



HAL
open science

Methods for Processing Languages with Ω

Yannis Haralambous, John Plaice

► **To cite this version:**

Yannis Haralambous, John Plaice. Methods for Processing Languages with Ω . International Symposium on Multilingual Information Processing '97, Nov 1997, Tsukuba, Japan. pp.115-128. hal-02101608

HAL Id: hal-02101608

<https://hal.science/hal-02101608>

Submitted on 16 Apr 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Methods for Processing Languages with Ω

Yannis Haralambous* John Plaice†

November 10, 1997

Abstract

In this paper we discuss general issues of multilingual typesetting and methods used by the Ω typesetting system. We describe the different levels of language processing with Ω , giving special emphasis to the level of virtual fonts; in particular we give a complete description of the configuration file which is used to create 16-bit virtual fonts out of 8-bit PostScript fonts containing the necessary glyphs.

1 Introduction

Typesetting in different scripts and languages is a problem which has been mainly solved today, either by big software industries which adapt their products to what they consider local typesetting specifications, or by individual users brewing their own limited but practical systems for their personal needs.

Typesetting in different scripts and languages by *keeping the quality* of traditional typography for every script, and having an open system which can be adapted to any language and script without loss of power or *quality*, is up to now an unsolved problem, and Ω aims to be a solution to precisely this problem.

If we want to have an open system, and “re-invent traditional typography”, for any script and language, just as Don Knuth did with \TeX for the English language, we need to analyze and compare the aspects of multilingual typesetting, and establish general methods to solve them.

The purpose of this paper is to show that efficiency and openness can be improved by *working on different levels*, each one adapted to a specific aspect of multilingual typesetting. These levels correspond to methods used in the Ω system, and we are going to describe them by giving concrete examples from Latin, Greek, Cyrillic, Arabic, Hebrew alphabet languages.

*Atelier Fluxus Virus, 187, rue Nationale, 59800 Lille, France. yannis@pobox.com

†plaice@acm.org

2 Ω Methods: an Overview

When developing an Ω typesetting system for a given language, one can work on the following levels:

1. The font level. A font is a container of glyphs, needed to typeset in a given language. These glyphs may or may not correspond to the graphical unities of a script, whether these are called “letters”, “ideograms” or otherwise; these glyphs may be parts of graphical unities, or combinations of them, or new independent symbols.

They must be provided to the screen previewer and to the printing engine. Ω itself is not concerned by them (just like \TeX Ω works only with metrics, and as we will see in the next item, these metrics are not necessarily the ones of the fonts actually containing glyphs).

The font level is the lowest development level, in the sense that glyphs are indivisible unities which can be used in other higher level structures but cannot be dynamically modified by Ω itself.

2. The virtual font level. Once we have the glyphs we need, we combine them to form what a language usually considers as a grammatically correct script entity (an ring accent alone is of no use, but \AA is part of the grammar of the Swedish language).

A virtual font character is a combination of glyphs taken from several “real” fonts, or of other virtual font characters. In the \AA example above, the glyphs of the ring and of the letter a can actually be taken from entirely different fonts, in different formats (bitmap, PostScript, TrueType, etc.).

Going from “real” fonts to virtual fonts is mainly a matter of optimisation of storage and memory: taking the seven Greek vowels, the three accents, two spirits, diaeresis, subscript iota and macron/breve diacritics (that makes 16 glyphs) we produce 336 (!) virtual glyphs. The description of such a character in a virtual font is hundreds of times smaller than the PostScript code describing an hypothetical similar glyph.

Virtual fonts are the ones used directly by Ω : they can have up to 2^{16} positions and 2^{32} kerning pairs. The files created by Ω can be processed (pre-viewed, printed) by utilities we have adapted; if the user has to use his/her own utilities which are not “16-bit clean”, there is a tool to “devirtualize” the files [i.e. replace virtual fonts by the underlying real fonts] and make them as standard as \TeX output files.

3. The Ω TP level. When Ω reads a document, it first tokenizes it and expands commands and then forwards the data to the typesetting engine. Between these two steps we introduce an arbitrary number of filters which we call Ω Typesetting Processes (Ω TPs). These are written in a Lex-like syntax, are

loaded while reading the document and can be activated and de-activated dynamically.

A typical example for an Ω TP is contextual analysis of Arabic. Of course, this operation can also be done by \LaTeX commands; but an Ω TP will do it much faster, it will avoid conflicts with other \LaTeX commands, and it will be much easier to create and configure.

Contextual analysis of Arabic is a typical example of a script property which is low level and should not use any macros or other high-level structures. In the case of the Latin script one couldn't possibly expect from a user to constantly think of 'fi' and 'ff' ligatures and place them manually, and it would be very bad strategy to use an "fi ligature command"; in Arabic this property is of the same nature, and hence should remain completely transparent.¹

But efficiency should not only be limited to speed of typesetting: sometimes it is very important to optimize also the configuration and customization time and effort. A typical example is the management of encodings, whether input or output: by using a universal encoding (we call it Unicode++, it is a superset of Unicode) as intermediate step of our Ω TPs, every new input encoding requires only a $\text{foo} \rightarrow \text{Unicode++}$ Ω TP and every new font encoding only a $\text{Unicode++} \rightarrow \text{foo}$ one; these significantly easier to make as if one had to rewrite all processing steps (including contextual analysis and other bells and whistles).

Ω TPs can also be used to hide things from other Ω TPs and make them re-appear on a later stage: imagine you want an Arabic word with one letter in red. Placing a $\backslash\text{red}$ command in front of that letter would normally break the contextual analysis. Not so if we define a range of "characters" which are considered as characters and not commands, and whose only purpose is to carry information (for example the color switch). These characters would not interfere with contextual analysis, and after the analysis is done, would be transformed into commands.

4. The hyphenation and sorting engines. Hyphenation and sorting rules are grammatical properties of a language: they have nothing to do with typographical aspects, input methods, font encodings, etc. They have to be performed on the Unicode++ level, which is the most abstract one. Once defined at this level, they can immediately be used with every input and output encoding. Ω hyphenation and sorting engines are still under development; for the time being, Ω doesn't sort, and hyphenates like \TeX , using the (virtual) font encoding.
5. The macro-command level. \LaTeX is a terrific construction, featuring commands on different levels: \TeX primitives, plain \TeX commands, internal

¹And undoubtedly this would be the case if the lingua franca of computer science was Arabic or Urdu, and not English...

\LaTeX commands, higher \LaTeX commands and environments, user-defined commands and environments. Our goal is to keep all script-dependent processing independent of the \LaTeX macro level; once we have reached this goal, every \LaTeX package will be useable in any script and language, and both Ω package and \LaTeX package developers will be able to work independently but still producing mutually compatible software.

In the following sections we will discuss these levels in more detail.

3 Level 0: Studying Typesetting Aspects of the Target Language

This level has little to do with Ω , but it is the most fascinating, and that's why we include it. Learning a language is a difficult process that can take years, especially if one is not in the adequate environment. But studying the typesetting aspects of a given language is a different task, which is humanly possible in a few months, and which is equally fascinating. After having learned the basics (the alphabet/writing system, punctuation, special symbols) one spends months browsing books and searching for exceptions, special cases, unwritten laws and conventions, and in general trying to feel the esthetics of a script.

Often one finds symbols unknown even to natives, and almost always, when looking in older books one discovers what has been lost, what was brought by the computer, what current convention is just a technical compromise and would horrify any traditional typographer.

Once all the facts have been gathered, one has to make a model of the language's writing system. For example, traditional Mongolian uses the same model as Arabic language (letters connected by a basic stroke, and having contextual forms), only in the vertical direction. Indic languages mostly use the same model: letters connected forming clusters, with parts on the left or the right, above or below. Cambodian uses a more elaborate model having consonants in the center of the cluster, subscript consonants centered or adjacent below them, vowels surrounding them, either centered or right justified and eventually accents. Urdu *Nastaliq* uses a very particular model: characters connected as in Arabic but written diagonally from upper right to lower left; furthermore the contextuality of *Nastaliq* letters is stronger than the one of ordinary Arabic *Naskhi* or *Kufi*, since some letters lose their distinctive features (for example, dots) or obtain different features when in the presence of specific other letters [imagine a rule saying “whenever an *i* is followed by an *s* or a *b*, the latter gets the dot of the former: is this possible?]

After establishing the model one has to decide what is general behaviour of the script and what is exceptional or should always be requested by the user. Take for example the Arabic ligature lam-lam-hah لله used in the word Allah الله ; this ligature is mandatory for writing the word ‘God’, and consists of two lams (the second one with shaddah—meaning it is a double lam, so that actually there are three lams—

and vertical fatha), and a hah. When we write these letters without ligature we obtain: ﷲ. We could establish a rule saying “whenever the user asks for three lams, a vertical fatha and a hah, he/she should get the ‘God’ ligature instead”. But is this rule 100% sure? Could’nt we find some other word having exactly those letters, and maybe some other ones? On the other hand, making a macro for the word Allah would be sub-optimal as well. In our transcriptions we have chosen an intermediate approach: using a special way of transcribing this ligature, avoiding any ambiguity with other uses of this character string than ‘God’. Of course, this is a solution we have given, the user can easily configure the transcriptions to fit his/her needs, taste and convictions about what is common sense and what not.

Often this modelization is a sequence of more-or-less succesful attempts, hopefully converging to some sufficiently intelligent and compomised system. And this may also change under the burden of real-life use: a typesetting system is not valid until several natives and non-natives haven’t done entire books with it, giving the developpers critical feedback, and having created their own transcriptions and conventions.

Enough with generalities, let us return to the different levels of language processing with Ω , starting with the lowest level, which is the font level.

4 Level 1: Dealing with glyphs

What we call a “glyph” is actually simply a character from a PostScript, TrueType or Metafont font; so why calling it that way? Because we want to make the distinction between the “images”, and the “combinations of images” which Ω will use as characters for typesetting. In the next section we will see how we combine glyphs; in this section we should normally speak about the design of glyphs, the esthetics of multi-script font design, the aspects of homogeneization, the design of characters which never have been designed before, and the risks of innovation. Discussion of these issues could fill entire books, so we will limit ourselves to a few very precise issues which we had to face when designing these fonts.

4.1 A Multi-script Font: Does it Make Sense?

Does it make sense to design a Cyrillic font with Roman esthetics? Or a Greek font with Cyrillic esthetics? Do these esthetics actually exist? After all we live in a period where most of the printed documents use Microsoft’s or Apple’s standard Times and Arial/Helvetica fonts which do share the same esthetics—if any—and the first multi-script font would only be the merging of all those fonts.

Let us be realistic: the number of documents using more than one script is very small (although some combinations may occur more often than others: for example, in Greek texts one finds relatively often foreign words written in the Latin alphabet, simply because every Greek can read that alphabet as well) and it certainly a more

serious goal to obtain high *quality* typesetting staying inside a given script than trying to combine more of them.

But still there are excuses: one could try, when going from one script to the other, to *keep* only those features which can be common without conflicting with a given script’s internal esthetics: width of strokes (fat/thin, curves, etc.), serifs (whenever appropriate), global dimensions (upper/lowercase letter height and depth, sidebearings, etc.).

One could also try to use common graphical elements, as designer secrets, so that the reader feels a certain homogeneity without exactly where it comes from.

Anyway, to give a direct answer to the question, yes, it makes sense, at least for the Ω project; for two reasons: first of all, there is no multi-script font yet (at least not in the public domain) and the typographical atrocities we have seen in the last ten years are already sufficient to motivate us, and secondly because we consider our fonts as starting points: we will try to give them the most local typesetting *quality* without losing homogeneity—people with more specific (and more limited geographically) needs will replace the parts they consider suboptimal with other fonts, optimal for a given combination of scripts (or a single script) but not for the others. We are giving them the framework and the tools to do this—and we try to sensitize them to this eventuality.

4.2 How About Arabic, Hebrew, Indic scripts?

In the previous section we were mainly speaking about Cyrillic, Greek and Latin: three scripts which are close relatives, having common—although far away chronologically—ancestors. When we extend our fonts to Semitic and Indic scripts, some problems vanish and other arise.

The problems that vanish are those of “false twins”: a lowercase Latin ‘a’ is designed to be identical to lowercase Cyrillic one, but is this an universal truth? Must a Cyrillic El ‘Й’ look different than a Greek Lambda ‘Λ’? The answer seems to be yes (although there are Cyrillic fonts with Greek-like EIs), but how about the Cyrillic Pe vs. Greek Pi? When comparing Arabic and Hebrew and Indic and Latin, you don’t face those problems.²

There is a very interesting counter-example to that assertion; try to identify the script of the following two words: Հս Պսսսսսննննի. Does it look familiar? It should do so, because these glyphs have been designed in absolute conformance to the design principles of Latin script. Guessed what script it is? Well, it is Armenian, and contrarily to our Berber experimentations, this typeface is authentic: it can be found in traditionally typeset books and doesn’t shock Armenians the least. Maybe it started as an experiment a hundred years ago, and has by now become a typeface standard.

But this case remains *unique*: although we have seen other scripts designed so that they imitate the Latin script these were only artistic experimentations and can

²Except perhaps for punctuation, digits and other general typographical symbols, since more and more scripts tend to use the Latin versions of these.

not be taken into consideration for general use. Scripts like the Arabic, Hebrew, Indic, South-East Asian ones, remain esthetically independent of LGC (Latin-Greek-Cyrillic). What can be done to give a multi-script document, an homogeneous overall impression?

First of all, the fact that lines are shared by different scripts (we suppose here that these scripts are either horizontal only or vertical only), means also that they share the same baseline *skip*. This at first glance harmless fact can give very annoying results for scripts with strong ascenders and descenders. For example, Khmer has sometimes double subscripts and accents: the baseline *skip* must be significantly higher than for LGC; but what if there are only a few Khmer words in a paragraph of several lines of LGC text? Does this mean that only those lines will have a bigger baseline *skip*? This is the first problem one has to solve (the solution depends on the combinations of script... there is no general solution).

The second factor to look at is the width of strokes; many scripts have distinctive fat and thin strokes (not all of them, for example Berber letters are pure geometric forms, drawn with a fixed width pen). If this is the case, one can start by checking if these can have the same values as LGC fat and thin strokes. This is of course a very simplistic approach, but it can serve as a good starting point.

We also check the punctuation is identical, or similar, to the one of LGC (for example, the Arabic semicolon is an inverted LGC semicolon, with bigger sidebearings). All scripts we know of use LGC parentheses, brackets and braces. These are a good common test to see in which extent a script is different from LGC, and how one can play with the new script to bring it closer to LGC.

If these methods fail, one can always use the—more scientific, and less intuitive—method of comparing typographical “grays”: one prepares fonts in different weights, and compares the global impression of gray one gets when looking the page from a certain distance, with the one of LGC text. For this it is best to have Metafont fonts, since these allow fine tuning of weight, without redrawing of characters. In fact, through this method, one can find both the right weight and baseline *skip*.

4.3 Can One Avoid Innovation?

Hardly. After all, we are fixing ourselves new goals: our fonts don’t have just to be nice-looking and to follow a certain number of traditions; we also want the different scripts to fit together, and we also want the typographical variations of one script to have equivalents in other scripts, or at least we want the entire character set of a script to exist in the same number of variations.

Up to now, an IPA font was considered “exotic”, and no one really expected it to have bold and italic forms. After all, as in mathematics bold and italic characters are different symbols, may be IPA can have no bold or italic version anyway. This is only partly true: some IPA characters are in fact “italic characters in upright position”, so that globally switching to italic would cause confusion. The same goes for small capitals. Not so for bold, though; a complete Unicode++ font family needs a

bold IPA part.

Problems don't end there: the IPA has been designed to work with serif typefaces; what about sans-serif ones? The letter 'a' of a sans-serif font may already look like an 'v'—this is the case for the Futura typeface—, how do we point out the difference? Same question for 'g' and 'ǵ' [one finds the latter even in some serif fonts]. The IPA character 'ɪ' (which is in fact a small-cap 'i') differs from the dotless 'i' through its serifs, but what if the serifs aren't there? Should we postulate that IPA cannot be written with sans-serif fonts? And what about semi-serif fonts, like Optima? IPA is maybe a very special case, but the same problems arise with African languages and Asian Cyrillic.

And what about Coptic and Slavonic? Just like ancient Greek can be written with the same typefaces as modern Greek (although there are typefaces specifically made for ancient Greek, like Porson or Greek Sans Serif), why not typesetting Coptic and Slavonic excerpts without breaking the homogeneity of the surrounding “modern” font. We have tried to do this for Slavonic (including not only letters and accents found in Unicode, but also all the other diacritics and symbols needed for Old Church Slavonic), and the example looks like this: Цѣвѣща ѣи ѣиць ѣ рече ѣи. Nothing really exciting to look at, but that is the point: it should fit with the surrounding text instead of looking “exotic”. This is innovation; is it really useful? Time will show.

4.4 How Much can be Done by Combining Glyphs?

As we will see in the next section, a big part of the Unicode++ table (at least for LGC) is obtained by combining a limited amount of glyphs. This corresponds to a grammatical reality: it is quite natural to assume that placing diacritics does not affect the shapes of whether the base character or the diacritic itself. Often this is true, but there are times when typographical quality requires special shapes. This is typically the example of Arabic alif with hamza: alif is a “high” Arabic letter, we have chosen it to have approx. the height of a capital LGC letter; the hamza diacritic is placed upon the alif; it often happens that an alif with hamza takes another diacritic: a vowel (fatha, or dammah) or even a nountated vowel (fathatan, or dammatan): the latter are already quite high. The combination of alif, hamza and for example, dammatan, would normally be much too high. In traditional Arabic typography one can choose smaller diacritics, and after all there many esthetic ligatures which are also quite high, so that the result is homogeneous in the end. But in a modern typeface like the one in the Ω Arabci system, some other solution had to be found. We have chosen to lower the alif letter, so that the alif+hamza has almost the same height as the alif alone: compare \aleph and \aleph^{h} , so that the alif with hamza and a vowel is not excessively high, as can be seen in the following word: أَجْرِي.

This method can be applied to several cases in other scripts: one can imagine an 'h' with a lower vertical stem, so that the accent of the Esperanto character 'ĥ' is not excessively high. This is finally a matter of taste, and we intend to include this kind of characters as variant forms in our font, and leave the decision on using

them or not to the user.

Another typical example is the one of the diaeresis over Greek vowels iota and epsilon, whether lowercase or uppercase; the distance between the two dots of the diaeresis depends on the letter beneath: compare $\ddot{\iota}$, $\ddot{\upsilon}$, $\ddot{\text{I}}$, $\ddot{\text{Y}}$, while in the Latin script the diaeresis (Umlaut, tré) always has the same width: $\ddot{\text{i}}$, $\ddot{\text{o}}$, $\ddot{\text{w}}$...

5 Level 2: Combining Glyphs in Virtual Fonts

A character of a virtual font can be a combination of characters of other fonts (whether real, in which case we actually have glyphs, or again virtual, in which case we have sub-combinations etc.), of black boxes of arbitrary height and width, and of PostScript code [the latter feature is not implemented in all DVI drivers, so one should refrain from using it, at least for the moment].

We have chosen to work intensely with virtual fonts for two reasons: first because by combining glyphs we can optimize space (and space management is crucial when you deal with 16-bit fonts), and second, to be able to use 16-bit fonts without re-writing all DVI drivers of the world (the underlying real fonts are 8-bit only, so that a devirtualized 16-bit font, becomes 8-bit and can be processed by any decent DVI driver).

5.1 The General Scheme

The idea is to have a few “big” virtual fonts, for purposes distinct enough so that one doesn’t need to switch fonts all the time. For the moment, the virtual fonts we have are `om1gc` (covering LGC, IPA, Armenian, Georgian, Left-to-right Tifinagh) and `omr2l` (covering Arabic and Hebrew alphabet languages, Right-to-left Tifinagh and forthcoming Syriac alphabets). Other “big” fonts will follow, for Indic languages, South-East Asian ones, CJK and dingbats.

Each “big” virtual font is accompanied by a certain number of “small” PostScript fonts. These have names using the following scheme: `om` (for Ω), `+se` or `ss` (for serif or sans-serif) + a two-letter code defining the font encoding (`la` for Latin, `gr` for Greek, `cy` for Cyrillic, `cx` for Extended Cyrillic, and so on) + nothing or `b` or `i` or `bi` (for plain, bold, italic, bold-italic).

To build a “big” virtual font out of these, we use a Perl utility, called `makeovp.pl`. This utility will read a configuration file (which usually has the same name as the font and the `cfg` extension), a special file containing a certain number of “generic” kerning pairs, all the AFM files required to build the font, and a few global parameters taken from the command line. `makeovp.pl` will build PL files for all PostScript fonts, and an `OVP` file for the “big” virtual font; once this operation is completed, the utilities `PLtoTF` and `OVP2OVF` will convert these files to the binary font files needed by the Ω system for typesetting. This process can be entirely batch (using Unix scripts, DOS scripts under Windows or AppleScript on the Mac, with the CMacTeX Ω implementation) and is needed only when the font is modified; an av-

erage user is not expected to modify the font, but nevertheless we give him/her, all the necessary tools to do so.

In the following sections we will describe the configuration file and the operation of `makeovp.pl` in general.

5.2 Files Read by `makeovp.pl`

The Perl utility `makeovp.pl` will read AFM files, the configuration file, and the *kerning* file.

Why having a separate *kerning* file, instead of relying on the *kernings* included in the individual AFM files? The answer is trivial: because no AFM file can contain *kernings* with characters from other AFM files, and this is exactly what we need in our big virtual font: for example, punctuation is contained in a “common” PostScript font, while the “Latin” one contains only Latin alphabet letters, the “Greek” one contains Greek letters and specifically Greek punctuation (Greek guillemets), and so on. To *kern* between Latin letters and common punctuation, or between Greek letters and common punctuation one needs a separate file, with *kerning* pairs from different PostScript fonts.

Otherwise the *kernings* file has the same syntax as any AFM file:

```
KPX foo goo -121
```

where `foo` and `goo` are characters and `-121` denotes a *kern* of 121 PostScript units to the left.

AFM files contain information on character metrics: the width of the box containing the character, and the bounding box of the actual shape of the character. This information is sufficient for building a virtual font, and we will see how.

Finally a certain number of global values are given to `makeovp.pl` through the command line, we will see these in the next section, when describing the different operators of the configuration file.

5.3 Structure of the Configuration File

The configuration file (`om1gc.cfg`, `omr2l.cfg`, etc.) contains one line for each character of the virtual font. This line consists of (a) a 4-digit hexadecimal number, which is the character’s position in the font, (b) an operator, (c) one or more character names, depending on the operator, (d) a certain number of optional operators and values.

We will describe one by one the different operators, giving concrete examples.

5.3.1 The operator `N`

`N` stands for “NAME”, and means simply that the string following the operator is the (PostScript) name of one or more characters in some of the fonts loaded. For example,

03A5 N Upsilon

means that at position 03A5 of the virtual font, we have placed the *Greek* letter capital Upsilon Υ . The PostScript names we have used try to be as standard as possible (of course, most of the time there is no standard, so we just have to invent names...)

The same glyph can be used for several font positions: for example, the Croatian Dze Đ has exactly the same shape as the Icelandic Eth Ð: we will use the same glyph, and hence the same PostScript name. Nevertheless we try to optimize the use of glyphs so that one can typeset in one script without necessarily loading the PostScript fonts for other scripts: for example, although the capital Latin ‘A’ has the same shape as the *Greek* capital Alpha, we will use two different glyphs in two different PostScript fonts, so that when typesetting *Greek* one can avoid loading the Latin PostScript font. These considerations will be irrelevant when we will be able to use 16-bit monobloc PostScript fonts; for the moment, this is not part of the Adobe Type 1 font specification.

There are several options we can use with this operator: #KRN, #KRNLEFT, #KRNRIGHT, #HOFFSET, #VOFFSET.

The first three concern *kerning*: we can state that a given character has the same *kerning* behaviour as some other character, which we give by name. For example, ç will be *kerned* exactly like the letter ‘c’: everytime there is a *kerning* pair in the *kernings* file using ‘c’ a new *kerning* pair will be defined, using ç instead of ‘c’ (and if there is a ‘c-c’ *kerning* pair, three new *kerning* pairs will be defined: ‘ç-c’, ‘c-ç’, ‘ç-ç’). Sometimes we *kern* a letter like some given letter on the right and like some other letter on the left: this is typically the case of ligatures: æ will be *kerned* as ‘a’ on the left, and as ‘e’ on the right; in that case, we use the #KRNLEFT and #KRNRIGHT operators:

```
OOC6 N AE #KRNLEFT=A #KRNRIGHT=E
OOE6 N ae #KRNLEFT=a #KRNRIGHT=e
O110 N Eth #KRN=D
```

The operators #HOFFSET and #VOFFSET will offset the glyph without affecting the box of the character.

5.3.2 The operators XHAC and CHAC

These operators will place diacritics over letters, which are considered to have the same height as the lowercase letter ‘x’ (x-height) or the one of an LGC uppercase letter (cap-height). The idea is that the height of letters can fluctuate: a round letter, like ‘o’, is slightly higher than a flat letter, like ‘z’, to counterbalance a well-known optical effect. The height of accents over these letters must be the same, even if they aren’t exactly of x-height, or cap-height: take for example ó and ú; if we would take the real height of these letters, the accent on the former would be slightly higher than the one on the latter.

How is this accent placed precisely? The Perl utility centers the bounding box of the accent over the bounding box of the letter, with a fixed distance of EPSILON

(a global value) between the lower boundary of the accent and either the x-height or cap-height of the font (those two are also global values we provide the utility with).

The options available are: #KRN, #KRNLEFT, #KRNRIGHT, #LETTERLIKE, #ACCENTLIKE and #LETTERREVLIKE.

The options #LETTERLIKE and #ACCENTLIKE allow us to use given glyphs, with the metrics of other glyphs. These options are extremely important in certain cases. A typical example is Vietnamese: the letter ‘o with hook’ *ơ* is significantly wider than the plain ‘o’, nevertheless accents have to be centered on the “o part” of the letter: *ô*, *ò*, *ơ*, *õ*. This trick allows us to correctly place an accent on the vertical stem of a ‘b’ or an ‘h’: *ḃ*, *ḥ*, *ḥ̈*; to obtain this result, we simply ask the accent to be placed as if the letter was an ‘l’:

```
0603 CHAC b dot #LETTERLIKE=1
0623 CHAC h dot #LETTERLIKE=1
0627 CHAC h dieresis #LETTERLIKE=1
```

We have the same functionality with accents: using the #ACCENTLIKE operator we can place an accent as if it was some other accent. The typical example is again Vietnamese, where there are combined accents ‘circumflex + grave’, ‘circumflex + acute’, which have to be centered with respect to the middle axis of the circumflex (and hence as if there were no acute or grave accent):

```
06D0 CHAC 0 circumflexacute #KERN=0 #ACCENTLIKE=circumflex
06D1 XHAC o circumflexacute #KERN=o #ACCENTLIKE=circumflex
06D2 CHAC 0 circumflexgrave #KERN=0 #ACCENTLIKE=circumflex
06D3 XHAC o circumflexgrave #KERN=o #ACCENTLIKE=circumflex
06D4 CHAC 0 circumflexhook #KERN=0 #ACCENTLIKE=circumflex
06D5 XHAC o circumflexhook #KERN=o #ACCENTLIKE=circumflex
06D6 CHAC 0 circumflextilde #KERN=0 #ACCENTLIKE=circumflex
06D7 XHAC o circumflextilde #KERN=o #ACCENTLIKE=circumflex
```

Another example is Slavonic, with letters such as *ѣ̈*, where the accent has to be placed on the right part of the ligature, as in *ѣ̈́*. This means that we should use the metrics of a given character, justified on the right of our box: this is the rôle of the #LETTERREVLIKE operator.

5.3.3 The operator VARAC

This operator has the same syntax (and takes the same options) as the XHAC and CHAC operators, but produces accents of arbitrary height: it was necessary for letters like ‘t’ which are neither x-height nor cap-height.

5.3.4 The operators XHLLAP and CHLLAP

These operator have the same syntax (and take the same options) as XHAC and CHAC but instead of centering the accent over the letter, justify it on the left. We have

used this technique for certain cases where the accent placement had to be very precise: for the lowercase Greek *iota*, which takes special accents (narrower than for the other Greek vowels).

5.3.5 The operators DBLXHAC, DBLCHAC, DBLVARAC, DBLXHLLAP and DBLCHLLAP

These operators have the same properties (and take the same options) as the operators without DBL prefix, but take two arguments: they will place two accents on the same letter, either by centering them or by justifying them to the left. The two accents will be placed an EPSILON distance apart. An optional parameter #ACCENTOFFSET is used to modify the vertical distance between the accents. We have used this parameter for Greek letters with breve and acute/grave accent, where the latter is brought a little closer to the former than calculated through their bounding boxes (a bounding box gives no information about the actual shape of a glyph).

5.3.6 The operators SUBZERODPAC and SUBFIXDPAC

These operators are used to place diacritics below characters. The first operator will consider the base character being of zero depth and will center the diacritic according to its bounding box, so that the upper boundary of the diacritic is at depth EPSILON. This is useful for accents such as the lower dot as in *ŋ*, the lower ring as in *ä*, etc.

The second operator will move the diacritic horizontally only, supposing that its depth is already correct. We use this operator for accents such as the cedilla: *ç*, *ł*, *Ń*, *ř*, etc. Notice that although the cedilla is placed on the virtual font level, the letters with *ogonek* are separate glyphs: this comes from the fact that the cedilla is connected to the letter by a straight stroke, which is never tangent to the character, and will neither distort it, nor be distorted by it: this allows us to use the same cedilla shape for all character+cedilla combinations. On the other hand, the *ogonek* sticks to the letter and becomes part of its design; every character+*ogonek* combination has been designed separately to give an optimal result.

5.3.7 The operators XHSUBZERODPAC, CHSUBZERODPAC, XHSUBFIXDPAC and CHSUBFIXDPAC

These are combinations of the previous ones, and the CHAC and XHAC operators: they allow us to place simultaneously diacritics over and under a character, like in the case of *Š*, *č*, *ř*, etc.

5.3.8 The operator ADJ

This operator allows us to concatenate two characters, using a box with width equal to the sum of the widths of the two boxes, \pm the eventual *kerning* between those characters, and height/depth the maximum height/depth of the two characters. We first wanted to use that operator for the Croatian digraphs ‘Lj’, ‘Nj’, etc. but the

decided that the whole idea of having code positions for these digraphs was so silly that we could very well do without. The operator nevertheless proved very useful for cases such as the Greek capital vowels with accent Α, Ε, Η, Ι, Ο, Υ, Ω.

This operator takes a special option: #MOVELEFT; this option allowed us to override a kern between two characters and move the second one horizontally, together with his box.

5.4 Conclusion

The operators described above allowed us to create virtual fonts for the Latin, Greek, Cyrillic, Arabic, Hebrew, Armenian and Berber script. The advantage of the Perl language is that it is platform-independent and that modifications can be quickly and plainlessly. We will pursue the development of virtual fonts adding new scripts, and adding new operators whenever this is necessary.

Forthcoming developments include operators which will use the PostScript code of given characters to obtain more information: for example the exact coordinates of the vertical stem of letters like ‘h’ or ‘d’, so that accent positioning is more accurate for an arbitrary font.

6 Levels 3 and 4: ΩTPs and L^AT_EX Macros

Once the structure of fonts is well organized, we use ΩTPs for low level script- or language-dependent operations and macros for higher level operations. It would lead us too far to start discussing general principles of ΩTPs and macros. Instead, we plan to publish a number of papers on the different language and scripts covered by Ω, to start with Multilingual Typesetting with Ω, a Case Study: Arabic, which the reader can find in this same volume of proceedings of the 1997 Tsukuba Conference on Multilingual Computing.

To obtain more informations on the Ω project, please visit our Web site:

<http://www.ens.fr/omega>