

Sorting Unicode Tibetan using a Multi-Weight Collation Algorithm

Robert R. Chilton

Technical Director,
The Asian Classics Input Project (ACIP)

acip@well.com

www.asianclassics.org

Why do we need a sorting algorithm for Unicode Tibetan?

- Tibetan script encoded in Unicode and ISO/IEC 10646
- Full support of Tibetan within a computer environment also requires:
 - Keyboard(s) or other input methods
 - Rendering: readable, printable display of the encoded Tibetan-script data (fonts, etc.)
 - Collation rules for generating culturally acceptable sorting

Previous efforts to sort Tibetan data

- Utilized a single-weight sorting model
- Generally adequate for sorting native Tibetan orthographies within a specific application/environment
- Not able to robustly handle foreign transcriptions and other non-standard orthographies

- Designed for Romanized or font-encoded Tibetan, not Unicode Tibetan
- Treat Tibetan-script sorting in an exclusive, special case fashion; such proprietary sorting methods are not likely to be widely implemented

Features of multi-weight sorting methodology

- Well-understood and widely implemented*
- Uses a **collation element table** to achieve culturally acceptable sorting
- Enables searching at different degrees of precision (e.g., case-sensitive searches)

*References:

- ISO/IEC 14651 (2001-04) Ed. 1.0
Information technology -- International string ordering and comparison -- Method for comparing character strings and description of the common template tailorable ordering.
- Unicode Technical Standard #10: Unicode Collation Algorithm (UCA).

Advantages for implementers and users of Unicode Tibetan

- Collation element table for Tibetan can “plug into” existing sort logic at the operating system level
- Robust searching and sorting of Tibetan data thus becomes automatically available to all compliant applications running within that operating system environment

- The same collation element table can be used across multiple platforms – resulting in consistent sorting of Tibetan data within different operating system environments

Moving from methodology to algorithm: an overview

- A look at dictionary sorting
- Understanding the multi-weight sorting (“international string ordering”) model and extending this model to Tibetan
- Determining the collation elements needed for sorting Unicode Tibetan

Dictionary sorting of Tibetan

- General agreement (since 1900) on dictionary order of native orthographies
 - Main exception: treatment of *wazur*
- Lack of consensus on details of sorting foreign-origin orthographies
- Universal agreement that all entries must appear under one of the 30 letters (ཀ་ to ཨ་)
 - Example of སྐྱལ་ (Sanskrit: “skandha”): sorts under the collation slot for སྐ

How foreign words are sorted in dictionaries generally

- Words from foreign languages are sorted according to the sort rules of the dictionary's language (and not the sort rules of the origin language)
 - a Danish word beginning with å appears after letter Z in a Danish dictionary
 - the same word is sorted under letter A in an English dictionary

- Extending this convention to Tibetan, all words in a Tibetan dictionary – including foreign words – are sorted under 30 letters
- Extending this convention still further, all vowel signs are treated in terms of the 5 standard Tibetan vowels
 - implicit vowel ཨ
 - 4 explicit vowel signs

The multi-weight sorting model for international string ordering

- Weights are generally assigned at three (or more) levels
- In Latin scripts these levels correspond to:
 1. alphabetic ordering = primary level
 2. diacritic ordering = secondary level
 3. case ordering = tertiary level

(Additional levels may be used for tie-breaking between strings not distinguished at the first three levels)

Examples in Latin script

- **role** and **rule** differ at the primary (alphabetic) level
- **role** and **rôle** differ at the secondary (diacritic) level
- **role** and **ROLE** differ at the tertiary (case) level
- **role** and **RÔLE** differ at both the secondary level and the tertiary level

Extending the multi-weight model to Tibetan

- ཏ་ and ཐ་ differ at the primary level
- ཏ་ and ཏ་_ཐ differ at the secondary level
- ཏ་ and ཏ་_ཐ differ at the tertiary level
- ཏ་ and ཏ་_ཐ differ at both the secondary level and the tertiary level

Determining collation elements for Unicode Tibetan: an overview

- Prescripts in Tibetan orthographies
- The Unicode model for encoding Tibetan script
- What is a collation element?
- 167 primary-weighted collation elements
- 9 secondary-weighted collation elements

Prescripts in Tibetan orthographies

- In English, all 26 letters always have primary weight; thus “at” sorts far away from “vat”
- In Tibetan, letters written before the radical letter have less than primary weight; thus ཀན, ལྷན, བཀན, and བལྷན sort relatively near to each other, under letter ཀ

■ 11 possible prescripts (or “pre-radicals”) might occur before the radical letter:

– 5 prefix letters: ༀ ༁ ༂ ༃ ༄

– 3 head letters: ༅ ༆ ༇

– 3 two-letter sequences of ༂ prefix followed by one of the head letters, i.e.: ༂྅ ༂྆ ༂྇

- Grammar rules define which radical letters can take which prescripts

– For letter ཀ་ there are 7 possible prescripts:

དཀ བཀ ཏྱ ལྱ སྱ བྱ སྱ

- No radical letter can take all 11 prescript forms (and some take none at all)

The Unicode encoding model for Tibetan script

- 193 distinct characters defined in Unicode
- The 30 letters (along with conjunct and reversed forms) are encoded twice: in nominal position and in orthographic-subjoined position
 - reflects the fact that the Tibetan script is written from top to bottom as well as from left to right
- Example of བརྟེན་ encoded as 6 characters:

བ + ར + ཏ + འ + ན + འ

0F56 0F62 0F9F 0F7A 0F53 0F0B

What is a collation element?

- A *collation element* enables clustering of multiple Unicode characters such that they can be treated together as a single item for determining sort weights
- Single characters also function as collation elements
- The weights assigned to the collation elements determine their sort (or collation) order relative to one another

Defining Tibetan pre-scribed radical sequences as collation elements

- For letter ཀ' we can define each of the 7 prescript + radical clusters (དཀ བཀ རྐ ལྐ སྐ བརྐ བསྐ) as a collation element (also called a “collation grapheme”)
- We can then assign sort weights to these collation elements such that they sort in a culturally acceptable relative order

Primary-weighted collation elements

- 30 nominal letters: ཀ་ to ཏ་ – which may be either radical letters (མིང་གཞི་) or suffix letters
- 103 multi-letter pre-scribed radical forms
 - In many of these 103 forms the prescript is written at the head line (encoded as 1 or 2 nominal characters) and the radical letter is encoded as a subjoined character
 - In the example of བརྟེན་, the radical letter ཏ་ is encoded in subjoined position

Defining the 4 explicit vowels as collation elements

- As collation elements, suffix letters cannot be distinguished from bare radicals*
- Because a nominal letter serving as a radical letter carries the implicit vowel ʔ , the 4 explicit vowels must be given primary weights; and must be weighted heavier than the nominal letters -- since a radical letter marked with an explicit vowel will sort **after** the same letter not marked by an explicit vowel

* in a stateless implementation

Defining the 30 post-radical letters as collation elements

- Post-radicals = the 30 letters in subjoined position (when not functioning as the radical letter in a pre-scribed radical form)
 - Requires maximum-length substring matching
- Only 4 post-radical (subscribed) letters occur in native Tibetan orthographies:

◁ ◡ ◢ ◣

- Remaining 26 are required to treat non-native orthographies in a consistent manner
- Must be given primary weights; and heavier than the 4 explicit vowels

Relative order of the 167 primary-weighted collation elements

- First: 30 nominal letters and 103 multi-letter pre-scribed radical forms (= 133 collation elements)
 - given sort weights such that the 103 pre-scribed radical forms are interleaved as appropriate with the 30 nominal letters
- Next: 4 explicit vowels
- Next: 30 post-radical letters (i.e., in orthographic subscribed position)
- Thus, total collation slots at the primary-weight level: $133 + 4 + 30 = 167$

Secondary-weighted collation elements

■ These 9 have no primary weight

– 4 combining marks:



– 5 signs:



The remaining 120 Unicode Tibetan characters

- $30 + 4 + 30 + 9 = 73$ of the 193 Unicode Tibetan characters have been treated above, leaving 120
- 59 of these 120 have a primary weight (in addition to a secondary and/or tertiary weight):
 - 19 can be decomposed into simple elements and thus need not be treated in the collation element table
 - 9 are variants (primary and tertiary weighted) of certain of the 30 nominal letters
 - 3 are variants (primary and tertiary weighted) of certain of the 4 explicit vowels

- 8 are variants (primary and tertiary weighted) of certain of the 30 subscribed letters
- 20 are the digits and half-digits
- The remaining 61 characters are punctuation marks and other symbols which generally have no impact on dictionary sort order and thus have no primary, secondary or tertiary weight

Appendices

- The Unicode (and ISO/IEC 10646) character-encoding chart for Tibetan
 - highlighting characters in example: བརྟེན
- An ordered list of collation elements for Unicode Tibetan

	0F0	0F1	0F2	0F3	0F4	0F5	0F6	0F7
0	 0F00	 0F10	 0F20	 0F30	 0F40	 0F50	 0F60	
1	 0F01	 0F11	 0F21	 0F31	 0F41	 0F51	 0F61	 0F71
2	 0F02	 0F12	 0F22	 0F32	 0F42	 0F52	 0F62	 0F72
3	 0F03	 0F13	 0F23	 0F33	 0F43	 0F53	 0F63	 0F73
4	 0F04	 0F14	 0F24	 0F34	 0F44	 0F54	 0F64	 0F74
5	 0F05	 0F15	 0F25	 0F35	 0F45	 0F55	 0F65	 0F75
6	 0F06	 0F16	 0F26	 0F36	 0F46	 0F56	 0F66	 0F76
7	 0F07	 0F17	 0F27	 0F37	 0F47	 0F57	 0F67	 0F77
8	 0F08	 0F18	 0F28	 0F38		 0F58	 0F68	 0F78
9	 0F09	 0F19	 0F29	 0F39	 0F49	 0F59	 0F69	 0F79
A	 0F0A	 0F1A	 0F2A	 0F3A	 0F4A	 0F5A	 0F6A	 0F7A
B	 0F0B	 0F1B	 0F2B	 0F3B	 0F4B	 0F5B		 0F7B
C	 0F0C	 0F1C	 0F2C	 0F3C	 0F4C	 0F5C		 0F7C
D	 0F0D	 0F1D	 0F2D	 0F3D	 0F4D	 0F5D		 0F7D
E	 0F0E	 0F1E	 0F2E	 0F3E	 0F4E	 0F5E		 0F7E
F	 0F0F	 0F1F	 0F2F	 0F3F	 0F4F	 0F5F		 0F7F

	0F8	0F9	0FA	0FB	0FC	0FD	0FE	0FF
0	 0F80	 0F90	 0FA0	 0FB0	 0FC0			
1	 0F81	 0F91	 0FA1	 0FB1	 0FC1			
2	 0F82	 0F92	 0FA2	 0FB2	 0FC2			
3	 0F83	 0F93	 0FA3	 0FB3	 0FC3			
4	 0F84	 0F94	 0FA4	 0FB4	 0FC4			
5	 0F85	 0F95	 0FA5	 0FB5	 0FC5			
6	 0F86	 0F96	 0FA6	 0FB6	 0FC6			
7	 0F87	 0F97	 0FA7	 0FB7	 0FC7			
8	 0F88		 0FA8	 0FB8	 0FC8			
9	 0F89	 0F99	 0FA9	 0FB9	 0FC9			
A	 0F8A	 0F9A	 0FAA	 0FBA	 0FCA			
B	 0F8B	 0F9B	 0FAB	 0FBB	 0FCB			
C		 0F9C	 0FAC	 0FBC	 0FCC			
D		 0F9D	 0FAD					
E		 0F9E	 0FAE	 0FBE				
F		 0F9F	 0FAF	 0FBF	 0FCF			

An Ordered List of Collation Elements for Unicode Tibetan

[*Note: a comprehensive Collation Element Table for Tibetan script will include additional collation elements, such as ཀླ, ཀླ, ལྷ, ལྷ, ལྷ, beyond those listed here.]

A. Primary Weighted Collation Elements

A.1. The 133 radical-initial sequences (also covers the suffix letters):

ཀ དཀ བཀ ཀླ ཀླ ཀླ བཀླ བཀླ ཁ མཁ འཁ ག དག བག མག འག ཁླ ཀླ
 ཀླ བཀླ བཀླ ང དང མང ང ལ ལ བང བལ ཅ གཅ བཅ ལྷ བལྷ ཆ མཆ
 འཆ ང མང འང ང ལ བང ཉ གཉ མཉ ལྷ ལྷ བྷ བལྷ ཏ གཏ བཏ ཏླ
 ལྷ ལྷ བྷ བལྷ བལྷ ཐ མཐ འཐ ད གད བད མད འད ང ལ ལ བང བལ
 བལ བ གབ མབ ལྷ ལྷ བྷ བལྷ བ དབ ལ ལ བ དབ འབ བ དབ འབ བ ལ
 ལ མ དམ ལ ལ ཅ གཅ བཅ ཅ ལྷ བཅ བལྷ ཆ མཆ འཆ ང མང འང
 ང བང ཐ ལ གལ བལ ལ གལ བལ འ ཡ གཡ ང བང(seen in བྷ) ལ བ
 གལ བལ ལ གལ བལ ཏ ལྷ ལ

Key:

Black for the 30 nominal letters. Note that whereas any of these 30 can serve as a bare radical, 10 of these can also appear in suffix position in native orthographies.

Blue for (relatively) unambiguous cases of prefixed and/or superscribed radical letters. Note that certain unavoidable ambiguities arise between native orthographies and transcriptions from foreign languages.

Red for ambiguous cases where a 3rd codepoint is required to distinguish the sequence as being a prefixed radical letter (as opposed to a root letter followed by a suffix). Note that certain cases (in Dzongkha) require a 4th codepoint in order to distinguish a case of a prefixed radical letter from a case of a suffix letter followed by a secondary syllable that involves a vowel (i.e., རྟི or རྟོ).

Magenta for an ambiguous case (in Dzongkha) where a 3rd (or possibly 4th) codepoint is required to distinguish the sequence as being a prefixed radical letter (as opposed to a suffix letter ད followed by a secondary syllable བ or བོ).

A.2. The 4 explicit vowels:

ི ཱ ེ ཻ

A.3. The 30 post-radicals:

ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ
 ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ ྱ

B. Secondary Weighted Collation Elements (have no primary weight)

B.1. The 4 combining marks:

◌̣ ◌̤ ◌̥ ◌̦

B.2. The 5 signs (used in transliteration):

ཧྲི ཨྲི ཨྲི ཨྲི ཨྲི

C. The 120 Remaining Unicode Tibetan Characters

The characters listed above (in items A and B) account for 73 of the 193 Tibetan characters defined in Unicode. This leaves 120 characters, of which 19 can be decomposed into simple elements and thus need not be treated in the collation element table. There is also no need to assign primary secondary or tertiary weights to the 61 characters that function as punctuation marks and other symbols since these generally have no impact on dictionary sort order. [Note that "Syllable OM" at U+0F00 is here treated as an ornamental symbol rather than as having any lexical value due to the fact that there is no canonical or compatibility decomposition specified for this character.]

The digits and half digits account for 20 further characters. The remaining 20 characters are variations (i.e., having both primary and tertiary weights) of certain of the 30 nominal letters, 4 vowels, and 30 subjoined post-initial letters listed previously.

9 Nominal letter variants:

འ ག ཁ ག ག ག ག ག (fixed form) ག

3 Vowel variants:

འི འི འི

8 Subjoined letter variants:

འི གི ཁི གི ཁི གི ཁི གི