

Proposal to Encode Combining Glagolitic Letters in Unicode

Aleksandr Andreev*, Heinz Miklas, Yuri Shardt

Section 1. Introduction

Glagolitic, also known as “Glagolitsa”, is an alphabetic writing system used to record Church Slavonic and other Slavic languages. Originating in the 9th century, it is the earliest known Slavonic alphabet. The creation of the alphabet is attributed to the younger of the teachers of the Slavs, St. Cyril.

Glagolitic writing may be found in mediæval manuscripts and in printed liturgical books, mostly of a Croatian origin. In Bulgaria, Glagolitic was gradually replaced by the Cyrillic alphabet, and this Cyrillic alphabet was subsequently used also by other Slavs. For its part, the Glagolitic script has been preserved by some communities in Croatia even up to the present. Extant Glagolitic texts are of enormous value to linguists, palæographers, and scholars of liturgy.

Support for Glagolitic in the Unicode standard is required for two purposes. First, contemporary specialists need to be able to typographically represent mediæval texts written in the Glagolitic script, both in printed matter (such as academic publications) and in an electronic format (for use with computer analysis, such as string comparison, wordlist generation and searching). To this end, computer fonts that contain the repertoire of Glagolitic characters must be created. Second, owing to the close relationship between the Cyrillic and Glagolitic writing systems, scholars have traditionally represented Glagolitic texts also in Cyrillic transcription. To facilitate the transliteration process, an encoding model that parallels the model for the Cyrillic script needs to be available for Glagolitic.

The base repertoire of Glagolitic characters has been included in the Unicode standard since version 4.0. Nonetheless, this repertoire is incomplete because it lacks combining Glagolitic letters. Such combining letters exist in the Glagolitic script and play a function that is analogous to their role in Cyrillic – that is, they are used in abbreviations that are either space saving devices (for example, commonly written words are often abbreviated) or in *nomina sacra*. For full support of the Glagolitic writing system in Unicode, as well as for proper interoperability between the implementations of the Glagolitic and Cyrillic scripts, we propose the encoding of these combining characters in an additional block entitled *Glagolitic Extended*.

Section 2. Proposed Characters

The following table contains examples of combining Glagolitic letters that occur in various Glagolitic manuscripts and in printed literature. We propose to encode the characters as one block, in the same codepoint order as the base Glagolitic letters encoded at U+2C00 and following. This allows for simple computer manipulation of Glagolitic characters, as well as leaving some encoding positions empty to be used in the unlikely instance that additional combining characters are discovered by researchers and need to be encoded. Note that since a Glagolitic Extension is not in the Roadmaps to Unicode, all of the indicated codepoints in this proposal are provisional codepoints in the Private Use Area (PUA).

* Corresponding author, aleksandr.andreev@gmail.com.

Name	Codepoint	Appearance	Location in Sources
Combining Glagolitic Letter Azu	U+E000	ⱦ ⱦ◌	<i>Srez.</i> , p. 42; <i>MissSin</i> , 13v18, 44r16-17, 45v20; <i>PsDem</i> , 27v6, 42r7, 118r5
Combining Glagolitic Letter Buki	U+E001	Ⱨ Ⱨ◌	<i>Srez.</i> , p. 254; <i>MissSin</i> , 38r13, 49v20
Combining Glagolitic Letter Vede	U+E002	ⱨ ⱨ◌	<i>Srez.</i> , p. 59; <i>MissSin</i> , 7r2, 22r10, 40v10, 52r17; <i>PsDem</i> , 104r8
Combining Glagolitic Letter Glagoli	U+E003	Ⱪ Ⱪ◌	<i>Srez.</i> , p. 254; <i>PsDem</i> , 117r1
Combining Glagolitic Letter Dobro	U+E004	ⱪ ⱪ◌	<i>Srez.</i> , p. 228; <i>MissSin</i> , 54v14; <i>PsDem</i> , 21v1, 42r8, 78r19, 131r9
Combining Glagolitic Letter Yestu	U+E005	Ⱬ Ⱬ◌	<i>Srez.</i> , p. 59; <i>MissSin</i> , 29r20, 46r12
Combining Glagolitic Letter Zhivete	U+E006	ⱬ ⱬ◌	<i>Srez.</i> , p. 84
Combining Glagolitic Letter Zemlja	U+E008	Ɱ Ɱ◌	<i>EuchSinV</i> , 103r16m
Combining Glagolitic Letter Izhe	U+E009	Ɐ Ɐ◌	<i>Srez.</i> , p. 82
Combining Glagolitic Letter Initial Izhe	U+E00A	Ɒ Ɒ◌	<i>PsSinV</i> , 177r18
Combining Glagolitic Letter I	U+E00B	ⱱ ⱱ◌	<i>Srezn.</i> , p. 82
Combining Glagolitic Letter Djervi	U+E00C	Ⱳ Ⱳ◌	<i>Srezn.</i> , p. 84; <i>MissSin</i> , 18v6
Combining Glagolitic Letter Kako	U+E00D	ⱳ ⱳ◌	<i>Srezn.</i> , p. 224; <i>PsDem</i> , 126r2
Combining Glagolitic Letter Ljudie	U+E00E	Ⱶ Ⱶ◌	<i>Srezn.</i> , p. 40, p. 42; <i>PsDem</i> , 5v1; <i>MissSin.</i> , 30v22, 51r15
Combining Glagolitic Letter Myslite	U+E00F	ⱶ ⱶ◌	<i>Srezn.</i> , p. 224; <i>PsDem</i> , 50v15, 105r19, 113v3; <i>MissSin</i> , 13r18
Combining Glagolitic Letter Nashi	U+E010	ⱷ ⱷ◌	<i>Srezn.</i> , p. 59; <i>PsDem</i> , 21v3; <i>MissSin</i> , 20v23
Combining Glagolitic Letter Onu	U+E011	ⱸ ⱸ◌	<i>PsDem</i> , 126r2.; <i>MissSin</i> , 43v14, 43v21-2
Combining Glagolitic Letter Pokoji	U+E012	ⱹ ⱹ◌	<i>Srezn.</i> , p. 42, p. 248; <i>MissSin</i> , 13r18, 13v10
Combining Glagolitic Letter Ritsi	U+E013	ⱺ ⱺ◌	<i>Srezn.</i> , p. 42; <i>MissSin</i> , 35v15
Combining Glagolitic Letter Slovo	U+E014	ⱻ ⱻ◌	<i>Srezn.</i> , p. 228; <i>PsDem</i> , 126r24, 128v7, 126r24, 126v18
Combining Glagolitic Letter Tvrido	U+E015	ⱼ ⱼ◌	<i>Srezn.</i> , p. 36, p. 42, 59; <i>PsDem</i> , 118r4; <i>MissSin</i> , 18v3, 17v12, 45v20, 46r19

Name	Codepoint	Appearance	Location in Sources
Combining Glagolitic Letter Uku	U+E016	꙱ ⦿	<i>MissSin</i> , 43v14
Combining Glagolitic Letter Fritu	U+E017	꙲ ⦿	<i>MissSin</i> , 25r13, 29r14, 33r23, 40r15, 22r10
Combining Glagolitic Letter Heru	U+E018	꙳ ⦿	<i>MissSin</i> , 19r12, 13v15, 17v14
Combining Glagolitic Letter Shta	U+E01B	ꙴ ⦿	<i>MissSin</i> , 58r6, 72r/v8
Combining Glagolitic Letter Tsi	U+E01C	ꙵ ⦿	Srezn., p. 59; <i>MissSin</i> , 19r12
Combining Glagolitic Letter Chrivi	U+E01D	ꙶ ⦿	<i>MissSin</i> , 53r5, 45r9-10, 39v13
Combining Glagolitic Letter Sha	U+E01E	ꙷ ⦿	Srezn., p. 224; <i>MissSin</i> , 21r12, 13r1, 17r18, 22v18, 24r23, 24(15)v8
Combining Glagolitic Letter Yeru	U+E01F	ꙸ ⦿	<i>MissSin</i> , 18v6
Combining Glagolitic Letter Yeri	U+E020	ꙹ ⦿	Srezn., p. 42
Combining Glagolitic Letter Yati	U+E021	ꙺ ⦿	Srezn., p. 248; <i>PsDem</i> , 83r16; <i>MissSin</i> , 40r15
Combining Glagolitic Letter Yu	U+E023	ꙻ ⦿	Srezn., p. 254
Combining Glagolitic Letter Small Yus	U+E024	꙼ ⦿	<i>EuchSinV</i> , 32v17, 51r11
Combining Glagolitic Letter Yo	U+E026	꙾ ⦿	Does not exist as a single character, but is a component of U+E029.
Combining Glagolitic Letter Iotated Small Yus	U+E027	ꙿ ⦿	Srezn., p. 248; <i>MissSin</i> , 19r12
Combining Glagolitic Letter Big Yus	U+E028	꙽ ⦿	Mansvetov, p. 362 (given by Mansvetov in Cyrillic transcription only)
Combining Glagolitic Letter Iotated Big Yus	U+E029	꙾ ⦿	<i>EuchSinV</i> , 5r10, 17r1, 17v18, 28v3, 29r18 etc.
Combining Glagolitic Letter Fita	U+E02A	ꙿ ⦿	<i>Assem</i> , 125v5, 149v25; (<i>EuchSinV</i> , 83v14, 84v15, 85v4 for /f/)

The following entries are proposed for addition to UnicodeData.txt (note that all codepoints are provisional):

```

E000;COMBINING GLAGOLITIC LETTER AZU;Mn;230;NSM;;;;;N;;;;;
E001;COMBINING GLAGOLITIC LETTER BUKI;Mn;230;NSM;;;;;N;;;;;
E002;COMBINING GLAGOLITIC LETTER VEDE;Mn;230;NSM;;;;;N;;;;;
E003;COMBINING GLAGOLITIC LETTER GLAGOLI;Mn;230;NSM;;;;;N;;;;;

```


Section 4. Justification for Encoding

In this section, we explain the rationale behind encoding combining Glagolitic letters in the Unicode standard. Our rationale can be summarized as follows. First, we show that combining Glagolitic characters are in fact distinct from – and in the presentation of text should be handled separately from – their respective base forms. Second, we demonstrate that other possible approaches for handling combining letters in the Glagolitic script – the use of *ad hoc* markup, encoding in the PUA, and the use of advanced font features – are insufficient or overly complex. Finally, we argue that because combining Cyrillic characters have already been encoded in Unicode, correct interoperability between the two scripts demands that combining Glagolitic characters be also encoded.

4.1 Distinction and Use of Combining Characters in Church Slavonic

As we stated in the Introduction section above, superscription in Church Slavonic is used in two instances: in abbreviations (for example, the word *милостѣ* (mercy) is often written as *млѣтѣ*) and in *nomina sacra* (for example, the spelling *гдѣ* (Lord) is used when it refers to God and the spelling *господѣ* is used when it refers to a secular ruler (a lord), much the same way that capitalization is used in many modern languages).¹ Thus, superscription is a required feature of writing Church Slavonic. Unlike in modern English and other languages where superscription is a stylistic embellishment (e.g., in writing “2nd” as opposed to “2nd”), the superscripted characters in Church Slavonic are combining characters that act like true diacritical marks. In particular, these characters (both in the Cyrillic and Glagolitic scripts) are non-spacing characters, while the “n” and “d” in writing “2nd” in English are spacing characters. Handling such non-spacing characters by positioning spacing characters over a base character would not be correct from the standpoint of text processing in the Unicode standard. In fact, as far as text processing is concerned, the combining characters in Church Slavonic are in no way different from any of the other diacritical marks already encoded.

In his review of L2/14-103, David Birnbaum concedes that “standard modern ChSl [Church Slavonic] orthography does require superscript letters in some words.” But then he goes on to write, “I would have regarded the use of the “wrong” letter as culturally incorrect but nonetheless informationally adequate.” We strongly disagree with this premise. In many instances, using the inline letter instead of the combining letter is not only “culturally incorrect” but also “informationally inadequate”. For example, the sequence *ѡг[ъ]* means “god” (a pagan deity) while the sequence *ѡѣ* (with the combining letter Ge) is an abbreviation for *ѡгосподиченъ* (“theotokion” – a type of liturgical hymn). While one could write *ѡѣ* by using markup-level superscription, this “r” is still a spacing character. One could use kerning at the font level to force the “r” to position over the “o”, but such an approach is not correct from the standpoint of text processing. Moreover, under such an approach, the correct meaning of a text stream would not only be determined by its characters but also by markup (or formatting) and by font-level attributes, which cannot be exchanged between users in a text-only format.

In addition, it is also incorrect to write all Church Slavonic words in their full, unabbreviated form (resolving all abbreviations). As we have seen above, *гдѣ* and *господѣ* have two different meanings; thus, writing *гдѣ* as *господѣ* is “informationally inadequate,” not just “culturally incorrect”.

¹ Throughout we present examples in the Cyrillic script, since it is more familiar to the reader. We will then demonstrate that all of the arguments also hold true for the Glagolitic script because of the relationship that exists between the two scripts.

It is true that strict orthographic conventions did not take shape in Church Slavonic until after 1700 (with the publication of the Elizabeth Bible in 1751) and that Church Slavonic writing of earlier recensions (especially before the advent of the printing press) demonstrates a greater degree of leeway in spelling. Nonetheless, wherever combining letters occur in Church Slavonic, they are always treated as non-spacing marks. Treating them differently in computer-encoded text would be an unreasonable limitation.

In fact, combining characters are already encoded in Unicode for a variety of writing systems. In addition to the combining characters used for Cyrillic, Unicode includes a variety of combining characters used in writing classical Arabic (for example, Honorifics and Koranic annotation signs; for a discussion of these, see L2/01-425); and the various combining Latin characters used for the representation of mediæval texts (see ISO/IEC JTC 1/SC 2/WG 2 N2266 and L2/06-027). In their proposal for combining Latin characters, Everson *et. al.* argue that these characters are necessary to make possible a representation of mediæval text that “does not entail the replacement or the distortion of the original character set.” Similarly, a distortion of the Glagolitic character set is undesirable and so combining Glagolitic characters are needed.

4.2 Alternative Approaches

In addition to encoding, three methods to handle combining Glagolitic characters could be contemplated: the use of an *ad hoc* markup language, encoding the characters in the Private Use Area (PUA), and the use of advanced typographic features available in some font technologies.

Use of *Ad Hoc* Markup

The use of *ad hoc* markup languages is the approach proposed by Ralph Cleminson in his response to L2/14-103. Cleminson writes: “[i]t was recommended ... that in those texts (the majority) where superscription is largely a matter of scribal whim, it should be encoded using markup.” Of course, given the subsequent encoding of a large number of combining Cyrillic characters in the Unicode standard, this recommendation – which is nowhere articulated in the technical documentation to the Unicode standard – has not been followed. Nonetheless, let us briefly consider this approach. Under this scheme, an *ad hoc* markup language is used to mark up text and to indicate that a certain letter is a superscript. In fact, this approach is not new; it was first developed for Church Slavonic with the creation of the HIP (Hyperinvariant Presentation) technology, a markup language used to represent Church Slavonic with an 8-bit codepage that claims to be a “platform-independent representation of Church Slavonic texts ... designed to record the texts in a readable form.” In HIP, the backslash character is used to indicate a superscript; thus, the above example is encoded as:

б0\г

This approach has a number of problems. Though HIP claims to be “readable” (in the sense that it is a set of mnemonic conventions), it is in fact *not* “readable” in the sense that it is not a final presentation form. For example, the text б0\г cannot be presented to an end user on a webpage or in a printed edition. The HIP format – and any markup approach – can only be used for storing text in a plain text format and requires the use of a processor program that converts the markup-encoded text into a format that can be finally presented. The reliance on a processor or converter means that the user of the stored text is limited to a specific set of software tools (those that support add-ons, like Microsoft

Office or those that can be scripted, like the various flavors of TeX).² It also means that in addition to the use of the markup language we must use one of the other two approaches (encoding in the PUA or the use of stylistic alternatives) to present the final representation. For example, the processor can convert the HIP-encoded text to a Unicode presentation where the superscript characters are encoded in the PUA. In practice, since a way to work with the final representation of the text without relying on markup now becomes available (the post-processor output), the markup language and the processor itself become obsolete. In other words, as soon as post-processed text free of markup becomes usable, no one will use the text with markup (this is in fact what happened to HIP).

Use of the Private Use Area

The second approach is to encode the combining characters in the Private Use Area of Unicode. At the font level, the positioning of the combining character over the base character is handled via the use of the mark-to-base positioning feature (*mark*) in OpenType. While this approach is attractive because of its simplicity, it presents several problems. First, if the combining characters warrant their own codepoints in the PUA, why do they not warrant codepoints in the body of the Unicode standard? While various attempts to standardize the PUA have been undertaken (for example, the Medieval Unicode Font Initiative), most of these attempts see the use of the PUA as a temporary solution until the necessary characters can be formally proposed for encoding into the standard. Second, one must keep in mind that the Unicode standard is more than encoding: it also defines, for example, character properties, line breaking, and collation. None of these is well-defined for codepoints in the PUA. Characters in the PUA have a `General_Category` property of `Co` (other, private use). In principle, vendors can agree to override the `General_Category` property; to do this, they must “exchange privately defined data which describes [*sic.*] how each private-use character is to be interpreted” (Unicode Standard, p. 558). But since “the Unicode Standard provides no predefined format for such data” (*ibid.*), in practice this means that such a scheme cannot be contemplated beyond a small, tightly-knit user community using a limited set of software. Certainly, this would not allow for easy exchange of Church Slavonic texts over the Internet in a way that is supported by all browsers and mobile devices. Furthermore, the lack of a specified collation table would preclude proper indexation of the text by major search engines and hamper other string manipulation operations. Finally, many software applications do not correctly support OpenType features (including mark-to-base positioning) for characters in the PUA, so the correct rendering of the text would fail. All in all, the use of the PUA can only be contemplated as a temporary measure until such a time as the characters can be formally accepted into the Unicode standard.

Use of Advanced Font Features

The other approach is the use of advanced typographic features, such as the use of stylistic alternatives (*salt*) or stylistic sets in OpenType or the use of custom features in SIL Graphite. For example, the main glyph *r* can be stored as the base character in the font and the superscripted glyph *ṛ* can be stored as *stylistic alternative 1*. The user can select the superscripted glyph by selecting *stylistic alternative 1* in a text rendering application. The positioning of the superscripted glyph is again handled via mark-to-base positioning.

² A further problem becomes that markup used for character presentation must be distinguished from internal markup used for text style. For example, how does one store information about font color together with the markup? Would it even be possible to have a universal format that can be used in all software – Notepad, TeX, the various Offices, web browsers, mobile devices, and so forth? (For example, the backslash character is a control character in TeX and its use in HIP makes the implementation of HIP in a TeX context a logistical nightmare).

While this approach looks powerful at first glance, it in fact fails spectacularly. First of all, not all software supports advanced typographic features (how does one select stylistic alternatives in Notepad or even in LibreOffice?) Second, even in software where stylistic alternatives can be selected, indicating that a certain character in a character stream has a certain property (such as being a stylistic alternative) usually interrupts the character stream, preventing the mark-to-base positioning feature from working properly, so the practical rendering of combining characters becomes impossible. Now it could be argued that implementation bugs in software need not concern the Unicode Technical Committee. But why design an encoding model that cannot be used in practice given the existing software?

The third issue with this approach is that it does not provide for any way to store information about the superscript letters in a software- and platform-independent setting. While information about stylistic alternatives can be stored in a specific file format (for example, within a Microsoft Office DOCX file), it can be shared across platforms and applications only via the use of a markup language. For example, one could record the sequence in question in an XML-like format as:

```
60<salt alt="1">r</salt>
```

In addition to all of the problems involving the use of markup languages discussed above, this approach suffers from further interoperability limitations. No standard for stylistic alternatives exists: font developers can choose to support the combining character as *stylistic alternative 1*; or, they can choose to support it as *alternative 2* and support something else (for example, the italic Russian *z*) as *alternative 1*; or, they can choose not to support stylistic alternatives at all. At least the PUA provides a limited set of codepoints so that some informal agreement between vendors could be reached; under the stylistic alternatives approach, if we also allow for SIL Graphite features (which have vendor-defined names), the number of ways to render the combining letters becomes non-countable and exchange of data across applications and systems in a standard way becomes impossible. Moreover, since stylistic alternatives are designed to be optional, the end user would not even see an error message when a given alternative is not available in a font.

4.3 Interoperability of Cyrillic and Glagolitic

In addition to the reasons presented above, there is one further reason why combining Glagolitic letters should be encoded in the Unicode standard. Namely, the Cyrillic and Glagolitic scripts are closely related; combining characters have already been encoded for Cyrillic; a methodology needs to exist for the simple conversion of texts between Cyrillic and Glagolitic; thus, combining Glagolitic letters should also be encoded.

As we have pointed out in the Introduction section, scholars believe that Glagolitic was the first script used to record Church Slavonic. After the introduction of the Cyrillic script, it gradually replaced Glagolitic in most Slavic cultures. Nonetheless, Cyrillic and Glagolitic scripts are often used interchangeably in the academic community, and it is quite common to publish Glagolitic textual sources in Cyrillic transcription (see Figure 14). Because of this interoperability of Cyrillic and Glagolitic, computer software needs to be able to unambiguously convert (transliterate) between the two writing systems. In the Table below, we present the transliteration scheme commonly used and due to Jagic (1879):

Glagolitic			Cyrillic		
Name	Codept.	Disp.	Name	Codept.	Disp.
Glagolitic Letter Azu	U+2C30	ⱦ	Cyrillic Letter A	U+0430	А
Glagolitic Letter Buki	U+2C31	Ⱨ	Cyrillic Letter Be	U+0431	Б
Glagolitic Letter Vede	U+2C32	ⱨ	Cyrillic Letter Ve	U+0432	В
Glagolitic Letter Glagoli	U+2C33	Ⱪ	Cyrillic Letter Ge	U+0433	Г
Glagolitic Letter Dobro	U+2C34	ⱪ	Cyrillic Letter De	U+0434	Д
Glagolitic Letter Yestu	U+2C35	ⱬ	Cyrillic Letter Ie	U+0435	Е
Glagolitic Letter Zhivete	U+2C36	Ɱ	Cyrillic Letter Zhe	U+0436	Ж
Glagolitic Letter Dzelo	U+2C37	Ɒ	Cyrillic Letter Dze	U+0455	З
Glagolitic Letter Zemlja	U+2C38	Ⱳ	Cyrillic Letter Ze	U+0437	З
Glagolitic Letter Izhe	U+2C39	ⱴ	Cyrillic Letter Ukrainian I	U+0456	І
Glagolitic Letter Initial Izhe	U+2C3A	ⱶ	Cyrillic Letter Iota	U+A647	Ї
Glagolitic Letter I	U+2C3B	ⱸ	Cyrillic Letter I	U+0438	И
Glagolitic Letter Djervi	U+2C3C	ⱺ	Cyrillic Letter Djerv	U+A649	Ј
Glagolitic Letter Kako	U+2C3D	ⱼ	Cyrillic Letter Ka	U+043A	К
Glagolitic Letter Ljudie	U+2C3E	ⱼ	Cyrillic Letter El	U+043B	Л
Glagolitic Letter Myslite	U+2C3F	ⱼ	Cyrillic Letter Em	U+043C	М
Glagolitic Letter Nashi	U+2C40	ⱼ	Cyrillic Letter En	U+043D	Н
Glagolitic Letter Onu	U+2C41	ⱼ	Cyrillic Letter O	U+043E	О
Glagolitic Letter Pokoji	U+2C42	ⱼ	Cyrillic Letter Pe	U+043F	П
Glagolitic Letter Ritsi	U+2C43	ⱼ	Cyrillic Letter Er	U+0440	Р
Glagolitic Letter Slovo	U+2C44	ⱼ	Cyrillic Letter Es	U+0441	С
Glagolitic Letter Tvrido	U+2C45	ⱼ	Cyrillic Letter Te	U+0442	Т
Glagolitic Letter Uku	U+2C46	ⱼ	Cyrillic Letter O Cyrillic Letter U	U+043E U+0443	У
Glagolitic Letter Fritu	U+2C47	ⱼ	Cyrillic Letter Ef	U+0444	Ф
Glagolitic Letter Heru	U+2C48	ⱼ	Cyrillic Letter Ha	U+0445	Х
Glagolitic Letter Otu	U+2C49	ⱼ	Cyrillic Letter Omega	U+0461	Ѡ
Glagolitic Letter Shta	U+2C4B	ⱼ	Cyrillic Letter Shcha	U+0449	Ѣ
Glagolitic Letter Tsi	U+2C4C	ⱼ	Cyrillic Letter Tse	U+0446	Ѥ
Glagolitic Letter Chrivi	U+2C4D	ⱼ	Cyrillic Letter Che	U+0447	Ѧ
Glagolitic Letter Sha	U+2C4E	ⱼ	Cyrillic Letter Sha	U+0448	Ѩ

Glagolitic			Cyrillic		
Glagolitic Letter Yeru	U+2C4F	ꙁ	Cyrillic Letter Hard Sign	U+044A	Ѣ
Glagolitic Letter Yeru, Glagolitic Letter Izhe	U+2C4F U+2C39	ꙁꙑ	Cyrillic Letter Yeru with Back Yer	U+A651	ѢІ
Glagolitic Letter Yeri	U+2C50	ꙁ̇	Cyrillic Letter Soft Sign	U+044C	ѣ
Glagolitic Letter Yati	U+2C51	ꙁ̆	Cyrillic Letter Yat	U+0463	Ѥ
Glagolitic Letter Yu	U+2C53	ꙁ̈	Cyrillic Letter Yu	U+044E	Ѧ
Glagolitic Letter Small Yus	U+2C54	ꙁ̇̆	Cyrillic Letter Little Yus	U+0467	ѧ
Glagolitic Letter Iotated Small Yus	U+2C57	ꙁ̇̆̆	Cyrillic Letter Iotified Little Yus	U+0469	ѧІ
Glagolitic Letter Big Yus	U+2C58	ꙁ̆̆	Cyrillic Letter Big Yus	U+046B	Ѩ
Glagolitic Letter Iotated Big Yus	U+2C59	ꙁ̆̆̆	Cyrillic Letter Iotified Big Yus	U+046D	ѨІ
Glagolitic Letter Fita	U+2C5A	ꙁ̆̇	Cyrillic Letter Fita	U+0473	Ѧ̇
Glagolitic Letter Izhitsa	U+2C5B	ꙁ̆̆̇	Cyrillic Letter Izhitsa	U+0475	Ѧ̇̆

The Table presents the standard transliteration scheme for the main letters of the Cyrillic and Glagolitic scripts. However, we contend that an automated conversion (transliteration) algorithm also needs to correctly handle diacritical marks and combining letters present in the two scripts.

In the encoding model for the Cyrillic script, the combining Cyrillic letters are already available. In their responses to L2/14-103, Cleminson and Birnbaum present an overall negative view of the use of combining characters for Cyrillic. Cleminson believes that only those combining Cyrillic letters that occur in modern Church Slavonic should have been encoded and Birnbaum writes that he “cannot now endorse the inclusion of any additional superscript Cyrillic or Glagolitic characters.”

In hindsight, we agree that the use of combining characters for Cyrillic superscription was not the best approach (although our criticism of this approach is based on different reasons than those presented by Cleminson and Birnbaum). In our view, this approach is too limiting and does not allow for the unambiguous representation of various complexities occurring in mediæval texts, such as, for example, the use of superscription over multiple base letters or the use of multiple combining letters over one base letter. Also, we regret that the standardization of Glagolitic and Cyrillic was not handled simultaneously, leading to discrepancies in the encoding schemes. Nonetheless, given Unicode's stability policy, there is no use now to criticize an existing implementation. Rather, we desire to make the existing implementation more useful. Thus, because combining Cyrillic letters are encoded, we propose that combining Glagolitic letters also need to be encoded so that a meaningful transliteration scheme can be designed. An approach where Cyrillic superscription is handled via combining characters while Glagolitic superscription is rendered using markup or stylistic alternatives would not make such a transliteration scheme possible. But, as we have pointed out, since it is common to present Glagolitic texts in Cyrillic transcription, the availability of such a scheme is necessary. The Table below presents the proposed transliteration scheme for combining letters. Note, again, that all codepoints for Glagolitic characters are provisional.

Glagolitic			Cyrillic		
Name	Codept.	Disp.	Name	Codept.	Disp.
Glagolitic Combining Letter Azu	U+E000	ⱦ	Cyrillic Combining Letter A	U+2DF6	ⱦ
Glagolitic Combining Letter Buki	U+E001	Ⱨ	Cyrillic Combining Letter Be	U+2DE0	Ⱨ
Glagolitic Combining Letter Vede	U+E002	ⱨ	Cyrillic Combining Letter Ve	U+2DE1	ⱨ
Glagolitic Combining Letter Glagoli	U+E003	Ⱪ	Cyrillic Combining Letter Ge	U+2DE2	Ⱪ
Glagolitic Combining Letter Dobro	U+E004	ⱪ	Cyrillic Combining Letter De	U+2DE3	ⱪ
Glagolitic Combining Letter Yestu	U+E005	Ⱬ	Cyrillic Combining Letter Ie	U+2DF7	Ⱬ
Glagolitic Combining Letter Zhivete	U+E006	ⱬ	Cyrillic Combining Letter Zhe	U+2DE4	ⱬ
Glagolitic Combining Letter Zemlja	U+E008	Ɑ	Cyrillic Combining Letter Ze	U+2DE5	Ɑ
Glagolitic Combining Letter Izhe	U+E009	Ɱ	Cyrillic Combining Letter Yi	U+A676	Ɱ
Glagolitic Combining Letter Initial Izhe	U+E00A	Ɐ	Cyrillic Combining Letter Iota	N/A ³	
Glagolitic Combining Letter I	U+E00B	Ɒ	Cyrillic Combining Letter I	U+A675	Ɒ
Glagolitic Combining Letter Djervi	U+E00C	ⱱ	Cyrillic Combining Letter Djerv	U+2DF8	ⱱ
Glagolitic Combining Letter Kako	U+E00D	Ⱳ	Cyrillic Combining Letter Ka	U+2DE6	Ⱳ
Glagolitic Combining Letter Ljudie	U+E00E	ⱳ	Cyrillic Combining Letter El	U+2DE7	ⱳ
Glagolitic Combining Letter Myslite	U+E00F	ⱴ	Cyrillic Combining Letter Em	U+2DE8	ⱴ
Glagolitic Combining Letter Nashi	U+E010	Ⱶ	Cyrillic Combining Letter En	U+2DE9	Ⱶ
Glagolitic Combining Letter Onu	U+E011	ⱶ	Cyrillic Combining Letter O	U+2DEA	ⱶ
Glagolitic Combining Letter Pokoji	U+E012	ⱷ	Cyrillic Combining Letter Pe	U+2DEB	ⱷ
Glagolitic Combining Letter Ritsi	U+E013	ⱸ	Cyrillic Combining Letter Er	U+2DEC	ⱸ
Glagolitic Combining Letter Slovo	U+E014	ⱹ	Cyrillic Combining Letter Es	U+2DED	ⱹ
Glagolitic Combining Letter Tvrido	U+E015	ⱺ	Cyrillic Combining Letter Te	U+2DEE	ⱺ
Glagolitic Combining Letter Uku	U+E016	ⱻ	Cyrillic Combining Letter Monograph Uk	U+2DF9	ⱻ
Glagolitic Combining Letter Fritu	U+E017	ⱼ	(Cyrillic Combining Letter Ef)	(U+A69E)	ⱼ
Glagolitic Combining Letter Heru	U+E018	ⱽ	Cyrillic Combining Letter Ha	U+2DEF	ⱽ

3 This character will be proposed for encoding by the authors in a separate document.

Glagolitic			Cyrillic		
Glagolitic Combining Letter Shta	U+E01B	Ṣ̌	Cyrillic Combining Letter Shcha	U+2DF3	Ш̣
Glagolitic Combining Letter Tsi	U+E01C	Ț̣	Cyrillic Combining Letter Tse	U+2DF0	Ц̣
Glagolitic Combining Letter Chri vivi	U+E01D	Č̣	Cyrillic Combining Letter Che	U+2DF1	Ч̣
Glagolitic Combining Letter Sha	U+E01E	Ṣ̌	Cyrillic Combining Letter Sha	U+2DF2	Ш̣
Glagolitic Combining Letter Yeru	U+E01F	Ț̣	Cyrillic Combining Letter Hard Sign	U+A678	Ѣ̣
Glagolitic Combining Letter Yeri	U+E020	Ț̣	Cyrillic Combining Letter Soft Sign	U+A67A	ѣ̣
Glagolitic Combining Letter Yati	U+E021	Ț̣	Cyrillic Combining Letter Yat	U+2DFA	Ѥ̣
Glagolitic Combining Letter Yu	U+E023	Ț̣	Cyrillic Combining Letter Yu	U+2DFB	Ѧ̣
Glagolitic Combining Letter Small Yus	U+E024	Ț̣	Cyrillic Combining Letter Little Yus	U+2DFD	Ѧ̣
Glagolitic Combining Letter Iotated Small Yus	U+E027	Ț̣	Cyrillic Combining Letter Iotified Little Yus	N/A ⁴	
Glagolitic Combining Letter Big Yus	U+E028	Ț̣	Cyrillic Combining Letter Big Yus	U+2DFE	Ѧ̣
Glagolitic Combining Letter Iotated Big Yus	U+E029	Ț̣	Cyrillic Combining Letter Iotified Big Yus	U+2DFF	Ѧ̣
Glagolitic Combining Letter Fita	U+E02A	Ț̣	Cyrillic Combining Letter Fita	U+2DF4	Ѧ̣

4.4 Collation

In his response, Ralph Cleminson writes that the approach of using combining characters “has serious drawbacks for the processing of electronic text.” He does not specify what these drawbacks are, but we take this to mean the difficulty in comparing strings. Under this encoding methodology, for example, the strings **НАШНХ** and **НАШ^ХН** are not equivalent at the codepoint level even though under certain circumstances one would like to treat them as the same string. However, the proper way to compare these strings is not via direct (codepoint by codepoint) comparison but rather via the use of the Unicode Collation Algorithm (UCA), which allows for multi-level comparison codes. Under an appropriate UCA collation table, the characters **Х** and **Х̣** could be given the same primary weights but different secondary weights (which is in fact what is implemented for Cyrillic and what we propose for Glagolitic). In a context where the above strings need to be treated as identical one can simply perform a comparison ignoring the secondary weights, while in the context where these strings need to be different, one can compare them at the secondary level. In any case, comparing any two Unicode strings codepoint by codepoint is not advisable from the standpoint of text processing.

4 This character will be proposed for encoding by the authors in a separate document.

Section 5. Examples


- 
 1. 300-655P330-65 300-655P33: 30+
 2. 3P333P3 330-65 333-655P3 333P3 333

Figure 1: Combining Glagolitic Letter Azu (boxed in red). Source: Srezn.

06. (верхнюю половину страницы занимает большая
 записка.)
 1. 333 333P3 333P3 333P3 333P3 333P3 333P3 333P3

 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3

Figure 2: Combining Glagolitic Letters Buky (boxed in red); Glagoli (boxed in black); Ljudije (boxed in blue); Tvrido (boxed in yellow); and Yu (boxed in green). Note the use of the Pokrytie. Source: Srezn.


15. 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3

 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3
 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3

Figure 3: Combining Glagolitic Letters Dobro (boxed in blue); Slovo (boxed in red) and Tvrido (boxed in green). Source: Srezn.

333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3

 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3
 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3

Figure 4: Combining Glagolitic Letter Kako (boxed in red). Note the use of the Pokrytie. Source: Srezn.


6.

 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3
 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3
 333P3 333P3 333P3 333P3 333P3 333P3 333P3 333P3

Figure 5: Combining Glagolitic Letter Heru (boxed in orange). Source: Srezn.

References

- Assem*: Codex Assemanianus (Vat. slav. 3 glag.), Macedonia, early 11th c. Available online (printed in *Asemanievo evangelie. Faksimilno izdanie*. Sofia: Nauka i izkustvo, 1981).
- EuchSinV*: Euchologium Sinaiticum, old part (Sin. slav. 37 & SPb. RNB, Glag. 3), Bulgaria, 11th c.; facsimile edition: NAHTIGAL, R. *Euchologium Sinaiticum. Starocerkvenoslovanski glagolski spomenik. I. Fotografski posnetek*. Ljubljana: AZU v Ljubljani, Učiteljska tiskarna, 1941.
- Everson: EVERSON, Michael *et. al.* "Proposal to add medievalist characters to the UCS". Working Group Document L2/06-027. 2006.
- Jagic: JAGIĆ, Vatroslav, *Quattuor Evangeliorum Codex Glagoliticus, olim Zographensis, nunc Petropolitanus*. Berlin: Apud Weidmannos, 1879.
- Mansvetov: MANSVETOV, I. D., "Evchologium Glagolski, Spomenik monastira Sinai brda. Izdao D-r Lavoslav Geitler. U Zagrebu, 1882. (Древнейший Славянский требник, изданный г. Гейтлером по рукописи, находящейся в Синайской библиотеке)". In *Прибавления к Творениям св. отцов*. 1883. vol. 32, part 1, pp. 347-390.
- MissSin*: Sinai Missal (Sin. slav. 5/N), Bulgaria, 11th c.; cf. MIKLAS, H. (ed.): *Glagolitica – Zum Ursprung der slavischen Schriftkultur* (ÖAW, Phil.-hist. Kl., Schriften der Balkan-Kommission, Philologische Abt. 41). Vienna: ÖAW, 2000, pp. 117-129.
- PsDem*: Psalter of Demetrius (Sin. slav. 3/N), Sinai, 11th c. Facsimile edition: *Psalterium Demetrii Sinaitici (monasterii sanctae Catharinae codex slav. 3/N), adiectis foliis medicinalibus*. Ad editionem phototypicam praeparaverunt M. Gau, D. Hürner, F. Kleber, M. Lettner, H. Miklas sub redactione Henrici Miklas. Cum praefationibus sacri monasterii atque Ioannis Tarnanidae (Glagolitica Sinaitica 1). Vienna: Holzhausen, 2012.
- PsSinV*: Psalterium Sinaiticum, old part (Sin. slav. 38), Sinai, 11th c., Facsimile edition: ALTBAUER, M.: *Psalterium Sinaiticum. An 11th Century Glagolitic Manuscript from St. Catherine's Monastery, Mt. Sinai*. Skopje: The Macedonian Academy of Sciences and Arts, Goce Delčev Publishing House, 1971.
- Srez.: SREZNEVSKIJ, I. I., *Древніе глаголическіе памятники, сравнительно съ памятниками кириллицы*, St. Petersburg: Типографія Императорској академіи наук, 1866.

Glagolitic Extended (Proposed)

Warning: all codepoints are provisional!

0	1	2
† ⦿ U+E000	Ɔ ⦿ U+E010	⦿ ⦿ U+E020
⦿ ⦿ U+E001	⦿ ⦿ U+E011	⦿ ⦿ U+E021
⦿ ⦿ U+E002	⦿ ⦿ U+E012	U+E022
⦿ ⦿ U+E003	⦿ ⦿ U+E013	⦿ ⦿ U+E023
⦿ ⦿ U+E004	⦿ ⦿ U+E014	⦿ ⦿ U+E024
⦿ ⦿ U+E005	⦿ ⦿ U+E015	U+E025
⦿ ⦿ U+E006	⦿ ⦿ U+E016	⦿ ⦿ U+E026
U+E007	⦿ ⦿ U+E017	⦿ ⦿ U+E027
⦿ ⦿ U+E008	⦿ ⦿ U+E018	⦿ ⦿ U+E028
⦿ ⦿ U+E009	U+E019	⦿ ⦿ U+E029
⦿ ⦿ U+E00A	U+E01A	⦿ ⦿ U+E02A
⦿ ⦿ U+E00B	⦿ ⦿ U+E01B	U+E02B
⦿ ⦿ U+E00C	⦿ ⦿ U+E01C	U+E02C
⦿ ⦿ U+E00D	⦿ ⦿ U+E01D	U+E02D
⦿ ⦿ U+E00E	⦿ ⦿ U+E01E	U+E02E
⦿ ⦿ U+E00F	⦿ ⦿ U+E01F	U+E02F

U+E000: Combining Glagolitic Letter Azu
 U+E001: Combining Glagolitic Letter Buki
 U+E002: Combining Glagolitic Letter Vede
 U+E003: Combining Glagolitic Letter Glagoli
 U+E004: Combining Glagolitic Letter Dobro
 U+E005: Combining Glagolitic Letter Yestu
 U+E006: Combining Glagolitic Letter Zhivete
 U+E007: <not assigned>
 U+E008: Combining Glagolitic Letter Zemlja
 U+E009: Combining Glagolitic Letter Izhe
 U+E00A: Combining Glagolitic Letter Initial Izhe
 U+E00B: Combining Glagolitic Letter I
 U+E00C: Combining Glagolitic Letter Djervi
 U+E00D: Combining Glagolitic Letter Kako
 U+E00E: Combining Glagolitic Letter Ljudie
 U+E00F: Combining Glagolitic Letter Myslite

U+E010: Combining Glagolitic Letter Nashi
 U+E011: Combining Glagolitic Letter Onu
 U+E012: Combining Glagolitic Letter Pokoji
 U+E013: Combining Glagolitic Letter Ritsi
 U+E014: Combining Glagolitic Letter Slovo
 U+E015: Combining Glagolitic Letter Tvrido
 U+E016: Combining Glagolitic Letter Uku
 U+E017: Combining Glagolitic Letter Fritu
 U+E018: Combining Glagolitic Letter Heru
 U+E019: <not assigned>
 U+E01A: <not assigned>
 U+E01B: Combining Glagolitic Letter Shta
 U+E01C: Combining Glagolitic Letter Tsi
 U+E01D: Combining Glagolitic Letter Chrivi
 U+E01E: Combining Glagolitic Letter Sha
 U+E01F: Combining Glagolitic Letter Yeru

 U+E020: Combining Glagolitic Letter Yeri
 U+E021: Combining Glagolitic Letter Yati
 U+E022: <not assigned>
 U+E023: Combining Glagolitic Letter Yu
 U+E024: Combining Glagolitic Letter Small Yus
 U+E025: <not assigned>
 U+E026: Combining Glagolitic Letter Yo
 U+E027: Combining Glagolitic Letter Iotated Small Yus
 U+E028: Combining Glagolitic Letter Big Yus
 U+E029: Combining Glagolitic Letter Iotated Big Yus
 U+E02A: Combining Glagolitic Letter Fita
 U+E02B: <not assigned>
 U+E02C: <not assigned>
 U+E02D: <not assigned>
 U+E02E: <not assigned>
 U+E02F: <not assigned>