

# NATURAL LANGUAGE PROCESSING: MALAGASY PART-OF-SPEECH TAGGING

RAKOTONDRAFARA Aina Nambinintsoa<sup>1</sup>, RANDIMBINDRAINIBE Falimanana<sup>2</sup>,  
RANDRIAMBOLOLONA Nivo Harisoa<sup>3</sup>, ROBINSON Matio<sup>4</sup>

<sup>1</sup> Student in Doctoral School of Science and Technic Engineering and Innovation, Laboratory of Cognitive Sciences and Application, High School of Polytechnical University in Antananarivo, Madagascar

<sup>2</sup> Professor, in Doctoral School of Science and Technic Engineering and Innovation, Laboratory of Cognitive Sciences and Application, High School of Polytechnical University in Antananarivo, Madagascar

<sup>3</sup> Doctor, in Doctoral School of Science and Technic Engineering and Innovation, Laboratory of Cognitive Sciences and Application, High School of Polytechnical University in Antananarivo, Madagascar

<sup>4</sup> Doctor, in Doctoral School of Science and Technic Engineering and Innovation, Laboratory of Cognitive Sciences and Application, High School of Polytechnical University in Antananarivo, Madagascar

## ABSTRACT

The objective of this work is the realization of a part-of-speech tagging (POS tagging) for the official Malagasy Language. The tagger uses a set of tags, a tagged learning corpus, various mathematical models, and an algorithm that implements the basic rules of the Malagasy grammar. We will describe our tagset in a two-level approach and the performance of our part of part-of-speech tagging in relation to unigram-based tagger, the bigram model and the Hidden Markov Model. We will then present detailed results of each simulation.

**Keyword:** part-of-speech tagging, POS tagging, official Malagasy Language, tagset, corpus, algorithm, Malagasy grammar, unigram model, bigram model, Hidden Markov model (HMM Model)

## 1. INTRODUCTION

The part-of-speech tagging consists of assigning a tag that represents a grammatical class to a word or group of words. During this work, we will create an official Malagasy tagger based on a supervised learning. Our tagger uses the NLTK tool which will permit us to independently tag a sentence or a corpus from a learning corpus.

## 2. PRESENTATION OF THE TAGSET

In order to obtain satisfactory results, we have decided to fix the principles of Malagasy grammar, and to avoid the usual mistakes when tagging the learning corpus, in accordance with the writer and teacher Régis RAJEMISA-RAOLISON'S work "Grammaire malgache". As far as the creation of the tagset for the Malagasy language is concerned, we relied on the tagset of the NLTK tool.

Our tagset lead us to adopt a two-level approach tagset in which the first level corresponds to the 13 tags representing the 13 grammar classes and the second level to the complete set of 42 tags.

## 2.2 First level of tags

**Table -1:** First level of tags

Tags	Description
CC	Conjunction
CD	Number
DT	Determiner
NN	Noun
JJ	Adjective
PR	Pronoun
VB	Verb
RB	Adverb
IN	Preposition
RP	Particle
FW	Foreign word
IJ	Interjection
.	Punctuation

## 2.3 Second level of tags: Determiner's case

Malagasy Language is composed of two types of determiner [1]:

- Definite determiner: **ny, ilay, ikala**
- Nominal determiner: **i, Ra, An, ry**

**Table -2:** Second level of tags for the determiner

Second level of tags	Description	Example
DTDEF	Defined Determiner	<b>Ny</b> lanitra sy <b>ny</b> tany <b>The</b> sky and <b>the</b> earth.
DTNOM	Nominal Determiner	Tsy tonga <b>ilay</b> mpandrafitra <b>The</b> mason hasn't come.

## 2.4 Second level of tags: noun's case

Malagasy Language is composed of two kinds of name: the proper noun and the common noun. The proper name characteristics remain the same as that of any other language such as French, English, etc. Nonetheless, the common name's have its particularity.

**Table -3:** Second level of tag for the noun

Second level of tags	Description	Example
NNP	Proper noun	Mahay ny <b>Barea</b> The <b>Barea</b> are talented.
NN	Common noun	Apetraho ny <b>penina</b> Put your <b>pen</b> down.
NNDVB	Common noun derived from the verb	Milalao ny <b>mpianatra</b> The <b>students</b> are playing
NNDJJ	Common noun derived from the adjective	<b>Hatsaram</b> -panahin' olona A person's <b>goodness</b> of soul.

## 2.5 Second level of tags: adjective's case

There are five kinds of adjectives in Malagasy Language: the adjective qualifier, the numerical adjective, the interrogative adjective, the demonstrative adjective, the indefinite adjective. The possessive adjective does not exist in Malagasy, it is supplemented by the personal pronoun or suffix: **ko (o)**, **nao (ao)**, **ny ...**; instead of saying: my house, we say in Malagasy: *ny tranoko* (the house of mine)[1].

**Table -4:** Second level of tag for the adjective

Second level of tags	Description	Example
JJ	Qualifier Adjective	Trano <b>kely</b> A <b>little</b> house
JJCD	Numerical adjective	Mpianatra <b>roa</b> ihany no afaka Only <b>two</b> students succeeded
JJI	Interrogative adjective	Olona iza io ? <b>Who</b> is this man?
JJD	Demonstrative adjective	Omeo ahy <b>itsy</b> boky <b>itsy</b> Give me <b>that</b> book.
JJIND	Indefinite adjective	Olona <b>maro</b> no tonga <b>Many</b> persons came here.

## 2.6 Second level of tags: pronoun's case

There are six kinds of pronouns in Malagasy Language: the personal pronoun, the demonstrative pronoun, the relative pronoun, the interrogative pronoun and the indefinite pronoun.

**Table -5:** Second level of tag for the pronoun

Second level of tags	Description	Example
PRP	Personal pronoun	Nahoana <b>aho</b> no hihemotra ? Why would <b>I</b> fall back?
PRD	Demonstrative pronoun	Aza manao <b>izao</b> ! Do not do <b>it!</b>
PRR	Relative pronoun	Hosazina <b>izay</b> tsy manaraka ny lalàna We will punish <b>those</b> who violate laws
PRI	Interrogative pronoun	<b>Iza</b> no nilaza an'izany ? <b>Who</b> said that?

## 2.7 Second level of tags: verb's case

According to traditional grammar, the verb is a word that expresses the process, which means the action the subject does or undergoes [2].

There are three voices in Malagasy:

- The active voice where the subject is the agent of the action.
- The passive voice where the subject is the object of the action
- The relative or circumstantial voice where the subject is circumstance of the action

**Table -6:** Second level of tag for the verb

Second level of tags	Description	Example
VBA	Active verb	<b>Manarona</b> ny sakafo amin'ny lovia <b>aho</b> ? <b>I cover</b> the meal with the plate.
VBP	Passive verb	<b>Saronako</b> amin'ny lovia ny sakafo The meal is <b>covered</b> with the plate.
VBR	Relative verb	<b>Anaronako</b> ny sakafo ny lovia The <b>plate</b> serves <b>me</b> to <b>cover</b> the meal
VBN	Participle	Faly aho rehefa <b>tafiditra</b> ny trano <b>Once back</b> home I feel happy.

## 2.8 Second level of tags: adverb's case

In Malagasy grammar, there are several adverbs which are:

- Adverbs of time which express either time, duration, or frequency;
- Adverbs of place which specify ordinary places or even demonstrative;
- Adverbs of order which express ordering;
- Adverbs of number;
- Adverbs of affirmation which express affirmation;
- Adverbs of negation;
- Adverbs of interrogation;
- Adverbs of doubt;
- Adverbs of quantity;
- Adverbs of opposition;
- Adverbs of defense

**Table -7:** Second level of tag for the adverb

Second level of tags	Description	Example
RBT	Adverb of time	<b>Taloha</b> mora ny fiainana Once life was easy
RBP	Adverb of place	<b>Mifanatrika</b> ny tsena ny tranonay Our house is <b>in front of</b> the shop.
RBO	Adverb of order	<b>Voalohany</b> mila voaloha ny karamany <b>First of all</b> , our wages must be paid.
RBN	Adverb of number	Miditra <b>tsiroroa</b> ny mpianatra Student enter <b>two by two</b>
RBM	Adverb of manner	Mila ovaina <b>tsikelikely</b> ny toetsaintsika ratsy We haveto change our way of thinking <b>little by little</b> .
RBA	Adverb of affirmation	<b>Eny</b> ramose ! <b>Yes</b> sir !
RBC	Adverb of negation	<b>Tsy mbola</b> tonga ny ekipa nationaly National team hasn't arrived <b>yet</b> .
RBI	Adverb of interrogation	<b>Aiza</b> no nametrahany ny peninany? <b>Where</b> did he put his pen?
RBD	Adverb of doubt	Mety tara <b>angamba</b> izy He <b>might</b> be late.
RBQ	Adverb of quantity	Raha dinihina <b>kely</b> ny nataony If we take a look <b>a little more</b> on what he has done.
RBPR	Adverb of defense	<b>Aoka</b> izay ! <b>Stop</b> it !

## 2.9 Second level of tags: preposition's case

In the Malagasy language, prepositions are classified according to the relationship between the word to be completed and the complement[1].

**Table -8:** Second level of tag for the adverb

Second level of tags	Description	Example
IN	Preposition for a possession	Tranon'i Rakoto The house <b>of</b> Rakoto
INA	Preposition for an attribute	Miasa <b>ho an'ny</b> vahoaka izahay We work <b>for</b> the people
INP	Preposition for a place	Eo ankavanan'ilay lehilahy no misy azy She is located <b>on the right of</b> this man.
INT	Preposition for a time	<b>Mandritra</b> ny lanonana <b>During</b> the event.
INM	Preposition for a manner	Tsy fantatra mazava tsara ny <b>momba</b> an'ilay zaza very We do not know much <b>about the</b> missing child.
INR	Preposition for a reason	Tsy niasa izy <b>noho</b> ny antony ara-pahasalamana He didn't attend work <b>because of</b> health issues.
INW	Preposition for a mean	Mipetraka <b>amin'ny</b> rainy izy He lives <b>with</b> his father.
INC	Preposition for a comparison	<b>Toa</b> ireny mpanakanto ireny ianao You <b>look like</b> one of those artists.
INE	Preposition for an exception	Nalainy daholo ny entany <b>ankoatra</b> ireto akanjo ireto He took all his stuff with him <b>except</b> those clothes.

## 2.10 Combination of tags: the suffix or linked personal pronouns case

Here are the rules of suffixation of these pronouns to the words of which they are complements:

- If the word does not end by **KA, TRA, NA** we add the pronouns **ko, nao, ny, ntsika, nay, nareo and ny**.
- If the word ends by **NA**, we look for this final **NA** and we add the same forms of pronouns as before.
- If the word is terminated by **KA, TRA** we use the forms **o, ao, ny, tsika, ay, areo, ny** so that we drop the A of **KA, TRA** in front of the vowels **o, eo, ay, areo**, and the whole syllables **KA, TRA** in front of the consonants **ny, tsika**

The table below shows the application of these rules:

**Table -9:** Example of rule of suffixation

	Word who doesn't end with ka, tra or na	Word who end with na	Word who end with ka	Word who end with tra
<b>Words</b>	Loha / Head	Tanana / hand	Soroka / shoulder	Tongotra / feet
	Lohako / <b>my</b> head	Tanako	Soroko	Tongotro
	Lohanao / <b>your</b> head	Tananao	Sorokao	Tongotrao
	Lohany / <b>his</b> or <b>her</b> head	Tanany	Sorony	Tongony
	Lohantsika / <b>our</b> head	Tanantsika	Sorontsika	Tongotsika
	Lohanay / <b>our</b> head	Tananay	Sorokay	Tongotray
	Lohanareo / <b>your</b> head	Tananareo	Sorokareo	Tongotrareo
	Lohany / <b>their</b> head	Tanany	Sorony	Tongony
<b>Tags</b>	<b>NNPRP</b>	<b>NNPRP</b>	<b>NNPRP</b>	<b>NNPRP</b>

### 2.11 Combination of tags: The case of the preposition with possession report

In English, a case in point is: The driver's house, in Malagasy *tranon'ny mpamily*

- *tranon'ny* is composed of the name **trano** and the preposition **n'**
- the combination therefore gives as tag: **NNIN**

### 2.12 Combination of tags: The case of the derived name and the determiner

In English, a case in point is: Christ's disciple, in Malagasy *Mpianatr'i Kristy*

- *Mpianatr'i* is composed of the noun derived from the verb **mianatra** and the nominal determiner **i**
- the combination gives the tag then: **NNDVBDTNOM**

### 2.13 Combination of tags: The case of the participle and the personal pronoun

In English, a case in point is: At last, my work is done, in Malagasy *Vitako ihany ny asako /*

- *Vitako* is composed of the verb **Vita** and the personal pronoun **ko**,
- The combination gives the tag: **VBNPRP**

### 3. PRESENTATION OF THE LEARNING CORPUS

The learning corpus includes 6 468 grammatical elements provided from newspaper articles, advertisements and dialogues taken from Malagasy forums.

**Table -8:** Learning corpus characteristics

Grammatical elements	Frequency
Noun	1636
Adverb	732
Verb	742
Determiner	592
Punctuation	655
Adjective	618
Preposition	337
Particle	369
Conjunction	380
Pronoun	191
Foreign word	150
Number	66

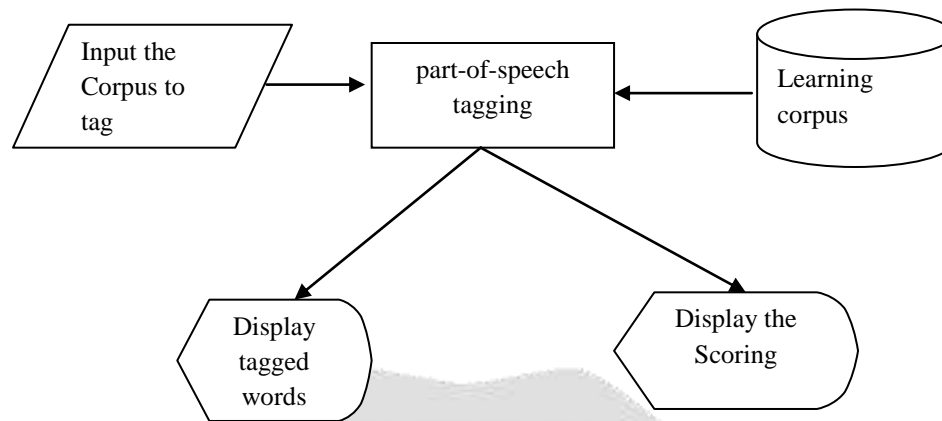
### 4. OPERATION OF THE MALAGASY PART-OF-SPEECH TAGGING

A part-of-speech tagger is a computer program that is able to recognize the grammatical nature of a word and assign a tag to that word. For our research, we used the NLTK tool.

The program operation itself can be divided into two stages:

- The first step is the part-of-speech tagging that consists of using our learning corpus and implementing our algorithm on how to assign a tag to a word according to the level of the desired set of tags.
- The second step is to tag each word of the test corpus and give a score on the percentage of correctly tagged words.





**Fig-1:** Our Part-of-speech tagger operation

With our learning corpus, there appear to be two major issues for our tagger:

- **An issue of precision:** our learning corpus is far from describing all the possible contexts of the Malagasy language, which means that even if the word is present in the corpus but understood in a context different from the corpus to be analyzed, it will be considered **unknown**.
- **An issue of coverage:** our learning corpus is very far from containing all the Malagasy words and the word order and tags possible, which means that a word or the word order that is not present in the learning corpus will be considered **unknown**.

Three models solve these problems but separately:

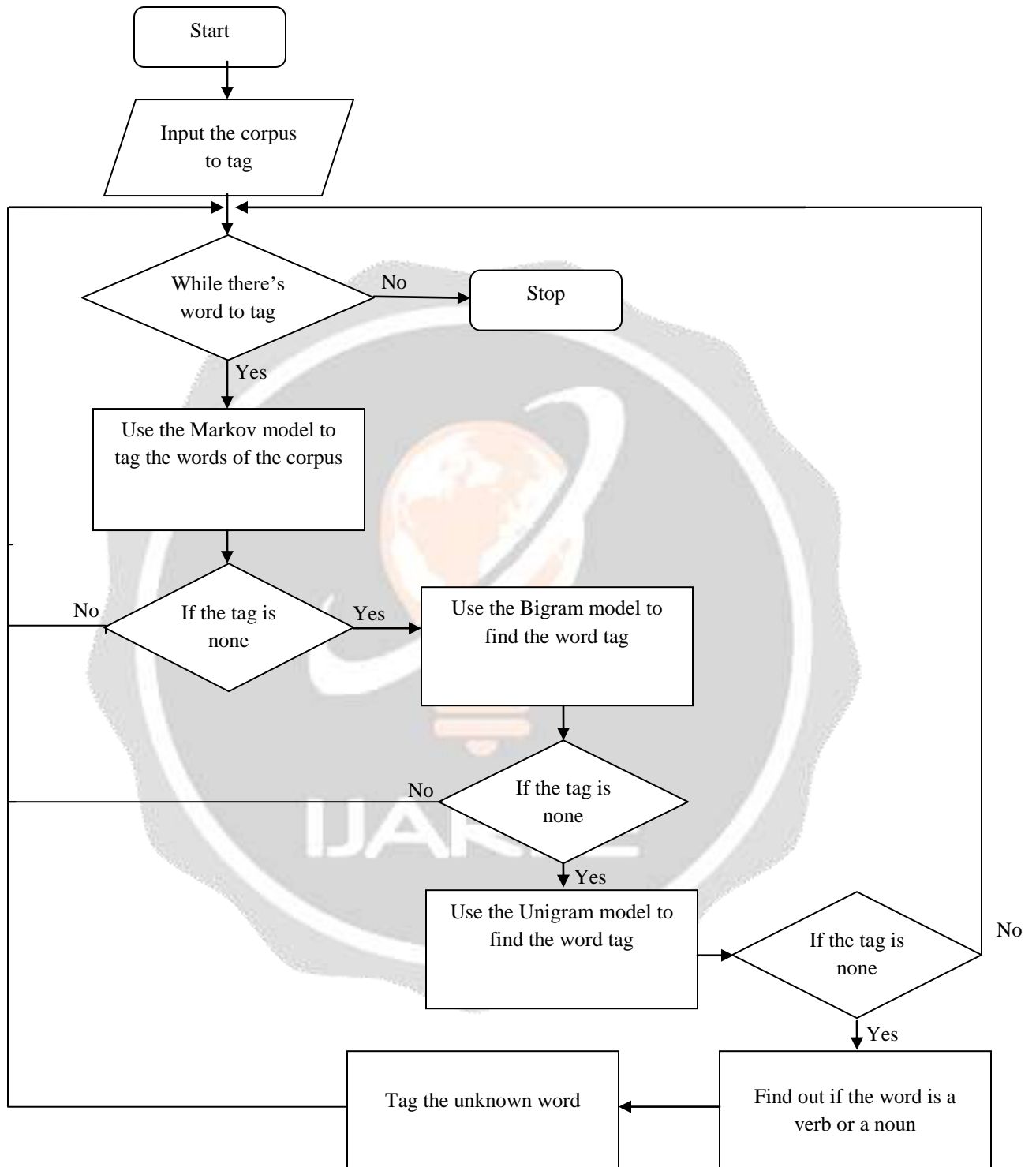
- The Unigram model, which solves most of the coverage problem but leaves a big flaw with regard to the precision problem. The reason is that this model relies on a simple statistical algorithm: for each word, it assigns the most probable tag for this word. For example, it attributes to the tag **JJ** any occurrence of the word *vaovao*, because *vaovao* is used more often as an adjective (for example, *hopitaly vaovao* / new hospital) than as a noun (for example, *Mijery vaovao* / Watch the news).
- The bigram model which resolves a part of the problem of precision and coverage. The reason is that this model is an improvement of the unigram model whose context is the current word associated with the tags of the word previously marked, so that we consider the previous word tagging in addition to the current word for tagging.
- The Hidden Markov Model (HMM) is the ideal model for maximizing accuracy in the case where our learning corpus is as rich and varied as Brown's corpus with 1,014,312 words, this model is a statistical mathematical model derived from Markov Chains, except that one cannot directly observe the sequence of states: the states are hidden. Each state issues "observations" which are observable. The purpose of the MMC is to associate with a sequence  $w = w_1 \dots w_i$  of words, a sequence  $t = t_1 \dots t_i$  of tags belonging to a set of tags.

Our tagger is based on a program that uses precision models, and then fills its coverage gaps with algorithms with more extensive coverage when needed. And for cases of unknown words, we have created a function that exploits the Malagasy grammar rules to recognize verbs or nouns.

The combination of these models works as follows:

- Use the tagger based on the NLTK Hidden Markov Model
- If the tagger based on the Hidden Markov Model cannot find a tag for the word, it refers to the tagger based on the NLTK bigram model.
- If the bigram tagger cannot find a tag for the word, it refers to the NLTK unigram tagger.
- If the unigram tagger is unable to find a tag for the word, it returns the unknown word to the function which tags the word with the first level of tag in verb or noun.

The flowchart below reflects this process:



**Fig-2:** The tagging process flowchart

## 5. RESULT AND INTERPRETATION

In order to test the effectiveness of our part-of-speech tagging, we have decided to make a simulation on four types of corpus of different nature and show the results obtained with a special focus on the quantity of word contained in our learning corpus and the different mathematical models.

**Table -9:** Result obtained according to the words contained in our corpus

	<b>Newspaper article</b> <b>4000 words</b>	<b>Advertising</b> <b>2000 words</b>	<b>Blog</b> <b>2700 words</b>	<b>Forum</b> <b>1700 words</b>
Words contained in our learning corpus	3875	1889	2008	950

**Table -10:** Results obtained with some mathematical models

<b>Tagger Precision</b>	<b>Newspaper article</b> <b>4000 words</b>	<b>Advertising</b> <b>2000 words</b>	<b>Blog</b> <b>2700 words</b>	<b>Forum</b> <b>1700 words</b>
Correct tag with Unigram Model	78,68 %	73,45 %	64,25%	55,35 %
Correct tag with Bigram Model	45,56 %	40,36 %	37,02 %	28,13 %
Correct tag with Hidden Markov Model	5,04 %	6,02 %	2,56 %	0,02 %
Correct tag with our POS tagging	80,03 %	75,63 %	67,25 %	57,85 %

### 5.1. Interpretation according to the nature of the analyzed corpus

For the case of the newspaper articles, the analyzed corpus comes from the newspaper “Midi Madagascar” and the type of the treated information concerns the politics and various facts. The table shows that the newspaper articles present better results than the other corpus which reason might be:

- More than 3/5 of our training corpus is composed of newspaper article which makes the majority of the words that compose the corpus is already known.
- Newspaper articles are written in official Malagasy without abbreviations
- The authors of the articles use a sustained register when the newspaper deals with information on the policy, and a current register when dealing with the various facts, both are identical to the one applied in our learning corpus.

For the case of the advertisements, the analyzed corpus comes from various ads in the newspapers “Midi Madagascar”. The results obtained do not differ too much from the results obtained with the articles. The reason might be:

- The use of official Malagasy
- The use of the same active verbs
- The non-respect of Malagasy grammatical structures because the purpose is to pass information quickly.
- The use of the current and sometimes familiar register
- Few use of technical term

For the case of the blogs, the corpus blog addresses to analyze young people's interest towards studies.

Those results might come from:

- Use of foreign words, usually French language but also English, to express technical terms
- The permanent use of a familiar register
- The blog deals with a particular topic, some words or vocabulary are not present in our learning corpus

For the case of the forums, the worst results are obtained with the forums in a Malagasy discussion group where the reasons might be:

- The amalgamation of Malagasy and French language
- The use of other Malagasy dialects except the official Malagasy
- The use of abbreviations and SMS languages

## 5.2. Interpretation according to the mathematical model applied

- The Unigram model shows high success rate in the chart which might be explained by the fact that the majority of the words that make up the different test corpus are contained in our learning corpus. The reason why the precision does not match lies in the numerous grammatical natures of a same word in Malagasy Language. Indeed, its nature depends on the sentence in which it is used.
- The Bigram model has a below average result which reason lies in the several concept (potential tag sequence) that are not present in our learning corpus.
- The Markov model presents the least satisfactory results which main cause lies in the incorrect tagging of unknown words. The Hidden Markov model is a statistical model which means that a mistagged word leads to a chain reaction and thus, all the following word / tag combination will be distorted.
- As far as our POS tagging is concerned, it was possible to obtain a result superior to the three models mentioned above. This is because our combination technique proves to be more efficient than the three individual models.

In order to obtain a more relevant interpretation through the mathematical model, we propose an example that accurately illustrates the difference between the results obtained for each model used.

The corpus to tag: *Ny Polisy, ny Bianco ary ny vahoaka*

In our example, the word “Bianco” is considered as an unknown word and is tagged with the tag **None**.

- Results obtained with the Unigram model:  
[(u'Ny', u'**DTDEF**'), (u'Polisy', u'**NN**'), (u', u'), (u'ny', u'**DTDEF**'), (u'Bianco', **None**), (u'ary', u'**CC**'), (u'ny', u'**DTDEF**'), (u'vahoaka', u'**NN**')] ]

As far as the Unigram model is concerned, despite the fact that this word “Bianco” is unknown, the remaining words are tagged correctly.

- Result obtained with the Bigram model:

[(u'Ny', u'**DTDEF**'), (u'Polisy', u'**NN**'), (u', u'.), (u'ny', u'**DTDEF**'), (u'Bianco', **None**), (u'ary', **None**), (u'ny', u'**DTDEF**'), (u'vahoaka', u'**NN**')] ]

Then for the Bigram Model, only the determiner followed by the name are tagged correctly because this sequence of tag exists in our learning corpora.

- Result obtained with the Hidden Markov model:

[(u'Ny', u'**DTDEF**'), (u'Polisy', u'**NN**'), (u', u'.), (u'ny', u'**DTDEF**'), (u'Bianco', u'**DTDEF**'), (u'ary', u'**DTDEF**'), (u'ny', u'**DTDEF**'), (u'vahoaka', u'**DTDEF**')] ]

Concerning the hidden Markov model, all the word that follows the unknown word “Bianco” is tagged with the tag **DTDEF** due to its recurrence in the corpus learning.

- Result obtained with our POS tagging:

[(u'Ny', u'**DTDEF**'), (u'Polisy', u'**NN**'), (u', u'.), (u'ny', u'**DTDEF**'), (u'Bianco', **NN**), (u'ary', u'**CC**'), (u'ny', u'**DTDEF**'), (u'vahoaka', u'**NN**')] ]

Due to our algorithm, the word “Bianco” is tagged as a noun and the rest of the word is tagged correctly.

## 6. CONCLUSIONS

In a nutshell, we can say that the accuracy of our POS tagging depends mainly on the right balance between the coverage of our learning corpus and the accuracy of the mathematical models used to deduce the nature of a word according to its place in the corpus.

As a perspective of future study, an improvement of the learning corpus could be set up. This improvement will carry a more assorted corpus in addition to the use of algorithms to recognize the SMS languages and the improvement of the algorithms of recognition of an extended grammatical item.

## 7. REFERENCES

- [1]. Régis RAJEMISA-RAOLISON. GRAMMAIRE MALGACHE. Fianarantsoa 1971
- [2]. Joro Ny Aina RANAIVOARISON. Modélisation de la morphosyntaxe du malgache : Construction d'un dictionnaire électronique des verbes simples. Université d'Antananarivo, 2014
- [3]. Abdelhamid EL JIHAD, Abdellah YOUSFI. Etiquetage morpho-syntaxique des textes arabes par modèle de Markov caché. Université Mohamed V, Rabat, Maroc, 2005
- [4]. Kristina Toutanova and Christopher D. Manning. Enriching the Knowledge Sources Used in a Maximum Entropy Part-of-Speech Tagger. In Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora (EMNLP/VLC-2000), 200