

himalayan linguistics

A free refereed web journal and archive devoted to the study of the
languages of the Himalayas

Himalayan Linguistics

Algorithmic description of the decomposition and checking of a Classical Tibetan written syllable

Elie Roux

Buddhist Digital Resource Center

Hélios Hildt

Buddhist Digital Resource Center

ABSTRACT

This document presents our research on the the correct formation of a Classical Tibetan syllable. It was triggered by attempts at defining the boundaries of well-formed syllables in Classical Tibetan for spell checking purposes.

Formalizing the formation of the syllable led us to inspect the small differences among grammar books, both in Western and Tibetan language. We then checked these differences against the Tibetan dictionaries we consider reliable, and also against the Kangyur.

Our inquiry finally led us to study the way to decompose a syllable, discussing the ambiguous cases, as well as the formation of the Dzongkha syllable.

KEYWORDS

Tibetan, spell checking, NLP

This is a contribution from *Himalayan Linguistics*, Vol. 17(1): 50-66.

ISSN 1544-7502

© 2018. All rights reserved.

This Portable Document Format (PDF) file may not be altered in any way.

Tables of contents, abstracts, and submission guidelines are available at
escholarship.org/uc/himalayanlinguistics

Algorithmic description of the decomposition and checking of a Classical Tibetan written syllable

Elie Roux
Buddhist Digital Resource Center

Hélios Hildt
Buddhist Digital Resource Center

1 Introduction

This document presents our research on the orthography of a Classical Tibetan syllable. The study of this question was triggered by attempts at defining the boundaries of well-formed written syllables in Classical Tibetan for NLP purposes. What we are trying to assess is the extend of the possibilities of which a modern Tibetan person would think as valid if they wanted to invent a new Classical Tibetan word. This is a necessary first step to be able to discriminate valid Classical Tibetan from other forms such as transliterated Sanskrit, or just mistakes. Of course, the syllable level is not enough for a full feature spell checker; but we hope our work will provide a basis for future work in this area, as well as a complete description of the various corner cases of syllable construction. This research has already been useful, allowing us to create

- a basic spell checker for Classical Tibetan using the *hunspell* library
- rules for collation of Classical Tibetan

Formalizing the formation of the syllable led us to inspect the small differences between grammar books, both in Western and Tibetan language¹. We then checked these differences against some chosen Tibetan dictionaries² and the Derge Kangyur³.

Our inquiry finally led us to study the way to decompose a syllable, discussing the ambiguous cases, as well as the formation of the Dzongkha syllable.

¹ Our main readings in Western languages were (Tournadre and Dorje, 2005) and (Beyer, 1992), but we also consulted (Bacot, 1946) and (Das, 1915); the Tibetan sources were (འཇུ་པ་མཚན་མོན་གྱི་འབྲུང་གནས།, XVIIIth c.) and (ཚུ་ཉན་ཞབས་བྱུང་།, 2003). The bibliography used in this document can be downloaded on <https://github.com/eroux/tibliography/>.

² The dictionaries we consulted are (Dorje et al., 2003), (མོན་ལམ།, 2016), (དུང་དཀར་ཆོ་བཟང་འབྲིན་ལས།, 2002), and in a lesser extent (Yisun, 1985), (ཚོས་སྐད་རྫོང་ཁ་ཚོགས་མཚན་མོ།, 2010) and (Negi, 1993). We also consulted (Hill, 2010).

³ (ཚོས་གྱུ་འབྲུང་གནས།, 1721), input by Esukhia, <https://github.com/Esukhia/derge-kangyur>

2 Elements of a Tibetan syllable

2.1 Description and vocabulary

2.1.1 Basic description

For the description of the different elements of a Tibetan syllable, we will take the vocabulary used in (Tournadre and Dorje, 2005). As an example, decomposing the syllable འཇུངས་ results in radical letter འ, prefix འ, superscribed འ, subscribed འ, vowel accent འ, first suffix འ and second suffix འ. Among these six categories, only the radical letter is mandatory, all the others are optional⁴. This constitutes a good framework for the description of the elements of a syllable but lacks a few elements.

2.1.2 Missing elements

First, some syllables (such as འཇུངས་) contain a wasur (འ). These are traditionally considered as subscribed letters, with the ability to combine with other subscribed (as in འཇུངས་). This last case breaks our initial description because it requires two optional subscribed letters instead of one, and the second can only be a wasur. Our proposal is to treat the wasur separately from the subscribed letters in order to keep things simple. So འཇུངས་ would be decomposed as radical letter འ, wasur, first suffix འ, second suffix འ; and འཇུངས་ as radical letter འ, subscribed འ and wasur.

A second problem is the འ “suffix” (as in འཇུངས་, “chapter”, /leu/, 2 syllables): it not considered as a suffix in traditional grammars (that are usually completely silent about it), and does not have normal suffix properties. We will call it a *special suffix*. At most one can appear in a syllable.

The third missing feature is affixed particles. These are appended to syllables with no suffixes or replace suffix འ. For instance འཇུངས་ is decomposed as radical letter འ and affixed particle འ. Two affixed particles can even be combined, as in འཇུངས་.

In order to describe the different elements, we will call *final part* what is written on the right side of the radical letter; *main stack* the radical and everything written above and below except vowel accents (so subscribed, superscribed letters and wasur).

We also would like to introduce the notion of *root*, which is constituted of everything preceding the vowel accent or final part. This notion will come handy in some parts of this document, and is crucial for collation purposes.

2.1.3 Formalization

Now we can propose a formalization of the elements of the syllable with some symbols coming from the regular expressions: “?” means an element that can either be omitted or appear once, and “|” simply means an exclusive “or”:

$$\begin{aligned}
 \text{syllable} &= [\text{root}][\text{vowel}]?[\text{final part}]? \\
 \text{root} &= [\text{prefix}]?[\text{main-stack}] \\
 \text{main-stack} &= [\text{superscribed}]?[\text{radical}][\text{subscribed}]?[\text{wasur}]? \\
 \text{final-part} &= ([\text{special-suffix}][\text{affixed-particles}]?) \mid [\text{suffixes}] \mid [\text{affixed-particles}] \\
 \text{affixed-particles} &= [\text{affixed-particle}][\text{affixed-particle}]?
 \end{aligned}$$

⁴ For the sake of simplicity, we will not consider that the “a” vowel accent or the first suffix འ are implied when no vowel accent or first suffix is present.

suffixes = [first-suffix][second-suffix]?

Note that although this is the way a syllable is built, some syllables can have multiple decomposition possibilities according to this scheme. For instance བར་ can be decomposed as either:

- radical བ + suffix ར (/bar/, “between”)
- radical བ + affixed particle ར (either /war/ nominalizer+oblique mark, or /bar/ “cow”+oblique mark)

This schema also doesn’t account for implied suffix འ: for instance if དཀའི་ did not have its affixed particle, it would be དཀའ་, and for some purposes it may also be useful to retain this information, but we chose the most simple schema for readability purposes.

2.2 Constraints on the construction of a syllable

In this part we will study the different possibilities of construction of these elements, first for each element taken independently, then the constraints of their relations.

2.2.1 Simple constraints

All consulted sources agree on all the constraints on the different elements except the special suffix. These constraints are:

- prefix can be ཀ ད བ མ OR འ
- first suffix can be ཀ ང ཅ ཆ ཇ ཈ ཉ ཐ ད དྷ ན པ ཕ བ བྷ མ ཙ ཐ འ OR ས
- second suffix can be ས⁵
- superscribed can be ར ལ OR ས
- subscribed can be ུ ཱ OR ི
- radical letter can be ཀ ཁ ག གྷ ང ཅ ཆ ཇ ཈ ཉ ཐ ད དྷ ན པ ཕ བ བྷ མ ཙ ཐ འ ཡ ར ལ ཤ ཥ ས OR མ
- vowel accent can be ི ུ ཱ ི OR ི
- wasur is ུ
- affixed particle can be འ ར ལ ཤ ཥ ས OR ས

When inspecting the different constraints on the relation between elements, we can see some consistent lists among all sources:

- ཀ ཁ ག གྷ ང ཅ ཆ ཇ ཈ ཉ ཐ ད དྷ ན པ ཕ བ བྷ མ ཙ ཐ འ ཡ ར ལ ཤ ཥ ས for superscribed ར + radical letter
- ཀ ཁ ག གྷ ང ཅ ཆ ཇ ཈ ཉ ཐ ད དྷ ན པ ཕ བ བྷ མ ཙ ཐ འ ཡ ར ལ ཤ ཥ ས for superscribed ལ + radical letter

⁵ Old Tibetan also has ད.

- སྐ སྐྱ སྐྲ སྐླ སྐྴ སྐྵ སྐྶ སྐྷ སྐྸ སྐྐྵ for superscribed ས + radical letter
- ལྱ ལྲ ལླ ལྴ ལྵ ལྶ ལྷ ལྸ ལྐྵ for radical letter + subscribed ལ
- སྐ སྐྱ སྐྲ སྐླ སྐྴ སྐྵ སྐྶ སྐྷ སྐྸ སྐྐྵ for radical letter + subscribed ལ
- ལྱ ལྲ ལླ ལྴ ལྵ ལྶ ལྷ ལྸ ལྐྵ for superscribed ས + radical letter + subscribed ལ
- ལྱ ལྲ ལླ ལྴ ལྵ ལྶ ལྷ ལྸ ལྐྵ for superscribed ལ + radical letter + subscribed ལ

2.2.2 Methodological choices

But other constraints are inconsistent among sources. Before we expose the constraints we see as the most relevant, we must expose our methodological choices.

First, the language we deal with is what is usually called Classical Tibetan; and we take it in a modern perspective: XXIth century Classical Tibetan. This leads us to consider that some old texts contain errors even if they follow what the rules of Classical Tibetan were at their period.

The syllables we are inspecting are what is sometimes referred to as གྲིམས་མཐུན་གྱི་ཚེག་བཟུང་།, literally “legal syllables”. While there is no clear cut definition of this concept, we exclude syllables appearing solely in 1. transliterated foreign languages, 2. proper nouns⁶, 3. misspellings, 4. words in other languages such as Zhang zhung language (འདྲ་ལྷོ་སྐད་) or Dzongkha; but we chose to include 1. onomatopoeia (སྐྱ་སྐད་), 2. regional orthographies (ཡུལ་སྐད་), 3. words coming from Old Tibetan (བརྗེ་རྒྱུ་བྱེད་), 4. words that are alternate spellings of other words. When a syllable comes from such a word, it will be suffixed by (OT) for Old Tibetan, (O) for onomatopoeia, (R) for regional words or (AS) for alternate spelling.

There is a whole spectrum of choices that can be taken for this kind of research, on one side of the spectrum we could have taken all the syllables in a purely Tibetan dictionary such as (Dorje et al., 2003) and consider that only these are legal syllables. On the other side we can take the most inclusive constraints on the formation of a legal syllable and state that all the combinations respecting these constraints are legal. We chose an in-between solution: quite strict constraints and listing exceptions. This has the inconvenience to imply some subjectivity in the words we include or exclude based on the supposed relevance of the sources, but we hope that the article will go in enough details to allow anyone to adapt this strategy to their own needs.

2.2.3 Constraints on the subscribed ས

The subscribed ས has the following possible combinations with radical letters in all the grammars we consulted: ལྱ ལྲ ལླ ལྴ ལྵ ལྶ ལྷ ལྸ ལྐྵ and ལྱ. But some sources list additional possibilities: ལྱ ལྲ ལླ ལྴ ལྵ ལྶ ལྷ ལྸ ལྐྵ and ལྱ. Among these, ལྱ and ལྲ are extremely rare and appear only in the following words: ལྱ་ལྱ་ (“mango”) and ལྱེག་ལྱེག་(O), so we consider the syllables ལྱ་ and ལྱེག་ as exceptions⁷. ལྲ is present in many words so we can count it as regular. ལླ is present only in ལླ་ (“ape”), ལླེས་ལླེས་ (“powder”) and ལླ་ལླ་ (name of a particular

⁶ Even the very common person names ལྱ་ལྱ་ and ལྲ་ལྲ་ are not considered “legal syllables”, so we consider other proper nouns shouldn’t. We do not consider the names of plants or lunar mansions as proper nouns when they do not come directly from a foreign language.

⁷ An “exception” here means that no other syllable can be built with these stacks, for instance ལྱ་ལྱ་ would not be considered legal.

design), so we consider the syllables ཏ and ཏེས as exceptions. We have not been able to find any word with the main stack ཏ except proper nouns, so we don't consider it as a possibility at all.

For combinations with radical letter and superscribed ས, all sources indicate ཏེས, ཏེས་ཀྱི་མཚན་ལོ་མཚན་ and ཏེས་ཀྱི་མཚན་ལོ་མཚན་. Some indicate ཏེས་ཀྱི་མཚན་ལོ་མཚན་ and ཏེས་ཀྱི་མཚན་ལོ་མཚན་ is present only in rare words such as ཏེས་ཀྱི་མཚན་ལོ་མཚན་ (two lunar mansions), and ཏེས་ཀྱི་མཚན་ལོ་མཚན་ (OT, “disordered”); ཏེས་ཀྱི་མཚན་ལོ་མཚན་ is only found in ཏེས་ (‘‘speak’’), ཏེས་ཀྱི་མཚན་ལོ་མཚན་ (OT, ‘‘speak’’) and derivatives: ཏེས་ཀྱི་མཚན་ལོ་མཚན་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་, ⁸ so we consider only these exceptions as valid for ཏེས and ཏེས་.

One of the sources not citing ཏ nor ཏེས as likely to be used in new words is (Beyer, 1992), invoking phonological reasons.

2.2.4 Constraints on the wasur

Wasurs explanation is even more different among sources. All the source for main stacks with wasurs give at least ཏ ཏེས ཏེས་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ and ཏེས་ཀྱི་མཚན་ལོ་མཚན་. It is also possible to find ཏ, ཏེས and ཏེས་ in some others. Our research showed that the word ཏེས་ (contraction of ཏེས་ཀྱི་མཚན་ལོ་མཚན་) can also be found in dictionaries.

The main problem with these descriptions is that they are too loose: among the many possibilities of syllables containing a wasur, so few exist that it seems easier to see them as exceptions. Also, the wasur does not seem to be used in the creation of new words⁹, so we chose to treat the syllable with wasurs as exceptions and list them. We have been able to list the following 25 syllables:

¹⁰

ཏེས་(O) ཏེས་(O) ཏེས་ (last three: O, AS for ཏེས་)
ཏེས་ ཏེས་ (in ཏེས་ཀྱི་མཚན་ལོ་མཚན་, OT for ཏེས་ཀྱི་མཚན་ལོ་མཚན་, “honey”) ཏེས་ཀྱི་མཚན་ལོ་མཚན་.

2.2.5 Constraints on prefix+radical

The following rules on prefixes are found in all grammars:

- ཏ can only be prefix of ཏ ཏེས ཏེས་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ OR ཏ
- ཏ of ཏ ཏེས ཏེས་ ཏེས་ OR ཏ
- ཏ of ཏ ཏེས ཏེས་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ OR ཏ
- ཏ of ཏ ཏེས ཏེས་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ OR ཏ
- ཏ of ཏ ཏེས ཏེས་ ཏེས་ཀྱི་མཚན་ལོ་མཚན་ OR ཏ

But (Tournadre and Dorje, 2005) gives other possibilities for prefix ཏ combined with ཏ ཏེས ཏེས་ and ཏེས་. This difference can easily be explained: in these cases, ཏ cannot be prefix of the radical letter when there is no superscribed (ex: ཏེས་ is impossible), it only can when a superscribed is present (ex: ཏེས་).

⁸ These possibilities are found in the various sources of (Hill, 2010).

⁹ (Tsering, 2006) does not list any word with wasur.

¹⁰ We did not include ཏེས་, ཏེས་ (in proper nouns), ཏེས་, ཏེས་ (Zhang zhung language), ཏེས་, ཏེས་, ཏེས་, ཏེས་ (Chinese), ཏེས་ (given only in ཏེས་ཀྱི་མཚན་ལོ་མཚན་, 1979)), ཏེས་ ཏེས་ nor ཏེས་ (Dzongkha).

3 Decomposing a syllable

3.1 Finding the main stack

What we propose here is to find the main stack in any valid syllable. Once the main stack is found, a set of rules described in a later section can be applied to find the different elements. Some rules found in (Tournadre and Dorje, 2005) describe a general way for doing that but are not precise enough to be implemented. For instance they would fail in cases such as བའི་, ལེན་, བའམ་, etc. Based on the information we gathered from our research, we propose to apply the following rules in this particular order¹⁷:

1. if the syllable contains a subscript, superscript or wasur then the main stack is what contains it (ex: དབྱེས་ → རྩ)
2. if a letter other than འ carries a vowel then it is the main stack (གཞི་ → འ)
3. if a vowel is carried by འ and འ is the first letter then འ is the main stack (འོད་ → འ)
4. if a vowel accent is carried by an འ which is not the first letter, the main stack is before the first འ with a vowel¹⁸ (མའིའོ་ → མ)

The following rules deal with syllables with no vowel, superscribed, subscribed nor wasur, only composed of letters that could be radical letters:

5. if the syllable has three or four letters and ends with འང or འམ, then the main stack is right before འང or འམ¹⁹ (བའམ་ → བ)
6. if the syllable is composed of only one letter this letter is the main stack
7. if the syllable contains two letters, then the first is the main stack (བར་ → བ)
8. if the syllable contains four letters, the main stack is the second letter (བཟབས་ → ཟ)

The following rules deal with syllables composed of a sequence of 3 simple letters:

9. if the final letter is not ས, then the main stack is the second letter (བཟའ་ → ཟ)
10. if the first letter cannot be a prefix to the second letter when it has no superscribed nor subscribed (see above for conditions on prefixes), then the main stack is the first letter (ཐགས་ → ཐ)
11. if ས cannot be second suffix after the second letter if it was a first suffix (meaning the second letter is not ཀ, ང, བ nor མ), then the main stack is the second letter (གནས་ → ན)

3.2 Ambiguous syllables

If the syllable doesn't match any of the above rules, then it is ambiguous. These syllables are those with no explicit vowel, no superscribed, no subscribed, no wasur, three letters, a final ས, the first letter can be a prefix to the second letter with no superscribed nor subscribed, and the second letter

¹⁷ This means that the first rule the syllable will match will be the one determining the main stack. This has the advantage to be transcribed easily into computer code.

¹⁸ This rule also works for Old Tibetan using འམ. The only exception to this rule is the very unlikely མའིའོ་(O), but the legality of this syllable is not clear.

¹⁹ This rule could be extended for Old Tibetan with འང.

is ག, ང, བ or མ. It is easy to find the 9 corresponding cases: མངས་, མགས་, དབས་, དངས་, དགས་, དམས་, བགས་, འབས་ and འགས་.

For each of these cases we have three possible structures:

1. ས is a second suffix (e.g. མངས་ can be decomposed as radical མ, suffix ང, second suffix ས, we note this decomposition $མ|ངས$)
2. ས is a suffix (e.g. མངས་ is radical ང, prefix མ, suffix ས, noted $མང|ས$)
3. ས is an affixed particle (e.g. མངས་ radical ང, prefix མ, suffix འ replaced by affixed particle ས, noted $མང|འ+ས$)

Case 2 and 3 have equivalent pronunciations but are important to distinguish because they will imply a different analysis in terms of lemmatization and part of speech tagging.

In order to decide where the main stack is, the only way is to take the decomposition with the highest probability according to our knowledge of the Tibetan language.

Let's review the different possibilities for the 9 cases, prefixing by * a form unattested in the dictionaries we consulted:

- མངས་: $མ|ངས$, $*མང|ས$, $མང|འ+ས$ → ambiguous, མ as main stack is more intuitive to our Tibetan informants
- འབས་: $*འ|བས$, $*འབ|ས$, $འབ|འ+ས$ → བ is the main stack
- མགས་: $*མ|གས$, $*མག|ས$, $*མག|འ+ས$ → ambiguous, མ as main stack is more intuitive to our Tibetan informants
- བགས་: $བ|གས$, $*བག|ས$, $*བག|འ+ས$ → བ
- འགས་: $*འ|གས$, $འག|ས$, $འག|འ+ས$ → ག
- དབས་: $*ད|བས$, $དབ|ས$ (OT), $དབ|འ+ས$ → བ
- དགས་: $*ད|གས$, $དག|ས$, $དག|འ+ས$ → ག
- དངས་: $ད|ངས$ (misspelling of དངས་²⁰), $*དང|ས$, $*དང|འ+ས$ → ད
- དམས་: $*ད|མས$, $དམ|ས$, $དམ|འ+ས$ → མ

3.3 Decomposition of the syllable

Once the main stack is found, the prefix, vowel accent and wasur are immediate to find. If a superscript or subscript is present, they can be immediately found with the rules exposed in “Simple constraints”.

Suffixes can be immediately classified between special suffix, first suffix, second suffix and affixed particle; the exceptions are ས and ར, plus the syllables ཀའི, ཟའོ, ཟླའོ. The latter are very rare and the context should make it obvious, but they are not decidable at syllable level.

²⁰ Attested in (Negi, 1993) (vol. 6) and (Duff, 2000).

| | བ | པོ | བཤམས | ལ | མ/མོ |
|---|---|----|------|---|------|
| ག | X | X | | | X |
| ང | X | X | | | X |
| ད | X | X | X | | X |
| ན | X | X | | | X |
| བ | X | | | | X |
| མ | X | X | | | |
| ལ | | | X | X | X |
| ར | X | X | X | | X |
| ལ | X | X | X | | X |
| ལ | X | X | X | | X |
| - | X | X | X | X | X |

At the time of the redaction of this article, some questions are still pending, like the possibility to affix one of these particles after the special suffix ལ, or the case of the following syllables found in the official lists: རྫོང་ལ་ རིལ་ལ་ རྫོང་ (we are not sure if these are errors or if other fusion with ལ or other syllables are possible).

4 Conclusion

We have described all the possible combinations of a Classical Tibetan syllable, listing all constraints and exceptions we found, resulting in a complete set of rules easy to implement in a computer language.

An immediate application has been to implement the rules in a spell checker running with the *hunspell* library,²³ freely available on <https://github.com/eroux/hunspell-bo>. Our formalization of the ལ and ལ endings allowed us to implement stricter rules in our spell checker and detect more potential mistakes.

We have also managed to build solid rules for the decomposition of a Classical Tibetan syllable, listing ambiguities and possible disambiguation.

REFERENCES

- Bacot, Jacques. 1946. *Grammaire du Tibétain littéraire*. Paris: Librairie d'Amérique et d'Orient.
- Beyer, Stephan. 1992. . Albany (NY): State University of New York Press.

²³ Hunspell (<https://hunspell.github.io/>) is the most popular spell checking library, used in all free software (LibreOffice, Firefox, etc.), but also in many closed source software such as the Adobe Suite, Mac OSX, etc. this makes our spell checker easily usable.

- Das, Sarat Chandra. 1915. An introduction to the grammar of the Tibetan language. Darjeeling: Darjeeling Branch Press.
- Dorje, Pema; and Hanfutun; and Drakpa. 2003. དག་ཡིག་གསར་བསྐྱོག། (The new spelling-dictionary). Xining: མཚོ་ཕོ་མོ་མི་རིགས་དཔེ་སྐྱུར་ཁང་།.
- Duff, Tony. 2000. The Illuminator Tibetan-English encyclopedia dictionary. Kathmandu: Padma Karpo Translation Committee.
- Dzongkha Development Commission. 2010. ཚོམ་སྐད་རྫོང་ཁ་ཚོག་མཛོད་ཆེན་མོ། (A comprehensive Chöke-Dzongkha dictionary), at <http://www.dzongkha.gov.bt/en/publications/title/a-comprehensive-ch-ke-dzongkha-dictionary>
- Dzongkha Development Commission. 2011. བསྐྱོག་ཡིག་གསལ་གྱི་ཨ་འོང་། (A Handbook of Dzongkha and Chöké Abbreviations), at <http://www.dzongkha.gov.bt/en/publications/title/a-handbook-of-dzongkha-and-ch-k-abbreviations>
- Dzongkha Development Commission. 2010. གཞི་རིམ་རྫོང་ཁའི་བརྗོད་གཞུང་། (Basic Level Dzongkha Grammar Textbook), at <http://www.dzongkha.gov.bt/dz/publications/title/basic-level-dzongkha-grammar-textbook>
- Hildt, Hélios. 2016. Towards describing Tibetan syntax: From word segmentation to rewrite rules through a semi-automatic workflow. *Himalayan Linguistics* 15: 78–112.
- Hill, Nathan. 2010. A lexicon of Tibetan verb stems as reported by the grammatical tradition. Munich: Bayerische Akademie der Wissenschaften.
- Kirtivajra, 1982. བོད་ཉེར་གྱི་བརྗོད་ཡིག་མིན་ཚོག་དོན་གསུམ་གསལ་བྱེད། (Tibetan-Mongol dictionary, the elucidation of the three configurations of words) [Reprint in: Four Tibetan-Mongolian Lexicon. New Delhi: Sharada Rani], at <https://www.tbrc.org/#!rid=W00KG09211>
- Negi, J.S. 1993. བོད་སྐད་དང་ལེགས་སྐྱུར་གྱི་ཚོག་མཛོད་ཆེན། (Tibetan-Sanskrit dictionary). Sarnath: Central Institute of Higher Tibetan, at <http://tbrc.org/link?RID=W1KG5422>
- Tournadre, Nicolas; and Dorje, Sangda. 2005. Manual Of Standard Tibetan: Language and civilization. Ithaca (NY): Snow Lion.
- Tsering, Tashi. 2006. Standardizing Tibetan terms of information technology. The China Tibetology Research Center, at http://digitaltibetan.org/images/e/e0/Standardizing_Tibetan_Terms_of_IT-China_Tibetology_Research_Center-2006.pdf
- Yisun, Zhang. 1985. བོད་རྒྱ་ཚོག་མཛོད་ཆེན་མོ། (Great Sino-Tibetan Dictionary). Beijing: མི་རིགས་དཔེ་སྐྱུར་ཁང་།, at <http://tbrc.org/link?RID=W29329>

- མཁམ་སྟོན་རྗེ་དབང་ལྷན་པུ་. 1979. དུས་གསུམ་རེའུ་མིག་ལྷ་མིའི་དགོངས་གཏེར། (The Treasure of Thomi's Insight, The Complete Verbal Forms). New Delhi: C.T. Khrto, at <http://tbrc.org/link?RID=W1KG14490>
- ཚེ་ཏན་འབས་བྱུང་།. 2003. བོད་གངས་ཚན་གྱི་སྐྱེ་རིག་པའི་བསྟན་བཅོས་ལེ་ཚན་འགའ་ཕྱོགས་བསྟུན། (A compilation of a few themes from the grammatical treatises of the Land of Snow). Xining: མཚོ་ཕྱོན་མི་རིགས་དཔེ་སྟེན་ཁང་།, at <http://tbrc.org/link?RID=W1KG25350>
- ཚས་ཀྱི་འབྲུང་གནས། (Ed.), 1721. བཀའ་འགྲུའ། (Tibetan Buddhist canon), ཟླ་དགེ། [Reproduced: Esukhia, 2012-2018], at <https://github.com/Esukhia/derge-kangyur>
- དུང་དཀར་སློབ་ཐབས་འཕེན་ལས།, 2002. དུང་དཀར་ཚིག་མཛོད་ཆེན་མོ། (Dungkar's great dictionary). Beijing: China Tibetology Publishing House, at <http://tbrc.org/link?RID=W26372>
- མོན་ལམ།. 2016. མོན་ལམ་ཚིག་མཛོད་ཆེན་མོ། (Monlam Grand Tibetan dictionary), at <http://www.monlamit.com/>
- ས་ཚན་ཀུན་དགའ་སྟོང་པོ།. 1092. ཡི་གེའི་བཞག་ཐབས་བྱེས་པ་བདེ་ལྷག་ཏུ་འཇུག་པ། (A reading manual, an accessible introduction for children) [in ས་སྐྱ་བཀའ་འབྲུམ། (The collected works of Sakya masters), ཟླ་དགེ. 1736. Reprint in: Dehradun. 1992], at [http://tbrc.org/link?RID=O01CT0026|O01CT00264CZ121847\\$W22271](http://tbrc.org/link?RID=O01CT0026|O01CT00264CZ121847$W22271)
- སི་ཏུ་པམ་ཆེན་ཚས་ཀྱི་འབྲུང་གནས།, XVIIIth Century. ཡུལ་གངས་ཚན་པའི་བརྡ་ཡང་དག་པར་སྦྱོར་བའི་བསྟན་བཅོས་ཀྱི་བྱེ་བྲག་སུམ་རྩ་བ་དང་ཉུགས་ཀྱི་འཇུག་པའི་གཞུང་གི་རྩམ་པར་བཤད་པ་མཁམ་པའི་མགུལ་རྒྱན་སྐྱེ་ཏེག་ཐེང་མངོས། (The beautiful rosary of pearls ornamenting the throat of scholars, a complete presentation of the thirty root verses and the guide to signs, two treatises showing the correct application of Tibetan grammar) [Reprint in: Xining: མཚོ་ཕྱོན་མི་རིགས་དཔེ་སྟེན་ཁང་།. 2001), at <http://tbrc.org/link?RID=W22787>

Elie Roux
elie.roux@telecom-bretagne.eu

Hélios Hildt
hhdрупchen@gmail.com

APPENDIX: The list of valid roots and exceptions

The following pages contain the complete lists of valid roots and exceptions, easily deducible from this article. They are in a simple form:

RootOrSyllable/PropertySuffix

Where PropertySuffix is:

- A if any vowel + suffix or affixed particle can appear after the root, with the exception of the suffix འ
- NB if any vowel + suffix or affixed particle can appear after the root, but at least one has to appear
- C if only affixed particles can appear after the root or syllable

Note that we take special suffixes into account separately.

As an example, \mathcal{A} means that all the following possibilities are valid:

ཀ, ཀལ, ཀལལ, ཀང, ཀངས, ཀད, ཀན, ཀབ, ཀབས, ཀམ, ཀམས, ཀལ, ཀའི, ཀའིའོ, ཀའོ, ཀའང, ཀའམ, ཀར, ཀས,
 ཀེ, ཀེལ, ཀེལལ, ཀེང, ཀེངས, ཀེད, ཀེན, ཀེབ, ཀེབས, ཀེམ, ཀེམས, ཀེལ, ཀེའི, ཀེའིའོ, ཀེའོ, ཀེའང, ཀེའམ, ཀར, ཀས,
 ཀུ, ཀུལ, ཀུལལ, ཀུང, ཀུངས, ཀུད, ཀུན, ཀུབ, ཀུབས, ཀུམ, ཀུམས, ཀུལ, ཀུའི, ཀུའིའོ, ཀུའོ, ཀུའང, ཀུའམ, ཀར, ཀས,
 ཀེ, ཀེལ, ཀེལལ, ཀེང, ཀེངས, ཀེད, ཀེན, ཀེབ, ཀེབས, ཀེམ, ཀེམས, ཀེལ, ཀེའི, ཀེའིའོ, ཀེའོ, ཀེའང, ཀེའམ, ཀར, ཀས,
 ཀོ, ཀོལ, ཀོལལ, ཀོང, ཀོངས, ཀོད, ཀོན, ཀོབ, ཀོབས, ཀོམ, ཀོམས, ཀོལ, ཀོའི, ཀོའིའོ, ཀོའོ, ཀོའང, ཀོའམ, ཀར, ཀས

\mathcal{NB} implies the following valid possibilities:

དཀའ, དཀལ, དཀལལ, དཀང, དཀངས, དཀད, དཀན, དཀབ, དཀབས, དཀམ, དཀམས, དཀལ, དཀའི, དཀའིའོ, དཀའོ, དཀའང, དཀའམ, དཀར, དཀས, དཀེ,
 དཀེལ, དཀེལལ, དཀེང, དཀེངས, དཀེད, དཀེན, དཀེབ, དཀེབས, དཀེམ, དཀེམས, དཀེལ, དཀེའི, དཀེའིའོ, དཀེའོ, དཀེའང, དཀེའམ, དཀར, དཀས, དཀུ,
 དཀུལ, དཀུལལ, དཀུང, དཀུངས, དཀུད, དཀུན, དཀུབ, དཀུབས, དཀུམ, དཀུམས, དཀུལ, དཀུའི, དཀུའིའོ, དཀུའོ, དཀུའང, དཀུའམ, དཀར, དཀས, དཀེ,
 དཀེལ, དཀེལལ, དཀེང, དཀེངས, དཀེད, དཀེན, དཀེབ, དཀེབས, དཀེམ, དཀེམས, དཀེལ, དཀེའི, དཀེའིའོ, དཀེའོ, དཀེའང, དཀེའམ, དཀར, དཀས, དཀོ,
 དཀོལ, དཀོལལ, དཀོང, དཀོངས, དཀོད, དཀོན, དཀོབ, དཀོབས, དཀོམ, དཀོམས, དཀོལ, དཀོའི, དཀོའིའོ, དཀོའོ, དཀོའང, དཀོའམ, དཀར, དཀས,
 དཀོམས, དཀོལ, དཀོའི, དཀོའིའོ, དཀོའོ, དཀོའང, དཀོའམ, དཀོར, དཀོས

And \mathcal{C} implies ཀེའི, ཀེའིའོ, ཀེའོ, ཀེའང, ཀེའམ, ཀེར, ཀེས.

These lists constitute the basis of the spell checker we built and are available on our git repository.

We have tested our spell checker against the data of (Hildt, 2016) and have only found expected discrepancies, due to the treatment of syllables built on ཅ, ས, སྐ, སྒ and སྔ as exceptions.

| | | | | |
|----------|-------|-------|--------|--------|
| Regular: | འཇུ/A | མང/NB | ལྟ/A | དལ/A |
| | འཇུ/A | ང/A | བལ/A | དལ/A |
| ཀ/A | ག/A | ལྟ/A | བལ/A | ལྟ/A |
| ཁ/A | ལྟ/A | ལྟ/A | བལ/A | ལྟ/A |
| ག/A | ལྟ/A | བང/A | མ/A | ལྟ/A |
| གྷ/A | ལྟ/A | བལ/A | མང/NB | ལྟ/A |
| དག/NB | དག/NB | ཅ/A | འཇུ/NB | ལ/A |
| དལ/A | དལ/A | གཅ/NB | ད/A | ལ/A |
| དལ/A | དལ/A | བཅ/NB | ད/A | ལ/A |
| བག/NB | བག/NB | ལྟ/A | གད/NB | འཇུ/NB |
| བལ/A | བལ/A | ཆ/A | བད/NB | འཇུ/A |
| བལ/A | བལ/A | མཆ/NB | མད/NB | འཇུ/A |
| བལ/A | མག/NB | འཆ/NB | འད/NB | བ/A |
| ཀྟ/A | མལ/A | ཇ/A | འད/A | ལ/A |
| ཀྟ/A | མལ/A | མང/NB | ང/A | ལ/A |
| ཀྟ/A | འག/NB | འང/NB | ལ/A | ལ/A |
| ཀྟ/A | འལ/A | ཇ/A | ལ/A | དབ/NB |
| ལྟ/A | འལ/A | ལྟ/A | བད/A | དལ/A |
| ལྟ/A | ག/A | བང/A | བལ/A | དལ/A |
| བཀྟ/A | ལྟ/A | ལ/A | བལ/A | འབ/NB |
| བལྟ/A | ལྟ/A | གལ/NB | ལ/A | འཇུ/A |
| བལྟ/A | ལྟ/A | མལ/NB | གལ/NB | འཇུ/A |
| ལྟ/A | ལྟ/A | ལྟ/A | མལ/NB | ལ/A |
| ལ/A | ལྟ/A | བལྟ/A | ལ/A | ལ/A |
| ལ/A | བལྟ/A | བལྟ/A | ལ/A | ལ/A |
| ལ/A | བལྟ/A | ད/A | བལྟ/A | ལ/A |
| མལ/NB | བལྟ/A | གལ/NB | བ/A | མ/A |
| མལ/A | བལྟ/A | བད/NB | ལ/A | ལ/A |
| མལ/A | ང/A | ལྟ/A | ལ/A | དལ/NB |
| འཇུ/NB | དང/NB | ལྟ/A | དབ/NB | དལ/A |

| | | | | |
|-------|-------------|------------------|-------|-------|
| མ/A | ར/A | འ ར ལ | ཧེ/C | ལུ/C |
| མུ/A | མ/A | | ཉེ/C | ལུ/C |
| མཱ/A | བམ/A | Specials suffix: | ཞེ/C | ལུ/C |
| མུ/A | པ/A | | ཏེ/C | ལུ/C |
| མ/A | པ/A | འཛེ/C | གཏེ/C | ལུ/C |
| གམ/NB | གམ/NB | མཛེ/C | ཏེ/C | ལུ/C |
| བམ/NB | བམ/NB | གལ/C | ཏེ/C | འཛེ/C |
| མ/A | མ/A | གེ/C | ཞེ/C | མུ/C |
| མ/A | མ/A | གེ/C | ལུ/C | མུ/C |
| བམ/A | མ/A | འཛེ/C | འཛེ/C | ལུ/C |
| བམ/A | གམ/NB | འཛེ/C | ལེ/C | མུ/C |
| མ/A | བམ/NB | ལུ/C | མུ/C | མུ/C |
| མམ/NB | བམ/A | ལུ/C | མུ/C | འཛེ/C |
| འམ/NB | བམ/A | ལུ/C | དེ/C | གཏེ/C |
| འ/A | ཏ/A | ལུ/C | དེ/C | མུ/C |
| མའ/NB | ཏ/A | གུ/C | མུ/C | མུ/C |
| འའ/NB | ཞ/A | འུ/C | དེ/C | མུ/C |
| འ/A | མ/A | འུ/C | ཞེ/C | མུ/C |
| བའ/A | | མུ/C | ཞེ/C | མུ/C |
| མ/A | Exceptions: | འགུ/C | མ/C | མུ/C |
| མ/A | | མུ/C | མུ/C | ཧེ/C |
| གམ/NB | དམེ | མུ/C | ཞེ/C | ཧེ/C |
| བམ/NB | མེ | མུ/C | དེ/C | གུ/C |
| མ/A | ཏ/C | མུ/C | ཞེ/C | མུ/C |
| མ/A | ཏེ | མུ/C | ཞེ/C | མུ/C |
| གམ/NB | མ/C | འུ/C | ཞེ/C | གུ/C |
| བམ/NB | མ/C | འུ/C | ཞེ/C | དེ/C |
| བམ/A | མེ | འུ/C | ཞེ/C | འུ/C |
| འ/A | མེ | གུ/C | ཞེ/C | འུ/C |
| ལ/A | མེ | གུ/C | འཛེ/C | འུ/C |
| གལ/NB | མེ | ལུ/C | འཛེ/C | གུ/C |

ཤེལ་/C ལྷ་/C
མིལ་/C ལྷག
མེལ་/C ལྷགས
མྱིལ་/C
མྱེལ་/C
བམེལ་/C
མལ་/C

Wasurs:

ཀ་/C
ཀའི་/C
ཀེ་/C
ཀུ་/C
ཀཱ་/C
ཀྱ་/C
ཀྲ་/C
ཀྲས
ཀྲགས
ཀྲེ་/C
ཀྲུ་/C
ཀྲེ་/C
ཀྲཱ་/C
ཀྲྀས
ཀྲཱ་/C
ཀྲཱེ་/C
ཀྲཱེ་/C
ཀྲཱེ་/C
ཀྲཱེ་/C
ཀྲཱེ་/C
ཀྲཱེ་/C
ཀྲཱེ་/C
ཀྲཱེ་/C