

# Typesetting Khmer

YANNIS HARALAMBOUS

*187, rue Nationale  
59800 Lille  
France*

---

## SUMMARY

Because of the complexity of Khmer script, up to now there has been neither a typesetting system nor standard encoding for the Khmer language. Presented in this paper are: (a) a complete typesetting system for Khmer based on  $\text{\TeX}$ , METAFONT and an ANSI C preprocessor, as well as (b) a proposal for an 8-bit encoding table for Khmer information interchange. Problems of phonic input, subscript and superscript positioning, collating order, spelling reforms and hyphenation are solved, and their solutions described. Finally an alternative solution using 16-bit output font tables is briefly sketched.

KEY WORDS Khmer  $\text{\TeX}$  METAFONT Computer typesetting

## 1 INTRODUCTION

Certain languages use characters or character combinations which change according to the context. A common example in English (but not in Portuguese and Turkish!) are the ‘fi’ and ‘fl’ ligatures. Everytime he/she encounters the combination of letters ‘f’, ‘i’, the typesetter has to replace it by the ligature ‘fi’. This practice, while remaining exceptional for the Latin script, becomes very important for certain Oriental scripts like Arabic (see [1], [2]), Indic scripts, Korean or Khmer.

Because of the repetition and transformation of the various shapes involved in this process, the best way of creating a font with strong contextual properties is to use a programming language, like METAFONT. Part of the contextual analysis can be done by  $\text{\TeX}$  (in simple cases, such as Latin or modern Arabic), otherwise one has to use an independent preprocessor.

In this paper we present a typesetting system for one of the most complicated scripts: Khmer. In this case, the  $\text{\TeX}$ /METAFONT/preprocessor approach is essential. Since there has been no standardization for Khmer information interchange yet, we also present a proposal for a Khmer 128-character table. This table has been submitted to ISO 10646 WG-2 for acceptance. Finally, the solutions to other typesetting problems such as hyphenation are also presented.

## 2 THE KHMER SCRIPT

The Khmer script is used to write Khmer, which is the official language of the Cambodian Republic and belongs to the Mon-Khmer group of Austroasiatic languages. It is a very old and beautiful script, and from the typesetter’s point of view, one of the most challenging and exciting scripts in the world.



We will call the combination of consonant and possible subscript consonant, second subscript consonant, vowel and diacritical mark, a *consonantal cluster*. Theoretically there can be 535 060 different consonantal clusters, but in practice less than 1% of them are really used. An analytic decomposition of A. Daniel's Khmer-French dictionary [3] has provided no more than 2 821 different consonantal clusters out of 25 000 entries; colloquial Khmer may require even fewer clusters.

Besides consonantal clusters there are also 14 'stand-alone' characters in the Khmer alphabet:

These carry neither subscript consonants, nor vowels, nor accents. They cannot be found in subscript form. Orthographical reforms of Khmer have in some cases replaced them by 'regular' consonantal clusters.

Inside a sentence, Khmer words are *not* separated by blank space. A blank space denotes the end of a sentence (or of part of a sentence: it acts like the period or the semicolon in Latin script).

Hyphenation occurs between *syllables*: a syllable consists of one or two consonantal clusters with the sole restriction that the second cannot have a vowel. When a word is hyphenated, a hyphen is used. Sentences are 'hyphenated' into words, but in that case, no hyphen is used. So from the typesetter's point of view, between two clusters hyphenation can be

1. forbidden (when the two clusters belong to the same syllable);
2. allowed and produce a hyphen (when the two clusters belong to the same word);
3. allowed without producing a hyphen (when the two clusters belong to different words in the same sentence).

This quick overview of the Khmer script has shown some of its particularities (see also [4], [5], [6]). To conclude, the author would like to underline the fact that the main difficulty in Khmer typesetting is the divergence between phonic and graphical representation of consonantal clusters (see Figure 1).

This paper is divided into five sections:

1. [the definition and discussion](#) of an 8-bit encoding table for information interchange and storage in the Khmer script. Consonantal clusters are encoded according to their phonic representation;
2. [the presentation of three Khmer font families](#), designed in the METAFONT language. These fonts correspond to the three main styles of Khmer type and provide sufficient *metaness*<sup>1</sup> to perform optical scaling, continuous interpolation from light to extra-bold weight and strong raster optimization;
3. [the description of the process](#) according to which the graphical representation of consonantal clusters is derived from the phonic one (this process being implemented in an ANSI C preprocessor);
4. [an overview of hyphenation and spelling reform rules](#) and their realization in the pre-processor;
5. [shortcomings of the Khmer typesetting system](#) and plans for future developments.

<sup>1</sup> In METAFONT lingo, *metaness* is the possibility of parametrized variation of the characters' shape, weight and style.

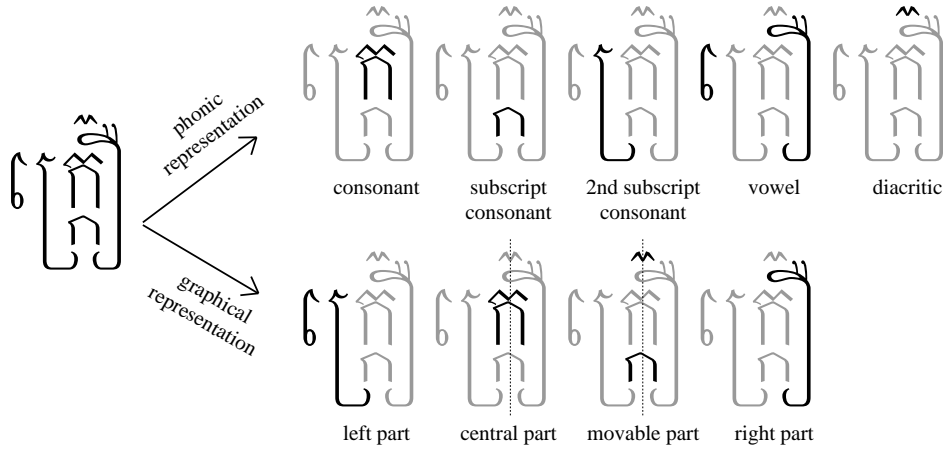


Figure 1. Decomposition of a Khmer consonantal cluster

### 3 AN 8-BIT ENCODING TABLE FOR THE KHMER SCRIPT

#### 3.1 Discussion

As mentioned in the introduction, Khmer language is written using consonantal clusters and stand-alone special characters. The collating order of consonantal clusters is given lexicographically according to the cluster components:<sup>2</sup>

Let  $C_1 = c_1 s_1 s'_1 v_1 d_1$  and  $C_2 = c_2 s_2 s'_2 v_2 d_2$  be two consonantal clusters, where  $c_1, c_2 \in \{\text{consonants}\}$ ,  $s_1, s_2 \in \emptyset \cup \{\text{subscript consonants}\}$ ,  $s'_1, s'_2 = \emptyset$  or  $\llcorner$ ,  $v_1, v_2 \in \emptyset \cup \{\text{vowels}\}$  and  $d_1, d_2 \in \emptyset \cup \{\text{diacritics}\}$ . Then

1.  $c_1 \succ c_2 \Rightarrow C_1 \succ C_2$ ;
2. if  $c_1 = c_2$  then  $s_1 \succ s_2 \Rightarrow C_1 \succ C_2$  (where  $\emptyset$  precedes any other element);
3. if  $c_1 = c_2$  and  $s_1 = s_2$  then  $s'_1 \succ s'_2 \Rightarrow C_1 \succ C_2$ ;
4. if  $c_1 = c_2, s_1 = s_2$  and  $s'_1 = s'_2$  then  $v_1 \succ v_2 \Rightarrow C_1 \succ C_2$ ;
5. if  $c_1 = c_2, s_1 = s_2, s'_1 = s'_2$  and  $v_1 = v_2$  then  $d_1 \succ d_2 \Rightarrow C_1 \succ C_2$ .

The table of 128 codes for Khmer characters presented below respects the collating order. Besides consonantal clusters and special characters, the following signs have been included in the 8-bit encoding:

1. digits: ០, ១, ២, ៣, ៤, ៥, ៦, ៧, ៨, ៩;
2. punctuation marks other than the ones borrowed from Latin script: ្ក (leikto) a variant form of the digit ២, indicating that the previous word is repeated (similar to Latin *bis*), ្ខ (khan) and ្គ (bariyatosan), equivalent to a full stop, ្ឃ (cammocpikuh) a graphical variant of the Latin colon, and the French *guillemets* « , »;
3. the currency symbol ៛ (rial);
4. the invisible code WBK (word-break) to indicate the word limits inside a sentence.

<sup>2</sup> The symbol  $\emptyset$  denotes an empty set.

The following have *not* been included in the table:

- the archaic characters 𑀓 and 𑀔 which were abolished about a century ago;
- the punctuation marks 𑀓 (cow's urine) and 𑀔 (cock's eye), used in poetry, divination and classical texts;
- the variant forms 𑀓, 𑀔 of 𑀓, 𑀔, used in [7].

These characters are nevertheless included in the  $\TeX$  output fonts and can be accessed via macros.

### 3.2 The Table

On Figure 2 the reader can see the table of codes 128–255 of the proposed 8-bit encoding for Khmer information interchange and storage. The 7-bit part of the table conforms to ISO 646 (standard 7-bit ASCII). Positions 0xCF and 0xDF are empty.

Codes 0x80–0x9F and 0xC0 represent consonants; the corresponding subscript consonants are offset by 32 positions: they are represented by codes 0xA0 – 0xBE and 0xE0. The consonant 0x9F does not have a corresponding subscript consonant. The practice of having subscripts 32 positions apart from primary consonants is similar to the 32-position offset of uppercase and lowercase letters in ISO 646 (7-bit ASCII).

Codes 0xC0–0xCE represent special characters. Digits have been placed in positions 0xD0–0xD9, vowels in 0xE1–0xF5 and diacritics in 0xF8–0xFF. Finally, 0xFA is the currency symbol, 0xDB–0xDE are punctuation marks and 0xBF is the word-break code WBK.

Because of the 128-character limitation, the following characters have not been included in the table: 𑀓, 𑀔, 𑀕, 𑀖, 𑀗, 𑀘, 𑀙, 𑀚, 𑀛:

They have to be represented by the following code pairs:

𑀓	=	0xE2	0xF4	𑀔	=	0xE4	0xF4	𑀕	=	0xE6	0xF4
𑀖	=	0xE9	0xF4	𑀗	=	0xEC	0xF4	𑀘	=	0xED	0xF4
𑀙	=	0xEF	0xF4	𑀚	=	0xEF	0xF4	𑀛	=	0xEF	0xF4

### 3.3 Requirements for Khmer script software

As in the case of Arabic and Hindi, software displaying Khmer text has to provide context-analytic algorithms. Below is an exhaustive list of the necessary context-dependent transformations:

1. When code 0xBA follows a code in the range 0x80–0x9E, 0xC0 then their glyphs must be permuted, for ex. 𑀓 + 𑀓 → 𑀔.
2. When code 0xBA follows a pair of characters  $\alpha\beta$ , with  $\alpha \in \{0x80–0x9E, 0xC0\}$ ,  $\beta \in \{0xA0–0xBE, 0xE0\}$  then the glyph of 0xBA must appear on the left of the glyphs of  $\alpha, \beta$ , for ex. 𑀓 + 𑀓 + 𑀓 → 𑀔𑀓.
3. When codes 0xE9–0xEC and 0xEF–0xF0 follow a combination of character codes  $\alpha, \alpha\beta, \alpha\beta\gamma$  where  $\alpha$  and  $\beta$  are as in the previous item and  $\gamma = 0xBA$ , then the glyph 𑀔 must appear on the left of the latter combinations. Example: 𑀓 + 𑀓 + 𑀓 + 𑀔 → 𑀔𑀓𑀓.

Hexa	80	90	A0	B0	C0	D0	E0	F0
0	ក	បំ	្ក	្ខ	អ	្ង	្ច	្ឆ
1	ខ	្ជ	្ឈ	្ញ	ត	្ឋ	្ឌ	្ឍ
2	ត	ត	្ណ	្ត	្ថ	្ទ	្ធ	្ន
3	ឃ	ឆ	្ប	្ផ	ឡ	្ភ	្ម	្យ
4	ង	ប	្រ	្ល	ឡ	្ល	្ល	្ល
5	ច	ត	្វ	្ឝ	ឡ	្ឝ	្ឝ	្ឝ
6	ឆ	ព	្ឞ	្ស	ប	្ឞ	្ស	្ស
7	ជ	ភ	្ហ	្ឡ	ប	្ហ	្ឡ	្ឡ
8	ឈ	ម	្អ	្ឣ	ព	្អ	្ឣ	្ឣ
9	ញ	យ	្ឤ	្ឥ	ព	្ឤ	្ឥ	្ឥ
A	ដ	រ	្ឦ	្ឧ	ជ	្ឦ	្ឧ	្ឧ
B	ប	ល	្ឨ	្ឩ	ព	្ឨ	្ឩ	្ឩ
C	ឌ	វ	្ឪ	្ឫ	ឡ	្ឪ	្ឫ	្ឫ
D	ឍ	ស	្ឬ	្ឭ	ឡ	្ឬ	្ឭ	្ឭ
E	ណ	ហ	្ឮ	្ឯ	ឡ	្ឮ	្ឯ	្ឯ
F	ត	្ឰ	្ឱ	WBK			្ឱ	្ឱ

Figure 2. Positions 0F80–0FFF of ISO 10646 (proposal)

- When codes 0xED and 0xEE follow a combination  $\alpha, \alpha\beta, \alpha\beta\gamma$  of codes as in the previous item, then their glyphs must appear on the left of these combinations.
- When code 0x89 (្ញ) is followed by a code in the range 0xA0–0xBE, 0xE0 then the variant glyph ្ញ must be used. Example: ្ញ + ្ក → ្ញ្ក. When the second code is 0xA9 then a variant glyph must be used for it as well: ្ញ + ្ខ → ្ញ្ខ.

These contextual transformations have been implemented by the author into a modified version of the Macintosh freeware text editor Tex-Edit by Tim Bender, included in the package. In Figure 3 the reader can see the effect of striking successively keys <ក>, <្ក> (<subscript modifier> followed by <ក>), <្ក្ក> (<subscript modifier> followed by <្ក>), <្ក្ក្ក>, to finally obtain the consonantal cluster ្ក្ក្ក.

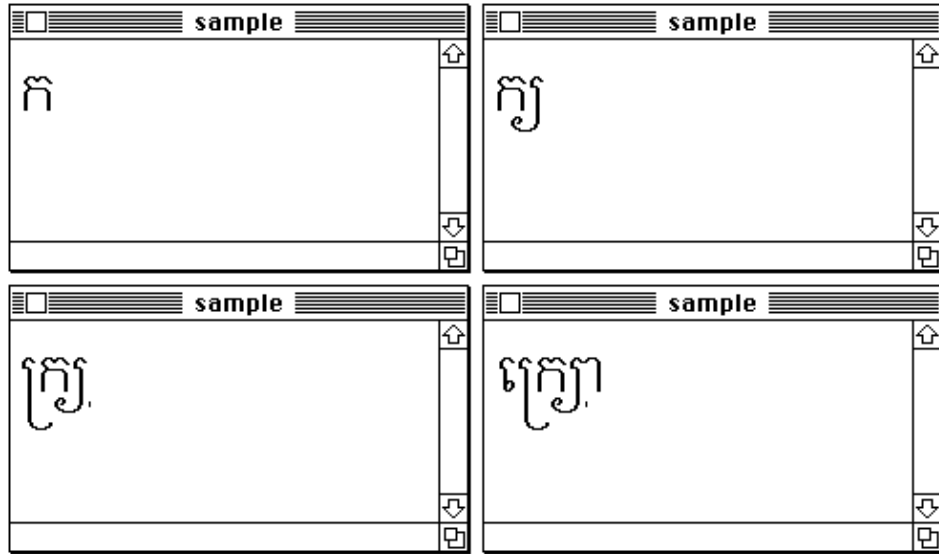


Figure 3. A text editor with Khmer contextual properties

#### 4 THE DESIGN OF KHMER FONTS IN METAFONT

##### 4.1 Font styles

There are three styles used in Khmer typesetting: standing (aksar ch-hor), oblique (aksar chrieng) and round (aksar mul). The last is virtually identical to inscriptions of the 12th and 13th centuries at Angkor Wat and is reserved for religious texts, chapter headings, newspaper headlines, inscriptions and other occasions where there is a desire to create a contrast with the oblique script, to add a touch of formality, or to provide variation of emphasis (see [5]).

The author has designed three METAFONT font families, corresponding to these styles; samples of these fonts in 14.4 point size follow hereafter; Figure 4 shows a sample headline using the round font.

##### Standing characters

នៅឆ្នាំ ១៩៦២ ប្រទេសកម្ពុជាបានចេញមក នៅសល់តែស្បែកនឹង-  
 ភ្លើង អំពីការប្រយុទ្ធនៃយុវអង្គ ដឹងជ័យោរ យោមួយ ប្រឆាំងនឹង-  
 ការដឹកនាំរបស់យួន ដឹងសៀម ។ ការឈឺចាប់នៃប្រជាពលរដ្ឋ-  
 យើងនៅពេលនេះ គឺហួសពីការគិតទៅហើយ ប៉ុន្តែពួកបស្ចឹមប្រ-  
 ទេស គេនៅតែពុំដឹងឮ ចំពោះការឈឺចាប់នេះទេ ។ ចំណែកខ្មែរវិញ  
 យើងមិន មានភ្នំសោះឡើយ ។

Oblique characters

នៅឆ្នាំ ១៩៦២ ប្រទេសកម្ពុជាបានចេញមក នៅសល់តែ-  
 ស្បែកគឺជំនួញ អំពីការប្រយុទ្ធនឹងយួរអង្កែង គឺជំងឺយោរ យោ-  
 មួយ ប្រឆាំងនឹងការជិះជាន់របស់យួន គឺជំងឺសៀម ។ ការឈឺ-  
 ចាប់នៃប្រជាពលរដ្ឋយើងនៅពេលនេះ គឺហួសពីការគិតទៅ-  
 ហើយ ប៉ុន្តែពួកបស្ចឹមប្រទេស គេនៅតែពុំដឹង ចំពោះការ-  
 ឈឺចាប់នេះទេ ។ ផំណែកខ្មែរវិញ យើងមិន មានភ្នួចសោះ-  
 ឡើយ ។

Round characters

នៅឆ្នាំ ១៩៦២ ប្រទេសកម្ពុជាបានចេញមក នៅសល់តែស្បែក-  
 និងន្លិច អំពីការប្រយុទ្ធនឹងយួរអង្កែង និងជំងឺយោរ យោមួយ ប្រ-  
 ឆាំងនឹងការជិះជាន់របស់យួន និងសៀម ។ ការឈឺចាប់នៃប្រ-  
 ជាពលរដ្ឋយើងនៅពេលនេះ គឺហួសពីការគិតទៅហើយ ប៉ុន្តែ-  
 ពួកបស្ចឹមប្រទេស គេនៅតែពុំដឹង ចំពោះការឈឺចាប់នេះទេ  
 ។ ផំណែកខ្មែរវិញ យើងមិន មានភ្នួចសោះឡើយ ។



Figure 4. Headline in round style

In contrast to systems like PostScript, in which fonts are interpreted during the printing process and where similar complexity can slow the process down, in the case of METAFONT (see [8]) fonts are compiled separately and stored on disk in a highly compacted



form. On powerful platforms, fonts can be created by METAFONT just before the printing process, stored on hard disk, and removed afterwards. On slower platforms (for example personal computers), fonts are stored permanently on the hard disk. A METAFONT font package takes much more space than a set of PostScript fonts; this disadvantage is counterbalanced by the fact that fonts created by METAFONT are under the complete control of the user: characters are already rasterized in an optimal and homogeneous way. Other advantages of using METAFONT, in particular for the design of Khmer fonts, are the following:

1. *Modularity.* Characters are designed in a modular way: descriptions of parts which are repeatedly used are stored as subroutines with an arbitrary number of parameters for adapting them to different situations where they can occur. A modular design makes the font more homogeneous and easier to modify: a change in a subroutine will affect the whole font.
2. *Metaness.* In Khmer, standing and oblique letters share the same design, except that certain curves of the latter are rounder than the corresponding curves of the former (for example, compare ឃ with *Ch*, or ឃ with */rKh/Ch*). To preserve the similarity between the two styles, standing and oblique fonts are generated using the same METAFONT code; only the values of slant, roundness and interletter spacing parameters are different.
3. *Raster optimization.* Vertical strokes of Khmer letters must always be of the same width, regardless of the resolution or of the output device. In the METAFONT “drawing space”, coordinates are given in pixels; this is possible because of the fact that output device characteristics are given at the beginning of the METAFONT run. The condition “two vertical strokes should have the same width” is given by a simple linear equation “*width of left stroke = width of right stroke*” with the sole precaution that the left edges of strokes fall on the pixel raster (this is obtained by the METAFONT primitive `round`).

To illustrate this we will take the example of the two vertical strokes of letter ឃ. As in [Figure 5](#) let  $\mu_s$  call  $z_1$  and  $z_2$  the points which are on the baseline and on the central paths of the two vertical strokes. Also let  $z_3, z_4$  be the upper extremities of central paths of the two vertical strokes. Let the straight segments  $[z_{1l}, z_{3l}], [z_{2l}, z_{4l}]$  be the left edges of the vertical strokes, and  $[z_{1r}, z_{3r}], [z_{2r}, z_{4r}]$ , their right edges. The fact that  $z_1, z_2$  are on the baseline can be expressed as  $y_1 = y_2 = 0$ , where  $y_*$  is the  $y$ -coordinate of  $z_*$ . In the same way, the fact that the strokes are vertical can be expressed by the couple of equalities  $x_1 = x_3, x_2 = x_4$  where  $x_*$  is the  $x$ -coordinate of  $z_*$ . Their precise location is given as a multiple of a global variable  $w$ , corresponding to the width of a generic character:  $x_{1l} = 1/14w, x_{2l} = 0.86w$ . As mentioned in the previous paragraph, the two strokes must have the same width (called *stem*). Since they are vertical, we can determine their width by using only  $x$ -coordinates. The width equality can be expressed as  $x_{2r} - x_{2l} = x_{1r} - x_{1l} = \text{stem}$ .

So far, so good. But let us consider an example in which things can go wrong. METAFONT does its calculations in pixels or fractional parts of pixels which are rounded afterwards to the closest integer value. Let us suppose that the two strokes are  $\text{stem} = 2.2$  pixels wide. Of course this should always be rounded to 2 pixels. Now suppose  $x_{1l} = 1/14w = 1.2$  and  $x_{2l} = 0.86w = 14.4$ . Values will be rounded in the following way:  $x_{1l} = 1.2 \rightarrow 1, x_{1r} = 1.2 + 2.2 = 3.4 \rightarrow 3$ , and hence  $x_{1r} - x_{1l} = 2$ ,

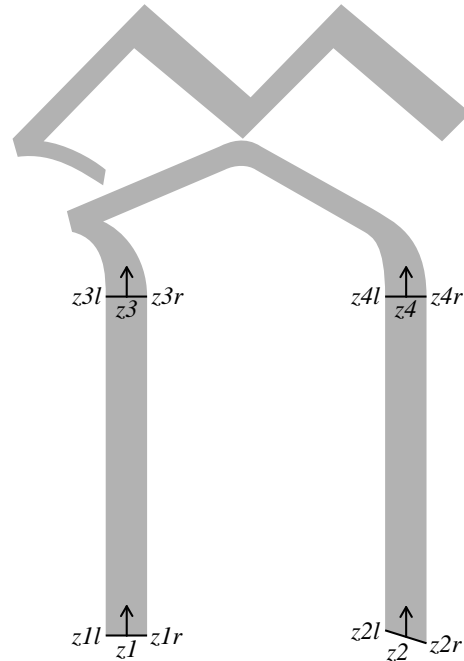


Figure 5. Raster optimization with METAFONT

while  $x_{2l} = 14.4 \rightarrow 14$ ,  $x_{2r} = 14.4 + 2.2 = 16.6 \rightarrow 17 \Rightarrow x_{2r} - x_{2l} = 3$  and so the right stroke is one pixel wider than the left one.

To prevent this, one simply instructs METAFONT to round values *while calculating point locations*. By writing  $x_{1l} = \text{round}(1/14w)$  and  $x_{2l} = \text{round}(0.86w)$ , both  $x_{1l}$  and  $x_{2l}$  will have integer values (in our example, 1 and 14); this implies that the fractional parts of  $x_{1r} - x_{1l}$  are the same  $x_{2r} - x_{2l}$  and hence both will get rounded in the same way (either to the right or to the left), and their values will remain equal after rounding.

The METAFONT code<sup>3</sup> that follows implements these operations; it is meant to illustrate the ease of raster optimization (*hinting*, in PostScript lingo) in this programming language. Lines starting with % are comments.

```
x1r-x1l=x2r-x2l=x3r-x3l=x4r-x4l=stem;
% strokes are of same width "stem"
x1r=x3r; x2r=x4r;
% and they are vertical
x1l=round(1/14w); x2l=round(0.86w);
% their left edges take integer pixel values
fill z1r--z3r--z3l--z1l--cycle;
fill z2r--z4r--z4l--z2l--cycle;
% fill the strokes with black
```

<sup>3</sup> This code is voluntarily kept simplistic; there are more elegant ways to program the same operations.

The reader may have noticed that the ‘pen position’  $z_2$  is actually oblique, while  $z_1$  is horizontal. This fact has not influenced rasterization, since all roundings done in this example are on the  $x$ -coordinate level. For the  $z_2$  pen position, a different kind of optimization can be performed: in low resolutions the straight segment  $[z_{2l}, z_{2r}]$  may look ‘broken’. This will mean that the angle being too small with respect to the pixel size, the segment will be displayed as a certain number of concatenated horizontal rows of pixels. The number of these rows is the rounded value of  $y_{2l} - y_{2r}$ . We can decide to replace the oblique segment  $[z_{2l}, z_{2r}]$  by a horizontal one, if this number is smaller than a certain value, for example 2. This will be written as:

```
if (round(y2l-y2r) <= 2): y2l:=y2; y2r:=y2; fi
```

where the assignment operator  $:=$  will change the values of  $y_{2l}, y_{2r}$ .

Special care has been taken for raster optimization, since output devices in Cambodia are mostly of very low resolution.

4. *Parametrization and optical scaling.* When type is scaled, widths of strokes are not necessarily scaled by the same factors. Large point sizes must be narrower and thinner proportionally to standard point size; small point sizes must be larger and with increased interletter space, to enhance readability. This problem is very well known for the Latin script and is solved in METAFONT-created font families like Computer Modern. Similar solutions have been adopted for Khmer.

There is a second advantage of optical scaling and parametrization of character shapes. In Cyrillic and Greek scripts one can define font families similar to Latin ones: there already exist Cyrillic and Greek Times, Helvetica, Courier, Garamond, Baskerville etc. The choice of a Khmer/Latin font combination is more delicate. Parametrization of character widths gives the user the possibility of changing the *grey density* factor of the Khmer font and adapt it to the Latin font he is using.

In Figure 6, the Khmer letter វ៉ has been reproduced 256 times, with different values of two parameters: the widths of ‘fat’ and ‘thin’ strokes. The central vertical symmetry axis represents ‘Égyptienne’-like characters, where the parameters have the same value. This classification can of course be refined and enables an arbitrarily precise choice of the font grey density.

## 5 TRANSLATING A PHONIC TO A GRAPHIC DESCRIPTION

In Section 3.3 we have given a quick overview of the minimal contextual analysis involved in displaying Khmer script on screen. The situation is much more complicated in the case of high-quality typesetting.

TeX is the ideal tool for typesetting in Oriental scripts like Khmer, because of the inherent fundamental concept of *boxes* (see [9], [10], [11]). As in mathematical formulas, elements of a consonantal cluster are moved to aesthetically correct positions and then grouped into a single and indivisible “box” which TeX treats as a single entity.

In this section we will see how the graphical representation of a cluster is constructed, using both the preprocessor and TeX.

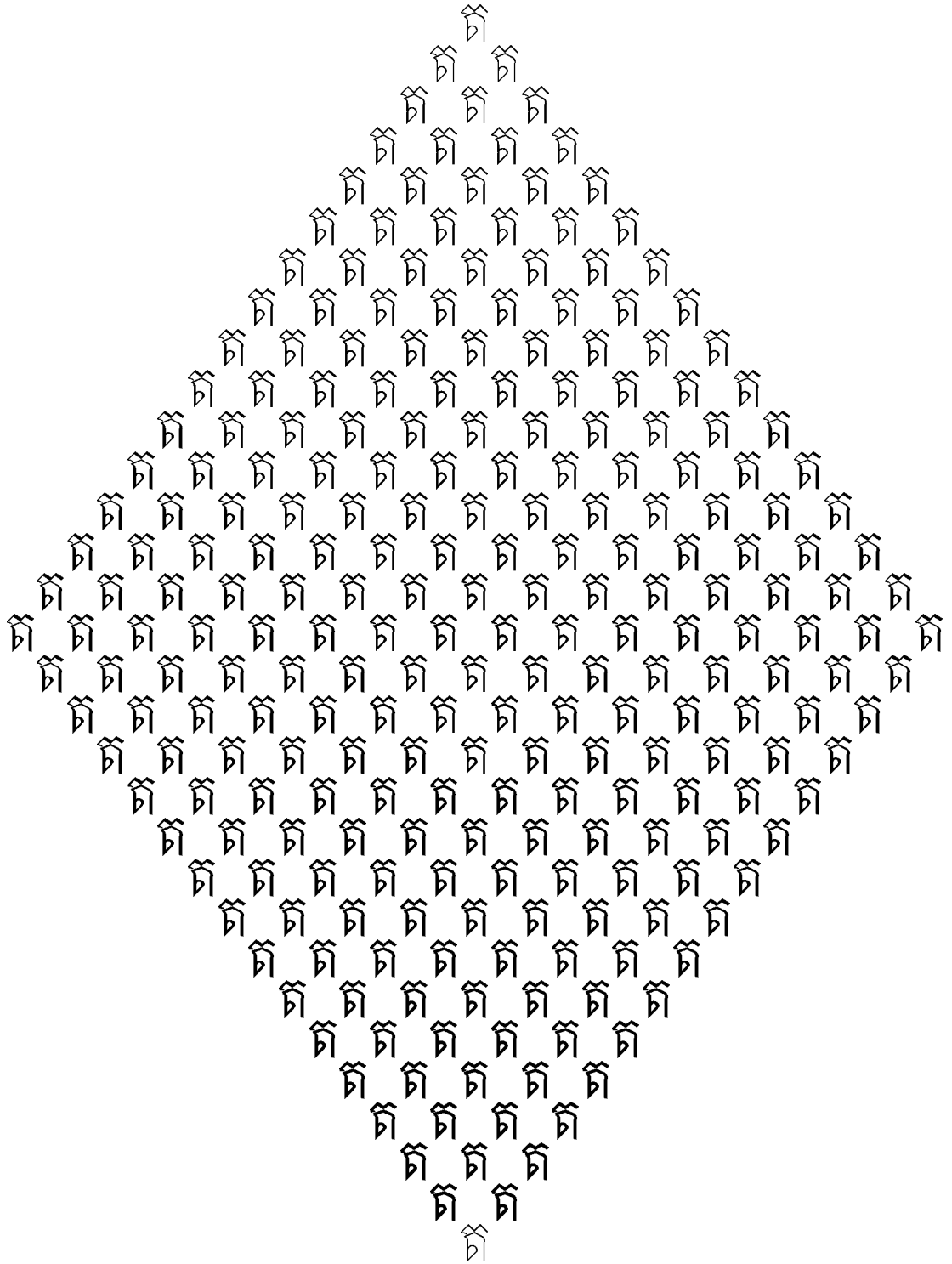


Figure 6. Test table for gray density fine-tuning of Khmer font

### 5.1 Graphical classification of Khmer cluster components

As already mentioned, there is a strong divergence between the phonic and graphical representation of a consonantal cluster: for example, is  $c = \text{𑄣𑄤}$ ,  $s_1 = \text{𑄣𑄤}$ ,  $s_2 = \text{𑄣}$ ,  $v = \text{𑄣𑄤}$ , then for the same cluster  $\text{𑄣𑄤𑄣𑄤}$ , the former representation is  $\langle c \rangle \langle s_1 \rangle \langle s_2 \rangle \langle v \rangle$  and the latter  $\langle v \rangle \langle s_2 \rangle \langle c \rangle \langle s_1 \rangle$  (left branch)  $\langle s_2 \rangle \langle c \rangle \langle s_1 \rangle \langle v \rangle$  (right branch).

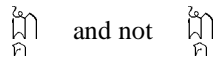
A thorough study of Khmer script and traditional typography, has resulted in the following classification of graphical components of a consonantal cluster:

1. **The ‘left part’.** Four elements which are placed on the left of a consonant:  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ .
2. **The ‘central part’.** All consonants:  $\text{𑄣}$ ,  $\text{𑄤}$  . . .  $\text{𑄭}$ . Also consonant + vowel  $\in \{ \text{𑄣𑄤}, \text{𑄣𑄤} \}$  combinations, whenever the vowel is attached to the consonant and not to a subscript:  $\text{𑄣𑄤}$ ,  $\text{𑄣𑄤}$ ,  $\text{𑄣𑄤}$  etc. but *not*  $\text{𑄣𑄤}$ .  
The difference between ‘left’ and ‘central’ part is that only the latter is taken into account when determining the symmetry axis of the cluster.
3. **The ‘movable’ part.** Subscripts and superscripts which are moved horizontally so that their symmetry axis coincides with the axis of the central part:  $\text{𑄣}$  . . .  $\text{𑄣}$ , and  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ .
4. **The ‘right’ part.** Elements placed on the right of the central part, and not involved in the determination of the cluster symmetry axis. In this category we have certain subscript characters:  $\text{𑄣}$  . . .  $\text{𑄣}$ , as well as selected subscript and superscript vowels and diacritical marks:  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ ,  $\text{𑄣}$ .

The effective graphical construction of a consonantal cluster by  $\text{\TeX}$ , is done in the following way: the preprocessor’s output replaces the phonic representation of a cluster (in the encoding described in Section 3.1) by a  $\text{\TeX}$  macro  $\text{\backslash KHcc1}$  with five arguments: the first is a 9-digit number representing the phonic representation of the cluster (and with the property that if  $N, N'$  are numbers representing clusters  $C, C'$  then  $C \succ C' \iff N > N'$ , where  $\succ$  is the collating order of clusters and  $>$  the usual ordering of integers); the remaining four arguments correspond to the four parts of the graphical decomposition of a cluster as described above. For example,

```
\KHcc1{050311501}{e/r}{gA}{/k}{'}
```

indicates a left part  $e/r$  ( $\text{𑄣𑄤}$ ), a central part  $gA$  ( $\text{𑄣𑄤}$ ), a movable part  $/k$  ( $\text{𑄣}$ ) and a right part  $'$  ( $\text{𑄣}$ ). This example illustrates the important fact that the symmetry axis of the central part is not necessarily the middle axis of the box containing the central part:



The difference is more than of just an aesthetic nature: in some cases the vertical alignment of elements within a cluster is necessary to determine the cluster itself. Take for example characters  $0x89$  ( $\text{𑄣𑄤}$ ) and  $0x96$  ( $\text{𑄣𑄤}$ ). When the latter is followed by a vowel  $\text{𑄣}$  it becomes  $\text{𑄣𑄤}$ , which is indistinguishable from the upper part of the former: it is the lower part  $\text{𑄣}$  that enables differentiation. But when both happen to carry the same subscript consonant, then this lower part vanishes. The difference will be found in the alignment of the



2. To prevent confusion between the letter ឃ followed by vowel ា, and the letter ម, the former combination of consonant and vowel is written មា. A variant of this letter is used in the presence of a subscript: ម + ា = មា.
3. When a cluster with ្រ contains vowel ា or ្រ, then the width of the primary consonant determines the depth of the vowel: ម + ្រ = ម្រ, but ម.ម + ្រ + ្រ = ម្រ.
4. The letter ម is not supposed to carry a subscript consonant; in some rare cases, it carries subscript ្រ: ម្រ.

### 5.3 Operations depending on collating order

As mentioned in the previous section, the  $\TeX$  command `\KHCC1`, obtained by the preprocessor, describes a cluster by means of five arguments. The last four arguments describe the cluster graphically: they correspond to the four parts of the graphical decomposition of a cluster, according to Section 5.1. The first argument corresponds to the phonic decomposition of the cluster; it is a 9-digit number  $N = c_1c_2s_1s_2s_3v_1v_2d_1d_2$  where

1.  $c_1c_2$  determines the primary consonant of the cluster:  $c_1c_2$  goes from 01 = ម, to 33 = ម;
2.  $s_1s_2$  determines the (first) subscript consonant:  $s_1s_2 = 00$  if there is no subscript consonant, otherwise  $s_1s_2$  goes from 01 = ្រ, to 32 = ្រ;
3.  $s_3 = 0$  if there is no second subscript consonant, 1 if there is a second subscript ្រ;
4.  $v_1v_2$  determines the vowel:  $v_1v_2 = 00$  if there is no vowel, otherwise  $v_1v_2$  goes from 01 = ា, to 28 = ្រ;
5.  $d_1d_2$  determines the diacritic mark:  $d_1d_2 = 00$  if there is no diacritic, otherwise  $d_1d_2$  goes from 01 = ្រ, to 08 = ្រ.

A complete list of characters, alphabetically ordered, is given in Section 2. Collating order rules mean that for clusters  $C, C'$  and their corresponding 9-digit numbers  $N, N'$ , we have

$$C \succ C' \iff N > N'$$

where  $\succ$  is the collating order of clusters. The numbers  $N, N'$  can easily be ordered since the collating order of clusters corresponds to their order as integers. This fact enables straightforward searching, sorting, indexing and other operations involving collating order.

## 6 HYPHENATION AND OTHER PREPROCESSOR FEATURES

### 6.1 Hyphenation

Hyphenation in Khmer obeys a very simple rule: words are hyphenated between syllables. Unfortunately this rule can hardly be implemented on a computer since there is no algorithmic way of detecting syllables: a syllable can consist of one *or two* consonantal clusters.

With the help of Prof. Alain Daniel, an empirical hyphenation mechanism has been developed out of several general rules and observations. Below is a first set of rules — there will be further refinement after thorough testing on bigger amounts of Khmer text.

Let  $C, C'$  be consonantal clusters. Hyphenation  $C-C'$  is possible whenever:

1.  $C'$  contains a vowel;
2.  $C$  contains a vowel such as ឃ្នះ, ឃ្នើះ, ឃ្នុះ, ឃ្នោះ, ឃ្នាះ, ឃ្នេះ, ឃ្នែះ, ឃ្នោះ, ឃ្នេះ, ឃ្នង, ឃ្ន្រ, ឃ្នាះ, ឃ្នេះ, ឃ្នាះ, ឃ្នេះ, ឃ្នាះ, ឃ្នេះ, ឃ្នាះ, ឃ្នេះ, or one of the diacritical marks ្ន, ្ន;

Hyphenation is always possible before or after special characters.

TeX provides an internal hyphenation mechanism based on *hyphenation patterns*. Unfortunately this mechanism cannot be used in the case of Khmer consonantal clusters, since these are enclosed in boxes and hence cannot be considered as characters by TeX. For this reason, the hyphenation algorithm is performed by the preprocessor; whenever one of the two above rules is satisfied, the TeX macro `\-` is included in the output. This command expands as

```
\def\-\{\discretionary{-}\{\}\}
```

so that a hyphen is obtained whenever a word is hyphenated. There is no algorithm yet for automatic decomposition of sentences into words: the user is asked to include WBK (word-break) codes between words inside a sentence. These codes are converted into `\KHwbk` commands by the preprocessor; `\KHwbk` expands into

```
\def\KHwbk{\discretionary}{\{\}\{\}\}
```

that is: a potential hyphenation point, *without* hyphen.

### 6.2 Decomposition of special characters and spelling reforms

The special characters (codes 0xC1-0xC6) are mostly historical residues and loans from other languages (Pali and Sanskrit). There have been many attempts by the Cambodian Ministry of Education to restrict their number, and eventually replace some of them by regular consonantal clusters.

This replacement can vary from word to word. Prof. Alain Daniel has established a list of reformed words and their replacements. This list is known to the preprocessor, which will output every special character as a TeX macro with a numeric argument indicating the potential replacement by some other special character or by a consonantal cluster. For example, depending on the surrounding word, ឃ្ន is output as `\KHao0`, `\KHao1`, `\KHao2`, `\KHao3` or `\KHao4`. If a certain boolean variable `\ifreformed` is false, then all five macros will always expand into ឃ្ន. On the other hand, if the boolean is true, then the first macro will expand into ឃ្ន, the second into អ្ន, the third into ឃ្ន, the fourth into អ្ន and the fifth into អ្ន.

Below is a first list of reformed words, known to the preprocessor. The special characters and their decompositions are set in bolder type.

<b>រំតិល</b> → រំអិល	<b>ឃ្ន្រស</b> → តិស	<b>ឃ្ន្រសាន</b> → តិសាន	<b>ឃ្ន្រស្រ</b> → តិស្រ
<b>រំតិល</b> → រំអិល	<b>ឃ្ន្រសឆ</b> → តិសឆ	<b>ឃ្ន្រស្រ</b> → តិស្រ	<b>ឃ្ន្រ</b> → អ្ន
<b>តិស្រការ</b> → ឃ្ន្រការ			



ក្រខ្ចី → ក្រអៅ	រព្យក → រលឹក	ព្យក្យ- → អៃក្យ-	ខ្ចង់ឡោង → អោងឡោង
ក្រាងក្រខ្ចី → ក្រាងក្រអៅ	ព្យ → លឹ	ព្យរវត → អៃរវត	ខ្ចង្កា → អង្កា
ខ្ចកា → ខ្ចកា	ព្យជ័យ → លឹជ័យ	ព្យស្វរ → អៃស្វរ	ខ្ចប្លស័ន → ខ្ចប្លស័ន
ខ្ចជី → អ្នជី	ព្យជិ → លឹជិ	ព្យស្វរ្យ → អៃស្វរ្យ	ខ្ចត្ត → អោត
ខ្ចន → អ្នន	ព្យលាស់ → លឹលាស់	ប្រខ្ចង់ → ប្រអោង	ខ្ចទ្យាន → ខ្ចទ្យាន
ខ្ចម! → ខ្ចម	ព្យសាយ → លឹសាយ	រខ្ចក → រអោក	ខ្ចបច្ចាតិក → ខ្ចបច្ចាតិក
ខ្ចរ → ខ្ចរ	ព្យរវណ → អៃរវណ	លខ្ចន → លអោន	ខ្ចវៃ → ខ្ចវៃ
ខ្ចរ្យ → ខ្ចរ្យ	ព្យ! → អៃ!	សខ្ចក → សអោក	ខ្ចសហ → ខ្ចសហ
ប្រ → រ	ព្យក- → អៃក-	ខ្ចវៃ → ខ្ចវៃ	ខ្ច → អៅ!
សប្ប្ច → សប្ប្ច		ខ្ច! → អៅ	ខ្ចក → ខ្ចក
រព្យក → រលឹក			

### 7 SHORTCOMINGS AND PLANS FOR FURTHER DEVELOPMENT

The system presented in this paper enables high-quality Khmer typesetting. It is the first Khmer typesetting system which solves problems such as text input in phonic order, positioning of subscripts and superscripts, optical scaling, hyphenation and replacement of special characters.

Nevertheless the graphical cluster-construction algorithm described in this paper has certain flaws; a few examples:

- if a consonant with subscript consonant carries the  $\square$  vowel, then the latter should be justified at the right edge of the *subscript*, which is not necessarily aligned with the right edge of the consonant. For example, in the (hypothetical) cluster  $\text{ក្រ}_{\text{ក}}$ , the  $\square$  is badly positioned;
- take a narrow letter (like រ, វ) which carries a large subscript (like  $\text{ក}_{\text{ក}}$  or  $\text{ន}_{\text{ន}}$ ) and suppose you are at the line boundary (either left or right); then contrarily to the normal use of subscripts, it is the subscript which should be used for line justification, and not the consonant.

These problems cannot be solved using the current mechanism (in which  $\text{T}_\text{E}_\text{X}$  considers that all subscripts and superscripts are of zero width). It could be possible to use subscripts with non-zero width, but (a) this would slow the process down, (b) it wouldn't solve the problem of the line boundary, since we are asking for contradicting properties: inside a sentence subscripts should not interfere in determining the distance between clusters, while at the line's boundary they should.<sup>4</sup> Furthermore, one could imagine a sentence ending with  $\text{ក}_{\text{ក}}$  and the next sentence starting with  $\text{វ}_{\text{វ}}$ . The blank space in between is hardly sufficient to prevent clusters from overlapping. Visually, the beginning of the sentence is lost.

Corrections to these problems can be performed manually (because these problems occur very rarely). However, a much more natural and global solution would be to treat consonantal clusters as individual codes in a 16-bit encoding scheme. As mentioned in the introduction, only 2 821 clusters (out of 535 000 theoretical possibilities) have been detected in the fairly complete dictionary of Prof. Alain Daniel, so a 16-bit table would be more than sufficient to cover them.

<sup>4</sup> Unfortunately, in  $\text{T}_\text{E}_\text{X}$  there is no such thing as a `\everyline` command.

This method of Khmer typesetting, is part of the  $\Omega$  project, undertaken by John Plaice (Université Laval, Canada) and the author. The first realization of  $\Omega$  is an extension of  $\text{\TeX}$  (and the two utilities  $\text{VPtoVF}$  and  $\text{DVICopy}$ ) to 16-bit fonts (allowing the use of 65 536 characters and 4 294 967 296 ligatures or kerning pairs). These fonts will be exclusively virtual: since the DVI file format allows up to 32-bit fonts there is no need to extend its specifications; DVI-files with 16-bit  $\Omega$  fonts will be ‘devirtualized’ through  $\text{DVICopy}$ : the 16-bit virtual fonts will be replaced by their 8-bit base fonts. In this way  $\Omega$  DVI files will be converted to standard 8-bit DVI files; no special DVI drivers will be needed (not even virtual font compatible ones). In the case of Khmer, the (unique) base font will contain the glyph descriptions (in  $\text{PK}$  or PostScript format) and the virtual font will contain the definitions of consonantal clusters. Since consonantal clusters will be treated by  $\text{\TeX}$  as individual characters, one will be able to define kerning pairs between them and solve the main problem of Khmer typesetting.

Text input could still be done using the 8-bit encoding of [Section 3](#); internal ligaturing will map the 8-bit description of consonantal clusters into their codes in the 16-bit table (a preprocessor can still be used to perform explicit construction, if for any reason they are not included in the table). This approach is similar to Kanji construction out of Kana characters in Japanese, or to Hangoul construction out of elementary strokes in Korean.

Other projects using  $\Omega$  concern vowelized Arabic, typesetting in Indic languages, Thai, Amharic without preprocessor, use of calligraphic fonts (such as Adobe’s *Poetica*), redrawing of Garamont’s *Greco du Roy* etc. First releases of  $\Omega$  projects are expected to take place in autumn 1994.

## AVAILABILITY

The METAFONT,  $\text{\TeX}$  and C sources of all software presented in this paper belong to the *public domain*. They constitute a proposal for a Khmer *\TeX* Language Package, submitted to the Technical Working Group on Multiple Language Coordination of the  $\text{\TeX}$  Users Group and will be released after ratification. The  $\alpha$  version of the package is currently being tested in Cambodia, and can be obtained from the author.

Khmer keyboard layouts using phonic input of consonantal clusters are currently being tested as well.

## ACKNOWLEDGEMENTS

The author would like to thank Prof. Alain Daniel (Institute of Oriental Languages and Civilizations, Paris) for his continuous support and encouragement and the Imprimerie Louis-Jean (Gap) in the person of Maurice Laugier, for having financed this project.

## REFERENCES

1. Daniel Berry and Johny Srouji, ‘Arabic formatting with  $\text{ditroff/ffortid}$ ’, *Electronic Publishing—Origination, Dissemination and Design*, 5(4), 163–208, (1992).
2. Yannis Haralambous, ‘Typesetting the holy Quran with  $\text{\TeX}$ ’, in *Proceedings of the 2nd International Conference on Multilingual Computing—Arabic and Latin script (Durham)*, (1992).
3. Alain Daniel, *Dictionnaire pratique cambodgien-français*, Institut de l’Asie du Sud-Est, Paris, 1985.
4. Alain Daniel, *Lire et écrire le cambodgien*, Institut de l’Asie du Sud-Est, Paris, 1992.

- 
5. Derek Tonkin, *The Cambodian Alphabet*, Transvin Publications, Bangkok, 1991.
  6. Akira Nakanishi, *Writing Systems of the World*, Charles E. Tuttle Company, Tokyo, 1980.
  7. វិចិត្រកម្មខ្មែរ, ការផ្តោយរបស់ពុទ្ធសាសនបណ្ឌិតយ, ព. ស. ២៥០៥ [*Dictionnaire Cambodgien*], Éditions de l'Institut Bouddhique.
  8. Donald E. Knuth, *The METAFONTbook*, Computers & Typesetting, Addison-Wesley, 1986.
  9. Donald E. Knuth, *The T<sub>E</sub>Xbook*, Computers & Typesetting, Addison-Wesley, 1986.
  10. Helmut Kopka, *L<sup>A</sup>T<sub>E</sub>X, eine Einführung*, Addison-Wesley, 1991.
  11. Helmut Kopka, *L<sup>A</sup>T<sub>E</sub>X, Erweiterungsmöglichkeiten*, 2nd edn, Addison-Wesley, 1991.