

ISO/IEC JTC/SC2/WG2

Universal Multiple—Octet Coded Character Set (UCS)

ISO/IEC JTC/SC2/WG2 N3919
2010.9.15

TITLE: Proposal to Encode Special Scripts and Characters in UCS for Uighur language
SOURCE: China
STATUS: NATIONAL BODY POSITION
ACTION: CONSIDERATION BY WG2 Preliminary Proposal, for collecting
comments

In order to deal with the incompleteness of information exchanges for Uighur languages , it is quite necessary to supplement **Eight** special scripts of Uighur letters in ISO/IEC 10646 /Unicode (table1-2, below).

Problame 1: 1) We have emphasize the importance of the nominal form of Uighur character “ئ” is a indivisible symbol, and that is do not represented by combining symbol “ل+ه=له ”, it does not consist with the writing regulations of Uighur letters. We emphasize again it is a completely single and a indivisible symbol . At the same time, the shape Uighur letter “ئ” is different to Arabic letter “ل”, thus it has to be separately encoded. Similarly, Uighur character “ئە”, “ئى”, “ئو”, “ئۇ”, “ئۆ” and “ئۈ” are indivisible symbols and cannot be represented by combining symbol with “ه”, or not dual-joining with U+0626 ARABIC LETTER YEH WITH HAMZA ABOVE. Other one You can see, they are different to Arabic letters “ه”, “ى”, “و”, “ۆ”, “ۇ” and “ۈ” in shapes and quantity. So These characters represented by sequence of two existing characters is a completely wrong action, such as NFS.
“05F7” = U+0626 ARABIC LETTER YEH WITH HAMZA ABOVE and U+0627 ARABIC LETTER ALEF
“05F8” = U+0626 ARABIC LETTER YEH WITH HAMZA ABOVE and U+06D5 ARABIC LETTER AE
“05FB” = U+0626 ARABIC LETTER YEH WITH HAMZA ABOVE and U+0648 ARABIC LETTER WAW
“05FC” = U+0626 ARABIC LETTER YEH WITH HAMZA ABOVE and U+06C7 ARABIC LETTER U
“05FD” = U+0626 ARABIC LETTER YEH WITH HAMZA ABOVE and U+06C6 ARABIC LETTER OE
“05FA” = U+0626 ARABIC LETTER YEH WITH HAMZA ABOVE and U+0649 ARABIC LETTER ALEF MAKSURA

Therefore, all of these letters need to be encoded separately.

2) The presentation forms of eight nominal forms for Uighur, languages are already encoded UCS (see FBEA -- FBFD in Table 133), but their nominal forms are not encoded till now.

Problame 2: There exist numerous nominal forms in Uighur, Kazakh , Kirghiz and Arabic 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 “ه” will be supplemented to Uighur language, because it has four variants: isolated form ("ه" and "ه"), final form ("ه" and "ه"). 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 "ه" 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 **Eight** 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 letter “ي” “ي” should be separately encoded suggestion code points to U+05FA.

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 “و” in Uighur language which has four variants: (“و”, “و”, “و” and “و”). 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 “و” separately, and we suggest to encode it at U+05FB. Some other letters such as U+06c7 “و”, U+06c6 “و” and U+0627 “و” have the same problem mentioned above. Therefore all of these letters should be encoded separately.

The Suggested code points for nominal forms of Uighur languages have to be in 05XX 05f7-05ff (may be arrange in other code points) .


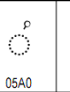
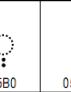




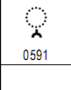
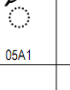
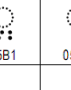
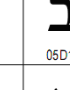
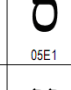
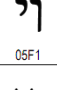

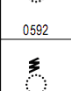
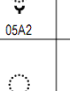
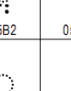
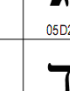
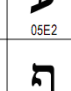
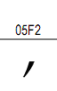



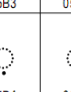
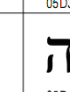












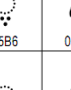
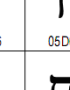
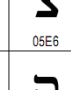


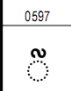
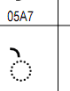
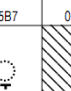
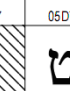
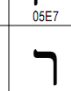



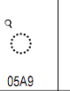





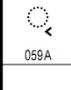
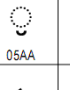
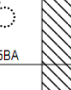
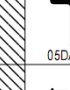



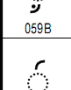

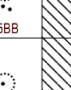














































Table 1-1: Some Uighur, Kazakh ,Kirghiz and Arabic Letters that produce to Double Ambiguous Codes in 10646

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	ا	ا	[a]	ا	ا						
		Uighur	ا	ئا	[a]	ا	ا			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	د	د				
		Uighur	ى	ئى	[i]	ى	ى	د	د	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	ۆ	ۆ	[v]	ۆ	ۆ						
		Uighur	ۆ	ئۆ	[θ]	ۆ	ۆ			ئۆ	ئۆ		

Table 1-2: Special Scripts to Be Supplemented in ISO/IEC 10646 /Unicode for Uighur languages

Language	Final form	Medial form	Initial form	Isolated form	Nominal form	No
Uighur	ا ئا FE8E FBEB			ا ئا FE8D FBEA	ئا 05F7 New	1
Uighur	ه مھ FEEA FBED			ئھ ه FBEC FEE9	ئھ 05F8 New	2
Uighur	پ ئپ FBE5 FBF7	پ FBD1 پ FBE7	پ ئپ FBE6 FBF8	پ ئپ FBE4 FBF6	ئپ 05F9 New	3
Uighur	ى ئى FEF0 FBFA	ئى FBD2 ئى FBE9	ى ئى FBE8 FBFB	ى ئى FEEF FBF9	ئى 05FA New	4
Uighur	و ئو FEEE FBEF			و ئو FEED FBEE	ئو 05FB New	5
Uighur	ۇ ئۇ FBD8 FBF1			ۇ ئۇ FBD7 FBF0	ئۇ 05FC New	6
Uighur	ۈ ئۈ FBDA FBF3			ۈ ئۈ FBD9 FBF2	ئۈ 05FD New	7
Uighur	ۋ ئۋ FBDC FBF5			ۋ ئۋ FBDB FBF4	ئۋ 05FE New	8

The Suggested code points for **Special Scripts** of Uighur languages have to be in 05XX 05f7-05ff (may be arrange in other code points) .

	059	05A	05B	05C	05D	05E	05F
0							
1							
2							
3							
4							
5							
6							
7							
8							
9							
A							
B							
C							
D							
E							
F							

FB50
Arabic Presentation Forms-A

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

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

ئا	05F7	Arabic	letter	ALEF	for	Uighur
ئە	05F8	Arabic	letter	AE	for	Uighur
ئى	05F9	Arabic	letter	YEH	for	Uighur
ئې	05FA	Arabic	letter	YEH	for	Uighur
ئو	05FB	Arabic	letter	WAW	for	Uighur
ئۇ	05FC	Arabic	letter	U	for	Uighur
ئۆ	05FD	Arabic	letter	OE	for	Uighur
ئۈ	05FE	Arabic	letter	OU	for	Uighur

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



Proposed eight special letters in Uighur words:

ئا : ئالىم, ئاستا, ئادالەت, ئىئانە, ئامراق
ئە : ئەمگەك, ئەقىل, ئەسئەت, ئەدەب
ئى : ئېيىق, سېسىق, دېتال, ئېرلان, ئېرا
ئى : ئىشك, ئىلمى, رۇقىي, رىۋايەت, كۆزى
ئو : ئوغاق, روشەن, ئوغلاق, ئويغاق, ئويلان
ئو : ئويغۇر, ئوچۇر, ئوزۇن, مەسئۇل, رۇسۇل
ئۆ : ئۆدەك, كۆل, گۆرۈ, ئۆي
ئۈ : ئۈزۈم, ئۈرۈمچى, كۈز, ئۈششۈك

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



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 to Encode Special Scripts and Characters in UCS for Uighur language
2. Requester's name: China
3. Requester type (Member body/Liaison/Individual contribution): Member body
4. Submission date: 2010-09-21
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:

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: Uighur
- b. The proposal is for addition of character(s) to an existing block: Yes
- Name of the existing block: Uighur
2. Number of characters in proposal: 8
3. Proposed category (select one from below - see section 2.2 of P&P document):
- | | | |
|----------------------------------------------------------------|------------------------------------------------------------------|-------------------------------------------------------------|
| A-Contemporary <input checked="" type="checkbox"/> | B.1-Specialized (small collection) <input 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 <input type="checkbox"/> | |
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?
Jpg files were provided by the Xinjiang University, China. Font will be provided later.
- b. Identify the party granting a license for use of the font by the editors (include address, e-mail, ftp-site, etc.):
chenzh@cesi.ac.cn wushour@xju.edu.cn
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.

C. Technical - Justification

1. Has this proposal for addition of character(s) been submitted before?	Yes
If YES explain	WG2n2820, WG2#45 WG2n3819, WG2#56
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.)?	Yes
If YES, with whom?	People in Xinjiang Autonomous Region, China.
If YES, available relevant documents:	
3. Information on the user community for the proposed characters (for example: size, demographics, information technology use, or publishing use) is included?	No
Reference:	
4. The context of use for the proposed characters (type of use; common or rare)	common
Reference:	For example in education: Uighur alphabet for children.
5. Are the proposed characters in current use by the user community?	Yes
If YES, where? Reference:	For example, Xinjiang Uighur Autonomous Region, China.
6. After giving due considerations to the principles in the P&P document must the proposed characters be entirely in the BMP?	No
If YES, is a rationale provided?	
If YES, reference:	
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)?	Yes
8. Can any of the proposed characters be considered a presentation form of an existing character or character sequence?	Yes
If YES, is a rationale for its inclusion provided?	Yes
If YES, reference:	See problem 1 in the proposal
9. Can any of the proposed characters be encoded using a composed character sequence of either existing characters or other proposed characters?	No
If YES, is a rationale for its inclusion provided?	
If YES, reference:	

10. Can any of the proposed character(s) be considered to be similar (in appearance or function) to an existing character?	No
If YES, is a rationale for its inclusion provided?	
If YES, reference:	
11. Does the proposal include use of combining characters and/or use of composite sequences?	No
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:	
12. Does the proposal contain characters with any special properties such as control function or similar semantics?	No
If YES, describe in detail (include attachment if necessary)	
13. Does the proposal contain any Ideographic compatibility character(s)?	No
If YES, is the equivalent corresponding unified ideographic character(s) identified?	
If YES, reference:	