

Old Cyrillic in Unicode*

Ivan A Derzhanski

Institute for Mathematics and Computer Science, Bulgarian Academy of Sciences
iad@math.bas.bg

The current version of the Unicode Standard acknowledges the existence of a pre-modern version of the Cyrillic script, but its support thereof is limited to assigning code points to several obsolete letters. Meanwhile mediæval Cyrillic manuscripts and some early printed books feature a plethora of letter shapes, ligatures, diacritic and punctuation marks that want proper representation. (In addition, contemporary editions of mediæval texts employ a variety of annotation signs.) As generally with scripts that predate printing, an obvious problem is the abundance of functional, chronological, regional and decorative variant shapes, the precise details of whose distribution are often unknown. The present contents of the block will need to be interpreted with Old Cyrillic in mind, and decisions to be made as to which remaining characters should be implemented via Unicode’s mechanism of variation selection, as ligatures in the typeface, or as code points in the Private space or the standard Cyrillic block. I discuss the initial stage of this work.

The Unicode Standard (Unicode 4.0.1) makes a controversial statement:

The historical form of the Cyrillic alphabet is treated as a font style variation of modern Cyrillic because the historical forms are relatively close to the modern appearance, and because some of them are still in modern use in languages other than Russian (for example, U+0406 “I” CYRILLIC CAPITAL LETTER I is used in modern Ukrainian and Byelorussian). Some of the letters in this range were used in modern typefaces in Russian and Bulgarian. Prior to 1917, Russian made use of *yat*, *fita*, and *izhitsa*; prior to 1945, Bulgaria made use of these three as well as the *big yus*.

(The last statement is incorrect, incidentally: and went out of use in Bulgarian a century earlier.)

In fact, the historical forms that are in modern use (or have been used in modern typefaces) are negligibly few, compared to the ones that aren’t (and haven’t); likewise, there is a plethora of modern forms with no historical prototypes (the ones invented expressly for non-Slavic languages). None the less, there are other benefits from treating the two forms of the script as style variations. There is the fact that oldstyle typefaces are often used for setting text in Modern Bulgarian, Russian and Serbian in an ‘Olde Cyrillick’ style, which may be an argument for designers of Old Cyrillic typefaces for cutting oldstyle forms of letters that are used now in these languages, though not in mediæval manuscripts and even in books such as the *Fish Primer* (Beron 1824), a book printed in an oldstyle script but with modern-style capitalisation. On the other hand, Gerov’s *Dictionary*, published at the turn of the 20th century (Gerov 1895–1908), is set in a modern typeface, but with two ‘historical’ letters that haven’t made it into Unicode’s Cyrillic plane and one incorrectly analysed there.

This implies a complex relationship between the two forms of the script. For example, mediæval manuscripts feature several shapes of the little jus (А А А А €), as well as the letter Ѧ (with its variant ѧ). In the *Fish Primer* the variety is reduced to one form of each letter. Gerov’s *Dictionary* gives the letters the shapes and I I (with connecting bars), respectively. In modern usage only the former letter survives; it is not readily recognised as a style variation of , though, so oldstyle Cyrillic typefaces can usefully contain both.

Several kinds of characters need to be accounted for:

* Work done in Prague in February–April 2005; especially discussions with Zdenka Ribarová. Most information on the historical usage of letters and variant forms is from Karskij 1928 and Andrejčin 1977; the inventory of letters is partly based on the appendix to Camuglia & Ribarov (ms).

- letters (majuscule and minuscule), including variant letter shapes;
- superscript letters;
- ligatures;
- diacritics;
- punctuation;
- annotation signs.

Of these, only the first category is relevant to both historical and modern Cyrillic; contrariwise, the last one (or perhaps the last three) can be presumed to be very similar to the needs of mediævalists working with other scripts (especially Greek and Roman).

It is a large part of the problem that variation comes in many kinds. It may be phonetic (if the letters stood for different sounds), functional (if they stood for the same sound but were associated with different positions in the word, or with individual words), chronological, regional or merely decorative. Many pairs of related shapes are obtained by reversal (ѡ Ѡ), by the use of diacritics (ѣ ѣ̇, ѡ ѡ̇, ѣ ѣ̇), by variation of width (ѣ ѣ̇), by enlisting Glagolitic letters (ѣ ѣ̇, ѣ ѣ̇); there is also much unclassifiable variation (ѣ ѣ̇). It stands to reason that similar cases should be treated in similar ways if possible.

For each shape, the following options are available:

- Find a matching shape to unify it with
 - in the Cyrillic plane,
 - in another public plane,
 - in the part of the private area used by the Medieval Unicode Font Initiative (MUFI).
- Relay it to Unicode's mechanisms for handling ligatures and diacritics, if it can be viewed as a compound character.
- Place it in a vacant part of the private area.
- Negotiate its inclusion into the primary Cyrillic space.

It is a general policy of Unicode to unify diacritics and punctuation across scripts but to split letters, so a Cyrillic letter shouldn't be unified with a Greek or Roman one (either from the public planes or from MUFI's part of the private area), but diacritics, punctuation and annotation signs may be.

For each pair of historical shapes which kind of are the same, but kind of aren't (e.g., ѣ and ѣ̇, both of which correspond – in function as well as history – to modern ѣ, U+0437), these options exist:

- Call them the same thing, and postulate that a font will contain one or the other in position U+0437, but never both. (No two shapes are ever identical in handwriting; the line needs to be drawn somewhere.) For this, it is a necessary condition that the shapes are not found in the same text (meaning that this is the obvious strategy for the treatment of chronological or regional variants); it may not be a sufficient condition, however (Unicode does separate some variant Greek letters, such as U+03B8 and U+03D1).
- Call them different glyphs for the same character, and postulate that if a font contains both (which it needn't do), then one will be generated by the character code alone (U+0437) and the other by the same code with a variation selector, a little-used part of the Unicode toolbox (U+VS_n). This is advisable if the two are in free variation in documents of the same tradition, but it is desirable to preserve the difference.
- Call them different characters, unify the first (which won't always be easy to identify) with the character code in the public plane (ѣ=U+0437), and

- find another code in the same public plane to unify the second with (the obvious choice for ζ would be U+04E1 , a letter used in Abkhaz for the affricate /ʒ/, in view of the shape), or
- use a code from a private area, and perhaps aim for inclusion into the primary space.

A tentative Old Cyrillic-minded interpretation of the Cyrillic Unicode plane

§1. Modern Cyrillic letters that never had Old Cyrillic forms, nor need any:

048A–048F	
0492–04CE	
04D2–04DF	
04E2–04F5	

The letters \ddot{a} and \ddot{u} (used respectively for \ddot{o} and \ddot{u} in Altaic languages), though they look like mere graphic modifications of a and u , in fact go back to Ӑ and ӑ , which were abolished in 1918 (Musaev 1982:22–23), but the former pair is separated by the sounds, and the latter by the shapes.

§2. Modern Cyrillic letters that weren't used in pre-Modern times, but can use oldstyle forms for the purpose of writing text in Slavic languages in an 'Olde Cyrillick' style:

0400	0450		
0401	0451		
0403	0453		
0408	0458		

0409	0459		
040A	045A		
040C	045C		
040D	045D		

042F	044F		
------	------	--	--

A Ӑ -like letter is found, in fact, in Russian cursive, but as a variant of i .

The letter ӑ superseded Ӑ (as Gerov 1895–1908 explicitly states), though graphically it is a descendant of both Ӑ and ӑ , which fell together in speedy (Russian) writing. Their separation (an artefact of Unicode) is perhaps advantageous, as oldstyle Ӑ letters have been cut, for the purpose of setting 'Olde Cyrillick' text in contemporary Bulgarian and Russian; one might also note that Dal' 1880–1882 lists ' ӑ ' as a separate entry in his *Dictionary*, as a 'nasalised [vowel] still heard in Polish and of interest only to the linguist'.

§3. Letters with regional or chronological variants, of which presumably one form should be chosen in a typeface. Codes of the form Kyr# and Kys# are assigned to letters not currently supported by Unicode:

0410	0430			Ӑ	or 'alpha' (similar to MUF1's U+F200 LATIN SMALL LETTER INSULAR A), a chronological semi-uncial variant.
0411	0431			ӑ	or 'tilted B' (90° clockwise), a Bosnian variant.
0412	0432			Ӓ	or 'loopy B' (a variant used in some texts), or one tilted (90° anticlockwise), frequent in 17c texts.
Kyr6	Kyr7			ӓ	or ditto with the crossbar going through the entire width of the letter.

041C	043C			М	or ditto with a deeper middle loop.
0420	0440			р/р	
0422	0442			т/т	Other variants include a 7-like shape.
Kys0	Kys1			ѣ/ѣ	Non-initial letter (as opposed to the digraph оѣ). Has a variety of shapes (ѣ ѣ; ѣ ѣ). The Latin Extended-B letters U+0222 LATIN CAPITAL LETTER OU and U+0223 LATIN SMALL LETTER OU are available in Code2000, but not in Arial Unicode MS, and at any rate they pertain to a different script.
0424	0444			Ѡ/Ѡ	
0426	0446			ѡ	or ditto with a shorter and more distant tail, or reversed.
0427	0447			Ѣ/Ѣ	
0460	0461			Ѡ/Ѡ	
0472	0473			Ѡ/Ѡ	

§4. Letters which tend not to appear in oldstyle Cyrillic typefaces:

040F	045F				Old Cyrillic form used in Beron 1824.
0490	0491				‘hard Г’, in South Russian/Galician/Ukrainian since 17c.
Kys4	Kys5				The Cyrillic letter Yn, used for ѣ, cooccurs in Roumanian with ѣ, used for ѣ, and so must be treated as a separate letter, not a local variant.

§5. More difficult cases. In the tables 0 indicates the primary variant (for which there may be more than one possible shape, the choice then being left to the designer of the typeface), and 1, 2 and 3 are optional variants, to be obtained with variation selectors:

				0	1	2	
0402	0452			Ѡ?			These two modern Serbian letters are descendants of the same Old Cyrillic letter, which is closer to Ѡ in shape, but was originally closer to Ѡ in sound, being voiced; its use for both voiced and voiceless sounds dates from the 14c.
040B	045B			Ѡ?			
0404	0454			Є			A broader shape was mostly used word-initially and a narrower one word-medially/finally. Allowing for a matching difference in sound (indeed, Є replaced Ѡ by the 14c), the unification of these two shapes with Ѡ (Ukrainian /je/) and Ѡ respectively seems appropriate, although the modern uppercase shapes differ in roundness, and the lowercase ones in openness, rather than width.
0415	0435			є	є		
042D	044D			ѡ		ѡ	A mediæval Bulgarian shape, originally interchangeable with Є or є (of which it was a mere graphic variation), but a separate letter in Modern Russian and most alphabets derived from it.

041E	043E		o	o	A broader shape was mostly used word-initially and a narrower one word-medially/finally.
047A	047B		o		
KyrA	KyrB		o		Used mostly, though not exclusively, instead of o in forms and derivatives of the word oko.
KyrC	KyrD		o		Used word-initially/postvocally from the 14c onwards: oko. (Perhaps a graphic variant of o?)
KyrE	KyrF		oo	oo	Used in дѡоѡ, ѡѡоѡ etc.; also in ѡѡѡѡ. (Even this is not the whole story: oo appears in one place in the word ѡѡѡѡ, and there is a ‘multiocular’ shape – appropriately in a word with this very meaning – made of 10 os in a 15c psalter, but that is by no means a regular part of the script.)
0460	0461		ω	ö	Both ω and ε were used word-medially instead of o and e for disambiguating between certain pairs of homonymous word forms, which is to say that both ω and o were to o as ε was to e, though in different ways (an argument in support of considering o a form of ω, as the Unicode description ROUND OMEGA suggests). The diæresis was decorative.
047C	047D		̄ω		The same thing as ω with titlo (and in fact no different from any other letter with this diacritic), as currently defined.
047E	047F		̅ω		The same thing as ω with superscript .
0421	0441		c	с	The wide variant is found very seldom, always word-initially. Thus the distribution of the variants of c parallels that of the variants of and , obviously because of the similarity of the forms of the letters, a phonetic difference being out of the question here.
0405	0455		ѕ		ѕ is sometimes an episemon (6), but ѕ is much more common in this function.
Kyr4	Kyr5		ѕ		
Kyr6	Kyr7		ѕ		Has the same sound value (/z/) as ѕ/ѕ, but isn’t graphically related to them, and is never used as an episemon.
04E0	04E1		ѕ?		The modern letter ѕ is used in Abkhaz for /z/, and has been proposed for /ð/ in some Samoyedic (Tereščenko 1986) and Tungusic languages. If it is considered too exotic as a counterpart of ѕ, the latter should be treated as a graphic variant of ѕ.
0417	0437		ѕ	ѕ?	
0406	0456		i	i?	Mostly an episemon (10). As a letter, used primarily as a prevocalic or line-final form; also in the name ic ‘Jesus’. The variant with one dot is very rare historically. In Beron 1824 the capital version has no dots, whereas the lowercase one has two unless it bears an accent mark.

0407	0457			ï	i?	Rare form, used word-initially or after и. Has a variant with a double acute. (Modern version used in Ukrainian for /ji/.)
Kyr8	Kyr9			ı		Only used in transcriptions of Glagolitic text.
0418	0438			н/н	н?	It could be postulated that if both shapes exist in a typeface, then +VS1 must generate a letter with a horizontal bar (the original one) and +VS2 one with a tilted bar (the modern one, used in Beron 1824), and what happens without a variant selector is up to the typeface.
041D	043D			н/н	н?	Similarly, if both shapes exist in a typeface, +VS1 must generate a letter with a tilted bar (the original one, used in Beron 1824) and +VS2 one with a horizontal bar (the modern one), and what happens without a variant selector is up to the typeface.
0423	0443			ѣ		The historical letter ѣ was a graphic variant of в; the modern letter was introduced in the place of ѣ.
0474	0475			ѣ		
0476	0477			ѣ		The double grave was decorative, as it appears.
Kys2	Kys3			ѣ		A variant Bulgarian Glagolitic letter (not found in Croatian Glagolitic), only used in one word: ѣЛЗМИ ‘βουνοι’. It is unknown whether this reflects an actual (i.e., phonetic) difference.
0466	0467			А/А		Word-initial/postvocalic letter in some traditions (as opposed to А). Postconsonantal in others (as opposed to А or Ѡ). Postconsonantal letter in Church Slavonic and in Beron 1824 (as opposed to Ѡ).
KysC	KysD			А/А		Word-initial/postvocalic letter in some traditions (as opposed to А). Postconsonantal in others (as opposed to А or Ѡ).
KysA	KysB			ѣ		A Glagolitic letter, only used in Cyrillic in one word: АѣΓ̄Λ̄З ‘ἀγγελος’. (The word is also spelt АΓΓ̄ΕΛ̄З, АН̄ѢΛ̄З and АН̄ЬІЄΛ̄З.)
046A	046B			Ѡ/Ѡ		The latter seems to be a palæographical variant.

§6. Letters with breves. They can be obtained by using U+0306 COMBINING BREVE; three of them, however, are available as precomposed characters:

04D0	04D1			А	Used in Beron 1824. (The modern version is used in Chuvash.)
0419	0439			н	Old Cyrillic form used in Beron 1824.
				ı	Found in the Sinai Psalter.
040E	045E			ѣ?	
				Ѡ	

References

Samuglia, M. & K. Ribarov, ms. 'Old Church Slavonic in Codes'.

Андрейчин, Л. Д. 1977. «Черковнославянска и гражданска азбука през възраждането». / *Из историята на нашето езиково строителство*, София: «Народна просвета», 151–156.

Берон, П. (Петъръ Х. Беровичъ) 1824. Бѣкварь съ различни поꙋченїа 'A Primer with Diverse Instructions' (popularly known as Рибен бѣквар 'Fish Primer').

Геров, Н. 1895–1908.

/ / - /
/ / ('Dictionary of the Bulgarian Language, with interpretation of the words in Bulgarian and Russian').

Даль, В. И. 1880–1882. *Толковый словарь живаго великорускаго языка* ('Explanatory Dictionary of the Living Great Russian Language').

Карский, Е. Ф. 1928. *Славянская кирилловская палеография*. Ленинград.

Мусаев, К. М. (отв. ред.) 1982. *Опыт совершенствования алфавитов и орфографий языков народов СССР*. Москва: «Наука».

Терещенко, Н. М. 1986. «Алфавит нганасанского языка; Алфавит энецкого языка». / П. Я. Скорик (отв. ред.), *Палеоазиатские языки*, Ленинград: «Наука», 45–47, 50–52.