

Universal Multiple-Octet Coded Character Set
 International Organization for Standardization
 Organisation Internationale de Normalisation
 Международная организация по стандартизации

Doc Type: Working Group Document

Title: Revised Proposal to encode Arabic characters used for Bashkir, Belarusian, Crimean Tatar, and Tatar languages

Source: See References

Authors: Ilya Yevlampiev, Karl Pentzlin, Nurlan Joomagueldinov

Status: Expert Contribution

Action: For consideration by JTC1/SC2/WG2 and UTC

Date: 20 May 2011

Supersedes: L2/11-138 "Proposal to encode Arabic characters used for Bashkir, Belarusian, Crimean Tatar, Karachay, Karakalpak, and Tatar languages"

This revision incorporates the decisions of UTC #127, as documented in WG2 N4065 = L2/11-201 »Proposed codepoints and properties for characters in "Proposal to encode Arabic characters used for Bashkir, Belarusian, Crimean Tatar, Karachay, Karakalpak, and Tatar languages"«. The three characters proposed in L2/11-138 which were not accepted by UTC #127 are no longer subject of this proposal. Therefore, the terms "Karachay" and "Karakalpak" were dropped from the title.

1. Introduction

The letters of Bashkir, Belarusian, Crimean Tatar, and Tatar Arabic orthographies that are not encoded in Unicode standard are presented here.

2. Proposed Characters

Annotations below address special issues for a character, or reference to figures where such special issues are discussed. (These annotations are not intended to be retained in the character list when copied into the standard.)

Proposed codepoint	Glyph	F M I	Proposed name	Linguistic comment
08AD	ﻱ	••ﻱ	ARABIC LETTER LOW ALEF • Bashkir, Tatar	<i>denotes the usage of back vowels in the word with dubious spelling in Yaña imlâ orthography</i>
08AE	ډ	ډ••	ARABIC LETTER DAL WITH THREE DOTS BELOW • Belarusian	<i>denotes the sound [dz]</i>
08AF	ڇ	ڇڇڇڇ	ARABIC LETTER SAD WITH THREE DOTS BELOW • Belarusian	<i>denotes the sound [ts]</i>
08B0	گ	گگگ	ARABIC LETTER GAF WITH INVERTED STROKE • Crimean Tatar, Chechen, Lak	<i>denotes the sound [j] in Crimean Tatar denotes the sound [k] in Chechen and Lak</i>
08B1	ﻭ	ﻭ••	ARABIC LETTER STRAIGHT WAW • Tatar	<i>denotes the sound [w] in Urta imlâ Alparov's orthography</i>

3. Encoding Considerations

Joining type and group for ArabicShaping.txt:

08AD; LOW ALEF; U; No_Joining_Group
08AE; DAL WITH 3 DOTS BELOW; R; DAL
08AF; SAD WITH 3 DOTS BELOW; D; SAD
08B0; KEHEH WITH STROKE BELOW; D; GAF
08B1; STRAIGHT WAW; R; STRAIGHT WAW

Unicode character Properties:

08AD;ARABIC LETTER LOW ALEF;Lo;0;AL;;;;;N;;;;;
08AE;ARABIC LETTER DAL WITH THREE DOTS BELOW;Lo;0;AL;;;;;N;;;;;
08AF;ARABIC LETTER SAD WITH THREE DOTS BELOW;Lo;0;AL;;;;;N;;;;;
08B0;ARABIC LETTER GAF WITH INVERTED STROKE;Lo;0;AL;;;;;N;;;;;
08B1;ARABIC LETTER STRAIGHT WAW;Lo;0;AL;;;;;N;;;;;

Collation:

Characters may be sorted after similar characters.

- U+08AD ARABIC LETTER LOW ALEF should have secondary or tertiary sorting after (its appearance before a word should not affect sorting, it depends on the letters after (U+08AD isn't a letter, but a sign before them). Only in the case when there are two words that totally coincide with each other with all letters and differs only where one has a U+08AD in the beginning and the second does not should the first word be the one with U+08AD (low alef) before it and the word without it be sorted afterward the second.
- U+08AE ARABIC LETTER DAL WITH THREE DOTS BELOW should sort after U+068D ARABIC LETTER DDAHAL
- U+08AF ARABIC LETTER SAD WITH THREE DOTS BELOW should sort after U+069D ARABIC LETTER SAD WITH TWO DOTS BELOW
- U+08B0 ARABIC LETTER GAF WITH INVERTED STROKE should sort after U+06AF ARABIC LETTER GAF
- U+08B1 ARABIC LETTER STRAIGHT WAW should sort after U+06CB ARABIC LETTER VE

Confusability

Confusability issues would only arise if any of these characters were decomposed.

4. Usage and Encoding Rationale

Additions for Bashkir, Tatar

The proposed letters were used in 1910s-1920s – in 1920s in Yaña imlâ (New Orthography) orthography and several experimental transition orthographies of 1910s colloquially named *Urta imlâ* (Middle Orthography) [1][2][4][5].

For the Tatar the majority of examples are stored here:

<http://fotki.yandex.ru/users/ievlampiev/album/64267>

<http://fotki.yandex.ru/users/ievlampiev/album/67014/> (also Bashkir)

<http://fotki.yandex.ru/users/ievlampiev/album/116321/>

See fig. 1 for Tatar alphabet.

Explanation of the Bashkir-Tatar Low Alef (based on ref [1][2]): Low Alef (؛; Tatar: qalınlıq bilgäse, калыңлык билгәсе, i.e. the hard sign) is a special character marking, that a word (or at least its first part) have back vowels. It was placed in the beginning of the words only.

Low alef isn't connected to other letters and occurs in initial and standalone positions only.

It works this way: the Tatar language has 5 pairs of vowels, one of the pair mates is a diphthong: ı - e [ы- э], i - iy (i) [и - ый], o-ö [о - ө], u-ü [у - ү], a-ä [a - ә]. In the Yaña imlâ alphabet, all of them used one vowel sign for each pair, except the last, i.e. Farsi yeh-damma Farsi yeh U+06CC (ﻯ), waw-damma U+06C7 (ﻮ), waw U+0648 (ﻮ), and paired alef U+0627 (ﺀ) and Arabic ae U+06D5 (ﺀ). Possibly, the last pair was signed by different characters due to many Arabic words that don't obey the vowel harmony law. Low alef was used in places where dubious reading might occur, for example, tor (stay) and tör (variety) without it. It was placed before the word that consisted of back vowels, i.e. in "tor" (ﺗﻮﺭ) but not for "tör" (ﺗﻮﺭ). See fig 3.

In some cases it was omitted: for words containing "q" (ق) or "ğ" (غ) - they usually have back vowels, or words, containing "a" (ا), as their other vowels are back due to vowel harmony law.

Rationale for encoding low alef as a separate character can be derived from 1) its form; 2) its usage.

- 1) Unlike the subscript alef (U+0656) low alef isn't placed below base characters, but is a standalone alef placed as subscript index. Moreover, subscript alef and low alef coexisted in the Tatar writing, but in the different orthographies (see fig. 4 & 5). Low alef appears in the last Arabic orthography Yaña imlâ, and the subscript alef appears in the transition Urta imlâ experimental orthographies, developed and used by many Tatar scholars until the Yaña imlâ was adopted. Subsequently they both should occur in the same font where they will have different glyph view. Also, subscript alef couldn't be used instead of low alef, as in any ordinary font it will be shown as haraka – for historical usage it is better for reader to see the "box" than incorrect view.
- 2) Low alef has the similar behavior to high hamza (U+0674) used for Arabic in Central Asian languages. It occurs only in beginning of the word and doesn't interact with other letters of the word. Low alef is not haraka, and not a letter, it has another sense, similar to Central Asian high hamza. See Fig. 3, 9.

Additions for Crimean Tatar, Chechen and Lak

The proposed letter was used in Crimean Tatar to designate [j] (fig. 14) in some orthographies and for Chechen and Lak to designate [k'] (fig 15, 16) in 1920s.

For Crimean Tatar usage examples see <http://fotki.yandex.ru/users/ievlampiev/album/67179/> .

Additions for Belarusian

The proposed letters were used in handwriting in the 16th-19th centuries for writing Belarusian and Polish among the Belarusian Tatars. In the 20th century it was used in printing among scholars studying this alphabet [3]. See fig 7, 10.

Examples of texts are found at:

<http://fotki.yandex.ru/users/ievlampiev/album/64978>

5. Acknowledgements

Many thanks to Tatarstan National Library workers for access for copying some texts. Special thanks to Lorna Priest for valuable help in composing this proposal.

6. References

- [1] See <http://fotki.yandex.ru/users/ievlampiev/album/116321/> (2nd and 3rd lines form title and <http://fotki.yandex.ru/users/ievlampiev/view/403389?page=0> Русьмактяблярэучюн татартлэдырсе леги, Казань, 1925)
- [2] Bashkir grammar (p .III; p. 9) found on: <http://fotki.yandex.ru/users/ievlampiev/album/67014/>
- [3] Антонович А.К. Белорусские тексты, писанные арабским письмом, и их графико-орфографическая система. Вильнюс, 1968. <http://fotki.yandex.ru/users/ievlampiev/album/64978/>
- [4] Курбатов, Халиф Рэхим улы. Татар теленең алфавит һәм орфография тарихы / Х.Р. Курбатов.—Казан: Татарстан китап нәшр., 1960.—131 б.; 23.—
<URL:http://z3950.ksu.ru/knigi/X_Qurbatov_Tatar_teleneng_alfavit_ham_orfografia_tarixi.djvu>.
- [5] Әхмәров К. З. Башкорт ызыуы тарихенән = Из истории башкирской письменности: (башкорт әз әби теленең алфавите һәм орфографияһы тарихе) / проф. К. З. Әхмәров; яуаплы ред. Ж. Ф. Кйекбаев.—Өфө: Башкортостан китап нәшр. , 1972.—132, [2] б.; 21.—Тит. л. парал. рус.— Библиогр.: с. 129-133 и в подстроч. примеч.<http://fotki.yandex.ru/users/ievlampiev/album/67013/>

7. Examples and Figures

Note: The figure numbers were not changed in doing this revision. As the revision no longer deals with some characters which were not accepted by UTC #127, some figures referring to these letters were dropped. In consequence, the sequence of figure numbers contains gaps.

Fig. 1. Tatar alphabet sorting from [1]

ئا، (ا)، ئه، (ه)، ب، پ، ت، ج، چ، د، ر، ز، ژ،
 س، ش، ع، ف، ق، ك، گ، گ، ل، م، ن، ئو، (و)، ئو، (و)،
 ۋ، ه، ئ، (ئ)، ئي، (ي).

Fig. 3. Usage of words with low alef and without it [1]

ئۇ - ۋ : كۈن. تۈن. بۇلن.
 ئۈچ. ئۈزن. تۈز. تۈز.

Fig. 6. Simultaneous usage of harakas (pink frames) and modifier dammas over yeh (blue frames). [See book title on <http://fotki.yandex.ru/users/ievlampiev/view/141416?page=6>]

۷۳ محمد (عەم) نىڭ وافاتى

محمد (عەم) بەھىللەشۈپ چەچىندە: «مىن بۇگۈن سىزنىڭ دىنىڭىزنى تۈگەللەندۈردىم، دېگەن
 سۆزىدە بىر آيەت ئوقۇغان ئىدى. بۇ آيەتتەن چەھزەت ئابوبىكر محمد (عەم) نىڭ ئوزاقلامى
 ئات بولسىن، آڭلاغان ئىدى: باشقا بىر آيەتتە «ھەر بىر پەيغەمبەر وافات بولغان كېيىن
 سۆزىدە، يا محمد، وافات بولۇرسىڭىز، دېگەن مەنىدا آڭلانغان ئىدى.
 ئاندىن باشقا غەلەمەتلەر بىلەن دە محمد (عەم) نىڭ ئەجەلى ياقىلانغان بىلەن ئىدى،
 ئۇنى كۆرۈش كۈنى كۈن آرتا بارا. سەھابەلەر ئۆزىنىڭ ئۆيى ئەھلىرىنى سولۇش

Fig. 9. Usage of low alef (marked red) in the Bashkir language [2]

بىگە بىيل عىلمى مەرکەز تارافىنان باشقۇرت تلىنىڭ سارف،
 نەحوەن ياذب ئۇلگرتو تاشقۇرلعاينى. ھەزرگە سارف ئۇلۇشن ياذب
 تەمام ئىتىدك. باشقۇرت تلىنە سارف ياذو بىرنسى تەزرىبە بولۇوى،
 ئول تورالا بر تۇرلىلە قوللانما، ماتىرىيال بولماوى ھەم باشقا سەبەبتەر
 تىلەكتى تولى كۆپى ئوتەوگە يول بىرمەنى.
 مەكتەبتەر ئۇسۇن تەقدىم ئىتىلگەن سارف كىتابى فائىتىكە

ھۇيلەم يەكى ھوز تۇركۇمدەرى بولھا ئىبتەشتەر ئىسنىنەن برەوھى
 بولماھا ئىكنسى برەوھى شول ھوزدى ئاگلاب، قالعان ئىبتەشتەرلنە
 تۇشۇندۇر ب ھۇيلەب بىرە ئالا. بر كىشى ياكى ئەذرلەگەن دەرىستەن

Fig. 10. Usage of proposed letters in Belarusian text [Кітабы — унікальная зьява ў беларускай мове/ В. І. Несьцяровіч http://www.pravapis.org/art_kitabl.asp].

تىرئفوق ایدی کوز نوری کوز مذدب • خروشا اوتمن ایت دجی ورمیدی •
 صمنا بلو اوج تطلانی بوز • قوز بیل صاباى نىرحال
 یفور یقاز بل اتردی ظلمات • سانس کیر تو پییدی قیامت •
 دوزخ باشو صمنا بلو • و صر ویران •
 — Адрывак лягенды «Мэрадж» з «Аль-Кітабу».

Fig. 14. Page of a Crimean Tatar primer for reading Qur'an [title is here
<http://fotki.yandex.ru/users/ievlampiev/view/158092?page=0>]

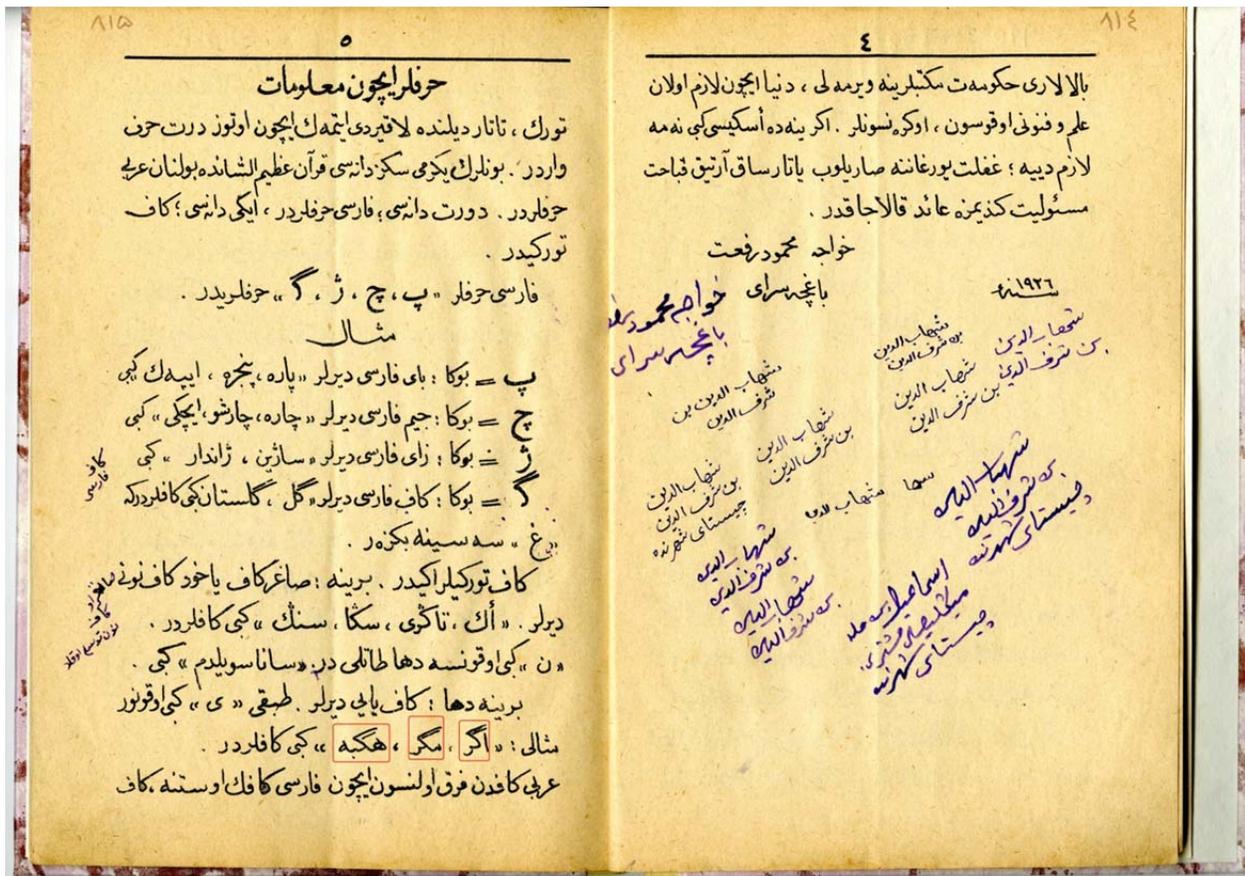


Fig. 15. Title of Chechen primer and its p. 42 [title is on the image,
<http://fotki.yandex.ru/users/ievlampiev/album/66437/>]

۴۳۶

<p>گور گور گومور گیر گیرا گهل گهیدی</p> <div style="display: flex; justify-content: space-around;">   </div> <p>گودول حاجگی</p> <p>دآدآ دآوردآ باخو وائی گاور گهل گهیدی یو ، عائشة نانا یورغو گودولایخی دالمش ، ما گه بنار دالمش ایخی ، دآورد گیرانا ماسا ده دو؟ وور ده دو ، بهتیق ماسا گیرا دو؟ دی گیرانا شیخ</p>	<div style="display: flex; justify-content: space-around;">   </div> <p>گازا بوچیگ دیگ</p> <p>عایشات - عائشة ! شون یآخان یاخانا گازانا - بوچیگنا یهئانا ، بوچیگینا یائان همیانا لوی گازا اووژا ، شورا لور احمدانا ، هوندا لوسور - آله احمدی شورا مالا .</p> <p>بوچیگ بو ---</p> <p>گ گ گ گ</p> <p>گا کاخ گانت گادیی</p>
--	--

Fig. 16. Lak Arabic alphabet table [Adyghe alphabet table, retrieved from Wikipedia 2010-05-08.

Accompanying information:

Lak arabic alphabet.JPG

English: Lak arabic alphabet from 1925 book

Source: Букварь на лакском языке. Primer auf Lak. Буйнакск, 1925 Buinaksk, 1925]

ا	ب	ت	پ	ج	ح	خ	چ
ج	د	ز	ر	س	ش	پس	ر
ط	ع	غ	ف	ق	ک	گ	گ
ک	ل	م	ن	و	ه	ی	او
ای	ث	ص	ض	ذ	ظ	ث	خ

(۱) وای قعور اگو حارف لاکو مازراو قانایستار عاراب
رایستارا چیچین داقا .
(۲) هارپا عاین جالاستا موقوو (۴) بیشایستاکوما مثلا : خوزو -
قعورو هارپا بوقواستا موقول یا لوکو (۵) دیشایستار مثلا : کارو - کارو .

**ISO/IEC JTC 1/SC 2/WG 2
PROPOSAL SUMMARY FORM TO ACCOMPANY SUBMISSIONS
FOR ADDITIONS TO THE REPERTOIRE OF ISO/IEC 10646¹.**

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

Please read Principles and Procedures Document (P & P) from <http://www.dkuug.dk/JTC1/SC2/WG2/docs/principles.html> for guidelines and details before filling this form.

Please ensure you are using the latest Form from <http://www.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html>.

See also <http://www.dkuug.dk/JTC1/SC2/WG2/docs/roadmaps.html> for latest Roadmaps.

A. Administrative

1. Title: *Revised proposal to encode Arabic characters used for Bashkir, Belarusian, Crimean Tatar, and Tatar languages*
2. Requester's name: *Ilya Yevlampiev, Karl Pentzlin, Nurlan Joomagueldinov*
3. Requester type (Member body/Liaison/Individual contribution): *Expert Contribution*
4. Submission date: *20 May 2011*
5. Requester's reference (if applicable):
6. Choose one of the following:
- This is a complete proposal: *Yes*
- (or) More information will be provided later: *No*

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 Extended-A*
2. Number of characters in proposal: *5*
3. Proposed category (select one from below - see section 2.2 of P&P document):
- | | | | | | |
|---------------------------------------|--------------------------|---|-------------------------------------|------------------------------------|--------------------------|
| A-Contemporary | <input type="checkbox"/> | B.1-Specialized (small collection) | <input checked="" type="checkbox"/> | B.2-Specialized (large collection) | <input type="checkbox"/> |
| C-Major extinct | <input type="checkbox"/> | D-Attested extinct | <input type="checkbox"/> | E-Minor extinct | <input type="checkbox"/> |
| F-Archaic Hieroglyphic or Ideographic | <input type="checkbox"/> | G-Obscure or questionable usage symbols | | | |
4. 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 P&P document? *Yes*
- b. Are the character shapes attached in a legible form suitable for review? *Yes*
5. Fonts related:
- a. Who will provide the appropriate computerized font to the Project Editor of 10646 for publishing the standard? *Lorna Priest, SIL International*
- b. Identify the party granting a license for use of the font by the editors (include address, e-mail, ftp-site, etc.):
6. 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*
7. 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*
Sorting and linguistic representations are discussed in the proposal

8. 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/UCD.html> and associated Unicode Technical Reports for information needed for consideration by the Unicode Technical Committee for

¹ Form number: N3702-F (Original 1994-10-14; Revised 1995-01, 1995-04, 1996-04, 1996-08, 1999-03, 2001-05, 2001-09, 2003-11, 2005-01, 2005-09, 2005-10, 2007-03, 2008-05, 2009-11)

inclusion in the Unicode Standard.

C. Technical - Justification

1. Has this proposal for addition of character(s) been submitted before? If YES explain	<i>No</i>
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? If YES, available relevant documents:	<i>Yes</i> <i>Linguists, librarians</i> <i>See examples in proposal</i>
3. Information on the user community for the proposed characters (for example: size, demographics, information technology use, or publishing use) is included? Reference:	<i>once used by 3 million people</i>
4. The context of use for the proposed characters (type of use; common or rare) Reference:	<i>Historically common</i>
5. Are the proposed characters in current use by the user community? If YES, where? Reference:	<i>Historical use</i>
6. After giving due considerations to the principles in the P&P document must the proposed characters be entirely in the BMP? If YES, is a rationale provided? If YES, reference:	<i>Yes</i> <i>Should be placed with similar characters</i>
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)?	
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:	<i>No</i>
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:	<i>No</i> <i>Discussion in proposal</i>
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:	<i>No</i>
11. Does the proposal include use of combining characters and/or use of composite sequences? If YES, is a rationale for such use provided? If YES, reference: Is a list of composite sequences and their corresponding glyph images (graphic symbols) provided? If YES, reference:	<i>No</i>
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)	<i>No</i>
13. Does the proposal contain any Ideographic compatibility character(s)? If YES, is the equivalent corresponding unified ideographic character(s) identified? If YES, reference:	<i>No</i>