

ISO/IEC JTC/SC2/WG2

Universal Multiple—Octet Coded Character Set (UCS)

ISO/IEC JTC/SC2/WG2N3819
2010.04.15

TITLE: Preliminary Proposal for Encoding Special Scripts and Characters in UCS for
Uighur, Kazakh and Kirgiz

SOURCE: China

STATUS: NATIONAL BODY POSITION

ACTION: Consideration by WG2, for collecting comments

In order to deal with the incompleteness of information exchanges for Uighur, Kazakh, Kirghiz and Arabic languages and occurrence of double ambiguous code, it is quite necessary to supplement Uighur, Kazakh, Kirghiz letters in table1-2.

We understand, the nominal forms with all the presentation forms, which other one different in both of their variant quantity and shape are not entirely identical, giving rise to double ambiguous code. In order to deal with ambiguity caused by letters in different language texts, which have the same shape in nominal forms, but have different quantities or shapes in deformed presentation forms, supplementation is necessary. There are some examples about it in UCS. Such as:

“ ﻱ ” (U+0649) in table 15 and “ ﻱ ” (U+06CC) in table 16;
“ ﻑ ” (U+0641) in table 15 and “ ﻑ ” (U+06A7) in table 16;
“ ﻫ ” (U+0647) in table 15 and “ ﻫ ” (U+06BE) in table 16;
“ ﺀ ” (U+0621) in table 15 and “ ﺀ ” (U+0674) in table 16;
“ ﺀ ” (U+0623) in table 15 and “ ﺀ ” (U+0622) in table 16;
“ ﺀ ” (U+0625) in table 15 and “ ﺀ ” (U+0673) in table 16;

The Nominal form shapes of these letters in corresponding pairs are same, but quantity of presentation forms or shapes are different. Thus, different nominal forms are assigned code points respectively for all of them. At present, there exist numerous nominal forms in Uighur, Kazakh and Kirghiz information processing and exchange, which cause ambiguity, since their presentation forms are different in quantity or shapes.

- 1) Table 1-1 lists Uighur, Kazakh, Kirghiz and Arabic letters for nominal form in UCS and nominal form in original shape, which are completely different in shape

and variant quantity. It certainly produces double ambiguous codes in UCS. Under serial number 2 in table 1-1, nominal form “**ئ**” will be supplemented to Uighur language, because it has four variants: isolated form (“**ئ**” and “**ه**”), final form (“**ه**” and “**ئه**”). While the nominal form of Kazakh and Kirghiz letter “**ه**”(06d5) has two variants: isolated form “**ه**” and final form “**ه**”. They are different in shape and variant quantity. If Nominal form of Kazakh and Kirghiz letter “**ه**”(06d5) replaces Uighur Nominal form “**ئ**” or isolated form “**ه**”, there appears ambiguity in information exchange as a result.

2) Likewise the nominal form of Arabic, Kazakh, Kirghiz language under serial number 1,2,3,4,5,6 of table 1-1 are different in shape comparing to nominal forms of Uighur language which have variant quantity and part of member in shapes, therefore five nominal forms for Uighur language need to be supplemented. For example, under serial number 3, Arabic letter YEH “**ي**” (U+0649) has two variants: (“**ي**” and “**ي**”), but its corresponding character in Kazakh and Kirghiz has four variants: (“**ي**”, “**ي**”, “**ي**” and “**ي**”). In Uighur, it has **Eight** variants. So Arabic letter U+0649 cannot be used to represent Uighur, Kazakh and Kirghiz letter mentioned above. Therefore, Uighur, Kazakh and Kirghiz letter “i” should be separately encoded.

We suggest to encode a character in U+076E for Kazakh and Kirghiz letter “**ي**”, and a character in U+0772 for Uighur “**ئ**”. In addition, U+0648 “**و**” is Arabic, Kazakh and Kirghiz letter, which has two variants (“**و**” and “**و**”), but it cannot be used to represent letter “**ئو**” in Uighur language which has four variants: (“**و**”, “**و**”, “**ئو**” and “**ئو**”). Although the nominal form of this letter is of Arabic, Kazakh, Kirghiz and Uighur languages, it has different variant quantity and shapes. So, it is necessary to encode Uighur letter “**ئو**” separately, and we suggest to encode it at U+076F. Some other letters such as U+06c7 “**ئو**”, U+06c6 “**ئو**” and U+0627 “**ئا**” have the same problem mentioned above. Therefore all of these letters should be encoded separately.

3) We here emphasize the importance of the nominal form of Uighur letter “**ئا**” which is an indivisible symbol and cannot be represented by combining symbol “**ئا = ئ + ا**”, it does not consist with the writing regulations of Uighur letters. Here we can see that the shape Uighur letter “**ئا**” is different to Arabic letter “**ا**”, thus it has to be separately encoded. Similarly, Uighur letters “**ئ**”, “**ئى**”, “**ئو**”, “**ئو**” and “**ئو**” are indivisible symbols and cannot be represented by combining symbol with “**ئ**”. You can see, they are different to Arabic letters “**ه**”, “**ي**”, “**و**”, “**و**” and “**و**” in shapes and quantity. Therefore, all of these letters need to be encoded separately otherwise they will cause double ambiguities.

4) The presentation forms of eight nominal forms for Uighur, Kazakh and Kirghiz languages are already encoded UCS (see FBEA -- FBFD in Table 133), but their nominal forms are not encoded.

Since 06XX and 07XX are already occupied or reserved for Arabic characters lately,

the supplemented code points for nominal forms of Uighur, Kazakh and Kirghiz languages have to be placed in 07XX in 076B-0772.

Table 1-1: Uighur, Kazakh and Kirghiz Letters that produce to Double Ambiguous Codes in 10646

Serial number	Code	Language	Nominal form in 10646	Nominal form -in original	Pronunciation	variants								
						Isolated form	Final form	Initial form	Medial form	Isolated form	Final form	Initial form	Medial form	
1	0627	Arabic Kazakh Kirghiz			[a]									
		Uighur		ا	[a]					ئا	ئا			
2	06D5	Kazakh Kirghiz	е	е	[e]	е	е							
		Uighur	е	ئە	[æ]	е	е			ئە	ئە			
3	0649	Arabic	ى	ى		ى	ى							
		Kazakh Kirghiz	ى	ى	[]	ى	ى	ر	ر					
		Uighur	ى	ئى	[i]	ى	ى	ر	ر	ئى	ئى	ئى	ئى	ئى
4	0648	Arabic Kazakh Kirghiz	و	و	[w. o]	و	و							
		Uighur	و	ئو	[o]	و	و			ئو	ئو			
5	06C7	Kazakh Kirghiz	ۇ	ۇ	[u]	ۇ	ۇ							
		Uighur	ۇ	ئۇ	[u]	ۇ	ۇ			ئۇ	ئۇ			
6	06C6	Kazakh	ۋ	ۋ	[v]	ۋ	ۋ							
		Uighur	ۋ	ئۋ	[θ]	ۋ	ۋ			ئۋ	ئۋ			

Table 1-2: Uighur, Kazakh and Kirghiz Letters to Be Supplemented

Language	Final form	Medial form	Initial form	Isolated form	Nominal form	No
Uighur	ا ما FE8E FBEB			ا ئا FE8D FBEA	ئا 076B New	1
Kazakh kirghiz	ه FEEA			ه FEE9	ه 076C New	2
Uighur	ه مە FEEA FBED			ه ئە FBEC FEE9	ئە 076D New	3
Kazakh kirghiz	ى FEF0	ه FBE9	د FBE8	ى FEEF	ى 076E New	4
Uighur	ى ئى FEF0 FBFA	ئى FBD2 ه FBE9	د ئە FBE8 FBFB	ى ئى FEEF FBF9	ئى 0772 New	5
Uighur	و ئو FEEE FBEE			و ئو FEED FBEE	ئو 076F New	6
Uighur	و ئو FBD8 FBF1			و ئو FBD7 FBF0	ئو 0770 New	7
Uighur	و ئو FBDA FBF3			و ئو FBD9 FBF2	ئو 0771 New	8

Suggested code points for letters which should be supplemented.

0750

Arabic Supplement

077F

	075	076	077
0	ي 0750	وي 0760	ي 0770
1	ث 0751	ي 0761	ي 0771
2	ب 0752	كي 0762	ي 0772
3	ت 0753	تي 0763	
4	ب 0754	كي 0764	
5	ب 0755	م 0765	
6	ب 0756	م 0766	
7	خ 0757	ن 0767	
8	ح 0758	ن 0768	
9	ط 0759	ن 0769	
A	د 075A	ل 076A	
B	ر 075B	ع 076B	
C	س 075C	ع 076C	
D	ع 075D	ع 076D	
E	ع 075E	ي 076E	
F	ع 075F	ي 076F	

FB50

Arabic Presentation Forms-A

	FB5	FB6	FB7	FB8	FB9	FBA	FBB	FCB	FBD	FBE	FBF
0	آ FB60	ا FB60	ق FB70	چ FB80	ک FB90	ن FBA0	م FBB0			و FBE0	ؤ FBF0
1	آ FB61	ا FB61	ق FB71	چ FB81	ک FB91	ن FBA1	م FBB1			و FBE1	ؤ FBF1
2	ب FB62	ث FB62	ج FB72	د FB82	گ FB92	ا FBA2				ؤ FBE2	ؤ FBF2
3	ب FB63	ث FB63	ج FB73	د FB83	گ FB93	ا FBA3			ك FBD3	ؤ FBE3	ؤ FBF3
4	ا FB64	ا FB64	ج FB74	ذ FB84	گ FB94	ة FBA4			ك FBD4	ي FBE4	ؤ FBF4
5	پ FB65	ث FB65	چ FB75	ذ FB85	گ FB95	ة FBA5			ك FBD5	ي FBE5	ؤ FBF5
6	پ FB66	ث FB66	ج FB76	ذ FB86	گ FB96	ه FBA6			ك FBD6	ا FBE6	ئي FBF6
7	پ FB67	ث FB67	ج FB77	ذ FB87	گ FB97	ه FBA7			ؤ FBD7	پ FBE7	ئي FBF7
8	پ FB68	ا FB68	ج FB78	ذ FB88	گ FB98	ه FBA8			ؤ FBD8	ا FBE8	ئي FBF8
9	پ FB69	ا FB69	ج FB79	ذ FB89	گ FB99	ه FBA9			ؤ FBD9	ه FBE9	ئي FBF9
A	پ FB6A	ث FB6A	ج FB7A	ذ FB8A	گ FB9A	ه FBAA			ؤ FBD A	ا FBEA	ئي FBFA
B	پ FB6B	ث FB6B	ج FB7B	ذ FB8B	گ FB9B	ه FBA B			ؤ FBD B	ا FBE B	ئي FBFB
C	پ FB6C	ق FB6C	چ FB7C	ر FB8C	گ FB9C	ط FBA C			ؤ FBD C	م FBE C	ي FBFC
D	پ FB6D	ث FB6D	چ FB7D	ر FB8D	گ FB9D	ط FBA D			ؤ FBD D	م FBE D	ي FBFD
E	ن FB6E	ق FB6E	ج FB7E	ک FB8E	ن FB9E	م FBA E			ؤ FBD E	ؤ FBE E	ي FBFE
F	ن FB6F	ق FB6F	ج FB7F	ک FB8F	ن FB9F	م FBA F			ؤ FBD F	ؤ FBE F	ي FBFF

Names and suggested code points for the eight letters to be supplemented:

ئا	076B	Arabic	letter	ALEF	for	Uighur
ە	076C	Arabic	letter	AE	for	Kazakh,Kirghiz
ئە	076D	Arabic	letter	AE	for	Uighur
ى	076E	Arabic	letter	YEH	for	Kazakh,Kirghiz
ئى	0772	Arabic	letter	YEH	for	Uighur
ئو	076F	Arabic	letter	WAW	for	Uighur
ئۇ	0770	Arabic	letter	U	for	Uighur
ئۆ	0771	Arabic	letter	OE	for	Uighur

Uighur alphabet for children (nominal forms of Uighur letters):



Proposed eight special letters in Uighur words:

ئا : ئالم , ئاستا , ئادالەت , ئىئانە , ئامراق
ئە : ئەمگەك , ئەقىل , ئەسئەت , ئەدەب
ئى : ئېيىق , سېسىق , دېتال , ئېرلان , ئېرا
ئى : ئىشك , ئىلمى , رۇقى , رىۋايەت , كۆزى
ئو : ئوغاق , روشەن , ئوغلاق , ئويغاق , ئوبلان
ئۇ : ئۇيغۇر , ئۇچۇر , ئۇزۇن , مەسئۇل , رۇسۇل
ئۆ : ئۆدەك , كۆل , گۆرۈ , ئۆي
ئۆ : ئۆزۈم , ئۆرۈمچى , كۆز , ئۆششۈك

Below are the pictures from Uighur alphabet for children (eight letters which should be supplemented):

ئا ا ا ئا

 ئارا	 ئاي
 ئاپئاق	 ئار

ئاپام ئايغا ئوخشايدۇ.

پ پ پ پ پ پ پ پ

 پىپىز	 پىتىز
 موشۇ كىپىق	 دىپىز

پىپىق ئورمانغا ئامراق،
ئىپىقتىن تۇرغىن يىراق.

ئە ە ە ئە

 دەرخ	 ئەينەك
 مەشئەل	 كۆڭلەك

— ئەينەك دېگەنە؟
— ئەينەك!
— يۇرمە ھەرگىز مەينەت.

ئە ە ئە ئە ئە ئە ئە ئە

 ئەككى	 ئەت
 گىلاس	 پىل

ئەككى خوراز سۆيۈشتى،
يۈز - كۆزىنى چوقۇشتى،
ئەككى ئايچاق دوست بولۇپ،
ئەك ئويلىدى چېپىشتى.

ئو و و ئو

 دۇپيا	 ئوكۇل
 رادىئو	 قول

ئون بارماقتا ئون تىرناق،
ئوك قولۇم، سول قولۇم،
ھەر قولۇمدا بەش بارماق.

ئو و و ئو

 دۇمباق	 ئوۋا
 مايىۋن	 ئاچقۇچ

ئۇيالىمۇ زىيەكشى،
بۇيالىمۇ زىيەكشى،
ياغلىقنىڭ مەلەم،
كەپىلىدىن ئەگەش.

ئو و و ئو

 كۈنلۈك	 ئۈزۈم
 گۈل	 بۈركۈت

ئۈزۈم ئۈزۈم ئۈزۈدۈم،
ئۈزۈمنى ئۈزۈم ئۈزۈدۈم.

ئو و و ئو

 تۆگە	 ئۆردەك
 بۆرە	 ئۆي

قاغا شاختا؛
— قاق - قاق - قاق؛
— سۇدا ئۆردەك؛
— غاق - غاق - غاق؛

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: Proposal for Encode Special Scripts and Characters in UCS for Uighur, Kazakh and Kirgiz

2. Requester's name: China

3. Requester type (Member body/Liaison/Individual contribution): National Body

4. Submission date: 2010.4.15

5. Requester's reference (if applicable): No

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): Yes

Proposed name of script: Uighur, Kazakh and Kirgiz

b. The proposal is for addition of character(s) to an existing block: No

Name of the existing block:

2. Number of characters in proposal: 8

3. Proposed category (select one from below - see section 2.2 of P&P document):

A-Contemporary Z B.1-Specialized (small collection) B.2-Specialized (large collection)

C-Major extinct D-Attested extinct E-Minor extinct

F-Archaic Hieroglyphic or Ideographic 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. Who will provide the appropriate computerized font (ordered preference: True Type, or PostScript format) for publishing the standard? Jpg files were provided by the Xiniang University, China. Font will be provided later.

If available now, identify source(s) for the font (include address, e-mail, ftp-site, etc.) and indicate the tools used: wushour@xju.edu.cn, chenzh@cesi.ac.cn, chenzh-zhuang@163.com

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)? No

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 inclusion in the Unicode Standard.

¹ Form number: N3152-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)

C. Technical - Justification

1. Has this proposal for addition of character(s) been submitted before? If YES explain	<i>WG2n2820, WG2#45</i>	<i>Yes</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>People in Xinjiang Autonomous Region, China.</i>	<i>Yes</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>For example, Table2 Orkhun-Yenisey Alphabet, Table-3 Old Turk(Orkhun)Alphabet, Table-4: "Iski Turk Yazitlari", Yusayin Namiq Orqun, Yurkiye, Istanbul</i>	<i>Yes</i>
4. The context of use for the proposed characters (type of use; common or rare) Reference:	<i>People through Interner</i> <i>For example ,Uyghur alphabet.</i>	
5. Are the proposed characters in current use by the user community? If YES, where? Reference:	<i>For example, Xinjiang Uighur Autonomous Region, China.</i>	<i>Yes</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>No</i>
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)?		<i>Yes</i>
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>See table1-1,1-2,FB50 Axabic Presentation from S-A</i>	<i>Yes</i> <i>Yes</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>
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>