

Unicode 4.1 and Slavic Philology Problems and Perspectives (II)

0. Introduction

In the first paper devoted to the current state of encoding Slavic characters into Unicode (KEMPGEN 2006)¹ we gave an introduction to the topic, outlined various achievements in the current version of Unicode (v. 4.1, 2005), mainly because not all of them are common knowledge or available in standard fonts like *Times* or *Times New Roman*, and then we discussed many of the historic or special characters of the Slavic languages that have not yet been encoded in Unicode. Space did not permit us to cover everything in the first article, so this second article will take up where the first paper left, and will also present more illustrative material for some of the topics from the first article. The reader might find it useful to be familiar with the first article before reading the present one, but knowledge of the first article is not a prerequisite.

As for the first paper, it should be noted that the purpose of this paper is *not* to claim that all characters mentioned in this article *should* be encoded in Unicode, but merely to point out those areas where further investigation is needed, where a common understanding of the principles and practice of treating characters should be developed among Slavists. The paper is also meant to be a contribution to formal proposals which will be submitted to Unicode, Inc. to have more characters encoded.

One important thing that is also worth repeating is that if a character, may it be a letter, an accent or some other sign, is defined in Unicode, that does not mean that a certain *font* contains an image of that character. Thus, it would be a misconception to think that because font *X* does not have character *Y*, *Y* is not available in Unicode. Most people today will, for example, use *Times New Roman* (from Monotype) on the PC (and now on the Macintosh, too), and *Times* (from Linotype) on the Macintosh. A comparison of these two common fonts (see Appendix) will show that, unfortunately, *Times New Roman* has less to offer to a Slavist than *Times*, but *Times* is not as widespread on the PC whereas on the Macintosh it is a standard since the invention of desktop publishing in the 80's (although of course only now under OS X with an extended Unicode character set).²

1 Available electronically from the 'Kodeks' server: <http://kodeks.uni-bamberg.de/>

2 It might be worth pointing out that version 3.05 of both *Times New Roman* and *Arial* for the Macintosh have the same character set as their PC counterparts. Thus, any data and file exchange that is based on these fonts will be completely without problems.

To make up for the deficiencies in the standard fonts mentioned above, the author of the present article has produced a font named *Kliment Std* that is available for free. This font is aimed especially at Slavic medievalists, and it features a lot of those characters not present in the above-mentioned reference fonts. ‘Kliment Std’ is available for download from

<http://kodeks.uni-bamberg.de/AKSL/Schrift/KlimentStd.htm>

and also from the ‘Repertorium’ website. It is being used in the present paper where the standard fonts do offer support for a character in question.

1. Unicode Blocks

For practical purposes, Unicode, a table consisting of 65.536 cells in which each character has its own unique number and also a name, is normally organized into “blocks”. Thus, for example, all Cyrillic characters for the Slavic languages form one block, Soviet additions to the Cyrillic alphabet form another block etc.

uni040	llocyril	djecyr	gjecyr	ecyrill	dzecyr	icyrill	vicyril	jecyril	ljecyr	njecyr	tshecy	kjecyr	uni040	lshort	dzhecy
È	Ë	Ђ	Ѓ	Є	Ѕ	І	Ї	Ј	Љ	Њ	Ћ	Ќ	Й	Ў	Ц
Acyril	Becyri	Vecyri	Gecyri	Decyri	Iecyril	Zhecyr	Zecyri	Iicyril	Iishort	Kacyri	Eleyri	Emcyri	Encyri	Ocyri	Pecyri
А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
Ercyri	Escyri	Tecyri	Ucyri	Efcyri	Khacyr	Tsecyr	Checyr	Shacyr	Shhac	Hardsig	Yericy	Softsig	Erever	Iucyri	Iacyri
Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
acyrill	becyri	vecyri	gecyri	decyri	iecyri	zhecyr	zecyri	iiicyri	iishort	kacyri	elcyri	emcyri	encyri	ocyri	pecyri
а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
ercyri	escyri	tecyri	ucyri	efcyri	khacyr	tsecyr	checyr	shacyr	shhac	hardsig	yericy	softsig	erever	iucyri	iacyri
р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
uni045	llocyril	djecyr	gjecyr	ecyrill	dzecyr	icyrill	vicyril	jecyril	ljecyr	njecyr	tshecy	kjecyr	uni045	lshort	dzhecy
è	ë	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	й	ў	ц
Omega	omega	Yatcyr	yatcyr	Eiotif	eiotif	Yuslitt	yuslitt	Yuslitt	yuslitt	Yusbig	yusbig	Yusbig	yusbig	Ksicyr	ksicyri
Ω	ω	Ѡ	ѡ	Є	є	А	а	Ѧ	ѧ	Ѩ	ѩ	Ѫ	ѫ	Ѭ	ѭ
Psicyr	psicyri	Fitacyr	fitacyri	Izhitsa	izhitsa	Izhitsa	izhitsa	Ukcyri	ukcyri	Omega	omega	Omega	omega	Otcyri	otcyri
Ψ	ψ	Ѳ	ѳ	Ѵ	ѵ	Ѷ	ѷ	Ѹ	ѹ	Ѻ	ѻ	Ѽ	ѽ	Ѿ	ѿ
Koppac	koppac	thousar	titlocy	palatal	dasiapr	psilipne	uni048	uni048	uni048	Iishort	iishort	Semis	semiso	Ercyri	ercyri
Ѕ	ѕ	Ѵ	ѵ	Ѷ	ѷ					Й	й	Б	б	Р	р

Fig. 1: Cyrillic Unicode Block (‘Lucida Grande’ font)

This pertains to the Latin accented characters used by today’s Slavic languages, and to the contemporary orthography of Cyrillic. The differences – and problems – however begin to start immediately beyond that point.

2. Notes on available characters and solutions

2.1. The Uk ligature

Among the historical additions to the Slavic block there is the ‘Uk’ digraph, i.e. Oy / oy. Sometimes this character is rendered as a digraph, and sometimes as a vertical ligature, i.e. 8 8 (see Fig. 1). However, one can either have the digraph or the ligature in the Slavic Cyrillic block, not both.

Just browsing through all characters available in a given font will, however, reveal the presence of 8 8 at another location: They are available as *Latin* letters in the Latin Extended-B block, character codes [0222] and [0223]. This is somewhat unexpected because it is well known to Slavists that these characters were introduced into the Cyrillic alphabet under the influence of the Greek alphabet and its writing conventions. In the Greek section of Unicode, however, these vertical ligatures aren’t available. The reason why the vertical ligature is not present in Unicode as a Greek character is because Greek official bodies vetoed its introduction on the basis that it’s use is ‘a mistake’, and that it will be corrected by teachers. However, the vertical ligature is being used today even in printed form, see Fig. 2 for an example in the name of a famous Greek biscuit company.



Fig. 2: Greek vertical OY ligature in use today

Anyway, from the Greek alphabet the ligature found its way into the alphabets of the Algonquin and Huron languages. Thus, the Greek vertical 8 8 became a *Latin* letter present in two Indian languages but the language it was borrowed from, Greek, does not have that character in Unicode because officials did not like it there, and the need of philologists wasn’t considered relevant enough to overturn the official vote.

Now from a user’s standpoint it wouldn’t make much of a difference if the 8 8 is a Greek, a Cyrillic or a Latin letter – if it’s there, it can be used. That’s indeed true, but it is also true that one will then have a mixture of Greek and Latin letters within one word, or of Cyrillic and Latin. This has negative side effects on sorting, spell checking, hyphenation etc. and such a mixture of characters should be avoided. For the very same reasons, Unicode has a Cyrillic ‘A’ even if already has a Latin ‘A’!

The same observation is also valid for other characters: Cyrillic characters not yet available in Unicode may look identical to Latin characters elsewhere in other Unicode blocks, and thus it is tempting to use them instead. However, as

we wanted to point out here, that is not the ideal approach, although certainly practical if no other solution is available at a given time.

2.2. Long and short vowels

Available in Unicode are precomposed characters for most long and short vowels:

Long vowels: ā ē ī ō ū ŷ
Short vowels: ǎ ě ĭ ǒ ů _

What's missing from this chart is a short y, which, however, can be composed from its parts: ŷ. One may be tempted to simply use the Byelorussian character instead, which looks identical. While for a printed text or an on-screen representation it makes no difference whether the ŷ is Cyrillic or Latin, this distinction is important for settings concerning the language of the text in a word-processor file, for hyphenation, for sorting etc., i.e. for text processing and text encoding – see above.

3. Problems in Unicode 4.1 – continued

3.1. Štokavian Accents – the missing double grave accent

Serbian and Croatian use four accent marks to denote the four possible combinations of long vs. short duration with rising vs. falling tone. The accepted norm uses acute (´) for long-rising, grave (`) for short-rising, an inverted bow (˘) for long-falling and a double grave accent (˝) for short-falling. These four diacritics go over any of the following six vowels: *a e i o u r*, totalling 6 x 4 = 24 different character combinations. Of these 24 characters, 23 are covered as precomposed characters by Unicode 4.1, while one has been forgotten: *r with grave accent*: ù. This omission was treated more fully in the first article. As one can see from the preceding sentences, it is sometimes necessary, when one writes about accents, to typeset them by themselves. This is why Unicode has among its blocks a section devoted to ‘spacing modifiers’, i.e. accents with their own positive width, and a section devoted to ‘combining diacritics’, i.e. accents which have zero-width, i.e. which will not move the insertion point in a word-processor, and which will be displayed instead above or under the preceding character. Such ‘combining diacritics’ are also known as ‘flying accents’.

While most accents required for Slavic phonetics are available as spacing accents as well as non-spacing accents, the one accent missing is the double grave accent. It is available in the ‘Combining Diacritical Marks’ block at [030F], but not in the ‘Spacing Modifiers’ section. There, however, we find a “modifier letter middle double grave accent” at number [02F5] among several other “middle” modifiers. They all differ from normal modifiers in their height:

they do not sit *above* characters but to the side of them: x`x. In other words: a normal *spacing modifier double grave accent* is missing from Unicode at present and should be added as a ‘spacing clone’ of the corresponding combining diacritic.

3.2. Sorbian and Polish orthography

In the first article on this subject, we already mentioned that for *Sorbian* several characters are missing. We will be presenting some additional illustrations here.

I. Lautlehre

§ 1. Buchstaben

Das Niedersorbische wird heute ausschließlich mit lateinischen Buchstaben geschrieben. Da bis 1933, wo der Nazismus jede niedersorbische Kulturarbeit erstickte, neben der erwähnten Schreibweise auch mit deutschen Buchstaben, letzteres sogar häufiger, geschrieben wurde, sei zum Verständnis der älteren Schriften auch diese Schreibweise mit angegeben.

a) Lateinische Schreibweise: a b ḃ c č d dź dż ź e ě f f' g h

b) deutsche Schreibweise: a b ḃ ċ ḋ ḋ ž̇ ė ě̇(ë̇) f' f' g h

ch i j k l l m ṁ n ṅ o p ṗ r ṙ (ř) s š t ś ć u w ẇ y z ž

ϥ i j k l l m ṁ n ṅ o p ṗ r ṙ (ř̇) ſ ẛ t ṥ ć̇ u w ẇ y z ž̇

Die Benennung der Buchstaben entspricht dem Deutschen; dazu merke: ć, dź, ś, ź nennt man ćej, dzej, śej, źej; č, ž, ě = čet, žet, ět, š = eš; ł = eł; ch, dż = cha, dža¹.

Fig. 3: Character table from Šwela (1952, 1).

The characters that are not available in Unicode are precomposed b' and f', which may not really be necessary given that they can be composed from their parts quite satisfactorily, and the *long s* (ſ) *with stroke*. See Fig. 5 for some more samples. The *long s with stroke* is an addition to Unicode that would be needed. Because this character was always and only printed using broken script fonts ('Fraktur'), the lowercase character is not the 'normal' s, but the long s (ſ). Implemented for a modern typeface, it would look like this:



Fig. 4: Serifed design of S and long s with stroke

- z, f weiches s (lesen) za – fa für (**S**ache), zły – fły böse.
 s, ß hartes s, ß (Grü**ß**e, mi**ß**): se – ße sich (Erb**ß**e), gus – guß die Gans (der **G**u**ß**).
 c, ç wie dt. z: cas – zaß die Zeit, cera – zera Linie (**z**erren).
 ź, ż ganz weiches dt. sch mit j verbunden: źiwy – żiwy wild.
 ś, śch dt. sch mit j verbunden: śele – śchela Kalb, śpa – śchpa Stube.
 ć, cź dt. tsch mit j verbunden = tś: ścina – cźcina Rohr, ścianiś – cźcianiś schillern.
 dź, dż = d + ź: zdźarżaś – dźarżasch erhalten.
 ż, żj volles sanftes sch, wie im Worte **j**ournal: żaba – żaba der Frosch, żywy – żywy lebendig.
 ś, śch härter als dt. sch (Pas**ch**a), śeść – śchescź sechs, śnarł – śchnarł der Goldammer (**s**chnarren), bęśćo – bęśchcjo ihr waret.
 ć, cź volles tsch (K**u**t**s**cher) = tś: laźcej – cźcej leichter.
 dź, dż = d + ź: ldzej – dźej leichter.
 ř, řch = sch, steht in manchen älteren Schriften nach k, p, t statt heutigem š: křej – kšej, třawa – tšawa.

Fig. 5: Character samples from Šwela (1952, 3).

In the history of *Polish* orthography, the work of Jakub PARKOSZ (ca. 1440) stands out as the first attempt to distinguish hard and soft characters. In the first part of the article on Slavic philology and Unicode, we already presented a scan from modern printed editions of his works. We will present here the handwritten original:

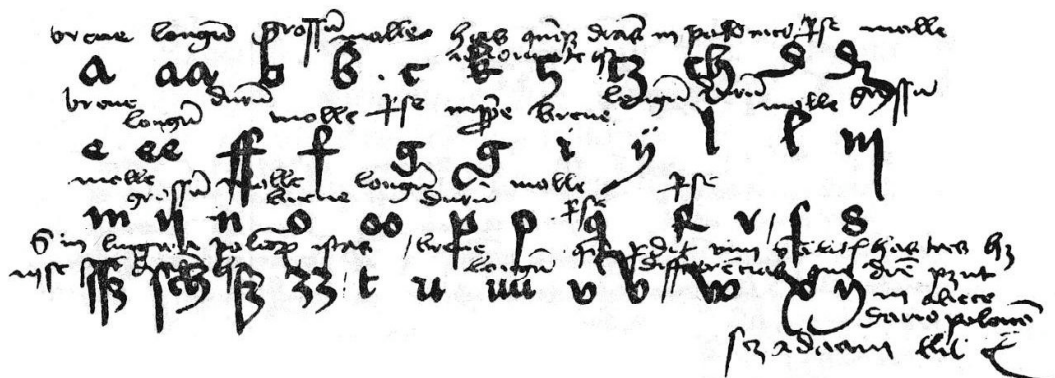


Fig. 6: Character samples from J. Parkosz (URBAŃCZYK/OLESCH 1985, 53)

Here, one can see the *square and round b* in the first line under the Latin descriptions ‘grossum’ (hard) and ‘molle’ (soft), and the *square and round p* in the third line, again with ‘grossum’ and ‘molle’ above them. As was already stated in the first article, Parkosz’s suggestions never became popular even

with his contemporaries because these distinctions were hard to realize and more a question of writing styles or font design than true graphemes.

We will present additional here material concerning the so-called *r rotunda*, the ‘round r’. Just as the Latin *s* is available in Unicode in two forms, the standard form and the *long s*, the *r rotunda* needs to be encoded as the second form of the *r*:



Fig. 7: *r*, *r rotunda*, *s*, and *long s* (‘Breitkopf’ font)

In contrast to German, in Polish the *r rotunda* carried a functional load: it helped to distinguish -rz- [ž] from -r-z- [rz] (see URBAŃCZYK/ OLESCH 1985, 36 on Murzynowski 1551). The *r rotunda* was also in use for Sorbian, see the second line of the table in Fig 8:

18

TABELLARISCHE UEBERSICHT DER NS. ALPHABETE.

Fabricius 1706	Fryco 1796	Zwahr 1847	Tešnar- Šwjela	Časopis M. S.	Mein Alphabet	Ober- sorbisch	Alt- slovenisch
p	p	p	p	p	p	p	p
r	r z	r	r	r	r	r	r
—	sch—fch	—	—	ř	—	ř	—
ff (ff)	ff	ss	ff	s	s	s	s
fch	ffch	sch	fch	š	š	š	š
fch	fch	schj	fch	ś	ś	—	—
t	t	t	t	t	t	t	t

Fig. 8: Alphabet chart from MUCKE (1891; repr. 1965, 18)

3.3. Unified Jers and Nasals

In Cyrillic texts, we find a so-called ‘unified jer’, which is written as a middle form between the two standard jers, and we also find so-called ‘unified nasals’:



Fig. 9: Unified Jer and unified nasals (‘Method’ font)

All these character pairs are not yet available in Unicode. Their very essence lies in the fact the scribes could not decide or did not know or were not aware of which one of each pair to write, and therefore they went for a middle form. This, of course, means that these middle forms cannot be identified with either one as being a stylistic variant. At this point, it should be noted that these ‘unified’ letters are not identical to the characters used in the so-called ‘one-jer-texts’ which exhibit only one of the jers, but where the form of this letter is clearly one of the two.

3.4. Transliteration of Glagolitic Nasals

Now that the Glagolica has been included in Unicode in its version 4.1, it is only logical to check if all characters needed for the *transliteration* of Glagolitic into Cyrillic are available. First problems in this area were noted in the first article – the case of the missing *Cyrillic Iota*.

The transliteration of the nasal vowels from the Glagolitic alphabet into Cyrillic letters presents another currently unsolved problem. Let’s have a look at the transliteration:

Ѣ	↔	А	✓
Ѧ	↔	А	¬
Ѣ	↔	ІА	✓
Ѣ	↔	Ѧ	✓
Ѣ	↔	ІѦ	✓
Ѣ	↔	‘YO’	¬

The Ѣ – nasal E (first line) – can be transliterated into Cyrillic. In the second line, there is special character that occurs in Zographensis and in Marianus: Ѧ; it is encoded in Unicode as a separate glyph even if it is recognized to be a variant of the glyph in the first line.

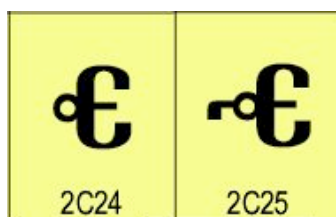


Fig. 10: Variants of EN in Unicode

According to TRUBETZKOY, the glyph € denotes a nasal O; however, because its shape is derived from the nasal E ($\text{€} < \text{€}$), it is usually transliterated using a modified Cyrillic nasal E ($\text{€} < \text{€}$) whose form actually occurs in Suprasliensis (but is not encoded as a separate letter in Unicode). Consequently, the € should be added to the Cyrillic block not only in its own right, but primarily because it is needed to transliterate € .

The jotated nasal E (€), the nasal O (€) and the jotated nasal O (€) can all be transliterated, as lines three to five show. However, sometimes the need arises to single out the first part of the jotated nasal O, i.e. € , for example, to write about it, and it also seems to occur in a non-connected form in actual manuscripts. Here, the situation becomes even less clear: it is unclear how to transliterate this part into Cyrillic; sometimes we simply find the Latin digraph ‘YO’ although it clearly does not represent the phonetic value correctly.

TRUBETZKOY’s OCS Grammar (1968, 22) presents two additional interesting transliteration characters:

22

N. S. Trubetzkoy

I. Gruppe (Einer)	II. Gruppe (Zehner)	III. Gruppe (Hunderter)	IV. Gruppe (Tausender)
1 € a	10 € € i_1	100 € r	1000 € €
2 € b	20 € i_2	200 € s	2000 € €
3 € v	30 € €	300 € t	3000 € €
4 € g	40 € k	400 € €	4000 € €
5 € d	50 € l	500 € f	5000 € € (€)
6 € e	60 € m	600 € €	6000 € €
7 € €	70 € n	700 € €	7000 € €
8 € €	80 € €	800 € €	8000 € €
9 € z	90 € p	900 € c	9000 € €

Fig. 11: TRUBETZKOY’s transliteration of Glagolitic

The transliteration for numbers “30” and “800” are unusual additions, derived from the shapes of the Glagolitic letters, it seems, and mixed into the set of Latin letters used for the rest of table. This seems to be a singular use of these nonstandard characters, so they may not merit inclusion into Unicode.

This, however, brings us to the next interesting area:

3.5. Transliteration of Glagolitic into Croatian (Latin)

While Glagolitic has been and usually is being transliterated into Cyrillic, we should not forget that in Croatia this is not the case: there, the Croatian or Square Glagoljica is transliterated using Latin characters.³ If we take a look at Fig. 12, we see one character that is not present in Unicode:

Ɱ	3
ⱮⱮ	ï
ⱮⱮ	za starije spomenike ġ, za mlade ŷ
Ɱ	ō
ⱮⱮ	ê
ⱮⱮ	č
Ɱ	ju
ⱮⱮ	ŷ
Ɱ	6
Ɱ	ě, za vrijednost »ja« — ê

Fig. 12: Croatian transliteration of Glagolitic (BRATULIĆ 1995, 98)

A not-so-common character is the c with circumflex (ĉ) which is available in Unicode at [0108] and [0109]; it can, of course, also be composed from its parts. The j with diacritic needs to be looked at more closely. While Fig. 12 seems to show a *j with circumflex* (which is available in Unicode at [0134] and [0135]), another figure from the same source clearly shows shows an inverted breve:

Ć, Ĵ

Fig. 13: Croatian transliteration of Glagolitic (BRATULIĆ 1995, 146)

3 Bulgarian or Round Glagolica and Croatian or Square Glagoljica are considered to be the same script in Unicode, and the difference between them to be stylistical.

A *j* with an inverted breve is not yet available in Unicode. Also, such a character cannot be made up from its parts because a ‘dotless *j*’ is also not yet available so that would be another character that needs consideration.

3.6 Croatian Glagoljica: variants or separate letters?

The encoding of the Glagoljica presents another interesting case: in the first drafts of the original submission, the two variants representing the back jer (first line) were both present individually, just as the two variants representing the front jer (second line).

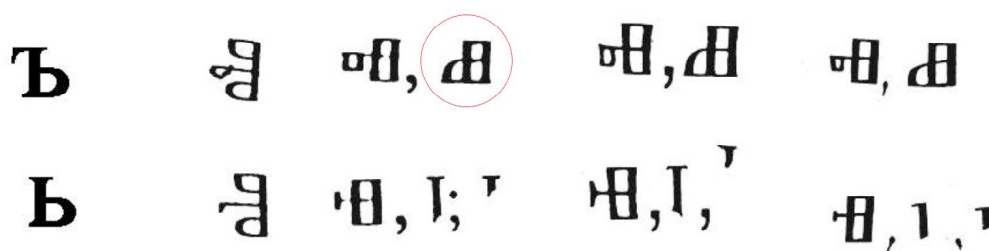


Fig. 14: Glagoljica variants of the Jers

In the final proposal, however, only the variants for the front jer survived, while the second variant for the back jer (encircled in Fig. 14) has not been encoded separately. While the two characters are surely related, one could still have preferred to have them both available.

3.7 Bosančica

This brings us to the one Slavic script that seems to have not yet been considered as such at all: The Bosančica.

đ	A	o	O
б	B	п	P
п	V	р	R
г	G	с	S
а, ђ, д	D	т	T
е	E	Ѹ	U
ѡ, ж, ѣ	Ž	ф	F
ѣ	DZ	х	H
з	Z	ω	OT
и	I	ψ	Ć, ŠT, ŠĆ
Ѡ, ѡ	Đ, Ć (đerv)	ч	C
к	K	ѵ	Č
л	L	ш	Š
љ	LJ	б	poluglas
м	M	ѣ, ѡ	JAT
н	N	ю	JU
њ	NJ		

Fig. 15: Bosančica alphabet (from Žubrinić 1996,70)

The question which of these characters are different enough to warrant separate encoding is not easy to answer. Several candidates could be singled out, primarily ‘V’, ‘D’, ‘Ć/Ď’, and ‘N’, with ‘Ć/Ď’ being the favorite.

3.8. Superscripts

Superscripts are another example where a solution has already been implemented for medievalists working with the Latin alphabet, while no comparable solution exists for Cyrillic. Nearly the complete Latin alphabet is currently already as superscripts; Fig. 16 shows a sub-set only:

0363	0364	0365	0366	0367	0368	0369	036A	036B	036C	036D	036E	036F
a	e	i	o	u	c	d	h	m	r	t	v	x

Fig. 16: Combining Latin superscripts in Unicode

These superscripts allow the writing of $\overset{e}{u}$ (= ü), $\overset{a}{a}$ (= ä), $\overset{o}{o}$ (= ö) etc., which is very important for medieval German.⁴ However, as we said, no superscripts are available for Cyrillic, even if some have very distinct shapes, see Fig. 17 for the ‘d’ and the ‘z’, maybe less so for the ‘x’.

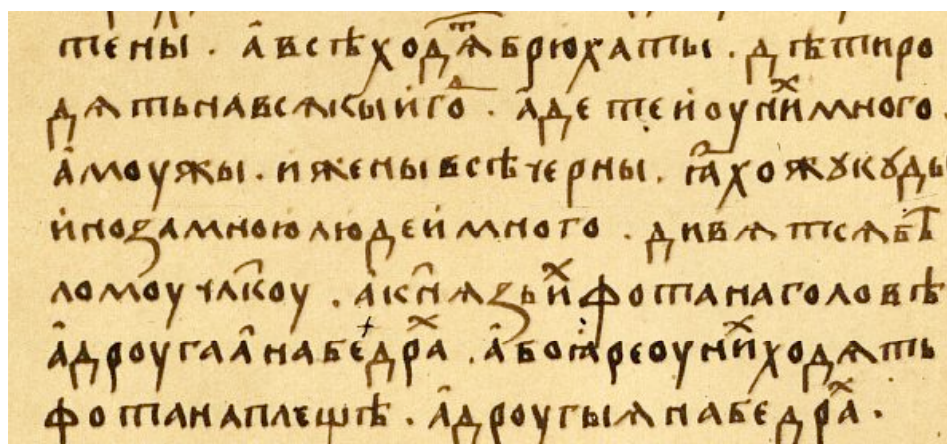


Fig. 17: Cyrillic superscripts in Afanasij Nikitin’s ‘Voyage Beyond the Three Seas’ (*Troickij spisok*)

4 German orthography is remarkable in that it always had clear rules of how to replace non-ASCII-characters by sequences of ASCII characters: ä > ae, ö > oe, ü > ue, ß > ss. These same replacements are also being used today in e-mail addresses (i.e. a person called ‘Müller’ will always choose the sequence ‘mueller’, never ‘muller’, for his/her e-mail address).

These rules do not seem to be well known abroad where the umlaut characters are usually simply replaced by the base character instead (i.e. ä > a, ö > o, ü > u). Very strange-looking is, however, the practice to substitute a Greek beta (β) if no German double s (ß) is available – this can sometimes be observed in Russian printing.

3.9. Ligatures

Ligatures are, by their very nature, a most pleasing subject, aesthetically speaking. The history of Latin typesetting knows much more ligatures than a normal user today will be aware of. Fig. 18 gives an impression of an extended set of ligatures in a modern ‘expert’ digital font.



Fig. 18: Expanded set of Latin ligatures in a modern typeface

In Unicode, ligatures are considered “presentation forms”. Quite a number of such forms are required for Arabic where the shape of a character depends on its position in the word (beginning, middle, end). For reasons of compatibility with legacy code pages, Unicode preserves some ligatures for Latin (see Fig. 20, left part). However, for Cyrillic and Glagolitic, there are no ligatures at all available in Unicode. Both script systems do know quite a number of ligatures. In Cyrillic, there are at least a dozen common ligatures (see Fig. 19), but there are hundreds of actual ligatures to be found in texts.



Fig. 19: Sample OCS Cyrillic ligatures (‘Method’ font)

The same is true for Glagolitic, and Fig. 20 (right half) shows a well-known table listing ligatures consisting of two and even three characters. It may be worth mentioning here that Greek manuscripts also show a large number of ligatures, which are also not present in Unicode.

FB0		XLVIII LITTERAE IUNCTAE FREQUENTER OCCURRENTES													
0	ff FB00	đc at	đu au	đnh ava	đuuu avy	đur ad	đur az	đur azd	đur ak	đur al	đur aa	đur ap			
1	fi FB01	đur va	đur vva	đur vvr	đur vd	đur vdti	đur vz	đur vzd	đur vz	đur vza	đur vzv	đur vzd	đur vzda	đur vzdv	đur vzvl
2	fl FB02	đur vtr	đur vtr	đur gza	đur gv	đur gvi	đur gd	đur gda	đur gz	đur gl	đur gla	đur glv	đur gla	đur gla	đur gla
3	ffi FB03	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr
4	ffl FB04	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr
5	flt FB05	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr
6	st FB06	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr	đur vtr

Fig. 20: Latin ligatures in Unicode, and a table of Glagolitic ligatures

3.10. Numbers

Fig. 21 shows the ‘Numbers’ block in Unicode, consisting of various fractions and then Roman numerals in uppercase and lowercase:

2150	2151	2152	2153	2154	2155	2156	2157	2158	2159	215A	215B	215C	215D	215E	215F
■	■	■	1/3	2/3	1/5	2/5	3/5	4/5	1/6	5/6	1/8	3/8	5/8	7/8	1/
2160	2161	2162	2163	2164	2165	2166	2167	2168	2169	216A	216B	216C	216D	216E	216F
I	II	III	IV	V	VI	VII	VIII	IX	X	XI	XII	L	C	D	M
2170	2171	2172	2173	2174	2175	2176	2177	2178	2179	217A	217B	217C	217D	217E	217F
i	ii	iii	iv	v	vi	vii	viii	ix	x	xi	xii	l	c	d	m

Fig. 21: Fractions and Latin (Roman) numbers in Unicode

In contrast to this, for Slavic no precomposed numbers are available. As a visual reminder of the Slavic writing conventions for numbers might serve Fig. 22. For Slavic, we have the ‘thousand’ symbol encoded, but nothing else.

СЛОВ'ЯНСЬКІ ЦИФРИ

·А·	·Б·	·Г·	·Д·	·Є·	·С·	·З·	·Н·	·Ф·	·І·
1	2	3	4	5	6	7	8	9	10
АІ	БІ	ПІ	ДІ	ЄІ	СІ	ЗІ	НІ	ДІ	
11	12	13	14	15	16	17	18	19	
К	Л	М	Н	З	О	П	Ч		
20	30	40	50	60	70	80	90		
Р	С	Т	У	Ф	Х	Ц	Ш	Щ	
100	200	300	400	500	600	700	800	900	
*А	*Б	*Г	*Д	*Є	*І				
1000	2000	3000	4000	5000	10000				

Fig. 22: Slavic numbers (source: web, from an unnamed 1977 almanac)

That is, for each number, the user is expected to type at least four parts: a half-high dot, then the letter, a titlo over the letter, and again a closing half-high dot. The problem here is that to successfully compose the Slavic numbers from their parts, we need a titlo that goes above *two letters* – see numbers 11 to 19, and such a titlo is not yet available, only one which goes directly above a letter. The best implementation for the longer titlo would probably be a nonspacing diacritic that is to be typed between the two letters, very similar in principle to the ‘combining double macron’ that is available in Unicode at [035E]. So, a ‘combining double titlo’ is what should be added to Unicode.

3.11. Balkan Philology

A broader perspective on Slavic writing systems is in order if we want to fully cover all their uses. First, there is the use of *OCS Cyrillic in Romania*, where it was used until the 19th century. From a manual of Romanian paleography, there is clear evidence of an additional accepted OCS letter (see Fig. 23).

LITERA Ѧ

Litera Ѧ, obținută prin modificarea grafică a literei Ж⁷⁸, în documentele muntenești ale secolului al XV-lea apare o singură dată, într-un act de întărire dat pentru mănăstirea Govora⁷⁹ de voievodul Radu cel Mare: **УТ ГЛАВ ПОИАНЕ ДОЛЕ ПО ѦН⁸⁰ ДОЛННЪ ЧОРЖЧЕК** „din capul poienei, în jos, prin (pe în) valea Cioriceii“ (1499 iul. 13, Arh. St. Buc., S. I., nr. 118).

Deși atestată într-un cuvânt românesc, totuși unicul exemplu în care găsim această literă nu ne dă posibilitatea de a ne referi la valoarea sau valorile acestei litere în general⁸¹. În acest unic caz, grafia ne determină să afirmăm că Ѧ notează nazalitatea vocalei care precedă pe Н⁸².

Fig. 23: OCS letter ‘IN’ in Romanian (DJAMO-DIACONIȚĂ 1971, 49)

This same character is also known from Slavic texts, representing a similar, but not identical sound.

It is well-known that the *Greek* script has been used on the Balkans for neighbouring languages: It was in use in Bulgaria during the First Kingdom before OCS was to become the official language with its own alphabet, and it was used for Macedonian in the 19th century before the language acquired an official status (after WWII). The Greek script was also used for Albanian, and in the sample shown in Fig. 24 we see some interesting character-diacritic combinations not known in Greek itself (epsilon with underline, kappa with over-dot, epsilon with gravis and underline, sigma with dieresis, pi with over-dot – all in the first two lines of the sample).

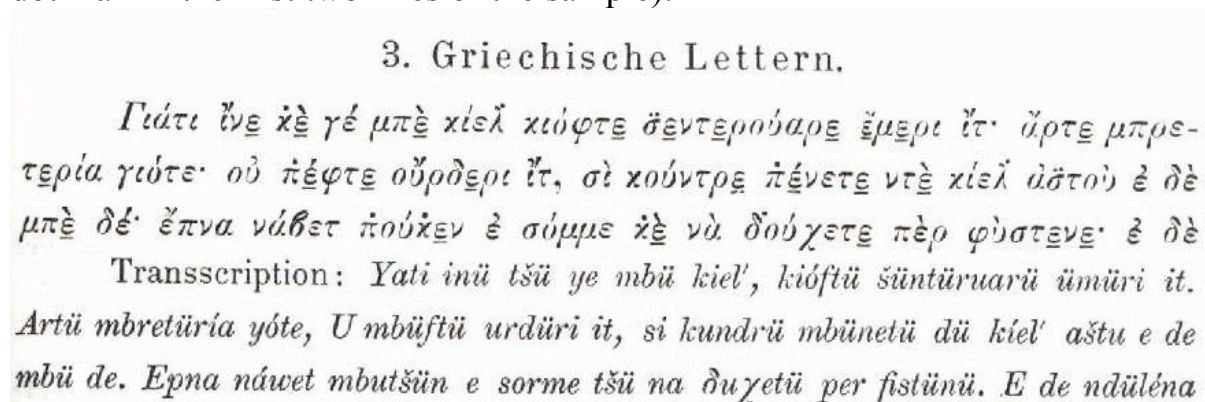


Fig. 24: Albanian written using Greek letters (after FAULMANN 1880)

At present there is no evidence that the Arabic script when applied to Macedonian yielded any additional characters, but its application to Byelorussian and Bosnian would have to be carefully researched to see whether this adaptation resulted in any characters or character plus diacritic combinations that are not available yet in Unicode or cannot be produce from available parts.

3.12. General Phonetics

Concluding our overview of areas that should be considered from a Slavist's point of view, general phonetics should no be completely forgotten. While most phonetic symbols are already available in Unicode, there are some symbols that were either not sanctioned by the IPA but still are in use or have been suggested at some point in time. Fig. 25 shows a list of various such additions not yet implemented in Unicode at present.



Fig. 25: Some phonetic symbols missing from Unicode

Many of these symbols are present in PULLUM/LADUSAW (1986), an excellent source on phonetic symbols, their history and meaning. It might be worth noting that one these symbols is being used by the Bavarian dialect atlas:



Fig. 26: open o-e ligature (PULLUM/LADUSAW 1986, 120)

4. Encoding Strategies and Questionable Characters

4.1 Ukrainian Ghe-upturn

Not every special character, however, that researchers have noted needs inclusion into Unicode. Let us look at such a case here.

TRUNTE (2001, 324–327) remarks that SMOTRYC’KYJ uses a special character in his 1619 grammar and he even typesets it. The character in question can be seen in Fig. 27 in the Greek loanwords *Grammátiki*, *Orfografía*, *Ety-mología*. TRUNTE calls this character an ‘allographe’ or ‘variant’ (2001, 325) of ɣ, and he says that it “stands for” Greek gamma (327).

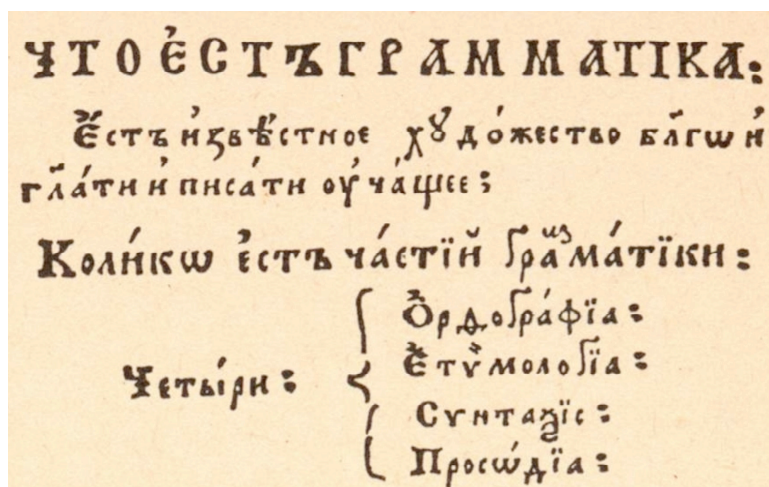


Fig. 27: Use of Ghe-upturn in SMOTRYC’KYJ (1619, a v)

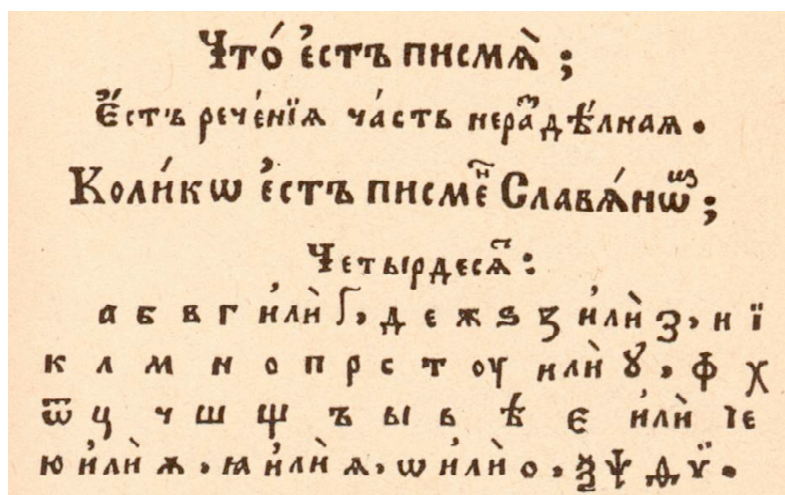


Fig. 28: Alphabet from SMOTRYC'KYJ (1619, a g)

Now what do these pictures really show? We see a “g” that does not fit the rest of Church Slavonic letters very well, and although it is reminiscent of the *capital* Cyrillic letter Г or Г, it is clearly a *lowercase* letter in both figures. Note that SMOTRYC'KYJ uses the ‘normal’ Г in the uppercase-only title question ‘What is Grammar?’ and, in the same word, but printed in lowercase letters, the Ghe-upturn in the fourth line of Fig. 27.

By pure chances, the author happened to have the label of a Greek wine bottle which allowed for some interesting observations: This label has the same letter in it which SMOTRYC'KYJ uses: again, it clearly is a lowercase letter (see the first word in the first line, where it is used twice, and the first word in the second line of the text). Now it becomes clearer that what SMOTRYC'KYJ uses is nothing but a *Greek lowercase Gamma* in the context of OCS letters. This explains the difference in design (it is lighter than the true OCS characters) and the special form.⁵



Fig. 29: Modern Greek wine label with script gamma

5 The same Greek gamma is also present in the table of Greek letters in Ivan Fedorov’s Ostrog ‘Azbuka’ from 1578 (see 1983, 2).

As the wine label shows, this form of the Greek lowercase gamma is still in use today when a font is required that should look a bit more traditional and handwritten.

If we now take a look at a Ukrainian grammar (RUDNYC'KYJ 1943, 2), we again see the same character in his alphabet table where he shows printed and handwritten forms of all characters:

А, а	<i>A a</i>		a	a
Б, б	<i>Б б</i>	<i>б</i>	be	b
В, в	<i>В в</i>	<i>в</i>	we	v, ѱ
Г, г	<i>Г г</i>		ha	h
Ґ, ґ	<i>Ґ ґ</i>		ge	g
Д, д	<i>Д д</i>	<i>д</i>	de	d
Е, е	<i>Е е</i>		e	e

Fig. 30 Ukrainian G and Ghe-upturn (RUDNYC'KYJ 1943, 2)

As we can see from these figures, the Greek Gamma has been adapted into the construction principles underlying the Cyrillic alphabet which features many lowercase characters that are “small versions” of the respective uppercase letter where the Latin alphabet uses special forms for the lowercase letters (which derive from handwriting). Just compare Cyrillic Т-т, М-м to Latin Т-t, М-m etc.⁶ That is, the special tall Greek lowercase Gamma which SMOTRYC'KYJ uses was later reinterpreted as being an uppercase character, and a small version was derived from it to form the lowercase character. This also explains why in hand-written form, the normal Cyrillic G and the Ghe upturn look different, although the printed forms are very similar to each other (see Fig. 30).

So, the remarks by TRUNTE could and probably should be worded a bit differently: In SMOTRIC'KIJ's Grammar we do see a new character, and because it stands for the sound [g], and not [h], it is a *phoneme* and not an *allophone*: if it is already considered to be a *Cyrillic* character, it isn't a tall variant of г, but of ґ! The very fact, however, that the tall Greek version of the Gamma is being

6 See KEMPGEN (1993) for some quantitative measurements of alphabet systems with regard to their internal structure.

used can be considered as a proof that this is still a borrowed Greek character and not yet a Cyrillic character. In contrast to this, one can indeed say, as TRUNTE does, that ѣ ‘stands for’ Greek theta (θ) in certain manuscripts, because both characters are different to each other and clearly belong to different script systems. But in our case we simply observe the origin of the Ukrainian ‘Ghe upturn’ character which began its life as a borrowed tall lowercase Greek script gamma.

If one wants to typeset SMOTRYC’KYJ today, there would be several solutions: one could simply use the modern Cyrillic *ghe upturn* character (i.e. г) or, if more fidelity with the look of the original is required, an alternate tall form of the lowercase *ghe upturn* could be used. Finally, another obvious solution would be to use a *Greek* font with a tall lowercase gamma, just as SMOTRYC’KYJ did for his ‘Grammar’. In any case, this is not a new character that would be missing from Unicode, and it seems a somewhat strange omission that TRUNTE does not identify the special character that he sees in SMOTRYC’KYJ’S ‘Grammar’ with the Ukrainian *ghe upturn* of today.

4.2. Encoding variants

This brings us to a more general consideration of how variants can be encoded within the Unicode framework. Let us take the Jery as an example (see Fig. 31). First, we have an old form (with the ‘hard sign’ as its first part) and a newer form (which has the ‘soft sign’ as its first part. This newer version also has two basic connected or ligature variants (see middle section): the connecting line either is somewhere at half height or at the top the character. And then we also have instances where the second part already has the dotted i as its second part.



Fig. 31: Variants of the ‘jery’

Basically, there would be four ways to encode variants:

- 1) use of the ‘private area’ for everything but the standard form;
- 2) separate fonts – one for the standard form, another one having variants;
- 3) use of font and software technology;
- 4) have Unicode, Inc. add them to the standard.

These solutions have their advantages and disadvantages:

The use of the ‘private are’ has the advantage of having all variants in a single font and the disadvantage of lack of compatibility between fonts and their users’ documents.

The use of different fonts for variants has the advantage that all variants have the same Unicode number as the basic glyph which means all software treats them as the same character with only different ‘looks’. An advantage is also that such a solution is immediately possible. The disadvantage here is that for each new variant of a letter, one needs a new font – possibly a dozen for some characters. As only certain characters do have variants at all, this means that many slots in these additional fonts would remain empty because they are not needed. This clearly is not the most economical or practical way to handle variants either.

Solution 3 can be demonstrated by taking a look at the Mac OS X version of ‘Lucida Grande’, the system font. This font has some Czech alternate characters built-in, but not in the private area, see Fig. 32.

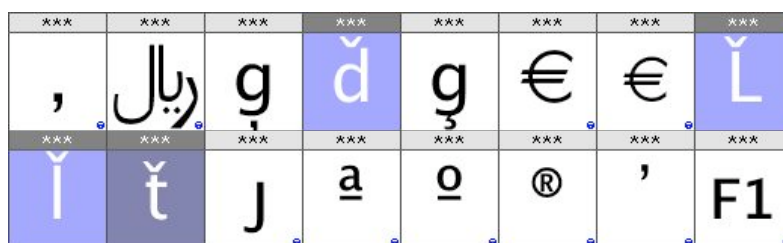


Fig. 32: Built-in variants in ‘Lucida Grande’ font

The clue to the non-use of the private area are the three asterisks above each of the character: characters in the private area have their own number which is not the case here. A user accesses these alternate shapes through functions built into his word processor. ‘Text Edit’, for example, OS X’s standard text editor, has a menu entry ‘Character Shapes’ > ‘Traditional form’. See also Fig. 33 for a possible user interface solution – it shows a part of the ‘character palette’ window where a user can simply point and click on a variant to select it.

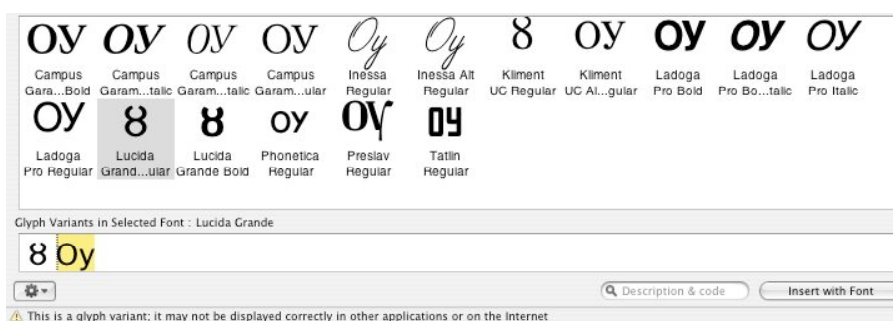


Fig. 33: User interface for selecting variants in OS X

The advantage of such a solution is that all variants are built into the basic font, no redundancy of empty slots in additional fonts is present, but there are also some important disadvantages: in this case the user must rely on software

vendors of word processors to implement a function to select ‘traditional forms’ in their software. Also, the correct display of these alternate shapes cannot be guaranteed with fonts other than the one that has originally been used, and cross-platform compatibility is even less clear. This, therefore, is the solution that requires the most coordinated effort by various software vendors (font vendors, word-processor vendors, and system software vendors) and as such is the least practical solution.

Solution 4 would be ideal if we could rely upon Unicode, Inc. to add all recognized variants to the standard in a timely manner and font vendors to follow up with expanded character sets in their fonts. However, the disadvantage here is that no immediate solution is possible if a solution is needed now and not in several years time. Also, more importantly, the success of such a strategy is unclear, even if some variants are already encoded as separate glyphs in Unicode:



Fig. 34: Sample variants available in Unicode

Another problem is the important question which shapes to treat as different characters and which ones to treat as variants, stylistic or other.

5. Conclusion and Outlook

As both papers have shown, there are quite a few areas in Slavic philology that need careful consideration which characters to submit to Unicode, Inc., for inclusion in a next revision of the standard. These amount to several dozen characters at the very least. If we compare the treatment that Latin has been given with the treatment of Cyrillic and Glagolitic, we clearly see that the Latin alphabet is at present much better supported with respect to the needs of medievalists, paleography etc., although even there many wishes remain unfulfilled at the moment. Much of the support an alphabet receives from Unicode and within Unicode depends on the initiative of individuals working in the field who are at the same time interested in computer technology and therefore willing to submit proposals to the Unicode consortium.

As of this writing (April 2006), Ralph CLEMINSON has submitted a proposal to Unicode, Inc. for inclusion of some 40 characters (most of them pairs). The author of the present article contributed various illustrations and suggestions apart from his oral presentation at the Sofia conference (which formed the basis for both articles). All characters included in the proposal are *Cyrillic* – Old Church Slavonic and its transliteration as well as Old Russian are the main

focus. The proposal does not include *all* Cyrillic characters mentioned in our two articles, but many of them, and *no Latin* characters. Consequently, more work needs to be done, in the area of Cyrillic characters as well as in the area of Latin characters, and even Glagolitic.

Abstract

The present paper is the second in a series presenting an overview of Slavic philology with respect to version 4.1 of the Unicode standard. Reviewing East, West and South Slavonic languages, their alphabets and writing systems, the first paper already revealed at least a dozen characters that need to be encoded in Unicode, among them several (soft) Cyrillic characters with a tail at the right side, the Cyrillic old-style ІѦ, a Cyrillic Iota (uppercase and lowercase), a Cyrillic dotless lowercase і, the Cyrillic Paerok (no distinction between uppercase and lowercase), number signs, accents etc. This second article presents even more characters, among them a Latin S with stroke, the ‘r rotunda’, problems in the transliteration of Glagolitic into either Cyrillic or Latin, a broader perspective on Balkan philology etc.

References

- Bratulić, J.:
1995 *Leksikon hrvatske glagoljice*. Zagreb.
- Cleminson, R.:
2006 *Proposal for additional cyrillic characters*. pdf file available from:
<http://www.unicode.org/~dwanders/06042-cleminson-cyrillic.pdf>
- Djamo-Diaconiță, L.:
1971 *Limba documentelor slavo-române emise în Țara Românească în sec. XIV și XV*. București.
- Faulmann, K.:
1880 *Illustrierte Geschichte der Schrift*. Wien–Pest–Leipzig.
- Kempgen, S.:
1993 Spezifika slawischer Schriften. In: ders. (ed.), *Slavistische Linguistik 1992*, München, 111–143.
- 2006 Slavic Philology and Unicode (I). To appear in: A. Miltenova et al. (eds.), *Proceedings of the ‘Azbuky net’ Conference, Sofia, Oct., Oct 24-27, 2005*. Sofia. 25 pages.
- Mucke, K.E.:
1891 *Historische und vergleichende Laut- und Formenlehre der niedersorbischen (niederlausitzisch-wendischen) Sprache*. Leipzig (Reprint 1965).

- Parkoszowicz, J.:
 1985 *Traktat o ortografii Polskiej Jakuba Parkosza*. Opracował Marian Kucała. Warszawa: PAN.
- Pullum, G.K., Ladusaw, W.A.:
 1986 *Phonetic Symbol Guide*. Chicago–London.
- Šwela, B.:
 1952 *Grammatik der niedersorbischen Sprache*. Leipzig (2nd ed.). First edition Cottbus 1906.
- Trubetzkoy, N.S.:
 1968 *Altkirchenslawische Grammatik. Schrift-, Laut- und Formensystem*. Hrsg. v. R. Jagoditsch. 2. Auflage. Graz–Wien–Köln.
- Trunte, N.:
 2001 *Slavenskij jazyk. Ein praktisches Lehrbuch des Kirchenslavischen in 30 Lektionen. Zugleich eine Einführung in die slavische Philologie. Band 2: Mittel- und Neukirchenslavisch* (Slavistische Beiträge, Bd. 370, Studienhilfen, Bd 9). München.
- The Unicode Consortium
 1991 *The Unicode Standard. Worldwide Character Encoding. Version 1.0*. Vols. 1-2. Reading, Mass.: Addison-Wesley Publ. Co.
- Unicode v. 4.1
 2005 Complete code tables are available in a single pdf file for download:
<http://www.unicode.org/Public/4.1.0/charts/Codecharts.pdf>
- Urbańczyk, S., Olesch, R. (eds.):
 1983 *Die altpolnischen Orthographien des 16. Jahrhunderts* (Slavistische Forschungen, Bd. 37). Köln–Wien: Böhlau.
- Žubrinić, D.:
 1996 *Hrvatska Glagoljica*. Zagreb.
- Федоров, И.:
 1983 *Азбука Ивана Федорова 1578*. Москва.

Appendix: Some fonts and their features for slavists

	Kliment Std v. 1.7 (Win/OS X)	Times / Helvetica v. 5.0d10e1 (OS X)	Times New Roman / Arial v. 3.05 (Win/OS X)
Basic Latin	✓	✓	✓
Latin-1 Supplement (= Western Europe)	✓	✓	✓
Latin Extended-A (= Eastern Europe & more)	✓	✓	✓
Latin Extended-B	some	most	some
Croatian Digraphs	✓	✓	---
Maced. Translit. (ǵ)	✓	✓	---
Štokavian Accents	✓	✓	---
Nasal o	✓	✓	---
Latin Extended Additional (256)	some	✓	ca. 1/3
Maced. Translit. (k)	✓	✓	---
Russ. Hist. Translit. (f̆, y̆)	✓	✓	---
Sorbian (m̆, p̆)	✓	✓	---
IPA – Phonetic	some	2/96	1/96
Spacing Modifiers	✓ 80/80	11/80	9/80
Translit. of Jers	✓	---	---
Combining Diacritics (= „Flying Accents“)	✓ 112/112	40/112	5/112
Greek			
Modern Greek	✓	✓	✓
Archaic Letters (Koppa, Stigma, Sampi...)	✓	---	---
Classical Greek	---	✓	---
Cyrillic			
Std. Russian & Slavic	✓	✓	✓
Macedonian Add. (è, ù)	✓	✓	---
Hist. Add. (Ѣ Ѧ ж ...)	✓	---	---
Ukrainian Ghe (Ґ ґ)	✓	✓	✓
Non-Slavic Cyrillic (ex GUS-Countries)	---	ca 1/2	ca. 1/10
Glagolitic	---	---	---
Transliteration into Cyrillic	✓	---	---

	Kliment Std v. 1.7 (Win/OS X)	Times/ Helvetica v. 5.0d10e1 (OS X)	Times NR/ Arial v. 3.05 (Win/OS X)
Armenian, Georgian, Hebrew, Arabic, Ethiopian	--	-- (supported by other fonts)	Hebrew, Arabic
General Punctuation	70/112	18/112	27/112
Superscripts/Subscripts (0...9)	30/46	--	--
Currency (Euro...)	1/48	3/48	6/48
Comb. Diacr. for Symbols (O)	✓	--	--
Number Forms	49/64		
Add. Fractions (2/3...)	✓	--	6/13
Roman Numerals	✓	✓	--
Arrows	6/112	-- (complete in Apple Symbols)	7/112 (complete in Wingdings)
Mathematical Operators (∏, ∫, ≠ ...)	45/256	12/256 (complete in Apple Symbols font)	15/256 (complete in other fonts)



Bibliographische Angaben / Bibliographical Entry:

Sebastian Kempgen: Unicode 4.1 and Slavic Philology – Problems and Perspectives (I). In: A. Miltenova, D. Radoslavova, E. Pancheva (eds.), *Computer Applications in Slavic Studies. Proceedings of Azbuky.net. International Conference and Workshop. 24–27 October 2005*, Sofia, Bulgaria. Sofia 2006, 131–159.

The original presentation given at the conference is available separately.



Bibliographische Angaben / Bibliographical Entry:

Sebastian Kempgen: Unicode 4.1 and Slavic Philology – Problems and Perspectives (II). In: T. Berger, J. Raecke, T. Reuther (Hgg.), *Slavistische Linguistik 2004/2005*, München 2006, 223–248.

The original presentation given at the conference is available separately.

Copyright und Lizenz / Copyright and License:

© Prof. Dr. Sebastian Kempgen 2021; ORCID: 0000-002-2534-9423
Bamberg University, Germany, Slavic Linguistics
<https://www.uni-bamberg.de/slaving/personal/prof-em-dr-sebastian-kempgen/>
<mailto:sebastian.kempgen@uni-bamberg.de>

License: by-nc-nd



January 2021, postprint