

Chapter III

P

art of Speech Tagset for Kashmiri

3.0 Introduction

Tagset development forms a foundation of any computational processing endeavor. It is generally accepted that, as a prelude to syntactic analysis of natural language by computers, a text must be annotated with tags indicating the POS. The first pre-requisite for automated POS tagging is a tagset that is a set of exhaustive categories into which any token of the language can be placed. While the nature of the language is that there will always be words that are hard to classify, or are ambiguous between two categories, the tagset categories should be designed in such a way so as to minimize such problems. The fundamental problems in POS tagging task stem from the fact that a word can take different lexical categories depending on its context. The tagger has to resolve this ambiguity and determine the best sequence for a sentence. POS tagging also known as morphosyntactic categorization or syntactic word class tagging (Halteren 1999) is the process of assigning a part of speech or other lexical class marker to each word in a corpus. Tags are also applied to punctuation markers, thus tagging for natural language is the same process as tokenization for computer languages, although tags for natural languages are much more ambiguous.

Part of speech tagging has been studied extensively in the past two decades and lot of work has been done in various European languages including many Indian languages like Hindi, Urdu, Sanskrit, Tamil and Kannada. Comparatively speaking, little work has been done in Kashmiri in this field. The only work done in Kashmiri so far in this domain is the development of a flat tagset following

ILMT(Indian Language Machine Translation) guidelines under the project “Development of Language Tools and Linguistic Resources for Kashmiri” at the Department of Linguistics, University of Kashmir.

3.1 General Framework for Kashmiri POS Tagging

For designing a Kashmiri tagset, apart from following the Eagles Guidelines and the Penn tree bank tagset, many other Indian tagging guidelines like IL-POST, ILMT and Sanskrit tagset were taken into consideration. The tagging schema for Kashmiri is designed taking into consideration both language features in general and the idiosyncratic features of Kashmiri. After careful consideration a hierarchical tagset was favored. The whole design of the tagset developed so far revolves around three distinct features into which the grammatical schema is distributed. The features are:

I. Category

II. Type

III. Attribute

Categories involve major grammatical categories like Nouns, Verbs etc. The type includes the type of those grammatical categories like Common Noun and Proper Noun for Noun category, Main Verb, Auxiliary Verb etc. for Verb category, and so on. The attribute level takes features within each type like Gender (masculine, feminine), Number (singular, plural), Case (dative, ergative, ablative, etc.), Tense and Aspect etc. into consideration. The category list includes all Kashmiri categories that can occur. The type list within a category includes all

types of the category that can occur. The attribute list includes all possible attributes of the type that can occur. A detailed account of the schema is given below:

Category	Types	Attributes
NOUN	Common Noun	Gender, Number, Case, Emphatic ¹ , Vocative ²
	Proper Noun	Gender, Number, Case, Emphatic
PRONOUN	Personal Pronoun	Gender, Number, Person, Distance ³ , Case, Emphatic, Honorific ⁴
	Demonstrative Pronoun	Gender, Number, Person, Distance, Case, Emphatic, Honorific
	Reflexive Pronoun	Gender, Number, Case, Emphatic
	Reciprocal Pronoun	Gender, Number, Case, Emphatic
	Relative Pronoun	Gender, Number, Case, Emphatic, Honorific
	Possessive Pronoun	Gender, Number, Person, Distance, Case, Emphatic, Honorific
	Indefinite Pronoun	Gender, Number, Case, Emphatic,
	Interrogative Pronoun	Gender, Number, Case, Honorific
VERB	Main Verb	Gender, Number, Person, Tense, Aspect, Finiteness, Honorific

¹ In the present tagset Emphatic is depicted both as a Type and as an attribute. As a Type it occurs as a free morpheme, and as an attribute it occurs as a bound morpheme as shown above.

² Like Emphatic, Vocative is also used as an attribute in as well as a Category.

³ 3rd Person Pronouns exhibit a two term distinction in terms of distance i.e. Proximate and Remote wherein Remote is further divided into Remote (within sight) and Remote (out of sight).

⁴ Honorific can be used as an Attribute as in Noun and Pronoun categories shown above, and also as a Type of the main category where it occurs as a free morpheme.

	Auxiliary ⁵	Gender, Number, Person, Tense, Aspect, Finiteness
	Conjunctive Verbs	Gender, Number, Finiteness
	Causative Verbs	Gender, Number, Tense, Aspect, Finiteness, First place, Second Place, Third Place
ADJECTIVE	Qualitative	Gender, Number, Case, Emphatic
	Quantitative	Number, Case, Emphatic
QUANTIFIER	Cardinal	Number, Case, Emphatic
	Ordinal	Number, Case, Emphatic
GERUND	—	—
ADVERB	Time	—
	Place	—
	Manner	—
POSTPOSITION	—	Case ⁶
INTENSIFIER	—	—
INTERJECTION	—	—

⁵ Auxiliary verbs also include modal thus while tagging all modal will be tagged as auxiliary with its attributes.

⁶Postpositions in Kashmiri take dative and ablative cases with them.

PARTICLE	—	—
VOCATIVE	—	—
EQUATIVE	—	—
HONORIFIC	—	—
NEGATIVE	—	—
EMPHATIC	—	—
MULTI WORD EXPRESSIONS	Compound Words	
	Reduplication ⁷	
	Echo Words ⁸	
	Part of Word	
CONJUNCTION	Coordinating	—
	Subordinating	
PUNCTUATION	—	—

⁷ Attributes for reduplication will be same as that of the first word, that is if the first word is a Common Noun then it will take Common Noun attributes and if it is a Verb then verbal attributes are used.

⁸ Attributes for Echo words will be same as that of the first word like that of Reduplication. The first word is tagged according to the category which it possesses and the second word is tagged as an Echo word with 'ECH' tag.

SYMBOL	—	—
UNKNOWN WORDS	—	—

Table 3.1: Categories, Types and their Attributes used in the Proposed Kashmiri POS Tagset

3.2 Attribute Description

The table given below gives a detailed description of attributes which are used in the present tagset for annotating the corpus.

Attribute	Value
GENDER	Masculine(mas) Feminine(fem) Neuter(null) ⁹
NUMBER	Singular(sng) Plural(plu) Dual(null) ¹⁰
PERSON	I Person(Ip) II Person(IIp) III Person(IIIp)

⁹ Neuter gender is not found in Kashmiri that's why it is kept as null in the present tagset. However its presence in the present tagset is because if the same tagset is used for any other language in which neuter gender is present then the same will be activated while using that particular language in order to maintain the cross linguistic usage of the Tagset.

¹⁰There are no special dual, trial or quadral forms in Kashmiri for any of the Pronouns and the reason for its presence in the present tagset is same as that of Neuter gender.

DISTANCE	III Person Proximate(prx)	
	III Person	within sight (rmn)
	Remote	out of sight(rmf)
CASE	Nominative(nom) Dative(dat) Ergative(erg) Ablative(abl) Genitive(gen)	
TENSE	Present(prt) Past(pst) Future(fut)	
ASPECT	Simple(smp) Progressive(prg) perfective(prf)	
FINITENESS	Finite(fn), Non-finite(nfn)	
EMPHATIC	Emph	
CAUSATIVE ATTRIBUTES	Ipl, IIpl, IIIpl	
HONORIFICITY	honorific(hon)	
VOCATIVE	Voc	

Table 3.2: Attributes and their values

3.3 Tag Structure

Using the above mentioned categories, types and attributes, the following tag structure is used for annotating the corpus.

<tag_cat><tag_types><attribute_list>

where **<tag_cat>** is used for the Main Category, **<tag_types>** is used for the Types of the Main Category and **<attribute_list>** denotes the sub features of the main category. The present tagset is hierarchical in nature and in order to maintain that hierarchy, the main category tags for types and for attributes are placed between ‘<’ and ‘>’ in order to differentiate between the Main Category, Type and Attribute. Furthermore there should be no white space in between the tags of a single word.

Consider *kaT'an* ‘boys-dat’. The tag structure of this token will be

Example	Tag structure
<i>kaT'an</i>	<N><NC><mas><plu><dat>
‘boys-dat’	

Table 3.3: Example of Tag structure

The above given example of Tag structure is explained as

Category	Type	Attribute
Noun(N)	Common Noun(NC)	mas,sng,dat

Table 3.4: Schema of Tag structure

3.4 Tag Description

A detailed description of each tag used in Kashmiri tagset is given below:

3.4.1 Noun

Nouns in Kashmiri are broadly classified as Proper and Common Nouns and this distinction is grammatically significant as Proper Nouns count as a different category for ergative and possessive marking. Common Nouns include person, place or a thing, Abstract Noun include emotions, ideas etc, Collective Noun includes group of things, animals, or persons, Countable and Non-countable Nouns. As Nouns (Common and Proper) inflect for number gender and case, in the present tagging scheme Nouns are tagged by using different tags according to their gender- number distinction. Some relevant examples are given below:

NOUN	Example	Tag structure	
Common Noun	<i>Kul</i>	<N><NC><mas><sng><nom>	
	'tree'		
	<i>kul'-an</i>	<N><NC><mas><plu><dat>	
	'trees- dat'		
	<i>ko:r-i</i>	<N><NC><fem><sng><abl>	
	'girl-abl'		

	<i>ko:r'-an hund</i>	<N><NC><fem><plu><gen>	
	'girls-gen'		
Proper Noun	<i>parimehal</i>	<N><NP><mas><sng>	
	'Parimahal'		
	<i>parimehl-an</i>	<N><NP><mas><plu><dat>	
	'Parimahal-dat'		
	<i>Ifat</i>	<N><NP><fem><sng><nom>	
	'Ifat'		
	<i>ift-an</i>	<N><NP><fem><plu><dat>	
'Ifat-dat'			

Table 3.5: Tag structure of Nouns

3.4.2 Pronouns

All Pronouns are tagged separately. In Kashmiri Pronouns are inflected for number, gender and case like Nouns but are distinguished from Nouns by having a category of person also. Thus all Pronouns will be tagged with their gender, number and person features. Some examples are given below

PRONOUN	Example	Tag
Personal Pronoun	<i>bi</i>	<P><PRP><mas><sng><Ip><nom>
	'I'	
	<i>Tse</i>	<P><PRP><mas><sng><IIp><erg>
	'you-erg'	
	<i>hum-is</i>	<P><PRP><mas><sng><IIIp><rmn><dat>
	'he-dat'	
	<i>yim-an</i>	<P><PRP><mas><sng><IIIp><prx><dat>
	'these-erg'	
	<i>as-i</i>	<P><PRP><mas><plu><Ip><dat>
	'we-dat'	
	<i>yem'</i>	<P><PRP><mas><sng><IIIp><prx><erg>
'this-erg'		
<i>hum-av</i>	<P><PRP><mas><plu><IIIp><rmn><erg>	
'they-erg'		

<i>tim-an</i>	<P><PRP><mas><plu><IIIp><rmf><dat>
'they-dat'	
<i>Me</i>	<P><PRP><mas><sng><Ip><erg>
'I-erg'	
<i>Tse</i>	<P><PRP><mas><sng><IIp><erg>
'you-erg'	
<i>hɔ</i>	<P><PRP><fem><sng><IIIp><rmn><nom>
'she'	
<i>sɔ</i>	<P><PRP><fem><sng><IIIp><rmf><nom>
'she'	
<i>Asi</i>	<P><PRP><mas><plu><Ip><dat>
'we-dat'	
<i>tuh'</i>	<P><PRP><fem><plu><IIp><nom>
'you-plu/hon'	
<i>tim-an</i>	<P><PRP><fem><plu><IIIp><rmf><dat>

	'they-dat'	
	<i>hum-av</i>	<P><PRP><fem><plu><IIIp><rmn><erg>
	'they-erg'	

Table 3.6: Tag structure of Personal Pronouns

Pronoun	Example	Tag structure
Demonstrative Pronoun	<i>hum'</i>	<P><PDM><mas><sng><IIIp><rmn><erg>
	'he-erg'	
	<i>Su</i>	<P><PDM><mas><sng><IIIp><rmf><nom>
	'he'	
	<i>hɔ</i>	<P><PDM><fem><sng><IIIp><rmn><nom>
	'she'	
	<i>sɔ</i>	<P><PDM><fem><sng><IIIp><rmf><nom>
'she'		

Table 3.7: Tag structure of Demonstrative Pronouns

Pronoun	Example	Tag
Reciprocal	<i>akhəkis</i>	<P><PRC>

Pronoun	'to one another'	
	<i>pa:nlvə:n'</i>	
	'amongst each other'	

Table 3.8: Tag structure of Reciprocal Pronouns

Pronoun	Example	Tag	
Reflexive Pronoun	<i>panun</i>	<P><PRF><mas><sng><nom>	
	'self's'		
	<i>panIn'-an</i>	<P><PRF><mas><plu><dat>	
	'self's-dat'		
	<i>panIn'-i</i>	<P><PRF><fem><sng><dat>	
	'self's-dat'		
	<i>panIn'-av</i>	<P><PRF><fem><plu><erg>	
'self's-erg'			

Table 3.9: Tag structure of Reflexive Pronouns

Pronoun	Example	Tag structure
	<i>m'ə:n-is</i>	<P><PPO><mas><sng><Ip><dat>

Possessive Pronoun	‘my-dat’	
	<i>cə:n-is</i>	<P><PPO><mas><sng><IIp><dat>
	‘your-dat’	
	<i>hum'-sund</i>	<P><PPO><mas><sng><IIIp><rmn><gen>
	‘his-gen’	
	<i>təm'-sund</i>	<P><PPO><mas><sng><IIIp><rmf><gen>
	‘his-gen’	
	<i>sa:n'-av</i>	<P><PPO><mas><plu><Ip><erg>
	‘our-erg’	
	<i>tuhund-i</i>	<P><PPO><mas><plu><IIp><abl>
	‘you-abl’	
	<i>human-hund</i>	<P><PPO><mas><plu><IIIp><rmn><gen>
	‘their-gen’	
<i>timan- hund</i>	<P><PPO><mas><plu><IIIp><rmf><gen>	
‘their-gen’		

<i>m'ə:n'</i>	<P><PPO><fem><sng><Ip><nom>
'my'	
<i>ca:n'-av</i>	<P><PPO><fem><plu><IIp><erg>
'your-erg'	
<i>humi-sInz</i>	<P><PPO><fem><sng><IIIp><rmn><gen>
'her-gen'	
<i>yem'-sInz</i>	<P><PPO><mas><sng><IIIp><prx><gen>
'this-gen'	
<i>sa:n-i</i>	<P><PPO><mas><plu><Ip><abl>
'our-abl'	
<i>tuhInz-av</i>	<P><PPO><fem><plu><IIp><erg>
'your-erg'	
<i>yiman-hinzI</i>	<P><PPO><mas><plu><IIIp><prx><gen>
'their-gen'	

	<i>human-hInzI</i>	<P><PPO><mas><plu><IIIp><rmn><gen>
	‘their-gen’	

Table 3.10: Tag structure of Possessive Pronouns

Pronoun	Example	Tag structure	
Interrogative Pronoun	<i>Kus</i>	<P><PIT><mas><sng><nom>	
	‘who’		
	<i>kImI</i>	<P><PIT><fem><plu><nom>	
	‘who’		
	<i>kɔsI</i>	<P><PIT><fem><sng><nom>	
	‘who’		
	<i>kam-av</i>	<P><PIT><mas><plu><erg>	
	‘who-erg’		

Table 3.11: Tag structure of Interrogative Pronouns

Pronoun	Example	Tag structure
	<i>Yus</i>	<P><PPR><mas><sng><nom>
	‘who’	

Relative Pronoun	<i>Yim</i>	<P><PPR><mas><plu><prx><nom>	
	‘these’		
	<i>yɔsI</i>	<P><PPR><fem><sng><nom>	
	‘who’		
	<i>yimI</i>	<P><PPR><fem><plu><prx><nom>	
	‘these’		

Table 3.12: Tag structure of Relative Pronouns

Pronoun	Example	Tag
Indefinite Pronoun	<i>kē h</i>	<P><PID><nom>
	‘something’	
	<i>kāh</i>	
	‘someone’	

Table 3.13: Tag structure of Indefinite Pronouns

3.4.3 Verbs

Verbal constructions in languages, apart from one word, may be composed of multi word sequences as well. A multi word verb group sequence contains a main verb and one or more auxiliaries (V AUX ...). In the present work after a careful study of Kashmiri peculiarities it was decided that it would be more

appropriate to use different tags for the main verb and the auxiliary verbs. As per the classification of verbs, they are tagged as main verb, auxiliary, causatives, conjunct, etc. in the present tagset. The features include tense, number, gender and person.

Verbs	Example	Tag	
Main Verbs	<i>Pakun</i>	<V><VM><fut><nfn>	
	‘to walk’		
	<i>pokus</i>	<V><VM><mas><sng><Ip><pst><prf> ><fn>	
‘walked’			
Auxiliary Verbs	<i>chus</i>	<V><VAUX><mas><sng><Ip><prt>	
	‘is-M’		
	<i>o:sus</i>	<V><VAUX><mas><sng><Ip><pst>	
	‘was-M’		
<i>ə:sIkh</i>	<V><VAUX><fem><sng><Iip><pst>		
‘was-F’			

	<i>chukh</i>	<V><VAUX><mas><sng><Ipl><prt>
	‘are-M’	
Causative Verbs	<i>kh'a:vna:vun</i>	<V><VCUS><Ipl><nfn>
	‘make someone to eat’	
	<i>pakna:vun</i>	<V><VCUS><Ipl><nfn>
	‘make someone to walk’	
Conjunctive Verb¹¹	<i>salah d'un</i>	<N><NC><mas><sng><V><VCNJ><n
	‘to give advice’	fn>

Table 3.14: Tag structure of Verbs

3.4.4 Gerunds

Gerunds are formed by adding an infinitive suffix *-unto* to verbs. Gerunds are tagged as VBG, for example,

¹¹Conjunctive Verbs consists of a Noun lexical item and a Verb lexical item. Thus the first part is tagged as a Noun with its type and attributes and the second part is tagged as a Verb with its type and attributes as given above.

Gerunds	Tag
<i>Yun</i>	<VBG>
‘coming’	
<i>gevun</i>	
‘singing’	
<i>pakun</i>	
‘walking’	

Table 3.15: Tag structure of Gerunds

3.4.5 Adjective

Adjectives are broadly classified into qualitative and quantitative Adjectives. All modifiers of quality like different colors (*vɔzʊl* ‘red’, *n'u:l* ‘blue’, *saphe:d* ‘white’, etc.), personal qualities (*ca:la:k* ‘clever’, *da:na:* ‘wise’ *buzdil* ‘coward’, etc.), physical qualities (*thod* ‘tall’, *tshoT* ‘short’, *v'oTh* ‘fat’, *zə:v'ul* ‘slim’ etc.), qualities of taste (*modur* ‘sweet’, *tsok* ‘sour’, *ToTh* ‘bitter’, etc.) fall under Qualitative Adjectives.

Numerals (cardinals, ordinals, fractions, multiplicatives), intensifiers (*kāh* ‘some’, *sə:ri:* ‘all’, *kam* ‘little’), Demonstrative Adjectives (*yu:t* ‘this much’, *t'u:t* ‘that much’) etc. are included in Quantitative Adjectives. According to the description given above, Adjectives are tagged as follows:

Adjective	Example	Tag structure
Qualitative	<i>asll</i>	<JJ><QL>
	'good'	
Quantitative	<i>kam</i>	<JJ><QN>
	'less'	

Table 3.16: Tag structure of Adjectives

3.4.6 Quantifier

Quantifiers are tagged as cardinals and ordinals.

Quantifier	Example	Tag	
Cardinals	<i>akh</i>	<QTF><CRD>	
	'one'		
	<i>ək'</i>	<QTF><CRD><erg>	
	'one-erg'		
Ordinals	<i>əkim</i>	<QTF><ORD>	
	'first',		
	<i>doyim-is</i>	<QTF><ORD><dat>	
	'second-dat'		

Table 3.17: Tag structure of Quantifiers

3.4.7 Cases

Cases in Kashmiri are primarily divided into nominative, dative, ergative, ablative and genitive. The word classes like Nouns, Pronouns, Adjectives, Verbs and Postpositions inflect for cases, and thus it was decided to include cases in the attribute list. To include cases in the present schema for tagging, the following attributes were used.

Cases	Example	Tag
Dative	<i>shi:l-as</i>	<dat>
	'sheela-dat'	
Ergative	<i>shi:l-an</i>	<erg>
	'sheela-erg'	
Ablative	<i>saku:l-I</i>	<abl>
	'school-abl'	
Genitive	<i>shi:l-un</i>	<gen>
	'shila-gen'	

Table 3.18: Cases as Attributes

Genitive markers like *hund*, *hInd'*, *hInz* and *hInzI* (Group III) are tagged as Postposition in the present tagging scheme as they behave as Postpositions. The

genitive markers like *un*, *uk*, *Ik*, *Ic* and *Ici* (Group I and II) are tagged as attributes.

Some examples of nouns with and without explicit case markers are given below.

Example	Gloss		Example with cases	Gloss
<i>Kul</i>	'tree-S'	Vs.	<i>kul-is</i>	'tree-dat.M.S'
<i>shur'</i>	'children-P'	Vs.	<i>shur'-an</i>	'children-dat.M.P'
<i>insa:n</i>	'human-S'	Vs.	<i>insa:n-as</i>	'human-dat.M.S'
<i>Kagaz</i>	'paper-S'	Vs.	<i>kagz-an</i>	'paper-erg.M.P'
<i>Pash</i>	'roof-S'	Vs.	<i>pash-I</i>	'roof-abl.M.S'
<i>ko:ri</i>	'girls-P'	Vs.	<i>ko:r'-av</i>	'girls-erg.F.P'
<i>me:z</i>	'table-S'	Vs.	<i>me:z-I</i>	'table-abl.M.S'
<i>Pash</i>	'roof-S'	Vs.	<i>pash-uk</i>	roof-gen .M.S

Table 3.19: Examples of Nouns with Cases

3.4.8 Particles

Particle is a function word that is not assignable to any of the traditional grammatical word class (such as Nouns, Conjunctions, etc.). Particles are uninflected and can float within a sentence and can be removed from the sentence without changing the grammaticality of the sentence. Expressions like *ti* 'also', *toti* 'still/even then' are tagged as Particles in this tagging scheme. The tag used for Particles is PRT.

3.4.9 Adverbs

Adverbs are words that modify part of speech like Adjectives and Verbs.

Adverbs are mainly divided into three types viz manner, place and time Adverbs.

The tag structure for Adverbs is given below.

Adverb	Example	Tag structure
Temporal Adverbs	<i>yeli</i>	<ADV><ATM>
	'when'	
Locative Adverbs	<i>yeli</i>	<ADV><ALO>
	'here'	
Manner Adverbs	<i>yithpə:Th'</i>	<ADV><AMN>
	'in this way'	

Table 3.20: Tag structure of Adverbs

3.4.10 Postpositions

Kashmiri has Postpositions. Apart from regular features, Postpositions in Kashmiri also express certain grammatical functions such as cases like Locative, Instrumental, Sociative and Allative. On the basis of their distribution Postpositions are divided into two groups.

Postpositions governing the dative case

Postpositions governing the ablative case

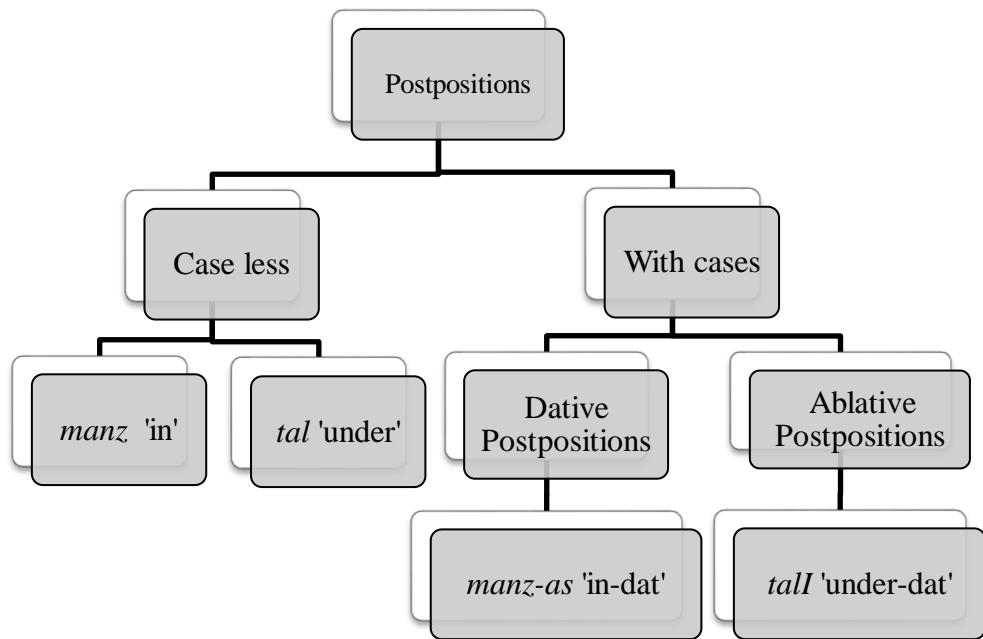


Fig.3.1: Postpositions with their attributes

Postposition	Example	Tag structure
Dative Postposition	<i>manz-as</i>	<PSP><dat>
	'room-dat'	
Ablative Postposition	<i>manz-I</i>	<PSP><abl>
	'room-abl'	

Table 3.21: Tag structure of Postpositions

3.4.11 Conjunctions

Conjunctions are divided into two types subordinating and coordinating Conjunctions. The tag 'CC' with 'SB' and 'CO' as its type tags is used for all Conjunctions, for example,

Conjuncts	Example	Tag Structure
Coordinating Conjunction	<i>Magar</i>	<CC><CO>
	'but'	
	<i>zantl</i>	
	'as if'	
	<i>hargah</i>	
	'in case'	
Subordinating Conjunction	<i>ti</i>	<CC><SB>
	'that'	
	<i>zi</i>	
	'that'	
	<i>ki</i>	
	'that'	

Table 3.22: Tag structure of Conjunctions

3.4.12 Intensifiers

The expressions intensifying an Adjective or an Adverb is tagged as INTF.

For example,

Intensifier	Tag
<i>va:rya:h</i>	<INTF>
‘very’	
<i>t'u:t</i>	
‘very/so much’	

Table 3.23: Tag structure of Intensifiers

3.4.13 Interjections

Interjections will be marked as INJ. As such very few expressions are used for Interjections in Kashmiri, for example,

Interjection	Tag
<i>o:ho:</i>	<INJ>
‘oh’	

<i>aha:</i>	
'ah'	

Table 3.24: Tag structure of Interjections

3.4.14 Negatives

Negation is expressed by the negative Particle *na*. The imperative *mI* 'don't', *maI* 'don't' and the negative conditional *nai* 'if not' is also used to express negation. Negative expression like *na* 'no', *kihinnI* 'nothing', *zIhIn'* 'never', etc. will be tagged as negatives and thetag used will be NEG.

Negation	Tag
<i>Na</i>	<NEG>
'no'	
<i>kihinn</i>	
'nothing'	
<i>mI</i>	
'don't'	

Table 3.25: Tag structure of Negatives

3.4.15 Question Words

All question words like *kya* ‘what?’ and *k'a:zi* ‘why?’ and Interrogative Adverbs like *kithI pə:Th* ‘how?’, *kati* ‘where?’ etc. are tagged as ‘QW’ in the present tagset.

3.4.16 Equatives

Equatives are used to show comparison between Nouns. To tag equatives ‘EQ’ is used, for example,

Example	Tag
<i>khətI</i>	<EQ>
‘than’	
<i>Nishi</i>	
‘in comparison of’	

Table 3.26: Tag structure of Equatives

3.4.17 Vocatives

Vocatives behave sometimes as Interjections but they have an independent function as well when used separately, thus it needs a separate tag. For vocatives ‘VOC’ is used as a tag. Vocatives can also be depicted as an attribute when they are attached to a word as a bound morph and the tag used is ‘VOC’.

Vocative Particle	Tag Structure
<i>hey</i>	<VOC><mas>
'hey'	
<i>hat.sa</i>	
'hello'	
<i>haye:</i>	<VOC><fem>
'hey (woman/girl)'	
<i>hatI bin'</i>	
'hey girl'	

Table 3.27: Tag structure of Vocative words

Vocatives	Tag Structure
<i>aslamo:</i>	<N><NP><mas><sng><voc>
'hey Aslam'	
<i>shi:lay</i>	<N><NP><fem><sng><voc>
'hey Sheela'	

Table 3.28: Tag structure of Vocatives as an attributes

3.4.18 Emphatic Words

There are no distinct Emphatic words that modify the Nouns, Pronouns, Verb, Adverbs and Adjectives. Emphatics are either bound morphemes attached to a word itself, or they can be free morphs also as in Emphatic Particles.

Noun	Emphatic marker	Emphatics
<i>bI</i>	+ <i>Iy</i>	<i>bIy</i>
‘I’	Emphatic	‘I-emph’
<i>shur'</i>	+ <i>i:</i>	<i>shuri:</i>
‘children’	Emphatic	‘children-emph’
<i>vɔzul</i>	+ <i>ui</i>	<i>vɔzlui</i>
‘red’	Emphatic	‘red-S.emph’

Table 3.29: Examples of Emphatics

Since Emphatics themselves fall into different grammatical categories, Emphatics are tagged according to the category they possess that is if a word is an Emphatic Pronoun then the word is tagged as a Pronoun with its types and attributes (here Emphatic is used as an attribute). Tag used for Emphatic Particles

is ‘EMPH’ and tag ‘emph’ is used while tagging Emphatics as attributes. Consider the following examples

Emphatic Particle	Tag structure
<i>pa:nay</i>	<EMPH>
‘by self’	

Table 3.30: Tag structure of full word Emphatics

Emphatics	Example	Tag Structure
Emphatic Pronoun	<i>bIy</i>	<P><PRP><Ip><mas><sng><nom><emph>
	‘me only’	
Emphatic Noun	<i>shur'nIy</i>	<N><NC><mas><plu><dat><emph>
	‘only by the kids’	
Emphatic Adjective	<i>vɔzI'avIy</i>	<JJ><QN><mas><plu><erg><emph>
	‘only by the reds’	

Table 3.31: Tag structure of Emphatics as an attribute

3.4.19 Honorific

This tag includes honorific free morphemes (words) like *hez*, *mahra* (honorific indicators), etc. which exist independently, as opposed to honorific attributes (which are bound) like *tuh* 'you- honorific'. Tag used is 'HON' for free morphemes and 'hon' for attributes.

3.4.20 Symbol

All special characters like #, @, &, \$, %, etc. are tagged as 'SYM'. This tag is similar to the Penn tag set tag 'SYM'.

3.4.21 Punctuation

All punctuation marks like ‘,’ ‘.’ ‘:’ ‘;’ ‘:’ are tagged as 'PUNC'.

3.4.22 Unknown

The unknown tag is assigned for those words where the category is not known or if there is a doubt regarding the category of a particular word. Tag used is 'UNK'. The tag used here is a temporary tag.

3.4.23 Multi Word Expressions

Multi word expressions include Echo words, Compounding and Reduplication, and Part of Words.

3.4.23.1 Echo Words

Indian languages have a highly productive usage of echo words. Echo words are very frequently used in Kashmiri and such words are tagged as 'ECH'. First part of the Echo Word is tagged according to the category which it possesses and second part of it is tagged as ECH.

Echo words	Tag
<i>ca:y sha:y</i>	<N><NC><fem><sng>
'tea and the like'	<ECH>

Table 3.32: Tag structure of Echo Words

3.4.23.2 Compounding

Compound words are made up of two constituents, each constituent either belongs to a Noun, Verb, Adjective or a Postposition. The compound itself may belong to the category Noun, Adjective or Verb. The compound is categorized on the basis of its head. A compound may be endocentric (when there is a definable head inside the group that has the same distribution as the compound as a whole) or it can be exocentric (i.e., when there is no definable head that relates to the compound as a whole). According to the present scheme of tagging, components of the compound words like Noun compounds, Adjectival compounds, Verbal compounds are treated as individual words and are tagged separately rather than giving a single tag to the whole word sequence (compound). That is, if a word is a compound then both of its components are tagged according to the category they

posses, that is, if a word is an Adjective- Noun compound then the first component is tagged as an Adjective (with its type and attributes) and the second component is tagged as a Noun (with its type and attributes) and if it is part of a Compound Adjective (Adjective-Adjective Compound) then the category ‘Adjective’ is used for both the components and so on and so forth. This decision was arrived at after great deal of thought and taking into consideration the inadequacy of handling multiple words (of a compound) under one tag (compound tag) without losing some intrinsic individual word (compound components) attributes and properties. Furthermore, it would be much more feasible to handle compounds at the Chunking level.

Some of the Kashmiri compound examples are given as under:

Compounds	Example		Tag
Compound Noun	<i>candI</i>	<i>tsu:r</i>	<N><NC><mas><sng><nom><N> ><NC><mas><sng><nom>
	pocket	thief	
	‘pickpocket’		
	<i>o:m</i>	<i>dɔd</i>	<JJ><QL><mas><sng>
	raw	milk	<N><NC><mas>
‘raw milk’			
Adjectival Compound	<i>Tsok</i>	<i>modur</i>	<JJ><QN><mas><sng><JJ><QN> ><mas><sng>
	sour	sweet	

	‘sour and sweet’		
Verb Compound	<i>she:run</i>	<i>pə:run</i>	<V><VM><mas><sng><nfn><V> ><VM><mas><sng><nfn>
	to make better	to decorate	
	‘to decorate’		
Adverb Compound	<i>ja:yi</i>	<i>ja:yi</i>	<ADV><LOC>
	place	place	
	‘everywhere’		
Pronoun Compound	<i>kāh</i>	<i>kāh</i>	<P><PID><mas><sng><nom>
	someone	someone	
	‘someone’		

Table 3.33: Tag structure of Compounds

3.4.23.3 Reduplication

When the same word is written twice with or without an internal change in it, and the resulting reduplicated form is grammatically or semantically significant, the process is known as reduplication. For such word combinations category, type and attributes of the first word is written followed by the tag ‘RDP’

Reduplication	Example		Tag
Pronoun Reduplication	<i>panIn'-an</i>	<i>panIn'-an</i>	<P><PRF><mas><plu><dat> <RDP>
	self-dat	self-dat	
	'of self's'		
Noun Reduplication	<i>ta:n-as</i>	<i>ta:n-as</i>	<N><NC><mas><sng><dat> <RDP>
	organ-dat	organ-dat	
	'every organ/body part'		
Verb Reduplication	<i>asa:n</i>	<i>asa:n</i>	<V><VM><mas><sng><prt><prog><fn> <RDP>
	laughing	laughing	
	'while laughing/easily'		
Adverb Reduplication	<i>jaldi:</i>	<i>jaldi:</i>	<ADV><AMN> <RDP>
	quickly	quickly	
	'very quickly'		
Postposition	<i>sI:t'</i>	<i>sI:t'</i>	<PSP>

Reduplication	with	with	<RDP>
	'close/together'		
Adjective Reduplication	<i>lɔkIT'</i>	<i>lɔkIT'</i>	<JJ><QN><mas><plu><nom> <RDP>
	little	little	
	'very little'		

Table 3.34: Tag structure of Reduplication

3.4.23.4 Part of Word

Part of word is used for those morphological words which are separated by an orthographic space and thus divided into two tokens as discussed below. For such words POW tag is used.

3.4.23.4.1 Word Boundaries (Token Division)

In Kashmiri, because of orthographic compulsion a single morphological word may be divided into two tokens i.e. a single word can be broken with a white space in between unlike in languages like English, Hindi, etc. In tagging every orthographic white space is considered as a word break even if it occurs within a lexical word.

Example:

Example	Gloss
دپی ز	'you should say'
<i>dəp'zi</i>	
میل ژار	'cooperation'
<i>miIltsa:r</i>	
ہتہ بورد	'hundreds'
<i>hatIbod</i>	
کھلہ ڈلہ	'spacious'
<i>khulIDull</i>	
نالہ موت	'embrace'
<i>nallmot</i>	

Table 3.35: Examples of Token division

Thus, the problem here was how to tag two tokens which make up a single morphological word. After considering various solutions to this problem it was decided to use two separate tags to follow the tagging rules i.e. for such words first part will be tagged according to the category which the whole word (the two word components together) possesses and the second part of the same word will be tagged as POW (part of the word). That means if we have a Noun or an

Adjective which is separated by an orthographic space then the first part of the word will be tagged as <N><...><...> or <JJ><...><...> and second part of the same word will be tagged as <POW>.

POW will also be used for those Nouns which take *-vo:l* as a suffix with them e.g. *gdɔdIvo:l* ‘milkman’, *sabzi:vo:l* ‘vegetable seller’. Such constructions will be tagged as ‘<N><type><attribute> [SPACE]¹²<POW>’. Some examples are given as under.

Word	مِلّے ژار	نالہ مۆت	ریتہ واد
	<i>milltsa:r</i>	<i>na:lImot</i>	<i>retIva:dI</i>
Tag	<N><NC><mas><sng> <POW>	<N><NC><mas><sng>< POW>	<N><NC><mas><plu> <POW>
Gloss	‘co-operation’	‘embrace’	‘months’
Word	اؤند پوک	بُتھی لاگے	کھلہ ڈلہ
	<i>Ondpok</i>	<i>buth'la:gay</i>	<i>khulIDull</i>
Tag	<ADV><ALO><POW>	<N><NC><fem><sng> <POW>	<JJ><QN><POW>
Gloss	‘surrounding’	‘defy’	‘spacious’
Word	دوڈ وول	سبزی وول	ماکان وول

¹²[SPACE] represents single white space

	<i>dɔdIvo:l</i>	<i>sabzi:vo:l</i>	<i>maka:nvo:l</i>
Tag	<N><NC><mas><sng> <POW>	N><NC><mas><sng><P OW>	N><NC><mas><sng>< POW>
Gloss	‘milkman’	‘vegetable seller’	‘Land Lord’

Table 3.36: Tag structure of POW

3.5 Proposed Tagset for Kashmiri

The final tagset which is proposed for Kashmiri is given in a tabular form below

S. No	Category	Tag name
1	Noun	N
1.1	Common Noun	NC
1.2	Proper Noun	NP
2	Pronoun	P
2.1	Personal Pronoun	PRP
2.3	Demonstrative Pronoun	PDM
2.4	Reflexive Pronoun	PRF
2.5	Reciprocal Pronoun	PRC
2.6	Possessive Pronoun	PPO
2.7	Relative Pronoun	PRL

2.8	Indefinite Pronoun	PID
2.9	Interrogative Pronoun	PIT
3	Adjective	JJ
3.1	Qualitative	QL
3.2	Quantitative	QN
4	Gerund	VGB
5	Verb	V
5.1	Main Verb	VM
5.2	Auxiliary Verb	VAUX
5.3	Conjunctive Verbs	VCNJ
5.4	Causative Verbs	VCUS
6	Adverb	ADV
6.1	Adverb of Location	ALO
6.2	Adverb of Time	ATM
6.3	Adverb of Manner	AMN
7	Particle	PRT
8	Postposition	PSP

9	Conjunction	CC
9.1	Coordinating Conjunction	CO
9.2	Subordinating Conjunction	SB
10	Interjection	INJ
11	Emphatic	EMP
12	Honirific	HON
13	Vocative	VOC
14	Question Words	QW
15	Equative	EQT
16	Intensifier	INTF
17	Negative	NEG
18	Quantifier	QTF
18.1	Cardinal	CRD
18.2	Ordinal	ORD
19	Reduplication	RDP
20	Echo Words	ECH
21	Part of word	POW

22	Symbol	SYM
23	Punctuation	PUNC
24	Unknown	UNK

Table 3.37: Proposed Tagset for Kashmiri

3.6 Conclusion

The aim of this tagset is to provide clear instructions for annotating the Kashmiri corpus. The tagset developed so far is hierarchical in nature as it is divided into main word categories, types of the categories and their sub features or attributes as discussed above. This tagset follows the Eagles guidelines, and here also word categories are obligatory, types are recommended and attributes are optional. The need for designing a hierarchical tagset is to capture the morphological richness of the language, and moreover it is effective in reducing confusions like ambiguities and inconsistencies that come up while annotating the corpus. A hierarchical design gives greater flexibility, extensibility and re-usability. Since hierarchical tagsets are more elaborative and comprehensive in nature, consequently the same tagset can be used at all levels - POS tagging, chunking, dictionary, morphological analysis and deep parsing. As far as the design of the present Kashmiri tagset is concerned, it is designed in a way where there is linguistic compatibility with other languages as well. However, the tagset designed here is developed keeping in view the morphosyntactic features of the language thus the tagset here is more compatible with Kashmiri language. However, there are certain feature(attributes) which are not present in

Kashmiri but are included in the above mentioned guidelines taking into consideration the cross linguistic usage of the tagset. In this way we can say the tagset developed here can be used for other languages as well with modifications and additions, wherever required.