

# Nelralec/Bhasha Sanchar Working Paper 2

## Categorisation for automated morphosyntactic analysis of Nepali: introducing the Nelralec Tagset (NT-01)

by Andrew Hardie, Ram Lohani, Bhim Regmi and Yogendra Yadava

July 2005

### 1. Introduction

This working paper presents a summary of the Nelralec Tagset, a categorisation system for the manual and automated analysis of morphosyntactic units in Nepali<sup>1</sup>.

This tagset has been compiled with reference to widely published grammars of Nepali. These grammars in the main describe standard Nepali. The tagset may not be sufficient for other dialects of Nepali<sup>2</sup>.

In sections 2 to 4, the most important issues underlying the construction of the categories in the tagset are discussed: the handling of case affixes/postposition, the handling of gender, and the model of Nepali verb inflection underlying the tagset categories. In section 5, the tagset itself is tabulated with examples.

### 2. Tokenisation and postpositions

The first step in both manual and automated tagging is the division of the text into *tokens*. Tokens are often described as “words” but in this context it may be better to think of them as “appropriately-sized units for morphosyntactic analysis”; such a unit may be larger than, or smaller than, an orthographic word (i.e. a string of letters bounded by white space or punctuation).

For the purposes of Nepali tagging, there is one major instance where words must be split apart to create appropriately-sized units for analysis: a group of morphemes which are variously considered to be clitics, case affixes, or postpositions.

While in many languages case, gender and number are best analysed together, in Nepali case and number are indicated by elements similar to (and tagged similarly to) postpositions, whereas gender is marked by an inflected affix. Thus, distinct considerations apply to gender and to case-number, (see below on gender).

In Nepali, there are a large number of forms that are considered to be *either* postpositions *or* case-suffixes. There is no 100% agreement on what is or is not a suffix. Some Nepali grammatical traditions consider the following forms to be suffixes: *haruu*, *ko/kii/kaa*, *le*, *laaii* (whose meanings are, respectively, plural/collective; genitive; ergative/instrumental; and accusative/dative), and

---

<sup>1</sup> It should be noted that this published document represents an abbreviation of the full tagset description used for manual and automatic tagging on the Nelralec project.

<sup>2</sup> When an example in the tagset definition is enclosed in [square brackets], that example has been reported to be limited to particular dialects.

consider all other postposition/suffix forms to be postpositions. However, in terms of part-of-speech tagging it is difficult to draw a consistent distinction between *haruu*, *ko/kii/kaa*, *le*, *laaii* and other postpositions such as *maa*, *baaTa*, *sanga*, *dekhi*, etc., because the combinations in which these morphemes can occur are too various to allow them to be classified using variations in the tag for the form that they are attached to<sup>3</sup>.

For this reason, *all* these forms are tagged separately to the noun, adjective, pronoun or other form to which they are attached. This means that postpositions are necessarily tokenised separately to the forms to which they are attached, and tagged separately.

### **3. Gender on nouns and adjectives**

Some Nepali categories mark gender by means of a three-way *o/ii/aa* distinction, where *o* is masculine, *ii* is feminine, and *aa* is “other” (it can indicate masculine plural, or feminine plural, or oblique case marking motivated by a following postposition, or honorific marking). There is also one postposition, *ko/kii/kaa*, which is marked for gender in this way.

However, this gender marking is not found on all members of the gendered categories. Many words are “unmarked” – that is, they have a single invariant form regardless of gender. For example, adjectives are a gendered category, speaking generally: they are marked for gender agreement with the noun they modify. But many, many adjectives are actually *unmarked* – that is, there is nothing in their morphology to indicate their gender.

Nepali nouns have *natural* gender rather than grammatical gender. Animate females are feminine, all other nouns are masculine. In most cases gender is not marked on nouns in the way that it is marked on some adjectives. For a minority of nouns, there are pairs of masculine and feminine nouns related through the *-o/ii* distinction, for instance *keTo* “boy” / *keTii* “girl”. But there are also numerous feminine nouns that end in *ii* without there being a masculine equivalent in *o* (e.g. *aaimaaii*, “woman”).

Thus the *o/ii/aa* distinction is ignored on nouns for the purpose of this system of POS tagging, but *not* on pronouns, adjectives, non-finite verbs, etc., where the distinction is motivated by agreement and is thus clearly inflectional rather than lexical-derivational, as is the case for nouns. Gender is not included in the tags for nouns because it is not an inflectional category in the same way that it is for adjectives.

### **4. Modelling Nepali verb inflections**

Nepali has a vast number of verb inflections. This is because compounding is an extremely productive process in the verb system of Nepali. Different combinations of various non-finite forms of the main verb and various inflected forms of secondary verbs creates a very large number of

---

<sup>3</sup> The following convention is used in this document when exemplifying tokens that result from the splitting-up of another token: when an example begins or ends in the symbol #, this indicates that the form being exemplified is a form that has been separated – i.e. it is either a morpheme which has been separated, or a form from which a morpheme has been separated.

tense-aspect-mood combinations<sup>4</sup>. If every distinction made in the verb system were to be indicated by a separate tag, then the tags for verbs would become entirely unmanageable. Thousands of tags would be required for the verbs alone. The simplified classification model of the Nepali verb that this tagset uses is described in this section.

To avoid the problem of compound verbs creating an unmanageable number of categories, the tagging system described here uses this rule for verbs:

**Every compound verb is tagged according to the last verb in the compound.**

That means that only the last identifiable verb in a compound verb is taken into account when deciding the tag. If there is only one identifiable verb, then the whole thing is taken into account.

As a further simplifying measure, certain aspects of the verb morphology system are ignored. These are passive, causative, and negative. This means that:

- Any passive verb is tagged the same as the corresponding active verb
  - E.g. tag *garinu* the same as *garnu*
- Any causative verb is tagged the same as the corresponding non-causative verb
  - E.g. tag *garaau~cha* the same as *garcha*
- Any negative verb is tagged the same as the corresponding positive verb
  - E.g. tag *garnechaina* the same as *garnecha*

No distinction is made in the tagging between auxiliary verbs and main verbs.

The tagset makes, for convenience of reference only, a distinction between finite (person-marked) verb forms, and non-finite (non-person-marked) verb forms. Many of the non-finite forms occur embedded at the start of a longer verb (e.g. **gardaithyo**, **garnuhuncha**). However, in accordance with the general rule, they are **not** tagged separately in this case. This means that non-finite tags are **only** used *if the non-finite form is written as a separate word, or if the non-finite form is at the end of the longer verb*.

The non-finite forms distinguished by tagset are as follows:

- The infinitive
- A group of participles<sup>5</sup>:
  - the *e(ko)*-participle, sometimes called the *perfect participle*
  - the *d*-participle (*do/dii/daa/dai*), which is used for three functions:
    - the converb/participle function (also known as the *conjunctive participle*, the *progressive participle*, and the *simultaneous converb*)
    - the modifier function
    - as an element of compound verbs (where it does not receive a separate tag)

---

<sup>4</sup> The term “compound verb” is used in a wide sense here, to refer to *any verb that consists of two or more elements that are identifiable as belonging to separate independent verb lexemes*. So, for instance, *garidiyo* would be a straightforward example of a compound verb; however, for the purposes of this description, words such as *garthyo*, *garcha*, *garirahyo*, and *garnecha* will be considered to be compound verbs as well. Compounds consisting of noun+verb or adjective+verb can be tagged using the same rule: only consider the second element, i.e. the last verb.

<sup>5</sup> They should probably not all be considered “participles” in the most precise grammatical sense. Some are “converbs” and some may be more precisely analysed as infinitival.

- the *ne*-participle (the *imperfect participle* or *infinitival participle*)
- the sequential converbs, also called *absolutive participles*, of which there are three, which all receive the same tag:
  - the *era*-participle
  - the *ii*-participle
  - the *iikana*-participle
- Two other forms:
  - The *e*-form (often referred to as *subjunctive* or *conditional*)
  - the *i*-form, sometimes referred to as the *passive root*, which is the form in which a non-final verb in a compound appears.
- Three command forms (also called the *imperative*)

The finite forms are modelled as follows:

There are six finite forms for each tense/mood/aspect combination, which are indicated in the third person singular by the following typical endings: *-yo*, *-thyo*, *-echa*, *-cha*, *-necha*, *-laa*. Many of these tenses follow similar patterns of inflection. An additional finite form, the *optative*, is given separate tags.

The following distinctions operate on finite verbs:

- Person: first, second, third
- Gender: masculine (default), feminine
- Number: singular, plural
- Honorific level: non-honorific, medial-honorific<sup>6</sup>

Therefore, in theory there should be  $3 \times 2 \times 2 \times 2 = 24$  tags. However, not all possible combinations of features have separate forms in Nepali. There is, for instance, no specific form in any tense-mood for the second-person non-honorific feminine plural. There is just a single second-person plural form.

In general, only second and third person singular verbs are marked for gender and honorific level. Plural verb forms are not marked for gender or for honorific level. First person verbs are not marked for gender or for honorific level. Similarly, at the medial-honorific level, there is no distinction between singular and plural, except for feminine verbs (which are listed separately below). This means that if we take only gender, number and honorific level into account, and if we treat alike those categories which are merged together, only ten tags are necessary.

In many cases gender is not marked on verbs. Even where it is marked, it is reported that the “masculine” verbs may be used in some varieties of Nepali with feminine subjects. For this reason, the tagset considers “masculine” to be the default option. So the six “masculine” tags listed in the table below, are in fact not defined as masculine: they are simply the usual, default forms. The four distinct feminine forms– the second person and third person singular, non-honorific and medial-honorific – are defined in contrast to the default: the tags for them are created by adding F to the tag.

There are separate tags for optative verbs, as these verbs behave differently in many ways to the other finite verbs (e.g. by taking a prefix to indicate the negative, rather than a suffix). However, the

---

<sup>6</sup> The higher honorific levels (high, royal) are conveyed through compound verbs.

person-number-honorific categories in the optative paradigm are directly parallel to those of the general finite paradigm. There are no feminine forms in the optative, and thus no tags for them.

## 5. The tagset

The following table gives a summary of the tagset. It contains 112 tags.

Some of the categories in this tagset are open, and some are closed. In terms of the tagset, this simply means whether or not an exhaustive listing of the word-forms that can take that tag is possible. A *closed category* is one where it is possible to list the contents exhaustively; an *open category* is one where it is not. The tags RR and UU, and any tags beginning with N, J, F (except FZ), or V, indicate **open categories**. The tags FZ, RD, RJ and RK, and any tags beginning with C, I, P, D, M, or Y, indicate **closed categories**. Examples given in *italics* represent a closed category. Not all closed categories are listed exhaustively in the table below.

Category definition	Examples (Latin)	Examples (Devanagari)	Tag
Common noun	keTo, keTaa, kalam	केटो, केटा कलम	NN
Proper noun	raam	राम	NP
Masculine adjective	moTo, raamro	मोटो, राम्रो	JM
Feminine adjective	moTii, raamrii	मोटी, राम्री	JF
Other-agreement adjective	moTaa, raamraa	मोटा, राम्रा	JO
Unmarked adjective	saphaa, dhanii, asal	सफा, धनी, असल	JX
Sanskrit-derived comparative or superlative adjective	uccatar, uccatam	उच्चतर, उच्चतम	JT
First person pronoun	<i>ma, haamii, mai#</i>	म, हामी, मै#	PMX
First person possessive pronoun with masculine agreement	<i>mero, haamro</i>	मेरो, हाम्रो	PMXKM
First person possessive pronoun with feminine agreement	<i>merii, haamrii</i>	मेरी, हाम्री	PMXKF
First person possessive pronoun with other agreement	<i>meraa, haamraa</i>	मेरा, हाम्रा	PMXKO
Non-honorific second person pronoun	<i>ta~, tai#</i>	तँ, तै#	PTN
Non-honorific second person possessive pronoun with	<i>tero</i>	तेरो	PTNKM

masculine agreement			
Non-honorific second person possessive pronoun with feminine agreement	<i>terii</i>	तेरी	PTNKF
Non-honorific second person possessive pronoun with other agreement	<i>teraa</i>	तेरा	PTNKO
Medial-honorific second person pronoun	<i>timii</i>	तिमी	PTM
Medial-honorific second person possessive pronoun with masculine agreement	<i>timro</i>	तिम्रो	PTMKM
Medial-honorific second person possessive pronoun with feminine agreement	<i>timrii</i>	तिम्री	PTMKF
Medial-honorific second person possessive pronoun with other agreement	<i>timraa</i>	तिम्रा	PTMKO
High-honorific second person pronoun	<i>tapaai~, hajur</i>	तपाईँ, हजुर	PTH
High-honorific unspecified-person pronoun	<i>yahaa~, wahaa~<sup>7</sup></i>	यहाँ, वहाँ	PXH
Royal-honorific unspecified-person pronoun	<i>sarkaar, mausuph</i>	सरकार, मौसुफ	PXR
Reflexive pronoun	<i>aaphuu</i>	आफू	PRF
Possessive reflexive pronoun with masculine agreement	<i>aaphno</i>	आफ्नो	PRFKM
Possessive reflexive pronoun with feminine agreement	<i>aaphnii</i>	आफ्नी	PRFKF
Possessive reflexive pronoun with other agreement	<i>aaphnaa</i>	आफ्ना	PRFKO
Masculine demonstrative determiner	<i>yasto, yatro</i>	यस्तो, यत्रो,	DDM
Feminine demonstrative determiner	<i>yastii, yatrii</i>	यस्ती, यत्री	DDF
Other-agreement demonstrative determiner	<i>yastaa, yatraa</i>	यस्ता, यत्रा	DDO
Unmarked demonstrative determiner	<i>yo, yas#, yi, yin#, yinii, yati, yatti</i>	यो, यस#,यी, यिन#,	DDX

<sup>7</sup> There is an alternative form for *wahaa~*: *uhaa~*, उहाँ. This form would also take PXH.

		यिनी, यति, यत्ति	
Masculine interrogative determiner	<i>kasto, katro</i>	कस्तो, कत्रो	DKM
Feminine interrogative determiner	<i>kastii, katrii</i>	कस्ती, कत्री	DKF
Other-agreement interrogative determiner	<i>kastaa, katraa</i>	कस्ता, कत्रा	DKO
Unmarked interrogative determiner	<i>ko, kas#, ke, kun, kati</i>	को, कस#, के, कुन, कति	DKX
Masculine relative determiner	<i>jasto, jatro</i>	जस्तो, जत्रो,	DJM
Feminine relative determiner	<i>jastii, jatrii</i>	जस्ती, जत्री	DJF
Other-agreement relative determiner	<i>jastaa, jatraa</i>	जस्ता, जत्रा	DJO
Unmarked relative determiner	<i>jo, jas#, je, jati, josukai</i>	जो, जस#, जे, जति, जोसुकै	DJX
Masculine general determiner-pronoun	arko	अर्को	DGM
Feminine general determiner-pronoun	arkii	अर्की	DGF
Other-agreement general determiner-pronoun	arkaa	अर्का	DGO
Unmarked general determiner-pronoun	aruu	अरू	DGX
Infinitive verb	garnu, garna, garna, nagarnu, nagarna, nagarna	गर्नु, गर्न, गर्ना, नगर्नु, नगर्नु, नगर्ना	VI
Masculine d-participle verb	gardo, nagardo	गर्दो, नगर्दो	VDM
Feminine d-participle verb	gardii, nagardii	गर्दी, नगर्दी	VDF
Other-agreement d-participle verb	gardaa, nagardaa	गर्दा, नगर्दा	VDO
Unmarked d-participle verb	gardai, nagardai	गर्दै, नगर्दै	VDX
e(ko)-participle verb	gae (as in gae saal), gare (as in garejati or gareko)	गरे	VE
ne-participle verb::	garne, nagarne	गर्ने, नगर्ने	VN
Sequential participle-converb	garera, gariikana, garii, nagarera,	गरेर, गरीकन, गरी,	VQ

	nagariikana, nagarii	नगरेर, नगरीकन, नगरी	
Command-form verb, non-honorific	gar, jaa	गर, जा	VCN
Command-form verb, mid-honorific	gara, jaau, jaao	गर, जाऊ, जाओ	VCM
Command-form verb, high-honorific	garnos, jaanos	गर्नोस्, जानोस्	VCH
Subjunctive / conditional e-form verb	gare, nagare	गरे, नगरे	VS
i-form verb	gari	गरि	VR
First person singular verb	gare~, garthe~, garina~, chu, hu~, garnechu	गरैँ, गर्थैँ, गरिनैँ, छु, हुँ, गर्नेछु	VVMX1
First person plural verb	garyau~, garthyau~, garenau~, chau~, hau~, garnechau~	गर्यौँ, गर्थ्यौँ, गरेनौँ, छौँ, हौँ, गर्नेछौँ	VVMX2
Second person non-honorific singular verb	garis, garthis, garinas, chas, hos, garnechas	गरिस्, गर्थिस्, गरिनस्, छस्, होस्, गर्नेछस्	VVTN1
Second person plural (or medial-honorific singular) verb	garyau, garthyau, garenau, chau, hau, garnechau	गर्यौँ, गर्थ्यौँ, गरेनौँ, छौँ, हौँ, गर्नेछौँ	VVTX2
Third person non-honorific singular verb	garyo, garthyo, garena, cha, ho, garnecha	गर्यो, गर्थ्यो, गरेन, छ, हो, गर्नेछ	VVYN1
Third person plural (or medial-honorific singular) verb	gare, garthe, garenan, chan, hun, garnechan	गरे, गर्थे, गरेनन्, छन्, हुन्, गर्नेछन्	VVYX2
Feminine second person non-honorific singular verb	garlis, ches, garthis	गर्लिस्, छेस्, गर्थिस्	VVTN1F
Feminine second person non-honorific singular verb	garthyau, chyau	गर्थ्यौँ, छ्यौँ	VVTM1F
Feminine third person medial-honorific singular verb	garina, garii, che, garthii	गरिन, गरी, छे, गर्थी	VVYN1F
Feminine third person medial-honorific singular verb	garin, garthin, garinan, chin	गरिन्, गर्थिन्, गरिनन्, छिन्	VVYM1F
First person singular optative verb	jaau~, garu~	जाऊँ, गरूँ	VOMX1
First person plural optative verb	jaaau~, garau~	जाऔँ, गरौँ	VOMX2



Second person non-honorific singular optative verb	gaes, gares	गएस्, गरेस्	VOTN1
Second person plural (or medial-honorific singular) optative verb	gae, gare	गए, गरे	VOTX2
Third person non-honorific singular optative verb	jaaos, garos	जाओस्, गरोस्	VOYN1
Third person plural (or medial-honorific singular) optative verb	jaauun, garuun	जाऊन्, गरून्	VOYX2
Adverb	raamrarii, ekdam, chiTo	राम्ररी, एकदम, छिटो	RR
Demonstrative adverb	yataa, utaa, tyataa; ahile, ahilyai, yahii~, yasarii, aba	यता, उता, त्यता, अहिले, अहिल्यै, यहीं, यसरी, अब	RD
Interrogative adverb	kataa, kahaa~, kahile, kasarii	कता, कहाँ, कहिले	RK
Relative adverb	jataa, jahaa~, jahile, jasarii, jaba	जता, जहाँ, जहिले	RJ
Postposition	agaaDi, pachaaDi, baaTa, dwaaraa, maa, maathi, saath, puurvak, tira, tarpha, vasa, sanga, binaa	अगाडि, पछाडि, बाट, द्वारा, मा, माथि	II
Plural-collective postposition	<i>haruu</i>	हरू	IH
Ergative-instrumental postposition	<i>le</i>	ले	IE
Accusative-dative postposition	<i>laaai</i>	लाई	IA
Masculine genitive postposition	<i>ko</i>	को	IKM
Feminine genitive postposition	<i>kii</i>	की	IKF
Other-agreement genitive postposition	<i>kaa</i>	का	IKO
Cardinal number	<i>ek, eu#, yau#, dui, tin, caar, paa~c</i>	एक, दुई, तीन, चार, पाँच	MM
Masculine ordinal number	<i>pahilo, dosro, tesro, cautho</i>	पहिलो, दोस्रो, तेस्रो, चौथो	MOM
Feminine ordinal number	<i>pahilii, dosrii, tesrii, cauthii</i>	पहिली, दोस्री, तेस्री, चौथी	MOF

		चौथी	
Other-agreement ordinal number	<i>pahilaa, dosraa, tesraa, cauthaa</i>	पहिला, दोसरा, तेसरा, चौथा	MOO
Unmarked ordinal number	<i>paa~cau~</i>	पाँचाँ	MOX
Masculine numeral classifier	<i>#To</i> [as in euTo]	#टो	MLM
Feminine numeral classifier	<i>(#)waTii, #Tii, #oTii, #auTii</i>	(#)वटी, #टी, #ओटी	MLF
Other-agreement numeral classifier	<i>(#)waTaa, #Taa, #oTaa, #auTaa</i>	(#)वटा, #टा, #ओटा, #औटा	MLO
Unmarked numeral classifier	<i>(#)janaa</i>	(#)जना	MLX
Coordinating conjunction	<i>ra, tathaa</i>	र, तथा	CC
Subordinating conjunction appearing <b>after</b> the clause it subordinates	<i>bhanne, bhani, bhanera</i>	भने, भनी, भनेर	CSA
Subordinating conjunction appearing <b>before</b> the clause it subordinates	<i>ki, yadi, yaddyapi, kinaki</i>	कि, यदि, यद्यपि, किनकि	CSB
Particle	<i>nai, caahi~, pani, hai, ra, re, kyaare, ho, khai</i>	नै, चाहिँ, पनि, है, र, रे, क्यारे, हो, खै	TT
Question marker	<i>ke</i>	के	QQ
Interjection	<i>oho, aahaa, hare</i>	ओहो, आहा, हरे	UU
Possessive reflexive pronoun without agreement	<i>merai</i>	मेरै	PMXKX
Non-honorific second person possessive pronoun without agreement	<i>terai</i>	तेरै	PTNKX
Medial-honorific second person possessive pronoun without agreement	<i>timrai</i>	तिम्रै	PTMKX
Possessive reflexive pronoun without agreement	<i>aaphnai</i>	आफ्नै	PRFKX
Unmarked genitive postposition	<i>kai</i>	कै	IKX
Sentence-final punctuation		? ! .	YF

Sentence-medial punctuation	, ; : :- / -		YM
Quotation marks	' "		YQ
Brackets	( ) { } [ ]		YB
Foreign word in Devanagari			FF
Foreign word, not in Devanagari			FS
Abbreviation	M.P.P.	म. पु. पु.	FB
Mathematical formula (and similar)	$e=mc^2$		FO
Letter of the alphabet			FZ
Unclassifiable			FU
Null tag: an element of the text which does not need a tag	<p>		NULL