

# Part-of-Speech Annotation Challenges in Marathi

Gajanan Rane, Nilesh Joshi, Geetanjali Rane, Hanumant Redkar,  
Malhar Kulkarni and Pushpak Bhattacharyya

Center For Indian Language Technology  
Indian Institute of Technology Bombay, Mumbai, India  
{gkrane45, joshinilesh60, geetanjaleerane, hanumantredkar,  
malharku and pushpakbh}@gmail.com

## Abstract

Part of Speech (POS) annotation is a significant challenge in natural language processing. The paper discusses issues and challenges faced in the process of POS annotation of the Marathi data from four domains *viz.*, tourism, health, entertainment and agriculture. During POS annotation, a lot of issues were encountered. Some of the major ones are discussed in detail in this paper. Also, the two approaches *viz.*, the lexical (L approach) and the functional (F approach) of POS tagging have been discussed and presented with examples. Further, some ambiguous cases in POS annotation are presented in the paper.

**Keywords:** Marathi, POS Annotation, POS Tagging, Lexical, Functional, Marathi POS Tagset, ILCI

## 1 Introduction

In any natural language, Part of Speech (POS) such as noun, pronoun, adjective, verb, adverb, demonstrative, etc., forms an integral building block of a sentence structure. POS tagging<sup>1</sup> is one of the major activities in Natural Language Processing (NLP). In corpus linguistics, POS tagging is the process of marking/annotating a word in a text/corpus which corresponds to a particular POS. The annotation is done based on its definition and its context *i.e.*, its relationship with adjacent and related words in a phrase, sentence, or paragraph. A simplified form of this is commonly taught to school-age children, in the identification of words as nouns, verbs, adjectives, adverbs, etc. The term ‘Part-of-Speech Tagging’ is also known as POS tagging, POST, POS annotation, grammatical tagging or word-category disambiguation.

In this paper, the challenges and issues in POS tagging with special reference to Marathi<sup>2</sup> have been presented. The Marathi language is one of the major languages of India. It belongs to the Indo-Aryan Language family with about 71,936,894 users<sup>3</sup>. It is predominantly spoken in the state of Maharashtra in Western India (Chaudhari *et al.*, 2017). In recent years many research institutions and organizations are involved in developing the lexical resources for Marathi for NLP activities. Marathi Wordnet is one such lexical resource developed at IIT Bombay (Popale and Bhattacharyya, 2017).

The paper is organized as follows: Section 2 introduces POS annotation; section 3 provides information on Marathi annotated corpora; section 4 describes Marathi tag set; section 5 explains tagging approaches, section 6 presents ambiguous behaviors of the Marathi words, section 7 presents a discussion on special cases, and section 8 concludes the paper with future work.

## 2 Parts-Of-Speech Annotation

In NLP pipeline POS tagging is an important activity which forms the base of various language processing applications. Annotating a text with POS tags is a standard low-level text pre-processing step before moving to higher levels in the pipeline like chunking, dependency parsing, etc. (Bhattacharyya, 2015). Identification of the parts of speech such as nouns, verbs, adjectives, adverbs for each word (token) of the sentence helps in analyzing the role of each word in a sentence (Jurafsky D. *et al.*, 2016). It represents a token level annotation wherein it assigns a token with POS category.

## 3 Marathi Annotated Corpora

Aim of POS tagging is to create a large annotated corpora for natural language processing, speech recognition and other related applications. Annotated corpora serve as an important resource in NLP activities. It proves to be a basic building block for constructing statistical models for the automatic processing of natural languages. The significance of large annotated corpora is widely appreciated by researchers and application developers. Various research institutes in India *viz.*, IIT Bombay<sup>4</sup>, IIIT Hyderabad<sup>5</sup>, JNU New Delhi<sup>6</sup>, and other institutes have developed a large corpus of POS tagged data. In Marathi, there is around 100k annotated data developed as a part of Indian Languages Corpora Initiative (ILCI)<sup>7</sup> project funded by MeitY<sup>8</sup>, New Delhi. This ILCI corpus consists of four domains *viz.*, Tourism, Health, Agriculture, and Entertainment. This tagged data (Tourism - 25K, Health - 25K, Agriculture - 10K, Entertainment - 10K, General - 30K) is used for various applications like chunking, dependency tree banking, word sense disambiguation, etc. This ILCI

<sup>1</sup><http://language.worldofcomputing.net/pos-tagging/parts-of-speech-tagging.html>

<sup>2</sup><http://www.indianmirror.com/languages/marathi-language.html>

<sup>3</sup><http://www.censusindia.gov.in/>

<sup>4</sup><http://www.cfilt.iitb.ac.in/>

<sup>5</sup><https://ltrc.iiit.ac.in/>

<sup>6</sup><https://www.jnu.ac.in/>

<sup>7</sup><http://sanskrit.jnu.ac.in/ilci/index.jsp>

<sup>8</sup><https://meity.gov.in/>

annotated data forms a baseline for Marathi POS tagging and is available for download at TDIL portal<sup>9</sup>.

#### 4 The Marathi POS Tag-Set

The Bureau of Indian Standards (BIS)<sup>10</sup> has come up with a standard set of tags for annotating data for Indian languages. This tag-set is prepared for Hindi under the guidance of BIS. The BIS tag-set aims to ensure standardization in the POS tagging across the Indian languages. The tag sets of all Indian languages have been drafted by Dept. of Information Technology, MeitY and presented as Unified POS standard in Indian languages<sup>11</sup>. Marathi POS tag-set has been prepared at IIT Bombay referring to the standard BIS POS Tag-set, IIIT Hyderabad guideline document (Bharati et al, 2006) and Konkani Tag-set (Vaz et. al., 2012). This Marathi POS Tag-set can be seen in Appendix A.

#### 5 Lexical and Functional POS Tagging: Challenges and Discussions

Lexical POS tagging (Lexical or L approach) deals with tagging of a word at a token level. Functional POS tagging (Functional or F approach) deals with tagging of a word as a syntactic function of a word in a sentence. In other words, a word can have two roles viz., grammatical role (lexical POS w.r.t. a dictionary entry) and functional role (contextual POS)<sup>12</sup>. For example, in the phrase ‘golf stick’, the POS tag of the word ‘golf’ could be determined as follows:

- Lexically it is a noun as per lexicon.
- Functionally it is an adjective as it is a modifier of succeeding noun.

In the initial stage of ILCI data annotation, POS tagging was conducted using the lexical approach. However, over a while, POS tagging was done using the functional approach only. The reason is that, by using the lexical approach we do a general tagging, i.e., tagging at a surface level or token level and by using the functional approach we do a specific tagging, i.e., tagging at a semantic level. This eases the annotation process of chunking and parsing in the NLP pipeline.

While performing POS annotation, many issues and challenges were encountered, some of which are discussed below. Table 1 lists the occurrences of discussed words in the ILCI corpus.

##### 5.1 Subordinators which act as Adverbs

There are three basic types of adverbs. They are time (N\_NST), place (N\_NST) and manner (RB). Traditionally, adverbs should be tagged as RB. Subordinators are conjunctions which are tagged as CCS.

However, there are some subordinators in Marathi which act as adverbs. For example, ज्याप्रमाणे (*tyApramANE*, like-

wise), त्याप्रमाणे (*tyApramANE*, like that), ह्याप्रमाणे (*hyApramANE*, like this), जेव्हा (*jevha*, when) and तेव्हा (*tevhA*, then). ज्याप्रमाणे (*tyApramANE*) and ह्याप्रमाणे (*hyApramANE*) are generated from pronominal stems viz., ज्या (*tyA*) and ह्या (*hyA*) hence they are lexically qualified as pronouns, however, they function as adverbs; hence to be functionally tagged as RB at the individual level. However, when these words appear as part of the clause then they should be functionally tagged as CCS.

This distinction was also observed by noted Marathi grammarian, Damle<sup>13</sup> (Damle, 1965) [p. 206-07].

##### 5.2 Words with Suffixes

There are suffixes like मुळे (*muLe*, because of; due to), साठी (*sAThI*, for), बरोबर, (*barobara*, along with), etc. When these suffixes are attached to pronouns they function as adverbs or conjunctions at a syntactic level. For example, words त्यामुळे (*tyAmuLe*, because of that), यामुळे (*yAmuLe*, because of this), यासाठी (*yAsAThI*, for this), ह्यामुळे (*hyAmuLe*, because of this), ह्याच्यामुळे (*hyAchyAmuLe*, because of it/him), यांच्यामुळे (*yAMchyAmuLe*, because of them), त्याचबरोबर (=तसेच) (*tyAchabarobara* (=tasecha), further) are formed by attaching the above suffixes to pronouns. These words which are formed are lexically tagged as PRP. However, functionally these words act as conjunctions at the sentence level; therefore, they should be tagged as CCD. Also, consider the words त्यावेळी (*tyAveLI*, at that time), ह्यावेळी (*hyAveLI*, at this time), ह्यानंतर (*hyAnaMtara*, after this), त्यानंतर (*tyAnaMtara*, after that). Here, wherever the first string/morpheme appears as त्या (*tyA*) and ह्या (*hyA*), the tag should be given as PRP, lexically. But functionally, all these words shall be tagged as N\_NST (time adverb).

##### 5.3 Words which are Adjectives

Adjectives are tagged as JJ. Consider the example below: त्याच्यामध्ये ही कला परंपरागत चालत आली आहे (*tyAchyAmadhya hi kala paraMparAgata chAlataAlIAhe*, this art has come to him by tradition). Lexically, the word परंपरागत (*paraMparAgata*, traditional) is an adjective, but, in the above sentence, it qualifies the verb चालत येणे (*chAlatayeNe*, to be practiced). Hence functionally, the word परंपरागत (*paraMparAgata*) should be tagged as an RB. Similarly, a word वाईट (*vAITa*, bad) has a lexical POS as an adjective (Date-Karve, 1932). But in the sentence मला वाईट वाटते (*mala vAITa vATate*, I am feeling bad), it functions as an adverb, as it is qualifying the verb and not preceding the pronoun मला (*mala*, I; me). Therefore, functionally word वाईट (*vAITa*) acts as adverb hence should be tagged as RB.

##### 5.4 Adnominal Suffixes Attached to Verbs

The adnominal suffix जोग (*jogaM*) and all its forms (जोगा, जोगी, जोगे, जोग्या; *jogA, jogI, joge, jogyA*) are always attached to verbs. For example, word करण्याजोग्या (*karaNyA-jogyA*, doable) is lexically tagged as a verb. However, word करण्या (*karaNyA*) is a Kridanta form of a verb करणे (*karaNe*, to do) and suffix जोग (*jogaM*) is an adnominal suffix attached to Kridanta form; hence, a verb with all the

<sup>9</sup> <https://www.tdil-dc.in/>

<sup>10</sup> <http://www.bis.gov.in/>

<sup>11</sup> [http://tdil-dc.in/tdildcMain/articles/134692Draft POS Tag standard.pdf](http://tdil-dc.in/tdildcMain/articles/134692Draft%20POS%20Tag%20standard.pdf)

<sup>12</sup> <https://www.cicling.org/2018/cicling18-hanoi-special-event-23mar18.pdf>

<sup>13</sup> <http://www.cfilt.iitb.ac.in/damale/index.html>

forms of जोगं (*jogaM*) should functionally be treated as adjectives. Therefore verbs with adnominal suffix should be tagged as JJ.

### 5.5 Words जसे (*jase*) तसे (*tase*)

As per Damle (1956), words जसे (*jase*, like this) and तसे (*tase*, like that) are tagged as adverbs. However, if they appear with nouns in a sentence, they are influenced by the inflection and gender property of that nominal stem. For example, words जसे (*jase*, like this) and तसे (*tase*, like that) have inflected forms like जसा (*jasA*, like him), जशी (*jashI*, like her), तसा (*tasA*, like this), तशी (*tashI*, like this), तसे (*tase*, like this), etc. All these words function as a relative pronoun in a sentence. Hence, the words and their variations should be functionally tagged as PRL.

### 5.6 Word तसेतर (*tasetara*)

A word तसेतर (*tasetara*, as it is seen) is the same as तसे पाहिले तर (*tase pAhile tara*, as it is seen). Lexically, it can be tagged as a particle (RPD) but since it has a function of conjunction; it should be tagged as CCD. For example, in a sentence तसेतर तणावामुळेही काळी वलय येतात (*tasetara taNAvAmuLehI kALI valaya yetAta*, as it is seen that black circles appear because of stress as well), word तसेतर (*tasetara*) functions as conjunction and hence should be tagged as CCD instead of tagging it as RPD.

### 5.7 Word अन्यथा (*anyathA*)

The standard dictionaries give POS of the word अन्यथा (*anyathA*, otherwise; else; or) as an adverb/indeclinable. For example, consider a sentence अन्यथा तो येणार नाही (*anyathA to yeNArA nAhI*, Otherwise, he will not come). Here, while annotating अन्यथा (*anyathA*) there is a possibility that annotator can directly tag this word as an adverb at a lexical level. However, it behaves like conjunction at the sentence level and hence it should be tagged as CCD.

### 5.8 Different Forms of कसा (*kasA*)

As per BIS Tag-set, words कसा, कशी, कसे (*kasA, kashI, kase*; how) shall be tagged as PRQ. However, the PRQ tag is only for pronoun category and the word कसा (*kasA*) is not a pronoun; it can behave as an adverb or as a modifier. Consider the examples below:

1. तो माणूस कसा आहे हे त्याच्याशी बोलल्यावरच कळेल (to mANUsA kasA Ahe he tyAchyAshI bolalyA-varacha kaLela, we will come to know about him only after talking to him) [adnominal]
2. सरकारी ठरावाने कायद्याचे कलम कसे रद्द होणार (sarakArI TharAvAne kAyadyAche kalama kase radda hoNARA, How can this clause of law be prohibited by Government Resolution?) [adverbial]

In the 1st case, word कसा (*kasA*, how) functionally acts as a pronoun, hence to be tagged as PRQ. While, in the 2nd case, it acts as an adverb, hence to be functionally tagged as RB.

## 5.9 Word मात्र (*mAtra*)

A word मात्र (*mAtra*) is very ambiguous in its various usages; it is difficult to functionally identify the POS of this word at a sentence level. Various meanings of word मात्र (*mAtra*) are given in Data-Karve dictionary<sup>14</sup>. Some of the different senses of मात्र (*mAtra*) are discussed here:

- When the word मात्र (*mAtra*) conveys the meaning of ही, देखील, सुद्धा (*hl, dekhlla, suddhA*; also) then it should be tagged as RB functionally.
- When a word is related to the preceding word तेथे (*tethe*, there) and its function is an emphatic marker च (*cha*) then it should be tagged as RPD functionally.
- When word मात्र (*mAtra*) appears in the form of conjunction then it should be marked as CC functionally.
- If the word is modifying the succeeding noun, then it should be tagged as JJ functionally.
- If the word is modifying the preceding word, then the tag will be RPD as a particle functionally.

Therefore, it is noticed that the word मात्र (*mAtra*) does not have one single POS tag functionally and it depends upon the appearance in a sentence. Hence, should be tagged as per the usage.

Token	Lexical	Functional	Occurrences
ज्याप्रमाणे	Pronoun	Adverb	94
त्याप्रमाणे	Pronoun	Adverb	180
ह्याप्रमाणे	Pronoun	Adverb	8
जेव्हा	Subordinator	Time adverb	1496
तेव्हा	Subordinator	Time adverb	1577
मात्र	Adverb	Post-position Conjunction Particle	426
तसेतर	Particle	Conjunction	8
कसा	Wh-word	Adverb	269
त्यामुळे	Pronoun	Conjunction	530
ह्यामुळे	Pronoun	Conjunction	424
ह्याच्यामुळे	Pronoun	Conjunction	8
त्यावेळी	Pronoun	Time adverb	244
ह्यावेळी	Pronoun	Time adverb	33
ह्यानंतर	Pronoun	Time adverb	71
त्यानंतर	Pronoun	Time adverb	298
परंपरागत	Adjective	Adverb	104
वाईट	Adjective	Adverb	246
अन्यथा	Adverb	Conjunction	24
जसे	Relative pronoun	Adverb	1007
तसे	Relative pronoun	Adverb	511
करण्यजोग्या	Verb	Adjective	97

Table 1: Occurrences of discussed words and lexical v/s functional tags assigned to these words

## 6 POS Ambiguity: Challenges and Discussions

Ambiguity is a major open problem in NLP. Several POS level ambiguity issues were faced by annotators while annotating the Marathi corpus. Following are some POS

<sup>14</sup><http://www.transliterator.org/dictionary/mr.kosh.maharashtra/source>

specific ambiguity problems encountered while annotating.

### 6.1 Ambiguous POS: Adjective or Noun?

Examples: वयस्कर (*vayaskara*, the aged)

- कुटुंबाच्या वयस्कर सदस्यांनी मतदान केले (*kuTuMbAchyA vayaskara sadasyAMnI matadAnakele*, all the aged members of the family voted).
- सर्व वयस्करांनी मतदान केले (*sarva vayaskarAMnI matadAna kele*, all the aged people voted).

In the above examples, the word वयस्कर (*vayaskarAMnI*) lexically acts as an adjective as well as a noun. However, at the syntactic level, in the first example, it is functioning as adjective hence to be tagged as JJ, while in the second example it is functioning as a noun hence to be tagged as N\_NN. This is one of the challenges while annotating adjectives appearing in nominal form. Annotators usually fail to disambiguate these types of words at the lexical level; therefore such words should be disambiguated at syntactic level. Hence, annotators need to take special care while annotating such cases.

### 6.2 Ambiguous POS: Demonstrators

While annotating demonstrators such as हा, ही, हे, तो, ती, ते (*(hA, hI, he)*, this), (*(to, tI, te)*, that) annotators often get confused whether to tag them as DMD or DMR. Simple guideline can be followed is, if the demonstrator is directly following noun, then tag it as DMD, otherwise tag it as DMR i.e., if the demonstrator is referring to previous noun/person.

### 6.3 Ambiguous POS: Noun and Conjunction

Example: word कारण (*kArANa*, reason; because). At semantic level, the word कारण (*kArANa*) has two meanings, one is 'a reason' which acts as a noun and another is 'because' which acts as a conjunction. Annotators have to pay special attention while tagging such cases.

### 6.4 Ambiguous Words: ते (*te*) and तेही (*tehI*)

The word ते (*te*) has different grammatical categories like pronoun (they), demonstrator (that) and conjunction (to). Examples:

- ३० ते ४० (*30 te 40*, 30 to 40)
- The word ते (*te*) lexically and functionally acts as conjunction, hence to be tagged as CCD.
- ते म्हणाले (*te mhaNAle*, they said)
- Here word ते (*te*) acts as personal pronoun, hence to be tagged as PR\_PRP
- ते कुठे आहेत? (*te kuThe Aheta?*, where are they?)
- Here word ते (*te*) acts as relative demonstrator, hence to be tagged as DM\_DMR
- राकेशने पोलीसांना फोन केला आणि ते दोन्ही चोर पकडले गेले (*rAkesane poliIsAMnA phona kela ANi te donhI chora pakaDale gele*, Rakesh called police and those two thieves got caught).
- Here, word ते (*te*) is modifying its succeeding noun चोर (*chora*, thief) so it is Deictic demonstrator, hence to be tagged as DM\_DMD.

- त्यांना हे कधीच पसंत नव्हते, त्यांच्या मुलाने संगीत शिकावे आणि तेही नृत्य (*tyAMnA he kadhIchapasaMtanavhate, tyAMchyAmulAnesaMgItashikAveANitehInRRitya*, He never wanted his son to learn music and that too the dance form)

Here, the word तेही (*tehI*) is an ambiguous word. It is modifying succeeding noun or previous context. Here, ही (*hI*) is a bound morpheme and conveys the meaning 'also'. Therefore word तेही (*tehI*) should be tagged as DM\_DMR.

### 6.5 Ambiguous word: उलटा (*ulaTA*)

Examples:

- उलटे टांगून सुकवले जाते (*ulaTe TAMgUna sukavale jAte*). Here, उलटे (*ulaTe*, upside down is behaving as manner, not a noun, hence to be tagged as RB.
- उलटे भांडे सुलटे कर (*ulaTe bhAMDe sulaTe kara*). Here उलटे (*ulaTe*) it is modifying succeeding noun, hence it is an adjective, hence to be tagged as JJ.

In the above examples, annotator should identify word behavior in the sentence and tag accordingly.

### 6.6 Ambiguous words: कितीही (*kitIhI*), ना का (*nA kA*) and असू दे ना का (*asU de nA kA*)

Examples:

- संगणक हा कितीही प्रगत किंवा चतुर असू दे ना का, तो केवळ तेच काम करू शकतो ज्याची विधी (पद्धत) आपल्याला स्वतः माहित आहे. (*saMgaNaka hA kitIhI pragata kiMvA chatura asU de nA kA, to kevalA techa kAma karU shakato jyAchi vidhi (paddhata) ApalyALA svata: mAhita Ahe*, The computer how much ever may be advanced and clever, it only does that work whose method we only know). Here, कितीही (*kitIhI*, how much) is a quantifier, hence to be tagged as QTF.
- In the phrase असू दे ना का (*asU de nA kA*), the token ना (*nA*) is a part of verb असू दे (*asU de*, let it be) and should be tagged as VM, hence the phrase should be tagged as VM, while the token का (*kA*) is acting as a particle in this phrase and not as a question marker, therefore का (*kA*) should be tagged as RPD.
- किती माणसे जेवायला होती? (*kitI mANase jevAyala hotI*, how many people were there for a meal?). Here, किती (*kitI*, how many) is a question so it should be tagged as DMQ.

### 6.7 Ambiguous word: तर (*tara*)

Examples:

- Conjunction: जर मी वेळीच गेलो नसतो तर हा वाचला नसता (*jara mI veLicha gelo nasato tara hA vAchalA nasatA*, if I had not gone on time he would have not survived).
- Particle: 'हो! आता मी जातो तर!' = 'मी अजिबात जाणार नाही' (*'ho! AtA mI jAto tara!' = 'mI ajibAta jANA-*

*ra nAhl*’, ‘yes! now I am leaving then’ = ‘I am not at all leaving’).

In the above sentences, the word तर (*tara*) is used as a supplementary or stressable word so somewhat special as to give meaning in the sentence. (Date-Karve, 1932). Hence it should be treated as CCD.

- तुम्ही तर लाख रुपये मागतां व मी तर केवळ गरीब पडलो (*tumhI tara lAkha rupaye mAgatAM va mI tara kevaLa garIba paDalo*, you are asking for lakh rupees and I am a poor person). In this sentence, the word तर (*tara*) indicates opposition with respect to meaning between two connected sentences. (Date-Karve, 1932). Hence, it should be treated as a RPD.

## 7 Discussions on Some Special Cases

- Words अम्लयुक्त (*Amlayukta*), मलईरहित (*malairahita*), मेदरहित (*medarahita*), दुष्काळग्रस्त (*duShkaLagrasta*) are combinations of noun plus adjective suffix such as युक्त (*yukta*), ग्रस्त (*grasta*) and रहित (*rahita*). In such cases, even though noun is a head string and adjective part is a suffix, the whole word shall be tagged as JJ.
- Before tagging अभंग (*abhaMga*, verses), ओव्या (*ovyA*, stanzas), काव्य (*kAvya*, poetry), etc., annotator shall first read between the lines; understand the meaning which it conveys and then decide upon the grammatical categories of each token. For example, in sentence कळवे तयासी कळे अंतरीचे कारण ते साचे साच अंगी (*kaLAve tayAsI kale aMtarIche kArAna te sAche sAcha aMgI*) the POS tagging should be done as साचे\V\_VM साच\N\_NN अंगी\N\_NN, etc.
- Doubtful cases of word कोणता (*koNatA*)

Examples:

- कोणता मुलगा हुशार आहे (*koNatA mulagA hu-shAra Ahe*)?
- वाहतुकीच्या दरम्यान कोणतीही हानी झालेली नाही (*vAhatukiChyA daramyAna koNatIhI hAnI jhAlell nAhl*).
- ह्यांच्या बोलण्याचा माझ्यावर कोणताही परिणाम झाला नाही (*hyAMchyA bolANyAchA mAjhyAvara koNatAhlI pariNAMA jhAlA nAhl*).
- शेतकऱ्यास कोणत्याही वर्षी पाण्याची कमतरता भासणार नाही (*shetakaryAsa koNatyAhlI var-ShI pANyAchI kamataratA bhAsaNARA nAhl*).

Here, in the 1st example, the word कोणता (*koNatA*, which one) undoubtedly is DMQ. In rest of the examples कोणतीही (*koNatAhlI*, whichever, whomever), कोणताही (*koNatAhlI*, whichever,

whomever), कोणत्याही (*koNatyAhlI*, whichever, whomever) are DM adjective (DMD).

## 8 Conclusion and Future Work

Marathi POS tagging is an important activity for NLP tasks. While tagging, several challenges and issues were encountered. In this paper, Marathi BIS tag-set has been discussed. Lexical and functional tagging approaches were discussed with examples. Further, various challenges, experiences, and special cases have been presented. The issues discussed here will be helpful for annotators, researchers, language learners, etc. of Marathi and other languages.

In future, more issues such as tagging for words having multiple senses; words having multiple functional tags will be discussed. Also, tagset comparison of close languages will be done. Further, the evaluation of lexical and functional tagging using statistical analysis will be done.

## References

- Akshar Bharati, Dipti Misra Sharma, Lakshmi Bai, Rajeev Sangal. (2006). AnnCorra : Annotating Corpora Guidelines For POS And Chunk Annotation For Indian Languages. *Language Technologies Research Centre, IIIT, Hyderabad*.
- Chitra V. Chaudhari, Ashwini V. Khaire, Rashmi R. Mur-tadak, Komal S. Sirsulla. (2017). Sentiment Analysis in Marathi using Marathi WordNet. *Imperial Journal of Interdisciplinary Research (IJIR)* Vol-3, Issue-4, 2017 ISSN: 2454-1362.
- Damle, Moro Keshav. (1965). Shastriya Marathi Vyakran. *A scientific grammar of Marathi*, 3<sup>rd</sup> edition. Pune, India: RD Yande.
- Daniel Jurafsky & James H. Martin. (2016). *Speech and Language Processing*.
- Edna Vaz, Shantaram V. Walawalikar, Dr. Jyoti Pawar, Dr. Madhavi Sardesai. (2012). BIS Annotation Standards With Reference to Konkani Language. *24<sup>th</sup> International Conference on Computational Linguistics (COLING 2012)*, Mumbai.
- Lata Popale and Pushpak Bhattacharyya. (2017). Creating Marathi WordNet. *The WordNet in Indian Languages*. Springer, Singapore, 2017. 147-166.
- Pushpak Bhattacharyya, (2015). *Machine Translation, Book published by CRC Press, Taylor and Francis Group, USA*.
- Yashwant Ramkrishna Date, Chintman Ganesh Karve, Aba Chandorkar, Chintaman Shankar Datar. (1932). *Maharashtra Shabdakosh*. Published by H. A. Bhave, Varada Books, Senapati Bapat Marg, Pune.

## Appendix A

### Marathi Parts of Speech Tag-Set with Examples

SI. No	Category	Label	Annotation Convention **	Examples
<b>Top level &amp; Subtype</b>				
1	<b>Noun (नाम)</b>	<b>N</b>	<b>N</b>	
1.1	Common (जातीवाचक नाम)	NN	N_NN	गाय\N_NN गोठ्यात\N_NN राहते.
1.2	Proper (व्यक्तीवाचक नाम)	NNP	N_NNP	रामाने\N_NNP रावणाला\N_NNP मारले.
1.3	Nloc (स्थल-काल)	NST	N_NST	1. तो येथे\N_NST काम करत होता. 2. त्याने ही वस्तू खाली\N_NST ठेवली आहे.
2	<b>Pronoun (सर्वनाम)</b>	<b>PR</b>	<b>PR</b>	
2.1	Personal (पुरुष वाचक)	PRP	PR_PRP	मी\PR_PRP येतो.
2.2	Reflexive (आत्म वाचक)	PRF	PR_PRF	मी स्वतः\PR_PRF आलो.
2.3	Relative (संबंधी)	PRL	PR_PRL	ज्याने\PR_PRL हे सांगितले त्याने हे काम केले पाहिजे.
2.4	Reciprocal (पारस्परिक)	PRC	PR_PRC	परस्पर
2.5	Wh-word (प्रश्नार्थक)	PRQ	PR_PRQ	कोण\PR_PRQ येत आहे?
2.6	Indefinite (अनिश्चित)	PRI	PR_PRI	कोणी\PR_PRI कोणास\PR_PRI हासू नये. त्या पेटीत काय\PR_PRI आहे ते सांगा.
3	<b>Demonstrative (दर्शक)</b>	<b>DM</b>	<b>DM</b>	
3.1	Deictic	DMD	DM_DMD	हे पुस्तक माझे आहे. तो\DM_DMD मुलगा हुशार आहे. हा\DM_DMD मुलगा हुशार आहे. ही\DM_DMD मुलगी सुंदर आहे. जेथे\DM_DMD राम होता तेथे\DM_DMD तो होता.
3.2	Relative	DMR	DM_DMR	हे\DM_DMR लाल रंगाचे असते.
3.3	Wh-word	DMQ	DM_DMQ	कोणता\DM_DMQ मुलगा हुशार आहे?
4	<b>Verb (क्रियापद)</b>	<b>V</b>	<b>V</b>	
4.1	Main (मुख्य क्रियापद)	VM	V_VM	तो घरी गेला\V_VM.
4.2	Auxiliary (सहाय्यक क्रियापद)	VAU X	V_VAUX	राम घरी जात आहे\V_VAUX.
5	<b>Adjective (विशेषण)</b>	<b>JJ</b>		सुंदर\JJ मुलगी
6	<b>Adverb (क्रियाविशेषण)</b>	<b>RB</b>		हळूहळू\RB चाल.
7	<b>Conjunction (उभयान्वयी अव्यय)</b>	<b>CC</b>	<b>CC</b>	
7.1	Coordinator	CCD	CC_CCD	तो आणि\CC_CCD मी.
7.2	Subordinator	CCS	CC_CCS	जर\CC_CCS त्याने सांगितले असते तर\CC_CCS हे काम मी केले असते.
7.2.1	Quotative	UT	CC_CCS_UT T	असे\CC_CCS_UT म्हणून\CC_CCS_UT तो पुढे गेला.
8	<b>Particles</b>	<b>RP</b>	<b>RP</b>	
8.1	Default	RPD	RP_RPD	मी तर\RP_RPD खूप दमले.
8.2	Interjection (उद्गार वाचक)	INJ	RP_INJ	अरेरे\RP_INJ ! सचिनची विकेट दापली.
8.3	Intensifier (तीव्र वाचक)	INTF	RP_INTF	राम खूप\RP_INTF चांगला मुलगा आहे.
8.4	Negation (नकारात्मक)	NEG	RP_NEG	नको, न
9	<b>Quantifiers</b>	<b>QT</b>	<b>QT</b>	
9.1	General	QTF	QT_QTF	थोडी\QT_QTF साखर द्या.
9.2	Cardinals	QTC	QT_QTC	मला एक\QT_QTC गोळी दे.
9.3	Ordinals	QTO	QT_QTO	माझा पहिला\QT_QTO क्रमांक आला.
10	<b>Residuals (उर्वरित)</b>	<b>RD</b>	<b>RD</b>	
10.1	Foreign word	RDF	RD_RDF	
10.2	Symbol	SYM	RD_SYE	\$, &, *, (, ),
10.3	Punctuation	PUN C	RD_PUNC	. (period), ,(comma), :(semi-colon), !(exclamation),?(question), :(colon), etc.
10.4	Unknown	UNK	RD_UNK	Not able to identify the Tag.
10.5	Echo-words	ECH	RD_ECH	जेवण बिवण, डोके बिके