# Corpus Management and Multi-lingual Lexical Database Tools

CDAC, Thiruvananthapuram

## Abstract

In this paper we will discuss about the tools for corpus management and Multi-lingual lexical database which were used for English to Malayalam translation system called AnglaMalayalam. AnglaMalayalam is an adaptation of AnglaBharti technology which was developed by IIT Kanpur.

## 1. Introduction

Corpus is a body of language. Corpus is an inevitable resource in Natural Language Computing. It is widely used to study the behaviour of Natural Languages, develop statistical language computing rules and building lexical resources. In the field of Machine Translation, a good quality corpus is very much essential for good quality translation. Thus corpus collection and its analysis form an essential task. In the following section we will discuss about the tools used for lexicon updation and domain specific lexicon entry for AnglaMalayalam system.

## 2. Corpus Collection and Analysis

The corpus is acquired from web resources, Journals, Magazines, text books, etc. The collected corpus had to be processed for analysis. The process includes cleaning, unique word Identification, sentence boundary marking, parsing, extracting word and POS, extracting bi-gram and tri-gram and extraction of NP and VP.

The corpus which is collected from web may contain certain junk characters, extra spaces after the words or sentences. Before analysing the corpus, such characters and spaces should be removed properly. The abbreviations and acronyms were extracted from the cleaned corpus based on certain heuristics. Accurate sentence splitting is an important building block of many NLP systems. Most Part of Speech taggers require input in the form of one sentence per line. We used a sentence tokenization tool developed in perl with a standard CPAN Module Lingua::EN::Sentence.

### 2.1 Unique Word Extraction

From the raw corpus, word and frequency count will be taken. Then the word is compared with the main lexicon. The compared words are subjected to the stemming process. Stemming is mainly done to remove the plural markers. After stemming process the words are again compared with the main lexicon. If the word is not found in the lexicon then it will be stored in to a separate file. Manual verification is required after the process to eliminate the nonsensical words.
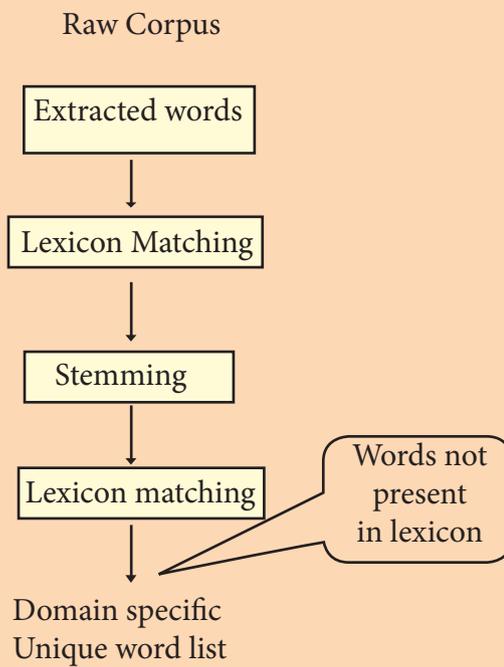
Raw Corpus

Extracted words

↓

Lexicon Matching

↓

Stemming

↓

Lexicon matching → Words not present in lexicon

↓

Domain specific
Unique word list

**Figure 1: Flow Chart for Unique Word Extraction**

DATA ANALYSIS

WordCount

Input file /folder: bestos_N_0807_org.txt   Browse

Output folder: Settings\nisha\Desktop   Browse

Find WordCount

1 files processed ;Total WordCount = 143573

LexComp

Input file: isha\Desktop\output.txt   Browse

Output file: a\Desktop\LexComp.txt   Browse

Compare Lexicon

3077 Words After Comparing with Lexicon

Stem Plurals & Compare with Lexicon

Input file: gs\nisha\Desktop\LexComp.txt   Browse

Output file: s\nisha\My Documents\Final.txt   Browse

Stem

Unique Words = 1825

**Figure 2: Screenshot for Unique Word Extraction Tool**

## 2.2 Multiword Extraction Tools

### 2.2.1 Sentence Boundary Detection

The Sentence Boundary Detection tool is developed in perl with a standard CPAN Module Lingua::EN::Sentence. The module is an implementation of standard sentence boundary detection algorithm.
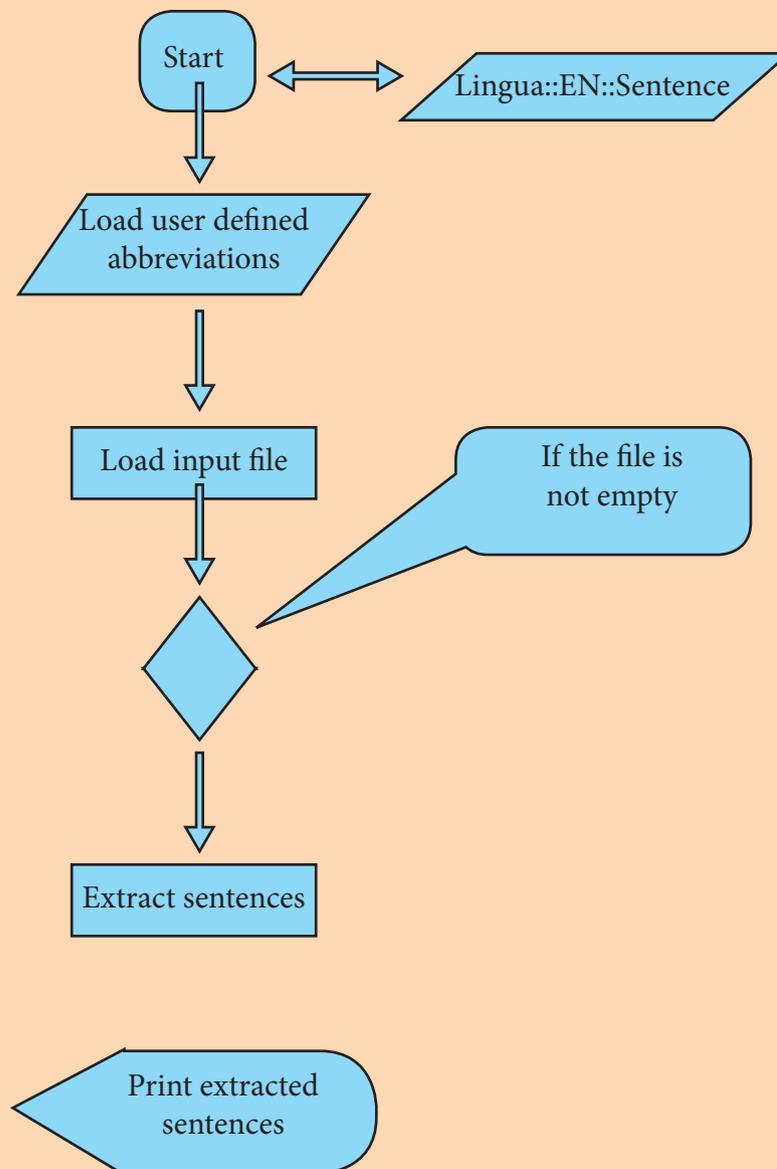


**Figure 3: Flow Chart for Sentence Boundary Detection**

## 2.2.2 Named Entity Extraction

The Named Entity extraction tool developed in Perl with a Standard CPAN (Comprehensive Perl Archive) Module called Lingua::EN::NamedEntity. It is a perl implementation of basic Named Entity Recognition Algorithm. The implementations can be installed in the system to use the algorithm in perl programs.
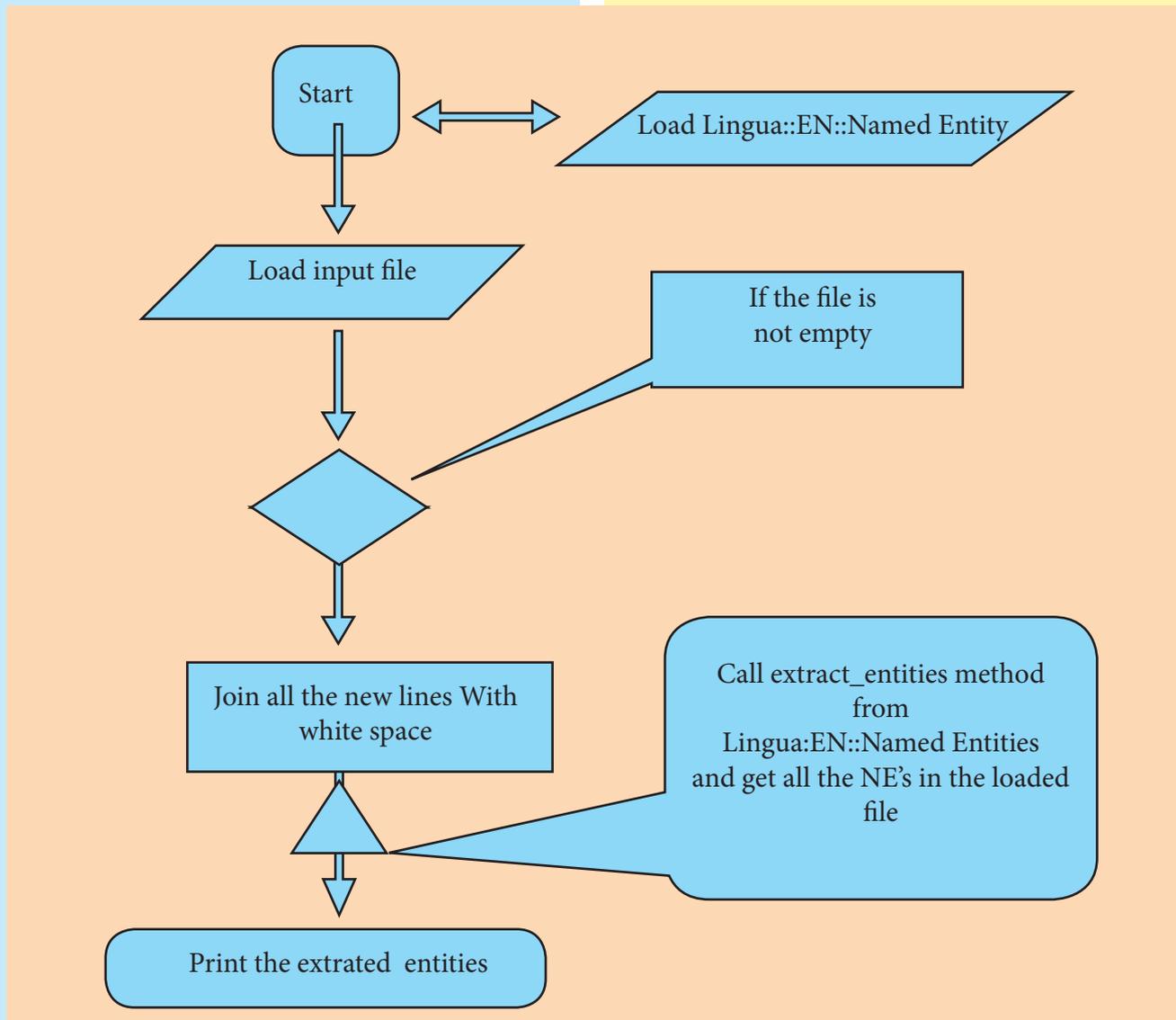
```
Start  <-->  Load Lingua::EN::Named Entity
  |
Load input file          If the file is
  |                        not empty
  <diamond>
  |
Join all the new lines With       Call extract_entities method
white space                        from
  |                           Lingua:EN::Named Entities
  |                          and get all the NE's in the loaded
Print the extrated entities        file
```

**Figure 4: Flowchart for Named Entity Recognition**

### 2.2.3 Parser

The tool is developed with Charniak parser and Perl programming language. Charniak parser is a statistical parser for English language. The parser can be downloaded from ftp://ftp. cs.brown.edu/pub/nlparser/ under the GNU GP license. The Perl program takes absolute path to sentence boundary marked corpus. It will take the files one by one. Processes and stores the output to the directory parsed.
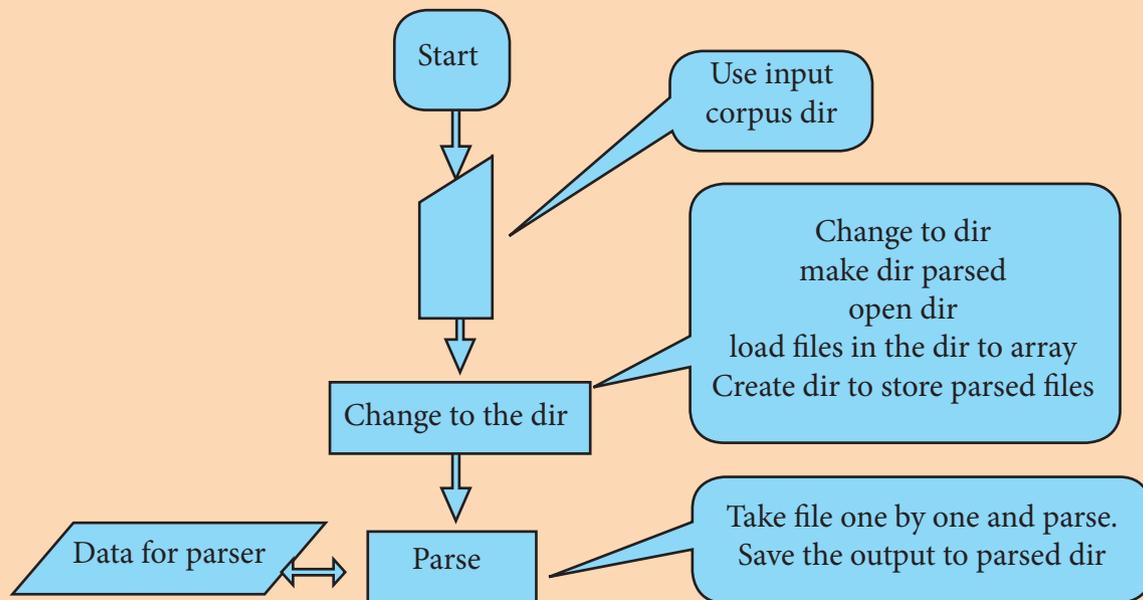
Start

Use input corpus dir

Change to dir
make dir parsed
open dir
load files in the dir to array
Create dir to store parsed files

Change to the dir

Data for parser

Parse

Take file one by one and parse.
Save the output to parsed dir

**Figure 5: Flowchart for the Parser**

## 2.2.4 MWE extraction

The program is developed in perl. It takes a java bigram and Trigram tool; which extracts bigrams and Trigrams along with POS from the parsed corpus. The tool is divided in to four different perl files. It recieves the absolute path to parsed corpus file and extracts NP word and POS, word and POS, bigram and trigaram and MWES.
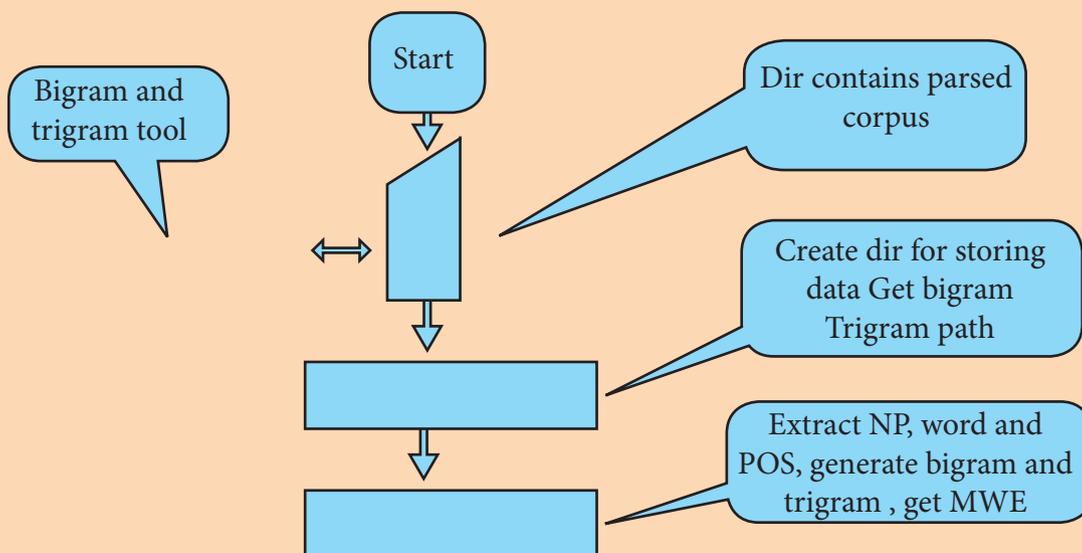
Bigram and trigram tool

Start

Dir contains parsed corpus

Create dir for storing data Get bigram Trigram path

Extract NP, word and POS, generate bigram and trigram , get MWE

**Figure 6: Flowchart for Multiword Extraction**

## 3. Lexicon Updation Tool

A lexicon is usually structured as a collection of English words along with their syntactic and semantic information. The syntactic information includes the POS, Tense, aspect,

Here the user can enter the meaning by using the online keyboard provided along with the tool or he can use the inscript keyboard. The paradigm number can be assigned by using the different combinations of inflections for a particular
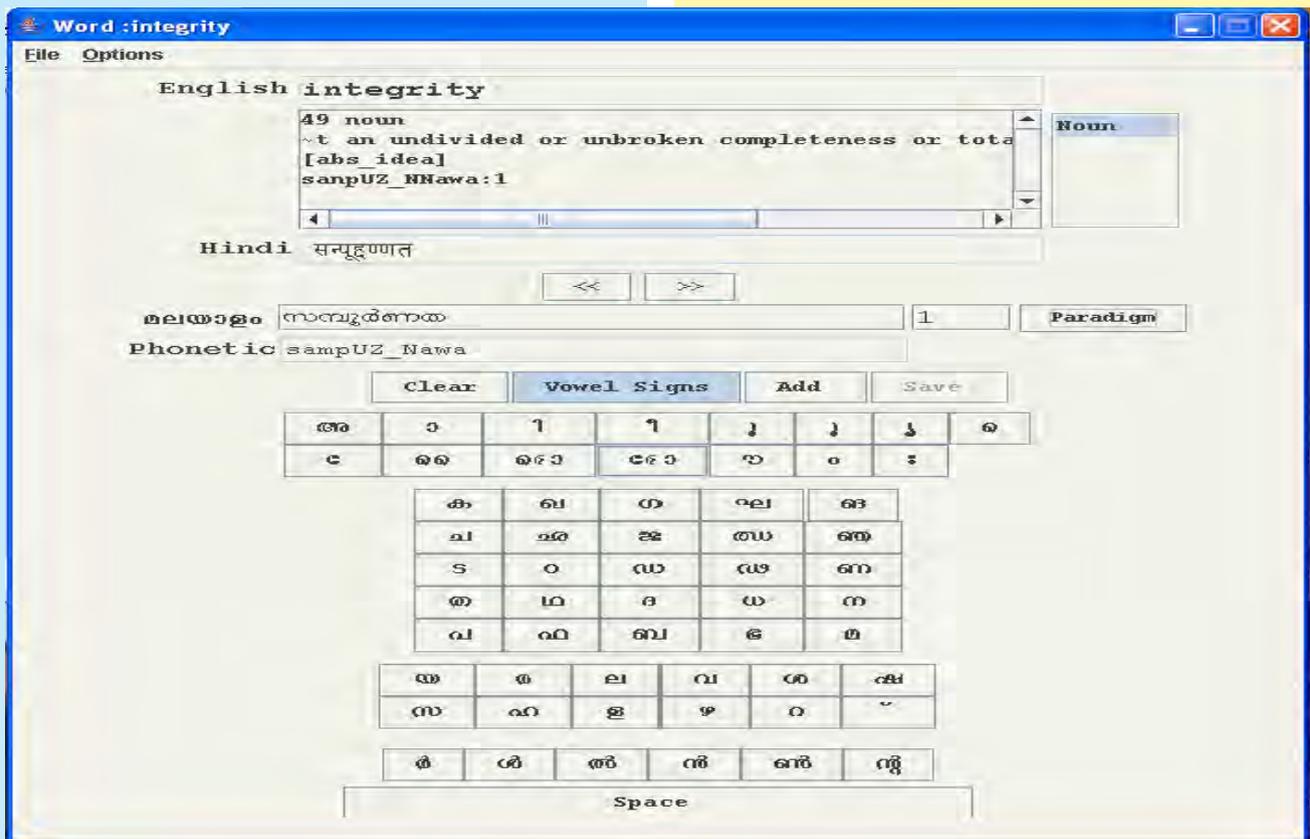


**Figure 7: Lexicon Updation Tool**

Modality, etc. Lexicon is one of the major parts in English to Malayalam translation process. We have used English – TL lexicon as a reference to develop English – Malayalam lexicon. Since the English – TL lexicon contains a round 45000 words. Manual entry will affect the dictionary structure. So a tool is developed to replace the

category of word provided by the system. We can choose the apt paradigm number through manual selection. The pop up coming upon clicking the Paradigm button is given below.
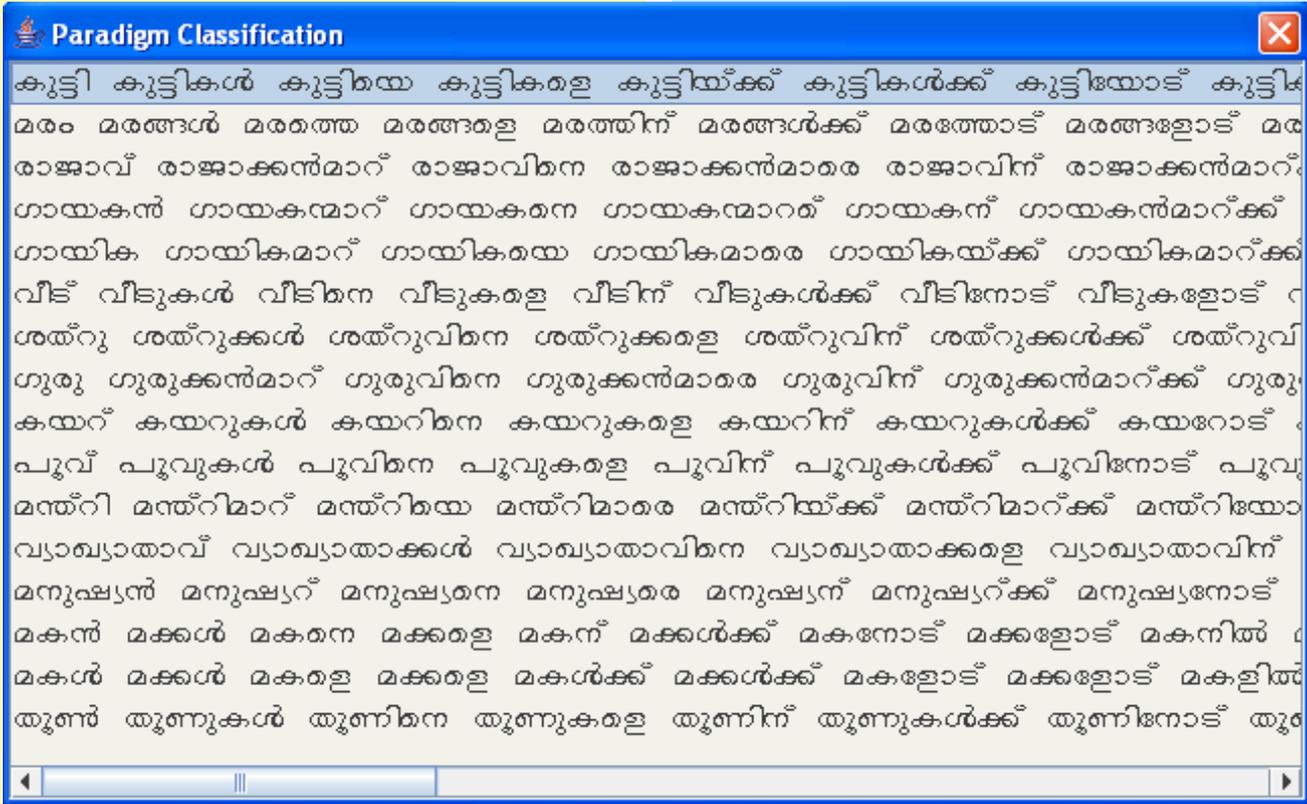
**Figure 8: Inflection for nouns**

## 4. Domain specific Lexical Entry

Lot of new words are extracted from the corpus and extracted domain specific words. Domain specific lexical entry tool was used to create domain specific lexicon. This tool can also be used to update the existing lexicon. The lexicon is having a specific format, which contains syntactic and semantic information. The tool is an XML based tool which has paradigm generators for Malayalam and English words. Since it is a domain specific lexical entry tool, there is a provision for selecting the domain like General, Health or Tourism. The other options are for entering the root word and POS tag of the word. Also, there were provisions for adding, editing, deleting and searching the desired root word and its information's.

In the Noun entry field, we have provided options for entering the details such as gender (masculine/feminine/neuter/don't care), number (singular/plural), person (first/second/third) and also the Target language details. In the field of Target language details, there were provisions for entering the meaning, semantic tag and paradigm number. For Adjective entries, there were provisions for entering the degrees of comparison (positive/comparative/superlative) and Target language details as mentioned above. For Verb entries, options were there for entering the tense forms (present/past/past participle/continuous), transitive, intransitive and bitransitive details. In the case of transitive verbs, enter the meaning, subject and object and paradigm number. In the case of intransitive verbs, enter the

meaning, subject and paradigm number and for bitransitive verbs, enter the meaning, subject, first object, second object and paradigm number. From the below figure (9), we can see that the English paradigm "16" has been generated based on the tenses, which is strictly based on certain

## 5. Summary

In this paper we have discussed about the tools for corpus management and multi-lingual lexicon database. We have seen the various tools used for corpus cleaning; identify unique words, sentence boundary marking, parsing, extracting
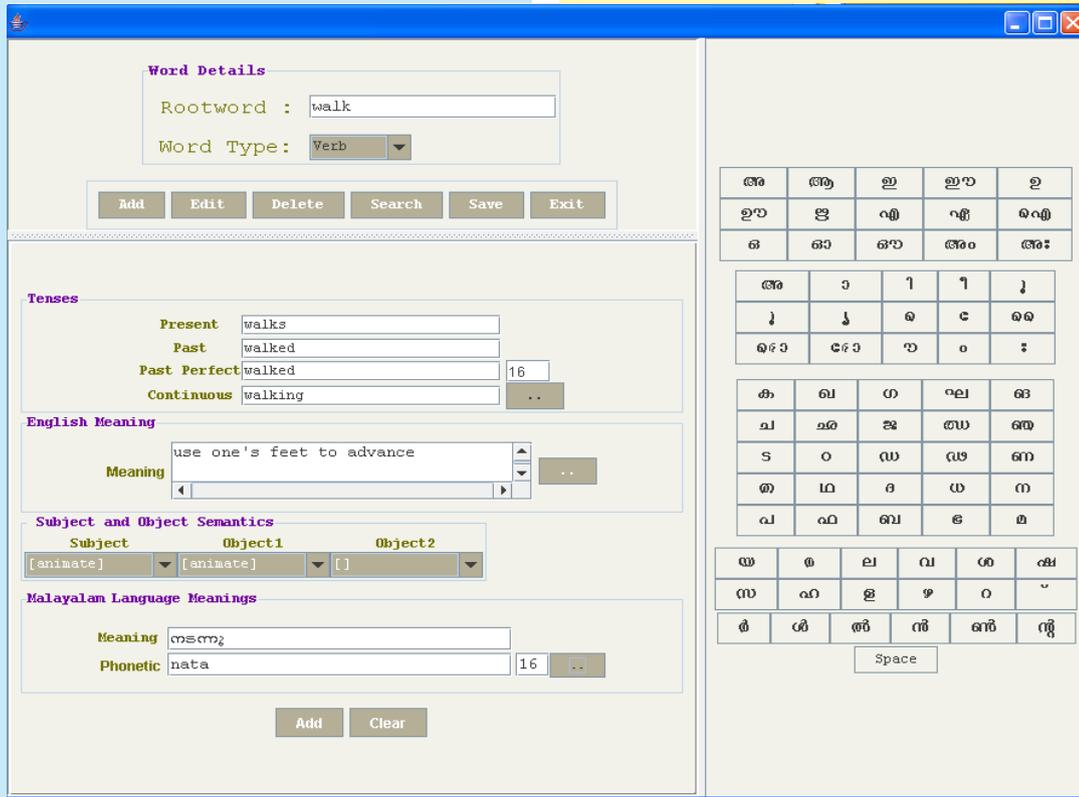


**Figure 9: Lexicon Entry Tool**

heuristics used in AnglaBharti technology. We have displayed available English meanings of the root word. The user can select the correct one which is similar in meaning to the root word. For this purpose, we have incorporated a dictionary of Word Net for getting the matching meanings corresponding to the words. The tool has the provisions to enter the Malayalam meaning using the online keyboard. The apt paradigm number for Malayalam can be assigned by using the different combinations of inflections for a particular category of word provided by the System.

word and POS, extracting bi-gram and tri-gram. We have also discussed about the tool for multi-lingual lexical database creation.

## Abbreviations / Acronyms

**NP:** Noun Phrase, **VP:** Verb Phrase, **PP:** Preposition Phrase, **POS:** Part of Speech, **XML:** Extended Mark-up Language, **CPAN:** Comprehensive perl archive, **NN:** Noun singular, **NNS:** Noun plural, **NNP:** Proper Noun Singular