

DATE: 2002-05-18

DOC TYPE: Expert contribution

TITLE: Comments on 02098-KhmerSrakOm_Aam (Proposal to add Om and Am)

SOURCE: Maurice Bauhahn

PROJECT: Khmer Unicode

STATUS: Comment on proposal

ACTION ID:

DUE DATE:

DISTRIBUTION: Worldwide

MEDIUM: PDF file

NO. OF PAGES:

A. Administrative	
1. Title	Comment on Proposal to add Om and Am
2. Requester's name	Maurice Bauhahn
3. Requester type	Expert
4. Submission date	18 May 2002
5. Requester's reference	ISO/IETC JTC 1/SC 2/WG2
6a. Completion	
6b. More information to be provided?	Only as required.

SUMMARY: The proposed **KHMER VOWEL SIGN OM** and **KHMER VOWEL SIGN AM** characters are in fact ligatures which reflect only a small part of vowel/sign combinations needed for languages expressed using the Khmer script. They should be rejected.

The Khmer script grew from Indic roots, developed over time, retained complexity, is susceptible to multiple interpretations, and is used for a variety of languages. Furthermore, Khmer reportedly has more vowels than any other language. Hence an official group of nine expert Khmer linguists brought together by the Cambodian Government deliberated over a period of days to come to studied conclusions on how to handle Unicode encoding. They issued a four page report dated 14 August 1996. These linguists understood the need (and many of the complications) of

expressing Sanskrit, Pali, and minority languages in the Khmer script. The following excerpt¹ shows the vowels they stated should be in Khmer Unicode:

ខ. ១ - សញ្ញាស្រ្ត:និស្ស័យមាន ១៦ ក្តី :

-ា ិ ឺ ឺ ឺ ុ ុ ុ េ េ្រ េ្រ េ- ែ- ែ- េ-ា េ-ា

The authors of the 02098 document have come to very different conclusions with regard to Khmer script vowels. They did not document their rationale, but evidences they could use include the list of vowels from the authoritative Chuon Nath dictionary as shown below² (note the two items proposed are circled):

យកក្ត អ ជាប្រធានដើមហើយបំបែកផ្សំជាមួយនិងស្រ្ត:និស្ស័យទាំង ២១ ក្តជា

អ អា អិ អឺ អី អឺ អុ អូ អួ អើ អឿ អៀ អេ ែ ៃ េ េា េា

អុំ អំ អាំ អះ

Furthermore these vowels are spelled with single syllable sounds. In addition, the derived Thai and Laotian scripts have registered such characters in their respective sections of Unicode. The NIKAHIT is not spelled separately in Khmer in these two instances (spelled OM and AM). And to top it all off: alphabetic ordering includes these in the sequence.

In the face of such obvious evidence, what rationale could there have been for a group of the leading Khmer linguists to not include these two items (or any of the four last *vowels* above) among the vowels in their recommendation for Khmer Unicode?

This document seeks to demonstrate that (1) NIKAHIT should be regarded as a distinct character, (2) that the combination of NIKAHIT with vowels would introduce ligatures into Khmer Unicode, (3) adding such ligatures would likely lead to inconsistent data entry/data store, and (4) adding such ligatures would lead to inadequate handling of non-Khmer language texts (especially Sanskrit).

I. NIKAHIT combines with vowels but is not essentially a vowel

1. National Higher Education Task Force. Decisions of the secretariat of the National Higher Education Task Force regarding encoding of the Khmer language into Unicode/computer. p. 2. Available at <http://www.bauhahnm.clara.net/Khmer/Page1.JPG Page2.JPG Page3.JPG Page4.JPG> and in translation TaskForce.html
2. Dictionnaire Cambodgien. Tome I. K.M. Cinquieme Edition. Phnom Penh: Éditions de l’Institut Boudhique, 1967. p. 15

2. The same Sanskrit grammar¹ extracted below identifies five signs of which the first is NIKAHIT.

វណិយុក្កិ ៥ យ៉ាង

១- អនុស្សាវៈ ៖ ។. កំ កំ kam

២- អនុនាសិកៈ ៗ ។. កំ កំ kañ

៣- វិសតិៈ : ។. កៈ បតិៈ កៈ បតិៈ kah

៤- វិរាមៈ ្រ ។. ក្ក កិ k, វាក្ក វិកិ vā

៥- អវគ្គិយៈ ៗ ។. តោ ៗ បិ លោ' ប

(លេខ ១១-១២-១៣-១៤-១៥)

3. This is not confined to Sanskrit, however. The authoritative Chuon Nath dictionary defines NIKAHIT (but does not give such an entry to dependent vowels).
4. The sign NIKAHIT typically has an 'm' or 'ng' sound. This sound is preserved in both of the proposed items. Although it appears to be lost when followed by KHMER LETTER NGO, that is largely because it is redundant with the sound of the following consonant. Hence NIKAHIT has strong consonant qualities as opposed to vowel ones.

1. Préah Sakyavong H. Tath, p. 3-4

- The Chuon Nath dictionary identifies ANUSVARA (another name for NIKAHIT) as derivative of Sanskrit and Pali in a separate entry¹ (there is no separate entry for dependent vowels).

អនុស្សាវ ឬ អនុស្សាវ (អៈនុស-ស្សាវៈ ឬ -សាវៈ) ន.

(សំ. ឬ បា.) និក្ខហិតឬជំលើ (°) : អនុស្សាវមានសំឡេង

ជាគរុ (ប្រើជា-រៈ ក៏បាន) ។

អនុស្សាវនក ឬ -កៈ (ម. ព. អនុស្សាវនា) ។

In a further entry for NIKAHIT², it is noted that the sign is virtually a consonant in Sanskrit.

និក្ខហិត (និក-គៈហិត) ន. (បា. និក្ខហិត ឬ និក្ខហិត)

ឈ្មោះគ្រឿងសម្គាល់មានរូបសណ្ឋានជាសូន្យមូលយ៉ាងនេះ

“o” រាប់ចូលជាព្យញ្ជនៈដែរ, មានសំឡេងដូចគ្នា “ង”

ប្រកបក៏មានដូចគ្នា “ម” ប្រកបក៏មាន ។ សំ. និង បា. រាប់

និក្ខហិតនេះចូលក្នុងពួកព្យញ្ជនៈសេសវគ្គផងដែរ អាចថា

អ័ង ឬ អ័ម ។ ឯទម្លាប់នៃជនបរទេសខ្លះ គេអាចថា អ័ង ឬ

អ័ម ឲ្យគ្រាន់តែសំឡេងដូចគេគ្រហឹមក្នុងច្រមុះប៉ុណ្ណោះ

ដោយគេយល់ថាអាចកុំឲ្យឮច្បាស់ថា អ័ង ឬ អ័ម ដើម្បីឲ្យ

និក្ខហិតនេះអាចរលាស់សំឡេងផ្សំនឹងស្រៈ ឥ, ខ ហើយនិង

ស្រៈឯទៀតខ្លះបានដោយងាយ សម្រាប់ភាសារបស់

គេ ... ។

- NIKAHIT is already recognised as a sign (not a vowel) in Khmer Unicode, set apart with the similar characters KHMER SIGN REAHMUK and KHMER SIGN YUUKALEAPINTU.
- Note that the last four ‘vowels’ in the Chuon Nath dictionary are of a different nature from those that precede: (a) They all involve combinations with signs (two with the unwritten inherent vowel), (b) They follow the last vowel which matches Sanskrit/Pali vowels, (c) they do not match Sanskrit or Pali vowels, and (d) they all have following consonantal sounds (ng, m, or h).

1. Dictionnaire Cambodgien, p. 1658.

2. Dictionnaire Cambodgien, p. 532.

II. The combination of NIKAHIT with vowels creates ligatures.

1. Not only is the NIKAHIT sound preserved in the two items proposed, the vowel sound is also preserved. Hence, two distinct characters retain their identity when used in combination. This is not the case in other Khmer vowels formed without signs that nevertheless share glyphs.
2. The Khmer and Laotian analogues which were brought into Unicode from legacy encodings have presented problems in their [insert thailaoam.tif] decompositions¹.

The categories discussed above have covered essentially all of the characters with compatibility decompositions defined in TUS 3.1, leaving only a dozen or less. I cannot give an exact count since it is so unclear (at least to me) how the remaining characters are best characterised that another person might have already included them in one of the above sets and perhaps have a different residue. I will just mention two that I have found hard to categorise: U+0E33 THAI

CHARACTER SARA AM “ ̂ ”, and U+0EB3 LAO VOWEL SIGN

AM “ ̂ ” Because these both decompose into a combining mark followed by a base characters, with the former combining with some preceding characters, I have not know which of the preceding categories, if any, to put them into.

3. By way of analogy there are many other combinations of dependent vowels with the signs KHMER SIGN REAHMUK and KHMER SIGN YUUKALEAPINTU which sort with the same weight as separate vowels in the Chuon Nath dictionary (see the document at <http://www.bauhahnm.clara.net/Khmer/KhmerSortingUnicodegamma.pdf>) but are not ‘officially’ recognised as vowels.

III. Adding these two items would introduce ambiguity into Khmer Unicode.

1. Were the two proposed items to be accepted, they would introduce alternative ways to enter the same material in addition to that originally intended. This, unfortunately, would lead to identical words entered by alternative means being treated as separate and distinct in binary searches and spell checks. A compatibility algorithm could, of course, resolve such difficulties...but given the speed required of these basic functions over a growing corpus of data, it is unlikely that a speed degrading compatibility algorithm will be introduced in the foreseeable future.
2. Whenever there is ambiguity, data entry would be slowed while the typist seeks to mentally resolve which alternative to choose.

IV. Adding such ligatures would lead to inadequate handling of non-Khmer language texts (especially Sanskrit).

1. Peter Constable. “A review of characters with compatibility decompositions”. Implementing Writing Systems: An Introduction. Melinda Lyons and NRSI Team, Editors. Preliminary edition. Dallas: SIL International, 2001. p. 212.

1. One of the rules of Khmer/Sanskrit/Pali is that at most one dependent vowel is allowed in a single consonantal cluster. By defining NIKAHIT as a vowel, we would encounter the awkward situation where combinations of NIKAHIT with ‘other’ vowels (and in the case of Sanskrit that could be any vowel) in a single cluster violates that grammatical rule.
2. In order to facilitate proper intra-cluster ordering of Khmer script characters, fonts and keyboard drivers should enforce a sequence: one consonant or independent vowel, zero or one register shifter, zero one or two COENG + consonant/independent vowel pairs, zero or one dependent vowel, zero or one sign, and zero or one vowel-like sign (this has not been standardised yet...but something like this would be the end result). It would damage categorisation to have NIKAHIT function in two different categories (i.e., dependent vowel and vowel-like sign).

In conclusion, NIKAHIT should remain in Khmer Unicode as a separate sign. Ligatures such as the proposed KHMER VOWEL SIGN OM/KHMER VOWEL SIGN AM should not be added. Nevertheless, the combinations of KHMER SIGN NIKAHIT, KHMER SIGN REAHMUK or KHMER SIGN YUUKALEAPINTU with inherent or dependent vowels may be treated as unique vowels in future Khmer sorting algorithms. To explicitly encode these two ligatures, however, would unnecessarily complicate data entry, code tables, and display algorithms. Encoding the proposed two characters would be inconsistent with the handling of the KHMER SIGN REAHMUK or KHMER SIGN YUUKALEAPINTU which function similarly in combination with dependent vowels.