

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 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. *dušich* vs Suprasliensis. *доуѣхъ* <*duxěx̥; Polish *wszak* vs Supr. *вѣакъ*, Ru. *всѣакъ* [уѣ] <*vьx-akъ), 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, *květъ, 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* <*svinьkъ and *полотѣнѣца* <*polьtьnъka), and those of *g are not unknown: Zaliznjak (2004: 47) admits that palatalised forms of the Germanic loan *кѣнѣзъ* <*kьneg- 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* < *qamış (> Kaz. *қамыс*)).

More importantly, as Galinskaja (op. cit.) points out, in all of the well-known Novgorodian forms of the pronoun *vьxъ 'all' which supposedly show a lack of PV3 by retaining both /x/ and back/hard desinences (e.g. fem. gen. sg. *вѣхѣ* <*vьxoĭĕ from letter №850), and which come from letters which otherwise correctly convey the jers (by writing <ѣ,e> for *ь and <о,ѡ> for *ѡ), 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 *vьxъ 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
	ǐ ɨ ǐ̇		ɣ ɣ̇ ɨ̇	
Mid	e ɛ		ɔ ɒ	
	ě ǣ			
Low	æ			
		a		

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

Labial		Dental		Palatal		Velar	
m		n		ɲ			
b	p	t	d	ħ	h̥	k	g
		s z		š ž		x	
				č			
		l		ʎ			
		r		ʀ			
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: /ḳ ɡ̣ ʃ̣ f̣ ü/, e.g. in respectively *кѣтъ* <*kítʰ, *ѣгемонъ* <*igemonʰ, *хитонъ* <*xítʰonʰ, *иосифъ* <*ijosifʰ³, and *мүро* <*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': *ѣгѣп̣т̣* vs *ѣгѣп̣т̣* vs *ѣгѣп̣т̣* vs *ѣгѣп̣т̣* vs *ѣгѣп̣т̣*⁴, etc.. One might also ask whether a separate <ʎ> letter for /ɡ̣/ (and the writing of <ʃ̣> 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 <ʁ, ɣ> or <ʀ, ɰ> 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.

Vowels

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. *рѡбѣ* <*orb-ɔ vs Uspenskij Sbornik nt. acc. sg. *рѡдѡ* <*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 /y/ and /ě/ are required to account for the split between North (East and West) and South Slavic forms of certain inflectional-endings: *y for the nom. sg. masc./nt. pres. act.

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

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

participle of certain verb-classes whose present-stem ends on a hard-consonant, which in South Slavic remains high and backed, e.g. Supr. **ꙗꙋꙋꙗ** <*zovy, Psalterium Sinaiticum **ꙗꙋꙋꙗ** <*stergy-**ꙗꙋ**, Codex Marianus **ꙗꙋꙋꙗ** <*jĀdy-**ꙗꙋ** (these forms lead Kortlandt (1979:260) to posit that some dialects of early OCS retained some kind of nasal character in this vowel and may even have developed the special "hooked" nasal letter <ꙗ> for it), but which in most of North Slavic lowered to /a/: Old Polish (Kazania Świętokrzyskie) has both *recꙗ* and *recꙋ* (with the special Old Polish letter for the merged reflex of *ę and *ꙋ) <*reky, and in other texts also *biorꙋ* <*bery; Russkaja Pravda **ꙗꙋꙋꙗ**, Uspenskij Sbornik **ꙗꙋꙋꙗ** <*dojdy, Ru.Ch.Sl. (Vita Methodii) **ꙗꙋꙋꙗ** <*vꙗxemogy-**ꙗꙋ**. Kortlandt's positing of a CS *ꙗ (which he writes as *aN) is far from universally accepted, and others consider these forms the result of various dialect-specific analogical process; see references and discussion in Olander (2015: 88-92). Whatever the truth of the matter, our *ꙗ is a convenient placeholder which allows all the relevant evidence to be retrieved.

The need for the retention of the / \bar{E} / archiphoneme, which represents merged Early Common Slavonic * \bar{e} * \bar{a} in the position after palatal consonants, up to this point of LCS, is explored in detail in Winslow (2022), but the same archiphoneme (along with its short counterpart / \bar{E} /) was explicitly posited by Kortlandt as far back as 1979 (p.266) as part of his ECS system. In short, a combination of:

3.) the evidence of certain modern Bulgarian dialects, which have reflexes of LCS *ě in words like *ж'ѣба* <*žĕba 'toad' (Stojkov 1954: 74–78),

5 Other annoying pre-LCS morphological isoglosses reflected in the texts include the masc./nt. instr. sg. *o- and *jo-stem endings *-ѣмь (N.Sl.) and *-омь (S.Sl.), which are most commonly (e.g. Olander 2015:168) thought to be analogical replacements of the original instr. sg. ending ECS *ā which is preserved in the adverb *въчера ‘yesterday’, and the *-тъ (N.Sl.) vs *-тѣ (S.Sl.) verbal endings of 3rd sg. and pl. present (plus its extension to 2nd and 3rd sgl. aorists like OCS НА҃҃҃҃҃҃҃҃҃҃҃҃, OR (Uspenskij Sbornik) Б҃҃҃҃҃҃҃҃҃҃҃҃, НА҃҃҃҃҃҃҃҃҃҃҃҃). Here I have no choice but to index them with dummy-symbols in the database: *-Omь for the instr. sg. ending and *-tQ for the verb-endings.

The need for both front and back $*r_{\text{f}}$ $*r_{\text{b}}$ is unambiguously shown by the East Slavic reflexes /er/ and /or/ (Ru. *смерть*, *морковь*), but $*r_{\text{f}}$ vs $*r_{\text{b}}$ is more complicated: PIE $*p_{\text{h}}nos$, $*w_{\text{h}}lk^{w}os$ > Lithuanian *pilnas* 'full', *wilkas* 'wolf' (LCS $*p_{\text{h}}n_{\text{h}}$, $*v_{\text{h}}lk_{\text{h}}$) vs Lith. *stulpas* (LCS $*st_{\text{h}}p_{\text{h}}$ 'pillar') suggests that Balto-Slavic had differentiated front/back variants of the PIE syllabic $*l$ (Bethin 1998: 69), but the ancestor to East Slavic backed all vowels preceding tautosyllabic /l/ (Ru. *молоко* < Proto-ESl. $*molko$ < LCS $*melko$ > OCS *млѣко*), and thus only has /ol/ reflexes here: Ru. *волк*, *столб*, *полный*. It's true that Polish has *wilk* and *milczeć* (< $*m_{\text{h}}l_{\text{h}}č_{\text{h}}\text{Æti}$), but the Polish reflexes are complicated and likely have more to do with the surrounding consonants: $*p_{\text{h}}n_{\text{h}}$ by contrast gives *pełny* with hardened /l/ and the Polish non-palatalising-/e/ reflex of $*r_{\text{b}}$, and the differing reflexes in *wierzech* < $*v_{\text{h}}r_{\text{h}}x_{\text{h}}$, *śmierć* < $*s_{\text{h}}m_{\text{h}}r_{\text{h}}t_{\text{h}}$ and *martwy* < $*m_{\text{h}}r_{\text{h}}t_{\text{h}}v_{\text{h}}j_{\text{h}}$ rule out any explanation based on the nature of the LCS syllabic-liquid alone (for more discussion see Bethin op. cit.: 73-75).

Consonants

the East Slavs inherited their writing system ultimately from the Urkirchenslavisch system designed for such a dialect, rather than one which had a clear way of writing <soft consonant> + <o>, is likely the reason that /o/ reflexes are so rarely detectable in the early texts, since <e> had to be used for both /e/ and /'o/, cf. the spelling ѠВШАНЪ of the Kipchak word /jovšan/ ‘wormwood’ in the Hypatian Codex, whose modern cognates (Turkmen *jowšan* /jowšan/, Kazakh *жусан* /žuwšan/, Azeri *yovšan*) unambiguously point to a Kipchak /o/, and the history of the East Slavic /o/ reflexes remains the subject of much disagreement, so it's simpler for everyone if I continue the traditional practice of writing LCS *e after palatals, even if that strictly speaking is inconsistent with my use of *Ē.

- 7 In the database I will have to use the single Unicode characters <ř ř ǀ ǁ>, rather than what's shown in my table, since the latter cannot actually be rendered without using the letters for /r ř l l/ plus the 'combining ring below' U+0325 symbol, which means searches for the consonantal liquids on their own will also return results containing syllabic liquids. The same problem affects /ǣ y/, which I will have to replace with <ǣ ŷ>.
- 8 To my mind the only evidence in support of a genuine jer + liquid stage comes from the paradigms of verbs like OCS *сѣтъти* < *sǣtǣti, where the syllabic /ǣ/ in the stem alternates with /br/ depending on the vocalicity of the following morpheme: the e.g. 3sg. pres. *sǣtǣrěť (Zogr., Supr. *сѣтърѣтъ*) or (one possibility of the) 3rd sg. aorist *sǣtǣre (Supr. *сѣтъре*) must have /bre/, while the 3rd pl. aorist *sǣtǣřę (Supr. *сѣтърша*) and the other possibility for the 3rd sg. aorist *sǣtǣ (Psal. *сѣтъѣ*, or with a different prefix Mar. *сѣтъѣ* < *otǣ), being word-final or pre-consonantal, must be syllabic /ǣ/. The same alternation occurs in the zero-grade forms of verbs like *umerti, as is clear from the Polish reflexes *umarł* < *umǣł vs *umrę* < *umǣrę. The argument could be that at some stage, before the LCS tendency towards Open Syllables became dominant, the stems in these paradigms were surely unitary /ǣr/, /mǣr/, i.e. 3sg. aor. /sǣ.ǣr/ vs 3rd. pl. aor /sǣ.ǣr.řę/, and that the latter's closed /ǣr/ syllable was only forced to open itself up by changing to /ǣ/ because of the Law of Open Syllables. Thus at least one source of the syllabic-liquids could be shown to have developed from a vowel + liquid stage, but that still doesn't prove that they all did, or that the change of /ǣr/ to /ǣ/ in these verb-forms was not merely a move to an already-existing syllabic-liquid phoneme.

Reflexes of the so-called jot-palatalisation are all written either as unitary palatal phonemes, or in the case of jot-palatalised labials as /v́ ḿ b́ ṕ/, rather than as sequences of consonant + /j/, hence /ń í ř/ for *nj *lj *rj. The ‘dejotated’ reflexes of *tj (and *kt+front-vowel) and *dj are denoted using the modern Serbian Cyrillic letters /ћ/ and /ђ/ respectively, because the commonly used alternatives, i.e. /t̚d̚/ (as used in e.g. Olander 2015) or /k̚ ġ/ (as used by me in Winslow 2022), or variations thereof, are visually too close to symbols used elsewhere in the system. /k̚, ġ/ are anyway already used in my system for foreign /k, g/ before front-vowels, and /t̚ d̚/ look too similar to the common denotations of secondarily-palatalised post-Jer Shift /t' d'/, as used in discussions of systems like Russian or Eastern Bulgarian where they arise.

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 <ꙗ ꙗ> for both sets (i.e. alongside attested ꙗꙗꙗꙗꙗꙗ < ḡȣēmŏn would have been **ꙗꙗꙗꙗꙗꙗ < *osq̣hēni, and alongside attested ꙗꙗꙗꙗꙗꙗ < *d̥shērĕ would have been **ꙗꙗꙗꙗꙗꙗ < κῆνσος⁹), 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.

There are convincing arguments for PV2/3 having preceded dejotation, at least in more central areas, most recently presented in e.g. Vermeer (2014: 197) and Wandl & Kavitskaya (2023 244-247), and therefore it could be objected that my system, which contains the dejotation reflexes /ħhǵńĺ/ but not the PV2/3 reflexes /c ś dž/, is ahistorical. However it should be emphasised that the primary goal of my LCS reconstructions is to act as an index which allows reflexes in texts to be found, not to be a historically realistic description of some actually-existing LCS dialect. The absence of PV2 in Novgorodian shows that it cannot have preceded dejotation everywhere in Slavic, and in any case the replacement of the sequences /tj dj nj lj rj/ by articulatorily distinct combined units, no longer associated by speakers with their /t/ and /j/ phonemes, is structurally completely irrelevant unless and until these new units merge with existing phonemes (or new sequences of dental + /j/ are introduced), as e.g. in the KF dialect where /tj/ merged with /c/ from PV2/3, or in ESL where it merged with /č/ from PV1. A language which had distinct Serbian-like palatal /c' dj'/ reflexes of *tj and *dj, and also no sequences of [tj, dj], could not convincingly be argued to have undergone dejotation at the phonemic level, as these new units would just be phonetic realisations of /tj, dj/. Analysed like that, the symbols /ħhǵńĺ/ in my system strictly speaking would really just be cover-symbols for the pre-jotation sequences, but such notation is preferable since it prevents searches for groups containing /j/ alone from returning results polluted by all the dejotation-groups. As I explored in my previous article (Winslow 2022), the status of /j/ as a phoneme in the earliest OCS texts is an intricate problem, so the ability to investigate the reflexes of *j in isolation from the dejotation-reflexes is important.

Word-initial *jĀ-/ *a-

9 Interestingly, this aspect of the hypothesised Urksl. orthographic system has rearisen in the modern Macedonian standard due to Turkish loanwords: *ќемер* < Tk. *kemer* ‘belt’, *ќе* < *[хъ]he[тъ]; *ѓон* < Tk. *gön* ‘leather’, *меѓу* < *meĥu.

Like Derksen, I assume that roots going back to PIE jot-less long *ē or diphthongal *oi-, e.g. the root for ‘to eat’, PIE *h₁ēd, all took prothetic *j and merged with *jĒ- from other sources, unlike Durnovo (1929: 54), who seems to think that such a development was limited to Bulgarian and Macedonian dialects, including those underlying OCS (where in the Cyrillic mss. we get regular ѣсти etc.). Isolated nominal forms like Ru. *яѣа* (which Derksen derives from a Balto-Slavic *oi-based on Lith. *aiža* and Old Prussian *eyswo*) suggest that *ě reflexes in the modern forms of verbs like Ru. *exамь*, Pol. *jeść* are later generalisations from prefixed forms like OR **ѣНѣСТН**, where no jot-prothesis could take place (cf. Schenker 1995: 88, Winslow 2022: 302 fn.14).

Difficulties arise though when deciding how to denote foreign sources of /ij/¹² which may or may not have been integrated into the native system as reflexes of $\hat{I}j$: words like $\mu\alpha\rho\iota\eta\alpha < \text{Μαρία}$, $\sigma\tau\alpha\delta\iota\eta\iota < \sigma\tau\acute{\alpha}\delta\iota\omicron\nu$, which are well-integrated into the morphological system as a fem. ja-stem and masc. jo-stem respectively, could either be reconstructed as consciously-foreign $*\text{marij}\bar{A}$, $*\text{stadij}\bar{b}$,

12 The sequence /ij/ is not totally banned from native words, since it appears to be preserved across morpheme-boundaries, such as in prefixed-verbs like *прийти* <*prijeti or long-form adjectives like masc. nom. pl. *други* <*drug-i, but within roots it does seem restricted to these post-LCS loanwords.

or as nativised *marъjĒ, *stadъjъ, but there are no occurrences of jer-spellings in these words in the OCS texts in TOROT. Other similarly-Greek words like ΔΗΛΑΒΟΛЪ (< διάβολος), however, do show up in OCS with jer-spellings: Supr. ДЫДВОЛА, Zogr. Luke 8 and Psal. Psalm 108 ѡѡДѡѡѡѡ, which (alongside the modern Macedonian *ѓавол* with the reflex of *ђ produced by the Macedonian so-called ‘new jotation’ of /d/ after the fallen jer brought it into contact with /j/) clearly suggest an early adaption of this foreign /ij/-group to native /ĭj/. Old Russian texts even show spellings of МАРИНА suggestive of full nativisation: Laurentian Primary Chronicle *мѣрьна*, *мѣрьно*, Zadonshchina *марѣна*, *марѣна*, as well as First Novgorod Chronicle gen. sg. *васнѣна* (jostem *васнѣнн* < Βασίλειος).

Word-initial *jɐ-/*ji-/*i-

I make an exception for certain forms of the personal-pronoun *jъ, however, and write *jimъ, *jima, *jixъ *jimъ and *jimi for the masc/nt. instr. sg. and dat./instr. dual/pl., because Czech here has *jim jich jimi*.

Prefixed forms like *do-jъti 'to come, arrive' for morphological reasons have to be distinguished from the class 4 verb *dojiti/dojiši/dojimъ etc. 'to breastfeed' (and its derived noun *dojidlika), a difference which is reflected in the modern Ukrainian *dіjmu* (<*dojъti with compensatorily-lengthened /o/ > /i/) vs *doimu*. Thus /i/ can follow /j/ when the former is part of a morpheme which just happens to be stuck onto a /j/-ending stem: I similarly allow words like *ščujika (щюица) and *vojinъ 'warrior' (воинъ, as opposed to *vojъnъ, the gen. pl. of *vojъna), or the loc. sg/pl. desinences of any jo-stem noun whose stem ends on /j/, e.g. Psal. *ѡбѣдѣ* <*žerbъji.

13 Spellings like Zogr. Mark 13:3 “πετρῶς ἰ ἀκωβῶς ἰ ἰωάννης.” “Peter and Jacob and John” would suggest that this initial *i- can get dropped after an /i/ of a preceding word, but whether this points to a dropping of the non-native *i-, simple deletion of a double /i i/ (haplology), or a native-like reflex of a weak-jer /*i j̥j̥Ēkovŋ/ > /i jakov/, is not really knowable, so indexing such words with a markedly foreign initial *ij- group is again the best way of allowing such difficult cases to be investigated.

14 Psalterium Sinaiticum contains just six occurrences of non-digraph <č>: *ѣдѣи*, *ѡѡѣи*, *ѡѡѣи*, *ѡѡѣи*, *ѡѡѣи*, and *ѡѡѣи*, according to Eckhoff's digitisation, five of which I've confirmed with Altbauer's (1971) facsimile of the manuscript. *ѡѡѣи* is from Psalm 151 in the newly-discovered part and thus not included in Altbauer's facsimile.

15 The leftover 14 are things like 1st. pres. dual. *ѡѡѣи* which Eckhoff's corpus wrongly lemmatises as **jъmati* instead of **jъmĕti*, and which thus get reconstructed as **jemlěvě* instead of **jъmavě*. At the time of writing only 3227/6862 Suprasliensis lemmas have been reconstructed, but those 3227 cover 89713/99194, or 90.4%, of the words.

metathesised *zork- root is never spelt <ЗРАК> and so seems to be tolerated, even though Diels cites prepositional forms like Supr. **БЕЗДРАЗУМА**, **БЕЗДРАЛА** which come from metathesised *orT-groups <*bez ◡*orzuma, <*bez ◡*ordla but do show inserted /d/. Such inconsistency is hard to explain unless the addition of /d/ has been partly morphologised as a variant of specifically the prepositions before /r/.

With such a sound-change that appears most often at morpheme or straight-up word-boundaries, there is a strong drive to restore the underlying shape of the constituent parts, hence the modern languages have mostly restored /zr/ groups in e.g. Russian *разрешить*, and there are traces of this even in Psalterium Sinaiticum: Psalm 48 **ЗѢБѢЮЩѢ** **ПѢЮЩѢ** (Diels 1963: 122). The Old Rus. Uspenskij Sbornik is pretty consistent in keeping prefixed verb-forms like **РАЗДРУШИТЬ** <*orzrušitъ, but by the time of the Laurentian Codex we get forms like **ВЗРАДУЕМ** and **НЕНЗРѸЕННОЕ**.

Because things like the *zr > *zdr, or the *ss > *s occur most frequently at transparent prefix-boundaries, and because of the clear tendency, even in the earliest texts, to undo them, I prefer to reconstruct them will strictly speaking illegal *zr and *ss groups, because that way an investigator can see for themselves the extent of the adherence to the expected phonological development vs restoration

јъskĕliti (Meillet 1965: 133)
обновити/оходити ошьль

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 *taty 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 чловѣколюбѣіе іако прѣданника разбойника 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 *pastyрь > /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 *zvěрь > /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 *матю*), 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.

Forms like **ТАТА**, then, though they frustrate our goal of reconstructing entire texts, do provide us some objective measure of 'linguistic distance' between stages of a language, because

//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 /řĚ/ (palatal /ř/ < *rj + the archiphoneme /Ě/) to occur, then a morphological change which replaces the sequence /ri/ with /řĚ/ 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* < *orzařĚti from the prefixed *разорити* *orzoriti is irrelevant, because LCS /orzařĚti/ does not violate the rule of LCS phonotactics: palatal /ř/ can be followed by /Ě/ because such a combination exists in the paradigms of wholly securely reconstructable jo-stem nouns, e.g. nt. gen. sg. *mořĚ (> Pol. *morza*, OCS *морѣ*, Ru. *морѣ*, etc.)

In the case of Supr. Gsg. masc. *звѣрѣ*, for an original i-stem (*звѣри* < *zvěri), a direct LCS ancestor for the attested form can still be given (*zvěřĚ), 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 /zvěř'/, 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,

-Mention the problems with my class “16” verbs in the morphology-section – i.e. , PAPs in /v/ aren’t very realistic for *žьŋq, bastard Suprasliensis has PPP *заклатъ* (a noisome foulness), etc.
-Could talk about the impossibility of dealing with *съмѣшѣ* deviances of S-aorist *съмѣшѣ* (vs. *съмѣтѣ* *съмѣтошѣ*), since unlike with nasal-stems, the deviance-slot here is taken up with the -ox-aorists, leaving no room for deviantly-RUKI’d S-aorists

-Could use the *овсяяныи* OR adjective as another example of derivational-morphology made possible only by the rise of soft /s’/ (if it can be confirmed that the *ѣнь adjectival suffix in e.g. OCS *оловънь* ‘leadene’ is LCS)

-*възлакати* would be a good one to use to talk about the unmetathesised groups like *old-, *olk- etc., because the “corpus-forms” table of my thing shows many examples of metathesised and unmetathesised forms

LCS Morphology and the Autoreconstructor

- Consonant-stems – with the *tel- suffix agent-nouns, I mostly follow people like Meillet (1965: 426) in taking consonant-stem endings in most of the plural, but the nom. pl. it’s difficult to agree with his positing of a plain /le/, as opposed to palatalised /Ě/ desinence (i.e. with the consonant-stem vowel on the jo-stem stem), because Zographensis and Suprasliensis are consistent in marking such forms with their palatalisation-diacritic.

-Talk about the pres. forms of *telhi and link back to the discussion about the difficulties with syllabic *l̥, saying that Derksen and the two Czech dictionaries cite *tl̥q forms, and that Zogr. mostly spells this group with <лѣ> as well, which would suggest an switch from e- to o-grade ablaut between the full-grade and zero-grade stems, but also that the issue is confounded by the existence of an o-grade form of the verb *tolhi suggested by the PPP form **ⱭⱮⱲⱰⱦⱤⱴⱶⱵⱹⱼⱾⱿⱽⱺ** <*protolčēnqj̑ in Psal. Psalm 138

Autoreconstructed forms are actually built out of a pre-dejotation stage, with dejotation applied as a post-processing step, because this greatly simplifies the inflectional morphology in places like the past-active-participle and 1sg pres. indic. of class IV (-iti) verbs: we can just use the desinences *-jь and *jǫ regardless of stem-consonant, and then apply dejotation later in a post-processing step that every word undergoes, rather than needing a whole set of consonant-mutation rules for these endings. Therefore it would be possible to allow searching based on pre-dejotation forms, but in the case of *sj, *zj wider Indo-European evidence is needed to distinguish their LCS /š ž/ reflexes from the identical outputs of PV1 (e.g. Gothic *siujan* confirms an ECS form *sjū-tei for the verb *šiti), which it is outside the scope of this project to consider, as the goal here is to enable investigations of actual texts, for which such ECS differences are irrelevant. I cannot therefore consistently offer pre-dejotation reconstructions, because stems containing reflexes of *sj and *zj are only ever reconstructed with š ž.