

# Unicode for Slavic Medievalists

David J. Birnbaum

University of Pittsburgh

djbpitt+@pitt.edu

Euro Summer School:  
Electronic Publishing for Cultural  
Heritage Studies

Sofia, 28 September 2002

# Outline

- Lumpers and Splitters
- Early Cyrillic Writing and Unicode
- Early Glagolitic Writing and Unicode
- Where to Get Information and How to Participate

# Lumpers and Splitters

- A few types of “e”: € € € € € € € €
- Handwritten letter forms *always* differ
- Approaches to classification
  - *Splitter*: If they’re at all different, they’re different, and should be encoded differently
    - Example: Retain difference in rendering
  - *Lumper*: If they’re at all similar, they’re the same, and should be encoded identically
    - Example: Conflate differences for querying, collation

# Early Cyrillic and Unicode

- “The historic form of the Cyrillic alphabet is treated as a font style variation of modern Cyrillic because the historic forms are relatively close to the modern appearance and because some of them are still in modern use ... .”
- Early Cyrillic “а” and modern Cyrillic “а” are both u+0430

# Source Set Rule

## Legacy Characters

<i>Character</i>	<i>Image</i>	<i>Decomposition</i>
u+0477	ṽ	u+0475 (v) u+030f (¨)
u+047f	Ṽ	u+0461 (w) u+0442 (̄) (superscript)
u+047d	Ṽ	u+0461(w) u+0483 (̄)

# Complementary Distribution (Non-Slavic)

<i>Language</i>	<i>Phoneme</i>	<i>Spelling (lower case only)</i>
Greek	/s/	ς / ___# (u+03c2)
		σ / elsewhere (u+03c3)
Hebrew	Final and nonfinal consonants	
Arabic	Initial, medial, final, and isolated consonants	

# Complementary Distribution (Slavic)

<i>Language</i>	<i>Phoneme</i>	<i>Spelling (lower case only)</i>
Rusian (some)	/ja/ ~ /҇a/ (~ /Cä/)	а / C __ ѧ / elsewhere
Rusian (some)	/o/	о / C __ ѡ / elsewhere
Pre-1918 Russian	/i/	і / __ V, ѣ (exc. мірѣ) и / elsewhere

# Visual Ambiguity

- Run-of-the-Mill Unicode Visual Ambiguity
  - u+002d HYPHEN-MINUS
  - u+2010 HYPHEN
  - u+2212 MINUS SIGN
- Super-Duper Early Cyrillic Visual Ambiguity
  - One character mapped to more than one glyph
  - One glyph associated with more than one character
  - u+0486 COMBINING CYRILLIC PSILI PNEUMATA (’)
  - u+0311 NON-SPACING INVERTED BREVE (ˆ)



# Jotation

<i>Character</i>	<i>Image</i>	<i>Notes</i>
u+0465	℥	
u+044f	ƒ	
u+0469	ƒ	
u+046d	ƒ	
u+044e	ƒ	Jotation plus u+043e (◊) (!)
(none)	ƒ	Jotation plus u+0463 (ƒ)
(none)	ƒ	Δ not present in Unicode Variant of u+0469 (ƒ)?

*/i/*

<i>Character</i>	<i>Image</i>	<i>Numerical</i>	<i>Notes</i>
u+0438	И	8	
u+0456	і	10	
u+0457	ї	10	
(none)	ı	10	Variant? Of what?
(none)	Ӏ	10 (?)	Invented for transcriptions of Glagolitic

# Cyrillic Palatal Glides

- U+0458 CYRILLIC SMALL LETTER JE  
(j)
  - Serbian, Macedonian
- U+0439 CYRILLIC SMALL LETTER  
SHORT I (й)
  - Russian, Ukrainian, Belarusian, Bulgarian

# Jers

<i>Character</i>	<i>Image</i>	<i>Notes</i>
u+044a	Რ	
u+044c	Ტ	
(none)	Უ	“Neutral” jer Variant? Of what?

# Jery

<i>Part</i>	<i>Images</i>
First	Ʒ, ʙ
Second	н, і, і̇, л, ʟ

- Non-ligated and ligated
- (Modern: u+044b ЫI)

# Other Special Problems

- Upper and Lower Case
  - Modern and early Cyrillic are the same script
  - Modern Cyrillic distinguishes case
  - Early Cyrillic typically does not distinguish case
- Ligation is productive
- Superscription may require different glyphs
  - “Recumbant r” (ѣ)

# /u/

<i>Character</i>	<i>Image</i>	<i>Note</i>
u+0443	ŷ	Modern image y
u+0475	ʋ	Variously /ü/ or /u/
u+0479	oŷ	Horizontal digraph
u+043e	o	Sequence of two characters
u+0443	ŷ	Second may alternatively be u+0475
(none)	8	Vertical digraph (ligated)

/o/

<i>Character</i>	<i>Image</i>	<i>Notes</i>
u+043e	o	Narrow
u+047b	○	Broad
(none)	⊙⊙∞	Ocular (also “polyocular”)

Cf. u+0461 ω (omega)

Cf. /e/



# Greek

<i>Greek</i>		<i>Cyrillic</i>	
<i>Character</i>	<i>Image</i>	<i>Character</i>	<i>Image</i>
u+03c8	ψ	u+0471	Ψ
u+03be	ξ	u+046f	Ξ
u+03b8	θ	u+0473	Θ
u+03b1	α	(none)	α
u+03b5	ε	(none)	ε

/e/

<i>Character</i>	<i>Image</i>	<i>Notes</i>
u+0435	e (modern)	CYRILLIC SMALL LETTER IE Most common modern <i>glyph</i>
u+0454	є (modern)	CYRILLIC SMALL LETTER UKRAINIAN IE Most common early <i>glyph</i>
u+044d	э (modern)	
u+0465	ѣ	
(none)	Є Ǝ Ɔ Ǝ ǎ ǎ	Variants? Of what?

# Palatal Consonants

<i>Character</i>	<i>Image</i>	<i>Notes</i>
u+0459	љ (modern)	CYRILLIC SMALL LETTER LJE
u+045a	њ (modern)	CYRILLIC SMALL LETTER NJE
u+04a5	ҥ	CYRILLIC SMALL LIGATURE EN GHE (Altay, Mari, Yakut)
(none)	Ӏ	Palatal /l/
u+04a5 (?)	ӆ	Palatal /n/
(none)	Ӈ	Palatal /d/
(none)	(not available)	Palatal /m/

# Old Church Slavonic and Russian

<i>Character</i>	<i>Image</i>	<i>Sound</i>	
		<i>Old Church Slavonic</i>	<i>Rusian</i>
u+044f	Ѧ	/e/ ~ /je/	/ä/
u+0469	Ѧ		
u+0467	Ѧ	/ja/	

# Front Nasal (Unicode)

<i>Character</i>	<i>Image</i>	<i>Notes</i>
u+0467	Ѧ	CYRILLIC SMALL LETTER LITTLE YUS
u+0469	ѧ	CYRILLIC SMALL LETTER IOTIFIED LITTLE YUS

# Front Nasal (Manuscripts) 1

<i>Manuscript</i>	<i>Nonjotated</i>	<i>Jotated</i>
Savvina Kniga	Ɱ, Ɱ (rarely)	Ɱ̣
Zograph Folia	Ɱ	Ɱ̣
Suprasliensis Šuck Psalter	Ɱ	Ɱ̣
Hilandar Folia	Ɱ̣	Ɱ
Ostromir Gospel	Ɱ̣	Ɱ̣̣
Preslav ceramic inscription	Ɱ	Ɱ̣̣

# Front Nasal (Manuscripts) 2

## Manuscripts with a Single Front Nasal Letter

<i>Manuscript</i>	<i>Image</i>	<i>Notes</i>
Undol'skij Folia	ʌ	
Cyrillic Macedonian Folium	Δ	ʌ twice (/C__)

ʌ may represent etymological front or back nasal

# Mixed Corpora: Geographic

Old Church Slavonic u+0467 (Ɑ) may correspond to:

<i>Character</i>	<i>Image</i>	<i>Recension</i>
u+044f	Ɑ	Rusian
u+0454	€	Serbian
u+046b	Ɱ	Middle Bulgarian

But: Not all Rusian Ɑ, Serbian €, and Middle Bulgarian Ɱ correspond to one another.



# Diachronic Paleography: Ѧ and ꙗ

<i>Image</i>	<i>Character</i>	<i>Period</i>
Ѧ	u+0467	Early
я	u+044f	Modern
ꙗ	(none)	Early

	<i>East Slavic</i>	<i>South Slavic</i>
Sound	я = ꙗ = Ѧ	я = ꙗ ( $\neq$ Ѧ)
Paleography	я < Ѧ	

# Glagolitic and Unicode

“The Unicode Standard regards Glagolitic as a *separate* script from Cyrillic, not as a font change from Cyrillic. This position is taken primarily because Glagolitic appears unrecognizably different from Cyrillic, and secondarily because Glagolitic has not grown to match the expansion of Cyrillic. The Glagolitic script is not currently supported by the Unicode Standard.”

# Proposed Glagolitic Inventory

- ተላላላብ ጸጋጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ  
ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ ጸጋጠጥ
- P2, simple back jus, stapić, square ot.
- Generic short and long titlo
  - Deprecate u+0483 COMBINING CYRILLIC TITLO

# Notes on Glagolitic 1

- Round and square glagolitic are different “typographic variants” (fonts)
- **ꙸꙸ** is two characters, not one
- Stapić
  - Historically a variant of ꙸ
  - Graphically distinct
  - Treated as distinct by those who work with square Glagolitic

# Notes on Glagolitic 2

- Round and square đerv
  - Historically the same letter
  - Different functions (đ and j)
  - Separate characters?
- Additional diacritics or punctuation?
- Upper and lower case?

# Commission

Special Commission to the Executive Council  
of the International Committee of Slavists for  
the Computer-Supported Processing of Slavic  
Manuscripts and Early Printed Books

# Commission Projects

- Unicode Early Cyrillic Proposal
- Unicode Glagolitic Proposal
- Agreement (private standardization) of part of Unicode PUA for Slavic medievalist community

# How to Participate

- **ОБЩЕЖИТІЄ Obštežitie Portal**
  - <http://www.ceu.hu/medstud/ralph/obsht.htm>
- **Mailing list for Early Slavic Written Sources**
  - [slav-mss-list@port.ac.uk](mailto:slav-mss-list@port.ac.uk)
  - Subscription information available at Obštežitie
- **Commission Web Site**
  - <http://clover.slavic.pitt.edu/~repertorium/commission/>  
(after 15 October 2002)
- **Webmasters:**
  - David J. Birnbaum: [djbpitt+@pitt.edu](mailto:djbpitt+@pitt.edu)
  - Ralph Cleminson: [ralph.cleminson@port.ac.uk](mailto:ralph.cleminson@port.ac.uk)