

Amazigh Language Desktop Converter

Fadoua Ataa Allah Jamal Frain Youssef Ait Ouguengay
CEISIC, IRCAM, Morocco
{ataaallah, frain, ouguengay}@ircam.ma

Abstract— Before the creation of the Royal Institute of Amazigh Culture, the Amazigh language has not have the ability to take the advantages of information and communication technologies, even to be written in its native writing system "Tifinaghe". To overcome these limitations, the Amazigh language has undergone a process of standardization and integration into information and communication technologies. This process is passing through several stages, mainly the encoding stage and the development of appropriate standards for the keyboard layout, in addition to the stage of computational linguistics. In this context and in the aim to save the Amazigh cultural heritage, a desktop converter, allowing the ANSI-Unicode transition and Arabic-Latin-Tifinaghe transliteration, is developed.

Keywords—Amazigh language; transcoding; transliteration; ANSI; Unicode; Tifinaghe; Arabic and Latin scripts .

I. INTRODUCTION

The integration of Amazigh language into Information and Communication Technologies (ICT) has become a necessity to promote the Amazigh language. Nevertheless, this integration was confronted by several challenges, especially those related to language planning and computer science standardization.

To let the Amazigh language supporting and conveying knowledge, firstly, a writing form and an alphabetic system have been established. Secondly, based on a linguistic description of the most widely spoken varieties of Amazigh language, a spelling system has been stabilized [1]. Then, a stage of character encoding has been undertaken. However, the difficulty in these steps is to achieve generic solutions in limited time to allow the integration of Amazigh into the Moroccan educational system in 2003. Thus, the native Amazigh writing system encoding went through two steps: ANSI then Unicode encoding [2].

Furthermore, the promotion of the Amazigh culture implies the maintenance and the conservation of the literary heritage, and the diffusion of Tifinaghe script on all media. To this end, an Amazigh desktop converter is developed. This later is based on character transcoder allowing the ANSI-Unicode conversion, and an Arabic-Latin-Tifinaghe transliterator allowing the transmutation of Amazigh document content from the Arabic or Latin alphabets to Tifinaghe.

The remainder of this paper consists of three sections. In Section 2, we present a history overview and the writing systems of the Amazigh language. Then, we describe, in Section 3, the elaborated tool. Finally, in Section 4, we conclude and present some perspectives.

II. AMAZIGH LANGUAGE

A. Amazigh language history

Amazigh is the native language of North Africa. It is also known by the name of "Berber", and the local name "Tamazight". This language is present from Morocco to Egypt passing through Algeria, Tunisia, Libya, Niger and Mali (cf. Fig. 1). It was spoken by tens of millions of people as non-standardized dialects [3].

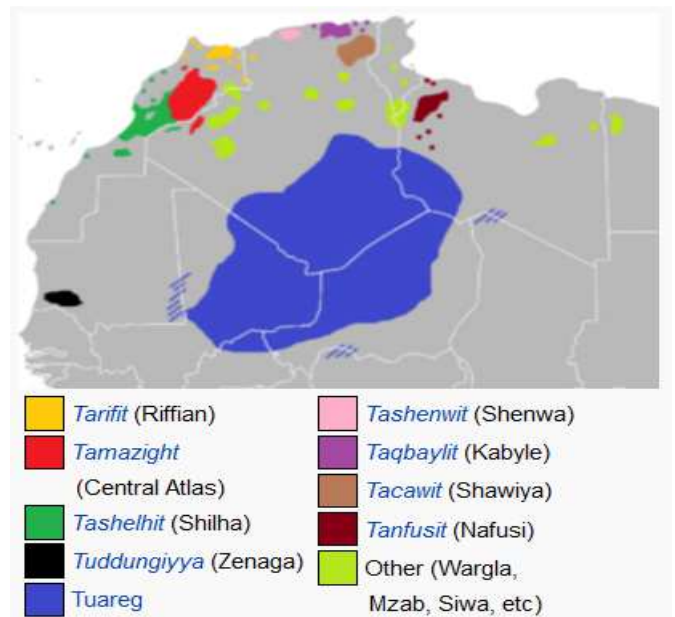


Fig. 1. Amazighophony Map¹

In Morocco, there are three main varieties of the Amazigh language: Tarifite in the North; Tamazight in the Center, the Middle Atlas and a part of High Atlas; and Tachelhite in the South, South-west of High Atlas, the Anti-Atlas and Sous. These varieties were primarily employed in oral communication. However, in order to preserve the Amazigh language, it was important to transit from orality to literacy and to upgrade the language, then to integrate the Amazigh language into the information and communication technologies.

Although the Amazigh language was primarily an oral tradition, the Amazigh language has, since antiquity, its own writing system called "Libyco-Berber" (Tifinaghe in Amazigh). This system dates back more than 40 centuries [4], [5].

¹ http://en.wikipedia.org/wiki/Berber_languages

However, the appearance form of its signs has undergone many modifications: since its inception "the Libyan" to the neo-Tifinaghe in the late sixties and Tifinaghe IRCAM-in 2001 [1].

Historically in ancient inscriptions, Amazigh was written horizontally from left to right, from right to left, vertically upward, downward or in boustrophedon (as illustrated in Fig. 2). Nonetheless, the orientation most often adopted in Amazigh language script is horizontal and from left to right, which is also adopted in IRCAM-Tifinaghe writing system.

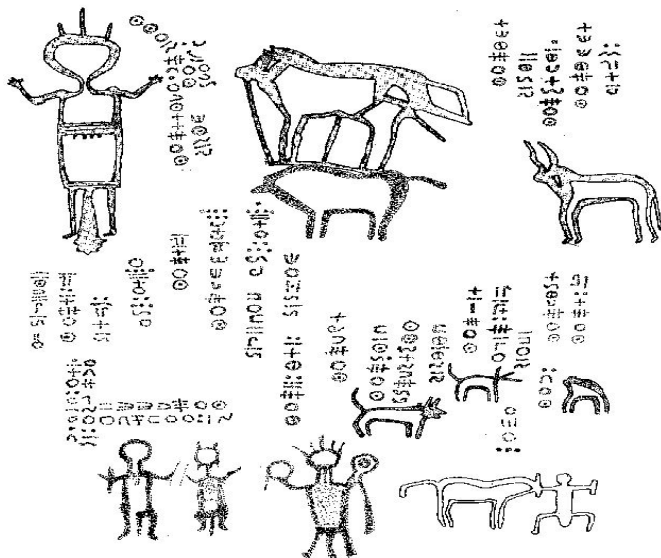


Fig. 2. Plate 9 Anou Elias, Mammanet Valley (Niger). Henri Lhote, *Oued Mammanet gravures. Les Nouvelles Editions Africaines. 1979*

B. Tifinaghe-IRCAM graphical system

Since 2003, Tifinaghe-IRCAM has become the official graphic system for writing Amazigh in Morocco [1]. This system contains:

- 27 consonants including: the labials (ⵍ, ⵍⵍ, ⵍⵍⵍ), dentals (ⵏ, ⵏⵏ, ⵏⵏⵏ, ⵏⵏⵏⵏ), alveolars (ⵙ, ⵙⵙ, ⵙⵙⵙ, ⵙⵙⵙⵙ), palatals (ⵛ, ⵛⵛ, ⵛⵛⵛ), velar (ⵝ, ⵝⵝ, ⵝⵝⵝ), labiovelars (ⵞ, ⵞⵞ, ⵞⵞⵞ), uvulars (ⵟ, ⵟⵟ, ⵟⵟⵟ), pharyngeals (ⵠ, ⵠⵠ) and the laryngeal (ⵡ);
- 2 semi-consonants: ⵣ and ⵢ;
- 4 vowels: three full vowels ⵍ, ⵎ, ⵏ and neutral vowel (or schwa) ⵓ which has a rather special status in Amazigh phonology.

No particular punctuation is known for Tifinaghe. IRCAM has recommended the use of the international symbols: “ ” (space), “:”, “;”, “,”, “.”, “?”, “!”, “...”, for punctuation markers; and the standard numeral used in Morocco (0, 1, 2, 3, 4, 5, 6, 7, 8, 9) for Tifinaghe writing system.

C. Amazigh encoding

Before the computer science standardization of Amazigh language, the Websites’ analysis showed that while there is an abundant Amazigh literature, there are very few Websites that use Amazigh language for communication. Furthermore, its presence on the web was especially based either on Arabic alphabet or the Latin alphabet enriched with some special

characters borrowed from the International Phonetic Alphabet (IPA). While, Tifinaghe characters appeared occasionally, and generally they were inserted as images. However, to allow the Amazigh people to communicate in their own language and follow at the same time the technological evolution, it was necessary to ensure a wide diffusion and a linguistic analysis of digital document content.

To this end, a digital transcription system has been created, and several efforts have been undertaken to encode Tifinaghe in Unicode/ISO 106461 [6]. Nevertheless, this process was progressing simultaneously with the integration of the Amazigh language into the Moroccan education system. So to resolve the lack of digital transcription system, while producing teaching materials in the official Amazigh script, the ANSI encoding has been used.

1) *ANSI encoding:* ANSI encoding is a slightly generic term used to refer to the standard code page on a system, usually Windows. It is more properly referred to as Windows-1252. This is essentially one of the American Standard Code for Information Interchange (ASCII) extension that includes all the ASCII characters with an additional 128 character codes. This difference is mainly due to the fact that "ANSI" encoding is 8-bit rather than 7-bit as ASCII is. But since there are many standard that differently use the 8th bit to define characters numbered from 128 to 255, we remark that the non-ASCII characters as accented ones (éËç) often display incorrectly.

2) *Tifinaghe ANSI encoding:* While expecting the Tifinaghe Unicode encoding, the Centre of Computer Studies Information Systems and Communication (CEISIC) has chosen to adapt the ANSI encoding for Tifinaghe characters as illustrated in the table of Fig. 3.

In this table, we remark the existence of the same Tifinaghe character in two different locations, which is due to the fact that each Latin character either uppercase or lowercase has been replaced by the same Tifinaghe character, since the Amazigh language does not support capital letter.

Once the encoding table is created, the ANSI encoding has allowed the manipulation of Tifinaghe digital documents, and the edition of publications and Amazigh scholar manuals. Moreover, the phonetic approach adopted by IRCAM in the choice of the graphical system and the correspondence between the Arabic script, Latin and Amazigh (see Table I) have greatly helped in accelerating the production of texts in Amazigh and facilitated the manipulation of multilingual texts with multiple writing systems. However, the scope of this IRCAM private encoding has been limited, because it is not considered as an appropriate Tifinaghe encoding.

3) *Unicode encoding:* The computer software internationalization and localization has inflected the definition of new character encoding standard that unifies all the existing character sets, and overcomes the limitations of incompatibility of old encoding standards. Thus, Unicode Consortium and the International Organization for

0000	0001	0002	0003	0004	0005	0006	0007	0008	0009	000A	000B	000C	000D	000E	000F	0010	0011	0012	0013	0014	0015	0016	0017	
NU	STX	SO	ETX	EO	ENQ	ACK	REL	BS	HT	LF	VT	FF	CR	SO	SI	DL	DC1	DC2	DC3	DC4	NAK	SYN	ETB	
0018	0019	001A	001B	001C	001D	001E	001F	0020	0021	0022	0023	0024	0025	0026	0027	0028	0029	002A	002B	002C	002D	002E	002F	
CAN	EM	SUB	ESC	FS	GS	RS	US		!	"	#	\$	%	&	'	()	*	+	,	-	.	/	
0030	0031	0032	0033	0034	0035	0036	0037	0038	0039	003A	003B	003C	003D	003E	003F	0040	0041	0042	0043	0044	0045	0046	0047	
0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?	@	⊖	⊕	⊗	⊘	⊙	⊚	⊛	
0048	0049	004A	004B	004C	004D	004E	004F	0050	0051	0052	0053	0054	0055	0056	0057	0058	0059	005A	005B	005C	005D	005E	005F	
⊜	⊝	⊞	⊟	⊠	⊡	⊢	⊣	⊤	⊥	⊦	⊧	⊨	⊩	⊪	⊫	⊬	⊭	⊮	⊯	⊰	⊱	⊲	⊳	
0060	0061	0062	0063	0064	0065	0066	0067	0068	0069	006A	006B	006C	006D	006E	006F	0070	0071	0072	0073	0074	0075	0076	0077	
⊴	⊵	⊶	⊷	⊸	⊹	⊺	⊻	⊼	⊽	⊾	⊿	Ⓚ	Ⓛ	Ⓜ	Ⓝ	Ⓞ	Ⓟ	Ⓠ	Ⓡ	Ⓢ	Ⓣ	Ⓤ	Ⓥ	
0078	0079	007A	007B	007C	007D	007E	007F	007F	20AC	20AD	20AE	20AF	20B0	20B1	20B2	20B3	20B4	20B5	20B6	20B7	20B8	20B9	20BA	
Ⓦ	Ⓧ	Ⓨ	Ⓩ	ⓐ	ⓑ	ⓓ	ⓔ	ⓕ	ⓖ	ⓗ	ⓘ	ⓙ	ⓚ	ⓛ	ⓜ	ⓝ	ⓞ	ⓟ	ⓠ	ⓡ	ⓢ	ⓣ	ⓤ	
0080	0081	0082	0083	0084	0085	0086	0087	0088	0089	008A	008B	008C	008D	008E	008F	0090	0091	0092	0093	0094	0095	0096	0097	
ⓥ	ⓦ	ⓧ	ⓨ	ⓩ	⓪	⓫	⓬	⓭	⓮	⓯	⓰	⓱	⓲	⓳	⓴	⓵	⓶	⓷	⓸	⓹	⓺	⓻	⓼	
0098	0099	009A	009B	009C	009D	009E	009F	009F	20AC	20AD	20AE	20AF	20B0	20B1	20B2	20B3	20B4	20B5	20B6	20B7	20B8	20B9	20BA	
⓽	⓿	Ⓚ	Ⓛ	Ⓜ	Ⓝ	Ⓞ	Ⓟ	Ⓠ	Ⓡ	Ⓢ	Ⓣ	Ⓤ	Ⓥ	Ⓦ	Ⓧ	Ⓨ	Ⓩ	ⓐ	ⓑ	ⓓ	ⓔ	ⓕ		
00A8	00A9	00AA	00AB	00AC	00AD	00AE	00AF	00B0	00B1	00B2	00B3	00B4	00B5	00B6	00B7	00B8	00B9	00BA	00BB	00BC	00BD	00BE	00BF	
Ⓠ	Ⓡ	Ⓢ	Ⓣ	Ⓤ	Ⓥ	Ⓦ	Ⓧ	Ⓨ	Ⓩ	ⓐ	ⓑ	ⓓ	ⓔ	ⓕ	Ⓠ	Ⓡ	Ⓢ	Ⓣ	Ⓤ	Ⓥ	Ⓦ	Ⓧ	Ⓨ	
00C0	00C1	00C2	00C3	00C4	00C5	00C6	00C7	00C8	00C9	00CA	00CB	00CC	00CD	00CE	00CF	00D0	00D1	00D2	00D3	00D4	00D5	00D6	00D7	
⓷	⓸	⓹	⓺	⓻	⓼	⓽	⓿	Ⓚ	Ⓛ	Ⓜ	Ⓝ	Ⓞ	Ⓟ	Ⓠ	Ⓡ	Ⓢ	Ⓣ	Ⓤ	Ⓥ	Ⓦ	Ⓧ	Ⓨ	Ⓩ	
00D8	00D9	00DA	00DB	00DC	00DD	00DE	00DF	00E0	00E1	00E2	00E3	00E4	00E5	00E6	00E7	00E8	00E9	00EA	00EB	00EC	00ED	00EE	00EF	
Ⓩ	ⓐ	ⓑ	ⓓ	ⓔ	ⓕ	Ⓠ	Ⓡ	Ⓢ	Ⓣ	Ⓤ	Ⓥ	Ⓦ	Ⓧ	Ⓨ	Ⓩ	ⓐ	ⓑ	ⓓ	ⓔ	ⓕ	Ⓠ	Ⓡ	Ⓢ	
00F0	00F1	00F2	00F3	00F4	00F5	00F6	00F7	00F8	00F9	00FA	00FB	00FC	00FD	00FE	00FF						
Ⓣ	Ⓤ	Ⓥ	Ⓦ	Ⓧ	Ⓨ	Ⓩ	ⓐ	ⓑ	ⓓ	ⓔ	ⓕ	Ⓠ	Ⓡ	Ⓢ	Ⓣ									

Fig. 3. Tifinaghe ANSI encoding table

Standardization (ISO) have combined their efforts to provide a universal character encoding scheme known as Unicode. This standard is designed to support the worldwide interchange, processing, and display of the technical disciplines of the modern world and the written texts of diverse languages even classical and historical ones.

Unicode can be implemented by different character encodings. The most commonly used encodings are UTF-8, UTF-16 and the now-obsolete UCS-2. UTF-8 uses one byte for any ASCII characters, which have the same code values in both UTF-8 and ASCII encoding, and up to four bytes for other characters. UCS-2 uses a 16-bit code unit (two 8-bit bytes) for each character but cannot encode every character in the current Unicode standard. UTF-16 extends UCS-2, using two 16-bit units (4 × 8 bit) to handle each of the additional characters².

4) *Tifinaghe Unicode encoding*: According to Mr Zenkour³ [7], the Unicode standard has constituted the main entrance of the Amazigh native writing system into the world of Information and Communication Technologies. Indeed thanks to this standard, nowadays, Tifinaghe is used on line based on Tifinaghe Webfonts [8] and internationalized URL [9], and it is integrated into the operating systems Linux and Windows.

Unicode covers almost all scripts in current use today, and reserves to each new writing system a set of blocks. Thus, it has reserved a block for Tifinaghe characters, where each character is currently determined by a single code point. Fig. 4 illustrates the Tifinaghe block and specifies the code points reserved for each of the four subsets of Tifinaghe characters: the basic set of IRCAM, the extended IRCAM set, other Neo-Tifinaghe letters in use, and modern Touareg letters. The first subset constitutes the set of characters chosen by IRCAM to arrange the orthography of the different Moroccan Amazigh varieties while preserving most characters of the historical Tifinaghe script. This subset is classified in accordance in the

range U+2D30...U+2D65, U+2D6F with the order specified in Tifinaghe-IRCAM alphabet. While, the Tifinaghe block is the range U+2D30...U+2D7F [6], [10]. For more details on the Unicode and the Information Technology component for adding Tifinaghe to Unicode one can see also [10].

	2D3x	2D4x	2D5x	2D6x	2D7x
0	Ⓩ	ⓐ	ⓑ	ⓓ	
1	ⓔ	ⓕ	Ⓠ	Ⓡ	
2	Ⓢ	Ⓣ	Ⓤ	Ⓥ	
3	Ⓦ	Ⓧ	Ⓨ	Ⓩ	
4	ⓐ	ⓑ	ⓓ	ⓔ	
5	ⓕ	Ⓠ	Ⓡ	Ⓢ	
6	Ⓣ	Ⓤ	Ⓥ		
7	Ⓦ	Ⓧ	Ⓨ		
8	Ⓩ	ⓐ	ⓑ		
9	ⓔ	ⓕ	Ⓠ		
A	Ⓢ	Ⓣ	Ⓤ		
B	Ⓦ	Ⓧ	Ⓨ		
C	ⓐ	ⓑ	ⓓ		
D	ⓕ	Ⓠ	Ⓡ		
E	Ⓣ	Ⓤ	Ⓥ		
F	Ⓦ	Ⓧ	Ⓨ		

- Basic Tifinaghe IRCAM
- Extended Tifinaghe IRCAM
- Other Neotifinaghe Letters
- Attested Touareg Letters
- Reserved for later Encoding

Fig. 4. Tifinaghe Unicode encoding block

² <http://en.wikipedia.org/wiki/Unicode>

³ The CEISIC's ex-director at IRCAM

5) *Amazigh language transcription*: Before adopting Tifinaghe as an official script in Morocco, like any oral language, Amazigh was writing by the graphic systems widely used in the country. Thus, the Arabic script was used for religion and rural poetry writing, while Latin supported by the International Phonetic Alphabet was used particularly by berberists since early works of missionaries.

III. AMAZIGH DESKTOP CONVERTER

Through its existence, the Amazigh language has known different forms of writing: Latin supported by the International Phonetic Alphabet, Arabic script, and Tifinaghe character based on ANSI and Unicode encoding. In the aim to allow users to read or write in a suitable form, and to save the Amazigh literature heritage in a standard unique form, a command-line converter has been developed [11]. This tool ensures an automatic shifting from one form to another. However, it has some limitations such as having menu-driven and graphical user interfaces, processing rich text format, and dealing with multilingual text especially when Amazigh language is writing with ANSI script. To overcome these shortcomings a desktop converter is developed.

A. Desktop converter technical architecture

The technical architecture of the desktop converter is based on an implementation of the Model-View-Controller (MVC) pattern based on technology .Net, which separates application data model and user interface views in separate components (cf. Fig. 5).

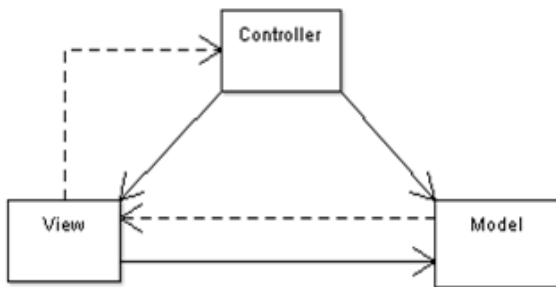


Fig. 5. Technical Architecture Diagram

B. Desktop converter functional architecture

As summarized in the diagram of Fig. 6, the desktop converter consists of two processes: transcoder and transliterator.

1) *Transcoder*: This process allows to shift the ANSI representation of Tifinaghe to Unicode representation for the content of a file in one of the following format: .txt, .rtf, .doc, and .docx (cf. Fig. 7), or the content of a text area (cf. Fig. 8).

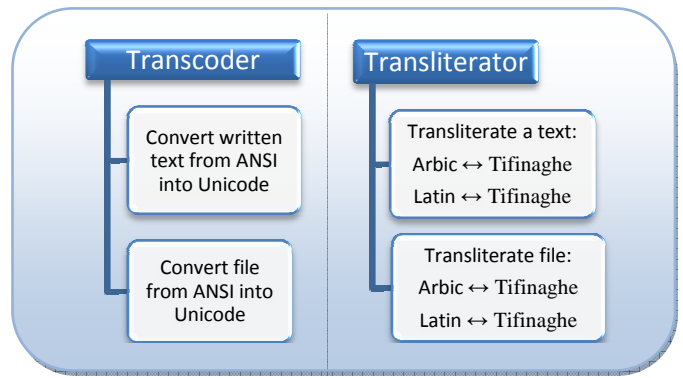


Fig. 6. Functional Architecture Diagram

In the case of file conversion, the transcoder allows browsing hard drives to choose the file to process, and choosing the input ANSI and the output Unicode fonts. Furthermore, it preserves the document layout and text formatting.

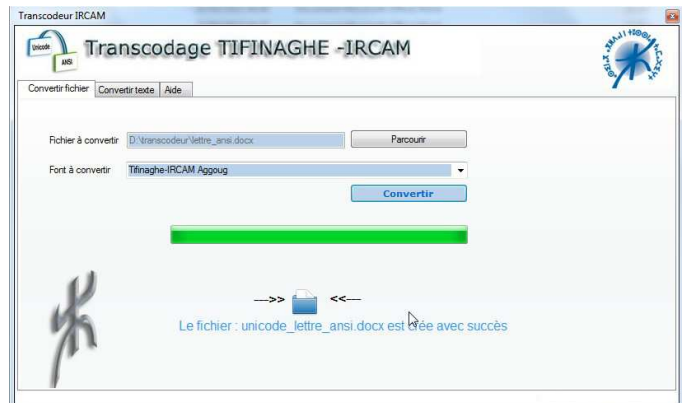


Fig. 7. Transcoder interface for converting a file

In the case of converting the content of a text area, the transcoder enables a real time encoding conversion. The input content could be pasted or typed, and the output could be copied or saved into a text file.

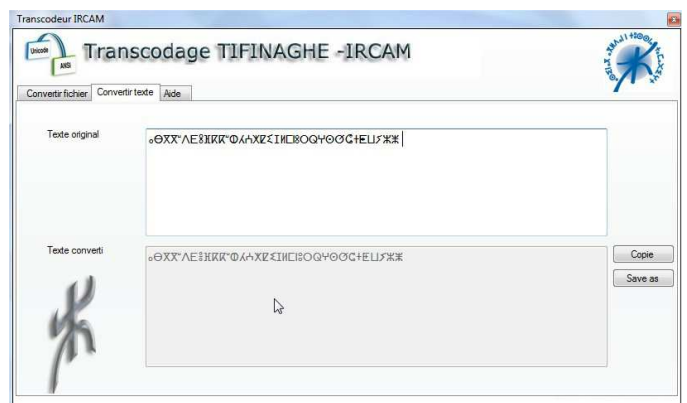


Fig. 8. Transcoder interface for the content of text area

IV. CONCLUSION

In the aim to promote the Amazigh language and to preserve its literary heritage, this paper has presented a new desktop converter enabling the ANSI-Unicode transcoding and the Arabic-Latin-Tifinaghe transliteration even in real time for typed text and rich text format for files. This work will be extended to an Amazigh toolbox that could integrate dictionary, conjugator, spelling checker and other computational linguistic tools.

REFERENCES

- [1] M. Ameer, A. Bouhjar, F. Boukhris, A. Boukouss, A. Boumalk, M. Elmedlaoui, E. M. Iazzi, H. Souifi, *Initiation à la langue amazighe*, IRCAM, 2004.
- [2] M. Moukhlis, L. Ouniam, *Bulletind'information de l'Institut Royal de la Culture Amazighe*, n° 5&6, June 2006.
- [3] S. Chaker, *Le berbère*, Actes de Les langues de France, 215-227. 2003.
- [4] M. Hachid, "Les premiers berbères", *Entre Méditerranée, Tassili et Nili*. Aix-en-Provence-Alger : Edisud-Ina-Yas, pp. 173-190.
- [5] A. Skounti, A. Lemjidi, E. M. Nami, *Tirra aux origines de l'écriture au Maroc*, IRCAM,2003.
- [6] P. Andries, *Unicode 5.0 en pratique : Codage des caractères et internationalisation des logiciels et des documents*. Dunod : Collection InfoPro, France, 2008.
- [7] L. Zenkouar, "Normes des technologies de l'information pour l'ancrage de l'écriture amazighe", in *Etudes et documents berbères*, n° 27, pp. 159-172, 2008.
- [8] P. Andries, "Demain encore plus de tifinaghes sur Internet", *Actes de colloque TICAM 2008*.
- [9] P. Andries, "Indiquer la langue, l'écriture, le pays dans des documents informatiques", *Actes de colloque TICAM 2008*.
- [10] L. Zenkouar, "L'écriture amazighe tifinaghe et Unicode", *Etudes et Documents Berberes*, 22, 2004: pp. 175-173.
- [11] F. Ataa Allah, and S. Boulaknadel, "Convertisseur pour la langue amazighe : script arabe - latin - tifinaghe", *The 2nd Symposium International sur le Traitement Automatique de la Culture Amazighe, SITACAM 2011*, 6-7 May 2011, Agadir, Morocco, 2011.
- [12] H. Stroomer, *Textes berbères des Guedmioua et Goundafa (Haut Atlas, Maroc)*, basés sur les documents de F. Corjon, J.-M. Franchi et J. Eugène, Edisud, 2001.
- [13] F. Ataa Allah, and S. Boulaknadel, "Toward computational processing of less resourced languages: Primarily experiments for Moroccan Amazigh language", in *Text Mining*. Rijeka: InTech, pp. 197-218, november 2012.
- [14] M. Chafik, *المعجم العربي الأمازيغي*, El Maarif Aljadida, Rabat, Morocco, 1993.



- b -

Fig. 9. Transliterator interface

Furthermore, the tool enables the user to set his/her proper correspondence mapping table that could be saved for another reuse (cf. Fig. 10). Once the new mapping table is saved, it will appear in the dropdown list.

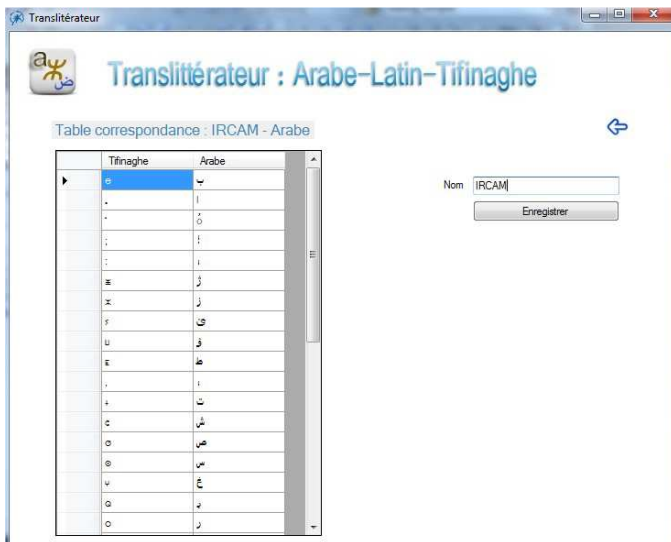


Fig. 10. Correspondence setting layout