



## Urdu Part of Speech Tagset

December 07, 2007

**CENTER FOR RESEARCH IN URDU LANGUAGE PROCESSING  
NATIONAL UNIVERSITY OF COMPUTER AND EMERGING SCIENCES, LAHORE  
PAKISTAN**

## Table of Contents

1	Introduction.....	4
2	Urdu Parts of Speech Classification.....	5

## Revision History

Name	Change Date	Version	Description of Changes
Hassan Sajjad	07-12-2007	1.0.0.0	Initial document

# 1 Introduction

Tagset of a language caters main parts of speech as well as morphological information of the language. A tagset may be consisted either of syntactic categories or it may be consisted of morpho-syntactic categories. Considering the efficiency in machine learning process and to reduce lexical and syntactic ambiguity, it was decided to concentrate on the syntactic categories of language.

There were three types of corpus available for analysis i.e. literature, news and poetry corpus. For the design of tagset, only literature and news corpus was analyzed. The corpus was based on the most recent available vocabulary used by local people.

## 2 Urdu Parts of Speech Classification

### Demonstrative:

Demonstratives are divided into four categories. All four categories of demonstratives have ambiguity with four categories of pronoun. Phrase level analysis was done to distinguish between demonstrative and pronoun. Following are some examples of demonstratives.

Personal demonstrative (PD)	This category includes the elements of demonstrative and personal demonstratives. Following is an example of it.
ہم، تم، آپ، یہ، وہ، اس	یہ <PD> مسجدیں <NN> ہماری <G> پہچان <NN> ہیں <VB>۔ <SM>
Relative demonstrative (RD)	جو <RD> لڑکا <NN> صبح <NN> آیا <VB> تھا <TA> وہ <PP> میرا <G>
جو، جن، جنہوں	دوست <NN> ہے <VB>۔ <SM>
Kaf demonstrative (KD)	کن <KD> لوگوں <NN> کو <P> آم <NN> اچھا <ADJ> لگتا <VB> ہے <TA>۔ <SM>
کن، کوئی	کمرے <NN> میں <P> کوئی <KD> لڑکا <NN> نہیں <NEG> ہے <VB>۔ <SM>
Adverbial demonstrative (AD)	میں <PP> ایسا <AD> کام <NN> نہیں <NEG> کر <VB> سکتا <AA>۔ <SM>
اب، تب، ادھر، یہاں	

### Nouns:

Nouns are divided into two categories. First category consists of simple nouns which are represented by NN in the tagset. However, there are other nouns that show adverbial nature like time, place, manner, etc. These are also catered under noun. The proper nouns are kept in a separate category. Following are some examples of different types of nouns.

Noun (NN)	یہ <PD> مسجدیں <NN> ہماری <G> پہچان <NN> ہیں <VB>۔ <SM>
جہاز، زمین، درخت، لڑکا، اوپر، اندر، سمیت، طرح، طرف	چھت <NN> کے <P> اوپر <NA> حامد <PN> ہے <VB>۔ <SM>
Proper noun (PN)	لاہور <PN> باغات <NN> کا <P> شہر <NN> ہے <VB>۔ <SM>
لاہور، پشاور، پاکستان	

### Pronouns:

Pronouns are divided into six categories based on their syntactic structure. Most of the categories are consistent with the types provided by Urdu grammarians. Following are some examples of the types of pronouns.

Personal pronoun (PP)	میں <PP> تمہارا <G> دوست <NN> ہوں <SM>_<VB>
میں، ہم، تم، آپ، یہ، وہ، اس	
Reflexive pronoun (RP)	میں <PP> اپنا <GR> کام <NN> خود <RP> کروں <VB> گا <SM>_<TA>
خود، آپ	
Relative pronoun (REP)	علی <PN> جو <REP> حامد <PN> کا <P> بھائی <NN> ہے <VB> میرا <G> دوست <NN> ہے <SM>_<VB>
جو، جن، جنہوں	

#### Adverbial pronoun (AP):

The adverbial pronouns occur at the place of nouns with adverbial nature and show the property of time, place, manner, etc. They are represented by AP in the tagset. Consider the following examples:

Example:	علی <PN> نے <P> اب <AP> کھانا <NN> کھایا <VB> ہے <SM>_<TA>
اب، تب، ادھر، یہاں	

#### Kaf pronoun (KP):

Kaf pronouns add interrogative property in the sentence. They are divided into two categories. Kaf pronouns, represented by KP, are used to ask question about a noun. The second category includes adverbial kaf pronouns which are used at the place of nouns with adverbial nature. Following are their examples:

Kaf pronoun (KP)	کمرے <NN> میں <P> کون <KP> ہے <SM>_<VB>
کون، کوئی، کن	
Adverbial kaf pro (AKP)	علی <PN> کدھر <AKP> گیا <VB> ہے <SM>_<TA>
کدھر، کب، کیسا	
Genitive reflexive (GR)	اپنا <GR> کام <NN> خود <RP> کرنا <VB> میرا <G> فرض <NN> ہے <SM>_<VB>
اپنا	
Genitives (G)	Consider the example of genitive reflexive
میرا، تمہارا، ہمارا، تیرا	

**Verb (VB):** At sentence level, any word showing action in any form is considered as verb. No further categorization is done. Consider the following examples of verb:

Example:	وہ <PP> روٹی <NN> کھا <VB> رہا <AA> ہے <SM>_<TA>
لکھنا، کھانا، جاتا، کرنا	

**Auxiliaries:**

Based on the syntactic nature of language, auxiliaries are divided into two categories. Aspectual auxiliaries always occur after main verb of the sentence. Tense auxiliaries are used to show the time of the action. They occurred at the end of the verb phrase. Consider the examples of aspectual and tense auxiliaries:

Aspectual auxiliary (AA)	Consider the example of verb.
رہا، کرنا، چکھ	
Tense auxiliary (TA)	Consider the above describe examples.
ہے، ہیں، ہوں، تھا، تھے، تھیں، گا، گی، گے، ہو، ہوں	

**Adjective (ADJ):**

Adjectives are catered as one category. The information related to the degree of adjective is not taken into account. Following are given some examples of adjectives.

ظالم، خوبصورت، کمزور، بیکار، سمجھدار، نفیس	حامد <PN> بہت <ADV> ظالم <ADJ> لڑکا <NN> ہے <VB>۔ <SM>

**Adverb (ADV):**

Adverbs are handled as one category in the tagset. Consider the following examples of adverbs.

Example:	وہ <PP> بڑا <ADV> محنتی <ADJ> لڑکا <NN> ہے <VB>۔ <SM>
بہت، نہایت، بڑا	

**Quantifier (Q):**

Consider following examples of quantifier:

Example:	سب <Q> لوگ <NN> تھوڑا <Q> انتظار <NN> کریں <VB>۔ <SM>
کچھ، چند، تمام، اتنے، سب، تھوڑا، تھوڑے، کئی، بعض، کل	

**Numerals:**

Numerals are divided into four categories based on their syntactic structure. Cardinal (CA), ordinal (OR), fractional (FR) and multiplicative (MUL) are types included in the tagset. Following are the examples of each category.

Cardinal (CA)	پہلے <OR> دو <CA> لڑکوں <NN> کو <P> بلاؤ <VB>۔ <SM>
ایک، دو، تین، چار بیالیس، افسٹھ، ننانوے، ہزار، دو ہزار	
Ordinal (OR)	Consider the example of cardinal.

پہلا، دوسرا، تیسرا، چوتھا، پانچواں، چھٹا، ساتواں، آٹھواں، آخری	
Fractional (FR)	<SM>_<VB> دینا <NN> دودھ <U> کلو <FR> ڈھائی
چوتھائی، ڈھائی، اڑھائی	
Multiplicative (MUL)	<SM>_<VB> ہے <ADJ> موٹا <MUL> دگنا <P> سے <PN> حامد <PN> علی
گنا، دگنا، دہرا، تہرا	

### Measuring unit (U):

They are frequently used with numerals. However, they have a different syntactic structure than numerals. Consider the example of fractional to see the occurrence of measuring units.

Example:	<SM>_<VB> دینا <NN> دودھ <U> کلو <FR> ڈھائی
پون، پائو، کلو، سیر	

### Conjunction:

Conjunctions are divided into coordinating and subordinating conjunctions. Following are their examples:

Coordinating (CC)	<NN> دوست <ADJ> اچھے <PN> علی <CC> اور <PN> حامد <SM>_<VB> ہیں
یا، اور	
Subordinating (SC)	<PP> مجھ <SC> کہ <VB> کہو <P> سے <PN> حامد <SM>_<VB> ملے <P> سے
کیونکہ، کہ	

### Intensifier (I):

There are only three words in this category. Consider their following examples:

Example:	<SM>_<TA> گا <VB> آؤں <PP> بھی <I> ا
ہی، بھی، تو	

### Adjectival particle (A):

This category includes only one word sa with its two inflection forms. This particle is normally used for comparison. Consider the following examples of adjectival particle.

Example:	<NN> میٹھک <NN> ایک <CA> عجیب <ADJ> سا <A> جانور <NN> <SM>_<VB> ہے
سا، سے، سی	



**KER particle (KER):**

These particles normally occur in verb phrase. There are only two entities in this class. Consider the following examples:

Example:	<SM>_<AA>دینا <VB>کر <NN>فون <KER>کر <VB>پہنچ <NN>گھر
کے، کر	

**Title:**

Titles are divided into two categories based on their pre and post occurrence around a proper noun. Consider their examples below.

Pre-title (PRT)	<NN>انسان <ADJ>اچھے <POT>صاحب <PN>سرمد <PRT>میاں <SM>_<VB>ہیں
حضرت، میاں	
Post-title (POT)	Consider the example of pre-title above.
جی، صاحب	

**Semantic Marker (P):**

Following are the list of particles included into this category. However, the entity سے is kept as separate category due to its ambiguous usage.

کا، کو، کی، کے، نے، میں، تک تک، پر	<SE> سے <NN>چھڑی <P>نے <PN>علی <P>کو <PN>حامد <SM>_<VB>مارا
<b>SE (SE): سے</b>	Consider the above example

**Wala (WALA):**

This category contains one word wala and its inflections. Consider its examples:

Example:	<SM>_<TA>ہے <VB>آیا <NN>آدمی <WALA>والا <VB>بیچنے <NN>پہل
والا، والی، والے	

**Negation (NEG):**

Consider the following examples of negation.

Example:	<SM>_<AA>سکتا <VB>کر <NEG>نہیں <NN>کام <AD>ایسا
نہ، نہیں	

**Interjection (INT):**

Interjections normally occur at the start of the sentence. They are kept as separate category in the tagset. Following are its examples:

Example:	<SM>_<TA> ہے <VB> کی <NN> بات <ADJ> اچھی <ADV> کیا <INT> واہ
واہ, سبحان اللہ, اچھا	

**Question words (QW):**

There are some words instead of kaf pronouns that are used for the interrogation in the sentence. However, these words cannot be replaced by a noun or pronoun. A separate category of question words has been formed for these words. Consider their examples below:

Example:	<SM>_<TA> گا <VB> جائے <NN> سکول <PN> علی <QW> کیا
کیا, کیوں	

**Punctuation marks:** In this tagset, punctuation marks are divided into two categories. Sentence markers mark the boundary of the sentence. Phrase markers are used inside the sentence but never used at the end of sentence. Consider their examples below:

Sentence marker (SM)	‘,’ , ‘?’
Phrase marker (PM)	‘,’ , ‘,’

DATE	2007, 1999
------	------------

**Expression (Exp):**

Any word or symbol which is not handled in this tagset will be catered under expression. It can be mathematical symbols, digits, etc.