

Khmer Sorting Analysis

Recent changes are in Red.

Sorting scheme for Khmer

Note that page references in this document are typically to Chuon Nath's Khmer-Khmer Dictionary, Japanese Reprint Edition with Arabic numbers at the bottom of the page.

Level 1 (Priority 1):

(Should Khmer numbers and signs precede the alphabet? Should U+17A3/U+17A4 precede the other letters of the alphabet?)

[U+1780..U+1793] The first 20 (of 33) Khmer consonants in the order they are encoded in Unicode: កខគឃងឝឥឦឧឬជ្ជដឝឧឝណតថទឝន

[U+1794] The next one (of 33) Khmer consonants in the order they are encoded in Unicode: ឝ It would probably be best to merge this and the next two entries under one heading, words with signs would list immediately after words with identical spelling without said signs. Is that acceptable?

[U+1794 U+17C9] A variant of the 21st Khmer consonant with 'p' pronunciation comes next (this is evident when marked as: ផ , however, there are hundreds of words in this section whose only distinction from a simple ឝ is their derivation)

[U+1794 U+17CA] A variant of the 21st Khmer consonant comes next (happily this is always marked as: ផ)

[U+1795..U+1799] An additional 5 (of 33) Khmer consonants in the order they are encoded in Unicode: ឝ្រឝ្លឝ្ឍឝ្តឝ្ម

[U+179A] An additional 1 (of 33) Khmer consonants in the order it is encoded in Unicode: ្រ It would probably be best to merge this and the next two entries under one heading (i.e., including ROBAT and the two independent vowels decomposed into ្រ and the appropriate dependent vowel). Is that acceptable?

[U+17CC] The ROBAT sign is (inconsistently in the Chuon Nath dictionary p. 465, 506, 538, 609, 750-1, 768, 1322, 1339-1340, 1633) treated for ordering purposes as an independent syllable. Should this be entered in phonetic order (as everything else is; I believe that would be appropriate)? What is its writing order when entered by a learned monk? It seems to fill

the roll of a superscript consonant and is not written stand-alone. If it is sorted as indicated here and not entered in phonetic order, there will have to be some mechanism to reorder it in the ordering algorithm.

This should probably be replaced by U+179A U+17D2

[U+17AB..U+17AC] These two independent vowels [ឃ្ម ឃ្ម្ម] are treated as consonants following 179A as they share a consonantal sound of 'r'

[U+179B] The next one (of 33) Khmer consonant: ឆ Should this and the following section be merged with decomposition of the following in U+179B plus the appropriate vowel?

[U+17AD..U+17AE] These two independent vowels [ឆ្ម ឆ្ម្ម] are treated as consonants following U+179B as they share a consonantal sound of 'l'

[U+179C] The next one (of 33) Khmer consonant: ឆ

[U+17AB..U+17AC] These two transliteration consonants [ឆ្ម ឆ្ម្ម] are treated as consonants following U+179C. They resemble the following Khmer consonant U+179F as they share a sound 's'. (Q: Are these two in the right order for sorting? Should they be integrated within the Khmer U+17DC for ordering purposes? None seem to be sorted in the Chuon Nate dictionary. Could we have examples of the characters they transliterate and the name of the script that character comes from? Have the glyphs and names been switched in Unicode?)

[U+179F..U+17A0] The next 2 (of 33) Khmer consonants: ឆ្ម ឆ្ម្ម

[U+17A1] The next 1 (of 33) Khmer consonants (this is separated because it is not available in a subscript form): ឆ្ម

[U+17A2, U+17A3..U+17AA, U+17AF..U+17B3] These characters are merged under one consonant (U+17A2) by means of decomposition into a glottal stop and a dependent vowel. For there to be a deterministic system this decomposition must be standardised. The resulting system (hopefully) will also sort transliterated Sanskrit/Pali text (note that Pali dictionaries sort the independent vowels first with separate sections for U+17A3 and U+17A4).

ឆ	ឆ	U+17A2
ឆ	ឆ	U+17A3 -> U+17A2 (?) ¹

¹There is a weak differentiation between short initial inherent vowel words (presumably U+17A3) and long inherent vowel words (presumably U+17A2) in the final section of the Chuon Nate Khmer dictionary. There is some controversy over the

អ	អ+ៗ	U+17A4->U+17A2 + U+17B6 (?)
ត	អ+ ្រ	U+17A5->U+17A2 + U+17B7
ឆ	អ+ ្រ	U+17A6->U+17A2 + U+17B8
ឧ	អ+ ្រ	U+17A7->U+17A2 + U+17BB ²
ឌ	អ+ ្រ (+ ក)	U+17A8->U+17A2 + U+17BB (+ U+1780) ³
ឍ	អ+ ្រ	U+17A9->U+17A2 + U+17BC
ណ	អ+ ្រ (+ វ)	U+17AA->U+17A2 + U+17BC (+ U+179C) ⁴
ដ	អ+ ្រ	U+17AF->U+17A2 + U+17C2 ⁵
ធួ	អ+ ្រ	U+17B0->U+17A2 + U+17C3
ឌី	អ+ ្រ	U+17B1->U+17A2 + U+17C4
ឍ	អ+ ្រ	U+17B2->U+17A2 + U+17C4 ⁶
ឌី	អ+ ្រ	U+17B3->U+17A2 + U+17C5

significance of U+17A3 and U+17A4 in Unicode. The linguist committee in Phnom Penh felt that there needed to be a distinction between the final Khmer consonant U+17A2 and the two independent Sanscrit vowels U+17A3..U+17A4. It would be good to clarify this issue if the particular Pali/Sancrit characters these are to represent could be shown.

²There are good examples of the equality of U+17A2 and the first part of the decomposed independent vowel on pages 1808-1850 (Arabic) of the Japanese reprint of Chuon Nath's dictionary.

³ The final Khmer consonant sound does not affect the ordering of this extremely rare and obsolete independent vowel. There will be some need of differentiating U+17A7 and U+17A8, but only at a higher level of sorting. This is referenced at the top of p. 1852 and p. 1877 of Chuon Nath's dictionary.

⁴ The final consonant U+179C does not figure in the sorting order, and is presented only for an understanding of the roots of the character. By this analysis there would seem to be an inconsistency on page 1851-1856, particularly with ឌី ... ឌីជំ ... ឌីឡុំ ... ឌី

... ឌី ... ឌីក If the Chuon Nath precedent were followed in this case it would seem to contradict the usage of decomposition for the other independent vowels that seem to separate into U+17A2 + x.

⁵ Note on p. 1860 the independent vowel in Chuon Nath's dictionary seems to have a secondary priority over the decomposition: ឌី ឌីអ

⁶ There are only two words which require the use of this character, the very common ឌី and the very rare U+17B2 U+1780 U+1789 U+17C0.

Level 1 (Priority 2)

◌̣	17C9	p. 195, 626 (in conjunction with 1794 higher level?), 1178
◌̣̣	17CA	p. 715 (in conjunction with 1794 higher level?), 1538-9, 1534-5

Level 2 (Priority 1): First subscript should include all the characters in Level 1 with the (possible) exception of a subscript form of ឡ which reportedly does not exist. However for sorting and display purposes it is assumed that any character in the range U+1780..U+17B3 could be a subscript. On the other hand only a subset of independent vowels are presently known to be subscripts (in addition to the consonant អ): ប្លង់ (ប្លង់ យ បង្កផន ផ្អផ)

Level 3 (Priority 1): Second subscript. Theoretically any of the characters under Level 2 may also sort in the same orders under Level 3. On the other hand in the Khmer language only about 9 are documented) ្រ ្រ ្រ ្រ ្រ ្រ ្រ ្រ ្រ

Level 4 (Priority 1): Vowels, 18 (Unicode: A committee of Khmer linguists voted to move three characters [U+17C6..U+17C8] from independent and combining forms of vowel to instead be signs as indicated in the Khmer Unicode section, reducing the number of dependent vowels that would need to be keyboarded.

The vowel/sign combinations which are known to exist using these are as follows:

◌̣̣̣	U+17B5	Short inherent p. 1583
◌̣̣̣̣	U+17B4	Long inherent
◌̣̣̣̣̣	U+17B6	
◌̣̣̣̣̣̣	U+17B6 U+17C7	p. 982, 1786, 1793
◌̣̣̣̣̣̣̣	U+17B7	
◌̣̣̣̣̣̣̣̣	U+17B7+U+17C7	p. 132, 1237, 1549
◌̣̣̣̣̣̣̣̣̣	U+17B8	
◌̣̣̣̣̣̣̣̣̣̣	U+17B7+17C7	p. 64, 251
◌̣̣̣̣̣̣̣̣̣̣̣	U+17B9	
◌̣̣̣̣̣̣̣̣̣̣̣̣	U+17B9+U+17C7	p. 760, 743-4, 1239, 1463

ㄹ	U+17BA	
ㄹᄇ	U+17BA+U+17C7	p. 246, 458, 597, 1887, 1808
ㄹ	U+17BB	
ㄹᄇ	U+17BB+U+17C7	p. 224, 542-3, 812, 1451, 1513, 1554
ㄹ	U+17BC	
ㄹᄇ	U+17BC+U+17C7	p. 1887
ㄹ	U+17BD	
ㄹᄇ	U+17BD+U+17C7	(Invalid? Not in Chuon Nath dictionary)
ㄹᄂ	U+17BE	
ㄹᄂᄇ	U+17BE+U+17C7	p. 743-4, 895, 1878-9
ㄹᄃ	U+17BF	
ㄹᄃᄇ	U+17BF+U+17C7	(Invalid? Not in Chuon Nath dictionary)
ㄹᄄ	U+17C0	
ㄹᄄᄇ	U+17C0+U+17C7	p. 748, 1242
ㄹᄅ	U+17C1	
ㄹᄅᄇ	U+17C1+U+17C7	p. 68, 215, 264, 689, 748 (but p. 1061)
ㄹᄆ	U+17C2	
ㄹᄆᄇ	U+17C2+U+17C7	p. 74, 142, 709, 761, 1475
ㄹᄇ	U+17C3	
ㄹᄇᄇ	U+17C3+U+17C7	(Valid? No example)
ㄹᄈ	U+17C4	
ㄹᄈᄇ	U+17C4+U+17C7	p. 76, 134-5, 142, 187
ㄹᄉ	U+17C5	
ㄹᄉᄇ	U+17C5+U+17C7	(Invalid? Not in Chuon Nath dictionary)
ㄹᄊ	U+17BB+U+17C6	
ㄹᄊᄇ	U+17BB+U+17C6 + U+17C7	(Invalid? Not in Chuon Nath dictionary)
ㄹᄋ	U+17C6	

កំ	U+17B6+U+17C6	
កំ៖	U+17B6+U+17C6+ U+17C7	(Invalid? Not in Chuon Nath dictionary)
ក៖	U+17C7	
កៈ	U+17C8	p. 413, 843, 1178, 1492, 1562, 1590, but lower priority to hyphen p. 1392-3!

Level 2 (Priority 2): Signs

កំ	U+17CE	p. 252, 542-3
!	(exclamation)	p. 1558
កំ	U+17CB	p. 119, 133, 148 (higher priority?), 177, 1178, 1544 (?)
-	(hyphen)	p. 1254, but why p. 1538-9
កំ	U+17D0	p. 119, 483, 681, 839, 1254
កំ	U+17CD	
កំ	U+17CF	
កំ	U+17D1	
—	(long hyphen)	p. 504, 1590, 1728, 1392-3
កំ	U+17D7	p. 252, 860

Level 3 (Priority 2) : Signs as above, relatively rare កំ៖ កំ៖ កំ៖

Test collation series

កំ	1	U+1780 ⁷	Single consonant
កំ	2	U+1780U+17CF	Single consonant and sign
កំកំ	3	U+1780U+1780	Consonant and next base consonant

⁷ When sorting ignore all spaces inserted into this column; they are purely for presentation/word-wrap purposes.

កក់	4	U+1780U+1780 U+17CB	Consonant and next base consonant and sign
កកវ	5	U+1780U+1780 U+179A	Could also be expressed with inherent vowels encoded U+1780U+17A5 U+1780U+17A5 U+179A (final consonant lacks vowel)
កកវ	5	U+1780U+17A5 U+1780U+17A5 U+179A	Identical to previous
កកត	6	U+1780U+1780 U+17B6U+178F	Vowel on second base resets cycling of third consonant
កកយ	7	U+1780U+1780 U+17B6U+1799	Third base consonant changes
កកេះ	8	U+1780U+1780 U+17C1U+17C7	Vowel on second base resets cycling, starting with no third base
កកែកកវ	9	U+1780U+1780 U+17C2U+1780 U+1780U+179A	ditto (presence of consonant in third base position follows absence of third base consonant)
កកែប	10	U+1780U+1780 U+17C2U+1794	Third base consonant cycle
កកោះ	11	U+1780U+1780 U+17C4U+17C7	Continuing to cycle through vowels on second base consonant
កក្រែក	12	U+1780U+1780 U+17D2U+179A U+17BEU+1780	Start cycling through subscript consonant on second base (reset cycling of vowel on second base)

កក្កា ក	13	U+1780U+1780 U+17D2U+17A2 U+17B6U+1780	Continue cycling through subscript consonant on second base (reset cycling of vowel on second base)
កក្កា ក	13	U+1780U+17B5 U+1780U+17D2 U+17A2U+17B6 U+1780	Identical to above (no implicit vowel when there is an explicit dependent vowel)
ខេតា ក	14	U+1781U+17C5 U+178FU+17B6 U+1780	Next consonant; cycling through vowel on first base
ខ្ញំ	15	U+1781U+17C6	Cycling through sign turned to vowel on first base
ខាំ	16	U+1781U+17B6 U+17C6	cycling through composed vowel on first base
ខាំង	17	U+1781U+17B6 U+17C6U+1784	Second base
ខះ	18	U+1781U+17C7	Cycling through sign turned to vowel on first base
ឃ្មោះ	19	U+178EU+17D2 U+1798U+17C4 U+17C7	
ឃ្មំ	20	U+178EU+17D2 U+1798U+17BB U+17C6	Composed vowel starts with subscript part first, then superscript.
ងោង	21	U+1784U+17C4 U+1784	
ងោង	22	U+1784U+17C4 U+17C9U+1784	Word with sign follows word without sign
នា	23	U+1786U+17B6	
នា ⁺	24	U+1786U+17B6 U+17CE	Sign follows vowel in entry order
នាៗ	25	U+1786U+17B6 U+17D7	Doubling sign indicates a consonant will follow (but weights as a sign)

ॐ	26	U+17A2U+17B4	Inherent vowel appear to have some affect
ॐ	27	U+17A2U+17B5	

The influence of inherent vowels in collation is a subject worth further investigation. For example, should words with a voiced inherent final vowel (Indic loanwords) be sorted before (or after!) words with final consonants lacking an inherent vowel? (Thanks to Kent Karlsson for raising this issue).

For corrections and suggestions please contact:
Maurice Bauhahn, 2 Meadow Way; Dorney Reach; MAIDENHEAD SL6 0DS;
U.K. Tel: +44(0)1628 626068; Email: bauhahnm@clara.net
5 December 2001 version 0.7gamma