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

Doc Type: Working Group Document

Title: Proposal to encode four Latin letters for Janalif Source: Karl Pentzlin, Ilya Yevlampiev (Илья Евлампиев)

Status: Individual Contribution

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

Date: 2008-11-03, revised 2009-03-16

Revision history: The revision of 2009-03-16 takes into account the code points (U+A790/U+A791) devised by UTC #117 for the n with descender. Moreover, it takes into account the name "Latin capital/small letter yeru" for the letter initially proposed as "Latin capital i with right bowl / Latin small letter dotless i with right bowl", as proposed by Michael Everson and continued by the German comments to PDAM7. Also, some sorting considerations were added for the Latin yeru, and fig. 6 was updated.

Additions for Janalif

N U+A790 LATIN CAPITAL LETTER N WITH DESCENDER

→ 04A2 cyrillic capital letter n with descender

N U+A791 LATIN SMALL LETTER N WITH DESCENDER

b U+A792 LATIN CAPITAL LETTER YERU

- → 042B cyrillic capital letter yeru
- → 042C cyrillic capital letter soft sign
- → 0184 latin capital letter tone six
- Ь U+A793 LATIN SMALL LETTER YERU
 - → 0131 latin small letter dotless i

Properties:

A790; LATIN CAPITAL LETTER N WITH DESCENDER; Lu; 0; L;;;; N;;; A791; A791; LATIN SMALL LETTER N WITH DESCENDER; Ll; 0; L;;;; N;;; A790;; A790 A792; LATIN CAPITAL LETTER YERU; Lu; 0; L;;;; N;;; A793; A793; LATIN SMALL LETTER YERU; Ll; 0; L;;;; N;;; A792;; A792

1. The Janalif alphabet (fig. 3, 4)

In 1908–1909 the Tatar poet Säğit Rämiev started to use the Latin alphabet in his own works. He offered the use of digraphs: ea for ä, eu for ü, eo for ö and ei for ı. But Arabists turned down his project. In the early 1920s Azerbaijanis invented their own Latin alphabet, but Tatarstan scholars set a little store to this project, preferring to reform the İske imlâ (en.wikipedia.org/wiki/iske_imla). The simplified İske imlâ, known as Yaña imlâ (en.wikipedia.org/wiki/yana imla) was used from 1920–1927. [1]

But Latinization was adopted by the Soviet officials and the special Central Committee for a New Alphabet was established in Moscow. The first project of the Tatar-Bashkir Latin alphabet was published in Eşce (The Worker) gazette in 1924. The pronunciation of the alphabet was similar to English, unlike the following. Specific Bashkir sounds were written with digraphs. However, this alphabet was declined. [1]

In 1926 the Congress of Turkologists in Baku recommended to switch all Turkic languages to the Latin alphabet. Since April of 1926 the Jana tatar əlifbasь/Yaña Tatar älifbası (New Tatar alphabet) society started its work at Kazan. [2]

Since 3 July 1927, Tatarstan officials have declared Janalif as the official script of the Tatar language, replacing the Yaña imlâ script. In this first variant of Janalif (acutes-Janalif), there weren't separate letters for K and Q (realized as K) and for G and Ğ (realized as G), V and W (realized as W). Ş (sh) looked like the Cyrillic letter Ш (she). C and Ç were realized as in Turkish and the modern Tatar Latin alphabet and later were transposed in the final version of Janalif. [1]

In 1928 Janalif was finally reformed and was in active usage for 12 years (see fig. 3, 4). This version of Janalif is the base of our proposal.

Some sources claim this alphabet having 34 letters, but the last was a digraph bj, used for the corresponding Tatar diphthong. [1] Another source states that the 34th letter was an apostrophe. They also give another sorting of the alphabet. (Θ after A, b after E) [2]

In 1939 Cyrillization of USSR was initiated. As was said, alphabet was switched to Cyrillic "by labor's request."

There are also several projects of Cyrillization. Ilminski's alphabet was already forgotten and it couldn't be used, due to its religious origin. As early as 1938 professor M. Fazlullin introduced an adaptation of the Russian alphabet for the Tatar language, without any additional characters. Specific Tatar letters should be signed with the digraphs, consisting of similar Russian letters and the letters \mathbf{b} and \mathbf{b} . [1]

In 1939 Qorbangaliev and Ramazanov offered their own projects that planned to use additional Cyrillic characters. Letters Θ, Θ, Y, h were inherited from Jaŋalif, but Җ and Ң were invented by analogy with Щ and Ц. Гъ and Къ should be used to designate Ğ and Q. By this project "ğädät" ("custom") was spelled as "гъэдәт", "qar" ("snow") as "къар". In Ramazanov's project W (Jaŋalif V) was marked by B before the vowel, and У, Y in the end of syllable. Jaŋalif: vaq - вак; tav - тау; dəv - дәү. In 5 May 1939 this project was established as official by the Supreme Soviet of TASSR. Surprisingly, "Tatar society disagreed to this project" and during 1940 July conference Cyrillic alphabet was finally standardized. 10 January 1941 this project was passed. According to this version, "ğädät" was spelled as "гадәт", "qar" as "кар". The principles were following: if ra/ro/ry/rь/ка/ко/ку/кы/ is followed by "soft syllable", containing "ə, e, e, и, ү" or soft sign "ь", they are spelled as ğä/gö/gü/ge/qä/qö/qü/qe, in other cases as ğa/go/gu/gi/qa/qo/qu/qı. rə/rə/ry/re/кə/кө/кү/ке are spelled as gä/gö/gü/ge/kä/kö/kü/ke. Similar practice were applied for e, ю, я, that could be spelled as ye, yü, yä and as yı, yu, ya. Examples: канәгать - qänäğät (satisfied); ел - yıl (year); ямь - yäm (charm). So, in Tatar Cyrillic soft sign hasn't sense of iotation, as in Russian, but a sense of vowel harmony. Unlike modern Russian, some words can end with ъ, to sign a "hard g" after the "soft vowel", as in балигь - baliğ (of the full legal age). [1]

All Russian words are written as in Russian and should be pronounced with Russian pronunciation.

In the 1990s some wanted to restore Janalif, or Janalif+W, as being corresponding to modern Tatar phonetics. But technical problems, such as font problems and the disuse of Uniform Turkic alphabet among other peoples forced to use "Turkish-based alphabet". In 2000 that alphabet was adopted by the Tatarstan government, but in 2002 it was abolished by the Russian Federation. [1]

2. The N with descender



Fig. 2 - Scan from [1]

The descender of the proposed letters U+A790/U+A791 LATIN CAPITAL (resp. SMALL) LETTER N WITH DESCENDER look like the descenders of e.g. U+2C67/U+2C68 LATIN CAPITAL (resp. SMALL) LETTER H WITH DESCENDER.

Therefore, the names proposed here were selected according to this example.

In current citations of Janalif texts, these letters are usually replaced by U+014A/U+014B LATIN CAPITAL (resp. SMALL) LETTER ENG, as these letters have a superficial but recognizable similarity to the correct Janalif letter, and as they are usually attributed to the same sound.

Also, the letter's usage was considered in 2000 Tatar Latin alphabet. Only some Tatar fonts use this glyph at the position of \tilde{N} .

Nevertheless, their form is distinctive and clearly different from the eng, which is also distinctive (even for the upper case eng of which all glyph variants concur in the form of their lower right appendage). The lower right appendage of the n with descender is always straight and placed right of the right n stem, while the lower right appendage of the eng is always a prolongation of the right n stem and bound

inwards.

Thus, the n with descender is no glyph variant of the eng.

If it were so, the letters U+0220/U+019E LATIN CAPITAL (resp. SMALL) LETTER N WITH LONG LEG also had to be regarded as glyph variant of the eng, as they in fact are more similar (the lower right appendage being straight but a prolongation of the right n stem).

Additional, the N with descender was used in parallel to the eng in the Latin alphabet used to the Khanty language about 1931-1936 (fig. 5).

Thus, it is a separate letter from eng in any case.

3. The Latin yeru



Fig. 1 - Scan from [1]

While the proposed U+A792 "LATIN CAPITAL LETTER YERU" (with its lower case counterpart U+A793 "LATIN SMALL LETTER YERU") looks like the Cyrillic letters U+042C/U+044C CYRILLIC CAPITAL (resp. SMALL) LETTER SOFT SIGN, it is by no ways a soft sign and never used as such in Janalif context.

In fact, it is a Latin equivalent to U+042B/U+044B CYRILLIC CAPITAL (resp. SMALL) YERU. Thus, it is an "i" variant by function, equivalent to the Turkish/Azerbaijani dotless i.

(The proposed naming does not prevent anybody from using the character as soft sign in nonstandard Cyrillic transcriptions or transliterations, as anybody is free to use any letters in any way.)

The letter is obviously different from the superficially similar U+0184/U+0185 LATIN CAPITAL (resp. SMALL) LETTER TONE SIX, where the vertical stem is terminated at the top by a distinctive slanted appendage, and where both capital and small form have cap-height and are distinguished by the lateral extension of the bowl.

Using the Cyrillic U+042C/U+044C as substitute in current citations of Janalif text (as it is in fact be done now due to the lack of an encoded Latin b/b), is as undesirable as having to use U+0420/U+0440 CYRILLIC CAPITAL (resp. SMALL) LETTER ER to denote the "p" in Latin text, as a substitute for a (hypothetically) not encoded U+0050/U+0070 LATIN CAPITAL (resp. small) LETTER P.

There also some points shall be noted which are similar to the situation of the Kurdish W/w [3], which was encoded at last (U+051C/051D). As pointed out above, Janalif is a stable alphabet, used for several years for several languages beyond Tatar, with a definitve sorting order: the yeru is the last letter in that alphabet after Z and Ξ (as long as the diphtong $\mathfrak b$ j is not considered). Since Tatar, over its history, is written in the Latin as well as in the Cyrillic alphabet, a multilingual wordlist cannot sort Kurdish correctly because the $\mathfrak b$ -looking letter (beyond its complete different function) cannot be in two places at the same time. (Sorting here means ordinary plain-text sorting, for instance of files in a directory.) Expecting Janalif users to have recourse to special language-and-script tagging software for these two letters alone is simply not a credible defense for the retention of the unification of two letters with complete different function.

4. References:

- [1] (Russian) М.З. Закиев. Тюрко-татарское письмо. История, состояние, перспективы. Москва, "Инсан", 2005
- [2] "Яңалиф". Tatar Encyclopedia. (2002). Kazan: Tatarstan Republic Academy of Sciences Institution of the Tatar Encyclopaedia.
- [3] Michael Everson et al., "Proposal to encode additional Cyrillic characters in the BMP of the UCS" (2007-03-21). Unicode document L2/07-003R; SC2/WG2 document N3194R.

5. Examples

A a	Вв	CC	Çç	Do
et a	86	80	80	20 d
೮	ب	₹	3	د
Еe	Эә	Ff	G g	0) 0
бе	9.	F. f	8 9	'
*.5 *.5	4€	ف	5	3
H h	l i	Ji	Kk	LI
H h	di	2 ;	Kh	21
da	ئي	ی (ثای)	ك	J
Mm	Nn	N, n,	0 0	0 е
M m	Nn		0.	0.
r	ن	ů.	ائۇ (زاق)	ۇ (رەك)
Pp	Qq	Rr	Ss	\$ \$
Pp	09	Ri	So	
ý	3	ر		ش
Yy	Tt	Uu	Vv	XX
34	Te	U u	or 2	X z
ئو (لەن)	ت	ائو (راق)	ر (ژ)	Ś
Zz	7 2	Ьь	Ьј ьј	
H z	X z	Бъ	bjbi	1000
;	ا ر	:5,	±,	

30 рос. Берлопперештон яналиф нятезеццо татар алфавиты («Яналиф», 1928, № 8) [Курбатов Х. Татар одоби теленец алфавит hом орфография тарихы.—Казац, 74

Fig. 3: Table of Jaŋalif, from [1]

	Приблизит, значение		Приблизит. эпачение
Aa	a	Nn.	**HL.,,
Вв	6	00	.0
Cc	ч	90	нак 1,0
Çç	дж	Pp	п
Dd	Д	Qq	,,К задненебное твердое
Ee	э	Rr	р
99	"а" широкое	Ss	С
Ff	ф	Şş	ш
Gg	Γ.	Tt	Т
Olo	,,Г ⁴⁴ фрикативное задненебное	Uu	у
Hh	неменкое ", ће	Vv	В
li	И	Хx	X
Jj	й	Уу	как помецкое "Д"
Kk	К	Zz	3
LI	л	Z₹	ж
Mm	M	Бь	ы
Nn	H.		

															1	9.	32	-19.	36														
a	В	c	d	e	Э	f	g	h	h	i	Ь	j	k	1	ł	ļ	ļ	m	n	ŋ	o	Θ	p	r	S	ş	s	t	u	v	Z	Z	Z,

Fig. 5: Table of the Latin alphabet used 1932-1936 for the Khanty language, showing the n with descender and the eng side by side as different letters.

*Retrieved 2008-10-31 from http://upload.wikimedia.org/wikipedia/commons/9/9c/Hanti_latin_alphabet.jpg

```
- <lta:var ty="variant" su="baku1926" ad="2007-04-18">
        <lt>
        <ld></ld>
        <l
```

<lta:pref>ky</lta:pref>
<lta:pref>sah</lta:pref>
<lta:pref>tk</lta:pref>
<lta:pref>tt</lta:pref>

<lta:pref>uz</lta:pref>

</lta:var>

Fig. 6: Entry in http://www.w3.org/2008/05/lta/lsr.xml (as of 2009-03-16). It shows the Latin yeru in a registry entry (Əlifbasь with transliteration Elifbasi, using the δ as well as the η as substitutes for the correct Janalif characters, as such a database is by nature confined to already encoded Unicode characters).



Fig. 7: Title page from a Kazhak newspaper from about 1937, showing all proposed letters. Retrieved 2008-10-25 from http://en.wikipedia.org/wiki/lmage:Sotsijaldy_qazaqstan.jpg.

The descender of the lower case n with descender shows a drop-like form here in the headline font, showing that the letter has developed some glyph variants during the time of its use.

Şavla, jahan, antal, jəş jerək, alqa, Qotqar miljon jənde ylemdən, Kyrhen jəm.e tormoş jəş salalar, Jalqan jərəgendən, qulundan. Kidən əzəlgə ser nəqrə menən, Şavlahan avaldar, qalalar. Saq jardamoqa jegərəp şəfqət menən Qotolhondar miljon salalar.

Bel jeldarda Bulat jarem agitator şaqir Bulep, Başlesa, asleqa qarşe kerəş ojoştorov esen agitatsejon şiqler həm nəberdər jada. Ləkin unen qajhe Ber əbərdərendə tatar Burzeva ədəBejətenen joqondoho urne-urne menən neq qena saqelep quja. Mibalqa unen Bel jeldarda jadqan əbərdərenən Rəsəj Jəş Kommunistar Sojuzenen Başqortostan Olkə Komi-

Fig. 8: Example from a Bashkir text of the Janalif era. While there are a lot of easy to find Latin yerus, some n with descender are encircled in red.

(The letters encircled in cyan are special Bashkir Latin letters which are unencoded yet but not subject of this proposal.)

Retrieved 2008-10-28 from

http://ru.wikipedia.org/wiki//Википедия:Проект:Внесение_символов_алфавитов_народов_России_в_Юникод Picture reference: http://ru.wikipedia.org/wiki//Изображение:Ваshqortalifba.jpg

Beldergənnən son 15 kennən də sonqa qalmışıca administratsiə "dublikat" digən jazu belən jana xezmət knəgəse birə.

13. Xezmət knəgələre biry ocen alınqan tyləvnen, şulaj uq xezmət knəgələren juqaltqan ocen alınqan ştraflarının boten suması dəvlət doxoduna kerə.

14. Xezmət knəgələrennən zakonsız fajdalanqan өсеп, alarnı ваşqа keşelərgə вігдәп өсеп, poddelka өсеп həm alarnı təzətep jazqan өсеп ugolovnı tərtiptə çəza вігеlə.

15. Predpriətielər həm ucrezdenielər xezmət knəgələren tieşle

narkomatlardan həm ucrezdenielərdən alalar.

16. "Xezmət isemlekləre turыnda" SSR Sojuzь Xalьq Komisarlarь Sovetының 1926 псь jыl 21 псе sentəвеr qərагь (SSSR zakonnar çіыптықы, 1926 псь jыl, № 66 statiə 502; 1929 псь jыl, № 35, statiə 315)—juqqa сыңағына.

SSR Sojuzb Xalbq Komisarlarb Sovetb Predsedatele V. MOLOTOF.

SSR Sojuzb Xalbq Komisarlarb Sovetb Eşləren Idarə İtyce I. BOL'ŞAKOF.

Məskəv, Kreml. 20 dekaвг, 1938 jыl.

На Татарск. яз.

Fig. 9: Scan from the workbook (Трудовая книжка - Xezmət knəgəse) from В.П. Емельянов, the grand-grandfather of one of the authors of this proposal (I.Ye.), about 1938.

This example shows many Latin yerus and some n with descender (e.g. the last letter of the second word of the first line). — By the way, this example also shows the use of U+0299 LATIN SMALL CAPITAL LETTER B as lower case counterpart for U+0042 LATIN CAPITAL LETTER B (see e.g. the first word in the second line), as it came into use for Janalif to make the b dissimilar from the Latin yeru.

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.

Α.	Administrative
_	Title: Proposal to encode four Latin letters for Jaŋalif Requester's name: Karl Pentzlin, Ilya Yevlampiev
	Requester's name: Karl Pentzlin, Ilya Yevlampiev Requester type (Member body/Liaison/Individual contribution): Individual Contribution Submission date: 2008-11-03, revised 2009-03-16
	Requester's reference (if applicable):
6.	Choose one of the following:
	This is a complete proposal: Yes
	(or) More information will be provided later:
	Technical – General
1.	Choose one of the following: a. This proposal is for a new script (set of characters): Proposed name of script:
	b. The proposal is for addition of character(s) to an existing block: Yes
	Name of the existing block: Latin Extended-D
2.	Number of characters in proposal:
3.	Proposed category (select one from below - see section 2.2 of P&P document): A-Contemporary B.1-Specialized (small collection) X 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? a. If YES, are the names in accordance with the "character naming guidelines" in Annex L of P&P document? Yes Yes
5.	b. Are the character shapes attached in a legible form suitable for review? Who will provide the appropriate computerized font (ordered preference: True Type, or PostScript format) for publishing the standard? If available now, identify source(s) for the font (include address, e-mail, ftp-site, etc.) and indicate the tools used: http://www.pentzlin.com/proposalfont.zip (more information in the info.txt file included in that archive)
	References: a. Are references (to other character sets, dictionaries, descriptive texts etc.) provided? 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
St that Ex inf Co re se	Additional Information: bmitters are invited to provide any additional information about Properties of the proposed Character(s) or Script at will assist in correct understanding of and correct linguistic processing of the proposed character(s) or script. amples of such properties are: Casing information, Numeric information, Currency information, Display behaviour ormation such as line breaks, widths etc., Combining behaviour, Spacing behaviour, Directional behaviour, Default illation behaviour, relevance in Mark Up contexts, Compatibility equivalence and other Unicode normalization ated information. See the Unicode standard at http://www.unicode.org/Public/UNIDATA/UCD.html and associated Unicode Technical Reports for information eded for consideration by the Unicode Technical Committee for inclusion in the Unicode Standard.

Proposal to encode four Latin letters for Janalif — 2009-03-16

 $^{^{1} \ \}text{Form number: N3152-F (Original 1994-10-14; Revised 1995-01, 1995-04, 1996-04, 1996-08, 1999-03, 2001-05, 2001-09, 2003-11, 2005-01, 2005-09, 2005-10, 2007-03, 2008-05)}$

C. Technical - Justification

Has this proposal for addition of character(s) been submitted before?	No
If YES explain	
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? One of the authors (I. Ye.) is himself a member of the user	community
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?	see text
Reference: see text	
Reference: See text 4. The context of use for the proposed characters (type of use; common or rare)	common
Reference: common within their context (see text)	
5. Are the proposed characters in current use by the user community?	historical
If YES, where? Reference: see text	
6. After giving due considerations to the principles in the P&P document must the proposed character	ers be entirely
in the BMP?	Yes
If YES, is a rationale provided?	Yes
If YES, reference: Keeping in line with other Latin characters	
7. Should the proposed characters be kept together in a contiguous range (rather than being scatter	ed)? Yes
8. Can any of the proposed characters be considered a presentation form of an existing	,
character or character sequence?	No
If YES, is a rationale for its inclusion provided?	
If YES, reference:	
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?	Yes
If YES, is a rationale for its inclusion provided?	Yes
If YES, reference: See text (in short: resembles a Cyrillic character in form but r	not in function)
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) prov	ided? <i>n/a</i>
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)	
True, document docum (mondad diladimination in modestical)	
13. Does the proposal contain any Ideographic compatibility character(s)?	No
If YES, is the equivalent corresponding unified ideographic character(s) identified?	740
If VES reference:	
ii TES, releience.	