

Collating the
Rus' primary chronicle
(*Povest' vremennyx let*)

David J. Birnbaum
Varna, 2014-09-15
djbpitt@gmail.com
<http://www.obdurodon.org>

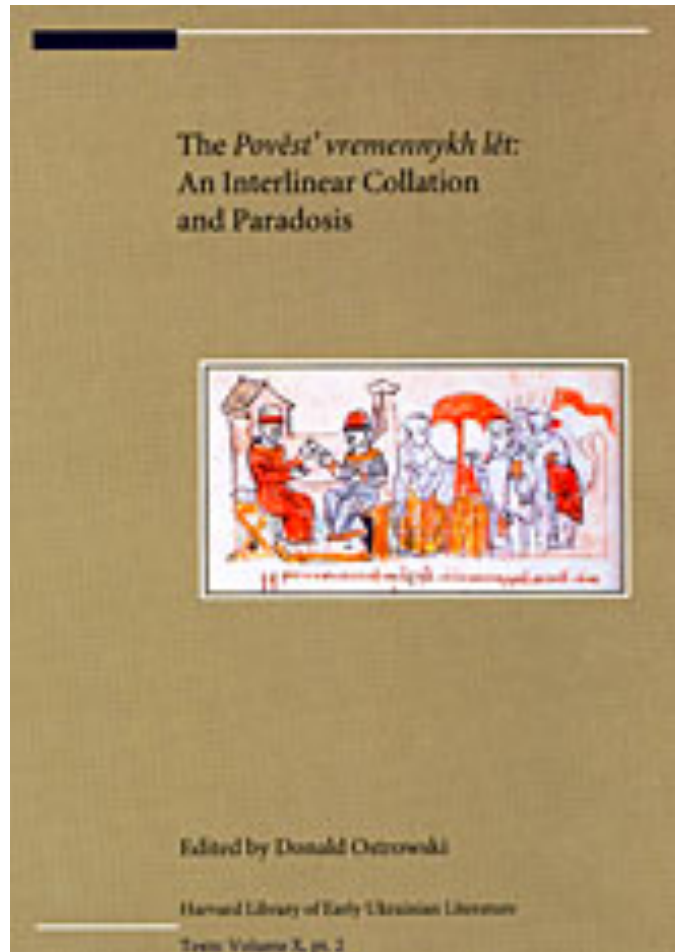
Why collate the *PVL*?

- Textual comparison
 - Relationships among the copies
 - Construction of a paradosis (alpha text)
 - History of transmission beyond alpha
- Linguistic comparison
- Orthographic comparison
 - Requires diplomatic transcription

Practical issues

- No funding
- Edition is under constant development
 - Donald Ostrowski, with David J. Birnbaum, Francis Butler, Inés García de la Puente
 - Must be able to rerun collation
 - Must be fully automated

Print edition



- *The Povest' vremennykh lět: 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 индикня в долготу

Trin: персида ватрь даже и до индикня в долготу

Radz: персида. ватрь. доже и до индикня. в долготу

Acad: персїда. ватрь. дон и до индикїа. в долготѣ |

Hypa: перьсидѣ. ватрь. доже и до индикня. в долготу

Khle: персида. ватрь. даже и до индикїа. въ | долготу

Bych: Персида, Ватрь, доже и до Индикня в долготу,

Shakh: Персида, Ватрь доже и до Индикня въ дълготу,

Likh: Персида, Ватрь, доже и до Индикня в долготу,

α : Персида, Ватрь доже и до Индикня въ дълготу,

Why interlinear?

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

How interlinear?

- Initially
 - Alignment by line (per Karskii 1926 edition)
 - Used by Müller in the *Handbuch* and Cross in the English translation (1930)
- Target
 - Alignment by word

Digital versions

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

Print edition workflow

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

First digital version

1, 4

Lav	персида. ватрь. тоже и до индикиѡа в долготу
Tro	персида ватрь даже и до индикия в долготу
Rad	персида. ватрь. доже и до ин ^А икиѡа. в до лготу
Aka	персїда. ватрь. до ^Ж и и до индикїѡа. в долготу
Ipa	перь сида. ватрь. доже и до инь дикиѡа. в долготу
Xle	персида. ватрь. даже и до индикїѡа. въ долготу
Vuř	Персѡда, Ватрь, доже и до Нѡдикня в долготу,
řax	Персѡда, Ватрь доже и до Нѡдикня въ дѡлготу,
Lix	Персѡда, Ватрь, доже и до Нѡдикня в долготу,
α	Персѡда, Ватрь доже и до Нѡдикня въ дѡлготу,

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
- Diplomatic transcription
 - Efficient comparison algorithms require exact string matching, which is rare in diplomatic transcription
 - Finding *closest* match requires a different—less computationally efficient—algorithm (method) than finding *exact* match

Word-aligned version

1,4

<i>Lav</i>	персида.	ватрь.	тоже	и	до	индикиа	в	долготу
<i>Tro</i>	персида	ватрь	даже	и	до	индикия	в	долготу
<i>Rad</i>	персида.	ватрь.	доже	и	до	индикиа.	в	до лготу
<i>Aka</i>	персїда.	ватрь.	дожи	и	до	индикїа.	в	долготу
<i>Ipa</i>	перь сида.	ватрь.	доже	и	до	инь дикиа.	в	долготу
<i>Xle</i>	персида.	ватрь.	даже	и	до	индикїа.	въ	долготу
<i>Byč</i>	Персида,	Ватрь,	доже	и	до	Индикня	в	долготу,
<i>Šax</i>	Персида,	Ватрь	доже	и	до	Индикня	въ	дѣлготу,
<i>Lix</i>	Персида,	Ватрь,	доже	и	до	Индикня	в	долготу,
<i>α</i>	Персида,	Ватрь	доже	и	до	Индикня	въ	дѣлготу,

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
 - 40000 ~ ᚠᚱ ~ 40 ТЫСЯЦЬ
 - РАЗУМЬН- ~ СЗМЫСЛЫН-

Adapting CollateX

- Preprocessing
- Collation (CollateX)
- Postprocessing

Preprocessing

- Normalized “shadow” copy
- Soundex simplification
- Collate on normalization, render 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
 - Consonants carry more information than vowels
- Information load may be sensitive to position
 - Beginning of word carries more information than end
 - Especially true for lexical (not morphological) searching in inflected languages

Adapting Soundex to Church Slavonic

- Neutralize variant spellings of initial vowel
 - ѿγ, γ, Ѯ = γ
 - ѡ, ѡ, Ѡ, ѡ = ѡ
- Case fold, neutralize consonantal variants
 - Not always one-to-one, e.g., ѱ = ѠѲ
- 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?
 - Ѣ҃҃҃҃, Ѣ҃҃҃, Ѣ҃ = ѢѠ҃҃҃ (Ѣ҃)
- Truncate or zero-pad (to what length?)

Soundex sample

- Ch397 и възвра|тить дьщерьше своѳе.
 - Ch384 и възвратит̂ дьщершоу свою.
 - Nbkm298 и възвратити братанитцѧ | своѧ
 - Berlin и въз̌вратити | братаницѧ свою.
-
- Ch397 и взвр дштр св
 - Ch384 и взвр дштр св
 - Nbkm298 и взвр брtn св
 - Berlin и взвр брtn св

Two types of normalization

- Collation
 - Find alignment points
 - Coarse adjustments
 - No harm in conflating, e.g., imperfect and aorist or 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 results
- Utilize forced matches
 - A B C
 - A D C
- Misses
 - Gap in alignment (no forced match)
 - Imperfect match
 - фраки ~ фраци
 - CollateX recognizes only perfect matches
 - Unable to recognize *closest match*
 - Computational complexity

3,5

3,5

<i>Lav</i>	гарѣмати	тавр[і] ани.	сирѣфыа.	фраци.
<i>Tro</i>	гарѣмати	тавриани	скуфиа	фраки
<i>Rad</i>	сармати	таврилни	скоуфиа	и фраци
<i>Aka</i>	сармати.	таврїани	скоуфїа.	и фра ци
<i>Ipa</i>	сармати.	тавриани.	скуфиа.	фраци.
<i>Xle</i>	сармати.	таврїани.	скѣфїа	фра ци.
<i>Byč</i>	Сарѣматн,	Тавриани,	Скуфна,	Фрацн,
<i>Šax</i>	Сарматн,	Тавриани,	Скуфня,	Фрацн,
<i>Lix</i>	Сарѣматн,	Тавриани,	Скуфна,	Фрацн,
<i>α</i>	Сарматн,	Тавриани,	Скуфня,	Фрацн,

Numbers

18,4

<i>Lav</i>						.[до]	ДѢДА.
<i>Aka</i>	.ѣ.ѣ.	а		исхо женїа	мѡисѣѡѡа	до	ДѢДА.
<i>Ipa</i>	ѣ.ѣ.		ѡ	исхожениа	мо исѣѡѡа.	до	ДѢДА.
<i>Xle</i>	.ѣ.ѣ.	а	ѡ	исхо*нїа	мѡѡсеѡѡа	до	ДѢДА
<i>Byč</i>	430;	а	отъ	нсхоженна	Монсѡѡа	до	Давнда
<i>Šax</i>	430;	а	отъ	нсхоження	Монсѡѡа	до	Давыда
<i>Lix</i>	430;	а	от	нсхоженна	Монсѡѡа	до	Давнда
<i>α</i>	430;	а	отъ	нсхоження	Монсѡѡа	до	Давыда

Problem areas

- Gaps in alignment
- No perfect match
- CollateX takes leftmost match
- 3,5
 - Orthography
 - скуфиа фра^ки (Tro)
 - скоуфиа и фра^ци (Rad)
 - Soundex
 - скф фр^к
 - скф и фр^ц

Postprocessing

- If there's a gap
 - If the column all matches, keep it
 - Else
 - Find all unique Soundex values in column and following
 - Move token to column with closest match
 - Edit distance (Damerau-Levenshtein)
 - » Insertion, deletion, substitution, transposition
 - Recursive

9,2

9,2

<i>Lav</i>	и	то	творать		мовенье	собѣ	а	не	мученье.
<i>Tro</i>	и	то	творять		мовенье		а	не	мученье
<i>Rad</i>	и	такѡ	творать	не	мытвѡу	собѣ		но	м ^{оу} ченіе.
<i>Aka</i>	и	такѡ	творатъ	не	мытвѡу	собѣ		но	моученье.
<i>Ipa</i>	и		творя ^т	не	мытву	себѣ	а	не	му ченье.
<i>Xle</i>	и		тво ^р а ^т	не	мытвѡу	себѣ		но	мжчєнїе,
<i>Byč</i>	н	то	творять		мовєньє	сѡвѣ,	а	нє	мучєньє”.
<i>Šax</i>	н	то	творять		мѣвєннє	сѡвѣ,	а	нє	мучєннє”.
<i>Lix</i>	н	то	творять		мовєньє	сѡвѣ,	а	нє	мучєньє”.
<i>α</i>	н		творять	нє	мытву	сєбѣ,		нє	мучєннє”.

In case of ties

- Thesaurus
- Most matches
- Length

Thesaurus

- Collect forced inexact matches
- Edit manually
- Use to break ties
- Close matches
 - ПОЛОМИША ~ ВЪЗЛОМИША
 - ПАМШ ~ ВЗАМ
- Non-matches
 - РАЗУМЪНЪ ~ СЪМЫСЛЪНЪ
 - РЗМНЪ ~ СМСЛЪ

Thesaurus

220,9

<i>Lav</i>	НАЛЕГША	первое	НА	СТОПОЛКА	И	ВЗЛОМИША
<i>Rad</i>	НАЛАГОША	первие	НА	СТОПОЛКА·	И	ПОЛОМИША
<i>Aka</i>	НАЛАГОША	первое	НА	СТОПОЛКА·	И	ПОЛОМИША
<i>Ipa</i>	НАЛАГОША	пе рвое	НА	СТОПОЛКА·	И	ВЪЗЛО МИША
<i>Xle</i>	НАЛЕГОША	пръвое	НА	СТОПОЛКА.	И	ВЪЗЛО МИША
<i>Byč</i>	НАЛЕГОША	первое	НА	Святополка,	Н	ВЗЛОМНША
<i>Šax</i>	НАЛЕГОША	първое	НА	Святопълка,	Н	ВЪЗЛОМНША
<i>Lix</i>	НАЛЕГОША	первое	НА	Святополка,	Н	ВЗЛОМНША
<i>α</i>	НАЛЕГОША	първое	НА	Святопълка,	Н	ВЪЗЛОМНША

ВЗЛОМИША

ПОЛОМИША

ВЗАМ

ПАМШ

What's next: many-to-one

141,11

<i>Lav</i>	а	прочѣхъ	вои	.ѿ.		и	поне	на	ѣтополка	нарекъ
<i>Rad</i>	а	прочѣхъ.		ѿѿ.		и	поне	[на]	ѣтополка.	на рекъ
<i>Aka</i>	а	прочѣхъ.		ѿѿ.		и	поне	на	ѣтополка.	нарекъ
<i>Ipa</i>	а	прочѣхъ	вои	.ѿ.	тысяць.	и	поне	на	ѣтополка.	нарекъ.
<i>Xle</i>	а	прочѣхъ	вои	.ѿ.	тысяцъ.		поне	на	ѣтополка	нарекъ
<i>Byč</i>	а	прочѣхъ	вой	40000,		н	поне	на	Святополка,	нарекъ
<i>Šax</i>	а	прочѣхъ	вон	40	тысяць,	н	поне	на	Святопѣлка,	нарекъ
<i>Lix</i>	а	прочѣхъ	вой	40,000,		н	поне	на	Святополка,	нарекъ
<i>α</i>	а	прочѣхъ	вон	40	тысяць,	н	поне	на	Святопѣлка,	нарекъ

Recent developments

- CollateX has been ported from Java to Python
 - Python module
 - Easier to modify
 - Closer integration of preprocessing and postprocessing
 - Fivefold improvement in processing time (so far)
- Decision tree architecture under development for CollateX core
 - Closer integration of postprocessing with core