

Universal Networking Language Based Analysis and Generation for Bengali Case Structure Constructs

Kuntal Dey and Pushpak Bhattacharyya

¹ **Abstract**– *Case structure analysis forms the foundation for any natural language processing task. In this paper we present the computational analysis of the complex case structure of Bengali- a member of the Indo Aryan family of languages- with a view toward interlingua based MT. Bengali is ranked 4th in the list of languages ordered according to the size of the population that speaks the language. Extremely interesting language phenomena involving morphology, case structure, word order and word senses makes the processing of Bengali a worthwhile and challenging proposition. A recently proposed scheme called the Universal Networking Language has been used as the interlingua. The approach is adaptable to other members of the vast Indo Aryan language family. The parallel development of both the analyzer and the generator system leads to an insightful intra-system verification process in place. Our approach is rule based and makes use of authoritative treatises on Bengali grammar.*

Keywords– *Universal Networking Language, Enconversion, Deconversion, Rule Theory, Bengali, kaarak (case), kriyaa (verb).*

1 Introduction

Bengali is spoken by about 189 million people and is ranked 4th in the world in terms of the number of people speaking the language ([2]). Like most languages in the Indo Aryan family, descended from Sanskrit, Bengali has the SOV structure with some typical characteristics. A motivating factor for creating a system for processing Bengali is the possibility of laying the framework for processing many other Indian languages too.

Work on Indian language processing abounds. *Project Anubaad* [18] for machine translation from English to Bengali in the newspaper domain uses the *direct translation approach*. *Angalabharati* [21] system for English Hindi machine translation is based on pattern directed rules for English, which generates a *pseudo-target-language* applicable to a group of Indian Languages. In MATRA [11], a web based MT system for English to Hindi in the newspaper domain, the input text is transformed into case-frame like structures and the the target language is generated by

parameterized templates. The *MANTRA* MT system for official documents uses Tree Adjoining Grammar (TAG) to achieve English Hindi MT ([1]). Project *Anusaaraka* [6] is a language accessor system rather than an MT system and addresses multiple Indian languages. Interlingua based MT for English, Hindi and Marathi [23] [22], that uses the UNL, transforms the source text into the *UNL representation* and generates target text from this intermediate representation. References to most of these works can also be found at [3]. Other famous MT systems are *Pivot* [20], *Atlas* [15], *Kant* [26], *Aries* [10], *Geta* [9], *SysTran* [17] etc.

The *Universal Networking Language (UNL)* has been defined as a digital meta language to describe, summarize, refine, store and disseminate information in a machine independent and human language neutral form. The information in a document is represented sentence by sentence. Each sentence is converted into a directed hyper graph having concepts as nodes and relations as arcs ([4]). Knowledge within a document is expressed in three dimensions:

1. Word Knowledge is expressed by Universal Words (UWs) which are language independent. These UWs are tagged using restrictions describing the sense of the word in the current context. For example, *drink(icl > liquor)* denotes the noun sense of *drink* restricting the sense to a type of *liquor*. Here, *icl* stands for inclusion and forms an *is-a* relationship like in semantic nets [7].
2. Conceptual Knowledge is captured by relating UWs through a set of UNL relations [14]. For example, *Humans affect the environment* in UNL is:

```
agt(affect(icl>do)).@present.@entry, human(icl>animal).@pl  
obj(affect(icl>do)).@present.@entry, environment  
(icl>abstract thing).@pl
```

agt means the *agent* and *obj* the *object*. *affect(icl > do)*, *human(icl > animal)* and *environment(icl > abstract thing)* are the UWs denoting concepts.
3. Speaker's view, aspect, time of event, etc. are captured by UNL attributes. 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 above discussion can be summarized using the example and it's UNL: *John, who is the chairman of the company, has arranged a meeting at his residence.*

¹ Kuntal Dey: u2ckuntal@yahoo.com

Pushpak Bhattacharyya: pb@cse.iitb.ac.in

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

```

;===== UNL =====
mod(chairman(icl>post).@present.@def,company
(icl>institution).@def)
aoj(chairman(icl>post).@present.@def, John
(icl>person))
agt(arrange(icl>do).@entry.@present.@complete,
John(icl>person))
pos(residence(icl>shelter), John(icl>person))
obj(arrange(icl>do).@entry.@present.@complete,
meeting(icl>event).@indef)
plc(arrange(icl>do).@entry.@present.@complete,
residence(icl>shelter))
;=====

```

In the expressions 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 detailed specification of UNL can be found at [5].

Our work is based on an authoritative treatise on Bengali grammar [8]. The strategies of analysis and generation of linguistic phenomena have been guided by rigorous grammatical principles.

2 EnConverter and DeConverter machines

The EnConverter (henceforth called *EnCo*) [25] is a language-independent parser, a multi-headed Turing machine [13] providing a framework for morphological, syntactic and semantic analysis synchronously using the UW dictionary and analysis rules. More details about the machine can be found at [25].

The machine has two types of *heads- processing heads* and *context heads*. The processing heads (2 nos.) are called *Analysis Windows (AW)* and the context heads called *Condition Windows (CW)*. The machine traverses the sentence back and forth, retrieves the relevant universal words (UW) from the lexicon and, depending on the *attributes* of the nodes under the AWs and those under the surrounding CWs, generates semantic relations between the UWs and/or attaches speech act attributes to them. The final output is a set of UNL expressions equivalent to a UNL graph.

The DeConverter (henceforth called the *DeCo*) [24] is a language-independent generator that produces sentences from UNL graphs. Like EnCo, DeCo too is a multi-headed Turing Machine. It does syntactic and morphological generation synchronously using the lexicon and the set of generation rules.

3 Rule theory

EnCo and DeCo are driven by *analysis rules* and *generation rules* respectively. These rules are *condition-action structures* that can be looked upon as *program* written in a

specialized language to process various complex phenomena of a natural language, both for analysis and generation. They have the following format:

```

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

```

Characters between double quotes are the predefined delimiters of the rule. The rules mean that

- **IF**
under the *left processing window* there is a node satisfying <COND1> and under the *right processing window* a node satisfying <COND2> attributes, and there are nodes that fulfill the conditions in <PRE>, <MID> and <SUF> in the order of left, middle and right sides of processing windows respectively,
THEN
the lexical attributes in processing windows are rewritten according to the <ACTION1> and <ACTION2> as specified in rule, and new attributes added if necessary. (By *processing window*, *analysis window* is meant for the *enconversion* and *generation window* for the *deconversion* process).
- The operations are done on the node-list depending on the <TYPE> of the rule. <RELATION1> describes the semantic relation of the node on right processing window to the node on left processing window and <RELATION2> describes the reverse [22].
- <PRIORITY> describes the interpretation order of the rules, whose value lies between 0-255. Larger number indicates higher priority. Matching rule with the highest priority is selected for multiple matching rules.

A sequence of such rules get activated depending on the sentence situation (the conditions of the nodes under the analysis/generation windows). These are the lexico-morpho-grammatical-semantic attributes of the words under processing. For example, for a sentence like *John laughs*, the *animate* attribute of *John*, the *verb* attribute of *laugh* and the *adjacency* of these two words under the analysis windows dictate with high probability establishing the *agt* (*agent*) relation between the corresponding two nodes in the UNL graph.

In order to adapt the UNL engines to *enconvert* the Bengali sentences into the UNL interlingua and to *deconvert* the UNL interlingua/graph into Bengali sentences, an *enconverter rule-base* and a *deconverter rule-base* have been written. The rules within the rule-base are compliant with the corresponding UNL engines and are focused to deal with the Bengali language structure.

4 Case Structure in Bengali: Kaaraks

In the Indian linguistic system- descended from Sanskrit- the *case constructs* are called *kaaraks* [12]. As in the traditional understanding, they denote the relationship of the nominals with the main verb of the clause except in the *genitive case* where two nominals are related to each other. The case structure in Bengali is complex. The *kaaraks* are broadly classified into 6 types [8], each having a finer categorization into sub-types. The correspondence between the Bengali *kaarak* system and the traditional linguistic concept of case [16] is shown by means of table 1. The *bibhakti signs* are the case markers. An exhaustive study of the *kaarak* system with a view to analyzing Bengali into UNL has been carried out. The foundation of this work is the *kaarak* theory [8]. Due to the word limitation, we exemplify the work with only the first *kaarak*, viz., the *kartri kaarak*.

Classical case	Corresponding Bengali kaarak	Bibhakti (Case Marker)
Nominative case	<i>Kartri kaarak</i>	None
Accusative case	<i>Karma kaarak</i>	<i>ke, re, ere</i>
Instrumental case	<i>Karan kaarak</i>	<i>dwaaraa, diye, diya, kartrik</i>
Dative case	<i>Sampradaan kaarak</i>	<i>janya, nimitta, ke</i>
Ablative case	<i>Apaadaan kaarak</i>	<i>theke, haite</i>
Genitive case	<i>Sambandha pad</i>	<i>r, er</i>
Case of time-place	<i>Adhikaran kaarak</i>	<i>e, te, ete</i>

Table 1: CASE-KAARAK CORRESPONDENCE

4.1 Kartri kaarak

Kartri kaarak denotes the *agent* of the action stated by the verb. It is divided into the following classes:

1. **Projojak kartaa** (প্রয়োজক কর্তা): The agent *causes* some event to take place, tending to compel the event to happen. The morphology of the verb is exploited and the extracted knowledge has the *causative* attribute. Example:
টম জনকে খেলাবে
tama janake khelaabe.
Tom John-to will-make-play.
Tom will make John play.
2. **Nirapekkha kartaa** (নিরপেক্ষকর্তা): Here there are more than one verb in the sentence with at least one *অসমাপিকা* (non-finite) verb and one *সমাপিকা* (finite) verb, and the *kartaas*, i.e., *agents* for these verbs are different or not related. The *kartaa* associated with the non-finite verb is called the *nirapekkha kartaa* (*nominative absolute* in English). As there is an *অসমাপিকা* verb involved, a *con* or *seq etc.* relation

is generated, also there is a possible generation of compound UW. Example:

টম খেলে জন খাবে
tama khele jana khaabe.
Tom if-eats John will-eat.
If Tom eats John will eat.

3. **Karmakartribaachyer kartaa** (কর্মকর্ত্বাচ্যের কর্তা): Here, the actual *kartaa* is not present, and hence the *karma*, i.e., the *object* acts as the *kartaa*. As a result, there is no *agt* or equivalent relation generated for conceptualizing an agent of the sentence, instead, an *obj* relation is realized. Example:
বালতি ভরেছে
baalti bhareche.
Bucket has-filled-up.
The bucket has filled up.
4. **Anukta kartaa** (অনুক্ত কর্তা): In cases of কর্মবাচ্য (*karma baachya*) and ভাববাচ্য (*bhaab baachya*) (which are variants of the passive voice), the *kartaa* is not emphasized on. Example:
টমের আজ খাওয়া হয় নি
tamer aaj khaoyaa hay ni.
Tom-of today eating not-happened.
Eating has not happened to Tom today.
5. **Sahajogi kartaa** (সহযোগী কর্তা): Two *kartaas* are present in the same sentence, co-acting with each other to perform the action specified by the verb. Example:
বাঘে গোরুতে খাচ্ছে
baaghe gorute khaacche.
Tiger cow eating.
Tiger is eating with cow.
6. **Bakyangsha kartaa** (বাক্যাংশ কর্তা): Here the noun phrase as a unit acts as the *kartaa*. A noticeable fact is that this noun phrase does not have any *সমাপিকা* (finite) verb. Example:
সৎপথে জীবনযাপন করা কঠিন কাজ
satpathe jiibanjaapan karaa kathin kaaj.
Honest-way-in leading-life hard work.
Leading a life in an honest way is hard work. (Note: Here *hard work* means *difficult*.)
7. **Upabakiya kartaa** (উপবাক্যীয় কর্তা): Here there is a noun clause in the sentence. This conceptually acts as the *kartaa*. However, to retain the *person* information present in the verb, a different term causing *agt* relation has to be introduced in the sentence during enconversion. The conceptual *kartaa* actually does not get identified as a *kartaa*, instead it is identified as something different (for example, *karma*). Example:
ভয় কাকে বলে জানি
bhay kaake bale jaani.
Fear to-whom call I-know.
I know what is called fear.
8. **Karta with 'e' bibhakti** (কর্তায় এ বিভক্তি): In spite of the presence of the *e* (এ) *bibhakti*, the *kartaa* has to be identified as an *agt* or equivalent relation. A

salient point to note is that the *e* bibhakti can be used with all other *kaaraks* as well, so appropriate analysis has to be done to identify its functionality. Often the context of occurrence of the word and the grammatical attributes available with the word from the lexical dictionary guide in identifying the *kaarak* in case of *e bibhakti*. Example:

ছাগলে ঘাস খায়
chaagale ghaash khaay.
 Goat grass eat.
 Goat eats grass.

(UNL relations generated for *kartri kaarak*: *agent (agt)*, *co-agent (cag)*, *partner (ptn)* etc.).

4.2 Other *kaaraks*

Five other *kaaraks* and two other related structures have been analyzed exhaustively as above.

1. *Karma kaarak* (6 subcategories): *Karma kaarak* is the person or thing on which the *kartri kaarak* executes the verb. (UNL relations for *karma kaarak*: *object (obj)*, *beneficiary (ben)*, *co-object (cob)*).
2. *Karan kaarak* (5 subcategories): *Karan kaarak* is the thing, tool or method by which the *kartri kaarak* of the sentence executes the specified action. (UNL relations for *karan kaarak*: *instrument (ins)*, *method (met)*).
3. *Sampradaan kaarak* (2 subcategories): *Sampradaan kaarak*s are cases where the agent (*kartri kaarak*) does or gives away something for or to someone. (UNL relations for *sampradaan kaarak*: *beneficiary (ben)*, *goal (gol)*, *purpose (pur)*, *reason (rsn)*).
4. *Apaadaan kaarak* (6 subcategories): This stands for the concept of sources of creation, location, position etc. All types of relations having the concept of *source* in some sense are eligible to come into this category. (UNL relations for *apaadaan kaarak*: *place-from (plf)*, *time-from (tmf)*, *from (frm)*, *source (src)*).
5. *Sambandha pad* (4 subcategories): If related to the next noun or pronoun, then the term having a *r* (র) or *er* (এর) *bibhakti* is called a *sambandha pad*. *Sambandha pad* always has some *bibhakti* (never *sunya bibhakti*). (UNL relations for *sambandha pad*: *modifier (mod)*, *possession (pos)*, *part-of (pof)*).
6. *Adhikaran kaarak* (8 subcategories): *Adhikaran kaarak*s are the ones that describe the place, time and topic of the action performed by the sentence. (UNL relations for *adhikaran kaarak*: *place (plc)*, *time (tim)*, *place-to (plt)*, *time-to (tmt)*, *to (to)*, *goal (gol)*, *virtual-place (scn)*, *objectified-place (opl)*).
7. *Sambodhan* (3 subcategories): *Sambodhan* (সম্বোধন) is the case where someone hails some other person and says something to this person. This act of hailing is captured by what is called সম্বোধন. This generates a *@vocative* attribute against the called person's appearance in the UNL graph.

Kaarak	Corresponding UNL Relations
<i>Kartri kaarak</i>	agt, cag, ptn, aoj, cao
<i>Karma kaarak</i>	obj, ben, cob
<i>Karan kaarak</i>	ins, met
<i>Sampradaan kaarak</i>	ben, gol, pur, rsn
<i>Apaadaan kaarak</i>	frm, src, plf, tmf
<i>Sambandha pad</i>	mod, pos, pof
<i>Adhikaran kaarak</i>	plc, plt, tim, tmt, to, gol, scn, opl

Table 2: KAARAKS VERSUS UNL RELATIONS

Table 2 summarizes the correspondence between Bengali *kaaraks* and the *UNL relations*.

The UNL relations that are not covered by the *kaaraks* in Bengali are: *and (and)*, *or (or)*, *quantity (qua)*, *proportion, rate or distribution (per)*, *content (cnt)*, *via (via)*, *condition (con)*, *sequence (seq)*, *co-occurrence (coo)*, *basis for expressing degree (bas)*, *duration (dur)*, *range: from-to (fmt)* and *manner (man)*.

5 *Kaarak* enconversion strategy

The basic idea follows. The non-verb primary (non-case [19]) words appearing in the sentences are one of the two types: (i) A word denoting a concept, which is a *kaarak* or *sambandha pad* or *sambodhan*, (ii) A word or *bibhakti* causing a conceptual relation to link two concepts. *Kaaraks*, *sambandha pads* and *sambodhans* get mapped to the UNL word concepts (UWs) after the analysis and appear in the UNL graph as **nodes**. The *bibhaktis* or conceptually relating words result in forming the **edges** of the graph which embed the logical relation between the two UWs. The lexical, morphological and semantic attributes in the dictionary entries of the UWs are also used to analyze the input. We illustrate with the following input to the enconverter:

কীর্তনে এবং বাউল গানে আমি মাতিয়ে রাখবো
 (kiirtane ebang baul gaane aami maatiye raakhbo)
 Kiirtan-by and baul song-by I enchant-will
 I will enchant with Kirtan and baul song

Strategy:

- Adding the *e* (এ) *bibhakti* to an abstract noun makes it a candidate for the *met* relation, so a *+MET* is added.
- Finally, a *met* relation is resolved when the node with the *MET* attribute and the verb become juxtaposed.

Salient rules:

- $+ \{N, Na, ABS, ^PLACE, ^CONCRETE, ^SCN, ^RSN, ^TIME, ^BLKINSERT: +MET, +MORADD, +eADD, +BLKINSERT::\} \{[e], NMOR, BLKINSERT::\} P30;$
- $> \{N, MET, ABS, ^V::met:\} \{V, ^METRES, : +METRES::\} P20;$

UNL:

met(enchant(icl>do):0T.@entry.@future,:01)
agt(enchant(icl>do):0T.@entry.@future,I(icl>person):0P)

and:01(song(icl>song):0K.@entry,kirwana(icl>song):00)
 mod:01(song(icl>song):0K.@entry,bAula(icl>song):0E)

The example gives a flavor of the procedure involved. Similar procedure is applied on all the various categories and sub-categories. (*Kirtan* and *baauk*: two Indian blends of songs.)

6 Verification

An exhaustive verification of the system has been carried out by writing a **UNL to Bengali Deconverter** (*i.e.* generator). This uses the same lexicon as the *Bengali enconversion* system and a set of *Bengali generation rules*. The enconverted input sentences have been re-generated from the UNL graphs and manually matched for conceptual equivalence. This is a form of intra-platform verification. It verifies both the preservation of information and meaning during enconversion and its wholesome retrieval during deconversion using the appropriate rule-bases. Examples follow. Many of the output sentences map back exactly to the same set of words and sentence structure as the input, with no divergence. However, to provide a more interesting delineation (within this short span) of the challenges faced, we mainly present instances of input output divergence.

1. Projojak karta (প্রযোজক কর্তা):

Input to enco: tama janake khelaabe
 Equivalent: টম জনকে খেলাবে
 Gloss: Tom John-to will-make-play
 Meaning: Tom will make John play.
 Output of deco: tama janake khelaabe
 Equivalent: টম জনকে খেলাবে
 Gloss: Tom John-to will-make-play
Remark: Exact match between input and output.

2. Nirapekkha karta (নিরপেক্ষকর্তা):

Input to enco: tama khele jana khaabe
 Equivalent: টম খেলে জন খাবে
 Gloss: Tom if-eats John will-eat
 Meaning: John will eat if Tom eats.
 Output of deco: jadi tama khaay jana khaabe
 Equivalent: যদি টম খায় জন খাবে
 Gloss: If Tom eats John will-eat
Remark: This is an interesting case where the *jadi* (*if*) clause has got introduced into the output of the deconverter while it was not explicitly present in the input to the enconverter. However, it is correct as these sentences have the same sense conceptually.

3. Upabakiyi karta (উপবাক্যীয় কর্তা):

Input to enco: bhay kaake bale jaani
 Equivalent: ভয় কাকে বলে জানি
 Gloss: Fear to-whom call I-know
 Meaning: (I) know what is called fear.
 Output of deco: aami jaani bhay kaake bale
 Equivalent: আমি জানি ভয় কাকে বলে
 Gloss: I know fear to-whom call
Remark: An explicit *aami* (*I*) has been introduced in the generated sentence.

4. Bakyangsha karma (noun phrase as an object) (বাক্যাংশ কর্ম):

Input to enco: aamtaa aamtaa kathaa balte bhaalobaasi naa

Equivalent: আমতা আমতা কথা বলতে ভালোবাসি না
 Gloss: Soft soft to-talk I-like not
 Meaning: (I) don't like to talk softly.

Output of deco: aami bhaalobaasi naa aamtaa aamtaa kathaa balte

Equivalent: আমি ভালোবাসি না আমতা আমতা কথা বলতে

Gloss: I like not soft soft to-talk

Remark: The structures and word-order of input and output differ, but conceptually they are same.

5. Karmer bipsaa (কর্মের বীঙ্গা) (Repetition in Karma):

Input to enco: kii kii caao bali
 Equivalent: কী কী চাও বলি

Gloss: What what you-want I-say

Meaning: (I)/(Let me) say what (you) want.

Output of deco: aami bali tomraa kii kii caao

Equivalent: আমি বলি তোমরা কী কী চাও

Gloss: I say you what what want

Remark: The input to enco has no default number information associated with the person. The output generates (by default implementation as per the rule base) a singular number output for the first person and a plural number output for the second person. Hence, *aami*, which means *I* (first person singular number) and *tomraa*, which means *you* (second person plural number), have been explicitly added to the output.

6. Karaner bipsaa (করণের বীঙ্গা) (Repetition in Karan):

Input to enco: taaraay taaraay bharaa raater aakaash
 Equivalent: তারায় তারায় ভরা রাতের আকাশ

Gloss: Star-with star-with filled night's sky

Meaning: (The) night's sky is filled with stars.

Output of deco: raater aakaash taaraader diye bhareche

Equivalent: রাতের আকাশ তারাদের দিয়ে ভরেছে

Gloss: Night's sky stars-with has-filled

Remark: Here the structural and morphological differences in the input and output is noticeable, although the conceptual meanings are same for both.

7. Sunya (no) bibhakti in karan (করণে শূন্যবিভক্তি):

Input to enco: gaadhaake haajaar caabuk maarleo se ghozDaa hay naa

Equivalent: গাধাকে হাজার চাবুক মারলেও সে ঘোড়া হয় না

Gloss: Donkey-to thousand whiplash in-spite-of-beating-with it horse become not

Meaning: In spite of thousand beatings with whiplashes a donkey does not become a horse.

Output of deco: jadi tomraa haajaar caabuk diye gaadhaake maaro tabuo se ghozDaa hay naa

Equivalent: যদি তোমরা হাজার চাবুক দিয়ে গাধাকে মারো তবুও সে ঘোড়া হয় না

Gloss: If you thousand whiplash with donkey-to beat yet it horse become not

Remark: The output is a complex sentence while the input is not, yet conceptually they mean the same.

8. **Asamaapikaa kriyaabaachak (infinite verb-related) apaadaan kaarak** (অসমাপিকা ক্রিয়াবাচক অপাদান কারক):
 Input to enco: aami marte bhiita nai
 Equivalent: আমি মরতে ভীত নই
 Gloss: I to-die afraid not
 Meaning: I am not afraid to die.
 Output of deco: maraar janya aami bhiita nai
 Equivalent: মরার জন্য আমি ভীত নই
 Gloss: To-die I afraid not
Remark: These differ by the *anusarga* (*janya* in the output), but the input means the same in Bengali as the output in spite of this difference in construction.
9. **Saamipyaa suchak (proximity-denoting) adhikaran kaarak** (সামীপ্যসূচক অধিকরণ কারক):
 Input to enco: tama darajaay daazDiye brishti dekhche
 Equivalent: টম দরজায় দাড়িয়ে বৃষ্টি দেখছে
 Gloss: Tom at-door standing rainfall seeing
 Meaning: Tom is seeing rainfall standing at the door.
 Output of deco: tama darajaay daazDiye daazDiye brishti dekhche
 Equivalent: টম দরজায় দাড়িয়ে দাড়িয়ে বৃষ্টি দেখছে
 Gloss: Tom at-door standing standing rainfall seeing
Remark: These two mean the same, although the word *daazDiye* comes twice in the deco output (to ensure the *coo* concept) in spite of it being present only once in the enco input.
10. **Bishayaadhikaran (topic denoting adhikaran) kaarak** (বিষয়াধিকরণ কারক):
 Input to enco: se taase pokta ebang futbale ostaad
 Equivalent: সে তাসে পোক্ত এবং ফুটবলে ওস্তাদ
 Gloss: He in-cards solid and in-football expert
 Meaning: He is solid in cards and expert in football.
 Output of deco: futbale ostaad ebang se taase pokta
 Equivalent: ফুটবলে ওস্তাদ এবং সে তাসে পোক্ত
 Gloss: In-football expert and he in-cards solid
Remark: This is an instance of free-format input natural language. The output structure varies significantly from the input structure, in spite of having the same meaning and hence being correct.

7 Conclusion

Systematic analysis of case structure forms the foundation for any natural language processing system. In this paper, we have described a system for the computational analysis of the Bengali case structure for the purpose of interlingua based MT using UNL. The complementary generator system too has been implemented, which provides the platform for intra system verification. Verification via cross system generation is being done using the Hindi generation system (also under development.) Apart from the case structure, computational analysis based on authoritative grammatical treatise, addressing complex phenomena involving verbs, adjectives and adverbs is under way.

References

- [1] <http://www.cdacindia.com/html/about/success/mantra.asp>.
- [2] <http://www.harpercollege.edu/~mhealy/g101ilec/intro/clt/cltclt/top100.html>.
- [3] <http://www.tdil.mit.gov.in/mat/ach-mat.htm>.
- [4] <http://www.unl.ias.unu.edu>.
- [5] <http://www.unl.ias.unu.edu/unlsys>.
- [6] Bharati A., Chaitanya V., and Sanyal R. *Natural Language Processing: A Paninian Perspective*. Prentice Hall India Private Limited, Oct 1996.
- [7] Woods W. A. *What's in a link: Foundations for semantic networks*. Morgan Kaufmann Publishers, Inc.
- [8] Chakrabarti B. *Uchchatarata Bangla Byakaran*. Akshay Malancha, Oct 1963.
- [9] Vauquois B. and Boitet C. Automated translation at grenoble university. acl.ldc.upenn.edu/J/J85/J85-1003.pdf <<http://acl.ldc.upenn.edu/J/J85/J85-1003.pdf>>.
- [10] Gonzalez J. C., Gon J. M., and Nieto A. F. *ARIES: A ready for use Platform for Engineering Spanish-Processing Tools*. London, October 1995.
- [11] Rao D., Mohanraj K., Hedge J., Mehta V., and Mahadane P. *A Practical Framework for Syntactic Transfer of Compound-Complex Sentences for English-Hindi Machine Translation*. International Conference on Knowledge Base Computer Systems, Mumbai, 2000.
- [12] Shastri C. D. *Panini Re-interpreted*. Motilal Banarasidass, New Delhi, 1990.
- [13] Hopcroft J. E. and Ullman J. D. *Introduction to Automata Theory, Languages and Computation*. Addison-Wiseley Publishing Company, 1989.
- [14] UNL Centre/UNDL Foundation. *The Universal Networking Language (UNL) Specifications*. November 2001.
- [15] Uchida H. Atlas. *MT Summit II*, pp. 152-157, 1989.
- [16] Fillmore C. J. *The case for case*. E Bach and R Harms (eds.), New York: Holt, Rinehart and Winston, 1968.
- [17] Hutchins W. J. and Somers H. L. *An introduction to Machine Translation*. London: Academic Press, 1992.
- [18] Dey K. *Project Anubaaad: an English Bengali MT System*. Bachelor of Engineering Dissertation, May 2001.
- [19] Dey K., Dubey S. K., and Bhattacharyya P. *Knowledge Extraction From Indo-Aryan Family of Languages Using A Rule Based Approach*. Dec 2002.
- [20] Muraki K. *PIVOT: Two-Phase Machine Translation System*. Hanakone, Japan, 1987.
- [21] Sinha R. M. K. *Machine Translation: The Indian Context*. International Conference on Applications of Information Technology in South Asian Languages, 1994.
- [22] Monju M., Dave S., and Bhattacharyya P. *Knowledge Extraction from Hindi Texts*. Knowledge Based Computer Systems, 2000.
- [23] Dave S., Bhattacharya P., and Girishbhai J. P. *Interlingua based English-Hindi Machine Translation and Language Divergence*. Sep 2002.
- [24] UNU/IAS. *De Converter Specifications*. UNU/IAS UNL Center, Nov 1998.
- [25] UNU/IAS. *En Converter*. UNL Centre/UNDL Foundation, May 2001.
- [26] Lonsdale D. W., Franz A. M., and Leavitt J. R. R. *Large-Scale Machine Translation: An Interlingua Approach*. Center for Machine Translation, Carnegie Mellon University.