



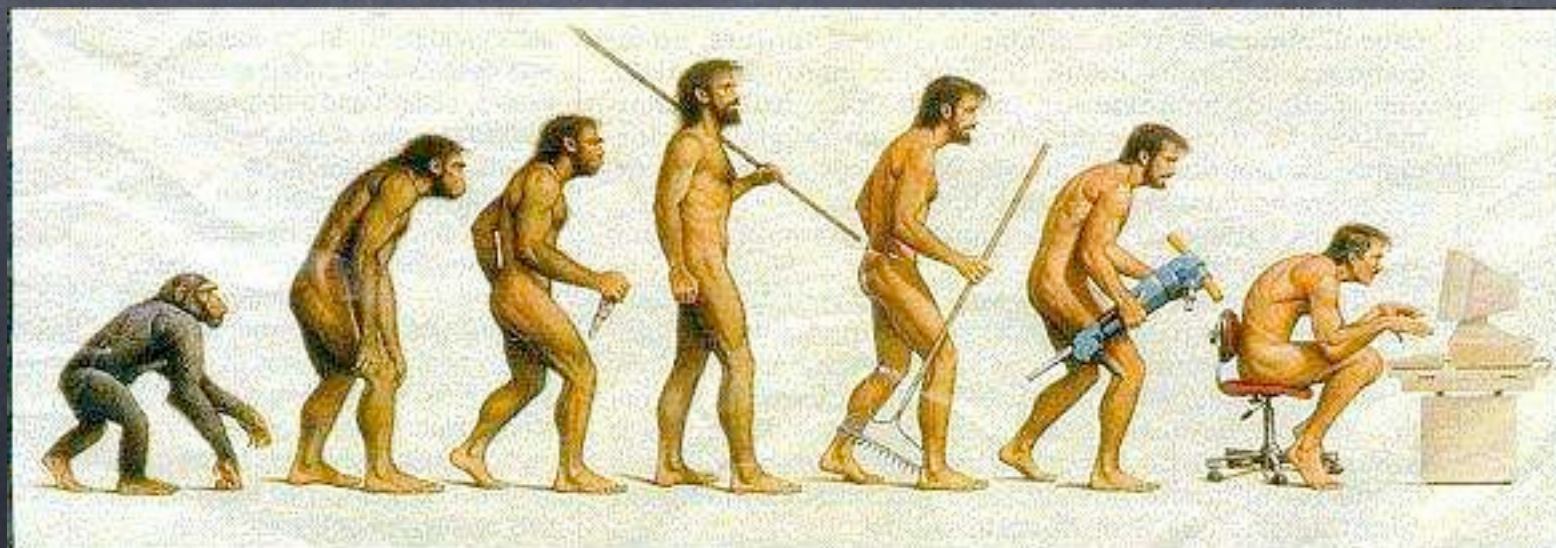
UNICODE 4.1 AND SLAVIC PHILOLOGY – PROBLEMS AND PERSPECTIVES –

Prof. Dr. Sebastian Kempgen
University of Bamberg, Germany



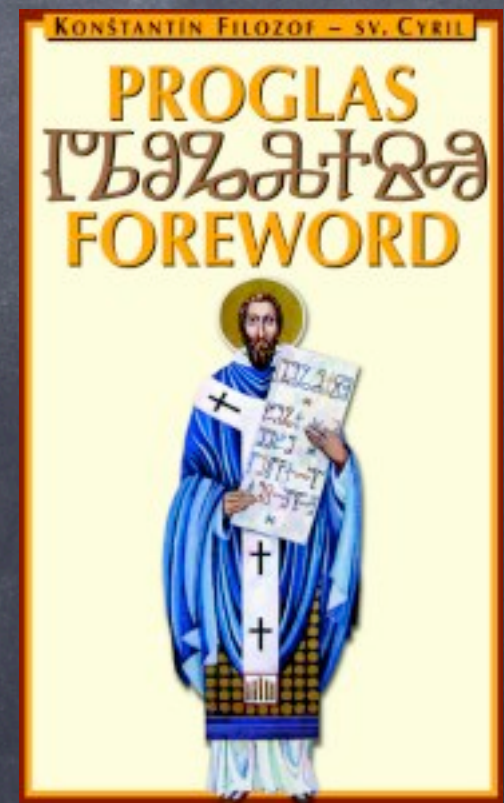
Overview

- 👁 Introduction
- 👁 Status Quo : What's Available
- 👁 Problems: Missing Pieces
- 👁 Perspectives: More Characters



Topics Covered

- Russian Hist. Orthography
- Russian Phonetics
- Polish Hist. Orthography
- Sorbian Hist. Orthography
- Croatian Accents
- Bulgarian Phonetics
- Old Church Slavonic
- Transliteration of Glagolitic
- Balkan Philology



1. Introduction

- ASCII: 2^7 characters = 128
- 'Code Pages': 2^8 characters = 256
- Unicode: 2^{16} characters = 65.536
- Versions: v. 1.0 1991, v. 4.1 2005
- Version 4.1: Glagolica ✓



Unicode Blocks

uni0401	locyril	Djecyr	Gjecyr	Ecyrril	Dzecyr	icyrilli	Yicyril	Jecyri	Ljecyr	Njecyr	Tshecy	Kjecyr	uni0401	Ushort	Dzhecy
È	Ё	Ђ	Ѓ	Є	Ѕ	І	Ї	Ј	Љ	Њ	Ћ	Ќ	Й	Ў	Ц
Acyril	Becyri	Vecyri	Geocyri	Decyri	Iecyri	Zhecyr	Zecyri	Iicyrill	Iishort	Kacyri	Elocyri	Emocyri	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	Shchac	Hardsig	Yericyr	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	iicyrill	iishort	kacyri	elocyri	emocyri	encyri	ocyri	pecyri
а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
ercyri	escyri	tecyril	ucyri	efcyri	khacyr	tsecyr	checyr	shacyr	shchac	hardsig	yericyr	softsig	erever	iucyri	iacyri
р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
uni0451	iocyri	djecyr	gjecyr	ecyrril	dzecyr	icyrilli	Yicyril	Jecyri	Ljecyri	Njecyr	tshecy	Kjecyr	uni0451	Ushort	Dzhecy
è	ё	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	й	ў	ц
Omegar	omegar	Yatcyri	yatcyri	Eiotifie	eiotifie	Yuslitt	yuslitt	Yuslitt	yuslitt	Yusbigi	yusbigi	Yusbigi	yusbigi	Ksicyr	ksicyri
Ω	ω	Ѡ	ѡ	Є	є	А	а	Ӑ	ӑ	Ӓ	ӓ	Ӕ	ӕ	Ӗ	ӗ
Psicyr	psicyri	Fitacyr	fitacyri	Izhitsa	izhitsa	Izhitsa	izhitsa	Ukeyri	ukcyri	Omegar	omegar	Omegar	omegar	Otcyri	lotcyri
Ψ	ψ	Θ	θ	Ʋ	υ	Ỳ	ỳ	Ɔ	ɔ	Ɔ	ɔ	Ẁ	ẁ	Ẃ	ẃ
Koppac	koppac	thousar	titlocyri	palatal	dasiapr	psilipne	uni048	uni048	uni048	Iishort	Iishort	Semis	semiso	Ercyri	ercyri
Ѕ	ѕ	Ѵ	ѵ	Ѷ	ѷ	Ѹ				Й	й	Ђ	ђ	Р	р



Unicode Blocks

- Basic Latin, Latin-1, Latin Extended-A, Latin Extended-B, Latin Extended Additional; Latin Extd C (5.0)
- Cyrillic, Cyrillic Extended (= Non-Sl.)
- Greek, Greek Extended (= Hist.)
- IPA, IPA Extensions + Supplement
- Spacing & Combining Diacritics



‘Private Area’

- 6.400 Slots
- Compatibility between fonts and documents not guaranteed
- Use at your own risk
- Already in use
- Use not coordinated between philologies



Unicode.org

- Submissions, Proposals: Technical and Philological Aspects!!
- Pipeline of characters and scripts
- Rejected scripts and characters
- Published Errata



2. Status Quo

- What is already available?
- What has been accomplished?
- A few appetizers



Russian Translit.

Transliteration for
historical letters:

Θθ: Ħħ (Irish Gaelic)

Vv: Ÿÿ ✓ (no info given)

(Latin Extd-Additional)



Russian Translit.

Miscellaneous phonetic modifiers

- 02B9 ' MODIFIER LETTER PRIME
- primary stress, emphasis
 - transliteration of mjagkij znak (Cyrillic soft sign: palatalization)
- 0027 ' apostrophe
→ 00B4 ´ acute accent
→ 02CA ´ modifier letter acute accent
→ 0301 ˆ combining acute accent
→ 0374 ρ greek numeral sign
→ 2032 ′ prime
- 02BA " MODIFIER LETTER DOUBLE PRIME
- exaggerated stress, contrastive stress
 - transliteration of tverdyj znak (Cyrillic hard sign: no palatalization)
- 0022 " quotation mark
→ 030B ˆˆ combining double acute accent
→ 2033 ″ double prime



(Spacing Modifiers)



Macedonian Translit.

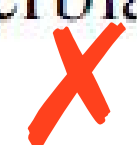
1E30 **ǲ** LATIN CAPITAL LETTER K WITH ACUTE
 ≡ 004B K 0301 6

1E31 **ǳ** LATIN SMALL LETTER K WITH ACUTE
 • Macedonian transliteration
 ≡ 006B k 0301 6



01F4 **Ǵ** LATIN CAPITAL LETTER G WITH ACUTE
 ≡ 0047 G 0301 6

01F5 **ǵ** LATIN SMALL LETTER G WITH ACUTE
 • Macedonian and Serbian transliteration
 ≡ 0067 g 0301 6



(Latin Extd-A, B)



Cassubian

(Didn't have a single CodePage)

Kaszëbsczé abecadlo:

aA ąĄ ãÃ bB cC dD eE éÉ ëË
fF gG hH iI jJ kK lL łŁ mM
nN ńŃ oO òÒ óÓ ôÔ pP rR
sS tT uU ùÙ wW yY zZ źŻ



Nasal Vowels

Ąą Ęę ĭ ĭ Ųų – Polish orthogr.,
Lithuanian

Qq – Sami, Old Icelandic

Q̄q̄ – Old Icelandic



ǫ ǣ – not precomposed



3. Missing Pieces

Oversights, mistakes,
problems, missing characters,
encoding problems etc.



Štokavian Accents

	rising	falling
long	/	⤿
short	\	//

Vowels: a e i o u r



Štokavian Accents

á é í ó ú ř • Latin 1

à è ì ò ù _ • Latin 1

â ê î ô û ŕ • Extended B

ă ă ĩ ồ ầ ử ử • Extended B

ř missing from UC!



Štokavian Accents

random examples...

Sr̂bija
hr̂vatski

	rising	falling
long	/	˘
short	\	˝



Štokavian Accents

Accents in Serbian...

брѣнати, брѣнѣм *уст.* боронѣть.
брѣндуша *ж.*, брѣнѣушка *ж.* бот. шафрѣн.
брѣнка *ж.* рѣжа, рѣжистое воспалѣние.
брѣња *ж.* 1) бѣлое пятнѣ (на морде живот-
ного); 2) козѣ с бѣлым пятнѣм на головѣ; лѣ-
шадѣ с бѣлым пятнѣм на мѣрде.

...same as in Croatian!



Štokavian Accents

Accents in Serbian...

бр̑нати, бр̑н̑ам *уст.* боронѣть.
бр̑ндуша *ж.*, бр̑нѣушка *ж.* бот. шафр̑ан.
бр̑нка *ж.* р̑ожа, р̑ожистое воспалѣние.
бр̑ња *ж.* 1) б̑ѣлое пятн̑о (*на морде живот-*
ного); 2) коз̑а с б̑ѣлым пятн̑ом на голов̑ѣ; л̑о-
шадь с б̑ѣлым пятн̑ом на морде.

...same as in Croatian!

...but not available in UC



Štokavian Accents

	rising	falling
long	/	˘
short	\	˝

Latin: 23 of 24

á é í ó ú ř • Latin 1
à è ì ò ù _ • Latin 1
â ê î ô û ř • Extended B
à è ì ò ù ř • Extended B

Cyrillic: nothing so far



Croatian Digraphs

matching Serbian Cyrillic

uhorn	Upsilonafr	Vhook	Yhook	yhook	Zstroke	zstroke	Ezh	Ezhrevers	ezhrevers	ezhtail	twostroke	Tonefive	tonefive	glottalinv	wynn
uʀ	Ϸ	U	Y	y	Z	z	З	Ʒ	Ʒ	Ʒ	Ʒ	5	5	ʈ	ƿ
clickdental	clicklateral	clickalveol	clickretrot	DZcaron	Dzcaron	dzcaron	LJ	Lj	Ij	NJ	Nj	nj	Acaron	acaron	Icaron
l	ll	≠	!	DŽ	Dž	dž	LJ	Lj	Ij	NJ	Nj	nj	Ǻ	ǻ	Ǽ
icaron	Ocaron	ocaron	Ucaron	ucaron	Udieresis	udieresis	Udieresis	udieresis	Udieresis	udieresis	Udieresis	udieresis	eturned	Adieresis	adieresis
ı	Ŏ	ŏ	Ŭ	ŭ	Ū	ū	Ů	ů	Ů	ů	Ů	ů	ə	Ä	ä
Adotmacr	adotmacr	AEmacron	aemacron	Gstroke	gstroke	Gcaron	gcaron	Kcaron	kcaron	Ogonek	ogonek	Ogonekma	ogonekma	Ezhcaron	ezhcaron
Ā	ā	Æ	æ	G	g	Č	č	Ǧ	ǧ	Q	q	Q̇	q̇	Ž	ž
jcaron	DZ	Dz	dz	Gacute	gacute	Hwair	Wynn	uni01F8	uni01F9	Aringacute	aringacute	AEacute	aeacute	Ostrokeacu	ostrokeacu
Ĳ	DZ	Dz	dz	Ć	ć	H	ƿ	Ŋ	ŋ	Å	å	Æ	æ	Ø	ø
Adblgrave	adblgrave	Ainvertedb	ainvertedb	Edblgrave	edblgrave	Einvertedb	einvertedb	Idblgrave	idblgrave	Iinvertedb	iinvertedb	Odblgrave	odblgrave	Oinvertedb	oinvertedb
À	à	Â	â	È	è	Ê	ê	Ì	ì	Î	î	Ò	ò	Ô	ô



(Latin–Extended B)



Cyrillic 'UK'

Ksicyrillie	ksicyrillie	Psicyrillie	psicyrillie
Џ	ѣ	Ѳ	ѳ
izhitsadblgrav	Ukeyrillie	ukeyrillie	Omegaroundcy
Ѹ	8	8	О
Koppacyrillie	koppacyrillie	thousandseyri	titlocyrillicon
Ҫ	ҥ	≠	~

Ksicyrillie	ksicyrillie	Psicyrillie	psicyrillie
Ќ	ќ	Ѳ	ѳ
izhitsadblgrav	Ukeyrillie	ukeyrillie	Omegaroundcy
Ѹ	OU	ou	О
Koppacyrillie	koppacyrillie	thousandseyri	titlocyrillicon
Ҫ	ҥ	≠	~

two variants:
ligature or digraph

(Cyrillic Hist.)



Cyrillic OU/ou

Ksicyrillie	ksicyrillie	Psicyrillie	psicyrillie
Џ	ѣ	Ѣ	ѣ
izhitsadblgrv	Ukeyrillie	ukcyrillie	Omegaroundoy
Ѣ	OU	ou	Ѣ
Koppacyrillie	koppacyrillie	thousandscyri	titlocyrillicon
Ѣ	ѣ	Ѣ	ѣ

0468	0469	046A	046B	046C	046D
Ѣ	ѣ	Ѣ	ѣ	Ѣ	ѣ
0476	0477	0478	0479	047A	047B
Ѣ	ѣ	OU	ou	Ѣ	ѣ
0484	0485	0486	0487	0488	0489
	ѣ	ѣ			

three forms needed

– just as for the Croatian digraphs!

UPPER CASE	MIXed	lower case
OY	Oy	oy



Cyrillic OU/ou











3 forms are needed;

but...

if only 2 are available,

they must be OU/ou,

not Ou/ou, because...

Ksicyrillie	ksicyrillie	Psicyrillie	psicyrillie
			
izhitsadblgrav	Ukcyrillie	ukcyrillie	Omegaroundcy
			
Koppacyrillie	koppacyrillie	thousandscyri	titlocyrillicon
			



Cyrillic OU/ou

3 forms are needed;

but...

if only 2 are available,

they must be OU/ou,

not Ou/ou, because...

Ksicyrillie	ksicyrillie	Psicyrillie	psicyrillie
Ѓ	ǣ	Ψ	ψ
izhitsadblgrav	Ukcyrillie	ukcyrillie	Omegaroundcy
ѐ	OU	ou	Ω
Koppacyrillie	koppacyrillie	thousandscyri	titlocyrillicon
Ғ	ғ	ҥ	Ү

this doesn't look too good

✗ ОѸГОДИТИ
оѸГОДИТИ



Cyrillic OU/ou

Mistake in Unicode 4.1 docs:

0478 Оу ~~✗~~ CYRILLIC CAPITAL LETTER UK
• basic Old Cyrillic uk is unified with
CYRILLIC LETTER U
→ 0423 У cyrillic capital letter u
0479 оу CYRILLIC SMALL LETTER UK

ИѦ	Оу	☉
0468	0478	0488
ИѦ	оу	☉
0469	0479	0489

Ksieyrillie	ksieyrillie	Psieyrillie	psieyrillie
Ѣ	Ѣ	Ѳ	Ѳ
izhitsadblgrav	Ukeyrillie	ukeyrillie	Omegaroundey
Ѣ	ОУ	оу	Ѡ
Koppacyrillie	koppacyrillie	thousandseyri	titloeyrilliecn
Ѣ	Ѣ	Ѣ	Ѣ



Russian Phonetics

just added in UC 4.1:

ŷ	ɀ
1D6B	1D7B
ɀ	ɀ
1D6C	1D7C
ɀ	ɀ
1D6D	1D7D

- ɀ LATIN SMALL CAPITAL LETTER I WITH STROKE
- used with different meanings by Americanists and Oxford dictionaries
- ɀ LATIN SMALL LETTER IOTA WITH STROKE
- used by Russianists



Russian Phonetics

just added in UC 4.1:

ŷ	Ƨ
1D6B	1D7B
Ƨ	Ƨ
1D6C	1D7C
đ	Ƨ
1D6D	1D7D

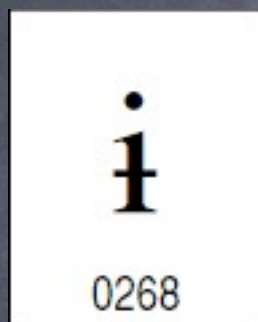
- Ƨ LATIN SMALL CAPITAL LETTER I WITH STROKE
- used with different meanings by Americanists and Oxford dictionaries
- Ƨ LATIN SMALL LETTER IOTA WITH STROKE
- used by Russianists

1) wrong
2) wrong

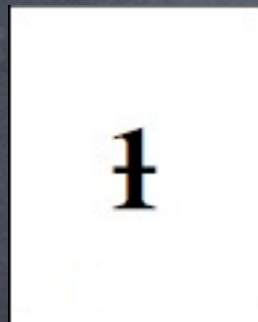


Russian Phonetics

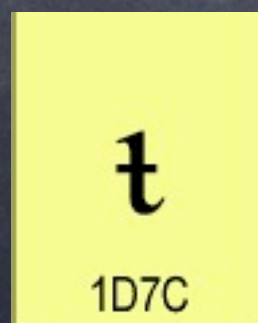
really needed for [ɨ]:



needed and available



needed but not available



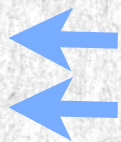
available but not needed



Russian Phonetics

1		2		3
nach API	nach Avanesov	nach API	nach Avanesov	
[ɪ]	[и ^е]	/e/ /a/ /o/ /e/ /a/ /o/ /i/ /e/ /a/ /o/ /e/ /a/ /o/	/e/ /a/ /o/ /e/ /a/ /o/ /и/ /е/ /а/ /о/ /е/ /а/ /о/	в лесу́ пяти́, часте́й село́ лесово́д, о ле́се пятиле́тка, у́часть мо́рем, перенесена́ би́ты шести́ дано́ гото́в шестикла́ссник вы́дан, кни́га, пу́ля на́ год, ле́том, мо́рем
[i]	[ь]			
[ɪ]	[ы]			
[ɪ]	[ы ^е]			
[ʌ]	[ʌ]			
[ə]	[ъ]			

dotless i, not iota!



K. Gabka (ed.), Einführung in das Studium der russischen Sprache.
Phonetik und Phonologie. Düsseldorf 1975, 34 & 164.



Russian Phonetics

СТИХОТВОРЕНИЕ М. Ю. ЛЕРМОНТОВА „ПАРУС“

па́рус

б'и^ел'ѐ^іѣт¹ па́рус / лд'инóкѣ² /
ф-тума́·н'ь мо́р'ѣ / гѣлубóм //
штó ѡш':ѣт-он / ф-стран'ѐ длл'·óкѣ^і /
штó к'ѡнул-он / ф-кря́·у раднóм //

игра́·і·ут во́лны / в'ѣт'ѣр с'в'ѡш':ѣт /
и-ма́·ч'ѣѣ / гн'·óѣѣ и-скрып'ѣѣ³ /
увѣ / он-ш':ѣс'т'иѣ⁴ н'ѣ-ѡш':ѣт /
и-н'ѣ-лѣ-ш':ѣс'т'иѣ⁴ / б'и^ежѣѣ /

па́д-н'ѣм стру́·ѣ / с'в'и^етл'ѐ^і лзѣ·р'и /
на́д-н'ѣм / лу́·ч'-со́нѣѣ зѣлѣтó·і /
л-óн м'и^ет'ѐжнѣ⁵ / прó·с'ѣт бѣ·р'и /
клѣ-бѣѣѣ⁶ в-бѣ·р'ѣѣ ѣѣс'т' па́кó·і ///

Avanesov
1972

Cyrillic +
Latin +
Greek +
Diacritics



Russian Phonetics

Cyrillic	абвгде...	✓
Latin/IPA	h i̯ j ə ʌ ɤ	✗
Greek	γ α	✓
Diacritics	ˆ ˊ ˋ ˙ ˚ ...	✗

Р.И. Аванесов, Русское литературное произношение, Москва 1972.



Russian Phonetics

Additional character 'Nasal': ∞

Самостоятельным, фонематическим признаком является „носовость вообще“ без локализации его образования. В словофонематической транскрипции носовую согласную фонему слабости по признаку дентальности-лабиальности обозначим знаком ∞. Вместе с тем отметим, что эта носовая фонема, слабая по признаку дентальности-лабиальности, одновременно является слабой и по признаку твердости-мягкости (так как перед [н'], [м'] она звучит мягко, а перед [н], [м] — твердо). Поэтому знак ∞ сопровождается цифрой 1 (справа ниже буквы), указывающей на неразличение этой слабой фонемы также по признаку твердости-мягкости. Исходя из сказанного приведенные выше слова в словофонематической транскрипции могут быть записаны следующим образом: |в'и∞₁т₂|, |ба∞₁т₂|, |ба́∞₁д₂α|, |кама́∞₁д₂α|, |в'й∞₁т'-ик₂|, |ба́∞₁т'ик₂|; |ра́∞₁п₂α|, |ла́∞₁п₂α|, |α∞₁ба́р|, |бо́∞₁ба́|, |баα∞₁б'йт'₂|.

Р.И. Аванесов, Русская литературная и диалектная фонетика, Москва 1974, 50.



Bulgarian Phonetics

Reviving OCS characters:

- # За означаване на дълга пауза в края на изречението:
и дъждове #.
- Ѹ За означаване на африкат дж: [Ѹоп], [ѸуѸè], [Ѹас].
- s За означаване на африкат дз: [sънкам], [сифт].
- l За означаване на изговора на средно [л] пред гласна
[lèсно].

Граматика на съвремения
български книжовен език. Том
I. Фонетика. София 1982, 29

s = dz ✓

Ѹ = dž ✗



Polish Orthography

breue longum groffum molle
a aa b

c has quinque differencias in Polonico idiomate habet
c k g cz ch

per se molle breue longum durum molle
d dz e ee ff f

per se inproprie breue longum durum molle
g g⁸⁷ i y l l

groffum molle groffum molle breue longum
m m n n o oo

durum molle per se per se
p p q R r

S in lingua Polonorum istas in se sex differencias habet
f s | sz ch z z
 breue longum

t u uu cum perdit vim vocalitatis,
 has tres habet differencias, que difference patent in abecedario Polonorum, scilicet adaam bił etc.
v v w x y ||

Jakub Parkosz, Traktat o ortografii polskiej.
 Warszawa 1985, 78. (Original: ca. 1440)



Polish Orthography

a á b **B** c ċ cz d dz dż e é e f g h i ü j ů y yi Ł
ł n m n ñ o ó p p r **r** rz s ss sz t v u w x.
X X

a á b **B** c ċ cz d dz dż e é e f g h i ü j ů y yi Ł
ł m m̄ n ñ o ó p p r **r** rz s ss sz t u u w w x.

St. Murzynowski, Orthographia Polska.
Königsberg 1551.



Polish Orthography

“r rotunda” – distinctive!

Um doppeldeutige Buchstabenverbindungen zu vermeiden, hat Murzynowski das Zeichen **7** eingeführt, da er die Verbindung $rz = r$ von $r+z/r+ż$ unterscheiden wollte in Beispielen wie *mierzy/mier-zi*, *skarżą/skarżą*. Der Buchstabe **7** wurde sehr häufig auch vom Schreiber des *Puławer Psalters* verwendet.

St. Murzynowski, Orthographia Polska.
Königsberg 1551.



Sorbian Orthography

Some Sorbian Characters:

ń p r ě ě ž ł m w



Ǻ ǻ ǣ ǿ Ȣ ȣ



(Latin Extd.–A, Extd. Additional)



Sorbian Orthography

18

TABELLARISCHE UEBERSICHT DER NS. ALPHABETE.

Fabricius 1706	Fryco 1796	Zwahr 1847	Tešnař- Šwjela	Časopis M. S.	Mein Alphabet	Ober- sorbisch	Alt- slovenisch
p	p	p	p	<i>p</i>	<i>p</i>	<i>p</i>	p
r	r >	r	r	<i>r</i>	<i>r</i>	<i>r</i>	r
—	sch — fch	—	—	<i>ṛ</i>	—	<i>ṛ</i>	—
ff (ff)	ff	ss	ff	<i>s</i>	<i>s</i>	<i>s</i>	s
fch	ffch	sch	fch	<i>ṣ</i>	<i>ṣ</i>	<i>ṣ</i>	š
fch	fch	schj	fch	<i>ṡ</i>	<i>ṡ</i>	—	—
t	t	t	t	<i>t</i>	<i>t</i>	<i>t</i>	t

K.E. Mucke, Historische und vergleichende Laut- und Formenlehre der niedersorbischen (niederlausitzisch-wendischen) Sprache. Leipzig 1891 (Reprint 1965,18)



Sorbian Orthography

Typographic challenge: Design



S with stroke; r rotunda



Czech Orthography

Latin, Medieval German etc.

$v = u = [u]$



Czech Orthography

Until 1849:
initial [u] = v; [v] = w
Needed:

ů

údolj “valley”

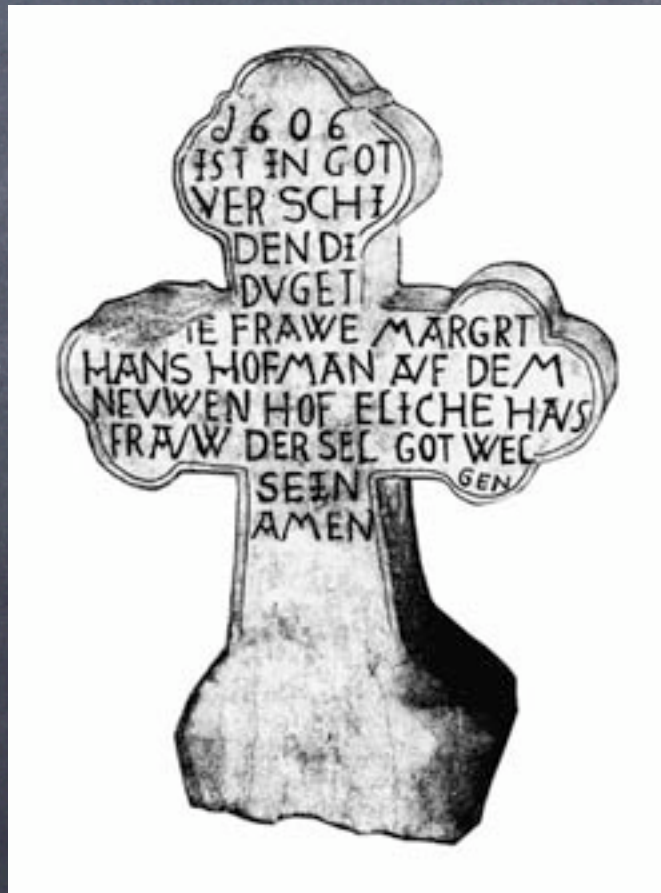
x



Czech Orthography

Needed for Czech: **ǚ** ‘ǚdolj’

Needed for German: **ǘ** ‘fǘhret’



ǚ

ǘ

Historical Cyrillic

The Basic Principle:

“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” (UC book, ch. 7)



Historical Cyrillic

The Basic Principle:

“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” (UC book, ch. 7)

– creates as many problems as it solves!



Historical Cyrillic

A Glimmer of Hope:

"If, at some future date, the old letterforms are adequately documented and the need for them demonstrated, then they can be added to this [the Cyrillic-Extended] block" (Unicode book v. 1, vol. 1, 45)



OCS, Old Russian

Application more or less trivial...

Ѧ ≈ А

Ѣ ≈ Б

Ѥ ≈ Щ

etc.



OCS, Old Russian

functional identity

Л ≈ Љ

Н ≈ Њ

Г̣ ≈ Ѓ

К̣ ≈ Ќ



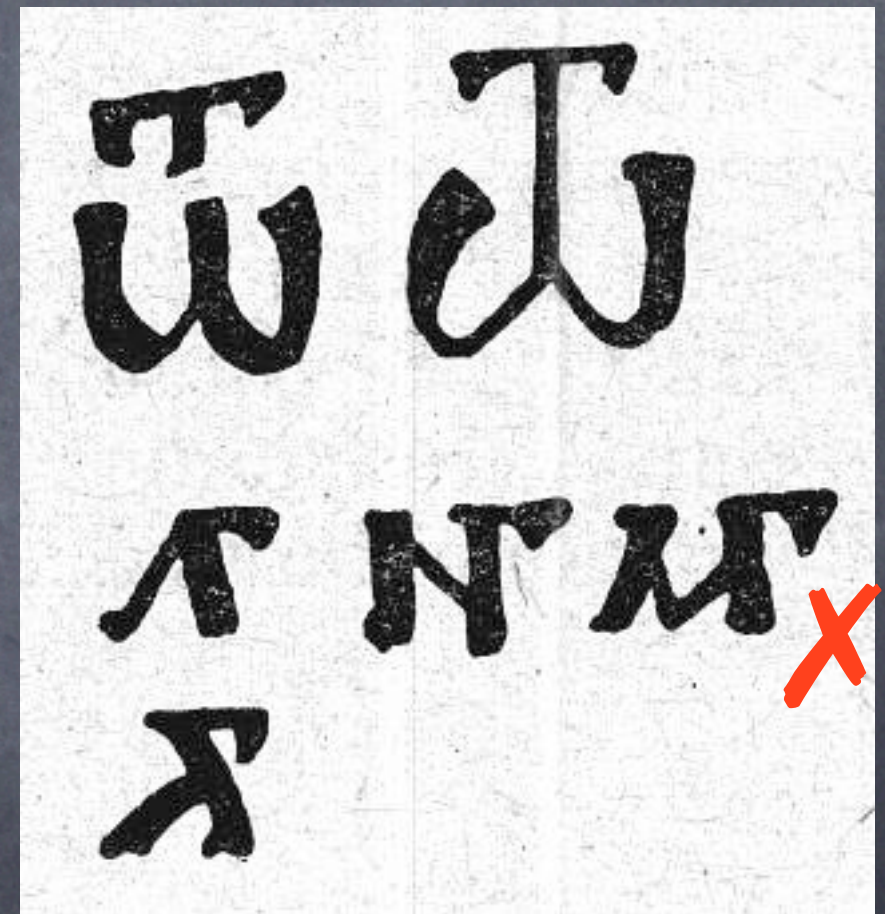
OCS, Old Russian

but then we find...



Мстиславово Ев. нач. XII в.

А.И. Соболевский, Славяно-русская
палеография. Изд. 2-е, СПб. 1908, 52.



В.Н. Щепкин, Русская палеография,
Москва 1967, 111.



OCS, Old Russian

functional identity no correspondence

Л ≈ Љ

Λ ≈ ?

Н ≈ Њ

Д ≈ Ђ

Ṗ ≈ ?

Ṛ ≈ Ṙ

Ḡ ≈ ?

Ṛ ≈ Ṙ

Ḳ ≈ ?



OCS, Old Russian

best-known problem...

Ѧ ≈ Я ?

Ѧ ≈ А ?



- Unicode does not have Ѧ
- being added to 'private area'



Early Modern Russian

	[jo]
18th century (AG-1802)	îô
Karamzin 1797	ë

x

✓

variation: îô, jô, ѣô, îô...



Early Modern Russian

	[jo]	[ja]
'Old Form'	Ѳѳ	Ѵѵ
'Modern Glyph'	ë	я

x

Unicode: status quo not satisfying;
too important for private area!



Old Russian

Jotated Jat'

Ѣѣ > ѣ ✕

Separate glyphs or variants of jat?



Old Russian

Jotated Jat'

Ѣѣ > ѣ ✕

Separate glyphs or variants of jat?

Ѧ-ѧ, Ѣ-ѣ, Ѥ-Ѭ, Ѧ-Ѩ



Old Russian

Not encoded yet:

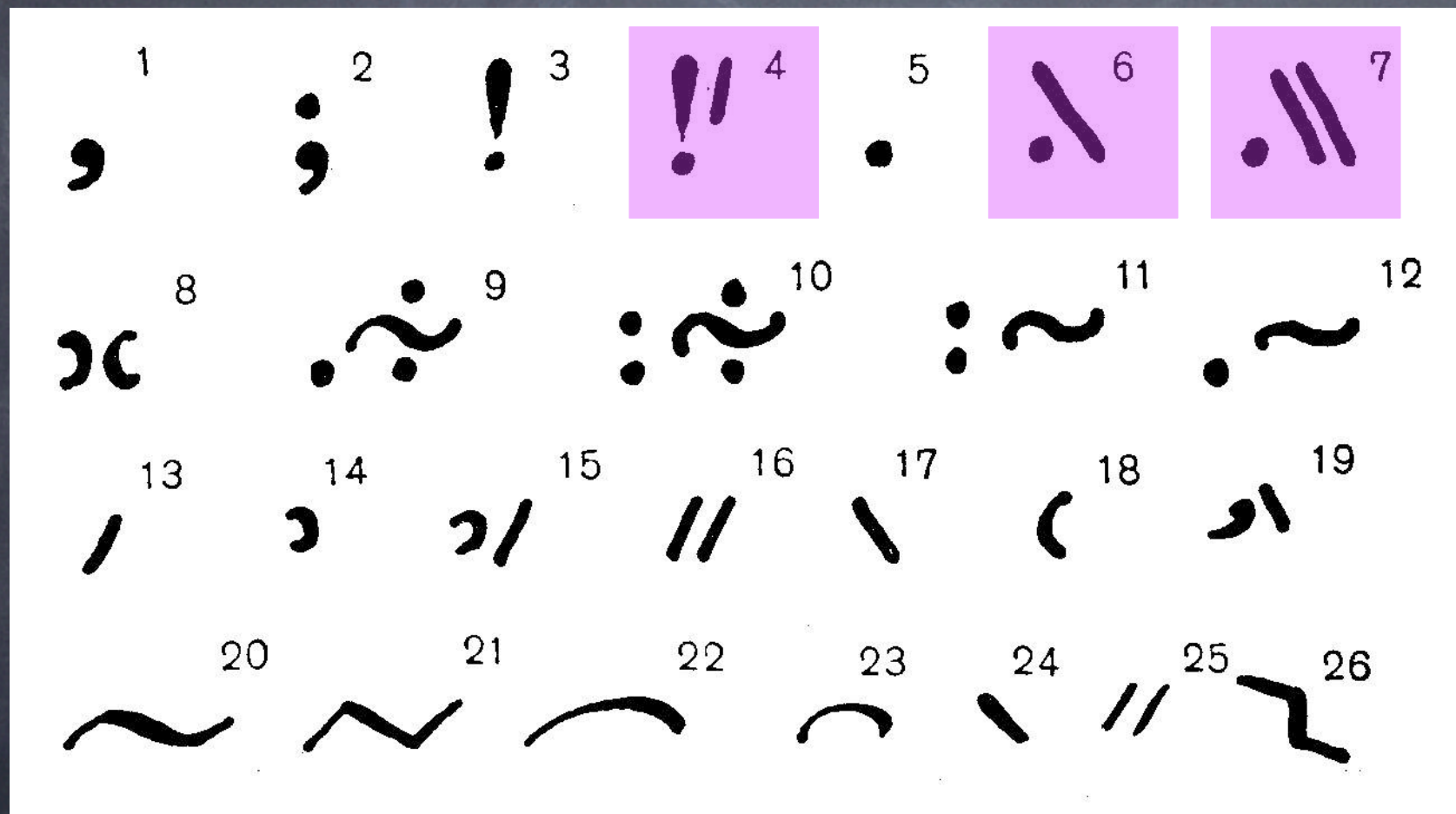
Ѧ Ѧ Ѧ Ѧ

Ѧ Ѧ Ѧ Ѧ Ѧ



Old Russian

znaki prepinanija 1



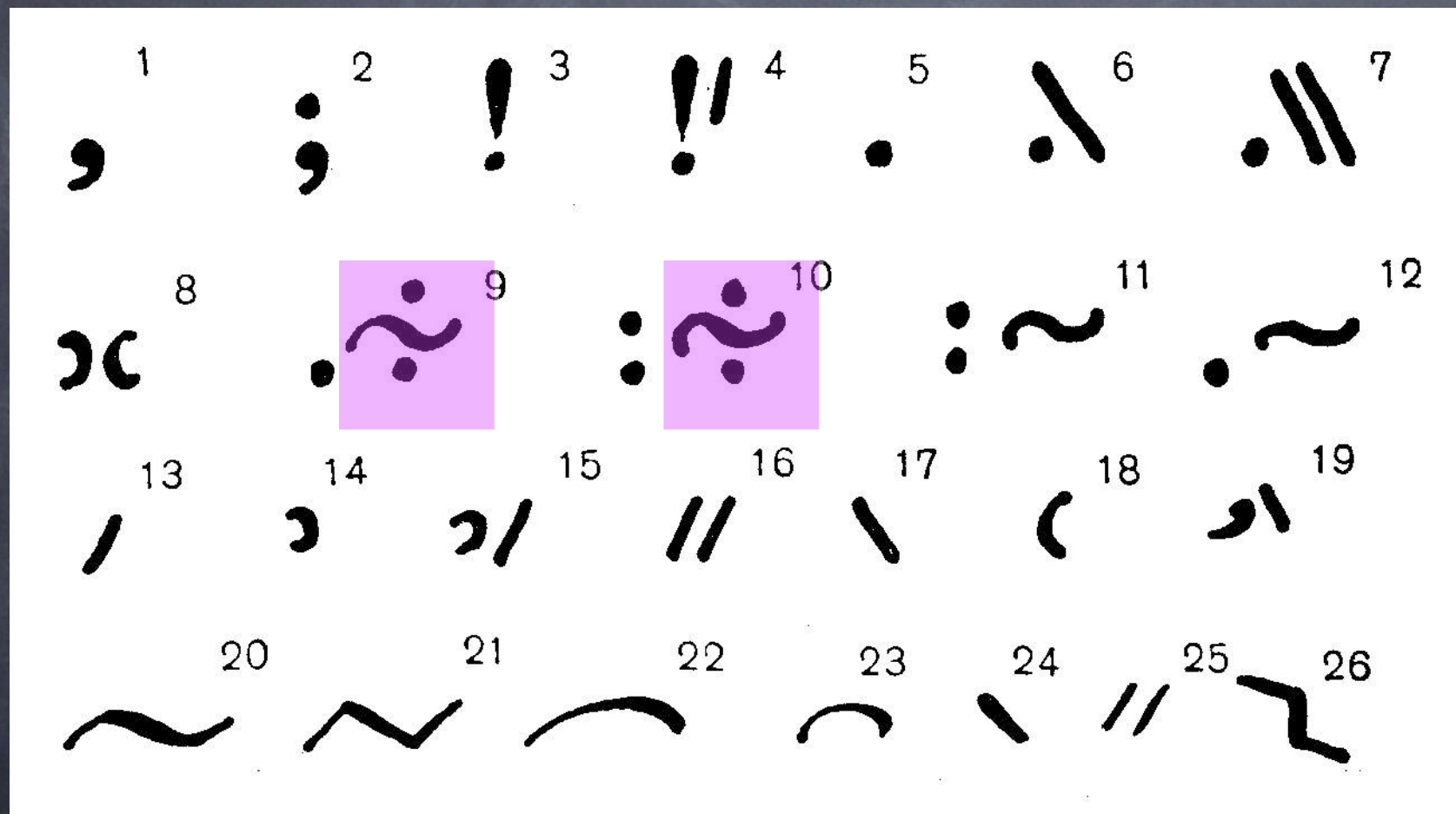
- 1 zapjataja
- 2 podstolija
- 3 pribyl'ca
- 4 pribavlenaja
- 5 točka
- 6 položitel'naja
- 7 kendema

Л.В. Черепнин, Русская палеография,
Москва 1956, 374–376.



Old Russian

znaki prepinanija 2



8 kavyka

9 slogija

10 statija

11 stišica

12 otrikal'

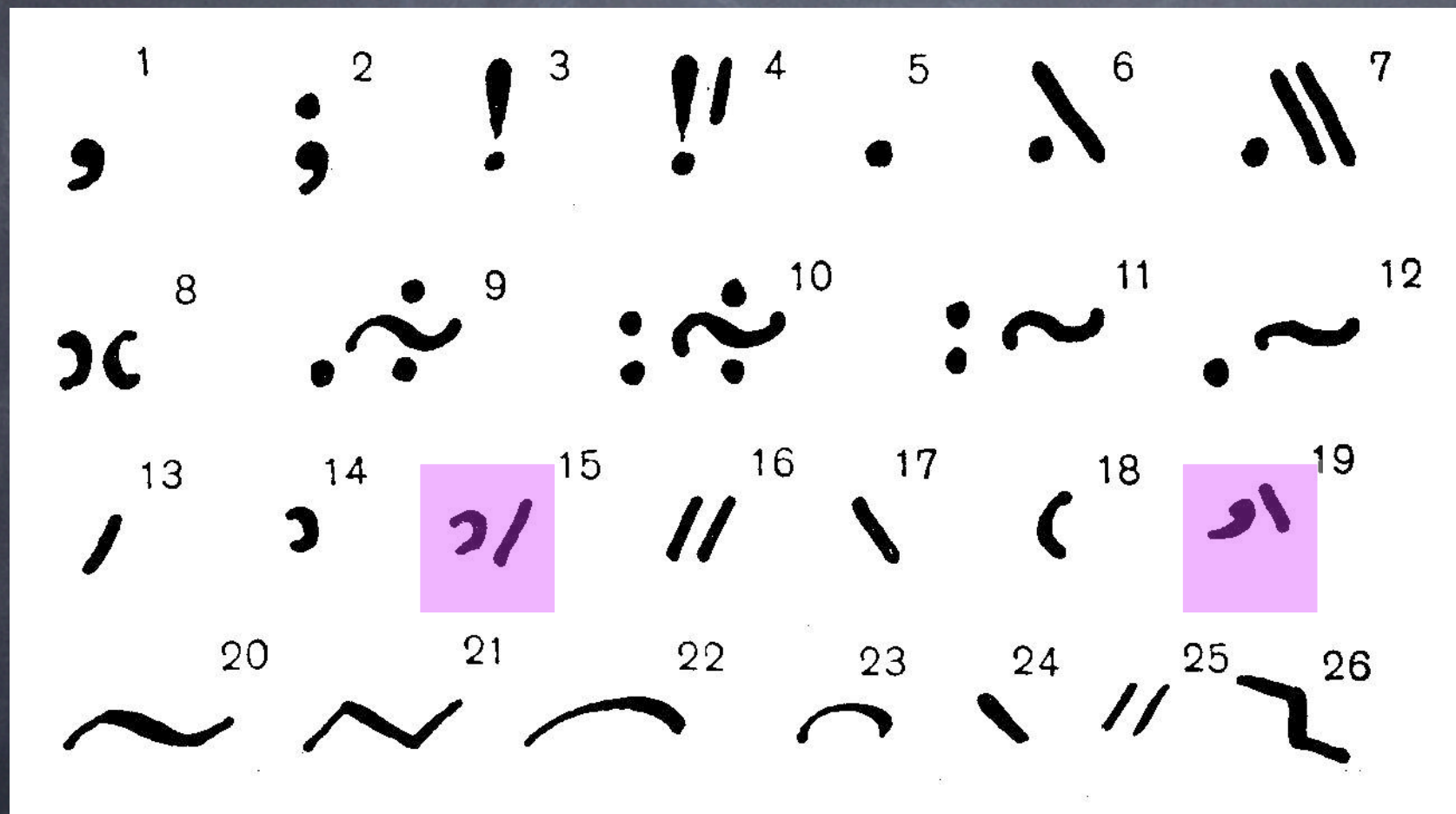
X

Л.В. Черепнин, Русская палеография,
Москва 1956, 374–376.



Old Russian

udarenija i pridyxanija



13 oksija

14 zvatelco/psili

15 iso/isso

16 okovavy

17 varija

18 dasija

19 apostrof

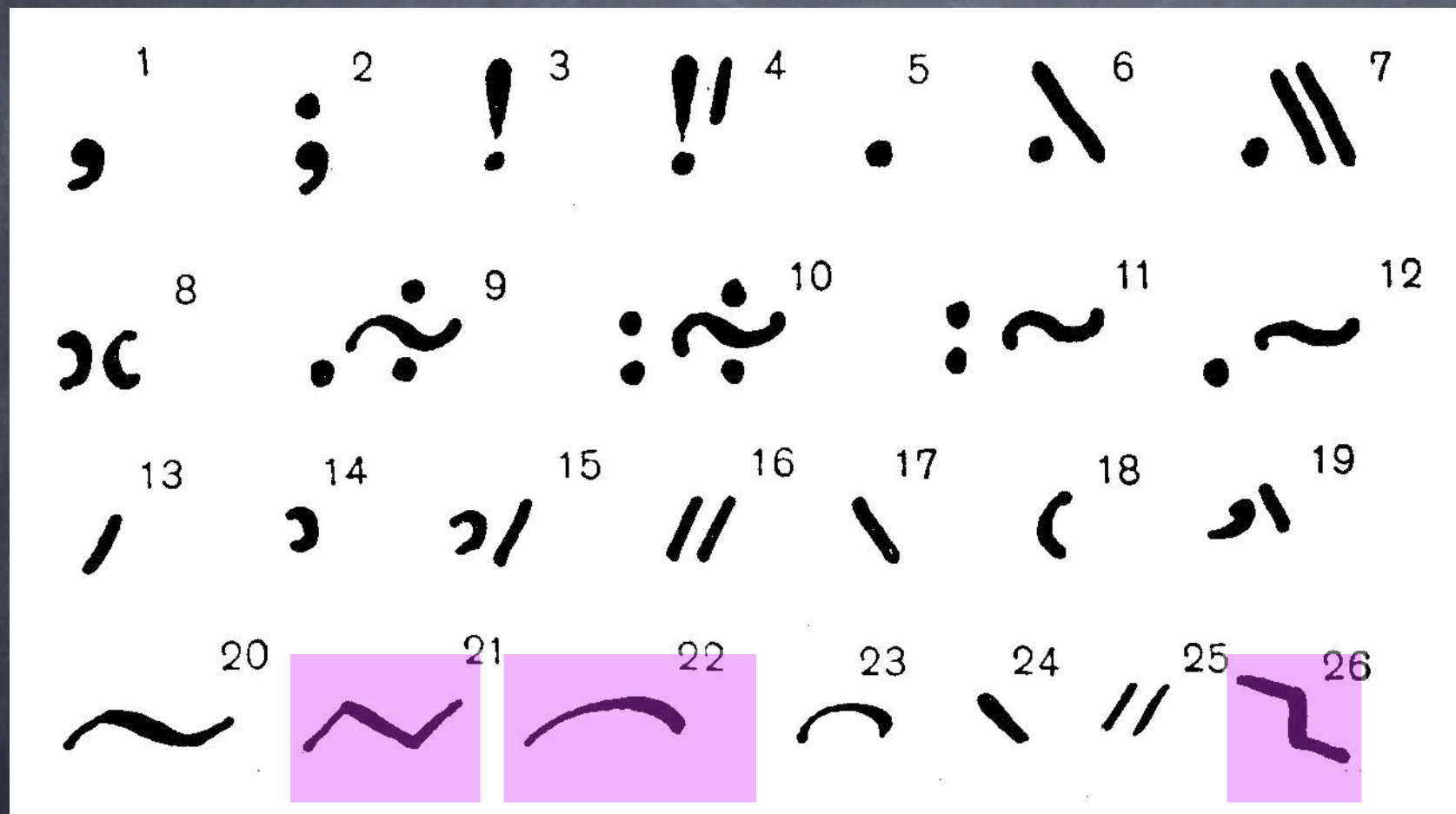
X

Л.В. Черепнин, Русская палеография,
Москва 1956, 374–376.



Old Russian

Titlo & Co.



20 titlo
21 vzmet
22 pokrytie
23 kamora
24 stjaga
25 smyček
26 erok, ertica X

Л.В. Черепнин, Русская палеография,
Москва 1956, 374–376.



Old Russian

problem: paerok

1. знак ' или \frown : мор'е, к'итъ, въплѣ.
2. знак ' или $\acute{}$: м'ного (мѣного), д'ва (дѣва),
 \times в'се (вѣсе), ч'то (чѣто).
3. титло \sim , \frown , $\bar{}$, \sim : снѣ (сынѣ), вѣка (владѣка),
ісѣ (исусѣ), глѣ (глава),
мѣцѣ (мѣсѣцѣ).



Old Russian

problem: paerok

2. знак ' или ' : М'НОГО (МЪНОГО), Д'ВА (ДЪВА),
X В'СЕ (ВЪСЕ), Ч'ТО (ЧЪТО).

- has two basic forms
- stands between characters!
- is **not** a diacritic!



Old Russian

problem: paerok

2. знак ' или ´ : М'НОГО (МЪНОГО), Д'ВА (ДЪВА),
X В'СЕ (ВЪСЕ), Ч'ТО (ЧЪТО).

- has two basic forms
- stands between characters!
- is **not** a diacritic!

≠ UC+033E "comb. vertical tilde": ´



Old Russian

part of the solution

02BC

,

MODIFIER LETTER APOSTROPHE

= apostrophe

,

- glottal stop, glottalization, ejective
- spacing clone of Greek smooth breathing mark
- many languages use this as a letter of their alphabets
- 2019 ' is the preferred character for a punctuation apostrophe
 - 0027 ' apostrophe
 - 0313 ́ combining comma above
 - 0315 ʹ combining comma above right
 - 055A ʻ armenian apostrophe
 - 2019 ' right single quotation mark

– letter!

– Mac./

Russ.

– one of
two forms
of paerok

‘apostrophe’ – Leskien






Old Russian

current Unicode solution?

	[jo] 0451	[ja] 044F	apostr. 02BC
'Old Form'	Ѡ	ѡ	Ѣ
'Modern Glyph'	ë	я	,



Old Church Slavonic

2150	2151	2152	2153	2154	2155	2156	2157	2158	2159	215A	215B	215C	215D	215E	215F
			$\frac{1}{3}$	$\frac{2}{3}$	$\frac{1}{5}$	$\frac{2}{5}$	$\frac{3}{5}$	$\frac{4}{5}$	$\frac{1}{6}$	$\frac{5}{6}$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{5}{8}$	$\frac{7}{8}$	$\frac{1}{1}$
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

Roman Numerals: Everything present 



Old Church Slavonic

СЛОВ'ЯНСЬКІ ЦИФРИ									
·Ѧ·	·Ѣ·	·҃·	·Д·	·Е·	·Ѕ·	·З·	·И·	·Ѧ·	·І·
1	2	3	4	5	6	7	8	9	10
ѦІ	ѢІ	҃І	ДІ	ЕІ	ЅІ	ЗІ	ИІ	ѦІ	
11	12	13	14	15	16	17	18	19	
Ѣ	Ѣ	Ѣ	И	З	О	П	Ч		
20	30	40	50	60	70	80	90		
Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ
100	200	300	400	500	600	700	800	900	
Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	Ѣ	
1000	2000	3000	4000	5000	10000				

Roman Numerals: Everything present X

Slavic Numerals: parts only X



Old Church Slavonic

Numbers, Numbers

10 T t'ma	100 T legion	1 Mio leodr	10 Mio voron	100 Mio koloda	1000 Mio t'ma tem
					

А.Х. Востоков, Грамматика церковно-
словенскаго языка изложенная по дрейвнейшимъ
онаго письменнымъ памятникамъ. СПб. 1863, 9
(Reprint Köln 1980)



Old Church Slavonic

encoded so far...

10 T t'ma	100 T legion	1 Mio leodr	10 Mio voron	100 Mio koloda	1000 Mio t'ma tem
(X)	X	X			



Glagolitic Translit.

Ѧ Ѧ Ѧ



И Ѧ Ѧ Ѧ



Glagolitic Translit.

Ѧ Ѧ Ѧ



И ѱ і і



Glagolitic Translit.

Samples: Cyrillic Iota

жаше его прѣдати. аште не би самъ хотѣлъ.
X ѿ того не можааше зърѣти. 21 ѿ гла-
дати. ꙗко свѣтильникомъ сжштемъ.
X ѿ свѣштамъ толнкамъ. се бо 22 ѿ при-
ѣ еванѣлистъ рече. ꙗко свѣтильники
свѣштъ ношахъ. ѿ тако его 27 ѿзиде же исъ
сатъ. ѿ нюда стоѣше съ ними. тѣ рекы

И.В. Ягичъ, Глаголическое письмо. В:
Энциклопедія славянской филологіи, вып. 3,
“Графика у славян”, СПб. 1911, 232.

Jagic, Edition of
Codex Marianus,
passim

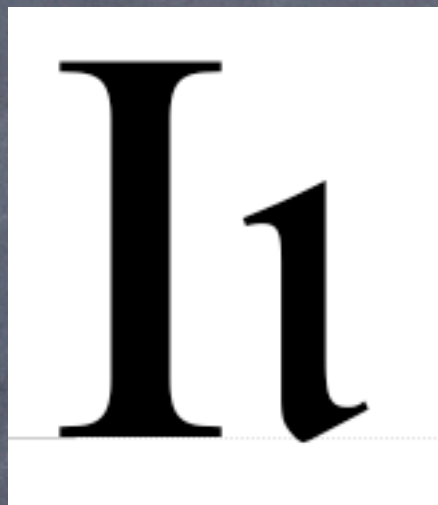


Glagolitic Translit.

Design: Cyrillic Iota



OCS
Black

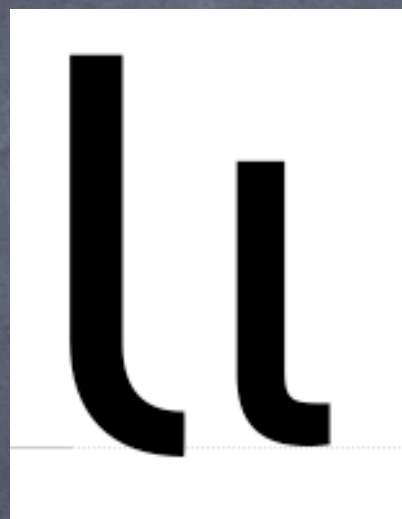


Greek
Serifed

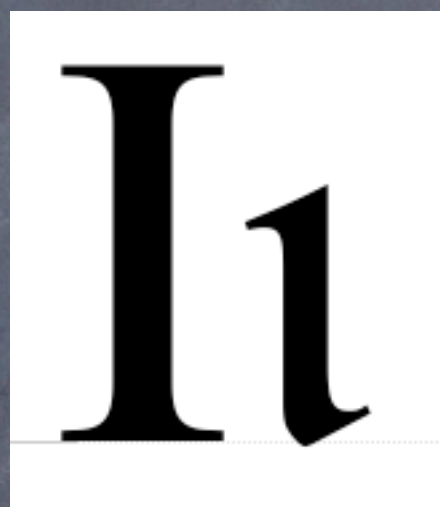


Glagolitic Translit.

Design: Cyrillic Iota



Cyrillic
Sans
Roman



Greek
Serifed
Roman



Cyrillic
Serifed
Roman



Glagolitic Translit.

Design: Cyrillic Iota

⋈ Ⲅ ⲡⲟⲛⲉ Ⲅⲉⲧⲏ ⲛⲉ Ⲅⲧⲏⲛ
⋈ ⲁⲡⲗⲧⲏ ⋈ ⲁⲗⲉ ⲓⲗⲁ ⲁ ⋈
Žalm 88. 6. Ⲓⲣⲟⲱⲕⲁⲗⲏ ⲛⲉⲥⲁ ⲕⲓⲟ ⋈

U1

X

ⲁ ⲓ ⲧⲏ ⲧ ⲁ Ⲅ
ⲛⲉ ⲛⲓ Ⲅⲉ ⲙ ⲙⲁ ⲓⲗⲁ ⲣⲡⲉ ⋈

Evangeliarum Assemani. Codex Assemani 3.
slavicus glagoliticus. Ediderunt Josef Vajs & Josef
Kurz. Tomus II. Edidit Josef Kurz. Pragae 1955, 3,
89



Glagolitic Translit.

Transliteration of Nasals

€ → A x



Glagolitic Translit.

Transliteration of Nasals

Ѣ → Ā X

Ѧ → A X



Glagolitic Translit.

Transliteration of Nasals

Ѣ → Ā x

Ѧ → A x

ѢѢ → X

ѦѢ → HX x



Glagolitic Translit.

Transliteration of Nasals

Ѣ → Ā X

Ѧ → A X

Ѣ → X

Ѧ → HX X

Ѧ → 'YO' X



Glagolitic Translit.

Ɱ	3
ⱮⱮ	ï
ⱮⱮ	za starije spomenike ġ, za mlade ĵ
Ɱ	ō
ⱮⱮ	ĉ ✓
ⱮⱮ	č
Ɱ	ju
ⱮⱮ	ž
Ɱ	6
Ɱ	ě, za vrije

Croatian
Norm

Ć, ĵ



Croatian Glagoljica

Ю ѡ ѡ, ѡ ✓ ѡ, ѡ ✗ ѡ, ѡ
Ь ѡ ѡ, ѡ; ' ✓ ѡ, ѡ, ' ѡ, ѡ, '

Different characters
or stylistic variants?



Bosančica

Ǧ	A	o	O
6	B	п	P
Π	V	р	R
л	G	с	S
А, Ѓ, Δ	D	т	T
Є	E	Ѹ	U
ƆЄ, Ж, ƆЖ	Ž	Ф	F
џ	DZ	х	H
З	Z	ω	OT
И	I	ψ	Ć, ŠT, ŠĆ
Ѓ, Ǧ	Đ, Ć (đerv)	ч	C
«	K	ѵ	Č
^	L	ш	Š
Ǧ^	LJ	б	poluglas
м	M	ѿ, Ъ	JAT
г	N	ю	JU
Ǧг	NJ		



Superscripts

	x	<	>	>^									
2	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
→													
*	***	***	***	***	***	***	***	***	***	***	***	***	***
	ˆ	ˆ	/	ˆ			ˆ	ˆ		ˆ	ˆ	ˆ	ˆ

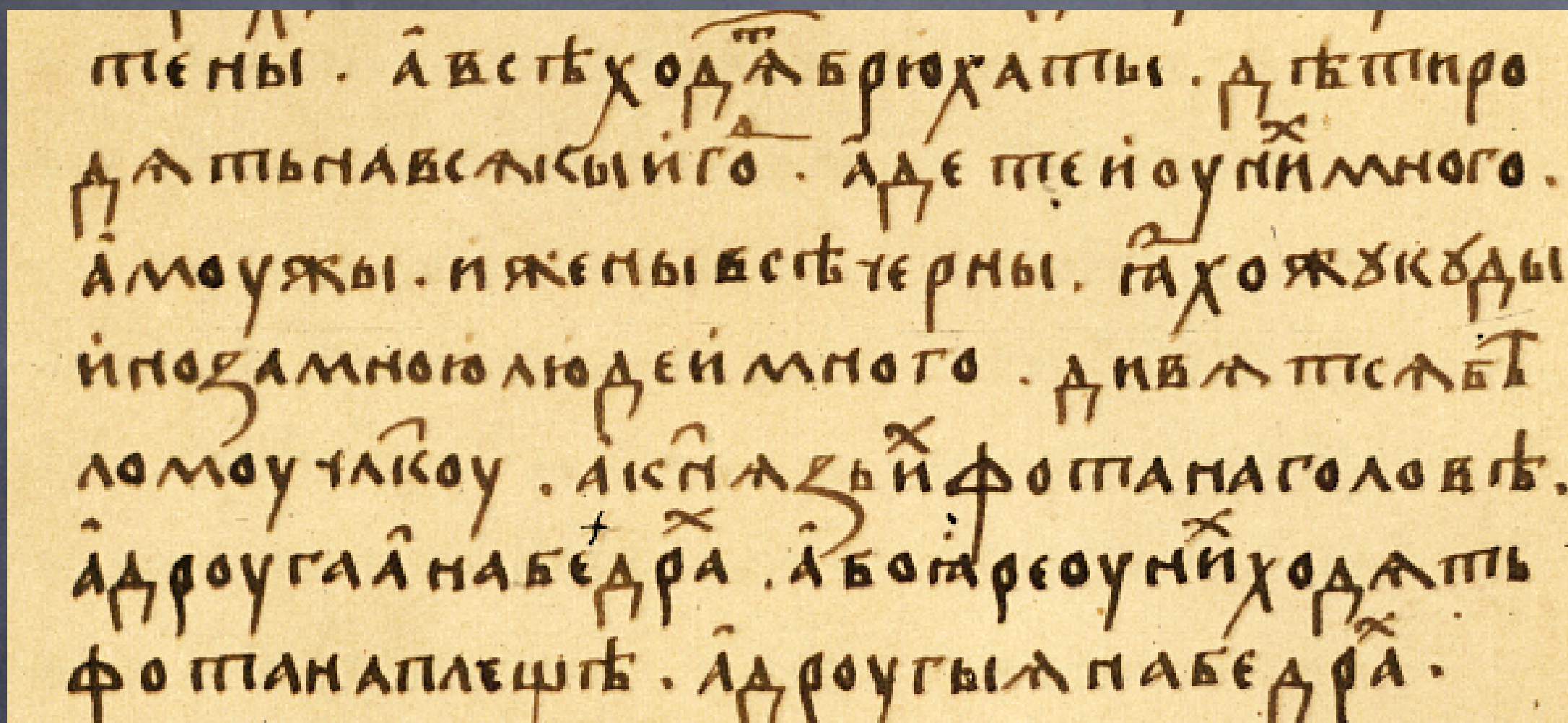
Latin: Many Superscripts (cmb)
for Medievalists (German!):
^eu (=ü) ^ea (=ä) ^eo (=ö); ^ov etc.



(Comb. Diacritical Marks)



Superscripts



Cyrillic: No Superscripts at all in UC



Афанасий Никитин, Хождение за три моря
Троицкий список



Ligatures

“Presentation forms”

- some for Latin
- many for Arabic

	FB0
0	ff FB00
1	fi FB01
2	fl FB02
3	ffi FB03
4	ffl FB04
5	ft FB05
6	st FB06



Ligatures

	FB0
0	ff FB00
1	fi FB01
2	fl FB02
3	ffi FB03
4	ffl FB04
5	ft FB05
6	st FB06

"Presentation forms"

- some for Latin
- many for Arabic
- none for Cyrillic
- none for Glagolitic

Ѣ Ѥ
Т т
Ѧ ѧ
Н н
Ш
Ѳ
Щ



Balkan Philology

- OCS-Cyrillic for Romanian

LITERA ♣

Litera **Л**, obținută prin modificarea grafică a literei **Ж**⁷⁸, 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 ЛН**⁸⁰ **ДОЛИНС ЧОРЖЧЕК** „din capul poienei, în jos, prin (pe în) valea Ciorîcei” (1499 iul. 13, Arh. St. Buc., S. I., nr. 118).

Deși atestată într-un cuvînt românesc, totuși unicul exemplu în care găsim această literă nu ne dă posibilitatea de a ne referi la valoarea sau valorile acestei litere în general ⁸¹. În acest unic caz, grafia ne determină să afirmăm că **Ț** notează nazalitatea vocalei care precedă pe **H** ⁸².

L. Djamo-Diaconiță: Limba documentelor
slavo-române emise în Țara Românească
în sec. XIV și XV. București 1971, 49



Balkan Philology

- Greek for Albanian

3. Griechische Lettern.

Γιάτι ἴνε **ἔ** γέ μπῆ κίε**λ** κιόφτε **σ**εντερούαρε ἔμερι ἴτ· ὄρτε μπρε-
τερία γιότε· οὐ **π**έφτε οὐρδερι ἴτ, σὶ κούντρο **π**ένετ ντε κίε**λ** ἀστοῦ **ε** δὲ
μπῆ δέ· ἔπνα νάθετ πού**κ**εν **ε** σόμμε **ἔ** νὰ δούχετε **π**έρ φῶστενε· **ε** δὲ

Transscription: *Yati inü tšü ye mbü kiel', kióftü süntüruarü ümüri it.*
Artü mbretüría yóte, U mbüftü urdüri it, si kundrü mbünetü dü kiel' aštu e de
mbü de. Epna náwet mbutšün e sorme tšü na duçetü per fistünü. E de ndüléna

Balkan Philology

- Greek for Albanian

3. Griechische Lettern.

Γιάτι ἴνε **ἔ** γέ μπῆ κίε^λ κιόφτε **ῥ**εντερούαρε ἔμερι ἴτ· ὄρτε μπρε-
τερία γιότε· οὐ **π**έφτε οὔροδερι ἴτ, σὶ χούντρο **π**ένετ ντῆ κίε^λ ἀστοῦ **ἔ** δὲ
μπῆ δέ· ἔπνα νάθετ πού^κεν **ἔ** σόμμε **ἔ** νὰ δούχετε **π**έρ φῶστενε· **ἔ** δὲ

Transscription: *Yati inü tšü ye mbü kiel', kióftü süntüruarü ümüri it.*
Artü mbretüría yóte, U mbüftü urdüri it, si kundrü mbünetü dü kiel' aštu e de
mbü de. Epna náwet mbutšün e sorme tšü na duçetü per fistünü. E de ndüléna

- Greek for Macedonian...

Balkan Philology

- Greek for Albanian

3. Griechische Lettern.

Γιάτι ἴνε **ἔ** γέ μπῆ κίε**λ** κιόφτε **ῥ**εντερούαρε ἔμερι ἴτ· ὄρτε μπρε-
τερία γιότε· οὐ **π**έφτε οὐρδερι ἴτ, σὶ κούντρε **π**ένεττε ντῆ κίε**λ** ἀστοῶ ἐ δὲ
μπῆ δέ· ἔπνα νάθετ πού**κ**εν ἐ σόμμε **ἔ** νὰ δούχεττε πέρ φῶστενε· ἐ δὲ

Transscription: *Yati inü tšü ye mbü kiel', kióftü süntüruarü ümüri it.*
Artü mbretüría yóte, U mbüftü urdüri it, si kundrü mbünetü dü kiel' aštu e de
mbü de. Epna náwet mbutšün e sorme tšü na duçetü per fistünü. E de ndüléna

- Greek for Macedonian...
- Arabic for Byelorussian & Bosnian...

4. Results



ll

Some Results...

llll

llll

lō ll ll

ll/ll/ll

llll

ll ll

ll ll ll

ll

ll
1D7C

ll

ll = ll

ll

ll
ll
ll

ll

ll ll ll ll

ll



Thank You!



N.B. The following two papers were the printed outcome of this presentation which was originally given at the corresponding conference:

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.

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.

Both papers are available online.



© Prof. Dr. Sebastian Kempgen 2021

ORCID: 0000-002-2534-9423

D-96045 Bamberg, University of Bamberg, Germany

sebastian.kempgen@uni-bamberg.de

<https://www.uni-bamberg.de/slaving/personal/prof-em-dr-sebastian-kempgen/>

