

A New Approach of Developing a Deconverting rules for Bangla Language

Aloke Kumar Saha¹, Muhammad F. Mridha¹, Kamal Kanti Biswas², and Jugal Krishna Das²

Department of Computer Science and Engineering¹

University of Asia Pacific, Dhanmondi, Dhaka¹

Department of Computer Science and Engineering²

Jahangirnagar University, Savar, Dhaka²

aloke71@yahoo.com, mdfirozm@yahoo.com, kamalbis@gmail.com, drdas64@yahoo.com

Abstract: The Universal Networking Language (UNL) is a worldwide generalizes form human interactive in machine independent digital platform for defining, representing, storing and dissipating knowledge or information among people of different nations. The theoretical and practical research associated with these interdisciplinary endeavour facilities in a number of practical applications in most domains of human activities such as creating globalization trends of market or geopolitical independence among nations. In this paper we will discuss the interlingua approach to machine translation. Here Universal Networking Language (UNL) has been used as the intermediate representation. In this thesis work, we have dealt with the language independent deconverter for the Bangla language it takes as input a UNL (Universal Networking language) expression. The system takes a set of UNL expression as input and with the help of language independent algorithm and language dependent data generates corresponding Bangla sentence. The process of deconversion involves syntax planning and morphology phase. The syntax planning phase is aimed at generation of proper sequence of words for the target sentence.

Index Terms: Universal Networking Language, morphology, morphological rules, Deconverter, Syntax analysis.

1 Introduction

Today the regional economies, societies, cultures and educations are integrated through a globe-spanning network of communication and trade. This globalization trend evokes for a homogeneous platform so that each member of the platform can apprehend what other intimates and perpetuates the discussion in a mellifluous way. However the barriers of languages throughout the world are continuously obviating the whole world from congregating into a single domain of sharing knowledge and information. The WWW represents a formidable tool for communication and information access. However, despite the abundance of information, languages very often cause problems. When most of the web pages today are written in few most commonly used languages like

English, French, Chinese etc, it becomes difficult for a person with insufficient knowledge of these languages to access and use this tool of communication and information. This has prompted the need to devise means of automatically converting the information from one natural language to another natural language, called Machine Translation. This process needs syntactic and semantic analysis of both source and target languages. Interlingua based machine translation has received a considerable attention because of economy of translation of effort and also additional attraction of the Interlingua providing a knowledge representation scheme.

As a consequence United Nations University/Institute of Advanced Studies (UNU/IAS) were decided to develop an inter-language translation program. The corollary of their continuous research leads a common form of languages known as Universal Networking Language (UNL). For the purpose of conversion we use Interlingua which follow the UNL specifications proposed by UNU/IAS Tokyo. UNL (Universal Networking language) is a language used to represent a semantic graph equivalent of a concept (contained in text document). The system takes a set of UNL expression as input and with the help of language independent algorithm and language dependent data generates corresponding Bangla sentence The Universal Networking Language Programme started in 1996, as an initiative of the Institute of Advanced Studies (IAS) of the United Nations University (UNU) [1, 5] in Tokyo, Japan. The mission of the UNU program is to allow people across nations to access information in Internet in their own languages. The core of the project is UNL, a language independent specification for serving as a common medium for documents in different languages. Researchers involved in this project from different countries have been developing UNL system for their respective native languages. The goal is to eliminate the massive task of translation between two languages and reduce language to language translation to a one time conversion to UNL. For example, Bangla corpora, once converted to UNL, can be translated to any other language given UNL system built for that language. The UNL system does this by representing only the semantics of a native language sentence in a hypergraph. Enconverter [12] converts each native language sentence to a UNL hypergraph and DeConverter [13] translates from hypergraph to any native language. The main aim of the UNL project is to overcome language barriers. This project currently includes 16 official languages. Bangla is not yet included. We have attempted to demonstrate that we can do similar tasks for Bangla as it has been done for other official languages. In this proposal we present a new approach of NLP through UNL for Bangla Language.

2 Benefits UNL

Once the information is converted to UNL form, it becomes language neutral and it can be converted to other different languages. Thus, it can be used for information exchange between languages. Information in a source language can be converted to UNL using source language Deconverter and then using Enconverter of target language, UNL can be enconverted in to that language. Since, UNL is in logical form, knowledge processing can be done unambiguously to produce useful and desired results.

It enables natural language phenomena to be expressed in formal semantic framework which enables computers to understand natural language. If the UNL is added to the network platforms, the communication status will be changed. UNL will make the communication among people through different Natural Languages possible, which will share information and provide a common educational environment as language is an essential part of the communication process. Communication between different nations will be easier since language barriers will be broken. Breaking language barriers, in turn, will result in, for example,

- a) Encouraging mutual understanding among different nations which is one of the ultimate goals of UNL. Sure, using foreign languages will make nations go through the risk of losing a big part of their culture; consequently, as time goes, their roots will be lost as well. With the existence of UNL this risk will not exist.
- b) Communication through UNL will make the mission of international organizations, like United Nations and UNESCO, easier as they are concerned about all people with different mother tongues; one of the main problems faced in the exchange of information between the organizations and different nations is the existence of language barriers.

3 UNL Structure

UNL is an artificial language that allows the processing of information across linguistic barriers [10]. This artificial language has been developed to convey linguistic expressions of natural languages for machine translation purposes. Such information is expressed in an unambiguous way through a semantic network with hyper-nodes. Nodes (that represent concepts) and arcs (that represent relations between concepts) compose the network. UNL contains three main elements:

- Universal Words: Nodes that represent word meaning.
- Relation Labels: Tags that represent the relationship between Universal Words i.e. between two nodes.

Tags are the arcs of UNL hypergraph. Relation: There are 46 types of relations in UNL. For example, agt (agent), agt defines a thing that initiates an action, agt(do, thing), agt(action, thing), obj(thing with attributes) etc.

- Attribute Labels: Additional information about the universal words.

4 UNL Deconverter

A "deconverter" is software that automatically deconverts UNL into native languages. It is also important that the basic architecture of the "deconverter" is widely shared throughout the world, in order to treat all languages with the same quality and precision standards. Technology developed for a language can be applied to other languages as long as the architecture is shared.

A tool called DeCo has been designed by UNU/IAS can deconvert both context-sensitive and context-free languages. It uses target language specific Word Dictionary, Co-occurrence Dictionary and Deconversion rules to generate the target language. So, developing a Deconverter for a language means developing dictionaries and writing deconversion rules, which are understood by the DeCo and these are language dependent. Each entry in

Word Dictionary includes native language Head Word, corresponding UW, and the attributes. Attributes include grammatical and semantic attributes. An example of an entry in Bangla Language Word Dictionary Attributes can be:

[পাখি] {} “bird(icl>animal>animate thing)”(N,ANI,SG, CONCRETE) <B,0,0>

[শহর] {} “city(icl>region)” (N, PLACE) <B,0,0>

Here, [পাখি] is the Bangla Head Word, book (icl>animal>animate thing) is UW and (N,ANI,SG, CONCRETE) is the attribute list. First, the deconversion rules are converted into binary format and then binary format rules are loaded. The UNL expressions are converted in to semantic net called Node-net. The UWs are replaced with corresponding native language Head Words. If it is not possible to unambiguously decide the correct Head Word for a given UW, Co-occurrence dictionary is used. Co-occurrence dictionary contains more semantic information for proper word selection without the ambiguity. But the use of Co-occurrence dictionary is optional. Node-net represents the hyper graph (a representation of UNL expressions) that has not yet been visited. Each node contains certain attributes initially loaded from the Language Dictionary and sometime generated by DeCo during runtime. These attributes can be read or deleted or new attributes can be added. This is governed by deconversion rules. Each node in the Node-net is traversed and inserted in to the Node-list. The result of rule application is operation on the nodes in Node-list like changing attributes, copy, shift, delete, exchange etc. and/or insertion of nodes from Node-net to Node-list. The rule application halts when either Left Generation Window reaches the Sentence Tail node or Right Generation Window reached the Sentence Head node. If post-editing is required the Deconverter will start applying post editing rules. Post editing rule has not been used for UNL Bangla Deconverter. At the end, the nodes in the Node-list represent the generated sentence [13].

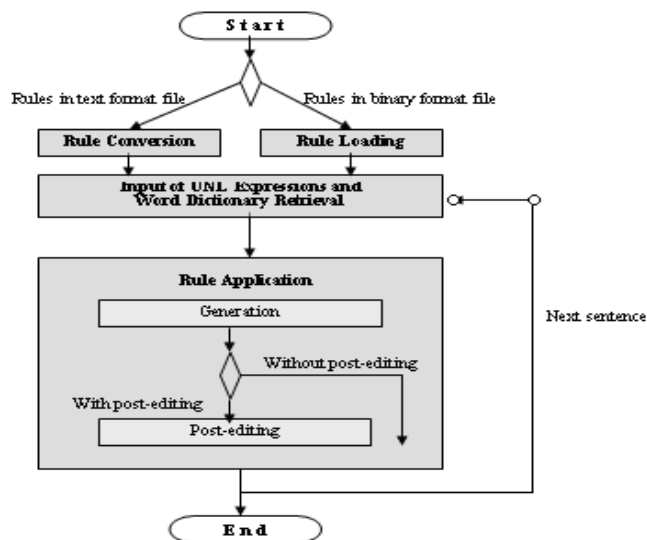


Figure 1: Flowchart of deconversion system

5 Proposed Model

In our proposed model a "Deconverter", which generates Bangla language from UNL, plays a core role in the UNL system. It is very significant that "deconverter" will be capable of expressing UNL information with very high accuracy. It will consist of word dictionary and conversion rules for a language. This will be language independent software that is applicable for any languages. This engine takes UNL expression as input and generates target language (Bangla) sentence with the help of various database files like lexicon files, morphological rule files [12].

In our proposed model we have used two phases to convert from UNL to Bangla as:

- 1) Syntax analysis phase
- 2) Morphological phase

This module is responsible for Bangla sentence formation by syntax analysis phase. The syntax analysis phase is aimed at generation of proper sequence of words for the Bangla language. In order to get the correct Bangla sentence as the output of the DeConverter system, all the rules should be applied in proper order. This module is also responsible for proper word formation through morphology generation. This module generates most of the words. This module handles noun, verb and adjective morphology generation. This module not only inflects the root words, but also introduces conjunctions, case markers and any other new words if necessary. The morphological rules are governed by UNL relations and attributes. Morphological rules due to UNL relations are called relation label morphology.

Syntax analysis is the process of linearizing the Semantic hyper-graph, i.e., it decides the word-order in the generated sentence. To make this process rule driven, we make several important assumptions:

The syntax analysis phase is aimed at generation of proper sequence of words for the target sentence. These phases first reads the input UNL file and convert it into semantic net like structure known as nodenet. We use lexicon files to map the UWs to target language worlds.

The process of deconversion involves syntax analysis phase and morphology phase. The syntax planning analysis is aimed at generation of proper sequence of words for the target sentence. These phases first reads the input UNL file and convert it into semantic-net like structure known as nodenet. Nodenet is a directed acyclic graph structure, which defines the sentence in the form of Directed Acyclic Graph. We use lexicon files to map the UWs to target language worlds. After generating a nodenet, the problem of the syntax plan generation get reduce to the problem of Directed Acyclic Graph traversal. Proper traversal of the node net generates the syntax plan of the target sentence. This syntax plan needs to be processed by the case-marking file, which apply proper case marker for each and every relations. This case-marking phase is next processed by the morphology phase. The morphology phase gives a final form of the target sentence.

6 Deconversion Rules

Deconversion (or “generation” in general) rules describe the conditions for rule application: the way of rewriting the attributes of nodes that satisfy those conditions, as well as the way of composing a native language sentence. DeConverter looks at, and operates on, the nodes in the Node-list through its windows, and the conditions and actions of a deconversion rule are matched to the windows. Each part of the rule expresses the conditions of, or actions on, the adjacent nodes in the Node-list in the order of the Left Condition Windows (LCW or PRE), the Left Generation Window (LGW), the Middle Condition Windows (MCW or MID), the Right Generation Window (RGW), and the Right Condition Windows (RCW or SUF).

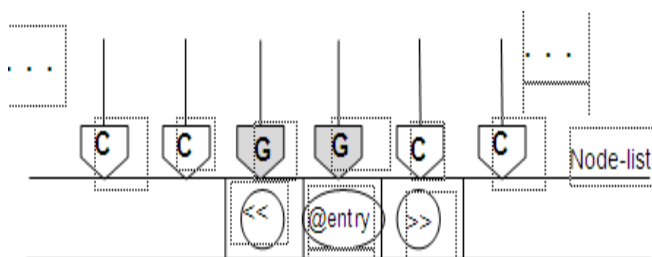


Figure 3: Initial state of the GW and the Node-list

DeCo can input either a string or a list of words form UNL to convert its native language. A list of entry node from UNL must be enclosed by <<> and >> [6]. When we input the word into DeCo, the Sentence Head (<<) will be on LGW, sentence texts/morphemes/words will be on RGW and the Sentence Tail (>>) will be on Right Condition Window (RCW) shown in figure 2. DeCo uses CWs for checking the neighbouring nodes on both sides of the GWs in order to judge whether the neighbouring nodes satisfy the conditions for applying a generation rule or not.

A deconversion rule has the following syntax:

```
<TYPE>
["("<PRE>")" ["*"]]...
{" | "" "" [ <COND1> ] ":" [ <ACTION1> ] ":" [ <RELATION1> ] ":" [ <ROLE1> ] " } | "" ""
    ["("<MID>")" ["*"]]...
{" | "" "" [ <COND2> ] ":" [ <ACTION2> ] ":" [ <RELATION2> ] ":" [ <ROLE2> ] " } | "" ""
["("<SUF>")" ["*"]]...
"P(" <PRIORITY> ");"
```

Some Proposed rule:

Subject Insertion:

Rule 1: : "HPRON,SUBJ,1SG:subj:agt" { V, ^IRG, ^pred:pred,1sg } P120;

Rule 2: : "HPRON,SUBJ,3SG:subj:agt" { V, ^IRG, ^pred:pred,3sg } P120;

Rule 3: : "HPRON,SUBJ,1SG:subj:agt" { V,IRG,@past,ED,^pred:pred,1sg } P120;

Rule 4: : "HPRON,SUBJ,1SG:subj:agt" { V,IRG,@complete,EN,^pred:pred,1sg } P120;

Rule 5: : "HPRON,SUBJ,1SG:subj:agt" { V,IRG,^@past,^@complete,^ED,^EN,^pred:pred,1sg } P120;

Rule 6: : "HPRON,SUBJ,3SG:subj:agt" { V,IRG,@past,ED,^pred:pred,3sg } P120;

Rule 7: : "HPRON,SUBJ,3SG:subj:agt" { V,IRG,@complete,EN,^pred:pred,3sg } P120;

Rule 8: : "HPRON,SUBJ,3SG:subj:agt" { V,IRG,^@past,^@complete,^ED,^EN,^pred:pred,3sg } P120;

Object insertion:

Rule 9: : { V,VDO,pred,^OBJ_inserted:OBJ_inserted }":obj:obj" P100;

Verb Insertion:

Rule 10: : { V,pred } "N,TIME::tim" P110;

Article insertion:

Rule 11: : "[a],ART" { N,^PRON,^TIME,^@pl,^VOW,^art:art } P100;

Rule 12: : "[an],ART" { N,^PRON,^TIME,^@pl,VOW,^art:art } P100;

7 Conclusions

This paper has described the development of UNL Bangla Deconverter, a Bangla language generator. Using the UNL system with its language components it has been proved to be a powerful environment for man machine communication. On the other hand, other machine translation systems will not be able to provide such environment for education and exchange of information as they are away from universality. They will never be inter-lingual. This will make their value limited to only the one or two languages involved in the translation. Consequently, communication and the distribution of information will be negatively affected. However one drawback can be that UNL has not yet conceived as a fully automatic machine translation system. The Bangla language could successfully be generated from UNL hyper semantic networks with a high degree of accuracy. The main skeleton of Bangla sentence structure has been handled however many problems remain unsolved such as generating passive structures, correct ordering of modifiers of the same type, selecting the correct word representing universal words which represents the main challenges of the future work.

REFERENCES

1. Muhammad Firoz Mridha, Mohammad Nurul Huda, Chowdhury Mofizur Rahman, Jugal Krishna Das, "Development of Morphological Rules for Bangla Root and Verbal Suffix for Universal Networking Language". 6th International Conference on Electrical and Computer Engineering, ICECE 2010, 18-20 December 2010, Dhaka, Bangladesh.
2. Muhammad Firoz Mridha, Manoj Banik, Md. Nawab Yousuf Ali, Mohammad Nurul Huda, Chowdhury Mofizur Rahman, Jugal Krishna Das, "Formation of Bangla Word Dictionary Compatible with UNL Structure," SKIMA'10, Paro, Bhutan, August, 2010.
3. H. Uchida, M. Zhu, and T. C. D. Senta, Universal Networking Language, NDL Foundation, International environment house, 2005/6, Geneva, Switzerland.
4. D. M. Shahidullah, "Bangala Vyakaran", Maola Brothers Prokashoni, Dhaka, August 2003, pp.110-130
5. Muhammad Firoz Mridha, Kamruddin Md. Nur, Manoj Banik and Mohammad Nurul Huda, "Structure of Dictionary Entries of Bangla Morphemes for Morphological Rule Generation for Universal Networking Language". International Journal of Computer Information Systems and Industrial Management Applications (IJCISIM) 2011
6. Muhammad Firoz Mridha, Kamruddin Md. Nur, Manoj Banik and Mohammad Nurul Huda, "Generation of Attributes for Bangla Words for Universal Networking Language(UNL)". International Journal of Computer Information Systems and Industrial Management Applications (IJCISIM) 2011.
7. H. Uchida, M. Zhu, "The Universal Networking Language (UNL) Specification Version 3.0", Technical Report, United Nations University, Tokyo, 1998
8. Muhammad Firoz Mridha, Manoj Banik, Md. Nawab Yousuf Ali, Mohammad Nurul Huda, Chowdhury Mofizur Rahman, Jugal Krishna Das, "Conversion of Bangla to UNL' ", SKIMA'10, Paro, Bhutan, August, 2010.
9. S. Dashgupta, N. Khan, D.S.H. Pavel, A.I. Sarkar, M. Khan, "Morphological Analysis of Inflecting Compound words in Bangla", International Conference on Computer, and Communication Engineering (ICCCIT), Dhaka, 2005, pp. 110-117
10. Bangla Academy, "Bengali-English Dictionary" Bangla Academy Dhaka 2007
11. <http://www.unl.ru/deco>
12. EnConverter Specification, Version 3.0, UNL Center, UNDL Foundation, Tokyo 150-8304, Japan 2002
13. DeConverter Specification, Version 2.7, UNL Center, UNDL Foundation, Tokyo 150-8304, Japan 2002.