postposition* and are strongly bound to the nouns. This allows Hindi to be a relatively free-word order language. English uses prepositional phrases as complements and qualifiers, and the order of the words is quite *fixed*.

The free word ordering, however, poses difficulties in the analysis of the Hindi sentences. In addition to the phrase and clause attachment problems, it also makes the task of distinguishing the clauses and phrases from the subject and object of the sentence difficult, as they all have case markers and can be placed anywhere in the sentence. For example,

H16. (A) राम ने चोरी करनेवाले लड़के को लाठी से मारा।

raam ne choree karanevaale ladake ko laathee se maaraa

Ram stealing boy-to stick-with hit.

(B) राम ने लाठी से चोरी करनेवाले लड़के को मारा।

raam ne laathee se choree karanevaale ladake ko maaraa

Ram stick-with stealing boy-to hit.

(C) चोरी करनेवाले लड़के को राम ने लाठी से मारा।

choree karanevaale ladake ko raam ne laathee se maaraa

stealing boy-to Ram stick-with hit.

(D) चोरी करनेवाले लड़के को लाठी से राम ने मारा।

choree karanevaale ladake ko laathee se raam ne maaraa

stealing boy-to stick-with Ram hit

(E) लाठी से राम ने चोरी करनेवाले लड़के को मारा।

laathee se raam ne choree karanevaale ladake ko maaraa

stick-with Ram stealing boy-to hit

E16. Ram hit with a stick the boy who had stolen.

Here, राम ने [raam ne] (Ram) denotes the agent, लड़के को [ladake ko](boy-to) the object and लाठी से [laathee se] the instrument. चोरी करनेवाले [choree karanevaale] is a clause qualifying लड़के

---

[ladakaa](boy). Relative positions of each of these phrases can be varied as is apparent from the sentences H16(A-E).

However, the postposition markers in Hindi always stay next to the nouns they modify and also have comparatively fixed roles. This partially compensates for the extra processing arising from the free word ordering.

English overloads the prepositions. For the UNL generation, not only the PP attachment but also the semantic relation of the PP with the noun or the verb should be determined. For example,

1. John ate rice with curd. → **cob(eat(icl>do).@entry.@past, curd(icl>food))**
2. John ate rice with a spoon. → **ins(eat(icl>do).@entry.@past, spoon(icl>thing).@indef)**
3. John ate rice with Mary. → **cag(eat(icl>do).@entry.@past, curd(icl>food))**
4. The Demon ate the rice with the goat. → **cob or cag??**

In the above sentences, the PPs starting with *with* have different roles. In the first sentence, the relation is *co-object*, in second it is *instrument* and in third it is *co-agent*. It is difficult to decide whether *goat* in the fourth sentence is a *co-object* or a *co-agent*! The system identifies these relations using the semantic attributes of the nouns placed in the lexicon. This analysis is explained in detail in (Parikh 2000).

Hindi is a Null-Subject language, while English is not. Null-Subject languages allow subjects to be dropped when the meaning is clear. H7 is an example of a Hindi sentence where the subject is dropped. Null-Subject languages do not have pleonastics. This phenomenon is discussed in section 6.1.4.

A very important feature of the Hindi language is that of *Conjunct and Compound Verbs* which are formed by combining two or more verbs or by combining a noun or an adjective or an adverb with verbs like कर (do) or हूँ (be). In the case of conjunct verbs, the first verb is usually the main one and the other is the subsidiary. All transformations of *Voice, Mood, Tense, Person, Gender* and *Number* affect the Subsidiary Verb only. The following sentences exemplify this:

H17. वह गाने लगी।  
vaha gaane lagee  
She singing started

E17. She started singing.

H18. हम गाने लगेंगे।  
ham gaane lagenge  
We singing will-start

E18. We shall start singing.

The following sentences show some of the interesting ways the verb जा [jaa](go) is used to emphasise or intensify the effects of the main verb. The literal translations show only the most common meanings of the constituent verbs:

H19. चले जाओ।  
chale jao  
walk go

E19. Go away.



---

H20. रुक जाओ ।

ruk jaao

stop go

E20. Stop there.

H21. झुक जाओ ।

jhuk jaao

bend go

E21. Bend down.

The phenomenon of compounding of verbs is a typical Indian language phenomenon. The **strategy** to deal with this, however, is quite simple. The presence of two verbs next to each other provides the clue that the second verb is the intensifier, and generally the UNL expression produced gets the attribute *@emphasis* attached to the first verb. For example,

**agt(go(icl>do).@entry.@imperative.@emphasis, you(icl>person))**

There are numerous lexical and syntactic differences between Hindi and English. Some of them are as under:

#### A. Number

Some words in English are always used in plural form, for example, *scissors*. This phenomenon does not occur in Hindi. It is impossible to determine from the second single sentence below whether the reference is to one or more *scissors*.

*The company manufactures scissors. (many)*

*The scissors are very sharp. (one or many)*

In Hindi, there are two different morphological forms for *scissors*- कैंची [kainchee] and कैंचियाँ [kainchiyaan](plural), and thus this problem does not arise.

In English, some words have a single meaning in the singular form and multiple meanings in the plural. For example, the word *premise* means an *assumption*, while the word *premises* mean *assumptions* or *the place that includes the building and the surrounding land*. Both these forms should occur in the UW dictionary. This leads to the problem of the correct UW selection when the word *premises* occurs in a sentence. For these words, for example, the lexicon needs to store the UWs *premise(icl>assumption)* and *premises(icl>place)*. The question of choosing the right sense of *premise* as in *clean the premises* will, however, arise, and this can be resolved only by using the lexical properties of the main verb and the surrounding words.

Hindi, like Japanese, has a special way to show respect. It uses plural forms of pronoun for this purpose. For example, आप [aap](all of you) is used instead of तू [too] (you) for a person when addressed with respect, and हम [ham](we) is used for मैं [main](I) to show one's own importance. English does not have such practices. Thus while translating from English to Hindi, the produced sentence may be unacceptable for a native speaker of Hindi. For example, तू [too] used instead of आप [aap] for *father* or a distinguished person will be frowned upon. We have explained in section 4.3 through the sentence H13 the strategy for dealing with this phenomenon.

#### B. Person

The *person* of a noun does not generally change in translating between Hindi and English. But there is one situation where this occurs, and this happens more with the spoken Hindi than with the written form. Hindi speakers often use the second person plural form instead of

---

the third person singular to describe a person who is being interviewed or is in focus of an event. For example,

H22. आप ने अमरिका से अपनी पी.एच.डी. की उपाधि प्राप्त की।  
aap ne amarikaa se apanee p.h.d. kee upaadhi praapt kee  
you (plural) America-from your (plural) Ph.D.-of degree obtained

E22. He/She obtained his/her Ph.D. degree from America.

It is not easy to deal with case. The fact that आप (aap) translates to *he/she* can be known only from the discourse and currently the UNL handles only single sentences, calls for post editing of the UNL expressions.

### C. Gender

Three gender forms are recognized in English- *masculine*, *feminine* and *neuter*, while in Hindi there are only two forms- masculine and feminine. This does not pose much of a difficulty in translation from Hindi to English and vice versa since the Language-UW dictionaries are different for the two languages. The gender attributes are language dependent. For the UW *child(icl>human)*, the English mapping *child* has the neuter gender, while the Hindi mapping बच्चा [bachcha] (child) has the Masculine Gender.

The other differences with respect to the gender occur with pronouns and the possessive case. Hindi does not have different pronouns for different genders. For example, there are *he* and *she* in the third person in English, but there is only a single pronoun वह [vah](he or she) in Hindi. The verb morphology helps identify the gender. In the Hindi enconversion, by default *he* is generated for वह [vah]. This mapping obviously is kept in the dictionary.

Gender specific possessive pronouns (*his*, *her* or *its*) are used in English, while in Hindi, उस [us](his or her) is used for both the genders. On the other hand, Hindi expresses the gender of the possessed entity by using different case markers. For example, in Hindi, उसका दोस्त [usakaa dost] (his/her he-friend) or उसकी दोस्त [usakee dost] (his/her she-friend) is used to refer to a boy friend or a girl friend respectively. In English the possessive preposition *of* is common for all genders, while in Hindi the corresponding case markers का [kaa](of-male) and की [kee](of-female) are used according to the gender of the possessed entity.

### D. Tense

There are irregular verbs in English, which require separate entries in the dictionary, since the irregular verbs cannot be morphologically derived in a simple way from the stems. In Hindi also, there are irregular transformations of verbs. For example, कर (do) and किया (did). An important distinction in terms of the tense is that English does not show any inflection from the stem for the future tense, but uses auxiliaries like *will* and *shall*. For example,

E23. He will read.                      He will write.

While in Hindi, the present continuous tense does not show any inflection.

H23. वह पढ़ रहा है।                      वह लिख रहा है।  
vah padh rahaa hai                      vah likh rahaa hai  
he reading-is                              he writing-is

Here, पढ़ [padh](read) and लिख [likh](write) are the base morphemes for all possible transformations with respect to tense and person. These phenomena are dealt with through the elaborate set of morphology rules in the analyser.

---

## 6 Language Divergence between Hindi and English

We have already described the major differences between Hindi and English. In this section, we discuss them in a more formal setting proposed in (Dorr 1993) which classifies various Language divergences and suggests solutions to them with respect to the Lexical Conceptual Structure (LCS).

Unlike LCS, UNL is based on the linking of word concepts in a semantic net like representation. We aim to show that most of the divergences described in (Dorr 1993) either do not affect UNL based translations or are comparatively easier to handle than in the LCS approach. Wherever possible, the examples from (Dorr 1993) are used.

### 6.1 Syntactic Divergence

Dorr (Dorr 1993) gives the following divergences arising from structural and syntactic aspects of German, Spanish and English languages:

- Constituent Order divergence
- Adjunction Divergence
- Preposition-Stranding divergence
- Movement divergence
- Null Subject Divergence
- Dative Divergence
- Pleonastic Divergence

In this section, we discuss the effect of each of these on the analysis of English and Hindi into the UNL form and also of generation from UNL into Hindi.

#### 6.1.1 Constituent Order divergence

Constituent Order divergence stands for the word order distinctions between English and Hindi. Essentially, the constituent order describes where the specifier and the complements of a phrase are positioned. For example, in English the complement of a verb is placed after the verb and the specifier of the verb is placed before. Thus English is a Subject-Verb-Object (SVO) language. Hindi, on the other hand, is a Subject-Object-Verb (SOV) language. The following shows the constituent order divergence between English and Hindi:

E24. Jim is playing tennis.

S V O

H24. जीम टेनिस खेल रहा है।

jeem Tennis khel rahaa hai

Jim tennis playing-is

S O V

Also, in Hindi, the qualifier of the complement succeeds the verb whereas in English, it succeeds the complement. For example,

E25. He saw a girl whose eyes were blue.

S V O Q

H25. उस ने एक लड़की को देखा जिसकी आँखें नीली थी।

us ne ek ladakee ko dekhaa jisakee aankhen neelee thee

He a girl-to saw whose eyes blue were

S O V Q

The UNL expressions generated from both English and Hindi are the same for these examples. In general, constituent order divergence does not affect the results of enconversion. But it does affect the *strategy of analysis*. The EnConverter system requires two UWs or Compound UWs *to be adjacent to each other* to generate a UNL expression between them. After every relation is generated, one of the participating UWs is deleted from the node-list and is made the child of the other UW in the semantic tree. For Hindi, the complement and its qualifier cannot be adjacent at any point of the analysis. Hence the SOV structure of the input sentence is converted in the intermediate steps into the SVO structure. The UNL expressions generated for the above example are:

```
[S]
aoj(see(icl>do).@past.@pred.@entry, he(icl>person))
obj(see(icl>do).@past.@pred.@entry, girl(icl>person))
pof(girl(icl>person), eye(icl>thing).@pl)
aoj(blue(icl>state),eye((icl>thing).@pl)
[/S]
```

### 6.1.2 Adjunction Divergence

Syntactic divergences associated with different types of adjunct structures are classified as *Adjunction divergence*. Hindi and English differ in the possible positioning of the adjective phrase. In the former, this phrase can be placed to the left of the head Noun. This is not allowed in English.

For example,

E26. \*the [living in Delhi] boy

H26. [दिल्ली में रहनेवाला] लड़का  
 [dillee mein rahanevaalaa] ladakaa  
 [Delhi-in living] boy

वाला [vaalaa] added to रहना [rahanaa](live) makes it an adjective phrase. This construction, in general, applies to only habitual actions. For example,

H27. (A) राम ने [मोहन को पसंद आनेवाला] तोहफा भेजा ।  
 raam ne [mohan ko pasand aanevaalaa] tohafaa bhejaa  
 Ram [Mohan-to like come-ing] gift sent  
 (B) राम ने वह तोहफा भेजा जो मोहन को पसंद आया ।  
 raam ne vah tohafaa bhejaa jo mohan ko pasand aayaa  
 Ram that gift sent that Mohan-to like came  
 (C) राम ने वह तोहफा भेजा जो मोहन को पसंद है ।  
 raam ne vah tohafaa bhejaa jo mohan ko pasand hai  
 Ram that gift sent that Mohan-to like is

E27. Ram sent the gift that mohan likes.

Sentences H27 (A) and (C) are equivalent. H26(B) cannot use वाला. The UNL expressions of the sentence H27 follow.

---

[S]  
 agt(send(icl>do).@entry.@past,Ram(icl>person))  
 obj(send(icl>do).@entry.@past,gift(icl>object).@indef)  
 aoj(like(icl>do).@present,Mohan(icl>person))  
 obj(like(icl>do).@present,gift(icl>object).@indef)  
 [/S]

The generator identifies an adjective clause by the two arrows coming into the noun node *gift(icl>object)* from the verb nodes *send(icl>do)* and *like(icl>do)*. It identifies the main verb of the sentence by the *@entry* attribute. It generates the sentence H27(A) if the verb *like(icl>do)* is in the present tense and the sentence H27(B) if the verb is in the past tense.

Another Divergence in this category is the *prepositional phrase (PP) adjunction* with respect to a verb phrase. In Hindi a PP can be placed between a verb and its object or before the object, while in English it can only be at the maximal level (*i.e.*, not between the verb and its object). For example,

E28. He called me [to his house.]

\*He called [to his house] me.

H28. (A) उसने मुझे [अपने घर] बुलाया ।

usne mujhe [apne ghar] bulaayaa

he to-me his house called

(B) उसने [अपने घर] मुझे बुलाया ।

usne [apne ghar] mujhe bulaayaa

he his house to-me called

The UNL expressions for both the sentences remain the same and the generator can produce any of the above Hindi sentences.

[S]  
 agt(call(icl>do).@past.@pred.@entry, he(icl>person))  
 obj(call(icl>do).@past.@pred.@entry, I(icl>person))  
 plt(call(icl>do).@past.@pred.@entry, house(icl>place))  
 [/S]

### 6.1.3 Preposition-Stranding Divergence

This divergence is accounted for by the choice of proper governors.

E29. Which shop did John go to?

H29. \*किस दुकान ज्होन गया में ?

kis dukaan john gayaa mein

[S]  
 agt(go(icl>do).@past,@pred.@entry, John(icl>person))  
 plt(go(icl>do).@past,@pred.@entry, shop(icl>place))  
 mod(shop(icl>place), which)  
 [/S]

H29 which is a literal translation of E29, is syntactically incorrect as the case marker में [mein](to) cannot be a proper governor for the noun phrase. In English, the preposition *to* is a proper governor for the trace. The case marker में [mein](to) is required to follow the noun which in this case is दुकान [dukaan](shop). The Hindi Generator does the syntax planning accordingly and produces the right case marker when it encounters *plt(go(icl>do), shop(icl>place))*.

---

### 6.1.4 Null Subject Divergence

In Hindi, unlike in English, the subject of the sentence can be left implicit. For example,

E30. Long ago, there was a king.

H30. बहुत पहले एक राजा था।

bahut pahale ek raajaa thaa

long ago one king was

Hindi allows dropping of the subject where the subject is obvious. For example, we repeat sentence H7:

H31. जा रहा हूँ।

jaa rahaa hun

going-am

E31. \*am going.

The subject मैं [main](I) is absent. Such omissions are permitted only in two situations. The first is that *a pleonastic is eliminated* and the second is when *a valid subject is omitted* as its implicit presence is reflected through the morphology of the predicate. The first case is discussed in the next sub-section. In the other case, the eliminated subject requires to be produced in the UNL expressions. This is done by examining the structure of the UNL graph during the analysis. *aoj* and *agt* are the only relations that relate the predicate with the subject of the sentence. The system takes care of this phenomenon by detecting the absence of *agt* or *aoj* relation with the main predicate in a non-passive sentence. If such a condition is detected then it inserts an appropriate UW, *I(icl>person)* in the above example, in the nodelist. The analysis of the sentence is then continued as usual. The UNL representation for H31 is:

```
[S]
agt(go(icl>do).@entry.@present.@progress, I(icl>person))
[/S]
```

A special kind of Null-Subject divergence is the Pleonastic Divergence.

### 6.1.5 Pleonastic Divergence

A Pleonastic is a syntactic constituent that has no semantic content. For example,

E32. It is raining.

*It* has no semantic role in the above sentence. Similarly in sentence E30, *there* does not have any semantic role. Frequently, pleonastics are linked to another constituent that carries the appropriate semantic content. If the UNL representation of the above sentence is done as follows then the Hindi generator will probably generate the sentence H32, which is stylistically incorrect.

```
[S]
aoj(rain(icl>do).@progress.@entry, it(icl>abstract thing))
[/S]
```

H32. ? यह बारीश हो रही है।<sup>7</sup>

yah baareesh ho rahee hai

---

<sup>7</sup> ? indicates that the sentence may be syntactically correct but its stylistic validity is questionable.

---

this rain happening-is

To deal with such problems, pleonastics are identified using semantic properties of the words in the sentence and they do not become part of the UNL expressions. For example, it has been observed that natural events like *rain*, *thunder*, *snow*, *etc.* make sentences using *it* as a pleonastic. Such words are given an attribute called *NATURAL-EVENT* in the lexicon, using which, the *it* in the sentence, as in E32, is eliminated from the UNL expressions. Now, the UNL representation of E32 is:

```
[w]
rain(icl>do).@entry.@pred.@progress
[/w]
```

Note that the UW *rain(icl>do)* is not related to any other word and the event is described by a single UW which means *rain is in progress*. This can be translated to a correct form of H32 as:

बारिश हो रही है।

Detailed information about detecting pleonastics can be found in (Parikh 2001).

## 6.2 Lexical Semantic Divergence

*Lexical-semantic divergence* (Dorr 1993)- arising from the properties of the entries in the lexicon- are of the following types:

- Conflational divergence
- Structural divergence
- Categorical divergence
- Head swapping divergence
- Lexical divergence

These are explained with examples along with their effect on the analyser and generator outputs.

### 6.2.1 Conflational Divergence

*Conflation* is the lexical incorporation of necessary components of meaning (or arguments) of a given action. This divergence arises from a variation in the selection of the word between the source language and the target language. For example,

E33. Jim stabbed John.

H33. जीम ने ज्हाँन को छूरे से मारा।

jeem ne john ko chhoore-se maaraa

Jim John-to knife-with hit

*stab* does not have a single word equivalent word in Hindi. We require the phrase छूरे से मारा [*chhoore se maaraa*] (*hit with a knife*). As a result, the UNL expressions generated from E32 and H32 vary. The HA produces:

```
[S]
agt(hit(icl>do).@entry, Jim(icl>person))
ben(hit(icl>do).@entry, John(icl>person))
ins(hit(icl>do).@entry, knife(icl>thing))
[/S]
```

---

However, the EA directly produces *stab(ict>do)*. But if the Hindi phrase छूरे से मारा [chhoore se maaraa] (*knife-with hit* - hit with a knife) is mapped to the UW *stab(ict>do)* in the Hindi-UW dictionary, the HA produces:

[S]  
agt(stab(ict>do).@entry, Jim(ict>person))  
ben(stab(ict>do).@entry, John(ict>person))  
[/S]

The EnConverter's property of picking up the longest lexeme has been exploited here. The expression is the same as the UNL expressions produced by the EA. Most cases of conflation divergence are handled this way. The opposite case of Hindi words being conflational has been discussed in section 2.2 for both noun (*devar*) and verb (*ausaanaa*).

### 6.2.2 Structural Divergence

E34. Jim entered the house.

H34. जीम घर में प्रवेश किया ।

jeem ghar mein pravesha kiyaa

Jim house-into entry did

The Hindi sentence diverges structurally from the English sentence, since the verbal object is realized as a noun phrase (*house*) in English and as a prepositional phrase (घर में [ghar mein]) in Hindi. In English, both *enter* and *enter into* will be allowed whereas in Hindi the prepositional phrase should strictly be used. The UNL expressions from both the English and Hindi sentences are the same:

[S]  
agt(enter(ict>do).@entry.@pred.@past, Jim(ict>person))  
plt(enter(ict>do).@entry.@pred.@past, house(ict>place))  
[/S]

If *into* is not present, the EA can generate *obj* between *enter* and *house*. This problem is solved by using the semantic attribute *PLACE* of the word *house* in the lexicon. This causes the generation of *plt* instead of *obj*. Thus, the lack of syntactic information (implicit prepositions) is compensated for by the semantic knowledge.

### 6.2.3 Categorical Divergence

*Categorical Divergence* arises if the lexical category of a word changes during the translation process. For example,

E35. They are competing.

H35. वह मुकाबला कर रहे हैं ।

vaha muqaabala kar rahe hai

They competition doing-are

Here, *competing* is expressed as a verb in English and as a noun-verb combination (*muqaabala kar*- competition do) in Hindi. This divergence is very common in English to Hindi MT- and in general in English to an Indian language MT. Hindi- like most Indian languages- forms *combination verbs* in which a noun is followed by a form of कर [kar] (*do*) or हो [ho] (*be*) to express the action suggested by the noun.



This phenomenon is handled by the HA by having two entries of such nouns in the lexicon- one as a noun and the other as a verb. The verb entry has an attribute *link* which indicates that a form of कर [kar](do) is to follow the noun. For the example in point, मुकाबला [muqaablaa] has the following two entries in the lexicon:

[मुकाबला] {} “competition(icl>action)” (N, NA, MALE, INANI, ABSTRACT);  
 [मुकाबला] {} “compete(icl>do)” (V, link);

Because of this, the UNL expressions for both the English and the Hindi sentences are the same:

[S]  
 agt(compet(icl>do).@entry.@pred.@present.@progress, they(icl>person))  
 [/S]

## 6.2.4 Head swapping divergence

### A. Demotional Divergence

Demotional divergence is characterized by the *demotion* (placement into a *lower down* position) of a logical head. In such a situation, the logical head is associated with the syntactic adjunct position and then the logical argument is associated with a syntactic head position.

For example,

E36. It suffices.

H36. यह काफी है ।

yaha kaafee hai

It sufficient-is

The word *suffice* is realized as the main verb in English but as an adjectival modifier काफी है [kaafee hai] in Hindi. The UNL expressions generated from the EA and the HA differ. The EA generates:

[S]  
 aoj(suffice(icl>do).@entry.@present, it)  
 [/S]

While, HA will generate the following UNL expressions:

[S]  
 aoj(sufficient.@entry.@present, it)  
 [/S]

The HG produces the sentence H36 from both the representations. This is because the Hindi-UW dictionary has *suffice(icl>do)* mapped to काफी है [kaafee hai](is sufficient). Hindi does not have any equivalent verb for *suffice*. Thus the divergence is handled in the lexicon with the following entry:

[काफी] {} “suffice(icl>do)” (V, VI);

### B. Promotional Divergence

Promotional divergence is characterized by the *promotion* (placement into a higher position) of a logical modifier. The logical modifier is associated with the syntactic head position and then the logical head is associated with an internal argument position.

E37. The play is on.

H37. खेल चल रहा है ।

---

khel chal rahaa hai

Play going-on-is

Here the modifier *is on* is realized as an adverbial phrase in English but as the main verb चल रहा है [chal rahaa hai](going-on-is) in Hindi. The UNL expressions generated by the EA are:

```
[S]
aoj(on(icl>state).@entry.@present, play(icl>abstract thing).@def)
[/S]
```

The HA on the other hand generates:

```
[S]
agtj(go on(icl>occur).@entry.@present.@progress, play(icl>abstract thing))
[/S]
```

The solution to this is same as that for demotional divergence. The dictionary entry in this case would be:

[चल] {} “go on” (V,Va);

### 6.2.5 Lexical Divergence

*Lexical* divergence means that the choice of a target language word is not a literal translation of the source language word. However, lexical divergence arises only in the context of other divergence types. In particular, lexical divergence generally co-occurs with conflation, structural and categorial divergences.

H38. ज्होन जबरजस्ती घर में घुस गया।  
john jabarjasti ghar mein ghus gayaa  
John forcefully house-in enter-go

E38. John broke into the house.

Here the divergence is lexical in the sense that the target language word is not a literal translation of the source-language word.

The EA and HA will both produce the following UNL expressions:

```
[S]
agt(enter(icl>do).@past.@entry.@force, John(icl>person))
plc(enter(icl>do).@past.@entry.@force, house(icl>home))
[/S]
```

It is clear how the HA can produce the above expressions. EA achieves this by mapping *break into* to *enter(icl>do)* in the English-UW dictionary. It also places an attribute *FORCED* into the lexicon which signals the generation of *@force* during analysis.

## 7 Experimental Observations

The English Analyser (EA), the Hindi Analyser (HA) and the Hindi Generator (HG) have been tested using the sentences in the United Nations Charter provided by the United Nations University. The corpus was designed to test the DeConverters of different languages all over the world. The corpus has around 180 sentences. It is in English and has been manually translated into Hindi for the HA. As the analysers are not yet equipped with Word Sense Disambiguation capability, inter-category word senses were manually disambiguated. As mentioned before, the analysers have intra-category or part of speech disambiguation capability. Approximately 80% of these sentences have been successfully converted to UNL expressions by the analysers without any change in the input sentences. The rest had to be pre

---

edited to a certain extent by simplifying the structure of the sentences and controlling the use of punctuations. The UNL expressions generated by the English and Hindi Analysers were given to the Hindi Generator. 95% of these UNL expressions were correctly converted into Hindi by the HG.

The Hindi analyser has also been tested on a huge Hindi corpus provided by the Ministry of Information Technology, Government of India. This corpus consisted mainly of stories from the political domain. The English Analyser too has been tested on documents like the *EnConverter Manual*, sentences from Brown Corpus and stock market stories downloaded from different web sites. We are continuously upgrading our system by testing on numerous corpora. The test base is currently considerable. The *Barcelona* corpus obtained from the multilingual information processing being conducted in Spain, sentences from the *Medline* corpus, *the agricultural corpora* from the Gujarat Government and such other corpora are being worked on. Thus the evaluation process is in progress.

Besides techno-scientific domains we have tested the analyser on literary works also. It is worth noting here that such sentences require more pre-processing than sentences from the technical domains. An example of a sentence not handled properly by the system from Wodehouse is:

*I loosed it down the hatch, and after undergoing the passing discomfort, unavoidable when you drink Jeeves's patent morning revivers, of having the top of the skull fly up to the ceiling and the eyes shoot out of their sockets and rebound from the opposite wall like racquet balls, felt better.*

However, with some obvious pre-editing as shown below the sentence is analysed accurately.

*I loosed it down the hatch and after undergoing the passing discomfort which is unavoidable when you drink Jeeves's patent morning revivers, of having that the top of the skull fly up to the ceiling and the eyes shoot out of their sockets and rebound from the opposite wall like racquet balls, felt better*

The verification of the analysis and generation processes have been carried out by converting Hindi sentences into UNL expressions and generating the sentence back. The results obtained are quite satisfactory in the sense that the generated sentences are in most cases the same as the source sentences. Sometimes the postposition markers are different while at other places a different word has been chosen. Yet other times, the structure of the generated sentence differs from the source sentence. However, in all cases the idea contained in the source sentence is conveyed in the generated sentence. A few examples are given below:

### **Example 1**

#### **Source Sentence**

H39. अध्ययन समूह उपकरण और सेवाओं से सम्बन्धित बहुत सारे मुद्दों को समाविष्ट करते हैं।

adhyayana samooha upakaran aur sevaaoM se saMbandhit bahoot saare muddoM ko samaavisht karate haiM.

E39. The Study Groups cover a wide number of issues related to equipments and services.

#### **UNL**

[S]

aobj(cover(icl>include):21.@entry.@present.@pred, Study Groups:00)  
obj(cover(icl>include):21.@entry.@present.@pred, issue(icl>important point):1R.@pl)  
mod(issue(icl>important point):1R.@pl, relate(icl>concerning):14)  
mod(issue(icl>important point):1R.@pl, wide number of(icl>very great):1F.@pl)

---

aoj(related>concerning):14, :01)  
and:01(service(related>assistance):0U.@entry.@pl, equipment(related>tool):0G)  
[S]

### Generated Sentence

H39'. अध्ययन समूह उपकरण और सेवाएँ सम्बन्धित बहुत सारे मुद्दों को समाविष्ट करते हैं।

adhyayana samooha upakaran aur sevaeM saMbandhit bahoot saare muddoM ko samaavisht karate haiM.

### Remark

Comparing the generated sentence with the source sentence we find that only the postposition marker of *sevaa* (*service*) has changed. The sentence is acceptable in Hindi and the meaning of course is conveyed.

### Example 2

#### Source Sentence

H40. अन्तर्राष्ट्रीय संस्था के रूप में आई टी यू सरकारों और गैर सरकारी संस्थाओं को दूरसंचार तन्त्र और सेवाओं के परिचालन के विस्तार और समन्वयीकरण हेतु कार्य करने के लिये और सभी देशों तक उनकी पहुंच को बढ़ावा देने के लिये एक साथ लाता है।

antarraashtriya saMsthaa ke roop meM aaii tii yoo sarakaaroM aur gair-sarakaarii saMsthaoM ko doorasaMchaar taMtra aur sevaeoM ke paricaalan ke vistaar aur samanvayiiikaraN hetu kaarya karane ke lie aur sabhii deshoM tak unakii pahuMch ko baDAvA dene ke lie eka saath laataa hai.

E40. As an international organization, ITU brings together governments and private sectors to work for expanding and coordinating the operation of the telecommunication networks and services, and to promote their access to all countries.

### UNL

[S]  
aoj(bring together(related>gather):6T.@entry.@present.@pred,  
ITU(related>International Telecommunication Union):0X)  
obj(bring together(related>gather):6T.@entry.@present.@pred, :01)  
pur(bring together(related>gather):6T.@entry.@present.@pred, :04)  
and:04(foster(related>nurture):69.@entry.@pred, work(related>do work):4J.@pred)  
obj:04(foster(related>nurture):69.@entry.@pred, access(related>approach):5X)  
scn:04(access(related>approach):5X, country(related> nation):5G.@pl)  
mod:04(access(related>approach):5X, those(related>pronoun):5R)  
aoj:04(overall(related>all):5B, country(related>nation):5G.@pl)  
pur:04(work(related>do work):4J.@pred, :03)  
mod:03(coordination(related>coordinating):3Y.@entry, operation(related>functioning):38)  
and:03(coordination(related>coordinating):3Y.@entry, expanding(related>expansion):3M)  
mod:03(operation(related>functioning):38, :02)  
mod:03(:02, telecommunication:2B)  
and:02(service(related>assistance):2Y.@entry.@pl, network(related>system):2N)  
and:01(institution(related>organization):1Z.@entry.@pl, government:16.@pl)  
aoj:01(private(ant>governmental):1K, institution(related>organization):1Z.@entry.@pl)  
aoj(ITU(related>International Telecommunication Union):0X, as:0L)  
obj(as:0L, institution(related>organization):0E)  
aoj(international(related>characteristic):00, institution(related>organization):0E)  
[S]

---

## Generated Sentence

H40'. सरकारों और गैर सरकारी संस्थाओं को अन्तर्राष्ट्रीय संस्था के रूप में आई टी यू दूरसंचार की तन्त्र और सेवाएं परिचालन के विस्तार और समन्वयीकरण के लिये कार्य करने और सभी देशों में उनकी पहुंच को बढ़ावा देने के लिये एक साथ लाता है।

sarakaaroM aur gair-sarakaarii saMsthaom ko antarraashtriiya saMsthaa ke roop meM aaii tii yoo doorasaMchaar kii taMtra aur sevaaeM paricaalan ke vistaar aur samanvayiiikaraN ke lie kaarya karane aur sabhii deshoM meM unakii pahuMch ko baDAvA dene ke lie eka saath laataa hai.

## Remarks

Here the phrase सरकारों और गैर सरकारी संस्थाओं को (governments and private sectors) has been placed at the start of the sentence. Being followed by के रूप में (*as*) this gives an impression initially that *ITU* is being qualified by the phrase. This, however, gets rectified as one reads ahead. The meaning is conveyed, but the source sentence is structurally better than the generated one. There are other minor changes like सेवाओं becoming सेवाएं and तक becoming में, which do not alter the meaning.

## Example 3

### Source Sentence

H41. यह उत्सव प्रदर्शनों का एक बड़ा कार्यक्रम और सांस्कृतिक क्रिया कलापों का एक विस्तृत क्षेत्र प्रदान करेगा जो पूरे 155 दिनों तक विश्व संस्कृतियों की सृजनात्मकता पर ध्यान केन्द्रित करेगा।

yah utsav pradارشanoM kaa ek badaa kaaryakram aur saaMskritik kriyaa-kalaapoM kaa eka vistrit kshetra pradaan karega jo poore 155 dinoM tak vishva saMskrtiyom kii srjanaatmakataa par dhyaan kendriwt karega.

E41. This Festival will offer a broad program of performances and a wide range of cultural activities that will focus on the creativity of world cultures over a period of 155 days.

## UNL

[S]

obj(provide(icl>do):2Q.@entry.@future.@pred, :01)  
aoj(provide(icl>do):2Q.@entry.@future.@pred, festival(icl>event):05)  
mod(festival(icl>event):05, this:00)  
aoj(focus(icl>concentrate):4W.@future.@pred, :01)  
and:01(range(icl>variety):2B.@entry, program(icl>performance):10)  
mod:01(range(icl>variety):2B.@entry, activity(icl>action):1Q.@pl)  
aoj:01(cultural(aoj>thing):1F, activity(icl>action):1Q.@pl)  
mod:01(program(icl>performance):10, performance(icl>abstract thing):0E.@pl)  
aoj:01(great(icl>characteristic):0U, program(icl>performance):10)  
tim(focus(icl>concentrate):4W.@future.@pred, day(icl>period):3H.@pl)  
scn(focus(icl>concentrate):4W.@future.@pred, creativity(icl>creativity):4D)  
mod(creativity(icl>creativity):4D, culture(icl>civilisation):3Y.@pl)  
aoj(world(mod<thing):3S, culture(icl>civilisation):3Y.@pl)  
aoj(around(icl>about):38, day(icl>period):3H.@pl)  
qua(day(icl>period):3H.@pl, 155:3D)

[/S]

## Generated Sentence

H41'. यह उत्सव सांस्कृतिक क्रिया कलापों के बड़ा प्रदर्शनों का एक कार्यक्रम और एक विभिन्न प्रकार प्रदान करेगा जो पूरे 155 दिनों में विश्व संस्कृतियों की रचनात्मकता पर ध्यान केन्द्रित करेगा।

---

yaha utsav saasMkritik kriyaa-kalaapoM ke badaa pradarshanoM kaa eka kaaryakram aur eka vibhinn prakaar pradaan karegaa jo poore 155 dinoM meM vishva saMskrtiyoM kii racanAwmakawaa par dhyaan kendriwt karegaa.

### Remarks

This illustrates changes of word as in (i) विभिन्न प्रकार (*range*) in place of विस्तृत क्षेत्र (another meaning of *range*) (ii) रचनात्मकता in place of सृजनात्मकता both meaning the same, *i.e.*, *creativity* and (iii) में (*in*) in place of तक (*over*). The reordering of phrases, however, is more serious as in सांस्कृतिक क्रिया कलापों के बड़ा प्रदर्शनों का एक कार्यक्रम (a program of a broad performance of cultural activities) replacing प्रदर्शनों का एक बड़ा कार्यक्रम और सांस्कृतिक क्रिया कलापों (a broad program of performances and cultural activities) where meaning alteration within that part of the sentence has taken place. The generated sentence, however, is not far in meaning from the source sentence.

The following example shows that though sentences in English and Hindi with identical meaning are represented as different sets of UNL expressions by the EA and the HA respectively, the HG generates the same output for both the representations.

The sentence is:

**UNEP has a mission to care for the environment.**

EA generated the following UNL expressions:

[S]  
aobj(have(icl>state):05.@entry.@present, UNEP(icl>United Nations Environment Programme):00)  
obj(have(icl>state):05.@entry.@present, mission(icl>duty):0B.@indef)  
pur(care(icl>do):0M.@present.@pred, environment(icl>state):0Z.@def)  
pur(mission(icl>duty):0B.@indef, care(icl>do):0M.@present.@pred)  
[/S]

The same sentence was manually translated to Hindi and input to the HA.

H42. यु एन ई पी का लक्ष्य पर्यावरण की देखभाल करना है।

U N E P kaa lakshya paryaavaran kee dekhabhaal karnaa hai

UNEP-of mission environment-of care-do is

E42. UNEP has a mission to care for the environment.

The output of the HA was:

[S]  
obj(care(icl>do):1I.@entry.@present.@pred, environment(icl>abstract thing):13)  
mod(mission(icl>duty):0W, UNEP(icl>United Nations Environment Programme):0H)  
aobj(care(icl>do):1I.@entry.@present.@pred, mission(icl>duty):0W)  
[/S]

The output of the HG for both the sets is:

H43. यु एन ई पी का लक्ष्य पर्यावरण का ख्याल रखना है।

U N E P kaa lakshya paryaavaran kaa khyaal rakhnaa hai

UNEP-of mission environment-of care-do is

E43. UNEP has a mission to care for the environment.

This lends credence to the capturing of the semantics by the UNL is a language independent way.

---

At this stage, it is difficult to compare the computational complexity of the analysis of Hindi and English sentences into UNL. However, we mention a few pointers in that direction:

1. UNL is based on a predicate centric framework. The analyser needs to know the predicate before it starts generating the UNL expressions. Because of the SOV structure of Hindi, in most case, the verb occurs at the end of the sentence. Thus the Hindi Analyser has to do a complete morphological analysis of the words on its way to the end of the sentence. There are examples in which the Hindi analyser completes the morphological analysis of words till the end of the sentence and then comes all the way back to the subject of the sentence. This normally does not happen in the case of the English Analyser. As soon as it encounters the predicate, it can start dealing with the complements and the prepositional phrases (PP).

The SOV structure also causes problems because of the computational model adopted. For example, the adjacency requirement of the logical units or constituents described in section 6.1.1, sometimes calls for manipulations like the exchange of syntactic constituents to change their order in the sentence.

2. Prepositions in English can be proper governors (Dorr 1993). Thus sentences like the following need to be dealt with:

**Which shop did John go to?**

The system is required to produce:

**plt(go(icl>do).@entry.@interrogation.@past, shop(icl>place))**

But because of the computational model adopted *to* is required to be adjacent to *shop*. This is achieved by exchanging *go* and *shop* when they are adjacent to each other in the node-list. Such computations can become very complex in the case of longer sentences with long distance dependencies. In Hindi the case markers cling to the noun they govern leading to simpler computation.

3. The problem of word sense disambiguation poses difficulties for both the analysers. UNL requires the analysers to generate an unambiguous word concept. Neither the EA nor the HA has any support for *sense disambiguation*. However, both perform very well for *part of speech disambiguation*. This helps prune options for a Universal Word.
4. Our experiments show that the number of rules fired is nearly the same for both English and Hindi analysis of most cases. This number is directly proportional to the number of lexemes. At least two rules- *shift* and *process*- are required for each morpheme. Hindi generally requires more morphological analysis. Thus the number of rules fired is **a bit** more than that of English. To illustrate this, the statistics for four sentences is given in Table 3. The sentences are:

E44. UNIFEM works to promote the economic and political empowerment of women.

H44. युनीफेम औरतों के आर्थिक तथा राजनैतिक अधिकार को बढ़ावा देने के लिए कार्य करती है।

yunifem ouraton ke aarthik tathaa raajanaitik adhikaar ko badhaavaa dene ke liye kaarya karatee hai.

UNIFEM women-of economic and political empowerment-to promote-give-for work-doing-is

E45. I know the lady who has worn a blue saree.

H45. मैं उस औरत को जानता हूँ जिsने नीली साड़ी पहनी है।

---

mai us ourat ko jaanataa hun jisane neelee saadee pahanee hai.

I that woman-to know-is who blue saree worn-has

E46. Uncle told us that Gita is removing dust from the kitchen with a broom.

H46. चाचा ने हम से कहा कि गीता रसोईघर में झाड़ू से धूल निकाल रही है।

chaachaa ne ham se kahaa ki geetaa rasoighar mein jhaadoo se dhool nikaal rahee hai

uncle us-to told that Gita kitchen-in broom-with dust removing-is

E47. With Lord Krips, his wife had also come and she wanted to buy a fine shawl from India for taking home.

H47. लर्ड क्रिप्स के साथ उनकी पत्नी भी आई हुई थी और वे भारत से स्वदेश ले जाने के लिए एक उमदा शाल खरीदना चाहती थी।

lord krips ke saath unakee patnee bhee aae huee thee or ve bhaarat se svadesh le jaane ke liye ek umdaa shaal khareedanaa chaahatee thee.

Lord Krips-with his wife also come-had and she India-from native land take-go-for one fine shawl buy-to wanted

Sen. No.	Type	No. of Lexemes		No. of Rules Fired	
		English	Hindi	EA	HA
40	Simple	22	30	54	64
41	Adjective Clause	20	20	46	55
42	Noun Clause	26	33	57	71
43	Compound	44	55	101	122

**Table 3: Statistical information for example sentences**

The difference in the number of rules fired can be accounted for from the fact of two rules used per lexeme. The other contributing factors are:

- Simple (E40, H40): The presence of the conjunction in the sentence. English requires looking ahead by several words to make sure it is not a compound sentence and is a simple conjunction of nouns. The morphology of Hindi helps in avoiding this processing.
- Adjective clause (E41, H41): The adjective clause requires the Hindi analyser to do extra processing as explained in section 6.1.1. This explains the 9 extra rules fired by HA.
- Noun clause (E42, H42): The difference here is exactly proportional to the difference in the number of morphemes.
- Compound (E43, H43): An extra rule fires in the case of the EA. This is for the look ahead processing of the compound sentence.



---

## 8 The Issue of Disambiguation

As has been mentioned at various places in the paper, our system currently does mainly *part of speech disambiguation* and a little bit of *sense disambiguation for postposition marker and wh-pronouns*. **The main instruments of disambiguation are the condition windows around the analysis heads and also the lexical attributes of the words.** This achieves the look ahead and look back necessary for disambiguation. We point out the specific sentences mentioned in the paper where disambiguation takes place. The words disambiguated are in bold.

1. *The soldier went away to the totally **deserted desert** to **desert** the house in the **desert**.*  
Part of speech disambiguation using the adverb *totally* which must precede an adjective which in turn must precede a noun.
2. *He went to my home **when** I was away.*  
POS disambiguation (adverb phrase) using the fact that *home* does not have *time* attribute.
3. *He met me at a time **when** I was very busy.*  
POS disambiguation (adjective phrase)- *when* can qualify a noun with *time* attribute.
4. The sentences with *se* in Hindi (sentences E3 to E6)  
Sense disambiguation using the lexical attributes of the preceding nouns.
5. The sentences 1, 2, 3 and 4 using **with** in section 5.  
Sense disambiguation using the attributes of the nouns in the sentences.

These examples throw light on the disambiguation capability of the analysers. However, more powerful lexical resources will have to be used for large scale WSD.

## 9 Conclusions and Future Directions

The criteria for deciding the effectiveness of an interlingua are that (a) the meaning conveyed by the source text should be apparent from the interlingual representation and (b) a generator should be able to produce a target language sentence which a native speaker of that language accepts as natural. A careful observer will notice that (a) and (b) are essentially the same. Still we put them down separately to emphasize the presence of a mechanical procedure in (b).

Keeping these criteria in view, our conclusions on the capability of the UNL *vis-à-vis* language divergence especially between English and Hindi are:

1. The UNL expressions generated from English and Hindi texts are mostly the same, as has been brought out in section 6.
2. When they differ, they do so mainly in the case of very overloaded constructs like *have* where the mechanical analyser does not capture the varied nuances.
3. The lexical-semantic divergence is actually handled in the L-UW dictionary. The generator primarily bears the burden of naturalness and idiomaticity in this case.
4. The syntactic divergence, on the other hand, is primarily tackled by the analysers. The capability is built into the rules.
5. The amenability to generation is being tested through at least another language, which is Marathi, a western Indian language, in our case. The results are approximately the same as in Hindi because of the similarity in structure between Hindi and Marathi.

There are several future directions. The L-UW dictionary has to be enriched enormously both in terms of the UW content and the semantic attributes so as to capture the word and world knowledge. The analysers need to be augmented with powerful word sense disambiguation modules. Hindi Generator needs to be thoroughly tested using the UNL expressions produced by the analysers for other languages. Investigation of the UNL as a knowledge representation scheme and the use of this knowledge for various purposes like text summarisation, automatic

---

hypertext linking, document classification, text-image consistency checking and such other knowledge intensive tasks should be carried out.

## 10 References

- Akshar Bharati, Vineet Chaitanya, Rajeev Sangal: 1996, *Natural Language Processing: A Paninian Perspective*, PHI.
- Arnold, D. and des Tombes, L.: 1987, 'Basic theory and methodology in EUROTRA,' in Nirenburg, S., editor, *Machine Translation: Theoretical and Methodological Issues*, Cambridge University Press, pp. 114-135.
- Arnold, Doug and Louisa Sadler: 1990, 'Theoretical Basis of MiMo,' *Machine Translation*, 5:3, pp. 195-222.
- Boitet, Christian: 1988, 'Pros and Cons of the Pivot and Transfer Approaches in Multilingual Machine Translation,' in Klaus Schubert and Toon Witkam (eds.), *Recent Developments in Machine Translation, Dan Maxwell*, Foris, Dordrecht.
- Carbonell, Jaime G. and Masaru Tomita: 1987, 'Knowledge Based Machine Translation, the CMU Approach,' In Sergei Nirenburg (ed.), *Machine Translation: The Theoretical and Methodological Issues*, Cambridge University Press, Cambridge, England, pp. 68-89.
- DeConverter Specification Version 2.0, UNU/IAS/UNL Centre, Tokyo 150-8304, Japan, March 2000.
- Dorr, Bonnie J.: 1992, 'The use of lexical semantics in interlingual machine translation,' *Machine Translation*, 7(3) pp. 135--93
- Dorr, Bonnie J.: 1993, *Machine Translation: a view from the lexicon*, The MIT Press.
- EnConverter Specification Version 2.1, UNU/IAS/UNL Centre, Tokyo 150-8304, Japan, March 2000.
- Guha, R. V., D. B. Lenat, K. Pittman, D. Pratt, and M. Shepherd: 1990, 'Cyc: A Midterm Report,' *Communications of the ACM* 33, no. 8.
- Hardt, S. L.: 1987, 'Primitives' in S. C. Shapiro (ed.), *Encyclopedia of Artificial Intelligence*, John Wiley and Sons, New York, NY, pp. 196.
- Katz, J.: 1966, *The Philosophy of Language*, Harper and Row, New York, pp-156.
- King, M. and Perschke, S.: 1987, *Machine Translation Today: The State of the Art*. Edinburgh University Press. EUROTRA.
- Senniappan, Kumavel and Bhattacharyya, Pushpak: 2000, 'Automatic Generation of Hyperlinks using Semantic Information,' in International Conference on Information Technology (CIT 2000), Bhubaneshwar, India.
- Landsbergen, J.: 1987, 'Isomorphic grammars and their use in the ROSETTA translation system.' In *Machine Translation Today: The State of the Art*. Edinburgh University Press, Edinburgh.
- Lytinen, Steven and Roger C. Schank: 1982, 'Representation and Translation,' Technical Report, Department of Computer Science, Yale University, New Haven, CT, 234.
- Martin, John C.: 1991, *Introduction to languages and the theory of computation*, McGraw Hill.
- Monju, M., Shilpa, T., Smita, D., Leena, G., Shachi, D., and Pushpak Bhattacharyya: 2000, 'Knowledge Extraction from Hindi Text', In Sasikumar M., Rao, Durgesh, Ravi Prakash, P. (eds), Proceedings of the International Conference on Knowledge Based Computer Systems (KBCS 2000), Mumbai, India.
- Muraki, K.: 1987, 'PIVOT: A Two-Phase Machine Translation System,' *Machine Translation Summit- Manuscripts and Program*, Japan, pp. 81-83.
- Muraki, K.: 1989, PIVOT: 'Two-phase machine translation system.' In Proceedings of the Second Machine Translation Summit, Tokyo. Omsa Ltd.

- 
- Nirenburg, Sergei, Victor Raskin and Allen B. Tucker :1987, 'The Structure of Interlingua in TRANSLATOR,' in Sergei Nirenburg (ed.), *Machine Translation: The Theoretical and Methodological Issues*, Cambridge University Press, Cambridge, England, pp. 90-113.
  - Nirenburg, Sergei, Jaime Carbonell, Masaru Tomita and Kenneth Goodman: 1992, *Machine Translation: A Knowledge Based Approach*, Morgan Kaufmann, San Mateo, CA.
  - Okumura, A., Muraki, K., and Akamine, S.: 1991, 'Multi-lingual sentence generation from the PIVOT interlingua.' In Proceedings of the Third Machine Translation Summit, Carnegie Mellon University.
  - Parikh, Jignashu and Bhattacharyya, Pushpak: 2001, 'Towards Realising *The Semantic Web*' sent for publication to *Journal of Natural Language Engineering*, Cambridge University Press.
  - Parikh, Jignashu, Soni, Trilok and Shah, Chirag: 2000, 'Conversion of English Language Texts to Universal Networking Language,' B.E. Dissertation, Dharamsinh Desai Institute of Technology, Nadiad.
  - Perschke, S.: 1989, 'EUROTRA project'. In Proceedings of the Second Machine Translation Summit, Tokyo. Omsa Ltd.
  - Rao, Durgesh, Mohanraj, Kavita, Hegde, Jayprakash, Mehta, Vivek and Mahadane, Parag: 2000, 'A Practical Framework for Syntactic Transfer of Compound-Complex Sentences for English-Hindi Machine Translation,' In Sasikumar M., Rao, Durgesh, Ravi Prakash, P. (eds), Proceedings of the International Conference on Knowledge Based Computer Systems (KBCS 2000).
  - D'Souza, Rayner, Shivakumar, G., Swathi, D., and Bhattacharyya, P.: 2001, 'Natural Language Generation from Semantic Net like Structures with Application to Hindi,' in Proceedings of International Symposium on Translation Support Systems, IIT Kanpur, India.
  - Gopinathan, S. and S. Kandaswamy (Eds) :1993, *anuvad ki samasayen [Problems of Translation]*, Lokbharti Prakashan.
  - Sinha, R. M. K.: 1994, 'Machine Translation: The Indian Context', in International Conference on Applications of Information Technology in South Asian Languages, AKSHARA'94, New Delhi.
  - Saudagar, Raziq A.: 1999, 'An automated Natural Language Generation for Hindi,' Technical Report, IIT Bombay.
  - Schank, Roger C.: 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 Roger C. Schank and K. M. Colby (eds.), *Computer Models of Thought and Language*, 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.
  - Schubert, K.: 1988, 'The architecture of DLT--interlingual or double direct.' In *New Directions in Machine Translation*, Floris Publications, Dordrecht, Holland.
  - Schutz, J., Thurmair, G., et al. (1991). 'An architecture sketch of Eurotra-II.' In Proceedings of the Third Machine Translation Summit, Carnegie Mellon University.
  - Tiwari, Bholanath and Naresh-Kumar: 1987, *Videshi bhashaon se anuvad ki samasayen [Problems of translation from various foreign languages]*, Prabhat Prakashan.
  - Vauquois, Bernard and Christian Boitet: 1985, 'Automated Translation at Grenoble University,' *Computational Linguistics* 11:1, 28-36.
-

- 
- Vauquois, Bernard: 1975, *La Traduction Automatique à Grenoble*, Dunod, Paris.
  - Witkam, T.: 1988, 'DLT--an industrial R&D project for multilingual machine translation,' In Proceedings of the 12th International Conference on Computational Linguistics, Budapest.
  - Uchida, H.: 1989, 'ATLAS-II: A machine translation system using conceptual structure as an interlingua.' In Proceedings of the Second Machine Translation Summit, Tokyo.
  - UNL: 1998, The Universal Networking Language (UNL) specifications version 3.0. Technical Report, United Nations University, Tokyo.  
URL: <http://www.unl.ias.unu.edu/unlsys/unl/UNL%20Specifications.htm>
  - Uchida, H., Zhu, M. and Della Senta, T.: 2000, *UNL: A Gift for a Millennium.*, The United Nations University.  
URL: <http://www.unl.ias.unu.edu/publications/index.html>
  - Wahlster, W.: 1997, *VERBMOBIL: Translation of Face-toFace Dialogs*, In Proceedings of the 3rd EUROSPEECH, pp. 29--38, Berlin, Germany, 1993. Lieske, Bos, Emele, Gamback, Rupp 4 In Proc. of EuroSpeech'97  
<http://verbmobil.dfki.de/verbmobil/>
  - Wilks, Yorick: 1972, *Grammar, meaning and the machine analysis of language*, Routledge & Kegan Paul, London.
  - Wilks, Yorick: 1987, 'Primitives' in S. C. Shapiro (ed.), *Encyclopedia of Artificial Intelligence*, John Wiley and Sons, New York, NY, 759-761.