

Multilingual Typesetting with Ω , a Case Study: Arabic

Yannis Haralambous*

John Plaice†

Abstract

In this paper we describe the internal structure of the Arabic script package for the Ω typesetting system, as well as the techniques and tools used for its development. This package allows typesetting using regular \LaTeX styles, in all Arabic alphabet languages: Arabic, Berber, Farsi, Urdu, Pashto, Sindhi, Uighur, etc.

We also give a description of the character codes added to Unicode, to obtain the Unicode++ encoding, used by the Ω system for typesetting purposes.

1 Overview of the Ω Arabic Script Package

Typesetting with Ω is a process similar to typesetting with \TeX : the user prepares a “source” file, containing the text of ^{his}_{her} document and a certain number of macro-commands for attribute changes of the text (font characteristics, language, case, etc.), references to figures (included in graphical format files on disk) and other material included in or accompanying the text.

Once this source file prepared, Ω is launched: it reads the file, expands the commands and typesets the text accordingly. To perform this task, Ω loads and executes several Ω TPs (Ω Translation Processes), which take care of low level properties of the document (contextual analysis of the script, case switching according to script and language, etc.). It also uses different fonts, most of which are virtual, in the sense that they themselves call other fonts. On a higher level, such a document uses \LaTeX packages, some of them modified to take advantage of the additional features of Ω vs. \TeX .

The leading idea of the Ω Arabic Script Package (as of any Ω language package) is that the low level properties of the script have to be separated from higher level typesetting commands. For example, contextual analysis of the Arabic script has to be completely independent of the \LaTeX command level, so that one can use Arabic text in any context (inside a table or a formula, or deeply nested inside several \LaTeX environments and commands, etc.) and under any circumstances.

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

†plaice@acm.org

There are two key aspects to Arabic script typesetting, unfortunately of unequal complexity: the first one is contextual analysis, that is the fact that Arabic letters change shape according to their position in a word, or according to the fact that they are part of an abbreviation, etc. This aspect can be handled easily and efficiently by Ω TPs. The second aspect is more global: it is the fact that Arabic script is written from right to left.

Two methods can be applied: the first one is to change the default direction of the whole document. This method is extremely efficient when the document is entirely in Arabic, or if left-to-right text excerpts are exceptional. Being global, this method applies also to page-level typesetting methods, such as the order of columns in a multicolumn environment, etc. Of course, mathematical formulas are not affected by this global direction change.

The second method is to keep left-to-right as default direction and to temporarily switch to right-to-left for every Arabic script sentence. This can be practical for a document where Arabic excerpts are exceptional.

2 Parts of the Ω Arabic Script Package

This package consists of the following elements:

1. The `OmegaSerifArabic` PostScript fonts: files `omsea1.pfb`, `omsea2.pfb`, `omsea3.pfb` and the corresponding AFM files. A Sans-serif font (`OmegaSansArabic`), as well as additional styles of the Serif font are under development.
2. The virtual font `omr1`: files `omr1.ovf`, `omr1.ofm`, `omsea1.tfm`, `omsea2.tfm`, `omsea3.tfm`.
3. The configuration file `omr1.cfg`, which is used by the PERL utility `MakeOVP` to create the virtual font out of the AFM files and other information.
4. A certain number of Ω TPs:
 - (a) `7arb2uni.otp`, 7-bit Arabic/Farsi transcription to Unicode;
 - (b) `7ber2uni.otp`, 7-bit Berber transcription to Unicode;
 - (c) `7urd2uni.otp`, 7-bit Urdu transcription to Unicode;
 - (d) `7pas2uni.otp`, 7-bit Afghanistani Pashto transcription to Unicode;
 - (e) `7pap2uni.otp`, 7-bit Pakistani Pashto transcription to Unicode;
 - (f) `7snd2uni.otp`, 7-bit Sindhi transcription to Unicode;
 - (g) `uni2cuni.otp`, contextual analysis, sending Unicode++ to cUnicode++ ('c' for 'contextual');
 - (h) `cuni2oar.otp`, cUnicode++ to `omr1` font.

These Ω TPs are available in human-readable and compiled binary format (**OCP**), the latter being loaded by Ω on runtime.

5. A \LaTeX style (`arabic.sty`) defining a command that will activate and deactivate the Ω TPs.
6. Documentation and test files (`testarab.tex`, `testsind.tex`).

3 Installation of the Ω Arabic Script Package

To use the Ω Arabic Script Package you must have Ω version 1.45 or higher installed on your machine. Place **OFM**, **OVF**, **TFM** and **OCP** files where the system expects to find them (if in doubt, consult the `texmf.conf` file). Keep the `arabic.sty` file somewhere where it can be found by Ω . Finally add the following few lines to the `psfonts.map` configuration file of `odvips`:

```
omsea1 OmegaSerifArabicOne </foo/omsea1.pfb
omsea2 OmegaSerifArabicTwo </foo/omsea2.pfb
omsea3 OmegaSerifArabicThree </foo/omsea3.pfb
```

where `/foo` stands for the absolute path of the directory containing the PFB files.

This is all you need to do: you can start already by launching Ω on files `testarab.tex` and `testsind.tex`.

In the following sections we will describe the use of the package, from the end users' point of view. We will assume that the user is familiar with the \TeX typesetting system and the \LaTeX macro package.

4 Basic Macros

Before starting a new document one has to choose if the “background language” is going to be an Arabic alphabet language, in other terms, if we expect pages and columns to be typeset from right to left, and the whole global page design to be right-to-left oriented.

If this the case, then the macro `\GlobalArabic[language]` has to be used in the document header, where the optional argument *language* is one of the following: `arabic` (by default), `farsi`, `urdu`, `pashto`, `sindhi`, `custom`.

This macro will switch the global typesetting direction of the document to right-to-left and will launch the Ω TPs necessary for the language chosen.

Inside the document, independently of the choice of background language, one can use \LaTeX environments `arabic`, `berber`, `farsi`, `urdu`, `pashto`, `pashtop`, `sindhi` to switch to the corresponding language, and `latin` or `greek` to switch to a Latin alphabet language or some flavour of Greek. It should be noted that these macros are only temporary and will be adapted to a more global language-switching scheme, currently being elaborated by the \LaTeX 3 and Ω working groups.

5 Input of Arabic Alphabet Text

5.1 You Have an Arabic Keyboard

If you have an Arabic Keyboard, containing sufficiently many keys for the language you want to typeset (for example, with a standard Arabic keyboard one can perhaps typeset Farsi, possibly Urdu but not Pashto and certainly not Sindhi), you need to configure Ω to your input encoding, by providing the appropriate input Ω TP by use of the `\ArabicInputEncoding` macro, which you have to place in the header of your document. We have already written such Ω TPs for three input encodings: Macintosh Arabic (`applemac`, covering Arabic, Farsi, Urdu), Windows Arabic (1256, covering Arabic and Farsi), MS-DOS Arabic ASMO (708, covering Arabic only) and ISO 8859-6 (`iso8859-6`, covering only Arabic). If your equipment is not in this list, go to section 6 to see how to write your own Ω TPs.

5.2 You Don't Have an Arabic Keyboard

In that case you can use a Latin transcription: we have prepared ASCII Latin transcriptions for each of the main Arabic-alphabet languages: Arabic, Berber, Farsi, Urdu, Pashto (Afghanistani and Pakistani), Sindhi. Here they are:

5.2.1 Arabic/Farsi Transcription

A	ا	p	پ	z	ز	`	ع	m	م	I	ى
'a	أ	j	ج	zh	ژ	gh	غ	n	ن	y	ي
'i	إ	H	ح	s	س	f	ف	'n	ن	'y	ئ
'A	آ	kh	خ	sh	ش	q	ق	-h	ه		ء
"A	آ	ch	چ	S	ص	v	ف	"h	ة	E	ع
b	ب	d	د	D	ض	k	ك	e	ة		
t	ت	dh	ذ	T	ط	g	گ	U	و	LLah	لله
th	ث	r	ر	Z	ظ	l	ل	'u	و	SLh	ﷲ

Remarks:

1. The tah marbutah \mathfrak{e} can be written in two ways: "h or "t.
2. The waw \mathfrak{w} can be written in two ways: w or U.
3. The hyphen in front of the transcription of \mathfrak{h} is only necessary to prevent confusion between cases such as kh (خ) and k-h (ك). We suggest you use it all the time.
4. **VERY IMPORTANT:** the duplication of consonants (shaddah) is obtained by writing the consonants twice. So for example, Dmm"h will produce ضمة and not ضمة; to obtain the latter, type Dm-m"h, as for example in the word تتحرك, which presents both cases, and which is typed t-tHrrk.

Vowels and other diacritics are obtained in the following way: (they are typed after the consonant to which they belong)

fatha	a
kasra	i
damma	u
soukoun	<>
vertical fatha	a
fathatan	aN
kasratan	iN
dammatan	uN

Example: it is a trivial task now to welcome you to this system of Arabic input, by saying

```
\begin{arab}
\Huge
'aahlAaN wa sahlAaN!
\end{arab}
```

أَهْلًا وَ سَهْلًا!

Example of vowelized Arabic:

لَأَنَّهَا أَلَانَ لَا تُفَكِّرُ فِي نَفْسِهَا، وَلَكِنَّهَا تُفَكِّرُ فِي أَخْوِيهَا
وَفِي أَلْخَطَرِ الَّذِي لِحِقَهُمَا.

transcribed:

```
li 'aannahaA "Al<>'Ana laA tufakkiru fiI naf<>sihaA,  
walakinnahaA tufakkiru fiI 'aakhaway<>haA  
wafiI "Al<>khaTari "AlladhiI laHiqahumaA.
```

5.2.2 Urdu Transcription

The Urdu transcription is similar to the Arabic/Farsi one described above, with a few additional characters, and one exception.

The additional characters are ٺ, ڙ, and ڻ, transcribed by 't, 'd, 'r. The exception concerns the two different uses of the hah glyph ه. In Urdu it can be used as the second part of a digraph, such as for example جھ, in which case we transcribe it as -h; it can also be the standard consonant hah, in which case we transcribe it by x. Notice the four forms of the latter in Urdu: ه, ڀه, ڃه, ڙه, while in Arabic the same letter is written هه ه.

Example:

ہماری طرف پرانے زمانے میں دستور تھا کہ اگر کسی شخص کو کاغذ پر کچھ لکھا ہوا گرا پڑا مل جاتا تو وہ اس پرزے کو احتیاط سے اٹھا کر کہیں رکھ دیتا یا پانی میں بہا دیتا تاکہ لکھے ہوئے حروف کی بے حرمتی نہ ہو۔

transcribed:

xmArI Trf prAnE zmanE my 'n dstUr t-hA kx Agr ksI
shkhS kU kAghdh pr kchh lk-hA xUA grA p'rA ml jAtA tU Uh
As prze kU AHtyAT sE A 't-hA kr kxy 'n rk-h dytA yA pAnI mI 'n
bxA dytA tAkx lk-hE xU'yE HrUf kI bE HrmtI nx xU.

5.2.3 Pashto Transcription

The Pashto transcription is similar to the Arabic/Farsi one described beyond, with a few additional characters and some exceptions. We are proposing two ΩTPs, using the same transcription, for the two flavors of written Pashto: Afghanistani and Pakistani.

1. Afghanistani Pashto

A	ا	'z	ع	'r	ر	D	ض	g	گ	-y	ي
b	ب	c	خ	z	ز	T	ط	l	ل	e	ې
p	پ	H	ح	zh	ژ	Z	ظ	m	م	ay	ئ
t	ت	kh	خ	'g	ږ	`	ع	n	ن	ey	ی
't	ټ	d	د	s	س	gh	غ	'n	ڼ		ء
's	ث	'd	ډ	sh	ش	f	ف	w	و		
j	ج	dh	ذ	x	ځ	q	ق	-h	ه	LLah	لله
ch	چ	r	ر	S	ص	k	ک	L		SLh	ﷲ

2. Pakistani Pashto

A	ا	'z	ع	'r	ر	D	ض	g	گ	-y	ے
b	ب	c	خ	z	ز	T	ط	l	ل	e	ے
p	پ	H	ح	zh	ژ	Z	ظ	m	م	ay	ے
t	ت	kh	خ	'g	گ	`	ع	n	ن	ey	ے
't	ٹ	d	د	s	س	gh	غ	'n	نر		ء
's	ث	'd	ڈ	sh	ش	f	ف	w	و		
j	ج	dh	ذ	x	خ	q	ق	-h	ه	LLah	لله
ch	چ	r	ر	S	ص	k	ك	L		SLh	ﷲ

Nevertheless, one should be aware that an automatic transcription from one glyph set to the other is not possible because, for example, a letter such as ځ is not used in Pakistani Pashto and can be replaced by خ or ش, depending on its pronunciation in a given word.

Example of Afghanistani Pashto:

که غواړئ چه د عقل یه زیان اوضرپوه شی دا و منی ء چه عقل هغه. قوتونه پهسپری
کښی وژنی ژه ژوندې ولاړدی. ژوندون پهعمل اواراده ولاړدی. غوښتنه لواراده دهریشتر
فت اصل اواساس دی. څومره چه عقل زیاتبز هغومره اراده ضعیفه کبزی.

and the same in Pakistani Pashto:

که غواړئ چه د عقل هه زيان اوضرپوه شه دا و مني ء چه عقل هغه. قوتونه پهسرټه
کڅه وژنه ژه ژوندې ولاړدې. ژوندون پهعمل اواراده ولاړدې. غوڅتنه لواراده دهرهسر
فت اصل اواساس دې. څومره چه عقل زياتبزه هغومره اراده ضعهفه کهبزه.

transcribed:

k-h ghUA 'ray chh d`ql yh zyAn AUDrrpUh shay dA
U mnI || chh `ql hghh. qUtUnh p-hs 'rI kxI wzhnI zhh zhUnde
wIA 'rdI. zhUndUn p-h`mI AUArAd-h wIA 'rdI. ghUxtnh lUArAd-h
d-hreysht ft ASl AUAsAs dI. cUmrh chh `ql zyAtebz hghUmrh
ArAd-h D`yf-h kebZI.

A variant form ږ of ږ is provided in the font. The user can change the ΩTPs (see 6) so that the former is used instead of the latter.

5.2.4 Sindhi Transcription

Sindhi being a language with many more letters than Arabic, and using Arabic letters in a way quite different than Arabic, it is not surprising that the Sindhi transcription is fundamentally different from the Arabic, Farsi, Urdu and Pashto ones. As a matter of fact we have tried to use as few non-alphabetic characters as possible, following a more-or-less rational scheme loosely based on the correspondence between Sindhi written in Arabic and in Devanagari script and the standard transcription of the latter. Since shadda is much more seldom in Sindhi than in Arabic, the “double consonant = consonant + shadda” convention is not valid in this transcription; instead we propose a transcription of the shadda diacritic: +.

A	ا	p	پ	dh	ڌ	sh	ش	kh	ڪ	y	ي
'A	آ	ph	ڦ	.=d	ڏ	.s	ص	.n	گ	'y	ئ
b	ب	j	ج	.d	ڍ	.z	ض	g	گ	meN	ڻ
=b	ڀ	=j	ڄ	.dh	ڊ	.t	ط	=g	ڳ	eN	ڻ
bh	ڀ	=n	ڃ	=z	ڙ	.z	ظ	l	ل		ء
t	ت	c	ڇ	r	ر	`	ع	m	م		
th	ٺ	ch	ڇ	.r	ڙ	gh	غ	n	ن		
.t	ٺ	.h	ح	z	ز	f	ف	'n	ڻ		
.th	ٺ	=kh	خ	zh	ڙ	q	ق	U	و	LLah	لله
=s	ث	d	د	s	س	k	ڪ	-h	ه	SLh	ﷲ

Remarks:

1. The transcription / is used for constructions such as به (b/), ته (t/), ڪ (kh/), etc.
2. The waw و can be written in two ways: w or U.

Example:

تنهن ڪري اسان کي پنهنجي ذهنن کي سجاڳ رکڻو پوندو ۽ پنهنجي جدوجهد ۾ ڏاهپ پيدا ڪرڻي. اهو به معلوم ڪرڻو پوندو ته سنڌ ۾ هر آئي وقت ڇا ڇا ٿي رهيو آهي ۽ دشمن اسان جي ۽ اسان جي جدوجهد جي ڪلاف ڪهڙا ڪهڙا گهٽا گهڙي رهيو آهي.

transcribed:

tn-hn kry AsAn khy pn-hnjy =z-hnn khy sjA=g rkh 'nU
 pUndU ||eN pn-hnjy jdUj-hd meN . =dA-hp pydA kr 'ny. AhU b/
 m`lUm kr 'nU pUndU t/ sndh meN hr 'A 'yy wqt chA chA thy r-hyU
 'Ahy ||eN dshmn AsAn jy ||eN AsAn jy jdUj-hd jy kh1Af k-h.rA
 k-h.rA g-hA.t g-h.ry r-hyU 'Ahy.

5.2.5 Berber Transcription

The Berber transcription is different from the previous ones because it is based on a tri-alphabetic system (Tifinagh, Latin and Arabic alphabets).¹ The goal of this transcription is to enable output in the three alphabets, out of the same code. In particular, since Latin alphabet has upper and lower case, it should be possible to distinguish these (and of course ignore the distinction when typesetting in Arabic or Tifinagh). In the table below, all transcribed letters are in lowercase ASCII, but can very well be written also in uppercase, producing the same result: Tifinagh, t ifinagh or TIFINAGH will all three produce تيفيناغ.

Tr.	Lat.	Ar.	Tif.	Tr.	Lat.	Ar.	Tif.	Tr.	Lat.	Ar.	Tif.
a	a	ا	·	.h	ħ	حs	ş	ص	Ɔ
b	b	ب	θ	i	i	ي	ɣ	t	t	ت	X
c	c	ش	G	j	j	ج	ɛ	.t	ţ	ط	E
gh	ɣ	غ	ː	k	k	ك	⇒	u	u	و	:
d	d	د	Λ	l	l	ل		x	x	خ	::
.d	ḍ	ض	E	m	m	م	⊃	z	z	ز	Ж
.e	ε	ع		n	n	ن	I	.z	ẓ	ز	#
f	f	ف	⊃	.n	ñ	.ن	≠	.i	i	ئ	ɣ
g	g	ق	X	q	q	ق	≡	--	-	-	-
.g	ğ	ج	I	r	r	ر	O				
h	h	ه	≡	s	s	س	⊃				

¹The reader can find more information in Un système T_EX berbère, Études et Documents Berbères, 11 (1994), La boîte à Documents/Édisud, Paris (France).

Remarks:

1. Letter و can also be transcribed w.
2. Letter ي can also be transcribed y.
3. The stroke _ is not to be confused with the graphical connecting stroke keshideh. It is placed between words and plays a grammatical role.
4. Duplication of consonants (shaddah) again is transcribed by writing the corresponding consonant twice.

Example:

تيفيناغ، دتيرا تيمزورا ن ئمازيغن. لانت دي تمورت نغ دات تيرا ن تا عرابت
 دتلاطينيت. تولفانت دات ئمير ن وقليد ماسينيسن. ئمازيغن ن ئميرن، تارون تنت
 غفيظرا، دق ئفران، غف ئقدورن، ماشا تيفتي غف يزكوان : تارون فل اسن ئسم ن
 ومتين، د ويت يلان، د واين يخدم دي تودرت يس اكن ور ت تتون ئناطفارن.

transcribed:

```
Tifinagh, d--tira timezwura n .imazighen.
Llant di tmurt--nnegh dat tira n ta .erabt d--tla .tinit.
Nnulfant--edd dat .imir n ugellid Masinisen. .Imazighen n
.imir--en, ttarun--tent ghefi .zra, degg .ifran, ghef .igduren,
maca tigg ti ghef i .zekwan : ttarun fell--asen .isem n umettin,
d wi--t--ilan, d wayen yexdem di tudert--is akken ur t ttettun
.ina .tfaren.
```

The same code will produce the following output in the Tifinagh alphabet:

```
XzJzI:, ^_XzO· XzLz=:O· I zL·zI. II·IX ^z XL:H_I: ^·X XzO·
I X·O·H ^_XII·EzIX. IV:II·IX_^^ ^·X zLzO I :XIIz^L C·MzIzMI.
zL·zI I zLzO_I, XX·O:I_XIX :Jz#O·, ^XX zJO·I, :J zXA·OI, C·C·
XzXz+z :J z#z=·I : XX·O:I JIIV_·MI zME I :EzXzI, ^ z_X_zI·I, ^
=z·SI S::^L ^z X:^H_zM ·z=I :O X XXXX:I zI·EJ·OI.
```

and the following one in the Latin alphabet:

```
Tifinay, d_tira timezwura n imaziyen. Llant di tmurt_nney dat tira n
taerabt d_tlatinit. Nnulfant_edd dat imir n ugellid Masinisen. Imaziyen
n imir_en, ttarun_tent yefizra, degg ifran, yef igduren, maca tigg ti yef
izekwan : ttarun fell_asen isem n umettin, d wi_t_ilan, d wayen yex-
dem di tudert_is akken ur t ttettun inatfaren.
```

6 Writing Your Own Transcription

We have developed and presented in this paper a certain number of Arabic alphabet language transcriptions for two reasons: first, to show the possibilities and power of Ω , and second, to give a starting point for the user to create ^{his}_{her} own transcriptions.

The process of creating a new transcription is twofold: the first part, which can be very difficult and painful, consists of finding the combination of letters, digits and ASCII symbols which will transcribe each character; the second one, which is straightforward (modulo some precautions) is to implement this in Ω by writing the appropriate Ω TP.

6.1 A Good Transcription: Is it Possible?

There are (at least) two goals for a good transcription:

1. It has to be readable and easily memorizable. In other words, AHmd is better than ' .hmd, for denoting احمد : although an apostrophe can be considered a logical choice for transcribing an alif and the period in front of the h may denote that it is an emphatic 'h' sound, taking an A for alif and a capital H for the emphatic h is more readable; also using rules such as "uppercase ASCII characters transcribe emphatic letters" is an easy way to memorize the transcriptions of ح, ط, ض, ص, ظ.
2. It has to be complete and avoid ambiguities. Of course all letters of the target language have to be covered, but having many letters to transcribe leads sometimes to ambiguities: for example taking h for ه, k for ك and kh for خ are perfectly logical choices; nevertheless there is a hitch: when you need to transcribe كه you are tempted to write simply kh and this will of course produce خ instead. The solution we have given to this problem is to type a hyphen between the letters which are not considered as a 'digraph', but this is only a compromise solution: the user must constantly be aware of this problem, and this is hardly the case when you are concentrated in your text...

It is clear that these two goals are contradictory: an accurate and unambiguous transcription has to be complicated and will be difficult to read and memorize; a friendly and easily readable transcription will be full of ambiguities.

An additional problem when making a transcription is to choose between (etymo)logical, phonetic and graphical representations of characters. A typical example is the standard Ω transcription of Greek: w is chosen for letter ω , this is a purely graphical choice: the 'w' looks like an omega, but has absolutely no other relation with, neither historical nor phonetic (the letter omega represents the sound 'o' in modern Greek); b is chosen for letter β , this is an etymological choice: the Latin 'B' derives from the ancient Greek 'B', otherwise β looks quite different than 'b' and is pronounced 'v' in modern Greek; finally, x is a phonetic transcription of letter ξ ; clearly they do not bear any resemblance, and historically it is not clear (at

least to the author) why ‘x’ should be derived from ξ (their positions in the alphabet is quite different as well, and this is an argument speaking against an etymological relation between the letters).

The reader may object that this distinction between etymological, phonetic and graphical representations is not relevant for Arabic alphabet transcriptions; actually this is only partly true: take for example bh for ب, this is an etymological transcription in the sense that it reflects the standard transcription of the Indic alphabet letter which corresponds to that Sindhi letter. Also ˘ for ayn is in some sense a graphical representation: it has been chosen because it resembles the IPA transcription of the ayn, which is ʕ. For the same reason, ˆ has been chosen for the hamza with carrier (in ˆ, ˙, etc.): the hamza’s IPA transcription is ʔ.

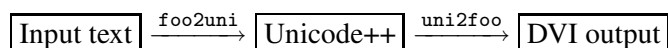
We hope to have convinced the reader that the making of a transcription is a difficult task, needing a lot of thought, compromises and tests. Once again, we would like to emphasize the fact that our transcriptions are only tentative proposals and should not be taken as standards of any kind; after all the power of Ω is that it can work with any input transcription without affecting further processing, be it contextual analysis, diacritic placement or esthetical ligaturing.

In the next section we will see how to implement a new transcription or change an existing one by writing/modifying an ΩTP file. But first some generalities on the ΩTPs used by the Arabic Ω system.

6.1.1 The ΩTPs used by the Arabic Ω system

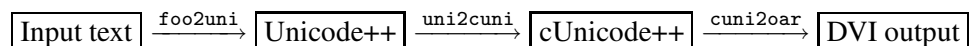
When Ω reads the text flow it places letters, digits and punctuation (whatever is not an escape or special character) into a buffer. When it encounters a special character it stops buffering and executes one after the other all currently active ΩTPs on the buffer. In theory, ΩTPs could be used to arbitrarily send character combinations to other combinations: one could very well imagine an ΩTP sending the string "Yannis" to "John" and "John" to "Yannis", or "Microsoft Word" to "μῆξΙΑε"; nevertheless, such an ΩTP would not be of general use...

Our development has mainly been focused in building ΩTPs in accordance to the following scheme:



where foo2uni sends text encoded in an arbitrary encoding into Unicode++ (Unicode++ is Unicode extended for the needs of Ω and typography), and uni2foo converts Unicode++-encoded data into the encoding of the output font. By this method we are able to keep completely separate input encoding and font encoding.

In the case of Arabic things are slightly more complicated since an additional step is needed: contextual analysis. This is where our scheme proves to be extremely efficient: by performing contextual analysis on the level of Unicode++, and hence obtaining the following new scheme:



we still remain independent of both the input and the font encoding. This means that if we need to adapt Ω to a new Arabic encoding we only need to indicate which code position corresponds to which Unicode character, and, on the other hand, if we want to adapt a new font to Ω , we only need to indicate which font position corresponds to which contextual form of which character, in cUnicode++.

In the next section we will partly describe the syntax of Ω TP files by giving examples of `foo2uni` cases.

6.2 Implementing a Transcription

The Ω TP files we will need for input encoding \rightarrow Unicode++ transformations use only part of the syntax of Ω TP files.² Such an Ω TP file is of the following form:

```
input: 1;
output: 2;
expressions:
...
...
```

where `input: 1; output: 2;` means that input is 8-bit while output is 16-bit, and `...` are lines of the following form:

```
before => after ;
```

where `before` is an expression before the transformation, and `after` after it. For example,

```
`a' => "o" ;
```

will transform all ‘a’s in the file into ‘o’s.

How do we describe characters and strings? On the left side of `=>` we can only put separate characters: they can be written either as “grave accent+ASCII character+apostrophe” or as `@"XYZT` where `XYZT` are hexadecimal digits: in this case we are not restricted to ASCII characters. The latter syntax can also be used on the right side. For example,

```
`i`j' => @"0133 ;
@"008E => @"00E9 ;
```

will send the string ‘ij’ to the Unicode++ character representing the Dutch ij ligature, and the 8-bit code 8E (a Macintosh ‘e’ with acute accent) to the Unicode++ character 00E9 (which is the Unicode ‘e’ with acute accent).

On the right side of `=>` we can also write complete strings, possibly containing Ω commands, which will be forwarded to the next Ω TP or to the typesetting engine of Ω . For example,

²The `uni2cuni` Ω TP file already needs more complicated constructions.

```
`~' => "\penalty10000" ;
```

sends the tilde character to the \TeX command of infinite penalty.³ We can also use ranges on the left side: for example, ``a' - `k'` means “all characters between a and k”.

By using parentheses and the vertical bar on the left side, we obtain the Boolean ‘or’ operator:

```
(`E' | `e') => ;
```

for example, will send both uppercase and lowercase letters ‘e’ to nothing (a transformation which would leave Percé’s book *La disparition* unchanged⁴).

This operator becomes even more useful by the fact that we can use on the right side the exact character matched on the left side: the commands `\1`, `\2`, ... , `\9` used on the right side stand for the first, second, ..., ninth character matched on the left side. For example:

```
`c' (`a' | `e' | `i' | `o' | `u') `t' => "m" \1 "p" ;
```

will send *cat*, *cet*, *cit*, *cot*, *cut* respectively to *map*, *mep*, *mip*, *mop*, *mup*.

We can go even further: Ω TP syntax allows us to add or subtract a fixed offset to the characters matched on the left side. For example:

```
`a' - `z' => #(\1 - @"0020) ;
```

will subtract **20** from the code position of the character found on the left side. The characters on the left side being precisely lowercase letters, this offset will turn them into uppercase ones.

6.2.1 Examples

The beginning of the Ω TP `7arb2uni`, used to send the ASCII transcription of Arabic to Unicode++, described in 5.2.1, to Unicode++, looks like this:

```
input: 1;  
output: 2;
```

expressions:

```
`L' `L' `a' `h'    => @"FDF2 ;  
`S' `L' `h'       => @"FDFA ;
```

³By this we obtain the same result as in \TeX but without turning tilde into an active character, a fact that \TeX users will surely appreciate.

⁴Although there are rumors that there is a single ‘e’ in that book... The authors were not able to find it yet.

Finally we send the triple hyphen to an m-dash ‘—’ and the single hyphen to nothing: its purpose is to prevent combinations of letters to be interpreted as digraphs: when reading k-h, Ω will not match it with kh: it will first match k with letter kaf, then send the hyphen to the vacuum of non-existence and when arriving to the h the k will already be matched so that it is too late to construct a kh digraph.

The period at the beginning of the last line is part of the Ω TP syntax we have not seen yet: it means ‘any character’. Since this is the last line of the file, we can interpret it rather like ‘any still not matched character’. This line simply sends any character not yet matched to itself.

6.3 Wrapping it up

Once the Ω TP file has been written or modified, one only needs to compile it (by using the `otp2ocp` utility) and place it where Ω expects to find it. On the \LaTeX command level, Ω TPs are loaded via the `\ocp` command, in a way similar to fonts: to load the file `foo2uni` one will write

```
\ocp\FooUni=foo2uni
```

Of course this is preferably done inside a \LaTeX package or style file: the final user should not need to deal with or understand this kind of code. Once the Ω TPs are loaded they are combined into lists. In this way we can push or pop simultaneously Ω TPs on/from a stack. This is useful because a language switch usually requires several Ω TPs to be changed at once. To define Ω TP lists we use the following syntax:

```
\ocplist\ArabicOCP=
\addbeforeocplist 100 \ArabUni
\addbeforeocplist 200 \UniCUni
\addbeforeocplist 300 \CUniArab
\nullocplist
```

The numbers (100, 200, 300) allow us to introduce additional Ω TPs, if necessary, between the already defined ones. Finally, to activate/desactivate an Ω TP list, we use the commands `\pushocplist` (followed by the name of the Ω TP list) and `\popocplist`. To take a real life example,

```
\ocp\ArabUni=7arb2uni
\ocp\UniCUni=uni2cuni
\ocp\CUniArab=cuni2oar
\ocplist\ArabicOCP=
\addbeforeocplist 100 \ArabUni
\addbeforeocplist 200 \UniCUni
\addbeforeocplist 300 \CUniArab
\nullocplist
\pushocplist\ArabicOCP
```

is sufficient to load all Ω TPs necessary for typesetting in the Arabic language.

7 Availability and Further Information

The Ω system is entirely in the public domain. It can be obtained from any CTAN server. The latest information on Ω and its Arabic system can be found on the Ω server:

`http://www.ens.fr/omega`

courtesy of the École Normale Supérieure de Paris.

8 Samples

Starting from next page, a few samples (Arabic, Berber, Sindhi). For these examples we have switched the background language to Arabic, so that even page numbers are in Arabic.

٨.١ أطفال الغابة

كان لأحد الملوك القدماء أخت تعيش معه في قصره، بعد أن ماتت زوجته، وتركت له من الأولاد ثلاثة: أميرين وأميرة. وقد ازداد حب الملك لأولاده، بعد وفاة والدتهم الملكة، وأحبهم حباً كثيراً؛ ليعوضهم ما فقدوه من عطف أمهم وحبها لهم، وتفكيرها فيهم؛ فكان يسأل عنهم كلما حضر، ويفكر فيهم كلما دخل، ويوصي بهم كلما خرج، ويطلبهم كلما جلس لتناول طعام الإفطار أو الغداء أو الشاي أو العشاء.

محنة أخيها لأولاده، وصممت فيما بينها وبين نفسها أن تعمل سراً كل وسيلة ممكنة لإبعادهم عن أبيهم والتخلص منهم.

وفي يوم من الأيام كان الأميران يلعبان مع أختهما الأميرة في حدائق القصر بعد خروج الملك، فشوقتهم عندهم وحببت إليهم الذهاب معها إلى الغابة للعب فيها، ووعدتهم أن تريحهم أشياء جميلة وألعاباً لذيذة سارة تحت الأشجار هناك.

فصدق الأميران والأميرة ما قالته عندهم، ولم يعرفوا ما تخفيه عنهم من الشر، وذهبوا معها للعب والرياضة في الغابة، ومشاهدة الأشياء الجميلة فيها، ورؤية الألعاب الغريبة تحت أشجارها.

وقد شعر الأطفال بسرور كثير عند ماخرجوا مع عندهم لهذه الرحلة. وأخذوا يمشون معها في الغابة حتى وصلوا إلى وسطها، فأحسوا بالتعب الشديد، وطهرت علاماته في مشيتهم، وعلى وجوههم بعد هذه الرحلة الطويلة المتعبة التي لم يجربوها من قبل. ولما شعرت العمة بشدة تعبهم، قالت لهم: ناموا هنا تحت هذه الشجرة حتى تحضر الحوريات لتلعب أمامكم ألعاباً لم تروها، وستجدون في مشاهدتها كل لذة وسرور.

٨.٢ الال ي وسقدش ن يضريس ن Ω دتامازيغت

ا د نسكن س وايس يف وسقدش ن Ω ي تيرا س توتلايت تامازيغت، اما س تيفيناغ، اما س يسكيلن يلاطائين. نويد تامازيغت ام، توتلايت يدرن (يتوارون س تيفيناغ تينين) : يزمر ومدان اد يسدو يال تيغورا ن وسودس ن تيرا، ي واراتن وسنانن، يتكنيكن نع ي ويد ن تسيكلا، ام ويد سخدامن ي وسمسارون تفرانسيست.

Ω، د امسلاي ن وسميهل ي وسودس ن تيرا. امواكن نزا، د اين ي د يتاكن يزويان وار تاقرارا ي وسقدش د وسيهرو، ماشا يسفك اد يلمد وقداش كرا توسنيوين. نونزاس، تارايت ن ولادا، نزمرا ت نسيفس س وسقدش ن يناقراون ن ورماس ن تيرا، يسغزانن ن وسميهشل ن واراتن، يتواسن اطاس (ويد يتنوزون، سراين غف ومدان، ويد يزمرن اد سخدامن تازمرت تاسمسيراوت ن كرا يناقراون يمهلائن ام ويد ن، Macintosh, Windows, Unix.

تاناضا تامزواروت ن Ω — غاس تين اي يتالاسن يسم ن Ω —، وس تلي اقرودم ي وقداش. ام ق يمسلاين ن وسميهل الك، اد يارو ومدان اهيل، دفير، ا ت يسفسو اكن ا ت يسوغال س انقال ن تماشيتن. دي Ω، اهيل د ارا ن وضريس (نطضن غورس كرا ن تسونضيوين ي وسيوني د تغسا تامزولت). اسفسو، د اسلكم ن واهيل Ω؛ انقال ن تماشيتن ارا د يفعن، د وين، ي د اقلام ن وسبتر اي يتوسودسن، يقيمد يميرون وسمسارو.

اكالايا، يزمر ا ت ياف يفرغ وين ينومن يسقداش ن يضريس غف، Macintosh, Windows، د ويض، ي دق اضريس ا د يفع دي تسمساروت اكن يلا ق وديل [اناقراوا يتواسن س يسميس يميوزيل، س تقليزيت «wysiwyg»، يشسغلاض كرا : اضريس ارا د يسوق وقداش، اد يلي غاس س تسدي ومي يساوض وقديل : اسقموض ارا د تسوق تسمساروت، يسمر اد يلي يوسر كرا.]

يوانگن اد يقيم وسقدش سراي ف ومدان، يزمر اد يسخدم اسغزان يتواسين د الين ي ورماس. تاغسا تامزولت ن وارا (يغفاون، تيفولا، تيسدارين، تيزميلين تينادان، تيميتار تينموداق، اسمل ن تكتابين)، ا ت يسيعال سي تيونيت ن وسغزانتي غر تسونضويون ن Ω. يميز، Ω، اد يسفسو انقالتني ا د يسوفغ اضريس يوقزن تاغسا تامزولت تامزاروت، ماشا تيرائيس اد يلينت ولاغنت وقار.

۸.۳ ڪتين ڪر مورثيا جڏهن

تنهن ڪري اسان کي پنهنجي ذهنن کي سجاڳ رکڻو پوندو ۽ پنهنجي جدوجهد ۾ ڏاهپ پيدا ڪرڻي. اهو به معلوم ڪرڻو پوندو ته سنڌ ۾ هر آئي وقت ڇا ڇا ٿي رهيو آهي ۽ دشمن اسان جي ۽ اسان جي جدوجهد جي ڪلاف ڪهڙا ڪهڙا گهڙا گهڙي رهيو آهي. اسان کي اها به خبر هجڻ گهرجي ته اسان جي آس پاس ۽ يسگردائي ۾ ڇا ڇا ٿي رهيو آهي. هندستان ۾ ڇا ٿي رهيو آهي، افغانستان ۾ ڇا ٿي رهيو آهي. عراق ۽ ايران ۾ ڇا ٿي رهيو آهي ۽ آمريڪا ۽ سوويت يونين ڇا ڇا سوچي رهيا آهن. جڏهن اسان سڄي دنيا جي سياست تي ۽ سڄي دنيا جي جدوجهد تي ۽ سڄي دنيا جي تبديلين تي نظر رکنداسون ۽ انهن تبديلين جي اثرن کي پنهنجي ملڪ، قوم ۽ عوام تي پوندي ڏسنداسون ته انهن تبديلين مان ڪهڙا منفي ۽ ڪهڙا مثبت اثر آهن. تڏهن ئي اسان پنهنجي جدوجهد کي بهتر به ڪري سگهنداسون ته چوٽڪاري وارو حل به ڳولي وينداسون.

رڳو ڳالهيو ڪندي ۽ نعرن هڻندي اسان جي قوم چاليه سال پيڙائون ۽ عذاب پوڳيا آهن ۽ انهن نعرن اسان جي قوم لاءِ وڌيڪ پيڙائون ۽ عذاب نازل ڪيا آهن. جيڪڏهن اسان ۾ اڄ قوم جي اميد پيدا ٿي آهي ته اها اسان جي عمل ۽ اسان جي بي لوٽ جدوجهد جي ڪري پيدا ٿي آهي ۽ ماڻهو اسان ڏانهن واجهائي رهيا آهن. ته اسان ئي آهيون جيڪي ڪجهه نه ڪجهه ڪنداسون. پر اسان کي ڏسڻو آهي ته دنيا جي اندر ڇا ٿي رهيو آهي ۽ اسان جو دشمن ڪيئن حالتن کي پنهنجن مفادن ۾ ڪتب آڻڻ جي ڪوشش ڪري رهيو آهي، انهيءَ جي لاءِ ضروري آهي ته اسان پاڻ ۾ ڏاهپ پيدا ڪريون ۽ پاڻ ۾ ڄاڻ جو هڪ وسيع خزانو پيدا ڪريون. ۽ اسان موجوده صورت حال کي سمجهڻ لاءِ روزمره جي ميڊيا ۽ دنيا جي اندر ٿيندڙ ڪاروائين تي گهري نظر رکون ته دشمن ۽ جارحيت پسند قوتون ۽ اسان تي قابض قوتون دنيا جي اندر ٿيندڙ تبديلين کي سندن حق ۾ ۽ سندن مفادن جي حق ۾، سنڌ تي جارحيت قائم رکڻ جي حق ۾، سنڌ کي مستقل قبضي ۾ ڪرڻ جي حق ۾ ڪيئن ڪتب آڻي رهيو آهن.