

Necromancing Diels: computerising the phonological analysis of early Slavonic texts using existing treebank data and a Late Common Slavonic computerised inflectional morphology

0. Introduction

Much progress has been made in the last twenty years in early Slavonic corpus linguistics as a result of the Old Church Slavonic part of the PROIEL project (Haug & Jøhndal 2008) and its subsequent expansion as the TOROT treebank (Eckhoff & Berdičevskis 2015), such that currently just over 240,000 words of canonical OCS have been manually lemmatised, part-of-speech and morphologically-tagged, and syntactically parsed. The focus of these projects, however, has been exclusively on the higher-level linguistic domains of syntax, semantics, and pragmatics: surface-morphology has been of only incidental concern, for example in investigations into differential-object marking (Eckhoff 2015, 2022). No inflection-class data is included in these corpora, and phonology has been totally ignored to the extent that some of the texts (esp. Kiev Folia, Codex Suprasliensis, and partially Codex Zographensis) contain quite severe typographical inconsistencies and errors that make them dangerous to use without reference to the manuscripts.

That being said, enough information is included in the lemmatisation and morphology-tagging that, with a few exceptions (e.g. comparatives), the morphological shape of the inflected text-forms can be predicted from just the tag-information, provided that inflection-class annotations are added to the lemmas. This means that the immediate Late Common Slavonic ancestors of surface-text forms can be generated by using a database of LCS inflectional-endings, reconstructing and inflection-class-marking the LCS stems of the lemmas, and then applying inflectional-endings to the stems according to the word's morphology-tag annotation¹. Such LCS reconstructions are an extremely useful form of 'phonological annotation', since theoretically all the information required to give rise to an attested form must be present in any correct reconstructed proto-form, and the complete regularity of the idealised LCS forms makes texts predictably searchable regardless of orthographic variability, abbreviations, or other irregularities in the surface-texts. When applied to whole texts, they make the exhaustive investigation of almost any phonological or orthographic question trivially easy compared to manually reading and extracting relevant forms, or using TOROT's existing lemmatisation and morphology-tagging to try to gather morphological categories which might contain the sound-groups one is interested in.

In the next section I will describe my computerised LCS inflectional-morphology in more detail, show how it can be used to "autoreconstruct" different OCS texts, and explain how difficulties caused by things like morphological innovations, badly-integrated foreign loanwords, or insufficiently-precise tagging-data can be overcome. (Possibly include here some demonstration of 'exhaustive investigation' of the autoreconstructed Marianus, since that is the highest-quality TOROT text and the only one virtually 100% covered by my lemmas?)

Since morphology-tagging and lemmatisation are a prerequisite for my method of automatic reconstruction, Section 2 will survey recent work on automating these tasks for early Slavonic texts. Thanks to modern deep-learning techniques and the large and growing amount of manually-produced training-data in Eckhoff's corpus, accuracies of 90%+ can easily be reached (depending on the target-text), and I will see how far up this can be pushed by better neural-network design and more careful and informed pre-processing of training and target-data.

As a test-case of "wholly automatic" phonological annotation, Section 3 will apply such methods to the Codex Assemanianus, an OCS lectionary containing most of the gospels which has been digitised in an ASCII-encoded format by Jouko Lindstedt but is not included in Eckhoff's corpus. Accuracy will be evaluated by comparing both the automatic tagging and lemmatisation, and the resulting LCS reconstructions, to 10 randomly-selected manually-annotated shorter sections.

¹ Morphological innovations and variations are detected by inspecting the text-forms and then applying 'alternative' endings as specified in the inflectional-endings database; see Section 1 for more detail.

Section 4 will then use the wholly-automatically-reconstructed Assemanianus as the basis for a short investigation into aspects of its phonological and orthographic system, which will be compared against existing treatments of this text in the literature, to see to what extent useful insights can be extracted even without any form of manual-annotation.

1. Auto-reconstructing texts using a computerised Late Common Slavic inflectional morphology

The premise of my chosen form of "phonological annotation" is that the earliest Slavic texts reflect languages which are **structurally** close enough to the broadly-agreed-upon system of Late Common Slavonic that the forms underlying the manuscript-spellings are more or less trivially derivable (by the application of sound-change rules) from their theoretical LCS ancestors.

By 'structurally' I am referring to the structure at the phonological level; structural changes at higher levels of analysis (i.e. inflectional morphology, derivational morphology) are of no concern unless they are **made possible only by intervening phonological changes**.

My contention is that before about 1100 not enough of these structural changes are in evidence in any Slavic text, and thus they can be relatively straightforwardly indexed using a well-chosen LCS system. Before giving examples of structural changes that are problematic for such an indexing-system, it's necessary to first lay out my LCS system in full:

1.1 Late Common Slavonic as a "phonological index"

In order to account for as much of the subsequently attested Slavic as possible, a point after the monophthongisation of diphthongs, but before the Second and Third Velar Palatalisations (PV2 and PV3) is chosen as the point of departure, because of the difference between the West Slavic /š/ and South/East /ś/ reflex of these two palatalisations of *x (Cz. loc. pl. *dušícĥ* vs Suprasliensis. *доуѣхъ* <*duxĕxĥ; Polish *wszak* vs Supr. *вѣсакъ*, Ru. *всѣхъ/ий* <*vĕx-akĥ), as well as the probable complete absence of PV2² in northern East Slavic (Old Novgorodian, see Zaliznjak 2004: 42-45 for the evidence), and the blocking of PV2 by an intervening *v in West Slavic (Pol. *gwiazda*, Cz. *květ* <*gvězda, *květĥ, etc.).

To be explicit, the native phonemes in my LCS system are given in the tables below:

2 The evidence regarding the possible absence of PV3 from Novgorodian is far less convincing: the Birchbark letters abound with examples of the PV3 reflex of *k (e.g. letter №439 from around 1200 has *свинѣцѣ* <*svinĕkĥ and *полотѣнѣца* <*polĕtnĕca), and those of *g are not unknown: Zaliznjak (2004: 47) admits that palatalised forms of the Germanic loan *кѣнѣзѣ* <*kĕnĕgĥ- are the rule, but considers this to be a "supradialectal" word originating outside of the Novgorodian dialect-area; Galinskaja (2014: 10) is less convinced and adduces the form *оуѣрѣзѣ* 'earrings' from letter №429 as a word of "вполне бытового характера" which thus supposedly shows a native Novgorodian reflex of PV3 of *g. (This is commonly assumed to be a Turkic loan, cognate with e.g. Kazakh *сырға*, but the fact that it appears in Slavic with front-vowels (Ru. *серьга*), unlike its back-voweled Common Turkic cognates, and the fact that it was borrowed early enough to undergo PV3 at all, suggests that Vasmer's derivation of it from "Old Chuvash" (i.e. some form of Oghur or Bulgar Turkic) is correct, and it thus belongs to an earlier layer of Turkic loans than those borrowed from the Kipchak dialects of the Polovtsians (e.g. Ru. *камыш* < *qamıŝ (> Kaz. *қамыс*)).

More importantly, as Galinskaja (op. cit.) points out, in all of the well-known Novgorodian forms of the pronoun *vĕxĥ 'all' which supposedly show a lack of PV3 by retaining both /x/ and back/hard desinences (e.g. fem. gen. sg. *кѣхѣ* <*vĕxojĕ from letter №850), and which come from letters which otherwise correctly convey the jers (by writing <ѣ,e> for *ĕ and <о,о> for *o), the weak-jer is always written with <ѣ,o>, unambiguously suggesting a /ь/ pronunciation. These forms therefore more likely point to a LCS doublet-form *vĕxĥ which would never contain the conditioning environment for PV3 anyway, and thus you can't use them as evidence of a lack of PV3 in Novgorodian (on the plausibility of such a doublet see Galinskaja (2014: 14), though cf. Zaliznjak's (2004: 54) less convincing explanation of the /ь/ in these words as an assimilation of original /ь/ to the back-vowels of the following syllable).

Table 2: LCS Vowels after the monophthongisation of diphthongs

	Front	Back	
High	i		y u
	ɨ ʏ ɪ	ʏ	ɯ ʊ
Mid	e ɛ		ɤ ɔ
	ě ě		
Low	æ		
		a	

*Table 1: LCS consonants before PV2/PV3
(adapted from Winslow 2022: 304)*

Labial	Dental	Palatal	Velar
m	n	ń	
b p	t d	ħ ǧ	k g
	s z	š ž	x
		č	
	l	ĺ	
	r	ř	
v		j	

In addition, the following symbols are used to represent phonemes of wholly foreign origin in order to represent badly-integrated foreign borrowings, whose level of integration into the native system we deliberately do not take a position on: /k̑ ġ ŋ f ü/, e.g. in respectively кѣтъ <*kitъ, ѿѣмонъ <*igemonъ, хитонъ <*χitonъ, иосифъ <*ijosifъ³, and мѹро <*müro. Almost none of the words containing these symbols would actually have existed in the language during Common Slavonic times, but they need to be included in the indexing-system because they often contain native Slavic elements (f.ex. inflectional endings). Normally they represent specific sounds in the source-language (usually Greek), so including them is useful for investigating the process of these sounds' integration into the native systems. For instance, the extent to which Greek /ü/ is integrated into either native /i/ or /u/ can be seen in variations in the OCS spellings of the word for 'Egypt': ѧѦѢрѥѡѧ vs ѧѢѤрѥѡѧ vs ѧѢѢѢрѥѡѧ vs егѢуѣтѣ vs ѧѦѢѢрѥѡѧ⁴, etc.. One might also ask whether a separate <Ѧ> letter for /ġ/ (and the writing of <ѣ> with the palatalisation-diacritic) could be linked to the inadmissibility in the native systems of soft [kʲ, gʲ] sounds, and whether their replacement with regular <к, г> or <ѣ, к> was more likely in systems with some level of native [kʲ, gʲ] (for instance, in Rus' after the so-called Fourth Velar Palatalisation, or in Novgorod due to the retention of native velars before front-vowels because of the non-action of PV2, etc.); in any case such questions are far easier to investigate if all relevant forms can be reliably retrieved by giving them even a consciously artificial LCS representation.

I have deliberately not included accentual information in my reconstruction of vowels, even though such information is in fact required to explain certain differing manuscript-reflexes, e.g. Russkaja Pravda fem. acc. sg. **ꙋꙋбѣ** < *orb-ŏ vs Uspenskij Sbornik nt. acc. sg. **ꙋꙋдѣ** < *ordl-o, because for too large a proportion of the vocabulary this information is not sufficiently securely and uncontroversially reconstructed to justify its inclusion, and anyway the (often post-LCS) derivational processes which are responsible for most of the actual words in the attested texts (and the inevitable accentual levelling processes likely to have occurred in the course of these derivations) complicate things even further.

The two extra nasal-vowels /ɣ/ and /ě/ are required to account for the split between North (East and West) and South Slavic forms of certain inflectional-endings: *ɣ for the nom. sg. masc./nt. pres. act. participle of certain verb-classes whose present-stem ends on a hard-consonant, which in South

3 Of course the sequence /jo/ violates LCS phonotactics as well.

4 Forms are given as they appear in the manuscripts; modern fonts mean that the misleading and unhelpful practice of transcribing Glagolitic into Cyrillic is no longer justified in any context.

Slavic remains high and backed, e.g. Supr. *zъzъl* <*zovy, Psalterium Sinaiticum *ⲁⲟⲩⲁⲗⲁⲩⲉⲧ* <*stergy-jъ, Codex Marianus *ⲁⲗⲁⲩⲉⲧ* <*jĀdy-jъ (as is clear from these forms, some dialects of early OCS retained some kind of nasal character in this vowel and may even have developed the special "hooked" nasal letter <ѣ> for it), but in North Slavic lowered to /a/: Old Polish (Kazania Świątokrzyskie) *recā* /r'eka/ <*reky, Ru.Ch.Sl. (Vita Methodii) *въсемогъи* <*vъxemogy-jъ. /ĕ/ is responsible for the NSl. /ĕ/ vs SSL. /ĕ/ shapes of jo-stem masc. acc. pl. and the ja-stem nom./acc. pl. and gen. sg. endings (Kortlandt 1979), which are reflected in respectively the post- and pre-revolutionary spellings of the Russian nom./acc. pl. long-adjective endings *-ѣ* <*yjĕ <*yjĕ vs *-ѣя* <*yja <*yjĕ <*yjĕ.

- 1.) the lack of any device in the Glagolitic alphabet to render /ja ná rá la/ sequences (for which Glagolitic texts must use the jat' <Ɑ> letter whose base-value is /ě/);
- 2.) overwhelming spellings of palatal-letter (<ⱭⱮⱲⱰ>) + jat' in the Kiev Folia (the oldest and therefore least distant ms. from the 'original' OCS, as first codified by Cyril and Methodius and for which the Glagolitic alphabet was devised) for the reflexes of LCS *č/š/ž/ħ + *Ē (e.g. ⱭⱮⱲⱰⱭⱮⱲⱰⱭⱮⱲⱰ <*ob-věhĒl, ⱭⱮⱲⱰⱭⱮⱲⱰⱭⱮⱲⱰ <*dušĒmi), as well as occasional traces of such spellings in later Glagolitic OCS (e.g. Psal. ⱭⱮⱲⱰⱭⱮⱲⱰ <*čĒše); and
- 3.) the evidence of certain modern Bulgarian dialects, which have reflexes of LCS *ě in words like ⱭⱮⱲⱰⱭⱮⱲⱰ <*žĒba 'toad',

The syllabic liquids /ʃ ʒ ʎ/ are included as unitary vocalic phonemes, following Schenker (1995: 94), rather than as combinations of /ʃ ʒ/ + /ʎ ʎ/, because these groups descend from PIE syllabic liquids and many descendant South Slavic dialects which retain syllabic liquids in this position (including most of those underlying canonical OCS) do not show any evidence of an intervening oral-vowel + liquid stage (such a view is shared by Bethin 1998: 71-72; cf. also Bulgarian dialectal

The need for both front and back $*\text{r}_\text{f}$ $*\text{r}_\text{b}$ is unambiguously shown by the East Slavic reflexes /er/ and /or/ (Ru. *смерть*, *морковь*), but $*\text{r}_\text{f}$ vs $*\text{r}_\text{b}$ is more complicated: PIE $*\text{p}_\text{f}\text{nos}$, $*\text{w}_\text{f}\text{lk}^\text{w}\text{os}$ > Lithuanian *pilnas* 'full', *wilkas* 'wolf' (LCS $*\text{p}_\text{f}\text{ln}_\text{b}$, $*\text{v}_\text{f}\text{lk}_\text{b}$) vs Lith. *stulpas* (LCS $*\text{st}_\text{f}\text{lp}_\text{b}$ 'pillar') suggests that Balto-Slavic had differentiated front/back variants of the PIE syllabic $*\text{r}_\text{f}$ (Bethin 1998: 69), but the ancestor to East Slavic backed all vowels preceding tautosyllabic /l/ (Ru. *молоко* < Proto-ESl. $*\text{molko}$ < LCS $*\text{melko}$ > OCS *млѣко*), and thus only has /ol/ reflexes here: Ru. *волк*, *столб*, *полный*. It's true that Polish has *wilk* and *milczeń* (< $*\text{m}_\text{f}\text{č}\text{Ēti}$), but the Polish reflexes are complicated and likely have more to do with the surrounding consonants: $*\text{p}_\text{f}\text{ln}_\text{b}$ by contrast gives *pełny* with hardened /l/ and the Polish non-palatalising-/e/ reflex of $*\text{r}_\text{b}$, and the differing reflexes in *wierzch* < $*\text{v}_\text{f}\text{ch}_\text{b}$, *śmierć* < $*\text{s}_\text{m}\text{r}_\text{f}\text{t}_\text{b}$ and *martwy* < $*\text{m}_\text{f}\text{tv}_\text{b}\text{y}_\text{b}$ rule out any explanation based on the nature of the LCS syllabic-liquid alone (for more discussion see Bethin op. cit.: 73-75).

Dejotation

The compelling hypothesis, first proposed by Durnovo (1929: 55-58) but most recently elaborated by Vermeer (2014: 209-214), and accepted by Mathiesen (2014: 197 fn. 22) and Winslow (2022: 310 fn. 25), according to which the Urkirchenslavisch reflexes of *h, ħ were close enough to foreign /g k/ before front-vowels that the original Glagolitic system used <Ѡ ѡ> for both sets (i.e. alongside attested Ѡ ѡ ѡ ѡ < ἡγεμὼν would have been **ѠѡѠѠѠѠ < *osŕĕjeni, and alongside attested Ѡ ѡ ѡ ѡ < *dĕherĕ would have been **ѡѡѡѡѡѡ < κῆνσοϛ⁷), does not prevent us from keeping the foreign sounds separate for our LCS stage, since clearly they differed enough in all the dialects underlying actually attested OCS to be written separately.

Therefore the only dejotation-reflexes which are not fully “undoable” under my notation are those of *sj and *zj and *kt, viz. /š ž ħ/, which are each of multiple origin.

7 Interestingly, this aspect of the hypothesised Ūrksl. orthographic system has rearisen in the modern Macedonian standard due to Turkish loanwords: *ќемер* < Tk. *kemer* 'belt', *ќе* < *[хъ]hē[тъ]; *ѓон* < Tk. *gön* 'leather', *меѓу* < *meħu.

some brief words about dejetated labials and epenthetic ɫ, which in my system is taken as the regular LCS outcome, removed only by later dialect-specific developments

The tendency for ECS *a- to have taken prothetic /j/ by LCS times (in accordance with the drive towards open syllables) can make it difficult to distinguish these groups from *jĀ- in the absence of wider Indo-European evidence. Normally I've followed Derksen (2009), or the ESSJA, but for certain lexemes, e.g. *ama 'pit', which in OCS is spelt overwhelmingly with амѣ or ѣмѣ , the single Greek cognate ἄμη adduced by ESSJa I p.70 in favour of jot-less *am- is not enough to categorically exclude the alternative *jĀma.

word-initial *jъ-, *ji-, Ukr. съкати, Derksen 2003 etc.

treatment of clusters like *ъс-щдъ* *ъс-кѣ*, which in OCS are simplified but kept separate by me, vs. *ot-xod*, *ob-xod*, which I reconstruct as simplified to *ox-*, *ox-*, issues of reformation with *jer-* containing prefixes like *отъ-* in many words

An example of morphological change contingent upon structural phonological change, leading to manuscript forms which preclude any valid reconstruction of their direct LCS-stage ancestors, is the replacement of i-stem endings with those of the corresponding jo- or jā-stems, in nouns whose stems end on labials or the subset of LCS dental consonants which lack palatal counterparts, viz. /d t s z/. Evidence for such a change is furnished by the Old Russian masc gen./acc. form **ТАТА** from the 1229 Treaty between Smolensk, Riga and Gotland (Version A). LCS *tāt is a masc. i-stem noun with genitive *tati, as it still appears in the Codex Suprasliensis translation of John Chrysostom's Homily for Holy Thursday (...то кажетъ владыкы ѡловѣколюбѣѣ ꙗко прѣданныка разбоѣника тати...), but in the dialect underlying the 1229 Treaty the rise of phonemically palatalised

/t'/ after the Jer Shift means that the stem (and the nom. sg. **ТАТЬ** /tat'/) of this noun now ends on the same class of "soft" consonants as original jo-stem nouns like *pastyř > /pastyr/, where the original LCS palatal *ř has fallen together with secondarily-palatalised /r/ from plain LCS *r before LCS front-vowels, in e.g. the original i-stem *zvěř > /zvěr/. This system thus no longer distinguishes between descendants of the original LCS palatals and the newly secondarily-palatalised consonants like /t'/, both are now together in the set of 'soft' consonants, opposed to their 'plain' or 'hard' counterparts, and so tend towards taking the same set of inflectional endings (in this case those of the original jo-stems). Consequently, a word like **ТАТЬ** has begun to take jo-stem endings, including the Old Russian /a/ reflex of LCS *Ā in the genitive/accusative singular. LCS /Ā/, though, by definition can only occur after LCS palatal consonants (see above), so a reconstruction *tatĀ is just nonsensical. In the case of the dat. sg. /u/-desinence (which isn't attested in our Treaty but it exists in modern Russian *татю*), we don't even have an LCS archiphoneme available to signal a preceding soft-consonant; there's simply no way of getting from LCS *tatu to Russian /tat'u/, because such a form was only made possible by the rise of phonemic /t'/, so our ability to index it with our LCS system is gone.

//this is just rough unstructured ideas, some of which may already have been incorporated into the text above

For example, if the phonotactic rules of our theoretical LCS system allow the sequence /řĀ/ (palatal /ř/ < *rj + the archiphoneme /Ā/) to occur, then a morphological change which replaces the sequence /ri/ with /řĀ/ is of no concern, because both are equally valid LCS. If, however, the same type of morphological change were to

For example, whether or not there actually existed at the LCS stage a mechanism for deriving secondary-imperfective verbs like OCS *разрѣти* < *orzařĀti from the prefixed *разорити* *orzoriti is irrelevant, because LCS /orzařĀti/ does not violate the rule of LCS phonotactics: palatal /ř/ can be followed by /Ā/ because such a combination exists in the paradigms of wholly securely reconstructable jo-stem nouns, e.g. nt. gen. sg. *mořĀ (> Pol. *morza*, OCS *морѣ*, Ru. *морѣ*, etc.)

In the case of Supr. Gsg. masc. *звѣръ*, for an original i-stem (*звѣри* < *zvěri), a direct LCS ancestor for the attested form can still be given (*zvěřĀ), because palatal /ř/ already exists in our LCS system, and one plausible explanation for this form is that the Eastern Bulgarian dialect underlying Suprasliensis developed secondary palatalisation of LCS plain *r before front-vowels, which was then phonemicised after the fall of word-final front-jers, and that newly-palatalised /r/ fell together with original LCS palatal /ř/, so that the nom. sg. *zvěř became /zvěr/, and its stem now ended on the same consonant /r/ as original ja- and jo-stems ending on LCS *ř like *морѣ* and *воуриѣ*, so it began to be inflected as a jo-stem masculine instead of an i-stem.

It should be emphasised that the historical reality of our reconstructions is only of concern at the phonological level, that is, phonemes and phonotactics; the plausibility of higher-level structures built out of these units,