

Network Working Group
Internet-Draft
Intended status: Informational
Expires: November 13, 2010

S. Sharikov
Regtime Ltd
D. Miloshevic
Afilias
J. Klensin
May 12, 2010

**Internationalized Domain Names Registration and Administration
Guideline for European languages using Cyrillic
draft-sharikov-idn-reg-05.txt**

Abstract

This document is a guideline for Registries and Registrars on registering internationalized domain names (IDNs) based on (in alphabetical order) Bosnian, Bulgarian, Byelorussian, Kildin Sami, Macedonian, Montenegrin, Russian, Serbian, and Ukrainian languages in a DNS zone. For completeness of the "European" languages, it also discusses the additional characters needed for Moldovan when it is written in Cyrillic script. It describes appropriate characters for registration and variant considerations for characters from Greek and Latin scripts with similar appearances and/or derivations.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on November 13, 2010.

Copyright Notice

Copyright (c) 2010 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents

(<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

This document may contain material from IETF Documents or IETF Contributions published or made publicly available before November 10, 2008. The person(s) controlling the copyright in some of this material may not have granted the IETF Trust the right to allow modifications of such material outside the IETF Standards Process. Without obtaining an adequate license from the person(s) controlling the copyright in such materials, this document may not be modified outside the IETF Standards Process, and derivative works of it may not be created outside the IETF Standards Process, except to format it for publication as an RFC or to translate it into languages other than English.

Table of Contents

1.	Introduction	4
1.1.	Similar Characters and Variants	5
1.2.	Terminology	6
2.	Languages and Characters	6
2.1.	Bosnian and Serbian	7
2.2.	Bulgarian	7
2.3.	Byelorussian	7
2.4.	Kildin Sami	7
2.5.	Macedonian	8
2.6.	Moldovan	8
2.7.	Montenegrin	8
2.8.	Russian	9
2.10.	Ukrainian	9
3.	Language-based Tables	9
4.	Table processing rules	9
5.	Table Format	10
6.	Steps after registering an input label	10
7.	Acknowledgments	11
8.	References	11
8.1.	Normative References	11
8.2.	Informative References	12
Appendix A.	European Cyrillic Character Tables	13
Appendix B.	Change Log	18
B.1.	Changes between -02 and -03 and comments about -03	19
B.2.	Changes between -03 and -04	19
B.3.	Changes between -04 and -05	19
	Authors' Addresses	19

1. Introduction

Cyrillic is one of a fairly small number of scripts that are used, with different subsets of characters, to write a large number of languages, some of which are not closely related to the others. When those languages might be used together in a zone (typical of generic TLDs (gTLDs) but likely in other zones both at and below the root), special considerations for intermixing characters may apply. Cyrillic also has the property that, while it is usually considered a separate script from the Latin (Roman) and Greek ones, it shares many characters with them, creating opportunities for visual confusion. Those difficulties are especially pronounced with "all of Cyrillic" is used rather than only the characters associated with a particular language.

This specification provides guidelines for the use of Cyrillic, as encoded in Unicode [[Unicode52](#)] with internationalized domain name (IDN) labels derived from most "European" languages that use the script (use of the term "European" is a convenience, since there is disagreement about the relevant boundaries for different purposes and, of course, much of Russia lies within geological Asia). Specifically it covers (in alphabetic order) Bosnian, Bulgarian, Byelorussian, the Kildin member of the Sami (often written "Saami") language family, Macedonian, Montenegrin, Russian, Serbian, and Ukrainian. Supplemental tables, based on information in the Unicode Standard and a recently-completed Montenegrin government standard [[MontenegrinChars](#)] are provided for use with Montenegrin. Moldovan is no longer in official use with Cyrillic script and no registrations are considered likely in Cyrillic, at least within the relevant ccTLD. Languages of Asia that use Cyrillic are not considered here and should be the subject of separate specifications.

While Cyrillic script is the primary one used for many of the relevant languages and countries, Latin script is often used instead of, or in combination with, it. Standard keyboards used in most of the countries have both Cyrillic and Latin characters. Therefore some registries could use Latin scripts for domain names registration in their zones. From time to time, some registries and users have claimed that there is a requirement for mixing Cyrillic and Latin characters in the same label. We strongly recommend against such mixing as user confusion is almost certain to result. In addition, registries that support many scripts will probably encounter the need to support labels in Greek or Latin scripts as well as Cyrillic and a large number of character forms are shared among those three scripts.

Because the DNS has no way for the end-user to distinguish among the languages that might have been used to inspire a particular label, it seems useful to treat the characters of a large number of languages

that use Cyrillic in their writing systems together, rather than trying to differentiate them. The discussion and tables in this specification should provide a foundation for developing more restrictive rules for zones in which only a single language is likely to be used, but it does not specify those language-specific rules.

Readers of this document should be aware that its recommendations are about use in DNS labels. The orthography for some of the languages involved, especially Kildin Sami, is not completely standardized and local usage sometimes permits substitution of Latin-based characters for their Cyrillic equivalents. Unless they are required by official orthographies, those substitutions should generally be avoided in DNS labels because of the risk of additional user confusion with the similar-appearing Latin characters.

1.1. Similar Characters and Variants

For some human languages, there are characters and/or strings that have equivalent or near-equivalent meanings. If someone is allowed to register a name with such a character or string, the registry might want to automatically register all the names that have the same meaning in that language. Further, some registries might want to restrict the set of characters to be registered for language-based reasons. In addition, IDNA [[RFC3490](#)] allows the use of thousands of non-alphanumeric characters, and some zone administrators will want to prohibit some or all of these characters.

So-called "variant techniques", introduced in [[RFC3743](#)] and generalized beyond East Asian language in [[RFC4290](#)], describe ways of registering IDN domain names to decrease the risk of misunderstandings, cybersquatting, and other forms of confusion.

The tables below (Appendix A) identify confusable characters in Latin and Greek scripts that might be easily confused with Cyrillic ones.

As with variant approaches for other scripts (e.g., see [RFC 4713](#) for the Chinese language [[RFC4713](#)] or [RFC 5564](#) for the Arabic language [[RFC5564](#)]), this document identifies sets of characters that need special consideration and provides information about them. A registry that handles names using these characters can then make a policy decision about how to actually handle them. The options for those policy decisions would include automatically registering all look-alike string to the same registrant, registering one such string and blocking the others, and so on.

1.2. Terminology

The terminology that follows is derived from [[RFC3743](#)] and [[RFC4290](#)], but this specification does not depend on them. All characters listed here have been verified to be "PVALID" under the recently-adopted IDNA2008 specification [[IDNA2008-Defs](#)].

A "string" is a sequence of one or more characters.

This document discusses characters that have equivalent or near-equivalent characters or strings. The "base character" is the character that has one or more equivalents; the "variant(s)" are the character(s) and/or string(s) that are equivalent to the base character.

A "registration bundle" is the set of all labels that comes from expanding all base characters for a single name into their variants.

A registry is the administrative authority for a DNS zone. That is, the registry is the body that makes and enforces policies that are used in a particular zone in the DNS. The term "registry" applies to all zones in the DNS, not only those that exist at the top level.

2. Languages and Characters

In the interest of clarity and balance, this document describes a "Base Cyrillic" set of twenty-three characters for use in comparing the character usage for Russian and Central European languages that use Cyrillic. The balance of this section compares the character usage of the individual languages in that group.

"Base Cyrillic" consists of the following Unicode code points (names associated with these code points and those below appear in [Appendix A](#)): U+0430, U+0431, U+0432, U+0433, U+0434, U+0435, U+0436, U+0437, U+043A, U+043B, U+043C, U+043D, U+043E, U+043F, U+0440, U+0441, U+0442, U+0443, U+0444, U+0445, U+0446, U+0447, U+0448.

In addition, modern writing systems that use Cyrillic do not have digits separate from the "European" ones used with Latin characters. For registries that permit digits to appear in domain name labels, the "Base Cyrillic" code point listed above should be considered to include U+0030, U+0031, U+0032, U+0033, U+0034, U+0035, U+0036, U+0037, U+0038, and U+0039 (Digit Zero, and Digit One, through Digit Nine). The Hyphen-Minus character (U+0029) may also be used.

It is worth noting that the EU top-level domain registry allows Cyrillic registrations using 32 code points [[EU-registry](#)]. That list

is sufficient for some of the languages listed here but not for others.

The individual languages that are the focus of this specification are discussed below (in English alphabetical order):

2.1. Bosnian and Serbian

Bosnian and Serbian have 30 letters in the alphabet and the additional seven characters to the base of 23 shared Cyrillic characters: U+0438, U+0458, U+0452, U+0459, U+045A, U+045B, U+045F.

2.2. Bulgarian

The Bulgarian alphabet has thirty characters, seven in addition to the basic twenty-three: U+0438, U+0439, U+0449, U+044A, U+044C, U+044E, U+044F.

2.3. Byelorussian

Byelorussian alphabet has 32 characters, i.e., nine characters in addition to the Base Cyrillic set of 23 characters: U+0451, U+0456, U+0439, U+044B, U+044C, U+045E, U+044D, U+044E, U+044F.

2.4. Kildin Sami

The phonetics of the Kildin Sami are quite complex and not easily represented in Cyrillic (see, e.g., [[Kert](#)]). The orthography is not standardized and the writing system may best be thought of as an attempt to transcribe the language phonetically (primary in Latin script in the 1930s but in Cyrillic more recently). Different scholars have reported different numbers of phonemes, further complicating the transcription process. Kertom identifies 53 consonants with long-short distinctions and, in many cases, hard-soft ones. He also identifies ascending and descending diphthongs and one triphthong as well as more common short and long vowels.

The primary reference for Kildin Sami that is apparently used by Sami language(s) experts in Scandinavian countries [[Riessl07](#)] and the references it cites, uses 56 characters, 33 of which do not appear in the basic set. Eight* of these characters have no precomposed forms in Unicode and hence must be written as a two-code-point sequence including U+0304 (Combining Macron). Using parentheses to make the two-code-point sequences more obvious, the additional characters are: (U+0430 U+0304)*, (U+0435 U+0304)*, U+0438, U+0439, (U+043E U+0304), U+044A, U+044B, (U+044B U+0304), U+044C, U+044D, (U+044D U+0304), U+044E, (U+044E U+0304), U+044F, (U+044F U+0304), U+0451, (U+0451 U+0304), U+0458, U+048B, U+048D, U+048F, U+04BB, U+04C6, U+04C8,

U+04CA, U+04CE, U+04D3, U+04E3, U+04E7, U+04ED, U+04EF, U+04F1, U+04F9.

- * These characters, CYRILLIC SMALL LETTER A with a COMBINING MACRON and CYRILLIC SMALL LETTER IE with a COMBINING MACRON, respectively, have the same visual appearance as LATIN SMALL LETTER A WITH MACRON (U+0101) and LATIN SMALL LETTER E WITH MACRON (U+0113). The combinations are not mapped to the Latin character sequences by NFC (or NFKC) normalization. Substitution of the Latin sequence for the second of these is specified by some sources including Riessler [[Riessler07](#)]. Substitution of the Latin character codepoint for the first sequence is not specified in any reference we found, but the relationship is obvious and may occur outside the user's control based on the keyboard or input functions in use. However, U+0101 and U+0113 are Latin Script characters so, if either is used, any tests on homogeneity of the script within a label need to be made with care. If some input systems produce U+0113 (or U+0101) and others produce the two-character combining sequence, a variant approach may be appropriate.

Similar issues may apply to other Kildin Sami characters constructed with combining sequences.

The key references in Russian [[Anto90](#)], [[Kert86](#)], [[Kuru85](#)] all propose slightly different character tables relative to each other and to Riessler's list. Because the latter list appears to be more comprehensive and to represent more recent scholarship, we have based the tables in this document on it. We recommend, however, that registries review these recommendations and the relevant papers should registration requests for Kildin Sami actually appear.

[2.5.](#) Macedonian

Macedonian has 31 characters in the alphabet. This is eight in addition to the basic set: U+0438, U+0458, U+0452, U+0459, U+045A, U+045C, U+045F, U+0491, U+0455.

[2.6.](#) Moldovan

Cyrillic is no longer in everyday use for Moldovan, so no IDN registrations are anticipated.

[2.7.](#) Montenegrin

According to the most recent, and now final, government specification [[MontenegrinChars](#)], Montenegrin has 32 characters in its alphabet, including two that have no precomposed forms in Unicode. This is

nine in addition to the basic set and two in addition to Bosnian and Serbian: U+0437 U+0302, U+0438, U+0441 U+0302, U+0452, U+0458, U+0459, U+045A, U+045B, U+045F.

See Bosnian, [Section 2.1](#), above.

[2.8.](#) Russian

The current Russian alphabet has 33 characters, consisting of the Base Cyrillic set plus an additional ten characters: U+0451, U+0438, U+0439, U+0449, U+044A, U+044B, U+044C, U+044D, U+044E, U+044F.

[2.9.](#) Serbian

See Bosnian, [Section 2.1](#), above.

[2.10.](#) Ukrainian

The character list for modern Ukrainian has apparently not completely stabilized. Some references claim 31 characters and therefore an additional 8 characters to the Base Cyrillic set of 23. Others claim 33, adding U+0438 and U+0439 and replacing U+044A (hard sign) with U+044C (soft sign), for a total of an additional 11 characters as compared to the Base Cyrillic set. Unless better information is available, the prudent registry should probably assume that all 34 characters are in use, i.e., the Base Cyrillic set plus U+0438, U+0439, U+0454, U+0456, U+0457, U+0491, U+0449, U+044A, U+044C, U+044E, U+044F.

[3.](#) Language-based Tables

The registration strategy described in this document uses a table that lists all characters allowed for input and any variants of those characters. Note that the table lists all characters allowed, not only the ones that have variants.

[4.](#) Table processing rules

The input to the process is called the "input label". The output of the process is either failure (the input label cannot be registered at all), or a registration bundle that contains one or more labels in A-label form.

5. Table Format

The table in [Appendix A](#) consists of four columns. The first and second identify the Cyrillic character and the third and fourth identify Latin or Greek characters that might be easily confused with them visually. If both a Latin and Greek character are present, the Greek one appears in the third and fourth columns on the subsequent line (with "..." in the first column to indicate more information about the character specified on the previous line). Variants needed only because of case folding are shown with "++" in the first column, as noted in the table.

Each character in the table is given in the "U+" notation for Unicode characters followed, in the next column, by its name as shown in the Unicode Standard. For easy reference, the characters are listed in the order in which they appear in the Unicode Standard.

The table does not, and any future revision MUST NOT, have more than one entry for a particular base character.

6. Steps after registering an input label

A registry has at least three policy options for handling the cases where the registration bundle has more than one label. These options, and their key implications, are:

- o Allocate all labels to the same registrant, making the zone information identical to that of the input label.

This option will cause end users to be able to find names with variants more easily, but will result in larger zone files. In principle, the zone file could become so large that it could negatively affect the ability of the registry to perform name resolution.

- o Block all labels so they cannot be registered in the future.

This option does not increase the size of the zone file, but it may cause end users to not be able to find names with variants that they would expect.

- o Allocate some labels and block some other labels.

This option is likely to cause the most confusion with users because including some variants will cause a name to be found, but using other variants will cause the name to be not found.

With any of these three options, the registry **MUST** keep a database that links each label in the registration bundle to the input label. This link needs to be maintained so that changes in the non-DNS registration information (such as the label's owner name and address) is reflected in every member of the registration bundle as well.

7. Acknowledgments

Support from Afiliias for a major portion of this work is appreciated.

The material on Kildin Sami would not have been possible without the efforts of Cary Karp for his help directly and his pointer to [\[Riessl07\]](#) and from Vladimir Shadrinov and Sergey Nikolaevich Teryoshkin for their own analyses and references to [\[Anto90\]](#), [\[Kert86\]](#), and [\[Kuru85\]](#) and partial translations from them. We are grateful for their efforts that facilitated treating it nearly the same way as other actively-used European languages that use Cyrillic script.

Careful reading of late drafts by Bill McQuillan and Alexey Melnikov identified a number of editorial problems, some of which might not have been caught otherwise.

8. References

8.1. Normative References

- [RFC3490] Faltstrom, P., Hoffman, P., and A. Costello, "Internationalizing Domain Names in Applications (IDNA)", [RFC 3490](#), March 2003.
- [RFC3491] Hoffman, P. and M. Blanchet, "Nameprep: A Stringprep Profile for Internationalized Domain Names (IDN)", [RFC 3491](#), March 2003.
- [Unicode52] The Unicode Consortium, "The Unicode Standard, Version 5.2.0", 2009.

Defined by: The Unicode Standard, Version 5.0, Boston, MA, Addison-Wesley, 2007, ISBN 0-321-48091-0, as amended by Unicode 5.1.0 (<http://www.unicode.org/versions/Unicode5.1.0/>) and Unicode 5.2.0 (<http://www.unicode.org/versions/Unicode5.2.0/>).

8.2. Informative References

- [Anto90] Antonova, A., "Sami primer", 1990.
- Published in Russian, no authoritative translation is known.
- [EU-registry] European Registry of Internet Domain Names (EURid), ".eu Supported Characters", January 2010, <<http://www.eurid.eu/en/eu-domain-names/technical-limitations/supported-characters>>.
- [IDNA2008-Defs] Klensin, J., "Internationalized Domain Names for Applications (IDNA): Definitions and Document Framework", January 2010, <<https://datatracker.ietf.org/drafts/draft-ietf-idnabis-defs/>>.
- [Kert] Kertom, G., "Kildin dialect of the Sami language".
- Published in Russian, no authoritative translation is known.
- [Kert86] Kertom, G., "Sami-Russian and Russian-Sami dictionary", 1986.
- Published in Russian, no authoritative translation is known.
- [Kuru85] Kuruch, R., "Sami-Russian dictionary", 1985.
- Published in Russian, no authoritative translation is known.
- [MontenegrinChars] Crna Gora Ministarstvo prosvjete i nauke (Ministry of Science and Education, Montenegro), "Pravopis Crnogorskoga Jezika I", 2009, <<http://www.gov.me/files/1248442673.pdf>>.
- In Montenegrin, no known English translation. See especially the table on page 8.
- [OmniglotSaami] Ager, S., "Sami (Saami)", 2009, <<http://www.omniglot.com/writing/saami.htm>>.
- [RFC3743] Konishi, K., Huang, K., Qian, H., and Y. Ko, "Joint

Engineering Team (JET) Guidelines for Internationalized Domain Names (IDN) Registration and Administration for Chinese, Japanese, and Korean", [RFC 3743](#), April 2004.

- [RFC4290] Klensin, J., "Suggested Practices for Registration of Internationalized Domain Names (IDN)", [RFC 4290](#), December 2005.
- [RFC4713] Lee, X., Mao, W., Chen, E., Hsu, N., and J. Klensin, "Registration and Administration Recommendations for Chinese Domain Names", [RFC 4713](#), October 2006.
- [RFC5564] El-Sherbiny, A., Farah, M., Oueichek, I., and A. Al-Zoman, "Linguistic Guidelines for the Use of the Arabic Language in Internet Domains", [RFC 5564](#), February 2010.
- [Riessl07] Riessler, M., "Kola Saami character chart (draft)", November 2007.

[Appendix A.](#) European Cyrillic Character Tables

These tables are constructed on the basis of the characters that can actually occur in the DNS, i.e., those that can be obtained by applying the ToUnicode operation of [RFC 3490](#) or the U-label transformation [[IDNA2008-Defs](#)] to an ACE-encoded label (A-label) as defined in those specifications. If the characters that can be mapped into those characters are to be considered instead, then the number of variants would increase considerably. For example, while Cyrillic Small Letter A and Greek Small Letter Alpha are readily distinguished visually, their capital letter equivalents are not, so, if the extended set of Nameprep [[RFC3491](#)] mappings are considered, the two small letters must be considered variants of each other.

These additional, possibly-required, variants are shown below with "+++" in the first column of the table.

Characters needed for European languages, other than Moldovan and Sami, written in Cyrillic.

+-----+-----+-----+-----+				
Cyrillic	Unicode Name	Variant	Unicode Name	
Char				
+-----+-----+-----+-----+				
U+0430	CYRILLIC SMALL LETTER	U+0061	LATIN SMALL LETTER	
	A		A	

+++		U+03B0	GREEK SMALL LETTER ALPHA
U+0431	CYRILLIC SMALL LETTER BE		
U+0432	CYRILLIC SMALL LETTER VE	U+0062	LATIN SMALL LETTER B
+++		U+03B2	GREEK SMALL LETTER BETA
U+0433	CYRILLIC SMALL LETTER GHE	U+0072	LATIN SMALL LETTER R
+++		U+03B3	GREEK SMALL LETTER GAMMA
U+0434	CYRILLIC SMALL LETTER DE		
+++		U+03B4	GREEK SMALL LETTER DELTA
U+0435	CYRILLIC SMALL LETTER IE	U+0065	LATIN SMALL LETTER E
+++		U+03B5	GREEK SMALL LETTER EPSILON
U+0436	CYRILLIC SMALL LETTER ZHE		
U+0437	CYRILLIC SMALL LETTER ZE		
U+0438	CYRILLIC SMALL LETTER I	U+0075	LATIN SMALL LETTER U
U+0439	CYRILLIC SMALL LETTER SHORT I		
U+043A	CYRILLIC SMALL LETTER KA	U+006B	LATIN SMALL LETTER K
...		U+03BA	GREEK SMALL LETTER KAPPA
U+043B	CYRILLIC SMALL LETTER EL		
+++		U+039B	GREEK SMALL LETTER LAMBDA
U+043C	CYRILLIC SMALL LETTER EM	U+006D	LATIN SMALL LETTER M
+++		U+03BC	GREEK SMALL LETTER MU
U+043D	CYRILLIC SMALL LETTER EN	U+0068	LATIN CAPITAL LETTER H (and SMALL LETTER H in some fonts)
+++		U+03B7	GREEK SMALL LETTER ETA
U+043E	CYRILLIC SMALL LETTER O	U+006F	LATIN SMALL LETTER O

...		U+03BF	GREEK SMALL LETTER
			OMICRON
U+043F	CYRILLIC SMALL LETTER	U+006E	LATIN SMALL LETTER
	PE		N
...		U+03C0	GREEK SMALL LETTER
			PI
U+0440	CYRILLIC SMALL LETTER	U+0070	LATIN SMALL LETTER
	ER		P
...		U+03C1	GREEK SMALL LETTER
			RHO
U+0441	CYRILLIC SMALL LETTER	U+0063	LATIN SMALL LETTER
	ES		C
U+0442	CYRILLIC SMALL LETTER	U+0074	LATIN SMALL LETTER
	TE		T
+++		U+03C4	GREEK SMALL LETTER
			TAU
U+0443	CYRILLIC SMALL LETTER	U+0079	LATIN SMALL LETTER
	U		Y
+++		U+03C5	GREEK SMALL LETTER
			UPSILON
U+0444	CYRILLIC SMALL LETTER	U+03D5	GREEK PHI SYMBOL
	EF		
+++		U+03C6	GREEK SMALL LETTER
			PHI
U+0445	CYRILLIC SMALL LETTER	U+0078	LATIN SMALL LETTER
	HA		X
...		U+03C7	GREEK SMALL LETTER
			CHI
U+0446	CYRILLIC SMALL LETTER		
	TSE		
U+0447	CYRILLIC SMALL LETTER		
	CHE		
U+0448	CYRILLIC SMALL LETTER		
	SHA		
U+0449	CYRILLIC SMALL LETTER		
	SHCHA		
U+044A	CYRILLIC SMALL LETTER		
	HARD SIGN		
U+044B	CYRILLIC SMALL LETTER		
	YERU		
U+044C	CYRILLIC SMALL LETTER	U+0062	LATIN SMALL LETTER
	SOFT SIGN		B
U+044D	CYRILLIC SMALL LETTER		
	E		
U+044E	CYRILLIC SMALL LETTER		
	YU		
U+044F	CYRILLIC SMALL LETTER		
	YA		

U+0451	CYRILLIC SMALL LETTER	U+00EB	LATIN SMALL LETTER
	IO		E WITH DIAERESIS
U+0452	CYRILLIC SMALL LETTER		
	DJE		
U+0453	CYRILLIC SMALL LETTER		
	GJE		
U+0454	CYRILLIC SMALL LETTER	U+03B5	GREEK SMALL LETTER
	UKRAINIAN IE		EPSILON
U+0455	CYRILLIC SMALL LETTER	U+0073	LATIN SMALL LETTER
	DZE		S
U+0456	CYRILLIC SMALL LETTER	U+0069	LATIN SMALL LETTER
	BYELORUSSIAN-UKRAINIAN		I
	I		
+++		U+03B9	GREEK SMALL LETTER
			IOTA
U+0457	CYRILLIC SMALL LETTER	U+03CA	GREEK SMALL LETTER
	UKRAINIAN YI		IOTA WITH DIALYTIKA
+++		U+00EF	LATIN SMALL LETTER
			I WITH DIAERESIS
U+0458	CYRILLIC SMALL LETTER	U+006A	LATIN SMALL LETTER
	JE		J
...		U+03F3	GREEK LETTER YOT
U+0459	CYRILLIC SMALL LETTER		
	LJE		
U+045A	CYRILLIC SMALL LETTER		
	NJE		
U+045B	CYRILLIC SMALL LETTER		
	TSHE		
U+045C	CYRILLIC SMALL LETTER		
	KJE		
U+045D	CYRILLIC SMALL LETTER		
	I WITH GRAVE		
U+045E	CYRILLIC SMALL LETTER		
	SHORT U		
U+045F	CYRILLIC SMALL LETTER		
	DZHE		
U+0491	CYRILLIC SMALL LETTER		
	GHE WITH UPTURN		
U+04C2	CYRILLIC SMALL LETTER		
	ZHE WITH BREVE		
+-----+			

Additional characters needed for Moldovan written in Cyrillic.

Cyrillic Char	Unicode Name	Variant	Unicode Name
U+04C2	CYRILLIC SMALL LETTER ZHE		
	WITH BREVE		
U+0437 + U+0302	Cyrillic Small Letter ZE with Acute		
U+0441 + U+0302	Cyrillic Small Letter ES with Acute		

Additional characters needed for Kildin Sami written in Cyrillic.

Cyrillic Char	Unicode Name	Variant	Unicode Name
U+0430 + U+0304	Cyrillic Small Letter A with Macron	U+0101	LATIN SMALL LETTER A WITH MACRON
...		U+03B0 + U+0304	Greek Small Letter Alpha with Macron
U+0435 + U+0304	Cyrillic Small Letter IE with Macron	U+0113	LATIN SMALL LETTER E WITH MACRON
U+043E + U+0304	Cyrillic Small Letter O with Macron	U+014D	LATIN SMALL LETTER O WITH MACRON
...		U+03BF + U+0304	Greek Small Letter Omicron with Macron
U+044B + U+0304	Cyrillic Small Letter YERU with Macron		
U+044D + U+0304	Cyrillic Small Letter E with Macron		
U+044E + U+0304	Cyrillic Small Letter YU with Macron		
U+044F + U+0304	Cyrillic Small Letter YA with Macron		
U+0451 + U+0304	Cyrillic Small Letter IO with Macron	U+00EB + U0304	Latin Small Letter E With Diaeresis and Macron

U+048B	CYRILLIC SMALL LETTER SHORT I WITH TAIL		
U+048D	CYRILLIC SMALL LETTER SEMISOFT SIGN		
U+048F	CYRILLIC SMALL LETTER ER WITH TICK		
U+04BB	CYRILLIC SMALL LETTER SHHA		
U+04C6	CYRILLIC SMALL LETTER EL WITH TAIL		
U+04C8	CYRILLIC SMALL LETTER EN WITH HOOK		
U+04CA	CYRILLIC SMALL LETTER EN WITH TAIL		
U+04CE	CYRILLIC SMALL LETTER EM WITH TAIL		
U+04D3	CYRILLIC SMALL LETTER A WITH DIAERESIS	U+00E4	LATIN SMALL LETTER A WITH DIAERESIS
U+04E3	CYRILLIC SMALL LETTER I WITH MACRON	U+016B	LATIN SMALL LETTER U WITH MACRON
U+04E7	CYRILLIC SMALL LETTER O WITH DIAERESIS	U+00F6	LATIN SMALL LETTER O WITH DIAERESIS
U+04ED	CYRILLIC SMALL LETTER E WITH DIAERESIS		
U+04EF	CYRILLIC SMALL LETTER U WITH MACRON		
U+04F1	CYRILLIC SMALL LETTER U WITH DIAERESIS		
U+04F9	CYRILLIC SMALL LETTER YERU WITH DIAERESIS		

[Appendix B.](#) Change Log

RFC Editor: Please remove this appendix.

B.1. Changes between -02 and -03 and comments about -03

- o Updated references to Unicode 5.2.
- o Updated information about Montenegrin and Kildin Sami.
- o Removed note about IDNA2003, inserted a comment about IDNA2008 verification, and changed terminology to reflect IDNA2008 where needed.
- o Corrected an error for Bulgarian.
- o Clarified role of this document vis-a-vis orthographies in use in various places.
- o Added text to clarify how the information in this document can be used.
- o Still reviewing Ukrainian (see [Section 2.10](#)) and mixed-script (see [Section 1](#)) requirements (see -04 changes immediately below).

B.2. Changes between -03 and -04

- o Revised text about mixed scripts slightly.
- o Updated material on Kildin Sami.
- o Improved the description of Ukrainian.

B.3. Changes between -04 and -05

- o Fixed several errors in the comparison table appendix.
- o Eliminated one residual case of IDNA2003 terminology.

Authors' Addresses

Sergey Sharikov
Regtime Ltd
Kalinina str.,14
Samara 443008
Russia

Phone: +7(846) 979-9039
Fax: +7(846)979-9038
Email: s.shar@regtime.net

Desiree Miloshevic
Afilias
Oxford Internet Institute, 1 St. Giles
Oxford OX1 3JS
United Kingdom

Phone: +44 7973 987 147
Email: dmiloshevic@afilias.info

John C Klensin
1770 Massachusetts Ave, #322
Cambridge, MA 02140
USA

Phone: +1 617 491 5735
Email: john-ietf@jck.com

