

Collating diplomatic transcriptions of manuscripts

David J. Birnbaum
Open philology workshop
Leipzig, 2014-07-14

djbpitt@gmail.com
<http://www.obdurodon.org>
<http://pvl.obdurodon.org>

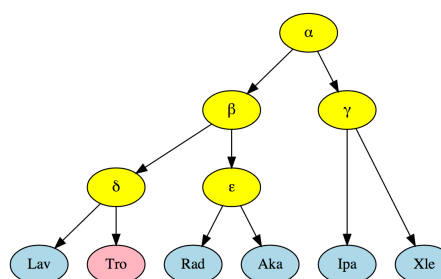
The Rus' primary chronicle

- *PVL* (*Повесть временных лет*)
- Historical chronicle of the East Slavs
 - Russia, Ukraine, Belarus
- Beginning from the creation (Hamartolos)
- Historical period (since 852) arranged by year
- Repeated edits and recompilations
- Fixed 1116
 - Incorporated into later chronicles

PVL textual tradition

- α Archetype
- β Northern branch
 - δ Laurentian (Lav, 1377) [Trinity (Tro; lost)]
 - ϵ Radziwiłł (Rad, 1490s), Academy (Aka, late 15th)
- γ Southern branch
 - Hypatian (Ipa, ca. 1425), Xlebnikov (Xle, 16th)
 - [Pogodin (Pog, early 17th, to supplement Xle)]
- Novgorod first
 - Commission (Kom), Academy (NAk), Tolstoj (Tol)

PVL stemma



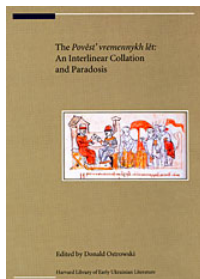
Why collate the PVL?

- Textual comparison
 - Relationships among the copies
 - Construction of a *paradosis* (alpha text)
 - History of transmission beyond alpha
- Why collate diplomatic transcriptions?
 - Linguistic comparison
 - Orthographic comparison
 - [More about diplomatic editions of manuscripts and critical editions of texts]

Practical issues

- No funding
 - Must reduce human effort
 - Especially repetitive human effort
- Edition is under constant development
 - Must be able to rerun collation
- Must be automated as much as possible

Print edition



- *The Povest' vremennykh let: An interlinear collation and paradosis*
- Donald Ostrowski, David J. Birnbaum, Horace G. Lunt
- Harvard UP, 2004
- 3 vv., 2368 pp.

Interlinear collation (print)

1,4:

Laur: персида, ватрь, тоже | н до индикниа в долготу
Trin: персида ватрь даже и до индикниа в долготу
Radz: персида, ватрь, доже н до индикниа, в долготу
Acad: персида, ватрь, дои н до индикниа, в долготу |
Hypa: персида, ватрь, доже н до индикниа, в долготу
Khle: персида, ватрь, даже и до индикниа, въ | долготу
Bych: Персида, Ватрь, доже и до Индикниа в долготу,
Shakh: Персида, Ватрь доже и до Индикниа въ долготу,
Likh: Персида, Ватрь, доже и до Индикниа в долготу,
 α : Персида, Ватрь доже и до Индикниа въ долготу,

Why interlinear?

- General
 - Variants are presented completely, not selectively
 - Ease of reading any individual copy
- Digital
 - Space, weight, cost are irrelevant
 - User can select witnesses
 - Searching on other than plain text
 - Lemma
 - Morphology

How interlinear?

- Alignment by line (per Karskii 1926 edition)
 - Ludolf Müller, *Handbuch zur Nestorchronik* (word index 1977)
 - Samuel Hazzard Cross, *The Russian primary chronicle Laurentian text* (English translation, 1930)
 - Performed manually for print edition
- Alignment by word
 - Too expensive to perform manually

Print edition workflow

- Typeset in troff
- Focus on producing print version
- Alignment is manual
 - Word-level alignment is impractical

Digital versions

- PDF of print edition
<http://hudce7.harvard.edu/~ostrowski/pvl/>
- HTML edition
<http://pvl.obdurodon.org>

First digital version

1, 4

Lav	персѣда, ватрь, тоже и до индикна в доаготу
Tro	персѣда ватрь, даже и до индикна в доаготу
Rad	персѣда, ватрь, доже и до индикна, в доаготу
Aka	персѣда, ватрь, до ^и и до индикна, в доаготу
Ipa	персѣда, ватрь, доже и до индикна, в доаготу
Xle	персѣда, ватрь, даже и до индикна, въ доаготу
Vuċ	Персѣда, Ватрь, доже и до Индикна в доаготу,
Šax	Персѣда, Ватрь доже и до Индикна въ доаготу,
Lix	Персѣда, Ватрь, доже и до Индикна в доаготу,
α	Персѣда, Ватрь доже и до Индикна въ доаготу,

First digital version

- Pro
 - Automated conversion from troff
 - Control over display
 - Fonts
 - Toggle individual witnesses on and off
 - Potential for annotation (lemma, morphology)
- Con
 - No support for word-level comparison

Why is collation difficult?

- Exponential complexity
 - Worst case: compare every word in every witness to every word in every other witness
 - Complicated by repetitions and transpositions
- Diplomatic transcription
 - Efficient comparison algorithms require exact string matching, which is rare in diplomatic transcription
 - Finding *closest* match requires a completely different (more computationally expensive) method than finding *exact* match

Word-aligned version

1,4

Lav	персѣда,	ватрь,	тоже	и до	индикна	в	доаготу
Tro	персѣда	ватрь,	даже	и до	индикна	в	доаготу
Rad	персѣда,	ватрь,	доже	и до	индикна,	в	доаготу
Aka	персѣда,	ватрь,	до ^и	и до	индикна,	в	доаготу
Ipa	персѣда,	ватрь,	доже	и до	индикна,	в	доаготу
Xle	персѣда,	ватрь,	даже	и до	индикна,	въ	доаготу
Vuċ	Персѣда,	Ватрь,	доже	и до	Индикна	в	доаготу,
Šax	Персѣда,	Ватрь	доже	и до	Индикна	въ	доаготу,
Lix	Персѣда,	Ватрь,	доже	и до	Индикна	в	доаготу,
α	Персѣда,	Ватрь	доже	и до	Индикна	въ	доаготу,

A more challenging passage

9,2

Lav	и то	творятъ	моуеньє	себѣ	а не	мученьє.
Tro	и то	творятъ	моуеньє	себѣ	а не	мученьє.
Rad	и тако	творятъ	не	мытву	себѣ	но м ^ѣ ученьє.
Aka	и тако	творятъ	не	мытву	себѣ	но мученьє.
Ipa	и	творя ^т	не	мытву	себѣ	а не мученьє.
Xle	и	творя ^т	не	мытву	себѣ	но мученьє.
Vuċ	и то	творятъ	моуеньє	себѣ,	а не	мученьє ^т .
Šax	и то	творятъ	моуеньє	себѣ,	а не	мученьє ^т .
Lix	и то	творятъ	моуеньє	себѣ,	а не	мученьє ^т .
α	и	творятъ	не	мытву	себѣ,	не мученьє ^т .

CollateX

- <http://collatex.net/>
- Interedition (Huygens Institute, the Hague)
- Advantage
 - Use someone else's collation algorithm and implementation
- Limitation
 - Requires exact string matching
 - Cannot find *closest* match
 - Cannot find logical matches that are not string matches
 - Digits vs words: 40000 ~ рд ~ 40 тысцѣ
 - Synonymy: разумъ – смыкълы-

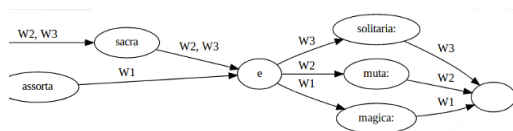
Recent CollateX developments

- Ported to Python (module)
- New collation algorithm
 - Non-progressive
 - Suffix arrays and LCP arrays
 - Not subject to order effects
- Remaining limitations
 - Exact string matching
 - Repetition
 - Transposition

Workflow (summary)

- Input is custom XML
- Convert to TEI, tokenize, add <w> tags
- Convert to JSON
- Preprocess: Enrich JSON with bespoke normalization
- Collate with CollateX, which creates variant graph
- Postprocess: Adjust rank in variant graph
- Generate JSON output
- Convert to custom XML
- Convert to HTML tables for rendering

Variant graph (excerpt)



Workflow (beginning)

- Input: Custom XML line blocks
- Convert to TEI with word (<w>) tags
 - Tokenizing mixed content
 - word1 wo<lb/>rd2 word3
 - <w>word1</w>
 - <w>wo<lb/>rd2</w>
 - <w>word3</w>
 - Not all whitespace represents token break
 - word1 <lb/> word2
 - <w>word <lb/></w>
 - <w>word2</w>
- Convert to JSON, enrich with normalizations

Normalization

- Create a normalized “shadow” copy
- Normalization based on Soundex, adapted for early Cyrillic writing
- Collate on normalization, return original

Soundex

- English-language surnames, 1918
- Algorithm (simplified)
 - Retain first letter
 - Delete other vowels; degeminate
 - Conflate other letters according to phonetic similarity (e.g., t/d = 3; m/n = 5)
 - Truncate or zero-pad to four characters
- Examples
 - Birnbaum B-651 (also ✓ Barenboim; also ✗ Brumble)

Soundex assumptions

- Character differences are not all equivalent with respect to information load
- Information load may be sensitive to position
- Beginning of word carries more information than end
 - Especially inflected languages
- Consonants carry more information than vowels
 - Except in short words

Adapting Soundex to Church Slavonic

- Neutralize variant spellings of initial vowel
 - оу, у, љ = у
 - ѡ, ѡ, ѡ, ѡ = ѡ
- Case fold, neutralize consonantal variants
 - Not always one-to-one, e.g., цѣ = шѣ
- Degeminate, delete other vowels, delete diacritics
 - Keep two letters of two-letter words
 - Higher information load
- Other conflations?
 - Knowledge based vs machine learning
- Expand abbreviations?
 - бѣа, бѣа, бѣа = бѣа (бѣ)
- Truncate or zero-pad (to what length?)

Soundex sample (*Bdinski sbornik*)

- Ch397 и възврѣ | титѣ дѣщерьше своею.
 - Ch384 и възвратиѣ дѣщершоу свою.
 - Nbkм298 и възвратити братаициѣ | своѣж
 - Berlin и възврати | братаициѣ свою.
-
- Ch397 и възвр дштр св
 - Ch384 и възвр дштр св
 - Nbkм298 и възвр брѣи св
 - Berlin и възвр брѣи св

Two types of normalization

- Collation
 - Find alignment points
 - Coarse adjustments
 - No harm in conflating grammatical forms
 - Imperfect and aorist; infinitive and supine
- Evaluation
 - Alignment points are already known
 - Finer comparisons
 - Many need to distinguish on the basis of small details

Collation after Soundex

- Greatly improved actual matches
- *Forced matches*
 - A B C
 - A D C
- Misses
 - Gap in alignment (no forced match)
 - Imperfect match
 - фраки ~ фраци (фрѣ ~ фрѣ)
 - CollateX recognizes only perfect matches
 - Unable to recognize *closest match*

3,5

3,5					
<i>Lav</i>	гарѣмати	таврѣ ани.	асуфѣа.		фраци.
<i>Tro</i>	гарѣмати	таврѣани	асуфѣа		фраци
<i>Rad</i>	сарѣмати.	таврѣани	асуфѣа.	и	фраци
<i>Aka</i>	сарѣмати.	таврѣани.	асуфѣа.	и	фраци
<i>Ipa</i>	сарѣмати.	таврѣани.	асуфѣа.		фраци.
<i>Xle</i>	сарѣмати.	таврѣани.	асуфѣа.		фраци.
<i>Vuč</i>	Сарѣмати,	Таврѣани,	Скуфѣа,		Фраци,
<i>Šax</i>	Сарѣмати,	Таврѣани,	Скуфѣа,		Фраци,
<i>Lix</i>	Сарѣмати,	Таврѣани,	Скуфѣа,		Фраци,
<i>α</i>	Сарѣмати,	Таврѣани,	Скуфѣа,		Фраци,

Numbers

18,4

<i>Lav</i>						[.ΔO]	ΔБДА
<i>Aka</i>	Ѓ.Ѓ.	Δ	исхо	женѡ	мнѡсѡвѡ	ΔO	ΔБДА
<i>Ipa</i>	Ѓ.Ѓ.	Δ	исхо	женѡ	мнѡсѡвѡ	ΔO	ΔБДА
<i>Xle</i>	Ѓ.Ѓ.	Δ	исхо	женѡ	мнѡсѡвѡ	ΔO	ΔБДА
<i>Буџ</i>	430;	Δ	отъ	исхоженѡ	монсѡвѡ	ΔO	ΔавѡΔΔ
<i>Šax</i>	430;	Δ	отъ	исхоженѡ	монсѡвѡ	ΔO	ΔавѡΔΔ
<i>Lix</i>	430;	Δ	отъ	исхоженѡ	монсѡвѡ	ΔO	ΔавѡΔΔ
<i>α</i>	430;	Δ	отъ	исхоженѡ	монсѡвѡ	ΔO	ΔавѡΔΔ

Problem areas

- Gaps in alignment
- No perfect match
- CollateX follows graph rank (leftmost match)
- 3,5
 - Orthography
 - скуѡфѡ и фѡрѡцѡ (Tro)
 - скуѡфѡ и фѡрѡцѡ (Rad)
 - Soundex
 - скф фрк
 - скф и фрц

Postprocessing

- Gap without perfect match on either side
 - Gap may span multiple columns
- Orthography
 - скуѡфѡ и фѡрѡцѡ (Tro)
 - скуѡфѡ и фѡрѡцѡ (Rad)
- Soundex
 - скф фрк
 - скф и фрц
- If there's a match, keep it
- Else
 - Find unique Soundex values in column and following
 - Move token to column with closest match
 - Damerau-Levenshtein edit distance
 - Insertion, deletion, substitution, transposition

9,2

9,2

<i>Lav</i>	и	то	творѡть	мѡвѡнѡ	совѡ	Δ	не	мѡченѡ.	
<i>Tro</i>	и	то	творѡть	мѡвѡнѡ	совѡ	Δ	не	мѡченѡ.	
<i>Rad</i>	и	такѡ	творѡть	не	мѡтѡу	совѡ	но	мѡченѡ.	
<i>Aka</i>	и	такѡ	творѡть	не	мѡтѡу	совѡ	но	мѡченѡ.	
<i>Ipa</i>	и		творѡть	не	мѡтѡу	совѡ	Δ	не	мѡченѡ.
<i>Xle</i>	и		творѡть	не	мѡтѡу	совѡ	но	мѡченѡ.	
<i>Буџ</i>	и	то	творѡть	мѡвѡнѡ	совѡ	Δ	не	мѡченѡ.	
<i>Šax</i>	и	то	творѡть	мѡвѡнѡ	совѡ	Δ	не	мѡченѡ.	
<i>Lix</i>	и	то	творѡть	мѡвѡнѡ	совѡ	Δ	не	мѡченѡ.	
<i>α</i>	и		творѡть	не	мѡтѡу	совѡ	не	мѡченѡ.	

In case of ties

- Thesaurus
- Most matches
- Length of match

Thesaurus

- Collect forced inexact matches
 - A B C
 - A D C
- Edit manually
- Use to break ties
- Close matches
 - поломѡшѡ – вѡзломѡшѡ
 - ламѡш – вѡзлѡмѡшѡ
- Non-matches
 - рѡзумѡнѡ – сѡмѡслѡнѡ
 - рѡзнѡ – сѡслѡнѡ

Thesaurus

220,9

<i>Lav</i>	налагѡша	первое	на	сѣопѡака	и	взаомнѡша
<i>Rad</i>	налагѡша	первие	на	сѣопѡака	и	пѡломнѡша
<i>Aka</i>	налагѡша	первое	на	сѣопѡака	и	пѡломнѡша
<i>Ipa</i>	налагѡша	первое	на	сѣопѡака	и	вѣзѡмнѡша
<i>Xle</i>	налагѡша	первое	на	сѣопѡака	и	вѣзѡмнѡша
<i>Vyē</i>	налагѡша	первое	на	свѣтопѡака	и	взаомнѡша
<i>Šax</i>	налагѡша	первое	на	свѣтопѡака	и	взаомнѡша
<i>Lix</i>	налагѡша	первое	на	свѣтопѡака	и	взаомнѡша
<i>α</i>	налагѡша	первое	на	свѣтопѡака, и	взаомнѡша	

взаомнѡша пѡломнѡша
вѣзѡмнѡша пѣрвѡша

What's next: many-to-one

141,11

<i>Lav</i>	а	прѡнѣ	вон	ѡ		и	пѡнде	на	сѣопѡака	нарекѣ
<i>Rad</i>	а	прѡнѣ		ѡ		и	пѡнде	на	сѣопѡака	нарекѣ
<i>Aka</i>	а	прѡнѣ		ѡ		и	пѡнде	на	сѣопѡака	нарекѣ
<i>Ipa</i>	а	прѡнѣ	вон	ѡ	тысяць	и	пѡнде	на	сѣопѡака	нарекѣ
<i>Xle</i>	а	прѡнѣ	вон	ѡ	тысяць	и	пѡнде	на	сѣопѡака	нарекѣ
<i>Vyē</i>	а	прѡнѣ	вой	40000		и	пѡнде	на	свѣтопѡака	нарекѣ
<i>Šax</i>	а	прѡнѣ	вон	40	тысяць	и	пѡнде	на	свѣтопѡака	нарекѣ
<i>Lix</i>	а	прѡнѣ	вой	40,000		и	пѡнде	на	свѣтопѡака	нарекѣ
<i>α</i>	а	прѡнѣ	вон	40	тысяць	и	пѡнде	на	свѣтопѡака, нарекѣ	

Acknowledgements

- Thanks to Ronald Dekker, lead developer of CollateX, for generous advice and consultation
- Thanks to Minas Abovyan, who implemented the PVL-specific Python code in our project

Thank you!

David J. Birnbaum
Open philology workshop
Leipzig, 2014-07-14

djbpitt@gmail.com
<http://www.obdurodon.org>