# Manipuri Morphological Analysis

**Y. Bablu Singh[1], Th. Mamata Devi[2] and Ch Yashawanta Singh[3]**
[1&2]Computer Science Department, Manipur University, Manipur, India
[3]Linguistics Department, Manipur University, Canchipur, Manipur, India
E-mail: yumnambablu@gmail.com, mamta_th@rediffmail.com, chungkhamyash@gmail.com

*Abstract* - **Morphological analysis is the basic foundation in Natural Language Processing applications including Syntax Parsing, Machine Translation (MT), Information Retrieval (IR) and Automatic Indexing. Morphological Analysis can provide valuable information for computer based linguistics task such as Lemmatization and studies of internal structure of the words or the feature values of the word. Computational Morphology is the application of morphological rules in the field of Computational Linguistics, and it is the emerging area in AI, which studies the structure of words, which are formed by combining smaller units of linguistics information, called morphemes: the building blocks of words. It provides about Semantic and Syntactic role in a sentence. It can analyze the Manipuri word forms and produces grammatical information, which is associated with the lexicon. Morphological Analyzer for Manipuri language has been tested on 4500 Manipuri lexicons in Shakti Standard Format (SSF) using Meitei Mayek Unicode as source; thereby an accuracy of 84% has been obtained on a manual check.**
*Keywords:* **Meitei Mayek, Morphological analysis; Machine Translation; Computational Morphology, Information Retrieval, SSF.**

## I. INTRODUCTION

The first step in natural language processing is to identify words in a sentence. The process is called morphological analysis. The Manipuri Morphological Analyzer is built using the methodology of finite-state compilers and algorithms, and the results are stored and run as finite-state transducers. Manipuri language also known as Meiteilon belongs to the Kuki-Chin [1] branch of the Tibeto-Burman language, sub-family of Sino Tibetan Language. It is an official language as well as a Lingua franca among the various speech communities [2]. Manipuri has been adopted as the medium of instruction and examination from the primary to the high school stage. Meiteilon has been the state language of Manipur since the 8th century A.D. It has been recognized as the 8th scheduled language in the Indian Constitution since 1992 [3]. Morphology consists of two branches: Inflectional morphology and Derivational morphology. Inflectional morphology is the study of those processes of word formation where various inflectional forms are formed from the existing stems [4].

Example:

> Plurals: boroI–>boroIsiQ
> Aspect: chai –>chaari

Derivational morphology is those processes of the word formation where new words are formed from the existing stems through the addition of morphemes. The meaning of the resultant new word is different from the original word and it often belongs to a different syntactic category [4].

Example:

> Adjective to Verb: acAb –>cAb
> Adjective to Adjective: acAb –>acAbgI
> Noun to Adjective: cArIb –>acAb

## II. PROBELEM STATEMENT

For the problem statement we will list five main objective questions as follows:

1. What are morphological categories in Manipuri Language?
2. What are computational morphological process in agglutinative language like Manipuri?
3. What are rules involved in morphological process?
4. What is the computational model for Morphological analysis in Manipuri Language?
5. What is the best approach for computational Morphological analysis?

## III. OBJECTIVES OF MORPHOLOGICAL ANALYSIS

Meiteilon being agglutinative language is highly inflectional language, which have the capability of generating more then thousands of words from a single root. Hence morphological analysis is vital for high-level applications to understand various words or lexeme in our language. So morphological analysis of Manipuri language forms the foundation for applications related with Natural Language Processing.

Agglutinative languages show high morpheme per word or head word ratio and have complex morphotactics structures, the absence of fusion at affix boundaries make the task of segmentation fluent once the architecture or the model for implementation of morphotactics is build. Here we are listing our main objectives of morphological analysis:

1. Identifying the morphological categories of Manipuri Language.
2. Identifying the number of prefix, Suffix as per category wise.
3. Identifying the morphological categories in Manipuri Language.
4. Identifying the computational process in agglutinative language, Manipuri.
5. Identifying the rules of morphological analysis in Manipuri Language.
6. Identifying the best approached for morphological analysis.

## IV. METHODOLOGY

The process goes in four modules i.e.

1. Data collection module.
2. Classification module.
3. Analysis module.
4. Implementation module.

In data collection process we start making paradigm table in category wise of word classification.

a).Lexical category can be classified in two class i.e. Inflected and Uninflected, later Inflected can be sub categorized in Major and Minor while Uninflected have indeclinables like Postposition, Adverbs, Conjunction, Interjections and Expressives.

Major word classes are:
1) Noun (singular/plural)          = n
2) Verb                                      = v
3) Adjective                              = adj

Minor word classes are:
1) Adverb                                 = adv
2) Pronoun                               = pn
3) Nlocative                             = nst

The possible values for Gender: m, f, n,mf,mn,fn, any .The possible values for Number: sg, pl, dual, any. The possible values for Person:1, 2, 3, any. The possible values for Case:d(direct),o(oblique). The possible case marker: dir, obl, ki, ku, ni, nu, lo, wo, yoVkkaetc. The possible feature structure for the word which is unknown to morph is <fsaf='word,unkn,,,,,,'>.

The possible feature structure for the punctuation mark which is unknown to morph is <fsaf='&sdot;,punc,,,,,,'>. The possible feature structure of the number is <fsaf='88,num,,,,,,'>. The possible values for case name: ex: nom, acc, dubi, etc.**or** 1, 2, 3

This file is read by morph in compiler mode during paradigm-data input.

## V. MORPHOLOGICAL ANALYZER

Morphological Analyzers perform morphological analysis. There are some important approaches for executing morphological analysis. But the two approaches, which are used widely, are:

1. Finite State Machines Based Approach
2. Machine Learning Approach

### A. Finite State Machines Based Approach

Section describes the Finite state machines approach used for building Finite State Transducers(FST) based morphological analyzers.

### B.Resources

The main goal of this approach is to list all the possible parses/analyses of an input word. In order to build a morphological parser using an FST based approach, the following resources are used in general:

### 1. Lexicons

Lexicon of a language is its vocabulary or the list of all words in Manipuri or particular languages. It is an explicit list of every word of the language. It is cumbersome to list every word in a language. Hence generally computational lexicons are used for this purpose. The Finite-state automaton (FSA) is generally used to model lexicons. A structured collection of the entire morpheme i.e. the root or headwords and morphemes or affixes of the words are collected [5].

TABLE I LEXICON TABLE

| Sl. No. | Lexicon | Numbers |
|---------|---------|---------|
| 1 | Headwords | 13463 |
| 2 | Suffixes of Noun | 37 |
| 3 | Suffixes of Pronoun | 39 |
| 4 | Numerals | 32 |
| 5 | Suffixes of Verb | 100 |
| 6 | Suffixes of Adverb | 23 |
| 7 | Suffixes of Adjective | 31 |
| 8 | Prefix | 9 |

### 2.Morphotactics

This explains the morpheme ordering ex: the plural morpheme follows the main noun morpheme. Example: $cars = car + N + pl$.

### 3. Orthographic rules

These rules are also known as spelling rules. They model changes when two morphemes combine. Example: fox + plural s = foxes. Here the *e*-insertion rule is applied. The number of paradigms per PARADIGM-INPUT-FILE is limited to one. More than one paradigm definitions in the

same file will lead to the particular file being rejected by morph for having improper number of word forms [6].

TABLE II NOUN PARADIGM TABLE

| yuM(root) | Category | Suffix |
|---|---|---|
| yuM-n | n.sg+nom(n) | -n |
| yuM-siQ-n | n+pl+nom | siQ-n |
| yuM-du-n | n.sg+Det+nom | -du-n |
| yuM-si-dgi-di | n.sg+Det+Abl+Spec | -si-dgi-di |

TABLE III PRONOUN PARADIGM TABLE

| əikhoi (root) | Category | Suffix |
|---|---|---|
| əikhoi -n | pro-pl+nom | -n |
| əikhoi-si-n | pro-pl+Det+nom | -si-n |
| əikhoi-n-di | pro-pl+nom+spec | -n-di |
| əikhoi-guMb | Pro-pl+Semb | -guMb |

TABLE IV VERB PARADIGM TABLE

| KNn (root) | Category | Suffix |
|---|---|---|
| KNn-ri | verb+dur | -ri |
| KNn -gni | verb+future | -gni |
| KNn -gL-li | verb+ habitual+asp | -gL-li |

TABLE V ADJECTIVE PARADIGM TABLE

| cetpə (root) | Category | Suffix |
|---|---|---|
| cetpə-n | adj+sg+nom | -n |
| cetpə-si-n | adj+det+nom | -si-n |
| cetpə-du-n | adj+Det+nom | -du-n |
| cetpə-si-bu-di | adj+Det+acc+Spec | -si-bu-di |

TABLE VI NUMBER PARADIGM TABLE

| təra (root) | Category | Suffix |
|---|---|---|
| təra-si-n | num+det+nom | -si-n |
| təra-du-n | num+Det+nom | -du-n |
| təra-du-di-n | num+det+spec+nom | -du-di-n |
| təra-du-n-di | num+det+conts+spec | -du-n-di |

TABLE VII ORDINAL PARADIGM TABLE

| ənisubə(root) | Category | Suffix |
|---|---|---|
| ənisubə-si-n | num+Det+nom | -du-n |
| ənisubə-si-n-di | num+det+nom+spec | -si-n-di |
| ənisubə-si-gi | num+dt+gen | -si-gi |
| ənisubə-si-dgi-n | num+det+abl+conts | -si-dgi-n |

The above paradigm tables are used to morphological database. These data are useful as a linguistics source to improve the Manipuri language. Morphological Analysis using Finite State Transducers in Manipuri language is given below as FST for lexicon formation of Manipuri words from one root.



cha, chaba, chabni, chagani......chabagidamak........  (280 words)

Fig.1 Formation of Lexicon from one root in more than 280 (might be more) words.

Flow Chart of Morph will show the overview of the program control from one module to the other. It shows the decision points in the program.
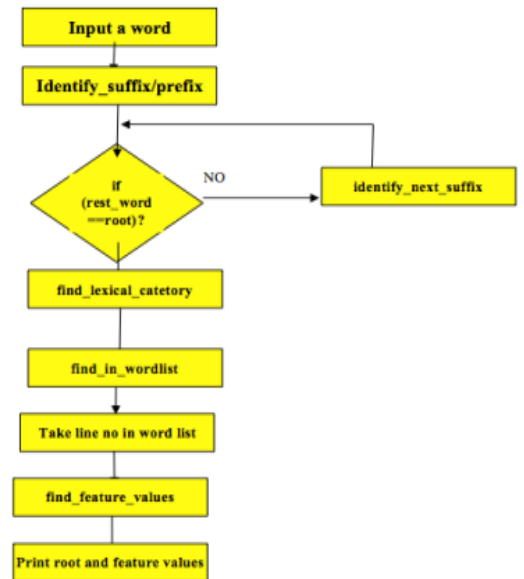


Fig. 2 Flow Chart for MMA

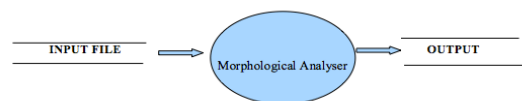DFD diagram showing how data are flow in MMA.



Fig.3 Zero Level DFD
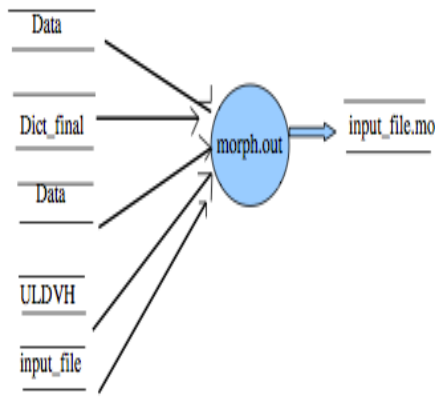
Data flow in another way:
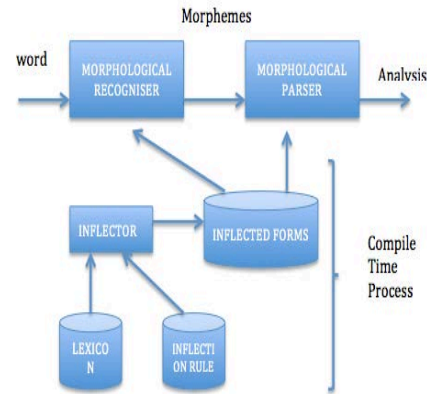
Fig.4 Zero Level DFD



Fig.7 Architecture of MMA

## VI. MORPHOLOGICAL ANALYSIS

In Morphological analysis we take words and we try to identify the suffix or prefix, and check weather this suffix is a valid suffix. If this suffix is present in the word-list then it is a vialed if not it's an unknown word for Morphological Analyzer. If the suffix is valid then check weather the stem is valid or not just, by converting it to root word, by adding and deleting is done. If the suffix and root words are valid then take the line number of the word in word-list then get the feature structure (like gnp, tam, case, case marker value). Add root word and feature structure to API- wrapper to print in the data tree.
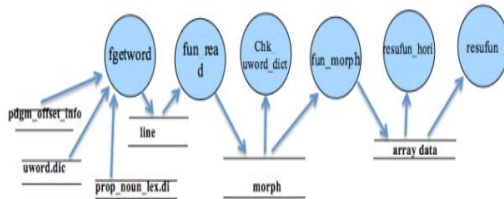


Fig.5 First Level DFD

Here the Figure 7 shown below DFDs is represented as per the programming concept. How data flows was shown from LEVEL 0 to 3.



Fig.6 second Level DFD

```
# This program is use as pre-processing module before tokenizer.

binmode(STDIN, ":utf8");
binmode(STDOUT, ":utf8");
while($line=<>)
{
        utf8::decode($line);
$line =~ s/\x{2018}/'/g; # <2018> ' is Replaced by single quote "'"
$line=~s/\x{2019}/'/g; # <2019> ' is Replaced by single quote "'"
$line=~s/\x{201C}/"/g; # <201C> " is Replaced by single quote "
$line=~s/\x{201D}/"/g; # <201D> " is Replaced by single quote "
$line=~s/\x{200D}//g; # <200D> is Removed
$line=~s/\x{200C}//g; # <200C> is Removed
$line=~s/\x{feff}//g; # <feff> is Removed
$line=~s/\x{0D}//g; #
is Removed
        print $line;
}
```
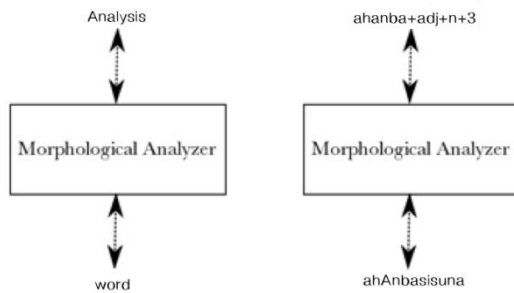
Fig. 8 Process of Morphological Analysis

Manipuri Morphology compiler is a program, which compiles and analyses words belonging to a Manipuri language (Meiteilon). It works in Manipuri language using Meitei Mayek; hence it can learn and recognize words from Manipuri language, which are written in Meitei script. Morph has two modes of operation, via the COMPILER mode and the ANALYSER mode. In the COMPILER mode it reads information about the words of a language from paradigm-input files and lexicon-data files and stores the processed information. This process is referred to as "compilation of data". Once this has been done Morph can recognize the words, which were present in the data. All the data need not be compiled in one go, fresh data can be fed to morph any time by running it in compiler mode. To recognize words one has to run morph in its second mode of operation, which is the ANALYSER mode.

It will recognize only those Manipuri words, which it has been "taught"i.e. trained data. It outputs all the descriptions of the given word it has read during compiler mode. It produces a diagnostic "Unknown word <given word>" to say that it does not find the word in its list. Compilation Manipuri data is done by running morph in compiler mode. While "compiling morph" means putting the morph source code through the C-compiler, which may be necessary especially to customize morph program to some specific need. Compilation of morph is not connected with "data compilation" or "compilation of data" in any manner. At times changes may demand recompilation of full data, while for some changes only recompilation of morph may suffice, and for some changes both may be called for.

*Input- output Specifications:*

> Input: TKN
> Output: Morphological Analysis

*Input specifications*

It requires that property TKN_ must be defined in the input SSF that is given to the Morphological Analysis

| ADD | TKN | CAT |
|-----|------|---------|
| 1 | rAmA | <UNDEF> |
| 2 | sIwA | <UNDEF> |

Input specifications require that property TKN_ must be defined in the input SSF that is given to the Morphological Analyzer.
Output specifications required that Morphological Analyzer as given below would define property of attribute feature. So the output SSF must contain Feature structures to all the valid values.

*Output:*

> An output SSF from Morph must contain all the four columns of SSF.
> The first column will have ADDR
> The second column will have TKN
> The third column will have CAT as UNDEF
> The fourth column will have the feature structure, *fs*
> The feature structure will be in the form of either abbreviated features, of and/or attribute-value pairs. If it has two feature structures then separate them with " **|** " character and between each column there is a tab that separates the fields.

The possible values of fs is listed below, for all possible POS

*A.Nouns*

A noun is analysed as root+suff+{features(such as gender, number,...)}.

The complete structure is presented below.

<fs af root = "Root of the word", lcat = "Lexical category of the root", gend ="Gender of the word", num = "Number coressponding to the word form", pers = "Person of the word", case = "Case ( Direct / Oblique )", vibh = "(cm / tam)" case_name= "case name", spec= "Specificity Marker", emph ="Emphatic Marker", dubi = "Dubitative Marker", interj = "Interjection Marker" conj = "Conjunction Marker" hon ="Honorific Marker" agr_gen ="Gender of the agreeing noun" agr_num ="Number of the agreeing noun" agr_per ="Person of the agreeing noun" suff ="Form of suffix representing all the above markers">

*B.Verbs*

The verb analysis structure is presented below.

<fsaf root = "Root of the word", lcat = "Lexical category of the root", tam ="Suffix indicating Tense Aspect Modality", gend = "Gender of the word", num = "Number corresponding to the word form", pers = "Person of the word", spec = "Specificity Marker", emph = "Emphatic Marker", dubi = "Dubitative Marker", interj = "Interjection Marker", conj = "Conjunction Marker"hon = "Honorific Marker", neg = "verb-neg Marker", voice = "Voice", caus = "Whether the verb form is causative or not(y/n)" finiteness = "Whether the verb form is finite or not (y/n)", suff = "Suff representing all the above markers">

### C.Adjectives

The feature structure for Adjectives is as follows:

<fsaf root = "Root",lcat = "Lexical category",gend = "Gender of the word",num = "Number",pers = "Person of the word", degree= "degree",-like = "like", dubi = "Dubitative",interr = "Interrogative", emph = "Emphatic", conj = "Conjunction Marker", ?spec = "Specific", suff = "complete suffix">

### D.Adverbs

The feature structuer for Adverbs is presented below.
<fs af root= "Root of the word", lcat= "Lexical category of the root", dubi= "Dubitative Marker", interr= "Interrogative", emph= "Emphatic Marker" conj= "Conjunction Marker", ?spec= "Specific", suff= "complete suffix">

### E.Noun Locative

The complete structure is presented below.
<fs af root = "Root of the word", lcat = "Lexical category of the root", gend ="Gender of the word", num = "Number coressponding to the word form", pers = "Person of the word", case = "Case ( Direct / Oblique )", vibh = "(cm / tam)", spec = "Specificity Marker", emph = "Emphatic Marker", dubi = "Dubitative Marker", nterj = "Interjection Marker", conj = "Conjunction Marker", case_name= "case name", hon = "Honorific Marker", agr_gen = "Gender of the agreeing noun", agr_num = "Number of the agreeing noun", agr_per = "Person of the agreeing noun", suff = "Form of suffix representing all the above markers">

### F.Number

A Numeral is analysed as root+suff+{features(such as gender, number,...)}.

The complete structure is presented below.

<fs af root = "numer", lcat ="num" gend ="", num ="", pers ="", case ="", vibh ="", suff ="" >

### G. Punctuation

A Punctuation is not analysed just give the NCR.

The complete structure is presented below.

<fs af root = "NCR" lcat ="punc" gend ="" num ="" pers ="" case ="" vibh ="" suff ="" >

### H.Unknown

A Unknown is not analysed and it is just repeatd.
The complete structure is presented below.
<fs af root = "word", lcat ="unk", gend ="", num ="", pers ="" case ="", vibh ="", suff ="">

## VII. IMPLEMENTATION

The code is separated into parts (subroutines) like command line parsing, program initialization, error handling and logging, and the main application. All the function/subroutines have been defined in the defn.h, struct.h struct1.h files. It would be better if, Functions/subroutines are defined in the separate program files, .cpp file.In the main program they must be called as subroutines
When a new paradigm file is added, or in an existing .p file # of lines is changed:

  a) p file should be placed in pc_data sub dir.
  b) Relevant info regarding the category, features & its values shouldbe entered in Ca, Ce& Fe files in test area.

Program to convert Meitei Mayek to WX mapping.

The given Perl programs will convert Meitei layout to wx mapping.

```
while ($line=<>){

        $line=~s/k([^AeiouO])/ka$1/g;
        $line=~s/s([^AeiouO])/sa$1/g;
        $line=~s/l([^AeiouO])/la$1/g;
        $line=~s/m([^AeiouO])/ma$1/g;
        $line=~s/p([^AeiouO])/pa$1/g;
        $line=~s/n([^AeiouO])/na$1/g;
        $line=~s/c([^AeiouO])/ca$1/g;
        $line=~s/t([^AeiouO])/ta$1/g;
        $line=~s/K([^AeiouO])/Ka$1/g;
        $line=~s/q([^AeiouO])/fa$1/g;
        $line=~s/q/f/g;
        $line=~s/T([^AeiouO])/wa$1/g;
        $line=~s/T/w/g;
        $line=~s/w([^AeiouO])/va$1/g;
        $line=~s/w/v/g;
        $line=~s/y([^AeiouO])/ya$1/g;
        $line=~s/h([^AeiouO])/ha$1/g;
        $line=~s/P([^AeiouO])/Pa$1/g;
        #$line=~s/a([^AeiouO])/aa$1/g;
        $line=~s/g([^AeiouO])/ga$1/g;
        $line=~s/J([^AeiouO])/Ja$1/g;
        $line=~s/r([^AeiouO])/ra$1/g;
        $line=~s/b([^AeiouO])/ba$1/g;
        $line=~s/j([^AeiouO])/ja$1/g;
        $line=~s/G([^AeiouO])/Ga$1/g;
        $line=~s/D([^AeiouO])/Xa$1/g;
        $line=~s/D/X/g;
        $line=~s/B([^AeiouO])/Ba$1/g;
        $line=~s/o/oV/g;
        $line=~s/i/i/g;
        $line=~s/I/i/g;
        $line=~s/A/A/g;
        $line=~s/e/eV/g;
        $line=~s/O/O/g;
```

```
$line=~s/u/u/g; #secondary u
$line=~s/U/u/g; #primary u
$line=~s/Y/E/g;
$line=~s/z/M/g;
$line=~s/x/tV/g;
$line=~s/d([^AeiouO])/xa$1/g;
$line=~s/d/x/g;
$line=~s/Q/fV/g;
$line=~s/C/kV/g;
$line=~s/L/lV/g;
$line=~s/M/mV/g;
$line=~s/F/pV/g;
$line=~s/N/nV/g;
print $line;
}
```

## VIII. CONCLUSION AND FUTURE WORK

In the present work, the development of a Manipuri Morphological Analysis has been described. The root dictionary stores the related information of the corresponding roots. The Analyzer can classify the word classes and sentence types based on the affix information. The verbs are under bound category. The verb morphology is more complex than those others. The distinction between the noun class and verb classes is relatively clear; the distinction between nouns and adjectives is often vague. Thus, the assumption made for word categories depend upon the root category and affix information. In the stripping of the morphemes the various morphemes pattern combinations are tested.

The Natural Language Processing tools need more text corpus with better transfer rules and techniques to achieve quality output. The performance of the various Manipuri NLP tools that have been developed in the present work need to be improved by experimenting with various machine-learning approaches with more training data.

Future works include the developments of automatic Morphological Analyzer using some machine learning algorithms. The exploration and identification of additional linguistics factors that can be incorporated into the Morphological Analysis to improve the performance is an important future task.

### Input

(Manipuri)

| ADDR_ | TKN_ | CAT_ |
|-------|------|------|
| 1 | BoroisIn | <NOUN> |
| 2 | pahabasisuna | <ADJECTIVE> |
| 3 | somnA | <PRONOUN> |

### Output

ADDR_   TKN_   CAT_           others_

| 1 | Boroi | <fsaf='Boroi,n,n,sg,,d,,'>|<fsaf='Boroi,n,m,pl,,d,,'> |
| 2 | pahabasisuna | <fsaf='pahaba,adj,n,sg,3,,su,su' spec='si_na'> |
| 3 | somnA | <fsaf='som,pron,n,sg,2,,na,su' spec= 'nA'> |

## REFERENCES

[1] Grieson, *Linguistics Survey of India,*Vol.III part III,1973.
[2] Singh Ch.Yashawanta , *Manipuri Grammar*, Rajesh Publications", New Delhi, 2000.
[3] A.Sarangi, *Language and Politics in India,* Oxford University Press, Ne Delhi, pp. 27, 2009.
[4] K.Nikhil, I.Abhilash and D.M.Sharma, "Hindi Derivational Morphological Analyzer. In *Proceedings of SIGMORPHON*, 2012.
[5] Yumnam Bablu Singh, *Web enabled Multilingual Manipuri Dictionary*,
[6] Aksher Bharathi and Rajeev Sangal, *et.al,* "*Natural Language Processing: A Paninian Perspective*".