

# First applications of Ω: Greek, Arabic, Khmer, Poetica, ISO 10646/UNICODE, etc.

Yannis Haralambous

Centre d'Études et de Recherche sur le Traitement Automatique des Langues  
Institut National des Langues et Civilisations Orientales, Paris.  
Private address: 187, rue Nationale, 59800 Lille, France.  
Yannis.Haralambous@univ-lille1.fr

John Plaice

Département d'informatique, Université Laval, Ste-Foy, Québec, Canada G1K 7P4  
plaice@ift.ulaval.ca

## Abstract

In this paper a few applications of the current implementation of Ω are given. These applications concern typesetting problems that cannot be solved by T<sub>E</sub>X (consequently, by no other typesetting system known to the authors). They cover a wide range, going from calligraphic script fonts (Adobe Poetica), to plain Dutch/Portuguese/Turkish typesetting, to vowelized Arabic, fully diacriticized scholarly Greek, or decently kerned Khmer.

For every chosen example, the particular difficulties, either typographical or T<sub>E</sub>Xnical (or both), are explained, and a short glance to the methods used by Ω to solve the problem is given. A few problems Ω *cannot* solve are mentioned, as challenges for future Ω versions.

## On Greek, ancient and modern (but rather ancient than modern)

**Diacritics against kerning.** It is in general expected of educated men and women to know Greek letters. Already in college, having used  $\theta$  for angles,  $\gamma$  for acceleration, and  $\pi$  to calculate the area of a round apple pie of given radius, we are all familiar with these letters, just as with the Latin alphabet. But the Greek language, in particular the ancient one, needs more than just letters to be written. Two kinds of diacritics are used, namely accents (acute, grave and circumflex), and breathings (smooth and raw) which are placed on vowels; breathings are also placed on the consonant rho.

Every word has at most one accent<sup>1</sup>, and 99.9% of Greek words have exactly *one* accent. Every word starting with a vowel has exactly one breathing<sup>2</sup>. It follows that writing in Greek involves much more accentuation than any Latin-alphabet language, with the obvious exception of Vietnamese.

How does T<sub>E</sub>X deal with Greek diacritics? If the traditional approach of the `\accent` primitive had been taken, then we would have practically *no* hy-

phenation (which in turn would result in disastrous over/underfulls, since Greek can easily have long words like  $\acute{\omega}\tau\omicron\rho\nu\lambda\alpha\rho\upsilon\gamma\gamma\omicron\lambda\omicron\gamma\iota\kappa\acute{\omicron}\varsigma$ ), *no* kerning, and a cumbersome input, involving one or two macros for every word.

The first approach, originated by Silvio Levy (1988), was to use T<sub>E</sub>X's ligatures ('dumb' ones first, 'smart' ones later on) to obtain accented letters out of combinations of codes representing breathings (> and <), accents ( ' , ' and ~ or =) and the letters themselves. In this way, one writes >'h to get ḥ. This approach solved the problem of hyphenation and of cumbersome input.

Nevertheless, this approach fails to solve the *kerning* problem. Let's take the very common case of the article τὸ (letter tau followed by the letter omicron); in almost all fonts there is a kerning instruction between these letters, obviously because of invariant characteristics of their shapes. Suppose now that omicron is accented, and that one writes τ'ο to get tau followed by omicron with grave accent. What T<sub>E</sub>X sees is a 't' followed by a grave accent. No kerning can be defined for these two characters, because we have no idea what may follow after the grave accent (it can very well be a iota, and usually there is no kerning between tau and iota). When the letter omicron arrives, it is too late; T<sub>E</sub>X has already forgotten that there was a tau before the grave accent.

A solution to this problem would be to write diacritics *after* vowels ("post-positive notation"). But

<sup>1</sup> Sometimes an accent is transported from a word to the preceding one:  $\acute{\alpha}\nu\theta\rho\omega\pi\acute{\omicron}\varsigma$  instead of  $\acute{\alpha}\nu\theta\rho\omega\pi\omicron\varsigma$ ,  $\tau\acute{\iota}\varsigma$ , so that typographically a word has more than one accent.

<sup>2</sup> With one exception: the letters ρρ are often written ῥῥ, when inside a word:  $\pi\omicron\acute{\rho}\acute{\rho}\acute{\omega}$ .

this contradicts the visual characteristics of diacritics in the upper case, since these are placed on the left of uppercase letters: "Έαρ could hardly be transliterated E>'ar. And, after all, T<sub>E</sub>X should be able to do proper Greek typesetting, however the letters and diacritics are input.

Ω solves this problem by using an appropriate chain of Ω Translation Processes (ΩTPs), a notion explained in Plaice (1994): as example, consider the word έαρ:

1. Suppose the user wishes to input his/her text in 7-bit ASCII; he/she will type >'ear, and this is already ISO-646, so that no particular input translation is needed. Another choice would be to use some Greek input encoding, such as ISO-8859-7 or EΛOT; then he/she may as well type >'εαρ or >έαρ (the reason for the absurd complication of being forced to type either >'ε or >έ to obtain έ, is that "modern Greek" encodings have taken the easy way out and feature only one accent, as if the Greek language was born in 1982, year of the hasty and politically motivated spelling reform). The first ΩTP will send these codes to the appropriate 16-bit codes in ISO 10646/UNICODE: 0x1f14 for έ, 0x03b1 for α, and 0x03c1 for ρ.
2. Once Ω knows what characters it is dealing with (Ω's default internal encoding is precisely ISO-10646), it will hyphenate using 16-bit patterns.
3. Finally, an appropriate ΩTP will send Greek ISO 10646/UNICODE codes to a 16-bit virtual font (see next section to find out why we need 16 bits), built up from one or more 8-bit fonts. This font contains kerning instructions, applied in a straightforward manner, since we are now dealing with only three codes: <έ>, <α> and <ρ>. No auxiliary codes interfere anymore.
4. x<sub>d</sub>vi<sub>copy</sub><sup>3</sup> will de-virtualize the dvi file and return a new dvi file using only 8-bit fonts, compatible with every decent dvi driver.

By separating tasks, hyphenation becomes more natural (for T<sub>E</sub>X one has to use patterns including auxiliary codes ', ', = etc.). Furthermore, an additional problem has been solved *en passant*: the primitives \leftthyphenmin and \rightthyphenmin apply to characters of \catcode 12. To obtain hyphenation between clusters involving auxiliary codes, we have to declare these codes as 'letter-like characters'. For example, the word έαρ, written as >'ear: the codes > and ' must be considered as letter-like (non-trivial \lccode) to allow hyphenation; but this means that for T<sub>E</sub>X, έαρ has 5 letters instead of 3, and hence even

<sup>3</sup> In the name of this program, which is an extended version of Peter Breitenlohner's dvi<sub>copy</sub>, 'x' stands for 'extended', not for 'X-Window'.

if we ask \leftthyphenmin=3, we will still get the word hyphenated as έ-αρ!! Ω solves this problem by hyphenating *after* the translation has been done (in this case >'e or >'ε or >έ → έ).

**Dactyls, spondees and 16-bit fonts.** Scholarly editions of Greek texts are slightly more complicated than plain ones<sup>4</sup>, one of the add-ons being a *third* level of diacritization: syllable lengths.

One reads in Betts and Henry (1989, p. 254): "Greek poetry was composed on an entirely different principle from that employed in English. It was not constructed by arranging stressed syllables in patterns, nor with a system of rhymes. Greek poets employed a number of different metres, all of which consist of a certain fixed arrangement of **long and short syllables**". Long and short syllables are denoted by the diacritics *macron* and *breve*. The diacritics are placed *between* the letter and the regular diacritic, if any (except in the case of uppercase letters, where they are placed over the letter while regular diacritics are placed to its left).<sup>5</sup>

The famous first two verses of the Odyssey "Ἄνδρα μοι ἔννεπε, Μοῦσα πολύτροπον, ὃς μάλα πολλὰ πλάγχθη ἐπεὶ Τροίης ἱερὸν πτόλιεθρον ἔπερσε form *hexameters*. These consists of six feet: four dactyls or spondees, a dactyl and a spondee or trochee (see again Betts and Henry 1989 for more details). One could write the text without accents or breathings to make the metre more apparent:

Ἄνδρα μοι ἔννεπε, Μοῦσα πολύτροπον, ὃς μάλα πολλὰ πλάγχθη ἐπεὶ Τροίης ἱερὸν πτόλιεθρον ἔπερσε  
 or one could decide to typeset all types of diacritics:  
 "Ἄνδρα μοι ἔννεπε, Μοῦσα πολύτροπον, ὃς μάλα πολλὰ πλάγχθη ἐπεὶ Τροίης ἱερὸν πτόλιεθρον ἔπερσε

Having paid a fortune to acquire the machine that sits between the keyboard and the screen, one could expect that hyphenation of text and kerning between letters remains the same, despite the constantly growing number of diacritics. Actually this isn't possible for T<sub>E</sub>X: there are exactly 345 combinations of Greek letters with accents, breathings, syllable lengths and subscript iota; T<sub>E</sub>X can handle at most 256 characters in a font. Therefore Ω is necessary for hyphenation of Greek text, whenever syllable lengths are typeset.

In this case, things do not run as smoothly as in the previous section: although 345 is a small number compared to 65,536 (= 2<sup>16</sup>), ISO-people decided that

<sup>4</sup> After all, scholars have been studying Greek text for more than 2,000 years now.

<sup>5</sup> These additional diacritics are also used for a different purpose: in prose, placed on letters alpha, iota or upsilon, they indicate if they are pronounced long or short (this time we are talking about *letters* and not about *syllables*).

there is not room enough for all combinations of accents, breathings and syllable lengths.<sup>6</sup>

Whenever ISO 10646/UNICODE comes too short for our needs, we use the *private zone*. Just as in the TV serial *The Twilight Zone*, in the private zone *anything* can happen. In the case of Ω all operations remain internal, so that we have absolute freedom of defining characters: according to the 1993-05-01 version of ISO-10646-1, the private zone consists of characters 0xe000 to 0xffffd (of group 0, that is the 16-bit part of ISO-10646-1), a total of 8,190 positions.

Ω will treat the input like in the previous section, but letters with macron and breve diacritics will occupy internal positions in the private zone. The rest of the treatment will be exactly the same. As for the transliterated input one could take ^ and \_ to denote macron and breve (after having changed their catcodes so that they do not interfere with math operators), or any other combination of 7-bit or 8-bit codes.

**A dream that may come true.** As stated by the first author in his Cork talk, back in 1990, his dream was—and still is—to draw a Greek font based on the famous « Grecs du Roi » typeface by Claude Garamont, graved in 1544–46 for the king François I. This typeface was designed after a manuscript of Ἄγγελος Βεργήχιος, a Cretan, calligrapher and reader of Greek at the French Court, in the beginning of the XVI century. There are 1327 different types, most of them ligatures of two or more letters (sometimes entire words). One can read in Nationale (1990) that “*this typeface is the most precious piece of the collection [of the French National Printing House]*”, certainly not the least of honors! Ω is the ideal platform for typesetting with this font, since it would need only an additional ΩTP to convert plain Greek input into ligatures. The author hopes to have finished (or, on a more realistic basis, to have brought to a decent level) this project in time for the International Symposium “Greek Letters, From Linear B to Gutenberg and from G to Ψ”, which will take place in Athens in late Spring 1995, organized by the Greek Font Society<sup>7</sup>.

<sup>6</sup> Nevertheless, they included lower and upper alpha, iota and upsilon with macron and breve, probably for the reasons explained in the previous footnote. However, combining diacritics must be used to code letters with macron/breve *and* additional diacritics.

<sup>7</sup> For additional information on the Symposium, contact Ἐταιρεία Ἑλληνικῶν Τυπογραφικῶν Στοιχείων, Ἑλληνικῶν 39-41, 116 35 Ἀθήνα, Greece, or Michael S. Macrakis, 24 Fieldmont Road, Belmont, MA 02178-2607, USA.

## Arabic, or “The Art of separating tasks”

**Plain Arabic, quick and clean and elegant.** Arabic typesetting is a beautiful compromise between Western typesetting techniques (finite number of types, repeated *ad infinitum*) and Arabic calligraphy (infinite number of arbitrarily complex ligatures). We can subdivide Arabic ligatures into two categories: (a) mandatory ones: letter connections (ه + م → هم) and the special ligature lām-alif (ل + ا → لا), and (b) optional ones, used for esthetic reasons.

The second category of ligatures corresponds to our good old ‘fi’ ‘fl’ ,etc. They depend on the font design and on the degree of artistic quality of a document. The first author has made a thorough classification of esthetic ligatures of the Egyptian typecase (see Haralambous 1992, reprinted in Haralambous 1993c). Here is an example of the ligaturing process of the word تحمل, following Egyptian typographical traditions:

- ت ح م ل (letters not connected);
- تحمل (only mandatory ligatures, connecting letters);
- تحمل (esthetic ligature between the first two letters);
- تحمل (esthetic ligature between the first three letters).

To produce more than 1,500 possible ligatures of two, three or four letters, three 256-character tables were necessary. Each ligature is constructed by superposition of small pieces. Once T<sub>E</sub>X knows which characters to take, and from which font, it only needs to superpose them (no moving around is necessary). The problem is to recognize the existence of a ligature and to find out which characters are needed. This process is highly font-dependent. A different font—for example in Kuffic or Nastaliq style—may have an entirely different set of ligatures, or none at all (like the plain font, in which the two words تيخ العربي are written, that is widely used in computer typesetting because of its readability); nevertheless, the mandatory ligatures remain strictly the same, whatever font is used.

Up to now, there are three solutions to the problem of mandatory Arabic ligatures:

- First, by K. Lagally (1992), is to use T<sub>E</sub>X macros for detecting and applying mandatory ligatures (we say: “to do the contextual analysis”). This process is cumbersome and long. It is highly dependent on the font encoding and the macros used can interfere with other T<sub>E</sub>X macro packages. All in all, it is not the natural way to treat a phenomenon which is a low-level fundamental characteristic of the Arabic script.

- Second, by the first author (Haralambous 1993b), is to use “smart”  $\TeX$  font ligatures (together with  $\TeX$ - $X\TeX$ , the bi-directional version of  $\TeX$ ); this process is philosophically more natural, since contextual analysis is done behind the scenes, on the very lowest level, namely the one of the font itself. It does not depend on the font encoding, since every font may use its own set of ligatures. The disadvantage lies in the number of ligatures needed to accomplish the task: about 7,000! The situation becomes tragic when one wants to use a dozen Arabic fonts on the same page:  $\TeX$  will load 7,000 probably strictly identical ligatures for each font. You need more than Big $\TeX$  to do that.
- Third, also by the first author (Haralambous 1993b), is to use a preprocessor. The advantage is that the contextual analysis is done by a utility dedicated to this task, with all possible bells and whistles (for example adding variable length connecting strokes, also known as ‘keshideh’); it is quick and uses only a very small amount of memory. Unfortunately there are the classical disadvantages of having a preprocessor treat a document before  $\TeX$ : one file may  $\backslash$ input another file, from any location of your net, and there is no way to know in advance which files will be read, and hence have to be preprocessed; preprocessor directives can interfere with  $\TeX$  macros; there is no nesting between them and this can easily produce errors with respect to  $\TeX$  grouping operations, etc.

None of these methods can be applied for large scale real-life Arabic production: in all cases, the Arabic script is treated as a ‘puzzle to solve’ and, inevitably,  $\TeX$ ’s performance suffers.

We use  $\Omega$ TPs to give a natural solution to the problem; consider once again the example of the word  $\text{تَحْمَل}$ :

- First,  $\text{تَحْمَل}$  is read by  $\Omega$ , either in Latin transliteration (tHm1) or in ISO-8859-6, or ASMO, or Macintosh Arabic, or any decent Arabic script input encoding.
- The first  $\Omega$ TP converts this input into ISO 10646/UNICODE codes for Arabic letters: 0x062a (ت), 0x062d (ح), 0x0645 (م), 0x0644 (ل);
- ISO 10646/UNICODE being a *logical* way of codifying Arabic letters, and not a graphical one, there is no information on their contextual forms (isolated, initial, medial, final). The second  $\Omega$ TP sends these codes to the private zone, where we have (internally) reserved positions for the combinations of Arabic characters and contextual forms. Once this is done,  $\Omega$  knows the form of each character.

- The third  $\Omega$ TP simply translates these codes to a 16-bit standard Arabic  $\TeX$  font encoding (this is a minor operation: the private zone being located at the end of the 16-bit table, we move the whole block near to the beginning of the table).
- If the font has no esthetic ligatures, we are done:  $\Omega$  will send the results of the last  $\Omega$ TP to the DVI file, and produce  $\text{تَحْمَل}$ . On the other hand, if there are still esthetic ligatures—as in  $\text{تَحْمَل}$ —then these will be included in the font, as ‘smart’ ligatures. Since the font table can have as many as 65,536 characters, there is plenty of room for small character parts to be combined.<sup>8</sup>

What we have achieved is that the fundamental process of contextual analysis is done by background machinery (just like  $\TeX$  hyphenates and breaks paragraphs into lines), and that the optional esthetic refinements are handled exclusively by the font (in analogy to Roman fonts having more ligatures than typewriter ones, etc.).

**Vowelized Arabic (things get harder).** In plain contemporary Arabic, only consonants and long vowels are written; short vowels have to be guessed by the reader, using the context (the same consonants with different short vowels, can be understood as a verb, a noun, an adjective etc.). When it is essential to specify short vowels, small diacritics are added over or under the letters. Besides short vowels, there are also diacritics for doubling consonants, for indicating the absence of vowel, and for the glottal stop (like in “Oh-oh”): counting all possible combinations, we obtain 14 signs. These diacritics can give  $\TeX$  a hard time, since they have to be coded between consonants, and hence interfere in the contextual analysis algorithm: for example, suppose that  $\TeX$  is about to typeset letter  $x$ , which is the last letter of a word, followed by a period. Having read the period,  $\TeX$  knows that the letter has to be of final form (one of the 7,000 ligatures has to be  $\langle x$  in medial form  $\rangle + \langle . \rangle \rightarrow \langle x$  in final form  $\rangle \langle . \rangle$ ). Now suppose that the letter is immediately followed by a short vowel, which in our case is necessarily placed between the letter  $x$  and the period.  $\TeX$ ’s smart ligatures cannot go two positions backwards; when  $\TeX$  discovers the period after the short vowel it is too late to convert the medial  $x$  into a final one.

Fortunately,  $\Omega$ TPs are clever enough to be able to calculate letter forms, whatever the diacritics surrounding them (which is exactly the attitude of the

<sup>8</sup> If these esthetic ligatures are used in several fonts, it might be possible that we have the same problem of overloading  $\Omega$ ’s font memory; in this case, we can always write a fourth  $\Omega$ TP, which would systematically make the ‘esthetic analysis’ out of the contextually analyzed codes.



All one needs to change is a macro at the beginning. Accomplishing this feature for Latin and Tifinagh was more-or-less straightforward; not so for Arabic. Unfortunately, that font has all the problems of plain Arabic fonts: it needs more than 7,500 ligatures to do the contextual analysis, and is overloaded: there is no longer room to add a single character, an annoyance for a language that is still under the process of standardization.

Another source of difficulties is the fact that the equivalences between Latin, Tifinagh and Arabic are not immediate. Some short vowels are written in the Latin text, but not in the Arabic and Tifinagh ones. Moreover, double consonants are written explicitly in Latin and Tifinagh, while they are written as a single consonant with a special diacritic in Arabic. And perhaps the most difficult problem is to make every Berber writer feel “at home”, regardless of the script he/she uses: one should not have the impression that one script is privileged over the others!

Finally, the last problem (not a minor one when it comes to real-life production) is that we need a special Arabic font for Berber, because of the different input transliteration: for example, while in plain Arabic transliteration we use ‘v’ for ف and ‘sh’ for ش, in Berber we are forced to use ‘g’ for the former and ‘c’ for the latter. There are two supplementary letters used for Berber in Arabic script: ج and ز; these letters are also used in Sindhi and Pashto, so that the glyphs are already covered by the general Arabic  $\TeX$  system; but in Berber, they have to be transliterated as ‘j’ and ‘z’, because of the equivalences with the Latin alphabet. This forces us to use a different transliteration scheme than the one for plain Arabic, and hence—because of  $\TeX$ ’s inability to clearly separate input and output encoding—to use a differently encoded  $\TeX$  output font. Imagine you are typesetting a book in both Berber and Arabic; you will need two graphically identical fonts for every style, point size, weight and font family, each one with more than 7,000 ligatures. And we are just talking of (esthetically) non-ligatured fonts!

$\Omega$  solves this problem by using the same output fonts for Berber and plain Arabic. We just need to replace the first  $\Omega$ TP of the translation chain: the one that converts raw input into ISO 10646/UNICODE codes. Berber linguists can feel free to invent/introduce new characters or diacritics; as long as they are included in the ISO 10646/UNICODE table we will simply have to make a slight change to the first  $\Omega$ TP (and if these signs are not yet in ISO 10646/UNICODE we will use the private zone).

**Comorian: African Latin versus Arabic.** A similar situation has occurred in the small islands of the Comores, between Madagascar and the African continent. Both the Latin alphabet (with a few additions taken from African languages) and the Arabic one are

used. Because of the many sounds that must be distinguished, one has to use diacritics together with Arabic letters. These diacritics look like Arabic diacritics (for practical reasons) but are not used in the same way; in fact, they are part of the letters, just like the dots are part of plain Arabic letters.

Once again, the situation can easily be handled by an  $\Omega$ TP. While it is still not clear what should be proposed for insertion into ISO 10646/UNICODE (this proposal, made by Ahmed-Chamanga, a member of the Institute of Oriental Languages in Paris, is now circulating from Ministries to Educational and Religious Institutions, and is being tested on natives of all educational levels), Comorians can already use  $\Omega$  for typesetting, and upgrade the transliteration scheme on the fly.

### Khmer

As pointed out in Haralambous (1993a), the Khmer script uses clusters of consonants, subscript consonants, vowels and diacritics. Inside a cluster, these parts have to be moved around by  $\TeX$  to be positioned correctly. It follows that  $\TeX$  *must* use  $\backslash$ kern instructions between individual parts of a cluster. Because of these, there is no kerning anymore: suppose that characters ๑ and ๓ need to be kerned; and suppose that the consonant ๑ is (logically) followed by the subscript consonant ๓, which is (graphically) placed under this letter: ๑๓. For  $\TeX$ , ๑ is not immediately followed by ๓ anymore, and so no kerning between these letters will be applied; nevertheless, graphically they still *are* adjacent, and hence need eventually to be kerned.

$\Omega$  uses the sledgehammer solution to solve this problem: we define a ‘big’ (virtual) Khmer font, containing *all currently known* clusters. As already mentioned in Haralambous (1993a), approximately 4,000 codes would be sufficient for this purpose. And of course one could always use the traditional  $\TeX$  methods to form exceptional clusters, not contained in that font.

As in Arabic, we deal with Khmer’s complexity by separating tasks. A first  $\Omega$ TP will send the input method the use has chosen, to ISO 10646/UNICODE Khmer codes (actually there aren’t any ISO 10646/UNICODE Khmer codes yet, but the first author has submitted a Khmer encoding proposal to the appropriate ISO committees, and expects that soon some step will be taken in that direction—for the moment we will use once again the private zone). A second  $\Omega$ TP will analyze these codes contextually and will send groups of them to the appropriate cluster codes. The separation of tasks is essential for allowing multiple input methods, without redefining each time the contextual analysis—which is after all a basic characteristic of the Khmer writing system.  $\Omega$  will

send the result of the second  $\Omega$ TP to the dvi file, using kerning information stored in the (virtual) font. Finally, `xdvi copy` will de-virtualize the dvi file and create a new file using exclusively characters from the original 8-bit Khmer font, described in Haralambous (1993a).

### Adobe Poetica

Poetica is a chancery script font, designed by Robert Slimbach and released by Adobe Systems Inc. Stating Adobe's promotional material, "*The Poetica typeface is modeled on chancery handwriting scripts developed during the Italian Renaissance. Elegant and straightforward, chancery writing is recognizable as the basis for italic typefaces and as the foundation of modern calligraphy. Robert Slimbach has captured the vitality and grace of this writing style in Poetica. Characteristic of the chancery hand is the common use of flourished letterforms, ligatures and variant characters to embellish an otherwise formal script. To capture the variety of form and richness of this hand, Slimbach has created alternate alphabets and character sets in his virtuoso Poetica design, which includes a diverse collection of these letterforms.*"

Technically, the Poetica package consists of 21 PostScript fonts: Chancery I-IV, Expert, Small Caps, Small Caps Alternate, Lowercase Alternates I-II, Lowercase Beginnings I-II, Lowercase Endings I-II, Ligatures, Swash Caps I-IV, Initial Swash Caps, Ampersands, Ornaments. The ones of particular interest for us are Alternate, Beginnings, Endings and Ligatures, since characters from these can be chosen automatically by  $\Omega$ . The user just types plain text, possibly using a symbol to indicate degrees of alternation. An  $\Omega$ TP converts the input into characters of a (virtual) 16-bit font, including the characters of all Poetica components. Using several  $\Omega$ TPs and changing them on the fly will allow the user to choose the number of ligatures he/she will obtain in the output. It will also allow us to go farther than Adobe, and define kerning pairs between characters from different Poetica components.

See Fig. 1 for a sample of text typeset in Poetica.

### ISO 10646/UNICODE and beyond

It is certainly not a trivial task to fit together characters from different scripts and to obtain an optically homogeneous result. Often the esthetics inherent to different writing systems do not allow sufficient manipulation to make them 'look alike'; it is not even trivial if this should be tried in the first place: suppose you take Hebrew and Armenian, and modify the letter shapes until they resemble one another sufficiently, to our Western eyes. It is not clear if Armenian would still look like Armenian, or Hebrew like Hebrew; furthermore one should not neglect the

psychological effects of switching between scripts (and all other changes the switch of scripts implies: language, culture, state of mind, idiosyncrasy, background): the more the scripts differ, the easier this transition can be made.

The only safe thing we can do with types from different origins is to balance stroke widths so that the global gray density of the page is homogeneous (no "holes" inside the text, whenever we change scripts).

These remarks concern primarily scripts that have significantly different esthetics. There is one case, though, where one can apply all possible means of uniformization, and letters can be immediately recognized as belonging to the same font family: the LGCAI group (LGCAI stands for "Latin, Greek, West/East Cyrillic, African, Vietnamese and IPA"). Very few types cover the entire group: Computer Modern is one of them (not precisely the most beautiful), Unicode Lucida another (a nice Latin font but rather a failure in lowercase Greek); there are Times fonts for all members of this group, but there is no guarantee that they belong to the same Times style, similarly for Helvetica and Courier. Other adaptations have been tried as well and it is to be expected that the success (?) of Windows NT will lead other foundries into "extending" their typefaces to the whole group<sup>10</sup>.

Fortunately,  $\TeX/\Omega$  users can already now typeset in the whole LGCAI range, in Computer Modern<sup>11</sup> (by eventually adding a few characters and correcting some others). The 16-bit font tables of  $\Omega$  allow:

1. hyphenation patterns using arbitrary characters of the group;
2. the possibility of avoiding frequent font changes, for example when switching from Turkish to Welsh, to Vietnamese, to Ukrainian, to Hawsa;
3. potential kerning between all characters.

But  $\Omega$  can go even beyond that: one can include different styles in the same (virtual) font; looking at the ISO 10646/UNICODE table, one sees that LGCAI characters, together with all possible style-dependent dingbats and punctuation, fit in 6

<sup>10</sup> As a native African pointed out to the first author, this will also mean that Africans will find themselves in the sad and paradoxical situation of having (a) fonts for their languages, (b) computers, since Western universities send all their old equipment to the Third World, but (c) no electricity for running them and using the fonts...

<sup>11</sup> *Definitely, sooner or later some institution or company will take the highly praisable initiative of sponsoring the development of other METAFONT typefaces; the authors would like to encourage this idea.*

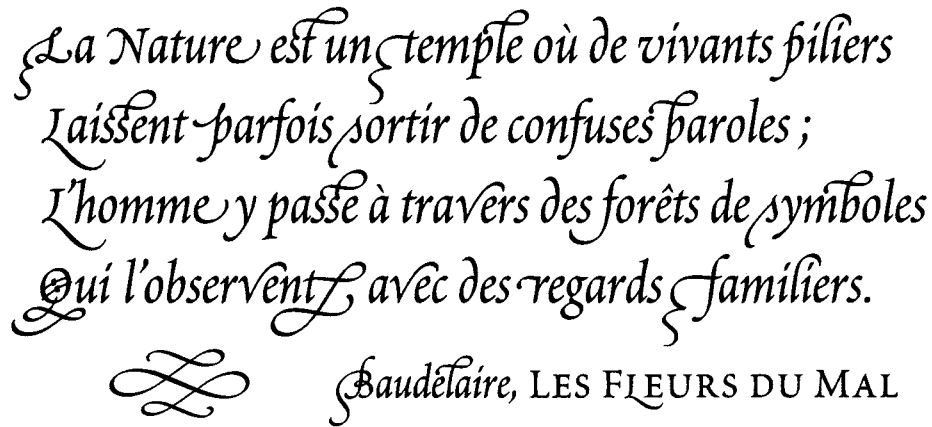


Figure 1: Sample of the Poetica typeface.

rows (1,536 characters). This means, that—at least theoretically—a virtual Ω font can contain up to 42 (!) style variations<sup>12</sup> of the whole LGCAI group, for example Italic, Bold, Small Caps, Sans Serif, Typewriter and all possible combinations [a total of 24 = (Roman or Italic) × (normal or bold) × (plain or small caps<sup>13</sup>) × (Serif or Sans Serif or Typewriter)]. Defining kerning pairs between all different styles will avoid the use and misuse of \ (italic correction) and give a better appearance to mixed-style text.

It will be quite an experience to make such a font, since many African and IPA characters have no uppercase or italic or small caps style defined yet; see Jörg Knappen's paper on African fonts (Knappen 1993) for the description of a few problematic examples and the solutions he proposes.

### Dutch, Portuguese, Turkish: the easy way

These three languages (and maybe others?) have at least one thing in common: they need fonts with a slightly different ligature table than the one in the Cork encoding. Dutch typesetting uses the notorious 'ij' ligature (think of the names of people very well known to us: Dijkstra, Eijkhout, van Dijk, or the name of the town Nijmegen or the lake IJsselmeer); this ligature appears in the Cork encoding (as well as in ISO 10646/UNICODE), but until now there was no user-transparent means of

<sup>12</sup> The authors would like to emphasize the fact that it is by no means necessary to produce all styles out of the same METAFONT code, as is done for Computer Modern. As a matter of fact the font we are talking about can very well be a mixture of Times, Helvetica, Courier; the advantage of having them inside the same structure is that we can define kerning pairs between characters of different styles.

<sup>13</sup> The figure is actually less than 18, since only lowercase small caps are needed...

obtaining it. In Ω you just need to place the macro `\inputtranslation{dutchij}` into the expansion of the Dutch-switching macro; according to Ω syntax (described in Plaice 1994) this ΩTP can be written as simply as

```
in: 1
out: 2
expressions:
'I' 'j' => @"0132;
'i' 'j' => @"0133;
.      => \1;
```

where @"0132 and @"0133 are characters 'IJ' and 'ij' in ISO 10646/UNICODE.

'f' + 'i' or 'f' + 'l'

Portuguese and Turkish typesetting do not use ligatures 'fi',... 'ffi' (in Turkish there is an obvious reason for doing this: the Turkish alphabet uses both letters 'i' and 'ı', and hence it would be impossible to know if 'fi' stands for 'f' + 'i' or 'f' + 'ı'). This is a major problem for T<sub>E</sub>X, since the only solution that would retain a natural input would be to use new fonts; and defining a complete new set of fonts (either virtual or real), just to avoid 5 ligatures, is more trouble than benefit. Ω solves that problem easily; of course, it is not possible to 'disable' a ligature, since the latter arrives at the very last step, namely inside the font. It follows that we must cheat in some way; the natural way is to place an invisible character between 'f' and 'i'; in ISO 10646/UNICODE there is precisely such a character, namely 0x200b (ZERO WIDTH SPACE); this operation could be done by an ΩTP line of the type 'f' 'i' => "f" @"200b "i" for each ligature. This character would then be sent to the Cork table's 'compound mark' character, which was defined for that very reason.

A still better way to do this would be to define a second 'f' in the output font table, which would not form a ligature with 'f', 'i', or 'ı'. This would give the font the possibility of applying a kern between



the two letters, and counterbalance the effect of the missing ligature (after all, if a font is designed to use a ligature between 'f' and 'i', a non-ligatured 'fi' pair would look rather strange and could need some correction).

### Other applications: *Ars Longa, Vita Brevis*

In this paper we have chosen only a few applications of  $\Omega$ , out of personal and highly subjective criteria. Almost every script/language can take advantage in one or another way of the possibilities of internal translation processes and 16-bit tables. For example, the first author presents in this same conference his pre-processor *Indica*, for Indic languages (languages of the Indian subcontinent (except Urdu), Tibetan and Sanskrit). *Indica* will be rewritten as a set of  $\Omega$ TPs; in this way we will eliminate all problems of preprocessing  $\TeX$  code.

*All in all, the development of 16-bit typesetting will be a fascinating challenge in the next decade, and  $\TeX/\Omega$  can play an important rôle, because of its academic support, openness, portability, and non-commercial spirit.*

### References

- Betts, G. and A. Henry. *Teach yourself Ancient Greek*. Hodder and Stoughton, Kent, 1989.
- Haralambous, Y. "Typesetting the Holy Qur'ân with  $\TeX$ ". In *Proceedings of the 2nd International Conference on Multilingual Computing—Arabic and Latin script (Durham)*. 1992.
- Haralambous, Y. "The Khmer Script tamed by the Lion (of  $\TeX$ )". In *Proceedings of the 14th  $\TeX$  Users Groups Annual Meeting (Aston, Birmingham)*. 1993a.
- Haralambous, Y. "Nouveaux systèmes arabes  $\TeX$  du domaine public". In *Comptes-Rendus de la Conférence «  $\TeX$  et l'écriture arabe » (Paris)*. 1993b.
- Haralambous, Y. "Typesetting the Holy Qur'ân with  $\TeX$ ". In *Comptes-Rendus de la Conférence «  $\TeX$  et l'écriture arabe » (Paris)*. 1993c.
- Knappen, J. "Fonts for Africa". *TUGboat* 14 (2), 104–106, 1993.
- Lagally, K. "Arab $\TeX$ ". In *Proceedings of the 7th European  $\TeX$  Conference (Prague)*. 1992.
- Levy, S. "Using Greek Fonts with  $\TeX$ ". *TUGboat* 9 (1), 20–24, 1988.
- Imprimerie Nationale. *Les caractères de l'Imprimerie Nationale*. Imprimerie Nationale, Éditions, Paris, France, 1990.
- Plaice, J. "Progress in the  $\Omega$  Project". Presented at the 15th TUG Annual Conference, Santa Barbara, 1994, and published in these Proceedings.