# A REVIEW PAPER ON PART OF SPEECH TAGGER FOR HINDI

Kartik Yadav, Saumya Mishra, Awadhesh Kumar Srivastava
Department of IT
K.I.E.T. Group of Institutions, Ghaziabad, U.P., India

**Abstract: Some segment of Part of Speech (POS) naming in hindi is so far an open issue. We regardless of everything miss the mark on a sensible strategy in executing a POS Tagger that uses hindi. At the present time delineate our undertakings to develop a POS Tagger which is based on Hidden Markov Model(HMM). We have used hindi Part of Speech mark set for the headway of this tagger. The Artificial Intelligence(AI) based approach for Named Entity Recognition(NER) is logically powerful and conservative and moreover requires less proportion of language dominance appeared differently in relation to manage based strategy. Among various AI systems HMM is one of the successful procedure to use and execute with various extra features discussed in before areas. Without a doubt, HMM has not been very used, and likewise the adequately developed work has not passed on reasonable precision, owing to this the ebb and flow inquire about work is given to HMM in NER for Indian vernaculars We have endeavored to achieve the most extraordinary exactness possible.**

**Keywords: Hidden Markov Model, Hindi Part of Speech Tag set, Part of Speech Tagger.**

## I.   INTRODUCTION

POS Tagging is the very basic stage for any Natural Language Processing Application. This is an effort which provides POS imprints to words which are mentioned in the sentence. The words which are present in the sentences are just a sequence of an action which need to be labelled. The previous word-name blend is used to provide new tagset to words in the sentences. The main application of POS tagging is parsing where words and its marks are changed according to previous tag-set combination. The taggers have been used in Machine Translation (MT) while designing an MT Engine based on commerce. According to a sentence structure, the tagger performs performs its word of labelling which is then further sent for parsing. In the case of NER, the word belongs to the name of an individual, region, affiliation, date and time etc.

The definition of components of speech has been primarily based on morphological and syntactic function; phrases that function in addition with appreciating to the affixes they take (their morphological properties) or with respect to what can occur nearby (their distributional properties') are grouped in training. While phrase classes are having inclinations toward semantic coherence (nouns frequently describe 'people, places or things', and adjectives regularly describe properties), this isn't always the case, and in semantic coherence isn't used as a defining criterion for elements of speech [1]. Parts of speech may be divided into two broad amazing categories: closed class types and open class types. Closed Class Types: Closed classes are those which have a relatively constant membership. For example, prepositions are a closed class due to the fact there is a hard and fast set of them in English; new prepositions are rarely delivered. Open Class Types: By assessment nouns and verbs are open classes because new nouns and verbs are continual brought or borrowed from different languages. It is likely that any given speaker or corpus will have exceptional open class words, but all audio systems of a language, and corpora that are massive enough, will likely share the set of closed class shared by Singh et al. (2013).

A large amount of work has been done on Part of Speech labelling. All the efforts can be directed into three parts which include: In the case of methodologies based on rules, human annotation is required, which is used for rule making for proper tagging of words or fact based approach is used where scientific details based on hybrid approach is used which is partially similar to rule based approach. Part of Speech taggers commonly utilize an AI based approach, but in Indian setting we don't have proper methodology. Right now we have HMM based tagger.

## II.   LITERATURE SURVEY

We have seen that the majority of the investigations on POS labeling on the South-Asian dialects has been finished utilizing stochastic labeling models like HMM, MEM and so on. A POS labeling approach dependent on Maximum Entropy Markov Model utilizing managed instruction. This model trains the utilization of a pre-labeled corpora and utilizations a list of capabilities anticipate the tag for a word. The capacity set incorporates POS labeling capacities, setting based highlights, word capacities, word reference capacities and corpus based thoroughly includes. The tagger surveys a precision of 89.34% on the advancement data of the NLPAI Machine Learning Competition 2006. To achieve such generally speaking execution, the tagger utilizes a pre-commented on

instruction corpus together with cycle 35,000 expressions commented on with a tag set counting 29 elites POS labels. As portrayed in Selvam at al.(2009) an HMM based POS tagger became developed which set up 83.41 accuracy.

The tagger was additionally prepared for a pre-clarified corpus including 40956 tokens. The tagger got to try on a clarified corpus, having 5967 tokens. The tagger sets up 75.25 accuracy when tried on a un-clarified check set, for example, 5129 tokens. However, every other HMM based absolutely tagger is characterized in Bharati at al.(2014), revealing a by and large execution of 76.49 curacy on tutoring and investigate information having about 15000 and 6000. expressions, separately. This tagger utilizes HMM in blend with chance models of certain relevant highlights. In Avinesh at al.(2015), the creators report a hybrid tagger for Hindi that sudden spikes in demand for two stages to POS label input content. In the primary stage, the HMM based TNT tagger is run on the untagged content to play out the underlying labeling. During this stage, a lot of change rules are initiated which are utilized later. In the subsequent stage, the arrangement of change rules learnt prior is utilized on the at first labeled content to address mistakes made in the principal stage.Be that as it may, the presentation of this tagger isn't in the same class as the other taggers detailed for Hindi. It utilizes a corpus of 25,000 words commented on with 24 labels, and the subsequent precision is 78.26% utilizing the TnT tagger. The creators recommend that the low score could be the consequence of the scantiness of the preparation information. The utilization of the arrangement of change governs in post preparing improves the general exactness to 82.74%.

For Bengali, Dandapat et. al.(2007) used the Hidden Markov model and Maximum Entropy. They used anatomical analysis for setting the tag set for corpus. They used pre-defined tagger and semi-rule based tagger and proposed an exactness of around 88%. They Support Vector Machine based taggers were also proposed. They showed the accuracy of 87.14%..Ekbal et. al.(2007) Additionally, built up a Conditional Arbitrary Fields based tagger. This tagger was made by using the prefix and postfix data of sentences with common labels.This provides an accuracy of more than 90%.They used this tagset to explain their corpus, and prepared their model subsequently. The tagset based on the support vector machine removed the phonetic data. Instead, five annotators were used to render the Part of Speech labeled set that finished outlining this structure in a quarter of a year's assignment.

For Malayalam, Manju et. had proposed a tagger based on HMM. Since they did not have a simple corpus, an anatomical analysis tool was used to produce a corpus which was then used to prepare the calculation for HMM. Another tagger to Malayalam was created by Shrivastava et. al. (2008) Employ Vector Support Machines. They used an SVM labeling tool which Gimé developed. Anthony et. To build up this tagger. Al. initially proposed a tagset that would be reasonable for Malayalam and subsequently made a comment on the corpus using the tagset. With their tagset, their tagset showed 89 per cent accuracy.

## III.    HMM BASED POS TAGGER FOR HINDI

We were first asked to transcribe a tagset based corpus in order to develop an HMM based tagger. We have used the tagset IL POS Modi et al.(2016) .To train our program we used 5,000 sentences (approx. 30,000 words) from the tourism domain.

A HMM-based POS tagger provides the best tag for a term by processing the marks forth and turning their probabilities around. Any probability of labeling change is determined by assessing repeat inspection of two names seen together in the catalogue, separated by repeat inspection of the former tag seen separately in the corpus. This is because we realize it is happening slowly. A run of the mill thing and not a relational word or pronoun will trail an enlightening word for instance.

To prepare our system, we have combined 4,500 sentences (approx. 20,000 words) from the travel industry room with the tourism domain. As this is a variant corpus, we have not invested in creating morphological analyzers in our energy amounts. Assume for example a thing or action, word or descriptive word or intensifier in the event that we have a word which is an open class word. There is a possibility at that stage that it will be assigned to various labels and we may be faced with the problem of vagueness. For example, we have an equivocal word designated for a thing and a word for action.

This representation environment is the revolutionary highlight of HMM which can pick the tag for a word by taking a close look at the tag of the past word and the tag of things to come. This wonder which is a conceptual design where there is a concealed simple generator of measurable occasions (the tag label's probabilities) and as many states can be seen this veiled generator. We'll probably be locating the states.

## IV.    MODULE DESCRIPTION

### A: Chunking

Chunking is distinguishing proof of limits of sensible substances which are fundamental for various NLP applications, for example, machine interpretation, adjusting equal writings, content to-discourse frameworks, programmed abstracting and named element acknowledgment. In this manner data on the named element limits is basic. Utilizing a full parser for distinguishing these limits alone, can frequently be a costly choice for this framework master has acknowledged that these restrictions on names substances are combined physically in the corpus. Chunking is breaking of a huge book into intelligent elements. In the setting of HMM framework named substances might be made out of an assortment of tokens. Consequently, for recognizable proof of every

one of these tokens which make single named substance, we have this preprocessing movement of NER which is known as Chunking. For instance 𝇊𝇊𝇊𝇊 𝇊𝇊𝇊𝇊𝇊 𝇊𝇊𝇊𝇊 𝇊𝇊𝇊𝇊𝇊 - 𝇊, is a solitary named element and complete name ought to be labeled as Organization (ORG) by Named Entity Recognition framework. In the current investigation, it has been expected that these sorts of words are as of now combined (chunked) into one unit either physically or by any accessible content chunkers, if any.

In the framework, a scramble (-) has been physically put in the middle of tokens for instance to make it single element else it won't be treated as single substances by the HMM framework and the framework will label it independently as follows CON /OTHER /OTHER /OTHER where CON speaks to the nation and OTHER speaks to NOT-A-Name To maintain a strategic distance from this sort of case, manual piecing action is required that ought to be performed on testing sentence and just as on the preparation corpora. Lumping should be possible naturally if such a content chunked framework exists with sensible accuracy. Therefore if testing sentence is Chunking ought to be performed first at that point continue further for next stage. Subsequent to piecing model sentence will look like in the accompanying structure and become single unit, one can expect that named substance recognizer is required to label them PER(PERSON), INST(Institute), and DEPT(Department) separately. PER, INST, DEPT are simply models, one can choose these labels as indicated by one's accommodation.

### B: Annotation

Annotation module is liable for improvement of assets, that is building up the preparation corpus from crude content. Since crude content can not be straightforwardly utilized by Viterbi. Calculation, in this manner information must be appropriately arranged before preparing. This arrangement requires labeling named substances of preparing corpus by their group of named elements. The annotation task itself is tedious and work escalated, including a lot of manual altering and focus. Furthermore, the recognizable proof of formal people, places or things and their classes may require a reasonable level of phonetic information on the annotator.

Like Indian dialects numerous normal dialects are asset poor, thus HMM framework gives an interface to build up these kinds of assets for any Indian dialects, and in any event, for any regular language with some additional human endeavors. These assets are valuable for choosing names substance labels for a sentence not present in preparing corpus. This bit of the structure will help the customer in making marked corpus by essentially clicking in a manner of speaking. This interface is developed with the objective that the tremendous corpus can be explained by the customer for name set picked without any other person with least undertakings and little language ability.

### C: Tokenizer

Tokenization is the act of splitting a sequence of strings into chunk like words, keywords, phrases, signs and other components called tokens. Tokens can be single words, saying or perhaps even entire sentences.

Since NER framework is relied upon to label all the named elements by their named substance class and different words by NOT-A-NAME tag, subsequently breaking of sentence in different tokens before named element labeling is required. Tokenization process is performed on the model sentence which has been lumped as of now. The model sentence is 𝇊𝇊𝇊𝇊𝇊- 𝇊𝇊𝇊𝇊𝇊 𝇊𝇊𝇊𝇊𝇊𝇊-𝇊𝇊𝇊𝇊𝇊𝇊𝇊𝇊 𝇊𝇊 𝇊𝇊𝇊𝇊-𝇊𝇊𝇊𝇊𝇊𝇊 𝇊𝇊Since NER framework is relied upon to label all the named elements by their named substance class and different words by NOT-A-NAME tag, subsequently breaking of a sentence in different tokens before the named element labeling is required. Tokenization process is performed on the model sentence which has been lumped as of now. And its tokenized structure is formed.

Sentence Extraction: This module extract the sentence from the corpus.
End structure: Punctuation as the delimiter.
Information: It is in Hindi.
Corpus Handling: The corpus is split into the sentences using end structure. The sentences are then stored into sentence agent table.

Word Tokenizer: This module extract words from the sentence.
End structure: Space as the delimiter
.delimiter. Data: Extracted Hindi sentence Preparing: Split the sentence as appeared by the space delimiter and store them into Lexicon table. Yield: Display the splitted Hindi words.

### D: Unknown Word Handling

Obscure words are characterized as the words which are not in the vocabulary. It has been told the best way to assess named substance label age probabilities of words happening in the corpus. Be that as it may, there can be numerous words in sentences which requires labels yet are not in preparing corpus. It is unimaginable to expect to list the entirety of the words in the dictionary. Likewise it is preposterous to expect to compose basic Guidelines which can list the named elements without over-age or under-age The event of typographical mistakes makes the issue much progressively entangled. It is hard to tell when an obscure word is experienced. This implies obscure words are a serious issue for taggers and practically speaking, contrasts exactness of different taggers over various corpora.

Recognizable proof for every obscure word is troublesome and might require embracing various methodologies. Diverse kinds of names elements must be recognized utilizing either substance or setting subordinate guidelines. Legitimate names have less substance, consistency thus recognizing them depends

more on relevant data. There is right now no good calculation for recognizing both obscure words and typographical blunders, yet specialists are independently chipping away at such sorts of issue. Since any setting subordinate choice about naming substance tag of obscure word will make the framework, language and area explicit and since just named elements must be recognized in this way the transliteration process has been attempted to distinguish named element tag of obscure words.

Transliteration Unknown Word Handling: Transliteration is the modification of content from one language to another. Transliteration might be characterized as changing over word by word or letter by letter starting with one language then onto the next. Transliteration can assume a vital job in settling the issue of obscure words in Named Entity Recognition. The transliterated obscure word dealing with, the framework will keep up transliterated message in different dialects for each named element in the corpus. When the obscure word is experienced in testing sentence (Observations) it is additionally transliterated in different dialects which are referenced in the current framework. The following stage is to discover transliterated content of an obscure word in the rundown of transliterated named elements in preparing corpus. For instance if framework is prepared for English name 'Smash' as PERSON and there is no sentence having Hindi name □□□ and Hindi □□□ exists in perceptions that is in trying sentence then in typical circumstance Viterbi calculation will produce mistake message for obscure word. In the proposed framework obscure word dealing with module will transliterate □□□ in English and create 'Smash' which is known as PERSON tag through preparing corpus. Proposed framework will allocate a similar tag to Hindi □□□ too.

Essentially it can deal with obscure expressions of different dialects names too. Hence, if there are three names in preparing corpus of three languages (for model three names written in English, Hindi, and Punjabi dialects) at that point this framework can recognize nine names of those three dialects. For instance English Ram, Hindi □□□□ and Punjabi □□□□ is labeled as PERSON in preparing corpus. By transliteration obscure word dealing with module we can get PERSON tag to Hindi and Punjabi proportional to 'Slam', English and Punjabi identical to '□□□' including Hindi and English proportionate to □□□□ (Punjabi) which was not prepared in preparing corpus.

### *E: Tagging*
This module name each Hindi word in the sentence with their linked tags like INTF, CCS, PRP, DMD, NN, QTC, PSP, VM, VAUX, PUNC and so on whose description (tag) are given in Table1. If the Hindi word isn't matching in any arrangement of the syntactic element tag by then name that word with "No_Tag". It in like manner perceives and show mark configuration like Start structure, Mid model, End plan. Data: Extracted Hindi sentence Handling: Tag every... Tag every declaration of information sentence. Yield: Display the mark yield.

Supervised POS Tagging: The directed POS labeling systems involve a pre-labeled corpora that is utilized for preparing to study data regarding the tagset,word-label quantities, principle sets and so forth [10]. The exhibition of the models by and large increment with the expansion in size of this corpora.

Unsupervised POS Tagging: Unlike the governed models, the unregulated POS labeling designs do not involve a pre-labeled corpora. Instead, they use advanced computational techniques such as the calculation Baum-Welch to function on targets, modify rules, and so on. Either they decide the probabilistic data needed by the stochastic taggers.

Rule based POS Tagging: The basic POS labeling models apply several transcribed guidelines and use logical data to relegate POS labels to words. Such principles are also regarded as the guidelines for setting outlines. For example, a setting outline rule can state something like: "If a Determiner precedes a questionable / obscure word Z and a Noun follows, mark it as an Adjective." Besides this, transformation based POS Tagging uses predefined...predefined set of standards that are used to create tagsets.

The term morphology is used in semantic analysis which suggests the way words are made from smaller units of importance which are called "morphemes". This tagging is utilized by certain models to remove ambiguity and provide disambiguated procedure. The standard based labelling models require proper preparation of predefined corpora. Still, large amount of work is used to accept those changes.

Stochastic POS Tagging: This method involves repetition, likelihood and insights. This methodology uses most frequent and most utilized tags for the given word which is used for preparing information. The data is used to label the word. Sometimes, the main problem is that it contradicts syntax rules of the language. Contrary to the word recurrence approach, a choice is known as the n-gram approach which establishes the probability of a given label succession. It determines the best tag for a word by measuring the probability that it will happen with the past n marks, where the value of n is set to 1,2 or 3 for purposes down to earth. These models are called the Unigram, Bigram and Trigram. The most well-known calculation for updating an n-gram method for the marking of new material is known as Viterbi Algorithm[8].

Conditional Random Fields POS Tagging: A Conditional Random Field (CRF) is a system of probabilistic model to fragment and mark a grouping of information. A restrictive model indicates the probabilities of conceivable mark arrangements given a perception grouping. The restrictive likelihood of the mark succession can rely upon self-assertive, non-free

highlights of the perception arrangement. The likelihood of a progress between marks may depend on the present perception, yet in addition on past and future perceptions [10]. The CRF model figures the likelihood dependent on certain highlights, which may incorporate the postfix of the present word, the labels of past and next words, the real past and next words and so on [9].

### F: Training

In phonetics a corpus is a huge organized arrangement of writings as a rule electronically put away and prepared. Corpora are assortment of the corpus. Since huge advancement can be made in content comprehension by endeavoring to naturally extricate data about language from an extremely enormous corpora. Hence, numerous assets have been created to help the learning task. There are assets like plain corpus which don't give any additional data about content however outright content and is usually named as crude corpus (untagged). There are other content assets which present a content with some additional data such content is a marked (explained) content which decide the class of name for which they are helpful.

Administered inquire about tagging has been the significant main impetus in ongoing NLP advancements. Computational Linguists use corpora in addition to other things, to watch (and propose) phonetic speculations (rules), to upgrade them and to at last assess them (or the methodologies dependent on those principles). This framework functions admirably for huge corpus size as well as for little corpus. These commented on corpora can be valuable for analysts working in these zones.

### G: Testing

This module provides for the customer the supply of testing discernment by changing sentence or by picking record having colossal testing corpus. Testing corpus supply each line of corpus consequent to tokenizing it as commitment to Viterbi Algorithm which will give perfect mark progression to up-and-comer sentence.

Table1: Description of Hindi POS Tagset

| S. No. | Tag | Description | Examples |
|---|---|---|---|
| 1. | NN | Common Nouns | □□□□, □□□□, □□□□, □□□□ |
| 2. | NST | Nourn Denotating Spatial and Temporal Expressions | □□□, □□□□, □□□□, □□□ |
| 3. | NNP | Proper Nourns (name of person) | □□□□, □□□, □□□□□ |
| 4. | PRP | Pronoun | □□, □□, □□□, □□□ |
| 5. | DEM | Demonstrative | □□, □□, □□ |
| 6. | VM | Verb Main (Finite or Non-Finite) | □□□□, □□□□, □□□□, □□□□, □□□□, □□□□ |
| 7. | VAUX | Auxiliary Verb (Any verb, present besides main verb shall be marked as auxillary verb) | □□, □□□, □□ |
| 8. | JJ | Adjective (Modifier of Noun) | □□□ □□ &□□, □□□□□□ , □□ □□□□ |
| 9. | RB | Adverb (Modifier of Verb) | □*□+, □□□□, □□□□ |
| 10. | PSP | Postposition | □-, □□, □□ |
| 11. | QC | Cardinals | □□, □□□, □□ |
| 12. | QF | Quantifiers | □□□□, □□□□, □□ |
| 13. | RP | Particles | □□, □□, □ |
| 14. | CC | Conjuncts (Coordinating and Subordinating) | □□ |
| 15. | WQ | Question Words | 2□3, 2□□, □□□ |
| 16. | QO | Ordinals | □□□□, □□□□□, □□□□□ |
| 17. | NEG | Negative | □□ |
| 18. | INJ | Interjection | □□□, □□□ |
| 19. | INTF | Intensifier | □□□□, □□□□, □□ |
| 20. | SYM | Symbol | ?, ; : ! |
| 21. | XC | Compounds | □□□-□□□□□ |
| 22. | RDP | Reduplications | □□□□+□□□□ □□+□□ |
| 23. | ECH | Echo Words | □□□-□□□, |
| 24. | UNK | Forigen Words | English |

## V. EVALUATION

Recall (R) and Precision (P) were used on a regular basis to measure the performance of data recovery and information derivation systems. Precision handles replacement and insertion errors while recall handles replacement and deletion errors. It is also important to provide a common performance indicator that covers all three forms of errors – simultaneous replacements, insertions and deletions. One merit figure, the F-measure, has been defined as a combination of P and R that is weighted.

Precision (P) = Number of correct tags given by system/Total tags provided by system.

OR

Precision (P) = Correct Responses / Correct +Incorrect + Missing Responses.

OR

Precision (P) = No. of correct POS tags assigned by the system/No. of POS tags assigned by the system.

Recall (R) = Number of correct tags given by system/ Number of possible responses.

<div align="center">OR</div>

Recall (R) = Correct Responses / Correct + Incorrect + Spurious Responses.

<div align="center">OR</div>

Recall (R) = No. of correct POS tags assigned by the system/No. of POS tags in the text.

F-Measure = 2xPxR/(P+R).

<div align="center">OR</div>

F -Measure=2PR/(P + R).

Our system has been tested for its performance by developing a test corpus, which comprises of 500 sentences (11,720words). The standard performance of a system is considered to be recall, precision and f-measure, therefore they are calculated.

The test scores given by our system are:
The number of POS tags assigned by the system correctly = 10652.
The number of POS tags alloted by the system = 11651.
The number of POS tags existing in the text = 11520.

The f-measure would be the same since recall and precision are same. Thus the accuracy of the system is 88.4%.

<div align="center">

VI.     CONCLUSION

</div>

In this paper, we have dealt with the POS Tagging using HMM approach. We have tried POS naming in hindi is order to solve the open issue. We have used Part of Speech mark for Indian language set for the headway of this tagger. We have found that the AI based approach for NER is logically powerful and conservative and moreover requires less proportion of language dominance appeared differently in relation to manage based strategy. Among various AI systems HMM is one of the successful procedure which we have used and execute with various extra features discussed in before areas. Without a doubt, HMM has not been very used, and likewise the adequately developed work has not passed on reasonable precision, owing to this the ebb and flow inquire about work is given to HMM in NER for Indian vernaculars. We have endeavored to achieve the most extraordinary exactness possible which is 85.40%.

<div align="center">

VII.     REFERENCES

</div>

[1] Singh, S., Gupta, K., Shrivastava, M., Bhattacharya, P., (2006) "Morphological Richness Offsets Resource Demand- Experiences in Constructing a POS Tagger for Hindi". In: COLING/ACL, pp. 779-786.

[2] Dalal, A., Nagaraj, K., Sawant, U., Shelke, S., (2006) "Hindi Part-of-Speech Tagging and Chunking: A Maximum Entropy Approach". In: NLPAI Machine Learning Competition.

[3] Dandapat, S., Sarkar, S., Basu, A., (2007) "Automatic Part-of-Speech Tagging for Bengali: An Approach for Morphologically Rich Languages in a Poor Resource Scenario". In: Association for Computational Linguistic, pp 221-224.

[4] Ekbal, A., Haque, R., Bandyopadhyay, S., (2007) "Bengali Part of Speech Tagging using Conditional Random Field". In: 7th International Symposium of Natural Language Processing(SNLP-2007), Thailand Pattaya, 13-15 December 2007, pp.131-136.

[5] Shrivastava, M., Bhattacharyya, P., (2008) "Hindi POS Tagger Using Naive Stemming: Harnessing Morphological Information Without Extensive Linguistic Knowledge". In: International Conference on NLP (ICON08), Macmillan Press, New Delhi.Manju K., Soumya S., Sumam, M. I., (2009) "Development of a POS Tagger for Malayalam - An Experience". In: International Conference on Advances in Recent Technologies in Communication and Computing, pp.709-713.

[6] Selvam, M., Natarajan, A.M., (2009) "Improvement of Rule Based Morphological Analysis and POS Tagging in Tamil Language via Projection and Induction Techniques". International Journal of Computers, 3(4).

[7] Dhanalakshmi, V., Kumar, A., Shivapratap, G, Soman, K.P., Rajendran, S, (2009) "Tamil POS Tagging using Linear Programming". International Journal of Recent Trends in Engineering, 1(2).

[8] Dhanalakshmi V, Anand kumar M, Rajendran S, Soman K P., (2009) "POS Tagger and Chunker for Tamil Language". Proceedings of Tamil Internet Conference 2009.

[9] Singh, J., Joshi, N., Mathur I., (2013) "Part of Speech Tagging of Marathi Text Using Trigram Method", International Journal of Advanced Information Technology, pp 35-41, Vol 3. No. 2.

[10] Bharati, A., Sharma, D.M., Bai, L., Sangal, R., (2014) "AnnCorra: Annotating Corpora Guidelines for POS and Chunk Annotation for Indian Languages", http://ltrc.iiit.ac.in/tr031/posguidelines.pdf.

[11] Avinesh, PVS, Karthik, G., (2015) "Part-Of-Speech Tagging and Chunking using Conditional Random Fields and Transformation Based Learning". In: NLPAI Machine Learning Competition.