

**ISO/IEC JTC 1/SC 2/WG 2**  
**PROPOSAL SUMMARY FORM TO ACCOMPANY SUBMISSIONS**  
**FOR ADDITIONS TO THE REPERTOIRE OF ISO/IEC 10646<sup>1</sup>**

Please fill all the sections A, B and C below.

(Please read Principles and Procedures Document for guidelines and details before filling this form.)

See <http://www.dkuug.dk/JTC1/WG2/docs/summaryform.html> for latest Form.

See <http://www.dkuug.dk/JTC1/WG2/docs/principles.html> for latest Principles and Procedures document.

See <http://www.dkuug.dk/JTC1/WG2/docs/roadmaps.html> for latest roadmaps.

**A. Administrative**

1. **Title:** Proposal to encode additional Arabic-script characters
2. Requester's name: INCITS/L2; Unicode Technical Committee; Jonathan Kew, SIL International
3. Requester type (Member body/Liaison/Individual contribution): Member; Liaison; expert contribution
4. Submission date: 2003-07-10
5. Requester's reference (if applicable): L2/02-274; L2/03-159; L2/03-168; L2/03-176; L2/03-210
6. (Choose one of the following:)  
 This is a complete proposal: Yes  
 or, More information will be provided later: \_\_\_\_\_

**B. Technical - General**

1. (Choose one of the following:)  
 a. This proposal is for a new script (set of characters): No  
 Proposed name of script: \_\_\_\_\_  
 b. The proposal is for addition of character(s) to an existing block: Yes  
 Name of the existing block: Arabic, and proposed Arabic Supplementary block
2. Number of characters in proposal: 30
3. Proposed category (see section II, Character Categories): A
4. Proposed Level of Implementation (1, 2 or 3) (see clause 14, ISO/IEC 10646-1: 2000): 2  
 Is a rationale provided for the choice? Yes  
 If Yes, reference: Includes combining marks
5. Is a repertoire including character names provided? Yes  
 a. If YES, are the names in accordance with the 'character naming guidelines  
 in Annex L of ISO/IEC 10646-1: 2000? Yes  
 b. Are the character shapes attached in a legible form suitable for review? Yes
6. Who will provide the appropriate computerized font (ordered preference: True Type, or PostScript format) for publishing the standard? Jonathan Kew, SIL International  
 If available now, identify source(s) for the font (include address, e-mail, ftp-site, etc.) and indicate the tools used:  
jonathan\_kew@sil.org  
TrueType font generated with FontLab 4.5
7. References:  
 a. Are references (to other character sets, dictionaries, descriptive texts etc.) provided? Yes  
 b. Are published examples of use (such as samples from newspapers, magazines, or other sources) of proposed characters attached? Yes
8. Special encoding issues:  
 Does the proposal address other aspects of character data processing (if applicable) such as input, presentation, sorting, searching, indexing, transliteration etc. (if yes please enclose information)?  
Yes: suggested Unicode character properties are included
9. Additional Information:

Submitters are invited to provide any additional information about Properties of the proposed Character(s) or Script that will assist in correct understanding of and correct linguistic processing of the proposed character(s) or script. Examples of such properties are: Casing information, Numeric information, Currency information, Display behaviour information such as line breaks, widths etc., Combining behaviour, Spacing behaviour, Directional behaviour, Default Collation behaviour, relevance in Mark Up contexts, Compatibility equivalence and other Unicode normalization related information. See the Unicode standard at <http://www.unicode.org> for such information on other scripts. Also see <http://www.unicode.org/Public/UNIDATA/UnicodeCharacterDatabase.html> and associated Unicode Technical Reports for information needed for consideration by the Unicode Technical Committee for inclusion in the Unicode Standard.

## C. Technical - Justification

1. Has this proposal for addition of character(s) been submitted before? If YES explain <u>(but characters approved for encoding at UTC #95, June 2003)</u>	<u>Not to WG2</u>
2. Has contact been made to members of the user community (for example: National Body, user groups of the script or characters, other experts, etc.)? If YES, with whom? <u>Local communities, linguists, NGOs working in S. Asia &amp; N. Africa</u> If YES, available relevant documents: <u>See §4 below</u>	<u>Yes</u>
3. Information on the user community for the proposed characters (for example: size, demographics, information technology use, or publishing use) is included? Reference: <u>See §3 below</u>	<u>Yes</u>
4. The context of use for the proposed characters (type of use; common or rare) Reference: <u>See §3 below</u>	<u>Various</u>
5. Are the proposed characters in current use by the user community? If YES, where? Reference: <u>South Asia, North African countries</u>	<u>Yes</u>
6. After giving due considerations to the principles in <i>Principles and Procedures document</i> (a WG 2 standing document) must the proposed characters be entirely in the BMP? If YES, is a rationale provided? If YES, reference: <u>Extensions to BMP Arabic repertoire, see §3 below</u>	<u>Yes</u> <u>Yes</u>
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)?	<u>Yes</u>
8. Can any of the proposed characters be considered a presentation form of an existing character or character sequence? If YES, is a rationale for its inclusion provided? If YES, reference: _____	<u>No</u> _____ _____
9. Can any of the proposed characters be encoded using a composed character sequence of either existing characters or other proposed characters? If YES, is a rationale for its inclusion provided? If YES, reference: _____	<u>No</u> _____ _____
10. Can any of the proposed character(s) be considered to be similar (in appearance or function) to an existing character? If YES, is a rationale for its inclusion provided? If YES, reference: <u>See §3.2 below</u>	<u>Yes</u> <u>Yes</u> _____
11. Does the proposal include use of combining characters and/or use of composite sequences (see clauses 4.12 and 4.14 in ISO/IEC 10646-1: 2000)? If YES, is a rationale for such use provided? If YES, reference: <u>Arabic-script vowels (see §3.1.2 below) are combining marks</u> Is a list of composite sequences and their corresponding glyph images (graphic symbols) provided? If YES, reference: _____	<u>Yes</u> <u>Yes</u> _____ <u>No</u> _____
12. Does the proposal contain characters with any special properties such as control function or similar semantics? If YES, describe in detail (include attachment if necessary)	<u>No</u> _____
13. Does the proposal contain any Ideographic compatibility character(s)? If YES, is the equivalent corresponding unified ideographic character(s) identified? If YES, reference: _____	<u>No</u> _____ _____

<sup>1</sup>Form number: N2352-F (Original 1994-10-14; Revised 1995-01, 1995-04, 1996-04, 1996-08, 1999-03, 2001-05, 2001-09)

### Submitter's Responsibilities

The national body or liaison organization (or any other organization or an individual) proposing new character(s) or a new script shall provide:

1. Proposed category for the script or character(s), character name(s), and description of usage.
2. Justification for the category and name(s).
3. A representative glyph(s) image on paper:  
If the proposed glyph image is similar to a glyph image of a previously encoded ISO/IEC 10646 character, then additional justification for encoding the new character shall be provided.  
**Note:** Any proposal that suggests that one or more of such variant forms is actually a distinct character requiring separate encoding, should provide detailed, printed evidence that there is actual, contrastive use of the variant form(s). It is insufficient for a proposal to claim a requirement to encode as characters in the Standard, glyphic forms which happen to occur in another character encoding that did not follow the Character-Glyph Model that guides the choice of appropriate characters for encoding in ISO/IEC 10646.  
**Note:** WG 2 has resolved in Resolution M38.12 not to add any more Arabic presentation forms to the standard and suggests users to employ appropriate input methods, rendering and font technologies to meet the user requirements.
4. Mappings to accepted sources, for example, other standards, dictionaries, accessible published materials.
5. Computerized/camera-ready font:  
Prior to the preparation of the final text of the next amendment or version of the standard a suitable computerized font (camera-ready font) will be needed. Camera-ready copy is mandatory for final text of any pDAMs before the next revision. Ordered preference of the fonts is True Type or PostScript format. The minimum design resolution for the font is 96 by 96 dots matrix, for presentation at or near 22 points in print size.
6. List of all the parties consulted.
7. Equivalent glyph images:  
If the submission intends using composite sequences of proposed or existing combining and non-combining characters, a list consisting of each composite sequence and its corresponding glyph image shall be provided to better understand the intended use.
8. Compatibility equivalents:  
If the submission includes compatibility ideographic characters, identify the equivalent unified CJK Ideograph character(s).
9. Any additional information that will assist in correct understanding of the different characteristics and linguistic processing of the proposed character(s) or script.

## 1. Proposed character additions (shaded cells)

	060	061	062	063	064	065	075	076	077
0	ا	آ		ذ	-	َ	ي	ف	
1	س	ع	ء	ر	ف	ّ	ث	ف	
2	م	ح	ا	ز	ق	ّ	پ	خی	
3	ص	ض	أ	س	ك	ّ	پ	شی	
4		ل	ؤ	ش	ل	ّ	ن	کی	
5		ط	إ	ص	م	ّ	ب	م	
6			ئ	ض	ن	ا	ن	م	
7			ا	ط	ه	ُ	ت	ن	
8			ب	ظ	و	ّ	چ	ن	
9			ة	ع	ی		ط	ن	
A			ت	غ	ي	ّ	د		
B		؛	ث		ّ	ّ	ر		
C	،		ج		ّ	ّ	ش		
D	،		ح		ّ		ت		
E	م	ٴ	خ		ّ		ت		
F	ع	؟	د		ّ		غ		



## 2. Names list for character additions

---

Listing of file *ArabicUTC95-names.txt*

---

```
; additions to the 0600 Arabic block
;
@@ 0600 Arabic 06FF
;
@ Punctuation
061E ARABIC TRIPLE DOT PUNCTUATION MARK
;
@ Other combining marks
;
; See Everson/Pournader's proposal L2/03-133R or N2581R2
; 0659 ARABIC ZWARAKAY
; * Pashto
;
065A ARABIC VOWEL SIGN SMALL V ABOVE
* African languages
065B ARABIC VOWEL SIGN INVERTED SMALL V ABOVE
* African languages
065C ARABIC VOWEL SIGN DOT BELOW
* African languages
;
; new Arabic Supplementary block
;
@@ 0750 Arabic Supplementary 077F
;
@ Extended Arabic letters
@+ These are primarily used in Arabic-script orthographies of African languages.
0750 ARABIC LETTER BEH WITH THREE DOTS HORIZONTALLY BELOW
0751 ARABIC LETTER BEH WITH DOT BELOW AND THREE DOTS ABOVE
0752 ARABIC LETTER BEH WITH THREE DOTS POINTING UPWARDS BELOW
0753 ARABIC LETTER BEH WITH THREE DOTS POINTING UPWARDS BELOW AND TWO DOTS ABOVE
0754 ARABIC LETTER BEH WITH TWO DOTS BELOW AND DOT ABOVE
0755 ARABIC LETTER BEH WITH INVERTED SMALL V BELOW
0756 ARABIC LETTER BEH WITH SMALL V
0757 ARABIC LETTER HAH WITH TWO DOTS ABOVE
0758 ARABIC LETTER HAH WITH THREE DOTS POINTING UPWARDS BELOW
0759 ARABIC LETTER DAL WITH TWO DOTS VERTICALLY BELOW AND SMALL TAH
* Saraiki
075A ARABIC LETTER DAL WITH INVERTED SMALL V BELOW
075B ARABIC LETTER REH WITH STROKE THROUGH
075C ARABIC LETTER SEEN WITH FOUR DOTS ABOVE
* Shina
075D ARABIC LETTER AIN WITH TWO DOTS ABOVE
075E ARABIC LETTER AIN WITH THREE DOTS POINTING DOWNWARDS ABOVE
075F ARABIC LETTER AIN WITH TWO DOTS VERTICALLY ABOVE
0760 ARABIC LETTER FEH WITH TWO DOTS BELOW
0761 ARABIC LETTER FEH WITH THREE DOTS POINTING UPWARDS BELOW
0762 ARABIC LETTER KEHEH WITH DOT ABOVE
* old Malay, preferred to 06AC
x (kaf with dot above - 06AC)
0763 ARABIC LETTER KEHEH WITH THREE DOTS ABOVE
* Moroccan Arabic, Amazigh
x (arabic letter ng - 06AD)
0764 ARABIC LETTER KEHEH WITH THREE DOTS POINTING UPWARDS BELOW
0765 ARABIC LETTER MEEM WITH DOT ABOVE
0766 ARABIC LETTER MEEM WITH DOT BELOW
* Maba
0767 ARABIC LETTER NOON WITH TWO DOTS BELOW AND DOT ABOVE
0768 ARABIC LETTER NOON WITH DOT ABOVE AND SMALL TAH
* Saraiki, Pathwari
0769 ARABIC LETTER NOON WITH SMALL V
* Gojri
```

---

### 3. Discussion of the proposed additions

The proposed characters can be considered in several categories: a collection of letters that have been used to extend Arabic script for various African languages; two *gaf* letters that merit individual discussion; several additional letters used in South Asian languages; and a punctuation character traditionally used when writing African languages in Arabic script.

#### 3.1 North African extensions to Arabic script

The principal source of information concerning many of the characters proposed for encoding here is Chtatou (1992). By way of general information on the languages involved, Chtatou writes:

##### *Fulfulde*

This language is also known by other names, mainly: Fula, Peul, Pular and Pulaar, and is spoken in many countries ... Fulfulde covers a larger geographical area than any other African language and, as a result, is considered as one of the principal languages of Africa: it is spoken by between 12 and 15 million people.

In the field of education, Fulfulde is used as a language of instruction in Guinea and Nigeria and there are pilot-projects considering its use in the Gambia, Mauritania and Niger. In addition, it is used as a language of instruction at secondary school level for the first two years in Guinea and at the university level in Nigeria. The language in question is used by the press and, as a result, two monthlies appear in it: one in Mali with a circulation of 500 copies and another in Niger with 3000 copies.

##### *Hausa*

Geographically speaking, Hausa is less spread out than Fulfulde; it is one of the principal languages of Africa, spoken by 40 million people ... An intense literary activity involving this language has been signaled mainly in the field of fiction and poetry. At the same time, the mass media have started using it more and more; in fact, there are in this country [Nigeria] 3 weekly newspapers published in this language. As for education, Hausa is optional at primary school level in Nigeria and obligatory at the secondary school level. In Niger, on the other hand, there is a pilot project to use Hausa as a language of instruction in secondary as well as in higher education.

##### *Songhoy*

This language is known as Zarma in Niger and Nigeria. Dendi, which is spoken in Benin, is considered as the same language as Songhoy because of mutual intelligibility.

Songhoy is the second language in Niger, in terms of the number of speakers ... in Mali it is the third. A lot of research on this language is being conducted in various countries ... in order to use it as a means to undertake literacy programmes in these areas. As for education, Songhoy is utilized as a language of instruction in the first three years of primary school. It is also taught at the university as an optional subject.

This language is also used in the mass media with a monthly newspaper selling 500 copies in Benin and three other monthlies in Niger with a circulation of 3,000 copies for one and 1,000 copies each of the remaining two.

##### *Wolof*

This language is currently spoken in Senegal, the Gambia and Mauritania. It is considered in all three countries as a community language. ... In all three countries where it is used, there are several pilot projects to use it as a language of instruction at the primary level of education; ... As for literacy, there are 43,000 people in the Gambia who have been initiated to this language, which is written in the Arabic script. There is a similar activity in Senegal, where the authorities use television to teach people to write the language.

Chtatou then goes on to discuss the desire in several African countries to develop Arabic-script orthographies for these languages:












...in the wake of decolonization, the interest in these languages was revived and governments set out to give them the status they deserve on a national level. ... These languages were also used to promote much-needed literacy programmes for people of different ages. As for some countries of Muslim Africa, their top priority was to devise for their languages an Arabic script in which they can be written so that the language can be used for the purpose of promoting literacy...

Keen on the development of such a script for their languages, African governments approached such international organizations as UNESCO and ISESCO ... The first workshops were organized by UNESCO through BRED (Bureau Régional d'Éducation pour l'Afrique), and later on ISESCO joined in the effort and sponsored other workshops on the same topic.

The bulk of Chtatou's work consists of reports on the writing conventions adopted at these workshops during the late 1980s. The charts he shows of existing and proposed transcription systems (either already in use or proposed as a result of the effort to standardize the writing systems) include a number of Arabic-script letters that are not supported in Unicode. This proposal, therefore, aims to extend the UCS repertoire to include the African characters documented in this study, as well as additional characters found in other African-language publications (see References).

### 3.1.1 Base (consonant) characters

The following 20 extended Arabic letters are found in the sources for African languages. All these proposed characters are of General Category Lo; Combining Class 0; Bidi Type AL. The suggested codepoints are those approved by the UTC in June 2003. (Gaps in the sequence of codepoints will be filled with characters discussed in subsequent sections of this document.) These characters are placed in a proposed new *Arabic Supplementary* block at U+0750, leaving the remaining spaces in the U+0600 block for Arabic-script diacritics, punctuation, etc.

<i>Glyph</i>	<i>Code</i>	<i>Character name</i>	<i>Shaping</i>	<i>See figures</i>
	0750	ARABIC LETTER BEH WITH THREE DOTS HORIZONTALLY BELOW	BEH	3, 7, 9
	0751	ARABIC LETTER BEH WITH DOT BELOW AND THREE DOTS ABOVE	BEH	4, 7, 9
	0752	ARABIC LETTER BEH WITH THREE DOTS POINTING UPWARDS BELOW	BEH	4, 7, 8, 10, 11
	0753	ARABIC LETTER BEH WITH THREE DOTS POINTING UPWARDS BELOW AND TWO DOTS ABOVE	BEH	5
	0754	ARABIC LETTER BEH WITH TWO DOTS BELOW AND DOT ABOVE	BEH	7, 8, 9
	0755	ARABIC LETTER BEH WITH INVERTED SMALL V BELOW	BEH	10
	0756	ARABIC LETTER BEH WITH SMALL V	BEH	10
	0757	ARABIC LETTER HAH WITH TWO DOTS ABOVE	HAH	2, 3, 7, 8
	0758	ARABIC LETTER HAH WITH THREE DOTS POINTING UPWARDS BELOW	HAH	8
	075A	ARABIC LETTER DAL WITH INVERTED SMALL V BELOW	DAL	11
	075B	ARABIC LETTER REH WITH STROKE THROUGH	REH	11
	075D	ARABIC LETTER AIN WITH TWO DOTS ABOVE	AIN	2, 3, 7, 9, 10, 11, 12, 13
	075E	ARABIC LETTER AIN WITH THREE DOTS POINTING DOWNWARDS ABOVE	AIN	6, 8
	075F	ARABIC LETTER AIN WITH TWO DOTS VERTICALLY ABOVE	AIN	8

ف	0760	ARABIC LETTER FEH WITH TWO DOTS BELOW	FEH	3, 7, 8, 9
ف	0761	ARABIC LETTER FEH WITH THREE DOTS POINTING UPWARDS BELOW	FEH	11
ك	0764	ARABIC LETTER KEHEH WITH THREE DOTS POINTING UPWARDS BELOW	GAF	11
م	0765	ARABIC LETTER MEEM WITH DOT ABOVE	MEEM	3, 7
م	0766	ARABIC LETTER MEEM WITH DOT BELOW	MEEM	12, 13
ن	0767	ARABIC LETTER NOON WITH TWO DOTS BELOW AND DOT ABOVE	NOON	1, 11, 13

3.1.2
Vowel signs

Several new signs have been used in African languages to represent vowel sounds not present in standard Arabic. The following three signs are proposed for encoding as combining characters in the Arabic block. The codepoints shown here are those approved by the UTC in June 2003. (Note that the ARABIC ZWARAKAY, as proposed in L2/03-133 (Everson & Pournader 2003), is to be encoded at U+0659.)

Glyph	Code	Character name	GC	CC	Bidi	See figures
◌ِ	065A	ARABIC VOWEL SIGN SMALL V ABOVE	Mn	30	NSM	12, 13
◌َ	065B	ARABIC VOWEL SIGN INVERTED SMALL V ABOVE	Mn	30	NSM	12, 13
◌ِ◌ْ	065C	ARABIC VOWEL SIGN DOT BELOW	Mn	32	NSM	14, 15, 16, 17

3.1.3
Samples showing African characters

The consonant and vowel characters listed above are illustrated in the following samples from Chtatou (1992) and other sources.

TABLE I

COMPARATIVE TABLE OF ARABIC AND SONGHOY CHARACTERS

Similar characters in Arabic and Songhoy		Specific Arabic characters		*Specific Songhoy characters	
Latin characters	Arabic characters	Latin characters	Arabic characters	Latin characters	Arabic characters
a	ا	θ	ث	p	ف
b	ب	ɸ	ح	ٴ	ن
t	ت	kh	خ	ny	ي

(19) cf. UNESCO/BREDA, Rapport général du séminaire atelier sur l'élaboration d'un système unifié de transcription du Songhoy en caractères arabes, du 14 au 19 mars 1987, Bamako, p.6.

**Figure 1:** Chtatou (1992), page 28. This also shows a BEH WITH TWO DOTS VERTICALLY ABOVE RIGHT SIDE, but it is unclear whether this needs to be distinguished from U+067A, and it is therefore not proposed for encoding at this time.

j	ج	z	ز	ŋ	غ
d	د	ʃ	ص	g	ع
r	ر	ɖ	ض		
z	ز	ɗ	ط		
s	س	ʒ	ظ		
ʃ	ش	ɾ	ع		
f	ف	ʁ	غ		
k	ك	q	ق		
l	ل				
m	م				
n	ن				
h	ه				
w	و				
y	ي				

**Figure 2:** Chtatou (1992), page 29.

d	ط	ز	ذ	ڨ	ن
f	ف	ر	ع	p	پ
k	ك	ق	غ	y	ي
l	ل		ق		
m	م				
n	ن				
w	و				
y	ي				

Fulfulde has adopted the same transcription rules as Songhoy (cf. Table I) except for the prenasalized consonants for which it devised new graphemes :

(17) nd	نڊ
mb	مب
ng	نگ
nj	نج

As for the sound /y/ which does not exist in Songhoy, it was transcribed as /يـ/ and the implosive /ɓ/ as /بـ/. Also the Arabic "tanwīn" or nunation "ـًـٍـِ" has been replaced in this transcription by the nasal /نـ/ at the end of the word.<sup>(21)</sup>

(21) "When the three vowel marks are written double at the end of a word, e.g. ـًـٍـِ, ـًـٍـِ and ـًـٍـِ they represent the three case endings, nominative, accusative and genitive of a fully declined, indefinite noun or adjective. The second vowel is pronounced "n"! Thus we have كَلْبٌ *kalbun*, a dog (nom.), كَلْبٍ *kalban*, a dog (acc.) and كَلْبِ *kalbin* a dog (gen.). This process of doubling the final vowel is called تَنْوِين *tanwīn*, or, by orientalisists, nunation, or "n'ing", from the Arabic name for the letter *n*. (cf. Cowan, (1958 : 6)).

**Figure 3:** Chtatou (1992), page 34. This also shows a BEH WITH THREE DOTS POINTING DOWNWARDS ABOVE RIGHT SIDE, but it is unclear whether this needs to be distinguished from U+067D, and it is therefore not proposed for encoding at this time.

#### ٣.٤.٤.٤. CONSONANTS

The 7 simple consonants that are proper to Pulaar have been transcribed in the following manner taking into consideration traditional transcriptions used for centuries by local scribes and religious scholars :

(24)	c	ش
	g	غ
	p	پ
	ɓ	ب
	y	چ
	ɟ	ڭ

**Figure 4:** Chtatou (1992), page 39. The positioning of the dots over/under the right end of the base form, rather than centrally, is considered to be an idiosyncrasy that does not merit separate encoding.

3.3.1. Hausa

The consonants that do not exist in the Arabic language have been transcribed in Hausa as follows :

3.3.1.1. Simple consonants

(26)	b	ب	bana	بنا
	c	٢	cocila	٢٢٢٢
			voiceless palato-alveolar fricative	
	ḍ	ط	ḍaki	طالبي house
	y	٢	yaruwa	٢رعو٢ parent

Figure 5: Chtatou (1992), page 42.

As for the labialized velar g<sup>w</sup>, it is transcribed by adding two dots to a normal *ghayn* /ع/ :

(28)	g <sup>w</sup>	ع٢	g <sup>w</sup> aba	ع٢با
------	----------------	----	--------------------	------

Figure 6: Chtatou (1992), page 43.

TABLE IV  
COMPARISON OF LETTERS PROPOSED BY MALI  
AND SENEGAL FOR PULAAR / FULFULDE

Latin characters	Proposed transcription		Examples	Gloss
	Mali	Senegal		
c	٢	ش	cakka	٢ك necklace
g	ع٢	غ	gerte	غر٢ ground-nut غر٢ peanut
p	٢ب	٢ب	paaka	٢اك knife ٢اك
ny	٢ن	٢ن	nyiiwa	٢ي٢ elephant
y	٢ب	ج	yiyal	٢ي٢ bone ٢ي٢
ŋ	ع٢	ق	ŋenyema	خ٢ earring ق٢
ḥ	٢م	٢ب	ḥoolde	٢ول٢ club ٢ول٢
nd	٢	ند	ndaayri	٢ا٢ ĩnenuphar ٢ا٢

Figure 7: Chtatou (1992), page 45. It seems clear that the writer is making a conscious distinction between three dots in a horizontal row and three dots in a triangle formation here.

TABLE V

COMPARISON OF LETTERS PROPOSED BY MALI AND NIGER  
FOR THE TRANSCRIPTION OF ZARMA / SONGHOY

Latin characters	Proposed transcription		Examples	Gloss	
	Mali	Niger			
c	ث	چ	ciiri	ثِير	salt
g	غ	غ	gaara	غَار	to solicit
ny	ن	ن	nyaamoy	نَامِي	
p	پ	پ	paate	پَات	
ŋ	ڭ	ڭ	naari	نَار	rice
o	و	و	koyra	كَيْر	village

Figure 8: Chtatou (1992), page 47.

TABLE VII  
PULAAR / FULFULDE STANDARDIZED ARABIC  
ALPHABET

Similar characters in Arabic and Pulaar/Fulfulde	Specific characters to Pulaar/Fulfulde	Pulaar/Fulfulde prenasalized sounds	Pulaar/Fulfulde vowels
ا ب ت ث ج د ر س ط ف ك ل م ن ه و ي	ث ب ت ق چ ن پ	مب ند نج نچ	و ا ي ه و د و

Figure 9: Chtatou (1992), page 50.

(34)

ا	without <i>hamza</i> for a, e, i, o, u
1	without <i>fatha</i> for aa
پ	Pesian and Urdu symbol for p (instead of پ)
ب	for b (instead of ب)
ث	for c (instead of ث)
ج	for y (no change)
ط	for d (no change)
گ	Persian and Urdu symbol for g (instead of غ)
ڭ	for ŋ (instead of ڭ)
و	for the vowel e

Figure 10: Chtatou (1992), page 54.

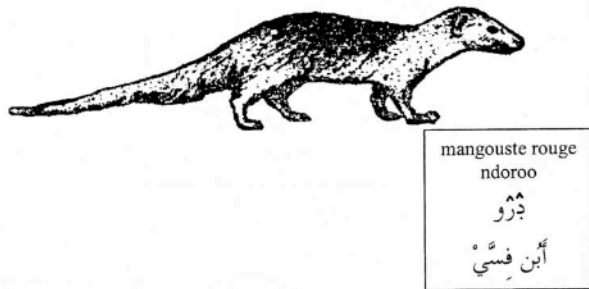
TABLE XI

### Consonants and semi-vowels

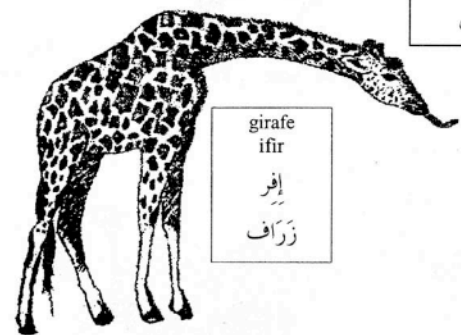
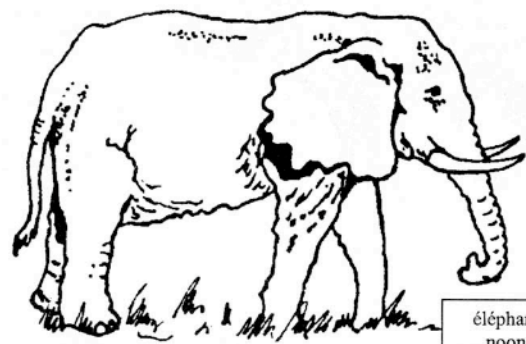
[illegible]

**Figure 11:** Chtatou (1992), page 58. Note that most of the characters shown in this chart are already encoded; only the indicated ones are relevant to this proposal.





6



7

Figure 12: Nodjindaina et al (2002), pages 6-7.

وَسَّ إِبْرَ حِدًّا نَا وَاقُنْ، إِبْرُقْ أَوْلِقُ حِدَ مِرُنْ  
كَأ، تَثْرِيَتْ تَمِي.

مَارِقُ حِدَ مِرُنْ تَ، كُرِيَا كَبِعَ نَا هَامَ دِرْت.  
كُغَالِقُ يَامَدَانُ، كُرِيَا كَبِعَ سُنْ أَلْعَ يَامَنْ دُم  
وَفَيْكَ، وَجَا كُجِنِ بَدَ زِت. كُغَالِقُ يَامَانُ  
جَا، كُرِيَا ثُور نَا هَامَ دِرْت.

أَسَالِقُ يَامَدَانُ، كُرِيَا ثُور سُنْ أَلْعَ يَامَنْ دُم  
وَفَيْكَ، وَجَا كُجِنِ بَدَ زِت. أَسَالِقُ يَامَانُ  
جَا، كُرِيَا أَتُق نَا هَمَ دِرْت.

ثُورِقُ يَامَدَانُ، كُرِيَا أَتُق سُنْ أَلْعَ يَامَانُ دُم  
وَفَيْكَ، وَجَا كُجِنِ بَدَ زِت. ثُورِقُ يَامَانُ  
جَا، إِذَا نَنِينُ أَشْبَقُ دُرْفُتَان.

Figure 13: Dahab et al (2002), page 19.

Mark	Description	Mandinka name	Roman equivalent
1.	dash above	sira tilindingo	ba, bē
2.	dash below	sira tilindingo duuma	bi, bē, bē
3.	half circle above	ngoo biringo	bo, bu
4.	circle above	sira murumurulingo	āb or consonant with- out following vowel
5.	large dot below left	tambi baa duuma	bē
<u>Nasals</u>			
6.	double dash above	sira tilindingo fula	bang, bēng
7.	double dash below	sira tilindingo duuma fula	bing, bēng, bēng
8.	double half circle above	ngoo biringo fula	bong, bung

Figure 14: Addis (1963), page 9: see item 5 in list of vowel marks.



	ḥ			ح	
c	c		ش	ش	ش
d	d	د			د
	dʰ			ط	ط
e	e		ـَـ	ـَـ	ـَـ
é			ـِـ		ـِـ
ë			ـِـ		ـِـ
f	f	ف			ف

Figure 15: Chtatou (1992), page 40: see transcriptions for /e/. Results of Workshop on the formulation of a standardized system of transcription of Wolof and Pulaar, 16–21 March 1987, Dakar, Senegal.

26

FULANI-HAUSA SCRIPTS

4. The **Short Vowels** are expressed by the signs in the third column :

	Hausa name	Sign	Value	Example	Fulani name
1	<i>wasalī bisā</i>	—	a	ب ba	masde dou
2	<i>wasalī ƙasa</i> (or <i>kisra</i> , or <i>kasāra</i> )	—	i	ب bi	masde les
3	<i>ḡigo</i> (or <i>gudā</i> ) <i>ƙasa</i> ; or <i>yamala</i>		e	ب be 	yamalēre (or yamalāre
4	<i>rufu'a</i>	—	u or o	ب bu	turnde

N.B.—In classical Arabic no. 4 has only the 'u' sound, whilst no. 3, i.e. 'e', does not exist; in consequence in many old Nigerian MSS. no. 2 will be found to represent both 'i' and 'e'.

Figure 16: Taylor (1929), page 26: see item 3 in table of vowel signs.

FULANI51

بِرْ بُولُورْ تَكِي دَغَمْ كَسَكَمْ دَغَمْ دَنْطْ دَا جَنْفَرْ دَمْ عَجِيحْ  
كَلْبِلْ عِيحْ جَمْ مَعْ عَجْمَرْ دَا دَسَارْ فُحْ : بَاوْمَنْ خَفِرْ د  
عَمْ جَمْ وَلَا كَلْبِلْ بَاوْمَنْ كَوَطْبِي سِي حِيْرْ عَمْبِلَرْ : عَمْبِلَرْ  
نَغَمْ قَحْمْ مِي يَمِيْنْ وَشَغْ عَسَرْ عَمْ فُحْ حَنْبِي كَ اَللهُ كِيَكِيْمْ  
دَرْ مَا كَ فُحْ كِيَكِي تَقْبَلْ نِي عَجْوَنْغْ : عَمْبُوْدَا عَطَارْ كَ  
مِي بِلْمَرْ وَنْطْ كِيَكِي تِي : بَاوْمَنْ قَحْمْ مِي يَمِيْنْ يَتُوْرْ طُوْبِي  
دَوْرَا كِيَكِي اَللهُ جِيْنْ بَلْمَرْ اِبْسَنْدْ مَنَغْ مَعْ عَمَّا جَوْنِيْمْ  
دَا بَا بَا جَوْنِيْمْ يَمِيْنْ فُحْ كُوْتَبْ نُوْلَامْ جَوْنِيْمْ سَكْ مَوْمَتْ  
اَللهُ بَسْمَدْ جَوْبِي مَعْ عَمْرِي مَعْ اَمِيْنْ هَذَا وَالسَّلَامْ :

I \*

Derēwol wurtake diga hā suka ma mo ardinḡa hā jangirde  
ma : ujineje | kofli [sc. kōfli] e yamgo njamu ma e njamu  
Dāda-sāre fuh. Bāwo man jangirde | ummi jam, walā ko fe'i  
bāwo ma kōdume, sei hairu e mo'ere. Andingo | ma fahin,  
mi timmini wittugo ngesa am fuh hande. Ko Allah hokki  
yam | nder māka fuh kabbe chappande nai e jowēgo. Am bō  
ndā ngedāri ko | mi yeḡi ma wonḡo, kabbe tati. Bāwo ḡon  
fahin mi yetti ma yettōre ḡunde | nde walā kempe, Allah juttin  
balḡe ma, o ḡesda mangu ma, gamā a heutini yam hā bāba  
heutinta ḡiyum fuh, kō bana nō lāmḡo heutinirta suka mūm. |  
Allah ḡesdu baḡde ma e darja ma, Āmīn. Hāzā wa's-salām.

you three bundles out of the largesse I am making. I am unend-  
ingly grateful to you : may God lengthen your days and grant you  
further honour, for you have helped me on as a father does his son,  
or as a chief advances his servant. May God increase your power  
and glory, Amen. This with weal.

D 2

Figure 17: Taylor (1929), page 51: showing use of dot for /e/ in running text (Fulani), and triple-dot punctuation mark (see section 3.4).



### 3.2 Jawi and Moroccan GAF characters

The Arabic letter *kaf*, nominally representing a /k/ phoneme, has several clearly distinct graphical forms. These derive from varying calligraphic traditions, and Arabic speakers understand them to be mere variations in the style of writing a single letter. As such, all can be represented by a single encoded character, U+0643 ARABIC LETTER KAF. Representative glyphs showing the three major forms of this letter, labeled A, B, and C, are shown in figure 18:



**Figure 18:** Three forms of the Arabic letter *kaf*

Form A (ك) is the form most commonly used in current text fonts, and is appropriately chosen for the representative glyph in the Unicode standard; it is also the form seen in typical charts of the Arabic alphabet, such as figure 19:

qāf	q, ƙ	[q]	100	ق	ق	ق	ق
kāf	k	[k]	20	ك	ك	ك	ك
lām	l	[l]	30	ل	ل	ل	ل

**Figure 19:** From chart of the Arabic alphabet, Daniels & Bright (1996), page 560.

However, where the Arabic script has been adopted for writing non-Arabic languages, variations in form that in Arabic were free variation or stylistic variants have sometimes been co-opted to make meaningful distinctions that merit encoding as separate characters. A clear example of this can be seen in Sindhi, where two forms of *kaf* are used as separate letters of the alphabet. The important phonemic distinction between /k/ (unaspirated) and /kʰ/ (aspirated) is represented by using form C (ﻙ) for /k/ and form B (ڪ) for /kʰ/, as shown in figure 20; the typical Arabic form A (ك) is not used in Sindhi.

q	[q]	ق	ق	ق	ق
k	[k]	ڪ	ڪ	ڪ	ڪ
kh	[kʰ]	ک	ک	ک	ک
g	[g]	گ	گ	گ	گ

**Figure 20:** From chart of the Sindhi alphabet, Daniels & Bright (1996), page 757.

To support the Sindhi usage (and other similar situations), we note that Unicode encodes the three *kaf* forms separately as distinct characters:

Code	Glyph	Name	Joining
0643	ك	ARABIC LETTER KAF	KAF
06A9	ڪ	ARABIC LETTER KEHEH	GAF
06AA	ﻙ	ARABIC LETTER SWASH KAF	SWASH KAF

**Figure 21:** Forms of Arabic *kaf* encoded in Unicode 4.0

(The name KEHEH used for U+06A9 is probably an attempt to transcribe the Sindhi name for the letter representing aspirated /kʰ/.) Figure 21 also shows that Unicode assigns these three characters to distinct joining groups, reflecting the fact that they are substantially different graphical forms and must each be shaped according to a different pattern.

A similar example of the disunification of Arabic glyph variants to become distinct characters when used for another language can be seen in the Urdu usage of the letter *heh*. Here, a contrast between ه and ه (and their related linking forms) is consistently used to distinguish the independent letter /h/, written with ه, from

aspiration of plosives and affricates, written with **هـ**. Arabic speakers would consider these variants of a single letter, but a clear distinction must be encoded for some other languages (and is therefore supported in Unicode).

Yet another example occurs with *yeh*: to an Arabic speaker, the form **يَ** is merely a calligraphic variant of the letter **ي** (or **ى**). But in Urdu, the form **يَ** has been adopted to write the vowel /e/, while the form **ی** represents /i/. This distinction must be encoded, and so Unicode includes U+06D2 **يَ** as a separate character.

So we see that where there are clearly distinct graphical forms in existence for an Arabic letter, it may well be appropriate to encode these forms separately. The fact that they originate as different calligraphic styles of a single letter in the Arabic language does not mean that this interpretation is adequate for all languages and regions.

Given that Unicode encodes these three forms of *kaf* separately, it seems appropriate to treat modified forms of *kaf* in a similar way. Where additional letters have been created by adding dots or other marks to an underlying *kaf*, this has often been done to one specific form of the letter (or in Unicode terms, to one of the three characters U+0643, 06A9, 06AA), and substitution of a different base form may not be at all acceptable. This is clear, for example, in the case of the Persian and Urdu /g/, written as U+06AF **گ**; the added 'bar' that creates the letter *gaf* can only be added to form B of the *kaf*.

The following two characters, representing *gaf* as written in different regions, are therefore proposed for encoding; they are discussed individually following the summary table. Both share similar Unicode properties: general category Lo; combining class 0; bidi type AL.

Glyph	Code	Character name	Shaping	See figures
ک	0762	ARABIC LETTER KEHEH WITH DOT ABOVE	GAF	22, 24
ک	0763	ARABIC LETTER KEHEH WITH THREE DOTS ABOVE	GAF	25, 26, 27, 28, 29

### 3.2.1 Jawi GAF

In the Jawi script (Arabic script used to write Malay), the /g/ sound is written using a *kaf* with one dot above. Such a character is encoded in Unicode at U+06AC (**ك**), with representative glyph based on form A of *kaf*. However, the Jawi *gaf* is properly based on form B, not form A:

ک	[k]	ق	ق	ق	ق
k	[k]	ك	ك	ك	ك
g	[g]	ک	ک	ک	ک
l	[l]	ل	ل	ل	ل

**Figure 22:** From chart of the Jawi alphabet, Daniels & Bright (1996), page 761.

Note in figure 22 that form B is used as the basis for the /g/ character, despite the fact that form A is used for /k/. This is a characteristic of Jawi writing, and not an artifact of this particular book's typography. Comparing the chart for Uighur, found on the previous page, we see that in some cases form A is used as the base for a modified *kaf*; the use of form B in the Jawi chart is no accident.

ق	[q]	ق	ق	ق	ق
k	[k]	ك	ك	ك	ك
g	[g]	گ	گ	گ	گ
ng	[ŋ, n]	ڭ	ڭ	ڭ	ڭ
l	[l]	ل	ل	ل	ل

**Figure 23:** From chart of the Uighur alphabet, Daniels & Bright (1996), page 760.

Figure 23 shows a form A *kaf* with three dots above, used for the /ŋ/ sound. Comparing this with the Jawi chart, we see that new letters based on *kaf* may involve a deliberate choice of one of the three forms of *kaf*, which may not be interchangeable or treated as glyph variants in this context.

All Jawi sources I have seen show the use of form B as the basis of the /g/, even though form A is commonly used for /k/. Figure 24 is taken from an introduction to the Jawi script published in Malaysia.



Huruf [ڤ] —huruf wau bertitik] dicipta dan diperkenalkan dalam tahun 1984 di Konvensyen Tulisan dan Ejaan Jawi di Kuala Terengganu untuk melambangkan huruf [v] dalam tulisan Rumi.

Huruf-huruf Jawi yang digunakan untuk menulis dan mengeja kata-kata dalam bahasa Melayu semuanya berjumlah 35 huruf yang dipinjam daripada huruf-huruf Arab (Huruf هجاء). Enam daripada jumlah huruf tersebut dicipta oleh orang Melayu sendiri bagi melambangkan bunyi-bunyi kata bahasa Melayu yang tidak terdapat dalam huruf-huruf Arab. Huruf-huruf tersebut dicipta berdasarkan bentuk-bentuk asal.

- (a) Huruf [ڤ] padanan huruf Rumi [c].
- (b) Huruf [ڤ] padanan huruf Rumi [ng].
- (c) Huruf [ڤ] padanan huruf Rumi [p].
- (d) Huruf [ڤ] padanan huruf Rumi [g].
- (e) Huruf [ڤ] padanan huruf Rumi [ny].
- (f) Huruf [ڤ] padanan huruf Rumi [v].

Figure 24: From Muhani (1998), page 6.

It is clear from the notes in the names list that the Unicode character U+06AC (ك) was encoded with the intent that it be used for Jawi /g/; however, we see that the representative glyph shown in the code charts and the joining group listed in *ArabicShaping.txt* are inappropriate for this purpose.

One possible response would be to change the glyph and the joining group of U+06AC to those expected in Jawi, thus making this character suitable for its originally-intended purpose. However, given the tendency, especially once modifying marks are added, for users to make a clear distinction between the different forms of *kaf*, not considering them merely as glyph variants of a single character, this comes dangerously close to changing the fundamental identity of the character. Moreover, there can be no assurance that a character having the specific form A with a dot above has *not* been deliberately used in some context, given that it exists in the current standard.

The fact that the Arabic joining group is considered a normative property of the Unicode character also weighs against the position that ك and ك could be considered variants of the same character; they must necessarily have different joining groups. It is therefore proposed that a new character ARABIC LETTER KEHEH WITH DOT ABOVE should be encoded, and a note added to U+06AC indicating that the new character is preferred for old Malay.

### 3.2.2 Moroccan GAF

Although standard Arabic does not write a /g/ sound, in Morocco the use of a form B *kaf* with three dots above is well established as the letter representing /g/. Published literature is generally in standard Arabic, and as such does not use this letter, but it is seen in other situations such as road signs, product labels, etc. It is also used in writing the Amazigh languages of Morocco. The following photographs show examples of this letter used in Moroccan Arabic:



Figure 25: Street name in Arabic and Latin scripts.



Figure 26: Petrol pump, labeled 'gasoil' in Arabic script.

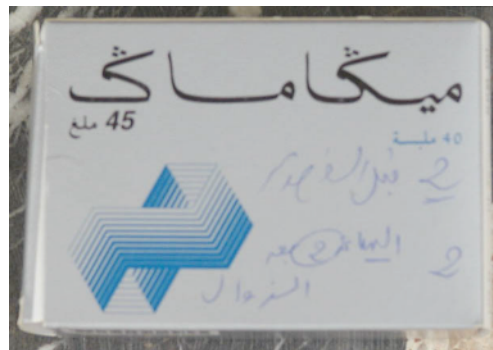


Figure 27: Medicine package 'Megamag', showing that the *gaf* is based on *kaf* form B.

Figure 28 is taken from the Royal Moroccan Academy's new Amazigh dictionary, and shows the letter in initial, medial, and isolated forms.

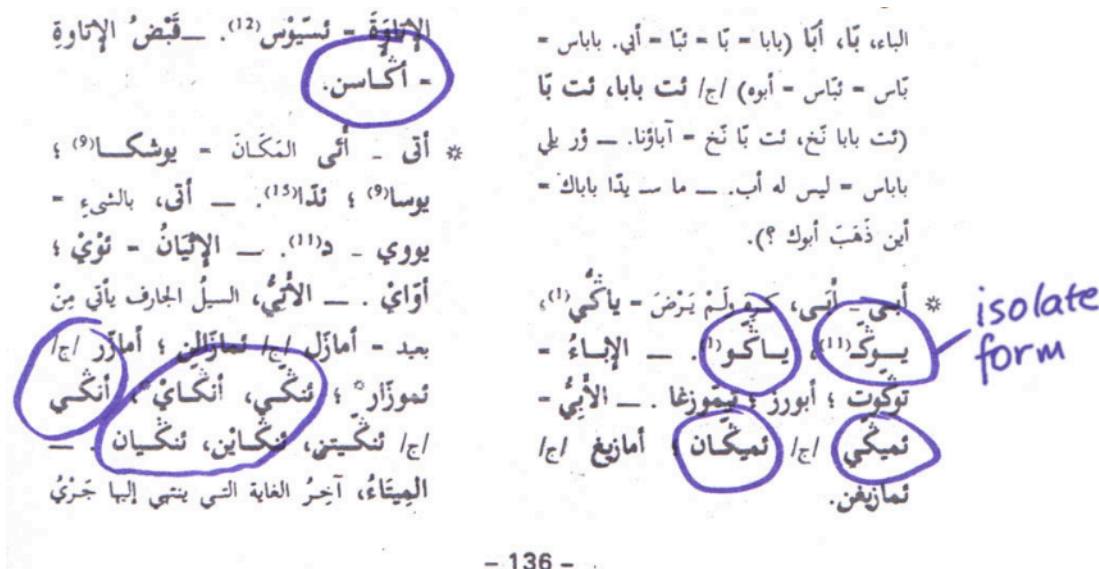


Figure 28: From Shafiq (1996), page 136; note the “makeshift” isolated *gaf*.

Close examination of figure 28 shows that the author did not have an appropriate character available, and therefore used a standard *kaf* and added the three dots by hand. The isolated form is of particular interest, as here



the author deliberately added a *kashida* (extender) character in order to achieve an approximation of *kaf* form B, in preference to simply adding the dots to the form A *kaf* which would otherwise have appeared.

A similar form is seen in Mauritanian texts, where a *kaf* written in shape B with three dots added is also used to represent /g/:

<i>ḍhar</i>	The 'black' way ( <i>jaamba</i> )	<i>gnaydiya</i>	The 'white' way ( <i>jaamba</i> )
ظهر	الطريق الكحلاء	الكنيديه	الطريق البيضاء
<i>kar</i>	<sup>1</sup> <i>entemaas</i>	<i>enweffal</i>	<sup>1</sup> <i>mekka mūsa</i>
كر	انتماس	انوقل	مك موس
	<sup>2</sup> <i>sayni kar</i> ( <i>ma yuḡarraṣ</i> )		<sup>2</sup> <i>el-faayez</i>
	سيني كر (ما يخرص)		الفايز
<i>faaḡu</i>	<sup>1</sup> <i>tenaččuuga</i>	<i>ššbaar</i>	<sup>1</sup> <i>ssruuzi</i> (' <i>arraay essruuza</i> )
فاغ	تنچوڤه	الشبار	السروزي (عراي السروزه)
	<sup>2</sup> <i>sayni faaḡu</i>	or <i>faaḡu lekbiir</i>	<sup>2</sup> <i>et-teḥraar</i> ( <i>el-hurr</i> )
	سيني فاغ	فاغ الكبير	التحرار (الحر)
<i>siññiima</i>	<sup>1</sup> <i>siññiimat hayba</i> <sup>3</sup>	('black') <i>el-mawṣṭi</i>	<sup>1</sup> <i>et-teḥzaam</i> <sup>4</sup>
سيمه	سيمه هيب	الموسطي	التحزام
	<sup>2</sup> <i>meqaččuuga</i> ( <i>Hawd</i> )		<sup>2</sup> <i>eññaama</i> <sup>5</sup>
	مقچوڤه		انيامه
	or <i>lebyaad</i>	('white') <i>menčalla</i>	<sup>2</sup> <i>rrbaabi</i> <sup>6</sup>
	لبياظ	منچله	الربابي
	or <i>etbaybi</i> ( <i>Trārza</i> )		(a) <i>leggetri</i> <sup>6</sup>
	اتببي		(b) <i>čaynna</i>
			(c) <i>liyyin</i>
			لين
<i>baygi</i> or <i>lebtayt</i>	<i>baygi</i>	<i>el-muḡaalef</i>	<i>el-'ittiig</i>
بيغي	بيغي	المخالف	العتيك
	intensive form <i>a'ḍḍaal</i>	or <i>baygi jḡraad</i>	
	اعضال	بيغي الجراد	
		or <i>baygi lekbiir</i>	
		بيغي الكبير	

<sup>1</sup> Entry (*ḍkuul* دخول) brought about through tightening the strings of the *tidīnit*.  
<sup>2</sup> Subsidiary (*rrdiif* الرديف) brought about through loosening the strings of the *tidīnit*.  
<sup>3</sup> Also called in western Mauritania *zzraag* الزراك.  
<sup>4</sup> Included in *faaḡu* by some *iggāwen*.  
<sup>5</sup> Not subsidiary to those who include *et-teḥzaam* in *faaḡu*.  
<sup>6</sup> A blending of white and black.

Figure 29: From Norris (1968), page 73. Note contrasting basic shape used for final *kaf* and *gaf*.

Here again, the printer has been forced to make do with a limited selection of available glyphs, and the results are instructive. The final and isolated forms of *gaf* are constructed by adding the 'tail' of a *kaf* (intended for

use in building form A) to a medial or initial *kaf* with three dots added, giving a result that resembles form B in having the added bar on top, but also has the ‘flourish’ typical of form A. The importance of the ‘form B-ness’ of this letter is evident when we note the trouble that has been taken to build it, in contrast to the simple form A used for *kaf* itself (highlighted in blue in figure 29).

In initial and medial joined forms, this letter would be visually identical to U+06AD ﻚ ARABIC LETTER NG. However, in final and isolated forms the difference is clear; Moroccan /g/ is consistently based on form B of *kaf*, even though form A is commonly used for the letter *kaf* itself. Given the clear distinction between these forms in users’ minds; the fact that the forms ﻚ and ﻜﻲ are not considered interchangeable, even where ﻚ and ﻜ are understood to be variants of the same letter *kaf*; and the different joining groups required, it is proposed that a new character ARABIC LETTER KEHEH WITH THREE DOTS ABOVE should be encoded in the UCS to represent the Moroccan/Amazigh *gaf*.

### 3.3 South Asian extensions to Arabic script

The characters listed in this section have been used in writing several languages in the South Asia region; examples of the use of each character in one or more languages are included, but these should not be interpreted as representing the entire scope of use. It also seems likely that as additional minority languages in the same areas establish orthographic conventions, they may adopt some of these characters.

Four characters from South Asian languages are proposed for encoding; they are discussed further in individual subsections following the summary table. All share similar Unicode properties: general category Lo; combining class 0; bidi type AL.

<i>Glyph</i>	<i>Code</i>	<i>Character name</i>	<i>Shaping</i>	<i>See figures</i>
طٲ	0759	ARABIC LETTER DAL WITH TWO DOTS VERTICALLY BELOW AND SMALL TAH	DAL	32, 33, 34
ش	075C	ARABIC LETTER SEEN WITH FOUR DOTS ABOVE	SEEN	30, 31
نٲ	0768	ARABIC LETTER NOON WITH DOT ABOVE AND SMALL TAH	NOON	32, 33, 34, 35, 36
نٶ	0769	ARABIC LETTER NOON WITH SMALL V	NOON	37, 38

#### 3.3.1 Retroflex DAL with two dots below

This character is used by Saraiki writers to represent an implosive retroflex /d/. Other Saraiki retroflexes are written with a small TAH above, following Urdu (e.g., U+0679, U+0688); other implosives are written with two dots vertically below, following Sindhi practice (e.g., U+067B, U+06B3). The natural tendency for a writer needing to write an implosive retroflex /d/, then, is to add the dots below a character like U+0688. Examples can readily be found in Saraiki newspapers and other books.

#### 3.3.2 SEEN with four dots above

Writers of the Shina language in Kashmir, needing to write a retroflex /s/ letter, have adopted the convention of using a letter based on SHEEN but with four dots in place of three.

As a pattern of four dots is not normally seen in the area (though it does occur in southern Pakistani languages such as Sindhi; see U+067F, U+0680, etc.), some writers have tended to replace the four dots with two horizontal lines. This is seen in some of the examples (below). However, as it is common in handwritten script to see horizontal pairs of dots written as a single line, it seems appropriate to consider this a glyph or stylistic variation, encoding the character as SEEN WITH FOUR DOTS and leaving it to font designers to determine whether to offer glyphs with the dots replaced by lines. (It does appear from Taj (1989) that some writers specifically choose to write the retroflex /s/ with horizontal lines, even when well-formed dots are used on other letters. Nevertheless, it seems clear that the form with four dots and that with two lines are glyph variants of the same underlying character.)

It may be noted that Akbar (1985) also shows a character with the form of HAH WITH FOUR DOTS ABOVE used by a Shina writer (to represent the phoneme /ts/). The information currently available suggests that this

character has not been as widely adopted as the others shown here, with other writers using U+0685 or U+0697 for this sound, so the case for adding it to the UCS repertoire may be less clear-cut, at least until further data is obtained.

### 3.3.3 Retroflex NOON with dot

A number of languages use this character to represent a retroflex /n/. The use of a small TAH as part of an Arabic-script letter to indicate retroflexion is derived from Urdu usage, but Urdu itself does not have a retroflex /n/.

Sindhi has this sound, and writes it with a NOON where the dot is replaced by the small TAH (rather than the small TAH being added, while retaining the dot as well). However, this is not an option in languages where the orthography is based on Urdu, as the initial and medial forms would be indistinguishable from the retroflex /t/ (U+0679). Writers of such languages have therefore devised this letter, which is not yet encoded in Unicode.

To reduce the likelihood of confusion with the Sindhi retroflex /n/ letter (encoded as U+06BB ARABIC LETTER RNOON), the suggested name explicitly mentions the presence of the dot in this character; a name such as NOON WITH SMALL TAH might be taken to imply that the small TAH replaces the dot instead of being an addition, which would result in a character identical to U+06BB.

### 3.3.4 NOON with small V

The Gojri community is found in both India and Pakistan, and there has been a significant amount of literature published in both countries from the 1980s onwards. While there has been some variation in orthographic conventions, there is widespread use of a NOON WITH SMALL V to represent the retroflex nasal consonant. (Gojri writers also use a LAM WITH SMALL V, as seen in the examples, but this is already encoded at U+06B5.)

Phonologically, this character is used (at least in the Gojri language) to write the same sound as Sindhi RNOON or the NOON WITH DOT AND SMALL TAH proposed here, but it clearly represents a distinct choice of extended-Arabic character for this sound; the SMALL V and SMALL TAH show two different conventions for the creation of new Arabic-script letters. The fact that the phoneme being written may be the same is irrelevant to the encoding of the written characters.

### 3.3.5 Samples showing South Asian characters

The examples shown here are drawn largely from books and newspapers published in Pakistan (an Indian example is also included), written in Shina, Saraiki, Pathwari, and Gojri. However, it can be expected that as literacy becomes more widespread among minority language communities in South Asia, many of these writing conventions may be “borrowed” by neighboring communities. The languages cited here serve to demonstrate the need to encode these characters, but should not be assumed to be the only users of them.

## بسکوپے حروف تہجی گہ حروف جار

رشتنا یا ش ای جیک خاص رسم الخط ایک حروف تہجی ایک نوش  
آپکو لکھو کے کار پس بو سے باک اردو ای حروف تہجی استعمال تھوئیں  
نا چیز سے کافی عرصہ بعد ریڈیو پاکستان ای خبر شعبہ انوکوم تھوئیں۔  
تو نے انوارہ انر بس رشتنا انر خبر سے لکھو انر کے بسکوپے  
حرفی گہنے حروف جار استعمال تھوئیں، اے انی ہن:  
بسکوپے حروف تہجی:

### حروف تہجی

اب پ ت ٹ ث ج چ ح خ  
د ذ ز ر ژ ت س ش ٹ ص من ط ظ  
ع غ ف ق ک گ ل م ن و ہ ع ی  
بھ چھ ٹھ ڈھ چھ ٹھ کھ گھ

حرف	لفظ	اردو معنی
پ	چلو	روشنی
ت	تھک	ٹھہرنا
ٹ	تھمو	آپ (جمع)
ث	تھا	بجائی
ج	شوٹو	گلا
چ	چھیلے	کپڑے

Figure 30: Akbar (1985), page 221 (Shina).

ش



شولسر

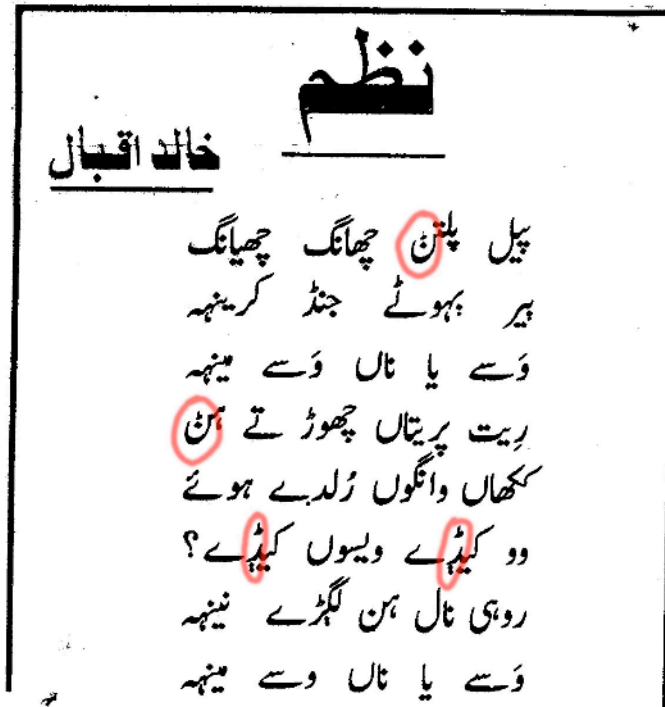


شولسر

Figure 31: Taj (1989), page 29 (Shina): examples of SEEN WITH FOUR DOTS ABOVE written using two lines in place of four dots (a practice also used for other four-dot letters in this book)

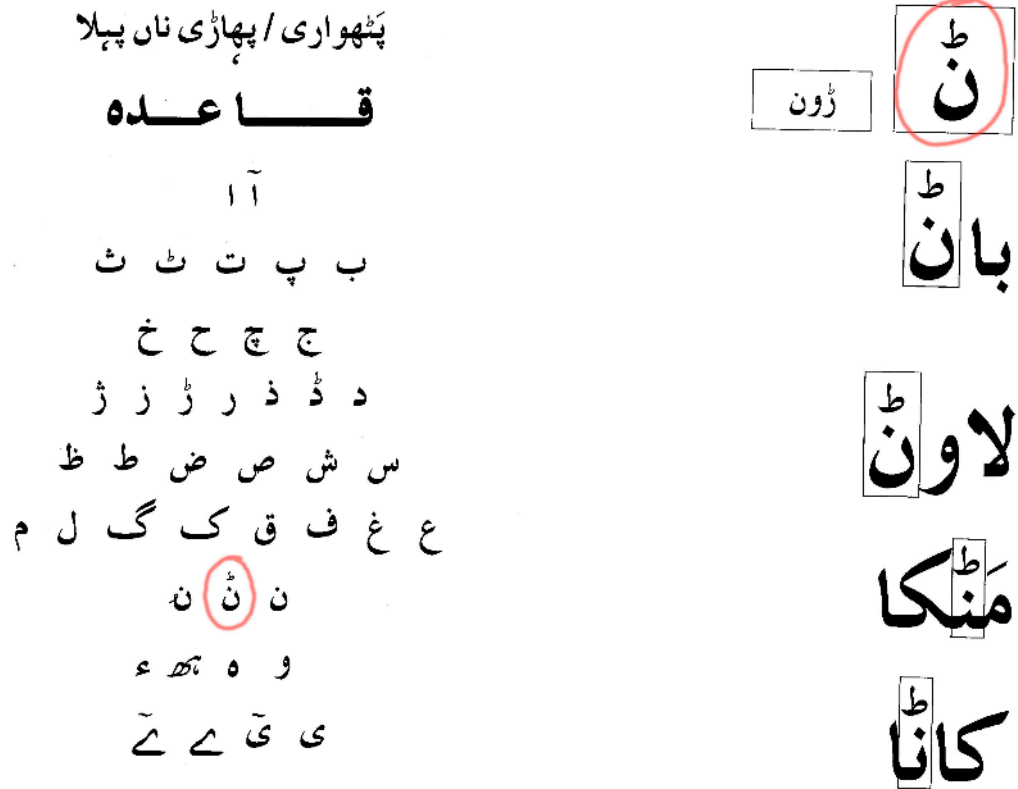
<p>آ ا ب پ ت ٹ ث ج چ ح خ د ڈ ذ ر ژ س ش ص ض ط ظ ع غ ف ق ک گ گ ل م ن ٹ و ہ ع ی</p>	<p>12 TREE وَن ٹ LEAF وَرَق و DEER ہرن ہ GUITAR یکتارا ی</p>
--	--

Figure 32: Mughal (1994), pages 12-13 (Saraiki).



**Figure 33:** Saraiki poem printed in *Daily Jhoke*, 29 January 2002. The small size of the added TAH and TWO DOTS on the special Saraiki letters arises because these letters are not supported in Urdu software used to produce this newspaper, so these have been added manually. A similar size anomaly can be seen in the dots under BEEH (U+067B) in line 2 of the poem; this character is also not available in the Urdu font used, so the Urdu BEH (with one dot) has been typed, and a second dot added manually.





**Figure 35:** Chitka (n.d.), inside front cover and page 20 (Pathwari). The alphabet chart also shows a second “modified NOON” character, but it is unclear whether this form can be considered well enough established for standardization.

گھڑے ناں رولا بجھی کے اکڑنی اکھ کھلی گئی۔ پرہن گدھڑ لپڑا بناسی !  
 ”ہاتے او مہاڑیے بے جی“ اکڑتے باکڑ ڈانگے۔  
 ”ہو کو، ہو کو، ہاتے ہاتے پاکڑ پٹیا“ مہاڑا کھڑ مہاڑا سیل“  
 تو ہاڑا کھڑ، تو ہاڑا سیل، اکڑتے باکڑ روٹیاں روٹیاں غتے نال پاکڑ کی تاڑیا  
 ”اوہ ساھڑا کھڑ“ ساھڑا سیل پاکڑ فٹوٹ بن بدلیا۔

**Figure 36:** Mehmood (2001), extract from page 41 (Pathwari).



جہنگل لہو کا دریا جاری کرے نگہ وؤہ اقدار تے ملک گیری کی ہوس برائیں۔  
 رت کے اٹھوئیں روان پرتل گئی ہے۔ تے کہ سے مہار قلم کار سچی اپنا  
 دور کا سوچتے بچ کا پارکھ بن سکتا تے کاٹے گل بن جاتی۔  
 ہوں خبر سے کن حالات ماں اپنا بکھاریاں کی کس دکھتی رگ پر ہتھ  
 رکھنا چاہوں۔ تے نہاں اس کے واسطے کس درجہ نفس شناسی کو اور اک حاصل  
 کر لو ہے۔ یاہ الگ گل ہے پر پھر بھی ہوں کہن تے بغیر نہیں رہ سکتو۔ کہ اپنا  
 زمانہ کا دکھاں درواں تے سی ہاں ناں نظر ملان تے بغیر جین آلا کچھ  
 ہو رہا ہوں تے پیا کہاویں پر شعر تے ادب کی بستی کا بسنیک بن کو ڈرامو  
 نہ چاویں۔  
 زیر نظر شہ ازہ ماں جن قلم کاراں کیس بکھت شامل ہیں۔ ہم  
 اُن کا شکر گزار ہاں پر جن اس جملہ پرانی نمائی ٹور ٹرن کو وقت گزر گئیو  
 ہے۔ ہن تے اوکھتاں تے الجھناں کے ناں جھن کو وقت آگیا ہے۔ لمیاں تے  
 کھکھیاں آباں ماں دل کا دیا بال کے کھن کی لوڑ ہے۔ سبجری تے نروئی  
 مسج کی رماں ناسوجاگ اکھاں ناں ہر کہے اکہن کی ضرورت ہے۔  
 اقبال عظیم

Figure 37: Azeem (1996), unnumbered page (Gojri).

محنت ناں ترقی کرے رکھ کے ماں ترقی کرے  
 محنت ناں ترقی کرے سالو سال ترقی کرے  
 نوں این بیس کے چور بھٹیا  
 منزل تیری دور نہ بھٹیا  
 تعلیم  
 بے ہوشو، ہن ہوش ماں آو ماشومان کو غم بھی کھاو  
 میری متو علم پڑھاو گدراں ناں اسکول چلاو  
 بے تعلیمی پود نہ چھوڑو!  
 اس بوٹا ماں پود نہ چھوڑو  
 علم کا دشمن اکلا لوک آن پڑھ ساد مراد لوک  
 وے نیہ رہیا زیادہ لوک علم پڑھاویں آج کہ لوک  
 دلا ناں بدلنو ہو سککو  
 ناں ہوا کے چیلو ہو کو  
 علموں باج نہ آگے چالو علموں باج نہ کھڑو بھالو  
 علموں باج لے بندو کالو علموں باج نہ رتہو پالو  
 علموں باج نجات تیہہ جوت  
 راضی رب کی ذات تیہہ جوت  
 ۵۷

Figure 38: Afaqi (n.d.), page 57 (Gojri).

### 3.4 Triple-dot Arabic punctuation mark

Traditional orthographic practice when writing African languages such as Hausa, Wolof, Fulani, Mandinka, etc., in Arabic script includes the use of a characteristic punctuation mark as a “stop”. This mark consists of three dots in a triangle, similar to the pattern of dots found on Arabic letters such as THEH and SHEEN, but used independently as a punctuation mark.

Although there is some suggestion that this mark is falling out of favor, with present-day writers tending to use more Latin-like punctuation, its widespread use in older texts appears to justify its encoding.

The following character is therefore proposed for addition to the UCS, with a suggested code value of U+061E:

Glyph	Code	Character name	GC	CC	Bidi	See figures
⋯	061E	ARABIC TRIPLE DOT PUNCTUATION MARK	Po	0	AL	17, 39, 40, 41

The character properties proposed here are analogous to those of other Arabic-script punctuation marks such as U+061B and U+06D4.

This character should be treated similarly to the other punctuation marks in collation; it is suggested that it be given a primary weight such that it sorts after U+06D4 ARABIC FULL STOP by default, although this is a rather arbitrary choice. The collation weight can of course be tailored as needed if a different behavior is desired in any particular language.

The following figures show examples of the use of this character.



The only punctuation mark that is used is equivalent to a full stop and is composed of three dots thus - ∴. One may occasionally encounter the Arabic interrogation mark ؟, but there is no other common form of punctuation, mainly because any punctuation as we know it can be easily confused with a letter or part of a letter in the Arabic alphabet. N.B. Writers today are tending to use the single dot full stop sign in preference to the three dot sign.

Figure 39: Addis (1963), extract from page 13.

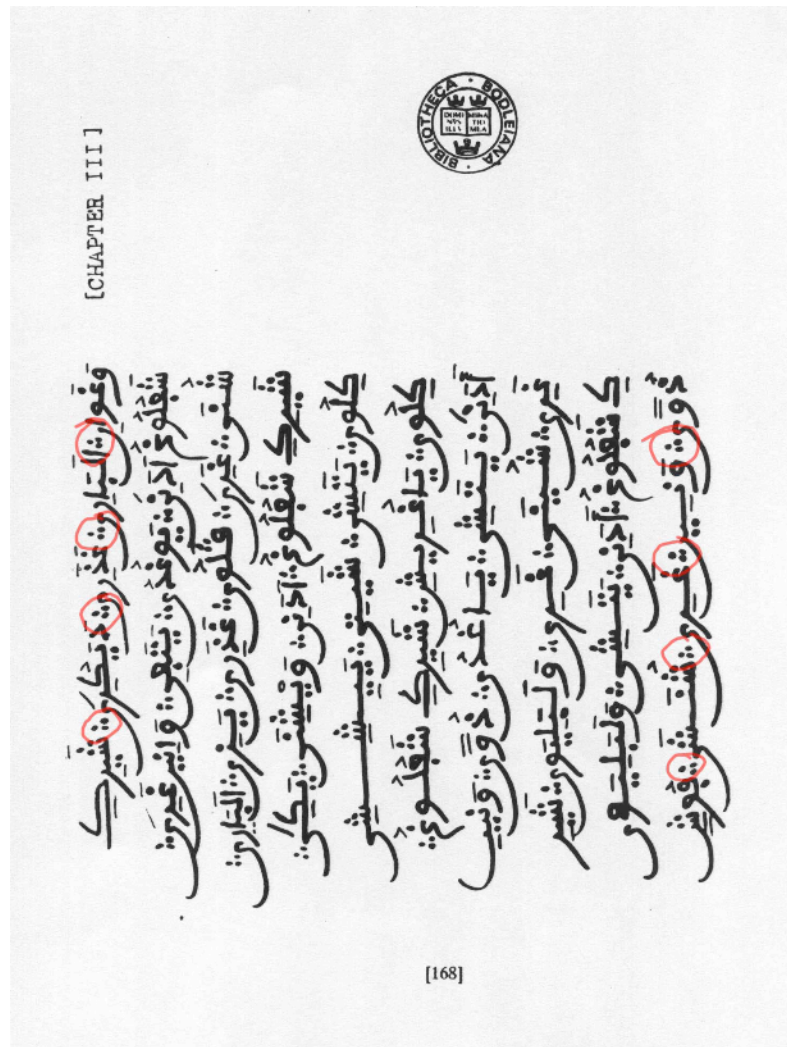


Figure 40: Piłaszewicz (1992), page 168. Many examples of the proposed character.

the syllable it is used alone in order to avoid two alifs : this is done in Arabic by the madda~, really an alif written horizontally : ex.

آى āya, the sign ◀ used as a stop

آية 'āyatun', is the same word in Arabic

الْقُرْآن Alkura'āni, and قُرْآن Kur'āni, the Koran

Figure 41: Taylor (1929), extract from page 33.

## 4. References

- Addis, R. T. 1963. *A study on the writing of Mandinka in Arabic script*. [SOAS library, London: call # LY Mandingo A/574203]
- Afaqi, Dr. Sabir. n.d. *Pegham-in-qalab*. Pakistan. [Gojri poetry]
- Akbar, Akbar Hussain. 1985. *Sumulo Rasuul (Holy Prophet)*. Islamabad, Pakistan: Modern Book Depot.
- Azeem, Iqbal Chaudhry, ed. 1996. *Sheeraza*. Srinagar: Jammu & Kashmir Academy of Art, Culture and Languages. [Gojri periodical]
- Bhaya, Bashir Ahmad. 1984, revised 1998. *Saraiki quaid te zubandani [Saraiki, father of languages]*. Bahawalpur, Pakistan: Saraiki Adabi Majlis.
- Chitka Committee. n.d. *Pathwari/Pahari nā pahala qaidah [First Pathwari/Pahari primer]*. Jhelum, Pakistan: Praala Publishers.
- Chtatou, Mohamed. 1992. *Using Arabic script in writing the languages of the peoples of Muslim Africa*. Rabat: Institute of African Studies. [Reports on a series of workshops held in several countries to work towards standardization of Arabic-script orthographies for major African languages.]
- Dahab, Abdoulay Ali, Abdoulay Issakha, Badour Adbelkerim and Evodie Zürcher. 2002. *Kitab aafe kiraa naa [Livret sur la santé]*. Abéché: Projet de développement de la langue maba. [Simple health booklet in the Maba language of Chad, Arabic script edition.]
- Daily Jhoke*. Multan, Pakistan. [Saraiki daily newspaper]
- Daniels, Peter T. and William Bright (eds). 1996. *The world's writing systems*. New York/Oxford: OUP.
- Kew, Jonathan. 2002. *Proposal for extensions to the Arabic block*. L2/02-274.
- . 2003. *Proposal to encode Arabic triple dot punctuation mark*. L2/03-159.
- . 2003. *Proposal to encode Arabic-script letters for African languages*. L2/03-168.
- . 2003. *Proposal to encode Jawi and Moroccan Arabic GAF characters*. L2/03-176.
- . 2003. *Draft chart showing UTC #95 additions to Arabic blocks*. L2/03-210.
- Mansour, Kamal. 2003. [http://www.bisharat.net/A12N/Afro-Arabic\\_Symbols.pdf](http://www.bisharat.net/A12N/Afro-Arabic_Symbols.pdf). [A collection of glyphs designed for African-language use; it is unclear in some cases whether actual use is established.]
- Mehmood, Tariq. 2001. *Sayana gidhar [Clever jackal]*. Jhelum, Pakistan: Praala Publishers.
- Mughal, Shaukat. *Saraiki qaidah [Saraiki primer]*. Multan, Pakistan: Saraiki Isha'ati Idarah.
- Muhani, Hj. Abdul Ghani. 1988. *Teman pelajar Jawi*. Petaling Jaya (Malaysia): Fajar Bakti.
- Nodjindaina, Jean-Bosco, Gami Ssane Mogaye, Mbanji Bawe Ernest, Susan Rose and Matt Day (illus.). 2002. *Erniye kadade-naanu (Ernime) [Les animaux sauvages de la brousse]*. 2ème édition révisée. N'Djaména/Abéché: Association SIL. [Picture book of animals and birds of the Maba area, Chad.]
- Norris, H. T. 1968. *Shinqiti folk literature and song*. Oxford: Clarendon Press.
- Piłaszewicz, Stanisław. 1992. *The Zabarma conquest of north-west Ghana and Upper Volta: A Hausa narrative "Histories of Samory and Babatu and others" by Mallam Abu*. Warsaw: PWN—Polish Scientific Publishers.
- Shafiq, Muḥammad (ed). 1990. *al-Mu'jam al-'Arabi al-Amazighi*. Rabat: Akadimiyat al-Mamlakah al-Maghribiyah [Royal Moroccan Academy].
- Taj, Abdul Khaliq. 1989. *Shina qaidah [Shina primer]*. Gilgit, Pakistan: Muhammad Book Stall.
- Taylor, F. W. 1929. *Fulani-Hausa readings in the native scripts: With transliterations and translations*. Oxford: OUP.