

# **Project Documentation**

## **Morphological Analyzer for Malayalam**

Jincy Baby

International Centre for Free and Open Source Software (ICFOSS)

## Project Overview

Morphological analyser is a program which compiles and analyses words belonging to a natural language. The pre-processing stage before the analysis, which splits the input word into morphemes, is the major part in an analyser. The current analyser developed uses a rule-dictionary based approach, where sandhi rules and pattern matching methods are considered for the splitting of words. This part is an extension of root extractor where the output is a list containing the root of the word and its respective suffixes that are stripped from the word. The suffixes to be stripped and tags to be assigned was fixed as per the discussions with a language expert.

## Introduction

Morphological analysis for Indian Languages is defined as the analysis of a word in terms of its lemma, gender, number, person, case, vibhakti, tense, aspect and modality. A major way in which morphologists investigate words, their internal structure, and how they are formed is through the identification and study of morphemes, often defined as the smallest linguistic pieces with a grammatical function. Consider the word "കാണുകയാണെന്ന്". We can break it into three morphemes: "കാണുക", "ആണ്", and "എന്ന്". "കാണുക" is called the stem, which is a base morpheme to which another morphological piece is attached.

## System Description

### Module 1 - Morpheme Generator

The module generates a list with the root and stripped suffixes of the given input word. The projects which needs both the root and suffixes can utilize this module rather than root extractor. The code is slightly different from root extractor but the concept is similar. The implementation method for morpheme is specified towards the end of the document.

The functions, dictionaries and suffix lists introduced in the morph analyser are almost similar to the root extractor. The only difference exhibited is the that an addition module is included for tagging the stem and suffixes after splitting the input into its appropriate stem and suffix list. The module make use of a tag database included specifically for tagging the suffixes. The stem is tagged into noun and verb based on the suffix list obtained.

### *stem*

Let  $x_1x_2...x_lX$  be the given input. The patterns stored in the database is stored as a list initially. The function checks whether the input ends with any pattern in the list and if it find such a pattern, it obtains the corresponding replacement string stored in the database. Assume that  $X$  is the pattern found and  $y$  is the corresponding replacement string to be attached to the word after removing the pattern. The function strips away  $X$  and replace it by  $y$  thus forming  $x_1x_2...x_ly$ . An additional column is available in the database which stores the suffix corresponding to  $X$  to be inserted in the output list.  $x_1x_2...x_ly$  is again returned to the function as input and the process is repeated until no such pattern is found. Finally

the last word is inserted into the list with suffixes before returning it as the output.

Algorithm(*stem*)

input: word, slist=list of patterns

output: list of word with no patterns left and suffixes

Steps:

1. Take the input word
2. Check the list of stopwords and return the word if found any.
3. For pattern in slist do
4. If word ends with pattern do
5. Replace the pattern with the corresponding replacement string from the database
6. Store the suffix in the output list
6. Re-initialize the word and go to step 1
7. Continue the process until no such pattern is found

The following shows a simple example which displays the detailed process of the function *stem*

Example1

Input Word: മാധ്യമപ്രവർത്തകൻ

Process:

Input: മാധ്യമപ്രവർത്തകൻ

pattern found : "ൻ"

Corresponding Replacement String: " " "

word obtained: മാധ്യമപ്രവർത്തകർ

list obtained: ["ൻ"]

Re-initialized Input: മാധ്യമപ്രവർത്തകർ

pattern found : "ർ"

Corresponding Replacement String: "ർ"

word obtained: മാധ്യമപ്രവർത്തകർ

list obtained: [" " ]

Output list: [മാധ്യമപ്രവർത്തകർ , ള്]

## ***verb***

A similar method is used in the function *verb*, with a major difference that the function is not recursive and it uses a verb dictionary for verification of the output. It converts an input word into verb form considering the pattern introduced in a database for the function and return the list of verb and suffixes.

Algorithm(*verb*)

Input: word

Output: word after replacing the pattern

Steps:

1. Take the input word

2. For pattern in list do
  4. If word ends with pattern do
    5. Replace the pattern with the corresponding replacement string from the database
    6. Insert the corresponding suffix into output list
    6. If the word obtained is in the verb dictionary, return the output list

A simple example to show the process of the function *verb* is discussed below.

### Example2

Input Word: പകർന്നു

Process:

Input: പകർന്നു

pattern found: "ർന്നു"

Corresponding Replacement String: "രക"

Corresponding suffix: "ന്നു"

word obtained: പകരക

Output list: [പകരക, ന്നു]

### *change*

*change* is a function which supports other functions mainly to form the suffixes. The function takes in the maatra and converts it into its appropriate Malayalam letter. It converts “ഓ” to “ആ”, “ഐ” to “ഇ”, “ഐ” to “ഇ”, “ഐ” to “ഉ”, “ഐ” to “ഉ”, “ഐ” to “എ”, “ഐ” to “എ”, “ഐ” to “ഒ” & “ഐ” to “ഓ”.

### *dvithva*

Another function *dvithva* checks the input word for any string in the list [“ക”, “പ”, “ത”, “ശ”, “ച”]. It takes ‘dvithva sandhi’, which occurs when the first part in a morpheme doubles when joined to another word, into consideration. Let  $x_1x_2\dots x_tx_{t+1}\dots x_n$  be the input where  $x_t = x_{t+1}$ . Then the function splits the word into  $x_1x_2\dots x_{t-1}$  and  $x_{t+1}\dots x_n$  by removing  $x_t$ . If it finds  $x_{t+1}\dots x_n$  in suffix list for dvithva, it returns  $[x_1x_2\dots x_{t-1}, x_{t+1}\dots x_n]$ .

Algorithm(*dvithva*)

Input: word, checklist=[“ക”, “പ”, “ത”, “ശ”, “ച”]

Output: list of splitted morphemes

Steps:

1. Take the input word
2. If any string in checklist is in word do
  3. for string in checklist do
    4. if word endswith the string, continue
    5. else do
      6. find index of the string as N
      7. if word[N+2:] in suffix list for dvithva do

4

13. Re-initialize the word and go to step 1
14. else do
15. if word in dictionary, insert it into output list
16. else take *verb(word)* and if non-empty, insert the obtained word and suffixes into output list
17. return the output list

A simple example is shown below to display the detailed process of the function *sp*

#### Example4

Input Word: ശേഖരിക്കുന്നതെന്ന

Process:

Input: ശേഖരിക്കുന്നതെന്ന

pattern found: "െ"

output list: [എന്ന]

word obtained after removing suffix: ശേഖരിക്കുന്നത്

list obtained using *stem*: [ശേഖരിക്കുന്ന,ത്]

output list: [ത്, എന്ന]

Reinitialize word: ശേഖരിക്കുന്ന

Input: ശേഖരിക്കുന്ന

pattern found: "ു"

output list: [ഉന്ന, ത്, എന്ന]

word obtained after removing suffix: ശേഖരിക്കു്

word obtained using *verb*: ശേഖരിക്കുക

Output list: [ശേഖരിക്കുക, ഉന്ന, ത്, എന്ന]

### *fn*

Another function *fn* checks the input word for any string in the list [“ഊ”, “ാ”, “ി”, “ു”, “ൂ”, “െ”, “േ”, “ൊ”, “ോ”]. The word upto the particular string is cut off and thus split into two. Let  $x_1x_2\dots, x_tx_{t+1}\dots x_n$  be the input word where  $x_t$  is the string found in the list. The word is then split into  $x_1x_2\dots x_t$  and  $x_{t+1}\dots x_n$ . If  $x_{t+1}\dots x_n$  is found in the suffix list or the dictionary, the system examines if  $x_1x_2\dots x_t$  is found in the dictionary else it passes  $x_1x_2\dots x_t$  through the function *sp*. After validation of both,  $[x_1x_2\dots x_t, x_{t+1}\dots x_n]$  or the output from *sp* after appending  $x_{t+1}\dots x_n$  is returned as the output.

Algorithm(*fn*)

Input: word, mlist=[“ഊ”, “ാ”, “ി”, “ു”, “ൂ”, “െ”, “േ”, “ൊ”, “ോ”]

Output: word stem

Steps:

1. Take the input word
2. If word is in dictionary, then return the word
3. For k in mlist do
4. if k is in word do
5. find all indices of k and store in a list

6. For N in list of indices do
7.   verbsuf=*verb*( $x_{N+1}...x_n$ )
8.   if  $x_{N+1}...x_n$  in suffix list or dictionary do
9.     if  $x_1x_2...,x_{N+1}$  in dictionary, return  $[x_1x_2...,x_{N+1},x_{N+1}...x_n]$
10.    else if  $x_1x_2...,x_N$  in dictionary, return  $[x_1x_2...,x_N,x_{N+1}...x_n]$
11.    else if  $sp(x_1x_2...,x_N + “\circ”)$  not empty, return  $sp(x_1x_2...,x_N + “\circ”) + [x_{N+1}...x_n]$
12.   if verbsuf[0] in suffix list or dictionary do
13.     if  $x_1x_2...,x_{N+1}$  in dictionary, return  $[x_1x_2...,x_{N+1}] + verbsuf$
14.     else if  $x_1x_2...,x_N$  in dictionary, return  $[x_1x_2...,x_{N+1}] + verbsuf$
15.     else if  $sp(x_1x_2...,x_N + “\circ”)$  not empty, return  $sp(x_1x_2...,x_N + “\circ”) + verbsuf$
16. if  $x_1x_2...,x_tx_{t+1}...x_n$  ends with k do
17.   if  $x_1x_2...,x_tx_{t+1}...x_{n-1}$  in dictionary, return  $[x_1x_2...,x_tx_{t+1}...x_{n-1}], change(k)$
18.   else if  $sp(x_1x_2...,x_tx_{t+1}...x_{n-1} + “\circ”)$  is not empty do
19.    return  $sp(x_1x_2...,x_tx_{t+1}...x_{n-1} + “\circ”) + [change(k)]$

The following shows a simple example which displays the detailed process of the function *fn*

#### Example5

Input Word: പകർന്നത

Process:

Input: പകർന്നത

pattern found: "ു"

list: [പകർന്ന , ത]

conversion of suffix using *verb*: [തക]

തക in dictionary

conversion of പകർന്ന using *verb*: [പകരക,ന്ന]

പകരക in dictionary

Output list: [പകരക,ന്ന,തക]

#### *ni*

*ni* searches the input for any string in the list [”ത”, ”ര”, ”ന”, ”ശ”, ”ഓ”]. If found, it splits the word upto the string and checks if the other half is in the suffix list or dictionary. The word derived is passed through *sp* and concatenated together to obtain the output. Let  $x_1x_2...x_{t-1}x_tx_{t+1}...x_n$  be the input. Let  $x_t$  be the string in the list.  $x_{t+1}...x_n$  is searched in the suffix list or dictionary and if found,  $sp(x_1x_2...x_t)$  is obtained and  $sp(x_1x_2...x_t) + [x_{t+1}...x_n]$  is returned as output.

#### Algorithm(*ni*)

Input: word= $x_1x_2...x_{t-1}x_tx_{t+1}...x_n$ , chlist=[“ത”, “ര”, “ന”, “ശ”, “ഓ”]

Output: output list

Steps:

1. Take the input word

2. for pattern in chlist:
3. if  $x_{t+1}...x_n$  in suffixlist or dictionary do
4. return  $sp(x_1x_2...x_t)+[x_{t+1}...x_n]$

The below example depicts the working of the function *ni*

#### Example6

Input Word: മഴയിൽനിന്ന്

Process:

Input: മഴയിൽനിന്ന്

pattern found: "ൽ"

suffix: നിന്ന്

word obtained after removing suffix: മഴയിൽ

word obtained using *sp*: [മഴ, ഇൽ]

Output list: [മഴ, ഇൽ, നിന്ന്]

### ***finallist***

The function obtains two lists. It removes the first element in the second list and add it with the first list to obtain the outputs to be returned.

### ***root***

The root module uses all the suffix stripping module to create a final list to be provided to the morphological analyser.

Algorithm

Steps:

1. Take the input word
2. Remove special characters and re-initialize the word
2. Stemlist = *stem*(word)
3. Store the first element in the stemlist in a list split.
4. Derive *ni*(split[0]) and if its not [split[0]], split=*ni*(split[0])
5. Derive *dvithva*(split[0]), split is set as *finallist*(*dvithva*(split[0]),split))
6. Derive *sp*(split[0]), split is set as *finallist*(*sp*(split[0]),split))
7. Derive *fn*(word), split is set as *finallist*(*fn*(split[0]),split))
8. Derive *verb*(word), split is set as *finallist*(*verb*(split[0]),split))
9. return *finallist*(split,stemlist))

### **Module2**



## NV

The function assigns appropriate tags to the root word taking into account the list of suffixes obtained from the finallist. The function considers two dictionaries verbsuffix which contains suffixes which can be attached with a verb and a noun dictionary noundict which contains suffixes that can be attached with a noun.

Algorithm(NV)

Input: root, clist=list of suffixes obtained from the finallist

Output: [morpheme,tag]

Steps:

1. if root in verb dictionary, return [morpheme, Verb]
2. if clist==[] or clist[0] in verbsuffix dictionary do
3. replace “ ~ ” by “ ള് ” and if the obtained word in dictionary, return [word,Verb]
4. if clist==[] or clist[0] in noundict do
5. return [word,Noun]

## tag

tag assigns appropriate tags to the input suffix. tags for the suffixes are obtained from the database introduced for tags. Each suffix is searched in the database and upon finding return the list with the suffix and tag.

## Morphological analyser

The main body of the morph analyser uses the root function to derive the list of stem and suffixes to be analysed further. The obtained stem is first searched to identify it as Noun or Verb utilizing the function NV. If the stem is verb, then the suffix list is checked in a database if the combination of suffixes are present . If not the suffixes are tagged one by one using the tag function. Algorithm(*Morph Analyser*)

Input: Word

Output: Analysis result

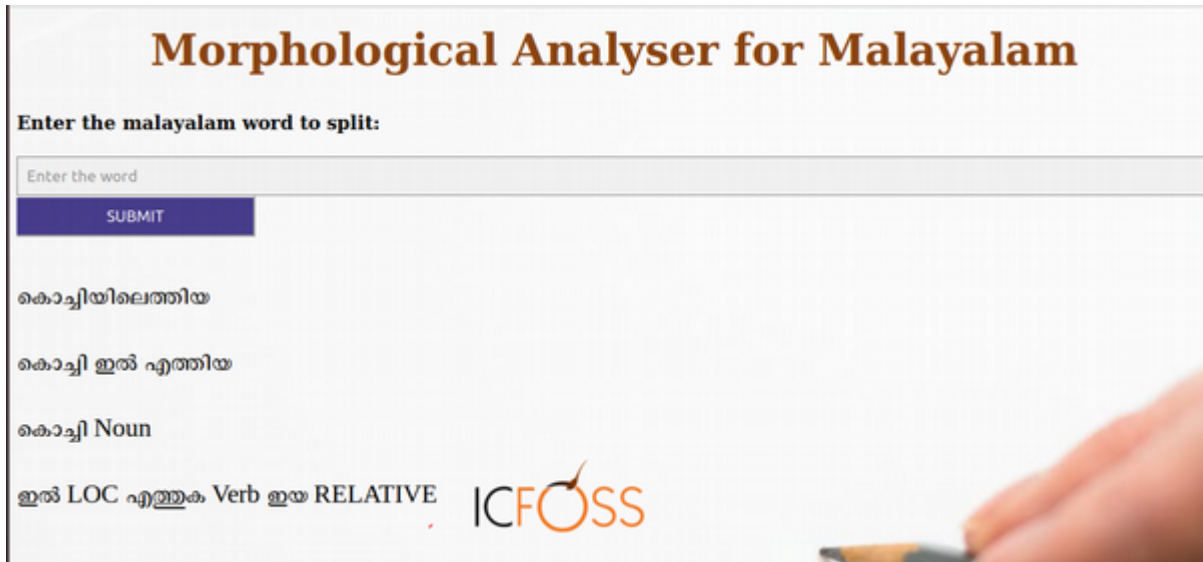
Steps:

1. Derive list of morphemes using *root*
2. Use NV to determine if the stem word is Noun or Verb
3. If verb do
4. if the list of suffixes is in the database, derive the tag of suffixes and give the output.
5. break
6. else, use tag function and return the result

## Results

Series of experiments has been conducted to improve and test the performance of the proposed algorithm. The tests were carried out with testing corpus from various sources. The extractor showed an increase in accuracy with updation of dictionary. The tests produced an average 92% accuracy for the analyser.

A total of 7000 words were tested out of which 6453 results produced accurate results. The updation of dictionary and tag list seems to increase the overall accuracy rate of the system.



**Morphological Analyser for Malayalam**

Enter the malayalam word to split:

Enter the word

SUBMIT

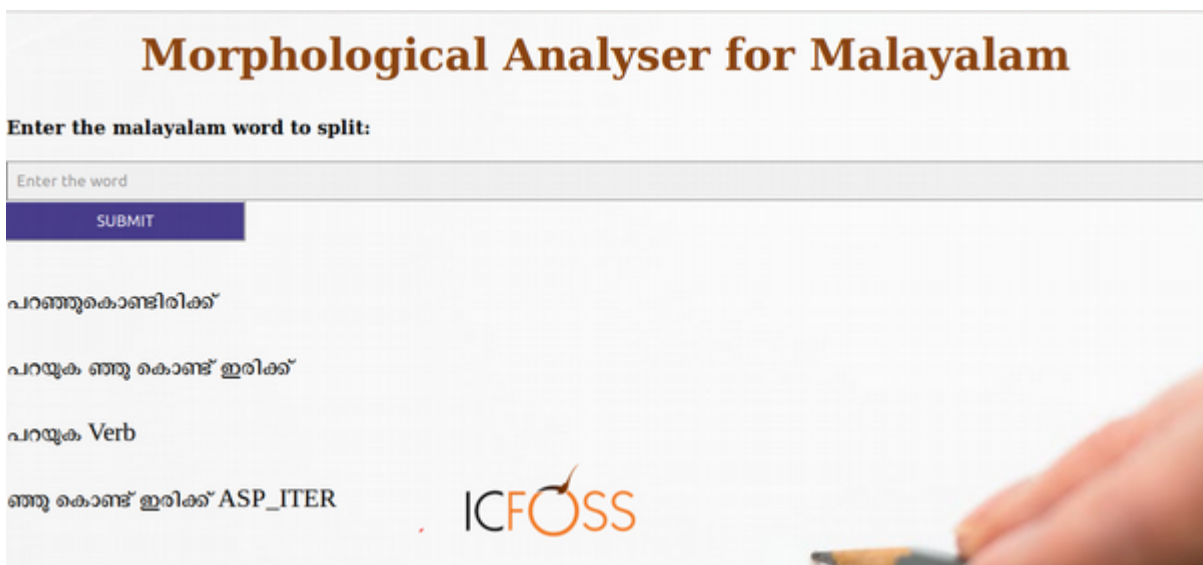
കൊച്ചിയിലെത്തിയ

കൊച്ചി ഇൽ എത്തിയ

കൊച്ചി Noun

ഇൽ LOC എത്തുക Verb ഇയ RELATIVE

ICFOSS



**Morphological Analyser for Malayalam**

Enter the malayalam word to split:

Enter the word

SUBMIT

പറഞ്ഞുകൊണ്ടിരിക്കു

പറയുക ഞ്ഞു കൊണ്ട് ഇരിക്കു

പറയുക Verb

ഞ്ഞു കൊണ്ട് ഇരിക്കു ASP\_ITER

ICFOSS

## Implementation of Morphological Generator

The extractor is in the form of a function and can be directly used in any NLP projects just by importing the module. The sentence obtained should be tokenized and the compound words obtained when treated with the morph() function will produce a list of morphemes for the input word.

```
In [11]: 1 word="നിയന്ത്രണത്തിലാക്കാനു്"  
2 FINAL=morph(word)  
3 for i in FINAL:  
4     print(i)
```

നിയന്ത്രണം  
ഇതി  
ആകു  
ആനി