

Towards a Computerization of the Lao Tham System of Writing

Grégory KOURILSKY
gregory@kourilsky.com

(1)
(1) INALCO
2, rue de Lille
75343 Paris Cedex 7, France
http://www.inalco.fr

Vincent BERMENT
Vincent.Berment@imag.fr

(1) (2)
(2) GETA-CLIPS (IMAG)
BP 53
38041 Grenoble Cedex 9, France
http://www-clips.imag.fr/geta/

ABSTRACT — Although the subject of under-resourced *languages* is a problem often taken into consideration, one omits to make the distinction with under-resourced *scripts*. The Tham script of Laos, used to write Buddhist Texts, is one of the two official writing systems used in Lao P.D.R. (Laos). But if the other one — the so-called “laic” Lao script — is now quite well computerized (numerous fonts, input software, word processors, Unicode area), the Tham script seems to have always been forsaken by modern technologies (typewriters and computers). And this phenomenon does not seem to be being reversed soon since the Unicode Standard does not integrate it in any zone. Understanding the sociological and technical reasons of this neglect, we present an approach to mend it.

RÉSUMÉ — Bien que le sujet des *langues* peu dotées informatiquement soit un problème régulièrement abordé, on omet souvent de distinguer le cas des *écritures* peu dotées. L’écriture tham du Laos, employée pour noter les textes bouddhiques, est l’une des deux écritures officielles de la République Démocratique Populaire Lao (Laos). Si l’écriture lao proprement dite (appelée parfois par opposition *écriture lao “laïque”*) est aujourd’hui relativement bien dotée informatiquement (nombreuses polices de caractères, logiciels de traitement de texte, zone Unicode), l’écriture tham semble avoir toujours été délaissée des techniques modernes de saisie (mécanique aussi bien qu’informatique). Ce phénomène ne semble pas devoir s’inverser puisque le consortium Unicode ne l’a pas intégrée dans son standard. Tout en évoquant les raisons sociologiques et linguistiques de ce délaissement, nous proposons des éléments pour y remédier.

KEYWORDS — Tham script, ThamWord, Under-resourced scripts, Unicode.

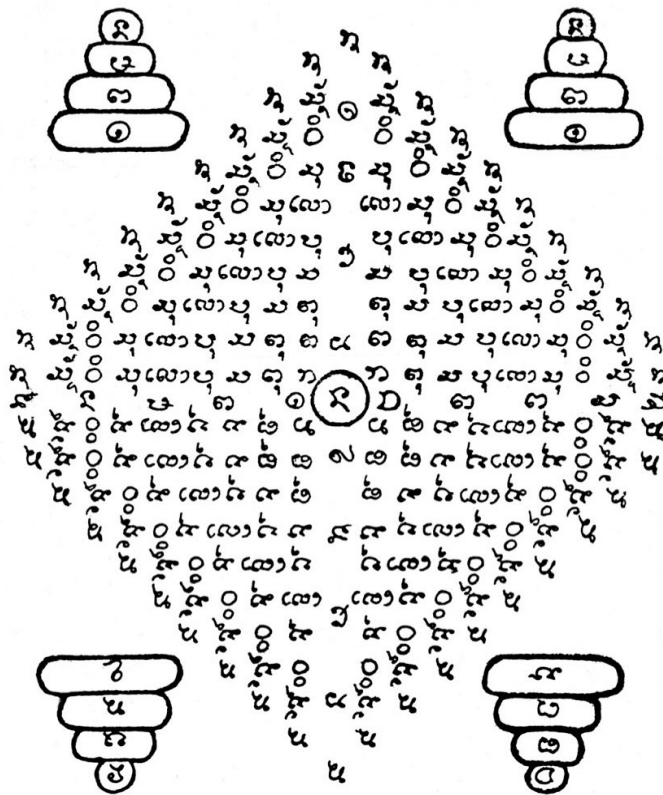
1. An Under-Ressourced Script

The Tham Script of Laos is the one of the two official writing systems used in Lao P.D.R. (Laos). The first one is the so-called “laic” Lao script used by population for administration and everyday life. The second one, Tham Script, is used to write Buddhist Texts, even if there are some examples of non-buddhist lapidary inscriptions like the famous stele of Dan Sai.

Contrary to the “laic” Lao Script which is now quite well computerized (numerous fonts, text input software, Unicode area), the Tham script seems to have been forsaken by modern technologies (typewriters and computers). We can find linguistical and sociological reasons to explain this phenomenon.

The linguistical reason is that the Tham Script is a complex script, never really codified, that allows many variant forms. Furthermore, it has the particularity to be used to write two different languages, Pali (language of Theravada Buddhism) and Lao, with different rules and specific characters for each of these languages (this aspect will be described further on).

The sociological reasons have several aspects. First of all, the Tham script is not well known among the Lao population: Principally high level Monks and few Scholars can read and write Tham. Secondly, the religious function of the Tham script makes it more than a simple language transcription system: as in Thailand, Cambodia and Myanmar, Pali letters in Laos are for an older time employed with a graphic consideration as disposition of the words, combinations, choice of graphic forms, that gives “an extra literary dimension to the reader consciousness” (BIZOT 1992:18). This aspect is well illustrated in the *yantra*, diagrams for protection made by grouping letters disposed in a non-linear way¹. This utilization of these dispositions of letters is perceptible not only in texts but on inscriptions on diverse supports as sheet of metal or wood, woven and even tatoos on bodies²:



Lao *gāthā*³ organized in a *yantra* with Tham letters⁴

It means that possibilities offered by text typing systems (mechanical typewriters and software) were a long time insufficient for the use of the Tham Script in Laos. Another explanation is, last but not least, Tham has begun to be obsolete since the influence of Thai Buddhism reform⁵ came deeper into Laos in the first decades of the twentieth century. Lao scholars like Maha Sila Viravongs who studied in pagodas in Thailand (especially a few years in Bangkok), were attached to Pali canonical texts from Ceylan and refractory to non orthodox practical as the use of letters evoked above. That ended in the Lao alphabet reform by the Buddhist Institute in adding missing characters to write Pali in the Lao “laic” alphabet⁶, so the Tham was not useful anymore

¹ See BIZOT 1996, 2001

² See BECCHETTI 1991.

³ *gāthā* means “Stanza’s heart”, and can be considered as magical formula (GABAUDE, 1988:246, note 6).

⁴ This *yantra* illustrates the cover of (PHONE PHRA NAO MONASTERY 1965).

⁵ King Mongkut (Rama IV), initiated in the first half of the 19th century a movement of reforms in Thai Buddhism that ended to the creation of a new *nikaya* called *dhammayut* (t. ธรรมยุตติกนิกาย) (GABAUDE 1988:31, note 73).

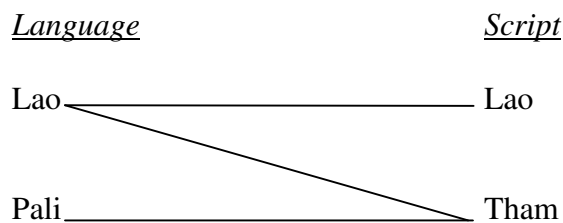
⁶ In the 1930's, scholars from the Buddhist Institute introduced fourteen letters in the Lao alphabet to complete the Pali consonantic system (see VIRAVONGS 1937): *gh*, *ch*, *jh*, *ñ*, *t*, *th*, *d*, *dh*, *n*, *dh*, *bh*, *l*, and

considering Pali transcription. In fact the magical function of Tham continued to be used in the open country but quite badly considered by official Buddhism in the cities, place of technology and computer development. But today it seems that we can observe a resurgence of Tham Script even in the Capital, and old magical uses begin to get out their hiding-place⁷.

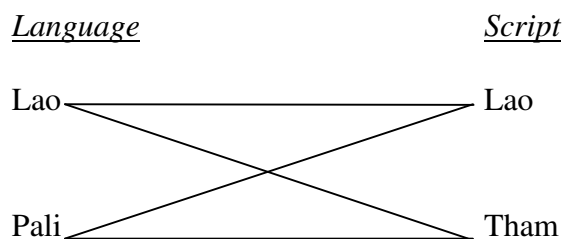
This phenomenon makes us think that the computerization of the Lao Tham system of writing is today a necessity more especially as modern technologies and computer softwares allow to render non-linear text and diagrams mentioned above. Far from analyzing the whole difficulties of computerization of the Tham Script system of writing, some important aspects will be reviewed.

2. Two “Sub-Writing Systems”

As we wrote, the Tham script has the particularity to transcribe two languages, Pali and Lao. Traditionally, relation between language and script in Laos can be illustrated in this schematic way:



Let’s notice that during the existence of the Buddhist Institute (1933-1975), the Pali complementary letters added⁸ in Lao “laic” alphabet made the relation completely symmetrical:



We can see that Tham is able (and employed) to write two languages, Pali and Lao, which are very different with regard to origin and phonology. Moreover, Phay Luang Maha Sena made a distinction in publishing two volumes of his manual to learn Tham Script: one volume concerns Tham to transcribe Lao, the second one teaches Tham transcribing Pali⁹.

The first form of Tham Script was probably created towards the 15th century by the Mon in Lanna to transcribe all Pali letters included those that does not exist in Thai languages (FERLUS 1995), like sounded or/and aspirate occlusives (*g, gh, j, jh, bh, etc.*), retroflexes (*ṭ, ṭh, ḍ, ḍh, ṇ*) and the *niggahita* sign *m̐* (l. *ນິກະຫິດ* /nikk^hahit/). Thus the Tham Script is related to Mon script and

niggahita ṁ *m̐*. The two sanskrit letters *ṣ* and *ṣ*, unknown in Pali, were added too. All these letters falled into disuse after 1975 with a new reform lead by Phoumi Vongvichit.

⁷ We can observe the same resurgence in Northern Thailand where Yuon (Lanna) script recently profited of new interest (but not necessarily for the same reasons).

⁸ Cf. note 6.

⁹ See SENA 1957, 1963.

ບ	ຜີ	ຜີ	ຜີ	ຜີ	ຜີ
/b/	/fɔ̃/ ¹³	/f/	/j/	/h/	/ɲ/
ບູ	ບູ	ບູ	ບູ	ບູ	ບູ
/ɲɔ̃/	/nɔ̃/	/mɔ̃/	/ɲɔ̃/	/lɔ̃/	/vɔ̃/

We could make the same distinction for the vowels and other particular signs like tone marks only used in Lao. In Pali, only 8 vowels are used (a, ā, i, ī, u, ū, e, o) but 29 in Lao (sometimes written in several manners), plus particular signs as ັ /a:/, ັ /ɛ:/, ັ /ru:/, etc¹⁴. But Pali uses independent vowels signs (ຳ, ຳ, ຳ, ຳ, ຳ, ຳ and ຳ) unknown in transcription of Lao.

Furthermore, rules and proprieties of writing differ when transcribing Pali or Lao. Thus it appears that a same word can be written in a different manner if written in a Pali text or a Lao text, for example a Lao word taken from Pali:

- (1) ຕາູ່, *tanhā*, « desire » (Tham transcribing Pali),
- (2) ຕາູ່, [ຕາູ່], /tanhā:/, « desire » (Tham transcribing Lao).

In a Pali transcription, a consonant following another consonant without inherent vowel is a subscript ((1) : ູ ັ), but in a Lao transcription it is the final or medial consonant that can be (but not systematically) subscript ((2) : ູ ັ).

Lastly, graphic forms of several subscript consonants differ if used in a Pali or in a Lao transcription:

value	Base Form	Pali subscript form	Lao subscript form
<i>k</i>	ກ	ກ	ກ / ັ
<i>ñ</i>	ຳ	ຳ	ຳ / ັ
<i>n</i>	ນ	ນ	ນ / ັ
<i>l</i>	ລ	ລ	ລ

For all these reasons, we've chosen to distinguish between two “sub-writing systems”, each working in its own manner:

- P-Tham (Tham transcribing the Pali language),
- L-Tham (Tham transcribing the Lao language).

We precise that this terminology is made here for the purpose of computerization of the Tham Script and is not known in Laos.

¹³ Lao high-class consonants are sounded, without tonemark, with a rising tone (INTHAMONE 1987:77) indicated with API sign ັ placed on the right of the consonant. Sign ັ represents any API vowel sign.

¹⁴For Tham Writing System (to write Pali and to write Lao), see (SENA 1957, 1963) (in Lao) and (KOURILSKY 2005) (in French).

3. Unicode and Input Methods

A simple way to mend the under-resourced problem of the Tham Script would be to create a Tham font which could be used on Windows system, the most popular operating system in the world. Considering the multi-utilization of several signs (for example the sign ၇ is used to note Pali independant vowel /a/, Lao vowel /ɔ:/ and Lao consonant /ʔ/; the sign ၆ is employed to write vowel /e:/ and many non-connex vowel as ၇ꠜ /ɔʔ/, ၇ꠞ /ao/, ၇ꠟ /jaʔ/, etc.), about 132 elementary signs are needed to write all Tham letters. Thus we created a font called *ThamStandard*¹⁵ that includes those 132 characters (plus a dozen of various forms). The problem is that much less than 132 characters are directly accessible on keyboards, so complex combinations are often necessary to type some letters (as *r* subscript ၇ has to be typed with [Ctrl] + [4], [E] ; ligature *bb* ၇ is typed with [CTRL] + [1], [o], etc.). Many Indic and derivated scripts fonts meet the same problem (Devanagari fonts, Khmer fonts¹⁶, etc.) and sometimes several fonts are needed to type all the letters of those alphabets. One solution to face this problem is a normalization by an integration into the Unicode Standard¹⁷. The aim of the Unicode Standard is first to normalize every script in the world to permit a perfect compatibility between different fonts and platforms in exchanges of computerized texts (like Internet). In Unicode, each character has a single code and any “Unicode font” has to respect it. For example A LATIN CAPITAL LETTER A has the code U+0041, \$ DOLLAR SIGN has the code U+0024, क DEVANAGARI LETTER KA has the code U+0915, ນ LAO LETTER KO has the code U+0E81, etc. But the useful point for our matter is a particularity of Unicode Standard that only *basic characters* are coded but not all *signs* (so-called *glyphs*). That means if क DEVANAGARI LETTER KA is coded U+0915, the various forms that can take this letter according to contexts (क, क्, कः, क्क, etc.) are not coded: it is the Windows Unicode Script Processor¹⁸ (called *Uniscribe*¹⁹) jointly with an OpenType font that will take over to display the appropriate glyph.

The Tham script is not in Unicode yet but the next few examples will illustrate that the distinction between basic character and glyph is pertinent for this script:

<i>Glyphs</i>	<i>Basic Characters</i>
၇ ၇ꠜ ၇ꠞ	Tham letter <i>k</i>
၇ ၇ꠟ ၇ꠠ	Tham letter <i>ñ</i>
၇ ၇ ၇	Tham Lao vowel /ɔ:/
၇ ၇	Tham vowel <i>ā</i>

¹⁵ In free access on <http://www.laosoftware.com>.

¹⁶ The Khmer font *Kdol*, created by ÉFEO (École Française d’Extrême-Orient) that includes numerous ligatures and various forms of letters regarding 17-19th centuries manuscripts, allows 438 signs (BIZOT 1992:19).

¹⁷ For more information about the Unicode Standard, <http://www.Unicode.org>.

¹⁸ Most of works concerning typing Unicode Asian fonts have been done for the Windows system. For more information about Unicode on Macintosh, see <http://apple.com/macosx/features/international> and http://www.alanworld.net/Unicode/fonts_macosx.html. For Linux, http://www.alanworld.net/Unicode/fonts_unix.html.

¹⁹ See <http://www.microsoft.com/typography/developers/uniscribe/default.htm>.

By the way, the Unicode Standard and the associated input methods are able to solve the problem of the computerization of Indic and derivated scripts with numerous letters, because only a limited number of characters is required on the keyboard whereas the Windows Unicode Script Processor displays the adequate glyph according to the context (letter in a final position, ligature, etc.). The matter to distinguish characters from glyphs is in practice not an easy task. The scripts implemented in Unicode solve their problems in different ways, not always conformable with general Unicode recommendations. A few cases taken among Indic and derivated Asian scripts (Devanagari, Tamil, Thai, Lao, Khmer, Burmese, Tibetan, etc.) will render illustrious the various choices available for an input method of the Tham script.

3.1 Consonants

Many Southeast Asian scripts use *subjoined*²⁰ forms for consonants. According to Unicode recommendations, subjoined forms have not to be coded and are supposed to be rendered by the system.

In Devanagari script, many various forms and ligatures occur, especially when a consonant lost its inherent vowel. The input method to render Unicode Devanagari introduces the special character 094D ◌ DEVANAGARI SIGN VIRAMA that points out that the preceding consonant has no inherent vowel. This character, according to the context, indicates to the system which form should be displayed for this preceding consonant:

0915 क + 094D ◌ + 092F य → क्य (क displays as क्; the sequence is read /kja/),

0915 क + 094D ◌ + 0024 क → क्क (क्क is a ligature; the sequence is read /kka/),

In some contexts the *virama* is explicite and the sequence occurs as क्, etc.

In Khmer script, the presence of a subscript consonant in a syllable also indicates that the preceding consonant lost its inherent vowel. The difference with the Devanagari script is that the second consonant of the cluster will change its shape (and not the first one like in Devanagari) in becoming a subscript. In the word ក្សា /ksa:/, the subjoined position ◌ of the consonant ស s indicates that we have to read /ksa:/ and not /kasa:/ (the consonant ក lost its inherent vowel /a:/). To inform rendering system that the consonant has to be rendered subjoined, an Unicode Khmer font has a character named *coeng*²¹ (17D2 ◌ KHMER SIGN COENG) that has to be typed before the consonant that has to be rendered subjoined:

1780 ក + $\underbrace{17D2 ◌ + 179F ស}_{◌}$ + 17B6 ា → ក្សា /ksa:/

The Unicode Standard specify that character *coeng* does not exist in the Khmer alphabet and “does not have a conventional visual form²² in Khmer as it is a control character to cause the formation of a subscript” (MICROSOFT CORPORATION 2004a:2).

²⁰ Terminology employed in Unicode’s publication to call non-basic form of a consonant. For exemple Khmer subscripts consonants are the *subjoined* forms of the basic consonantic characters.

²¹ That means “foot” in Khmer, and in a wider sense designs subscript consonants in the Khmer alphabet.

²² The sign ◌ which is the representation of the character U+17D2 SIGNE KHMER COENG in the Unicode’s publications is an arbitrary form employed only for the needs of writing.

Two particular cases are :

- The subjoined 𑄣 *r* placed before consonant but coded (and typed) after it in “logical” order²³:

$$1785 \text{ 𑄣 } + \underbrace{17D2 \text{ 𑄣 } + 179A \text{ 𑄣 }}_{\text{𑄣}} \rightarrow \text{ 𑄣 } /cra:/$$

- The *Robat* sign 𑄣 *r*, today not pronounced, is at the origin phonetically placed before the consonant above which it is placed but must be typed after it (KHMEROS 2004:13) in graphical order:

$$1796 \text{ 𑄣 } + 178E \text{ 𑄣 } + 17CC \text{ 𑄣 } \rightarrow \text{ 𑄣 } /poa/, \text{ etym. } poarna.$$

The case of the *Robat* illustrates that even a script so respectful of the Unicode recommendations as Khmer script can depart from the logical order.

The Tibetan Unicode works in a particular manner in coding twice the consonants: one basic form and one subjoined form (so in total contradiction with the Unicode recommendations). Example:

0F42 ། TIBETAN LETTER GA

0F92 ། TIBETAN SUBJOINED LETTER GA

The reason put forward are there is no *virama* in Tibetan²⁴ and the necessity to reduce needs for memory while the use of subscripts is more than frequent. Without discussing this choice, this example shows that the Unicode’s rules are not univocal.

3.2 Vowels

In most of the Asian scripts deriving from the Brahmi one, numerous vowels are formed by several graphemic signs (*non-connex* vowels) as in the few following examples:

𑌀	Tamil vowel /o/,
𑌁	Malayalam vowel /o/,
𑌂	Lao vowel /ja?/,
ໂ	Thai vowel /o?/,
𑄣	Khmer vowel /ua:/,
།	Tibetan vowel /i:/,
𑄣	Tham vowel /ua:/.

We can observe that two families of Unicode input methods could be distinguished: one “Indic” model and one “Thai-Lao” model.

The “Indic” model strictly follows the Unicode Standard’s recommendations and codes vowels as phonemes (or known as phonemes in the descriptive rules of those languages). Examples:

²³ The *logical order* agrees to phonetic order and in the most of case to typing order (UNICODE 2003:17).

²⁴ The letter *rok-me*, Tibetan equivalent to Indic *virama*, is employed only for transcribing Sanskrit.

- 0BCA ௮ TAMIL VOWEL SIGN O
- 0D4C ൠ MALAYALAM VOWEL SIGN AU
- 0BCA ្រ KHMER VOWEL SIGN YA
- 0F73 ི TIBETAN VOWEL SIGN II

We precise that from a general point of view the encoding method corresponds to the way of local spelling. That means that a single key has to be typed (after have typed the consonant) to display those vowels:

1780 ក + 0BCA ្រ → ក្រ /kɯə:/

This “Indic” vowel treatment implies to repositionning glyphs to avoid cross-over:

178F ត + 17B6 ្រ + 1780 ក + 0BCA ្រ → តក្រ /ta:kɯə:/

In this example the sylleme ក្រ has to be moved on the right to not cross-over with ្រ. The reason put forward for this “logical” order is to facilitate alphasorting and speech rendering (SOLÁ 2004:3) in always having a consonant coded first²⁵. There are exceptions as Khmer non-connex vowels not coded as ្រ, ្រ, ្រ, ្រ and ្រ, which have to be rendered with two characters (្រ + ្រ, etc.)²⁶.

In the opposite manner, Thai and Lao Unicode proceed differently in coding vocalic signs which are not systematically phonemes but often only part of them. Examples:

- 0E42 ໊ THAI CHARACTER SARA O,
- 0E30 ໋ THAI CHARACTER SARA A,

But the Thai vowel ໊໋ /o?/ is not coded by itself but needs the typing of the two characters 0E42 ໊ THAI CHARACTER SARA O and 0E30 ໋ THAI CHARACTER SARA A before and after the consonant to be displayed:

0E42 ໊ + 0E15 ๓ + 0E30 ໋ → ໊๓໋ /to?/

- 0EC0 ໄ LAO VOWEL SIGN E,
- 0EB1 ໊ LAO VOWEL SIGN MAI KAN,
- 0E8D ໆ LAO SEMIVOWEL SIGN NYO,

The Lao vowel ໄ໊ /ja?/ is not coded and has to be displayed in typing the three characters 0EC0 ໄ, 0EB1 ໊ and 0E8D ໆ, the first one before the consonant and the two others after it:

0EC0 ໄ + 0E9A ໘ + 0EB1 ໊ + 0E8D ໆ → ໄ໊ໆ /bja?/

Furthermore, scripts encoded following the “Indic” method have their vowels encoded in “logical order”, which means that the vowel character always follows the consonant character, even if the first is shaped before the consonant sign in the writing system²⁷. For example, the Devanagari

²⁵ Let's notice that in Tham script computerized alphasorting will need in all manner an electronic dictionary because of numerous *contractions* [L. 𑄓𑄗𑄚] /k^hamɲə/ like: 𑄓𑄗 /anvə:/ contraction of 𑄓𑄗𑄚, 𑄓𑄗 /di:li:/ contraction of 𑄓𑄗𑄚, etc.

²⁶ Nevertheless those vowels are directly accessible on the keyboard with one single key (KHMEROS 2004:10). That means that keyboard does not correspond to encoding and more again special treatments have to be done on keyboard.

²⁷ This order correspond roughly to typing order (ANDRIES 2002:68). See also note 22.

vowel sign ि *i* is placed before the consonant sign. If one wants to type the syllabe पि *pi* with a Devanagari Unicode font, one should type the consonant character प *p* before the vowel character:

092A प + 093F ि → पि *pi*

It is the Windows Unicode Script Processor that will proceed to the reorganization of the elements and should ensure corresponding between the logical order (characters) and the visual order (glyphs) (UNICODE 2003:260).

The “Thai-Lao” model goes differently in coding characters in graphic order. Thus prevowels signs are coded and typed before the consonant (the same way as handwriting) following a graphic rule. Examples:

0E41 ມ + 0E19 າ → ມາ /mɛ:/,

0EC4 າ + 0E81 ກ → າກ /kaj/.

The Myanmar Unicode system is between the two families in encoding prevowels in logical order but in not coding non-connex vowels:

1000 က + 1031 ဝ → ကဝ /ke:/,

1000 က + 1031 ဝ + 102C ဝ → ကဝဝ /kə/.

Thus the Myanmar consonant ဝဝ /ə/ is considered as an equivalent sequence of the two characters 1031 ဝ MYANMAR VOWEL SIGN E + 102C ဝ MYANMAR VOWEL SIGN A.

This “in-between” is quite strange but we can suppose a will of simplification according to Khmer Unicode in not encoding non-connex vowels.

In conclusion we observe that the implementation of Unicode concerning Indic and derived scripts presents differences in the consonants treatment (between Devanagari, Khmer, Tibetan) as well as in vowel treatment (between Thai-Lao, Khmer, Myanmar). Nevertheless it appears that the “Thai-Lao” input method is the most simple for implementation and utilization, especially for vowels. But because Thai nor Lao scripts use subjoined consonant, we have to choose a system inspired from another script that goes for Tham.

4. Proposal of an input method for Tham Script

4.1 P-Tham consonants

The P-Tham has a common system with Khmer in the use of subscripts: a subscript consonant indicates that the preceding consonant lost its inherent vowel. In the word တၢ်တၢ် *tanhā*, the subjoined form ဝ of တ *h* indicates that the consonant တ *n* has to be read without inherent vowel, so /n/ and not /na/. Thus a system similar as Khmer Unicode can be applied in adding a rendering character which we could call “Tham theoretic virama” (F04B ၵ) that will indicate to the system that the following character has to be rendered subjoined. This character finds here its linguistic justification in “killing” the inherent vowel of the preceding consonant and goes as a *virama*. As this character does not exist in Tham alphabet, we call it “theoretic”.

Example:

$$F00F \text{ ᨧ } + F013 \text{ ᨣ } + \underbrace{F04B \text{ ᨣᨣᨣᨣ } + F01E \text{ ᨣ}}_{\text{ᨣ}} \rightarrow \text{ᨣᨣᨣ}$$

Two particular cases are:

- $\overset{\circ}{\text{ᨣ}}$, subjoined form of ᨣ *n* occurs above the consonant and is a final (i.e. without inherent vowel). The rendering method is the same but has no linguistic justification because the preceding consonant hasn't lose its inherent vowel. To type the word ᨣᨣᨣ *saṅgha* (with ᨣ *s* and ᨣᨣ *gh*):

$$F01D \text{ ᨣ } + F003 \text{ ᨣ } + \underbrace{F04B \text{ ᨣᨣᨣᨣ } + F004 \text{ ᨣ}}_{\overset{\circ}{\text{ᨣ}}} \rightarrow \text{ᨣᨣᨣ}$$

- ᨣ , subjoined form of ᨣ *r*, occurs before the consonant. To type the word ᨣᨣ *bra* we have to choose between “logical order” or “graphic order”:

$$F016 \text{ ᨣ } + \underbrace{F04B \text{ ᨣᨣᨣᨣ } + F01A \text{ ᨣ}}_{\text{ᨣ}} \rightarrow \text{ᨣᨣ} \quad \text{Logical order}$$

$$\underbrace{F04B \text{ ᨣᨣᨣᨣ } + F01A \text{ ᨣ}}_{\text{ᨣ}} + F016 \text{ ᨣ} \rightarrow \text{ᨣᨣ} \quad \text{Graphic order}$$

The implementation in logical order has a linguistic justification (lost of preceding consonant's inherent vowel) and is conform to Unicode recommendations (and to other Unicode Southeast Asian scripts). The graphic order has the advantage to not need any special treatment for reordering and may be easier for users. We think better to choose the graphic order because of the similarity with Tham prevowels (see *infra*, 4.4).

4.2 L-Tham consonants

Two aspects differentiate L-Tham subjoined consonants from P-Tham's.

The first one is the context where subjoined occur. In L-Tham, subscripts occur:

- In a final position of a syllable: ᨣᨣ [ᨣᨣ] /ka:n/

Note that subscript is not systematic in a final position but depends on consonant and vowel: in fact it is graphical criterions that will determine to use or not a subjoined. For example in the word ᨣᨣᨣ [ᨣᨣᨣ] /mōk/ the final /k/ is written in its basic form because of an unaesthetic appearance of the word if written with a subscript as ᨣᨣᨣ .

This first aspect doesn't rule out the subjoined rendering system used in p-Tham:

$$F000 \text{ ᨣ } + F035 \text{ ᨣ } + \underbrace{F04B \text{ ᨣᨣᨣᨣ } + F013 \text{ ᨣ}}_{\text{ᨣ}} \rightarrow \text{ᨣᨣ}$$

Here the P-Tham linguistic justification (lost of preceding consonant's inherent vowel) disappears because the character $F04B \text{ ᨣᨣᨣᨣ}$ follows a vowel (in all manner there is no

inherent vowel in Lao). But the most important is to be able to type L-Tham and users will not take offense for linguistical rules in computer typing.

- In a medial position of a syllable: ລູງ [ລວາ] /lwa:/

One more time, the P-Tham rendering system can be applied:

$$F01B \text{ ລ } + \underbrace{F04B \text{ ັ} + F01C \text{ ອ } + F035 \text{ ງ}}_{\text{ອ}} \rightarrow \text{ລູງ}$$

The second aspect is differences between several subjoined forms used in P-Tham and L-Tham as we evoked in chapter 1.

For example if we take the Lao word ພຼັງ [ພລັງ] /phlan/, a problem appears because the glyph ັ needed for that Lao word is not the glyph ັ used in P-Tham. Thus if we use the character F04B ັ, we'll have:

$$F016 \text{ ັ } + \underbrace{F04B \text{ ັ} + F01B \text{ ລ } + F048 \text{ ັ}}_{\text{ັ}} + F01A \text{ ງ } \rightarrow \text{ັງ}$$

The same problem happens with glyphs ັ (subjoined ັ /k/) and ັ (subjoined ັ /ŋ/), that occur only in L-Tham.

Two solutions have been considered:

- Display glyphs automatically according to context.

It is not possible in Tham because of words written in different manners like ດີຫຼີ /di:li:/ that can be written ຊີ or ຊີ້, ຫອກ /nə:k/ that can be written ທຸ or ທຸ້, etc. The orthography is often let to the choice of the copyist. The use of contractions (l. ຄຳຫຼັ້ /kʰamɰə/) is frequent but also not obligatory. For example Lao word ເອົາ /ao/ can be written ງົງ or ເົງ, word ອັນວ່າ /anvə:/ can be written ວີ or ວີ້, etc.

- Display glyphs with a second subjoiner.

As character F04B ັ THAM THEORIC VIRAMA is employed for all P-Tham subjoined and all common L-Tham/ P-Tham subjoined, one complementary character could be used to display specific L-Tham subjoined forms, called F04C ັ THAM VARIANT SUBJOINER. Thus a user will be allowed to choose between two rendering characters depending on the glyph he wants to display.

+	F000 ັ	F004 ັ	F01B ັ	F01D ັ
F04B ັ	ັ	ັ	ັ	ັ
F04C ັ	ັ	ັ	ັ	ັ

Example with ັ /ŋ/ subjoined rendered ັ or ັ:

$$F00F \text{ ᨧ} + F038 \text{ ᨧ} + \underbrace{F04B \text{ ᨧ} + F004 \text{ ᨧ}}_{\text{ᨧ}} \rightarrow \text{ᨧ}$$

$$F00F \text{ ᨧ} + F048 \text{ ᨧ} + \underbrace{F04B \text{ ᨧ} + F004 \text{ ᨧ}}_{\text{ᨧ}} \rightarrow \text{ᨧ}$$

4.3 Various forms for subjoined consonants

A rendering method with two subjoiner characters allows an advantage in respect of the non-codified aspect of the Tham Script. In manuscripts many various forms occur especially for the subjoined consonants. It is because the Tham Script stills archaic and was never really strictly codified by any national nor religious institution²⁸. Sometimes three or four forms are known for one subscript consonant²⁹. The non-canonic aspect evoked (*supra*, 1) gives importance to graphic forms for physical appearance of texts. For example subscript form of ᨧ *m* can occur as ᨧ, ᨧ or ᨧ,

ᨧ *th* can occur as ᨧ or ᨧ, etc.

As the character THAM THEORIC VIRAMA will be used to render canonical Pali subjoined forms according to Maha Sena’s manuals, the THAM VARIANT SUBJOINER will allow to display L-Tham particular subjoined forms *and* variants (P-Tham and L-Tham). In that manner every consonant has two possibilities for subjoined by using one or the other of those two characters. Of course a system able to render all variants is not possible here (the user who wants to dispose numerous variants should have recourse to several fonts), but having two variants in a single font offers enough possibilities to be quite conform with many texts.

Examples:

+	F00C ᨧ	F018 ᨧ	F012 ᨧ	F01C ᨧ	F01C ᨧ	F01C ᨧ	F01C ᨧ
F04B ᨧ	ᨧ	ᨧ	ᨧ	ᨧ	ᨧ	ᨧ	ᨧ
F04C ᨧ	ᨧ	ᨧ	ᨧ	ᨧ	ᨧ	ᨧ	ᨧ

We precise that the fact that two forms of subjoined are available does not go against Unicode recommendations because only base form consonants and the two subjoiners are coded while the particular glyphs will be included in the font (OpenType or other).

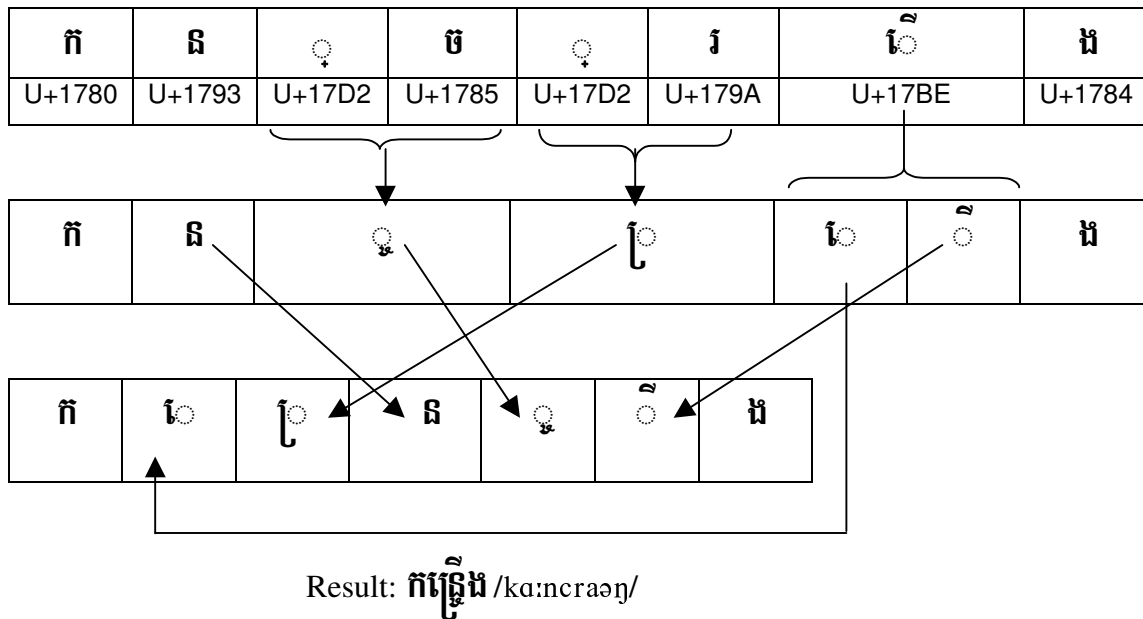
4.4 Vowels (P-Tham and L-Tham)

We’ve seen that two ways of encoding are possible for Unicode, one “Indic” input method and one “Thai-Lao” input method. If the “Indic” method could be in a first approach tempting in a reason of a linguistic logic (a vowel is before all a sound without consideration of the way it is written), some problems appear very quickly when one try to implement with complex words, in particular those

²⁸ The two volumes of Phay Luang Maha Sena respect graphic forms from a lapidary inscription dated from 1819 (GAGNEUX 1983:77). But those graphics were not really layed down as a canonical form by the Buddhist Institute.

²⁹ See (GABAUDE 1979:83-90). The author counted not less than three hundred graphic variations between two ten pages manuscripts.

contain at one and the same time a subjoined consonant and a non-connex vowel. For example to render the word **កំរ្រង** /ka:nkraɲ/, one has to type (example from (BAUHANH 2002:15-16)):



We can see the complexity of the operations that have to be done by the system:

- Transformation of consonants **រ** and **្រ** in their subjoined form **្រ³⁰** and **្រ**,
- Reordering of the glyph **្រ**,
- Split up of **ំ** in **ំ** and **្រ**,
- Reordering of **ំ** (just before **្រ**).

We can observe the misleading operation that consists in splitting up **ំ** in **ំ** and **្រ**, although Unicode insist to code the character **ំ** separately (U+17BE). This ends to important difficulties to make a Khmer Unicode Font: After years and work of numerous computer scientists and linguists, only one set of Khmer Unicode font is now available³¹ (to our knowledge) and it is not yet sure it gives entire satisfaction. The point of view of the users does not convince much more: Unicode admits that the typing is not intuitive and an education has to be delivered to users (UNICODE 2004:277).

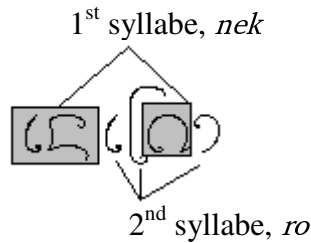
With the example of a word like **កំរ្រង**, which could have equivalent in Tham, we are able to conclude that if we can retain the Khmer Unicode's rendering process for subjoined consonants, the input method for the « Indic » model vowels does not convince. Neither does the “in-between” method of the Myanmar Unicode which does not seem to be more successful than the Khmer Unicode regarding to local users³².

³⁰ In this context the consonant **រ** subscript has to be rendered as **្រ**, more elongated as normal glyph **្រ**.

³¹ The font KhmerOS for Windows is available on http://www.Khmeros.info/Khmeros_download.html. It seems that no Unicode font are available for Macintosh yet.

³² “Although character codes for Myanmar Languages has been allocated in UCS/Unicode (U+1000-U+109F), lack of implementation makes unavailable to local end users. Much effort had been made to develop Myanmar character codes and fonts by many international experts and local experts.” (KO KO, YOSHIKI 2005:9).

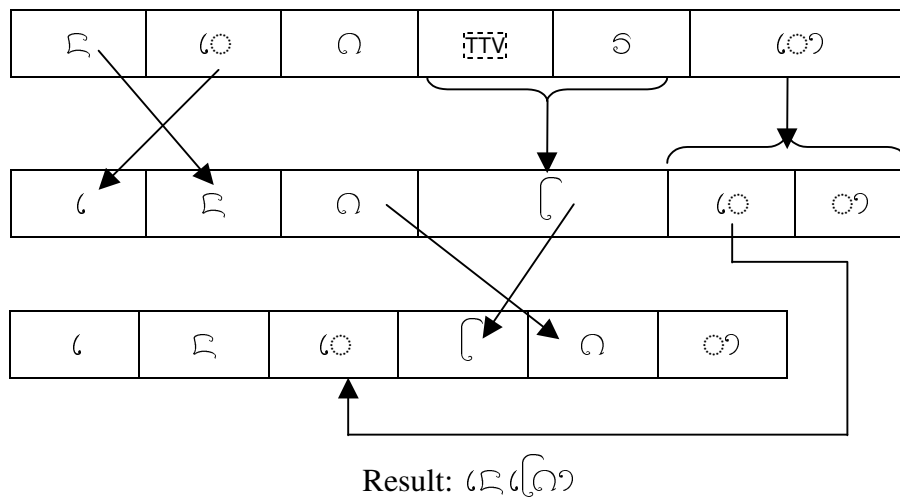
Let's take an example in Tham. The Pali word 𑀅𑀺𑀢𑀺 *nekro*, includes a pre-vowel (𑀺 *e*), a subjoined consonant (𑀢 *r*) and a non-connex vowel sign (𑀺 *o*). The ordering of the disposition of letters does not correspond to the ordering of the reading:



Representation with latin letters: $e n o r k^{33} o$

We can see that the first vocalic graphem (𑀺 , here part of vowel 𑀺 *o*) of the second syllabe is placed before the final consonant (𑀢 *g*) of the first syllabe. That means that visually speaking, the two syllabes cross-over (as in the latter Khmer exemple) and each does not form a monolithic group. To render the word in the Unicode manner, we have to determine how to display the good positioning for the non-connex vowel 𑀺 *o* regarding to the subjoined consonant 𑀢 *r* (for example not to write 𑀅𑀺𑀢𑀺). For this reason, “Indic” method seems to be too hard to implement non-connex vowels as single characters.

In “logical” order, one must type:



The operations executed by the system are:

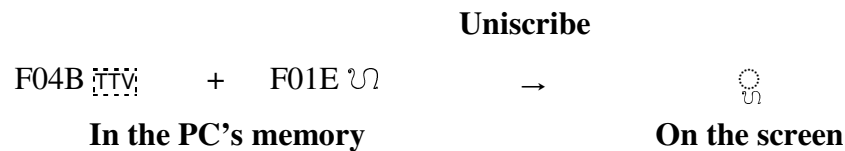
- Transformation of $\text{𑀢𑀺} + \text{𑀺}$ in 𑀢 ,
- Permutation of 𑀺 𑀢 in 𑀺 ,

³³ In a final position, 𑀢 is sounded /k/.

5. ThamWord: An Interim Tool for Tham Unicode Input and Rendering

5.1 Uniscribe: an obstacle to the implementation of a Tham Unicode system

Though the “Thai-Lao” input method we propose is much simpler to implement than the “Indic” one, some processing will need to be applied to render the text properly due to some of the options we took (e.g. to have the same key for the normal and for the subjoined forms, only differing by the \u2728\u2729 prefix). We mentioned *Uniscribe*³⁷ several times and said that it has this role of transforming character chains into its desired graphical form:



Unfortunately, the Tham script is not handled by Unicode yet so \u2728\u2729 will not appear on the screen when we will type the two characters or have them in the PC's memory. In a word, Uniscribe is not working for Tham yet.

5.2 ThamWord: An interim tool to bypass the obstacle

Until the days where Uniscribe will handle Tham, we propose ThamWord, an interim tool that will simulate its behavior. To do that, we first have to design a font in the Unicode private area. We selected the range F000-F0FF as suggested in (BERMENT 2004:153) for this font that we called *ThamUnicode*. In addition to the code points described in (KOURILSKY 2005:174-194), this font includes other points for the subjoined characters. It actually covers all the 132 characters (plus variants) as does the *ThamStandard* font so all the Tham glyphs can be displayed through a code. The extra codes are taken from the code points of the range that remained unused.

Then, when typing for example F04B \u2728\u2729 and then F01E \u2728 , the tool will replace the two codes by a third one which represents \u2728\u2729 . Note that this processing is different from what Uniscribe would do, as Uniscribe would only compute the graphical form associated to the two codes and ThamWord actually replaces the two codes themselves in the memory.

5.3 Virtual keyboards and technical solution

If the rendering of the text is probably the most complex issue in the absence of a Tham-capable Uniscribe, an interim tool has also to provide a Tham virtual keyboard. Several techniques have been described for that purpose in (BERMENT 2004). We chose to implement this function in a hook procedure that we derived from LaoWord, a Word add-in adapted to the Lao language that contains such a hook for its virtual keyboards (BERMENT 2003).

³⁷ The Unicode Script Processor embedded inside Windows.

The main features of ThamWord are:

- Providing an *ad-hoc* virtual keyboard for Tham,
- Transforming the newly typed characters and their environment according to the contextual rules described before.

The two features are included in the hook procedure which is launched by a Word add-in during its initialization. This add-in is a dynamic library (DLL) that is mapped in the memory space of Word when this application starts. It provides a human-machine interface through a toolbar added to the Word's toolbars.

The virtual keyboard is modifying the keycodes sent to Word following a table of correspondence defined both by the keyboard mapping and by the Tham Unicode zone³⁸. At this level, the font used doesn't need to have any glyph but the ones defined by a code in the Tham Unicode zone.

After this modification of the code, the hook procedure analyses the environment around the insertion point. If a rule has to be applied for the entered code in the current context, another modification is done, possibly with a deletion of some characters around the new one as for F04B $\text{ᨧᩢ᩠ᨦ} + \text{F01E } \text{ᨧ} \rightarrow \text{ᨧ}$ where ᨧᩢ᩠ᨦ is removed from the memory when ᨧ appears in the hook queue and is processed. In order to explore the environment around the insertion point, the hook procedures has to dialog with Word. This is done thanks to the Component Object Model (COM) interface.

5.4 Preparing the future

We can hope that the Tham writing system will soon be part of the Universal standard that Unicode aims to become and that the operating systems will be able to natively process complex Tham fonts (OpenType...). When this will happen? Who knows. At that time, the texts made in between with ThamWord or other equivalent software will constitute an existing corpus that will need to be slightly modified for fully complying the standard that was defined. Reverse rules like $\text{ᨧ} \rightarrow \text{ᨧᩢ᩠ᨦ} + \text{F01E } \text{ᨧ}$ will have to be applied to go back to the original input. This mechanism will soon be added to ThamWord.

³⁸ Actually it is still in the private area, as the submission to Unicode of a Tham zone is ongoing. See (KOURILSKY 2005).

Sources

- ANDRIES, P. (2002), *Introduction à Unicode et à l'ISO 10646*, Québec, Document numérique, Volume 6, n° 3-4, p. 51-88.
- BAUHANH, M. (2002), "Rendering the World's Complex Scripts: A Case Study in Khmer", 21th International Unicode Conference, Dublin, 20 p.
- BECCHETTI, C. (1991), *Le mystère dans les lettres*, *ปริศนาธรรมในอักษร*, Édition des Cahiers de France, Bangkok, 116 p.
- BERMENT, V. (2003), *LaoWord's Word Processing Functions*, Pan Asia Networking All Partners 2003, March 3 - March 10, 2003, Vientiane Novotel, Vientiane, Laos, 10 p.
- BERMENT, V. (2004), *Méthodes pour informatiser des langues et des groupes de langues « peu dotés »*, Grenoble, P.H.D. Thesis (Université Joseph Fourier), 277 p.
- BIZOT, F. (1992), *Le Chemin de Laïkā*, Paris, ÉFEO, Textes bouddhiques du Cambodge n° 1, Paris, 352 p.
- BIZOT, F. (1996), *ບູລິມະສິດ, La pureté par les mots*, Paris, Phnom Penh, Chiang Mai, Vientiane, ÉFEO, Textes bouddhiques du Laos, 275 p.
- FERLUS, M. (1995), "Les circonstances de l'introduction de l'alphabet Tham lanna", in: *La Thaïlande des débuts de son histoire jusqu'au XV^e siècle*, Premier Symposium Franco-Thai, 18-24 juillet 1988, Silpakorn University, p. 101-109.
- GABAUDE, L. (1979), *Les cetiya de sable au Laos et en Thaïlande. Les textes*, Paris, ÉFEO, PÉFEO 118, 96 + 338 p.
- GABAUDE, L. (1988), *Une herméneutique bouddhique contemporaine de Thaïlande : buddhadasa bhikkhu*, Paris, ÉFEO, PÉFEO CL, 690 p.
- GAGNEUX, P.M. (1983), "Les écritures lao et leur évolution du XV^e au XIX^e siècles", ASEMI, vol. XIV, 1-2, p.75-95.
- INTHAMONE, L. (1987), *Je lis et j'écris lao*, Paris, INALCO, 207 p.
- KO KO, W. YOSHIKI, M. (2005), "Languages of Myanmar in Cyberspace", Dourdan (France), TALN Conference, June 2005, 10 p.
- KOURILSKY G. (2005), *Éléments pour un traitement informatique de l'écriture tham du Laos*, Master Thesis (INALCO), 330 p.
- PELTIER, A.R. (2000), *ເຂົ້າໂຮມ, นางเฝ้าหมอก, La fille aux cheveux parfumés*, Vientiane, Institut de Recherches sur la Culture, 490 p.
- PINAULT, G.J. (2001), "Les écritures de l'Inde continentale", in *Histoire de l'écriture, de l'idéogramme au multimédia* (dir. Anne-Marie Christin), Paris, Flammarion, p.93-121.
- SENA, P.L. (MAHA SENA) (1957), *ແບບ ຮຽນ ໄວ ຮຽນ ອ່ານ ຫຼັງສື່ທັມ ຂຽນເປັນພາສາລາວ [apprendre rapidement à lire les caractères Tham dans les textes Lao]*, Bangkok, Kramol Tirannasur, 80 p.
- SENA, P.L. (MAHA SENA) (1963), *ແບບ ຮຽນ ໄວ ຮຽນ ອ່ານ ຫຼັງສື່ທັມ ຂຽນເປັນພາສາປາລີ [apprendre rapidement à lire les caractères Tham dans les textes Pali]*, Bangkok, Kramol Tirannasur, 51 p.
- SOLÁ, J. (2004), "Displaying Khmer", KhmerOS Initiative, <http://www.khmeros.info/download/DisplayingKhmerScript.pdf>, 4 p.

VIRAVONG, S. (MAHA SILA) (1937), *ບາລີເຊຍາກອດ ອັກຂະວິຣິກາທິນິງ ສັມບູລາກິສານແລະສັນຮິ*, [Grammaire Pali, Akkharavidhi, première partie, Samaññībhīdāna - Sandhi], Vientiane, Institut Bouddhique, 68 p.

Collectives :

KHMEROS (2004), "Creating and supporting OpenType fonts for the Khmer script", <http://www.khmeros/download/KhmerUnicodeTyping.pdf>, 18 p.

MICROSOFT CORPORATION (2004a), "Creating and supporting OpenType fonts for the Khmer script", <http://www.microsoft.com/typography/otfntdev/khmerot/default.htm>.

MICROSOFT CORPORATION (2004b), "Other encoding issues: Khmer OpenType specification", <http://www.microsoft.com/typography/otfntdev/khmerot/other.htm>.

UNICODE (2003), Unicode 3.1 annoté, <http://iquebec.hapax.com>.

UNICODE (2004), The Unicode Standard 4.0, <http://www.Unicode.org>.

PHONE PHRA NAO MONASTERY [ວັດໂພນພຣະເນົາ] (1965), *ຫັມມະຄາຖາ* [dhammagāthā], Vientiane, 88 p.