

DISAMBIGUATING TIBETAN VERB STEMS WITH MATRIX VERBS IN THE INDIRECT INFINITIVE CONSTRUCTION

EDWARD GARRETT, NATHAN W. HILL, ABEL ZADOKS
SOAS, University of London

Introduction

A great deal of digitized Tibetan texts are now available online, including the entire Derge Kanjur¹ and over 60 Dunhuang documents.² Nonetheless, there is a dearth of tools to access and process these data efficiently. To better exploit the profusion of available material, a research project at the School of Oriental and African Studies (SOAS), University of London, is currently working to create a part-of-speech tagged corpus together with an automatic word segmenter and part-of-speech tagger.³ Garrett et al. (forthcoming a) describes the part-of-speech tag-set used in this project. Garrett et al. (forthcoming b) describes Version 1.0 of the rule-based tagger, designed for Classical Tibetan materials. In the long term, the rule-based tagger will be combined with a statistical tagger to achieve improved results.

Our project began by hand-tagging an initial 17,522 words of the *Mdzañs blun*. We developed the initial part-of-speech tag set during this phase. In the next phase, covering the next 26,937 words of the *Mdzañs blun* and the first 32,083 words of the *Mi la ras pañi rnam thar*, we developed the rule-based tagger through an ad hoc process of trial and error. As of 8 June, 2014, the hand-annotated corpus contains 47,927 words of the *Mdzañs blun*, 32,083 words of the *Mi la ras pañi rnam thar*, and 10,069 words of the *Bu ston chos ḥbyun*. The rule-based tagger operates across this corpus with an accuracy of 0.998 and an ambiguity of 1.360.

In addition to the training corpus of hang-tagged data, there is a larger test corpus includes all of these three texts and a number of other

1 Available at <http://www.thlib.org/encyclopedias/literary/canons/kt/catalog.php#cat=d> on 10 June 2014.

2 Available at otdo.aa.tufs.ac.jp on 10 June 2014.

3 The project 'Tibetan in Digital Communication' is funded by the U.K's Arts and Humanities Research Council.

texts digitized by the University of Otani,⁴ including the *Sa-skya legs-bśad*, *Hgro-ba bzañ-moñi rnam-thar*, *Sa-paṅ-gyi rnam-thar gsuñ-sgros-ma*, *Gśen ñi-mañi rnam-thar*, *Rwa-loñi rnam-thar*, and *Mar-pañi rnam-thar*. This larger corpus permits the testing of hypotheses against a larger data set than the comparatively small hand-tagged training corpus.

The rule-based part-of-speech tagger uses the grammatical rules known to any first year Tibetan student to preclude unlikely part-of-speech interpretations. For example, the tagger knows that the syllable *-so* when it occurs after a verb stem that ends in *-s* and before the *śad* punctuation mark is not the noun 'tooth'. Unfortunately, for certain verb forms it is not possible in all cases to specify an unambiguous tense analysis.⁵ The circumstances giving rise to tense ambiguity are best illustrated with an example. The verb *gśegs* 'go' is invariant across all four tenses. Often syntactic cues disambiguate the correct tense (e.g. *gśegs śig* must be the imperative), but in other contexts disambiguation is not univocal. In the phrase *mi gśegs* the verb *gśegs* is either a present (cf. *mi byed*) or a future (cf. *mi bya*), but cannot be understood as a past. Thus, the tag [v.fut.v.pres] species that in this and comparable contexts it is impractical to decide between [v.fut] and [v.pres]. Finally, there are contexts such as *gśegs śiñ* and *gśegs so*, in which it is only possible to say that *gśegs* is not the imperative (cf. *byed ciñ*, *bya žiñ*, *byas śiñ*; and *byed do*, *byaño*, *byas so*). Rather than tagging such contexts with the lengthy [v.fut.v.past.v.pres] we instead employ the tag [v.invar].⁶ Following these protocols the rule-based tagger disambiguates verb stems in those places where they can be disambiguated with certainty, but leaves them ambiguous in cases where the interpretation is not clear-cut.

One of the main purposes of the creation of the corpus and tagger is to yield new insights into Tibetan grammar. Consequently, it would be foolish to rest content with those facets of verb stem distribution well known to the first year Tibetan student. Instead, the tools created so far should be harassed to discover new patterns, patterns which in turn can be fed back into refining the corpus and tagger. The current version of

4 Available at <http://web.otani.ac.jp/cr/twrrp/project/otet/> on 10 June 2014.

5 In this paper the term 'tense' is used to refer to the distinct four principal parts of verbs used in the indigenous grammatical tradition. This terminology is not intended to imply that the morphosyntactic categories recognized by the indigenous tradition correspond semantically to 'tense' (as opposed to 'aspect' or 'mood') as it is used in linguistic typology.

6 One must bear in mind that use of the tag [v.invar] is not a positive claim that a verb is (morphologically or otherwise) invariant, but rather is the negative claim that the stem of this verb in this context cannot be more precisely stated.

the tagger already includes one such new pattern, namely the prohibition of the future stem before the converb *-nas*. Although we introduced this rule on the basis of anecdotal evidence, it appears to be quite robust. For example, all instances of *bya nas* in the Derge Kanjur appear to involve either *bya* 'bird' or *nas* 'barley'. The goal of the current investigation is to see whether verb stems may be disambiguated before the terminative converb in the indirect infinitive construction (e.g. *sloñ-du hoñs* 'came to ask').

To explore this question we gathered evidence by searching the test corpus for the following types of part-of-speech patterns (cf. appendix).

[v.pres] [cv.term]
 [v.fut] [cv.term]
 [v.fut.v.pres] [cv.term]
 ...

We subsequently collated the results of these searches according to the rightward context to list the subcategorization patterns attested with each matrix verb, both in its stem form (e.g. *byed*) and its nominalized form (*byed-pa*).

The subcategorization patterns of positive matrix verbs

This section discusses the various patterns of subcategorization found with unnegated verbs. Verbs appear to fall into three classes, those that take future and present stems with equal preference, those that select for the present, and those that select for the future. Although the corpus does provide examples of subordinated past stems, this phenomenon appears to be unsystematic.

Eight matrix verbs occur with subordinate present and future stems and with stems that are ambiguous between present and future.

bcug(-pa) (24 v.pres, 20 v.fut, 6 v.fut.v.pres), *gzug(-pa)* (5 v.pres, 1 v.fut, 4 v.fut.v.pres), *hjug* (3 v.pres, 1 v.fut), *chug* (1 v.pres, 1 v.fut.v.pres)
soñ(-ba) (14 v.pres, 4 v.fut, 12 v.fut.v.pres), *hgro* (11 v.pres, 2 v.fut, 1 v.fut.v.pres), *phyin(-pa)* (9 v.pres, 3 v.fut, 2 v.fut.v.pres)
btan(-ba) (20 v.pres, 8 v.fut, 12 v.fut.v.pres), *gtoñ* (1 v.pres, 1 v.fut.v.pres)
byon(-pa) (5 v.pres, 5 v.fut, 10 v.fut.v.pres), *hbyon-pa* (1 v.pres)
gyur(-pa/-ba) (3 v.pres, 3 v.fut, 9 v.fut.v.pres), *hgyur* (3

v.fut.v.pres)
byuñ (2 v.pres, 1 v.fut, 8 v.fut.v.pres)
doñ-ba (1 v.pres, 1 v.fut, 1 v.fut.v.pres), *ḥdoñ* (1 v.fut.v.pres)
smra-ba (1 v.pres, 1 v.fut), *smras* (1 v.fut.v.pres)

Six matrix verbs appear to select only for the present stem. In the corpus there are examples of these matrix verbs selecting presents and selecting verb stems that are ambiguous between present and future, but there are no examples of these six verb selecting unambiguous future stems.

ḥoñs(-pa) (12 v.pres, 7 v.fut.v.pres), *ḥoñ(-ba)* (11 v.pres, 5 v.fut.v.pres), *yoñ-ba* (4 v.fut.v.pres), *yoñs-pa* (1 v.pres, 1 v.fut.v.pres)
mdzad(-pa) (1 v.pres, 13 v.fut.v.pres)
byas(-pa) (2 v.pres, 2 v.fut.v.pres)
btags-pa (1 v.pres, 2 v.fut.v.pres)
bzuñ-pa (1 v.pres, 1 v.fut.v.pres), *ḥdzin-pa* (1 v.fut.v.pres)
mchi (1 v.pres, 1 v.fut.v.pres)

One must however bear in mind that it is always possible the absence of unambiguous future stems in the subordinate clause might be an accident gap in the corpus, rather than a structural fact about the Tibetan language. It is only the verb *ḥoñ/yoñs* which is well attested enough that it seems quite likely that the occurrence of future stems in the subordinate clause can be precluded.

There are also six matrix verbs that appear to select only the future stem. In the corpus there are examples of these matrix verbs selecting futures and selecting verb stems that are ambiguous between future and present, but there are no examples of these six verb selecting unambiguous present stems.

gsol (19 v.fut, 7 v.fut.v.pres)
med(-pa) (7 v.fut, 8 v.fut.v.pres)
grags(-pa) (2 v.fut, 1 v.fut.v.pres)
yod-pa (2 v.fut, 2 v.fut.v.pres)
ruñ-ba (1 v.fut, 1 v.fut.v.pres)
gśegs-pa (1 v.fut, 1 v.fut.v.pres)

The fact that these six verbs show no stem changes lends credence to their interpretation as a structural class. The presence of both *med* and *yod* in this category cannot be a coincidence. In contrast, in light of the fact that motion verbs general appear among those verbs that take the present and future with equal preference, it seems likely that a lack of *gśegs* selecting a present stem is an accidental gap.

Turning to the past stems in subordinate clauses, there are some surprises. Although past stems are the most common stem in the corpus overall, only four verbs occur more than once as matrix verbs selecting the past stem: 1. *gsol(-ba)* (9), 2. *soñ* (2), *phyin* (1), *hgro-ba* (1), 3. *hoñ* (2), *hoñs-pa* (1), 4. *yod(-pa)* (3). All four have been discussed previously. The first (*gsol*) and fourth (*yod*) generally selects for the future. The second (*hgro/phyin/soñ*) selects for either present or future. The third (*hoñ*) selects the present. If one looks at the actual occurrences of these four verbs when they appear to select the past stem, the past stems in question turn out to be misspellings of the future stem.

bzuñ (gzuñ) du gsol (5)
bcad (gcad) du gsol (1)
bcad (gcad) du hoñs-pa (1)
bsad (gsad) du hoñ (1)
byuñ (hbyuñ) du soñ (1)
byuñ (hbyuñ) du hgro-ba (1)
byuñ (hbyuñ) du yod-pa (2)

This misspelling of futures as pasts is not unexpected, since the spoken languages lose the future and often the etymological future is homophonous in a given dialect with an etymological past because of sound change. In some cases the sandhi makes clear that a future was intended.

gyur du (past requires *-tu*) *hoñ* (1)
bsgrubs tu (past requires *-su*) *gsol-ba* (1)
bsgrubs tu (past requires *-su*) *yod* (1)

All of these cases merit more detailed philological discussion; here it suffices to conclude that there is no solid evidence for a class of matrix verbs that select the past stem in the subordinate clause.

For the purposes of verb stem disambiguation, one may conclude that if a verb has any interpretation other than [v.past] and it occurs as the subordinate verb of this infinitive construction, then the analysis of this verb as [v.past] may be removed. One may more tentatively conclude that a verb that has interpretations other than [v.fut] (and [v.past]) should have [v.fut] removed, if the matrix verb is one of the six verbs that selects the present, although perhaps this stipulation is only safely applied to *hoñ* and its forms. Because there are few futures in the corpus in general the lack of a future before the other verbs seven present selecting verbs may be a coincidence.

Reciprocally, if the matrix verb is one of the six future selecting verbs, then the analysis of the subordinate verb as [v.pres] can be precluded. In this case, the overall prevalence of presents over futures in the corpus means that the lack of presents before these verbs can be safely understood as a systematic gap. Furthermore, the fact that all future selecting verbs are invariant in their stem morphology, and that *yod* and *med* both occur in this category, bolsters the security of the analysis. Nonetheless, because *gšegs* is the only motion verb in this group, and most motion verbs can select either present or future stem, it would be a mistake to forbid the analysis of a verb subordinate to *gšegs* as a present.

The subcategorization patterns of negative matrix verbs

Because negated verbs are less frequent than positive verbs, the evidence of negated matrix verbs in the indirect infinitive construction is sparse. However, the investigation of subcategorization patterns with positive verbs provides a framework to analyze these data, i.e. if the evidence of negated verbs confirms the patterns seen with the unnegated verbs then the pattern in question more secure.

Four of the eight verbs that select for both present and future are attested with negation.

hjug-pa (1 v.pres), *bcug(-pa)* (2 v.pres, 1 v.fut.v.pres)
byuñ (1 v.pres, 1 v.fut)
phyin (1 v.fut.v.pres)
gtoñ-ba (1 v.fut.v.pres)

These data support the perspective that these verbs select both present and future, although they suffice to conclude this only in the cases of *byuñ*.

Turning to the verbs that in the positive appear to select for the present, the evidence from negated matrix verbs is very thin.

byas (1 v.pres)
yoñ (1 v.fut), *hoñs* (1 v.fut.v.pres), *hoñ* (1 v.fut.v.pres)

The occurrence of *yoñ* with an unambiguous subordinate future weighs against the theory that this verb selects for the present. The deviant examples is *bsgrub tu ma yoñ-ba* 'did not come to accomplish'. One may perhaps suspect a spelling mistake of *bsgrub* for *sgrub*.

The verb *ruñ* is the only verb that selects for future for which negated examples occur in the corpus. In fact, this verb occurs far more frequently negated than in the positive form. Most of the negated forms

support the analysis of this verb as selecting a future stem.

ruñ (11 v.fut, 4 v.fut.v.pres, 2 v.fut.v.past)

However, there are two problem cases. In one of these cases *ruñ* selects the past stem *byuñ*, which is possibly a simple misspelling for *ḥbyuñ*. The other case is a genuine problem. The passage is *gcod du mi ruñ* with an unambiguous present. So, paradoxically although (*mi*) *ruñ* selecting the future is among the most robust patterns in the corpus, *gcod du mi ruñ* is also one of the most clear-cut counter-examples to any pattern identifiable in the corpus.

A number of matrix verbs that predominantly appear in the negative were not noted in the discussion of positive matrix verbs. The verb *chud* appears to only select for the present.

chud (1 v.pres, 1 non-negated v.pres, 7 v.fut.v.pres)

Three verbs appear to select the present and future with equal ease.

btub(-pa) (2 v.pres, 1 v.fut, 10 v.fut.v.pres), *gtub-pa* (1 v.pres)

ḥdod(-pa) (1 v.pres, 3 v.fut.v.pres, 1 v.fut.v.past)

snañ (1 v.pres, 1 non-negated v.fut)

Concerning the verb *tshud*, one can only say that it does not select the past, which is not surprising since no verbs appear to select the past.

tshud (3 v.fut.v.pres)

The verb *gnañ* has a strange pattern of subcategorization, appearing to select the past in two cases (once selecting *phyuñ* and once *byuñ*).

gnañ (1 v.fut, 2 v.past)

The verbs *snañ*, *tshud*, and *gnañ* are not well enough attested in the corpus to yield any reliable conclusion.

Conclusions

With a certain fuzziness around the edges, the data surveyed here suggest that past tense verbs do not occur as the subordinate verbs of indirect infinitives and that the matrix verbs *gsol*, *med*, *grags*, *yod*, *ruñ* select the future tense in the subordinate clause. It is possible that one group of verbs selects the present tense whereas others are equally happy to select the present and the future, but the overall rarity of future stems in the corpus makes the line between these two categories difficult to draw.

REFERENCES

- Garrett, Edward, Nathan W. Hill, Adam Kilgarif, Ravikiran Vadlapudi, Abel Zadoks (forthcoming a). "The contribution of corpus linguistics to lexicography and the future of Tibetan dictionaries." *Revue d'Etudes Tibétaines*.
- Garrett, Edward, Nathan W. Hill, Abel Zadoks (forthcoming b). "A Rule-based Part-of-speech Tagger for Classical Tibetan." *Himalayan Linguistics*.

APPENDIX: MATRIX VERBS IN INDIRECT INFINITIVE CONSTRUCTION

[v.pres] [cv.term]

<i>bcug(-pa)</i> (24), <i>gzug(-pa)</i> (5),	<i>chud</i> (1)
<i>hjug</i> (3), <i>chug</i> (1)	<i>mchi</i> (1)
<i>son(-ba)</i> (14), <i>hgro</i> (11), <i>phyin(-pa)</i> (9)	<i>gñag-ba</i> (1)
<i>hoñs(-pa)</i> (12), <i>hoñ(-ba)</i> (11),	<i>btags-pa</i> (1)
<i>yoñs-pa</i> (1)	<i>doñ-ba</i> (1)
<i>btañ(-ba)</i> (20), <i>gtoñ</i> (1)	<i>drañs-pa</i> (1)
<i>byon(-pa)</i> (5), <i>hbyon-pa</i> (1)	<i>smra-ba</i> (1)
<i>gyur(-pa)</i> (3)	<i>btsugs</i> (1)
<i>hchad(-pa)</i> (2), <i>bśad</i> (1)	<i>brtsams-pa</i> (1)
<i>tshor-ba</i> (2)	<i>brtsis-pa</i> (1)
<i>bsgyur-pa</i> (1), <i>bsgyur-ba</i> (1)	<i>mdzad</i> (1)
<i>byuñ</i> (2)	<i>brdzañs-pa</i> (1)
<i>byas(-pa)</i> (2)	<i>bzuñ-pa</i> (1)
<i>byiñ(-ba)</i> (2)	<i>len</i> (1)

[v.fut] [cv.term]

<i>bcug-pa</i> (20), <i>gzug-pa</i> (1), <i>hjug-pa</i> (1)	<i>grags(-pa)</i> (2)
<i>gsol</i> (19)	<i>yod-pa</i> (2)
<i>son(-ba)</i> (4), <i>phyin-pa</i> (3), <i>hgro</i> (2)	<i>doñ-ba</i> (1)
<i>btañ(-ba)</i> (8)	<i>nus</i> (1)
<i>med(-pa)</i> (7)	<i>snañ</i> (1)
<i>byon(-pa)</i> (5)	<i>byuñ</i> (1)
<i>gyur(-ba)</i> (3)	<i>smra</i> (1)
<i>bsdus-pa</i> (2)	<i>zugs</i> (1)
	<i>ruñ-ba</i> (1)
	<i>gśegs-pa</i> (1)

[v.fut.v.pres] [cv.term]

<i>son</i> (12), <i>phyin</i> (2), <i>hgro</i> (1)	<i>gyur(-pa)</i> (9), <i>hgyur</i> (3)
<i>btañ(-pa)</i> (12), <i>gtoñ</i> (1)	<i>hoñs(-pa)</i> (7), <i>hoñ(-ba)</i> (5), <i>yoñ-</i>

ba (4), *yoñs-pa* (1)
bcug(-pa) (6), *gźug(-pa)* (4),
chug (1)
mdzad(-pa) (13)
byon(-pa) (10)
gsol (8)
byuñ (8)
med(-pa) (8)
bkod(-pa) (4), *hgod-pa* (1)
cha-ba (3), *chas-pa* (1)
phebs(-pa) (3)
bzuñ (1), *hđzin-pa* (1)
doñ-ba (1), *hdoñ* (1)
khrid-pa (2), *hkhrid* (1)
mkhyen (2)
ñal (2)
btags (2)
byas (2)
yod-pa (2)
bskyañs (1)
khol-ba (1)
grags (1)
rgyas (1)
bgyi-ba (1)
rgyugs-pa (1)

ñes (1)
mchi (1)
hchos-pa (1)
ñe (1)
rtogs-pa (1)
mthoñ (1)
btub (1)
ltem (1)
dod (1)
phul (1)
phed (1)
hbab (1)
smras (1)
btsal (1)
stsol-ba (1)
hđzud (1)
u (1)
bzo (1)
ruñ-ba (1)
śes-pa (1)
gśegs (1)
lhags-pa (1)
lhuñ-ba (1)

[v.past] [cv.term]

gsol(-ba) (9)
soñ (2)
hoñ (2), *hoñs-pa* (1)
yod(-pa) (3)
mid-pa (1)
ñe-ba (1)
lhuñ-ba (1)
byin-pa (1)
bsdus-pa (1)
byas-pa (1)
gsuñs-pa (1)

byon-pa (1)
skye-ba (1)
snyoms-pa (1)
mñah (1)
hgril (1)
bcad (1)
sbyar (1)
bcug (1)
hkhruñs (1)
spyod (1)

[v.fut.v.past] [cv.term]

gsol (11)
med(-pa) (2)

yoñ (1), *yoñs-pa* (1)
śog (2)

btuñ (1)*bzugs* (1)*byuñ* (1)*[v.pres] [cv.term] [neg]**ħjug-pa* (1), *bcug(-pa)* (2)*byas* (1)*btub* (2), *gtub-pa* (1)*byuñ* (1)*ruñ* (1)*snañ* (1)*dbañ* (1)*chud* (1)*ħdod* (1)*ster-ba* (1)*[v.fut] [cv.term] [neg]**ruñ* (11)*zad-pa* (1)*gnañ* (1)*yoñ* (1)*btub* (1)*byuñ* (1)*bzad-pa* (1)*[v.fut.v.pres] [cv.term] [neg]**btub(-pa)* (10)*phyin* (1)*chud* (7)*e* (1)*ruñ* (4)*stsal* (1)*tshud* (3)*dad-pa* (1)*ħdod(-pa)* (3)*gtoñ-ba* (1)*ħoñs* (1), *ħoñ* (1)*bcug-pa* (1)*[v.past] [cv.term] [neg]**gnañ* (2)*ltuñ-ba* (1)*ruñ* (1)*[v.fut.v.past] [cv.term] [neg]**ruñ* (2)*re* (1)*ħdod-pa* (1)