

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* [uŭ] <*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 <ь,е> for *ь and <о,ѣ> for *ѣ), the weak-jer is always written with <ь,о>, 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).

	Front	Back
High	i	y u
	ɨ ʏ	ɯ ʊ
Mid	e ɛ	ɔ o
	ě ě	
Low	Æ	
		a

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

Vowels

The two extra nasal-vowels /ȳ/ 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.

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-*ja*, Codex Marianus ꙗꙋꙋꙗ <*jĀdy-*ja* (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 <ѥ> for it), but which in most of North Slavic lowered to /a/: Old Polish (Kazania Świętokrzyskie) has both *recą* and *recø* (with the special Old Polish letter for the merged reflex of *ę and *o) <*reky, and in other texts also *biorø* <*bery; Russkaja Pravda ꙗꙋꙋꙗ , Uspenskij Sbornik ꙗꙋꙋꙗ <*dojdy, Ru.Ch.Sl. (Vita Methodii) ꙗꙋꙋꙗ <*vъxemogy-*ja*. Kortlandt's positing of a CS **y* (which he writes as *a_N) 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 **y* 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'еба* <*žĀeba '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 syllabic liquids /ǃ ǂ ǃ Ǆ/ are included as unitary vocalic phonemes, following Schenker (1995: 94), rather than as combinations of /ǃ ǂ/ + /ǃ Ǆ ǃ Ǆ/, 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 evidence in Stojkov 1954: 130-131, where hard consonants precede reflexes of the LCS /ǃ Ǆ/ even in dialects with secondarily-palatalised consonants before fallen weak LCS /ǃ ǂ/).

While most OCS shows no sign at all of a front-back distinction in the syllabic-liquids and writes the reflexes of these groups overwhelmingly with <рз> and <лз>, the Kiev Folia, which is the only OCS text that reflects a pre-Jer Shift stage and is very nearly flawless in its etymologically correct rendering of the jers, also spells *r̥_j, *l̥_j and *j̥ as one would expect: рѣръѣтъ - ѡвѣрьѣтъ - ѡбѣлъѣвъ - ѡбѣльѣвъ - ѡбѣльѣвъ (<*r̥_j, <*l̥_j> and <j̥>) (Winslow 2022: 313), and even Zographensis spells all 5 occurrences of *vľk- ‘wolf’ with вѣльк-/вѣльъв- and all 15 instances of its *-mlǫč- root with -мѣльчъ- (e.g. мѣльчъ + ꙗже). Therefore, taken as a whole the Slavic evidence pretty securely points to front and back variants of both syllabic liquids, and for searching purposes it’s far preferable to denote them with separate symbols⁷ rather than as the sequences /ɾ ʀ ɻ ɺ ɽ/⁸.

Dejotation

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 ѣти 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 /ĭj/: words like *μαρινα* < Μαρία, *stadii* < στάδιον, which are well-integrated into the morphological system as a fem. ja-stem and

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* <*drugi-ji, but within roots it does seem restricted to these post-LCS loanwords.

masc. jo-stem respectively, could either be reconstructed as consciously-foreign *marijĀ, *stadijĀ, or as nativised *marĭjĀ, *stadĭjĀ, but there are no occurrences of jer-spellings in these words in the OCS texts in TOROT. Other similarly-Greek words like ΔΗΑΒΟΛЪ (< διάβολος), however, do show up in OCS with jer-spellings: Supr. ДѢАВОЛА, Zogr. Luke 8 and Psal. Psalm 108 ѡѡДѡѡѡѡ, which (alongside the modern Macedonian *ѓавол* with the reflex of *ĥ 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 МАРИА suggestive of full nativisation: Laurentian Primary Chronicle *МѢРѢА*, *МѢРѢЮ*, Zadonshchina *МАРѢА*, *МАРѢА*, as well as First Novgorod Chronicle gen. sg. *ВАСИЛѢА* (jo-stem *ВАСИЛѢНН* < Βασίλειος).

Since we can’t ever be sure of the precise timing or route by which these late borrowings entered the various Slavic dialects, or of the extent of their adoption by Slavs beyond a tiny and often Greek-knowing scribal-class, the best solution is to set all such foreign /ij/ groups apart from the native vocabulary by using an *ij reconstruction, even where we can be pretty sure that early nativisation to reflexes of *ĭj occurred: *diĭĀvolъ, *vasilijъ, *marijĀ etc.

Word-initial *jъ-/*ji-/*i-

With native Slavic word-initial *ji-/*jъ-, I follow Derksen's (2009: 16) practice of writing *jъ-, even though Derksen himself (2003) has argued for a split between *ji- and *jъ- conditioned partly by accentological factors (which, as stated above, I have chosen not to consider). Most of the modern languages reflect these groups as just /i-/, except for Czech and Ukrainian: forms like Cz. *jdou* and Ukr. (after vowels) *йдуть* appear to have treated the weak-jer in *jъdŭť just like any other and retained the /j/, and Ukr. *ськати* <*jъskati (with the restricted meaning ‘look for nits/fleas in someone's hair’ after the base-meaning ‘seek’ was transferred to the Polonism *шукати*) shows the expected Ukr. softening of the /s/ after fallen weak-jer in *ъsk groups (cf. *польський*).

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*.

In badly-integrated clearly post-LCS foreign words, such as Biblical names like *ИАКОВЪ* (borrowed via Gk. *Ἰακώβ*), or *ІАЕМОНЪ* (< *ἡγεμών*), I keep a bare initial *i-, though this is rather an arbitrary choice and done partly as a way of marking such words as non-native¹³ (cf. my treatment of foreign initial *e- below). An exception is made for *ИСУСЪ* < Gk. *Ἰησοῦς*, which I have as *jisusъ, because of the greater likelihood that Slavs will have heard of Jesus even before the first biblical translations, and because spellings like Zogr. *ѡѡ Іѡѡ* suggest that it causes the same /Ŷ/ archiphoneme reflex of *ъ before *j as you get in e.g. native Mar. *ѡѡ Иѡѡѡ* < *vъ ѡ ѡѡstinŭ (see above).

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 *доїму* (<*dojъti with compensatorily-lengthened /o/ > /i/) vs *доїму*. 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 *šujika (*шюица*) and *vojinъ ‘warrior’ (*ѡѡѡнъ*, 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. *ѡѡѡѡѡѡ* <*žerbъji.

Word-initial *je-/*e-

No Glagolitic text makes any effort to distinguish /je/ (after vowels or word-initially) from post-consonantal /e/, writing both with <ѡ>, unlike the situation with the reflexes of *ję vs *ę, where in Zogr. and Mar. and partially in Assem. (Velcheva 1981: p.168) the full front-nasal digraph <ѡѡ> is

13 Spellings like Zogr. Mark 13:3 “*пѡѡѡѡ ѡ ѡѡѡѡ ѡ ѡѡѡѡ*.” “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 *jъjĀkovъ/ > /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.

reserved for *jē, while just the second ‘nasalising component’ <ɛ> is used for post-consonantal *e, e.g. Mar. 3rd pl. aorist **ѡѣсѣ** <*jesę, as opposed to KF **ꙗѣти** vs **ꙗѣли** <*prijeti vs *vъzeli¹⁴. Glagolitic evidence alone therefore would suggest that foreign borrowings with word-initial /e-/ were simply adapted to whatever the reflex of native LCS *je was. Suprasliensis, though, which uses the jotted <іє> letter, does in fact make an extremely consistent spelling distinction between foreign borrowings and native Slavic words: of the 157 occurrences of the 13 foreign lemmas I have so far reconstructed with word-initial *e/*je- which appear in Supr. (*episkupъ*, *evanġelъje*, *eġŭpъtъ*, *elisavetъ*, *elinъ*, *evanġelistъ*, *eġŭpъtъskъ*, *elinъskъ*, *episkupъstvo*, *evreјskъ*, *eliseјъ*, *етѣмаузъ*, *etiјоръskъ*), the only spellings with <іє> are ѡисеи, ѡппъ, ѡлини, and ѡлина, i.e. 4/157 or 2.5%. By contrast, of the 3172 native Slavic words in Suprasliensis which I Autoreconstruct as starting with *je- (not all of whose lemmas start with *je-, e.g. forms of *byti), just 88 are written with initial <ɛ>, vs 3070 with <іє>¹⁵. Thus 97.2% of native word-initial *je- in Suprasliensis is spelt with <іє>, while 97.5% of the occurrences of the clearly post-LCS Greek-mediated foreign borrowings listed above instead use plain <ɛ>, suggesting that some sort of difference was felt, at least by the scribes of Suprasliensis, and that we probably shouldn’t index these with the same *je- as used for native forms. I therefore use non-jotted *e- for such foreign borrowings, and the extent to which they take prothetic *j- and fall together with the native vocabulary is left as something for investigators to determine based on the evidence of each manuscript.

Prefixes

The last particularity of my LCS indexing-system worth mentioning relates to the handling of consonant-clusters in prefixes: as exhaustively exemplified by Diels (1963: 121-125), Common Slavic permitted only a restricted set of consonant-combinations in the syllable onset, generally either combinations of the continuants **s/*z* plus obstruent or sonorant (except **r*, see below), or of obstruents plus sonorant (with some curiosities such as the seeming tolerance of **bn* but not **pn*: OCS *ръивѣти* <**gyb-n̥ti* but *оучѣвати* <**usъr-n̥ti*, cf. 3sg. aor. *оучѣ*). Geminate consonants were banned and either simplified (*иѣшати* <**jъs-sekti*) or dissimilated (*процвѣсти* <**prokvit-ti*).

The ban on *sr/*zr is dealt with by insertion of *t and *d respectively, but the commonly-cited examples of *str <*sr (цестра, ерговѣ, остръ) all concern root-internal *sr where insertion of *t is common also to the Germanic and sometimes Baltic cognates. The examples given by Meillet (1965: 136) include: (for ерговѣ) Lith. dial. *srauja* next to Latvian *strauja*, then Germanic *straum- (> Eng. *stream*, Old Norse *straumr* etc.); (for остръ) Lith. *aštrus*, Gk. ἄκρος (here the *s is from PIE *k̑). As Meillet says, “*ce n' est pas un developpement germano-balto-slave ; d'une part, le developpement d'un -t- dans le groupe sr est chose naturelle et se retrouve ailleurs (fr. pop. castrole de casserole) et, d'autre part, le developpement de t en ces conditions n'est pas general en baltique: str est regulier en lette, mais sr subsiste couramment en lituanien.*”, so we can't really be sure when the Slavic change took place or whether it was still active during our LCS stage. The only indication of its activity in OCS is the single Psal. ~~срѣдѣ~~ <*sorm-omъ spelling cited by Diels (p. 122); otherwise new /sr/ from metathesised *sErC groups is tolerated unchanged.

New occurrences of *zr, on the other hand, are regularly generated in the language right up to OCS times, not only in the derivational-morphology because of the verb-prefixes *orz-, *vъz-, *jъz- (e.g. Supr. 3sg. aor. **ѡѡѡѡѡѡ** ‘roared’, from *vъz-ruti), but also because of the clitic prepositions *jъz and *bez, which form one phonological word with whatever follows them and thus cause OCS spellings like Mar. Luke 1 **ѡѡѡѡѡѡѡѡ** <*jъz ѡѡѡѡѡѡ. Meillet (p. 136) also cites the Old Polish adverb *zdręki* <*jъz ѡѡѡѡѡѡ, which proves that the phenomenon is not limited to SSL or OCS. Curiously, though, despite this overwhelming evidence of a synchronic /zr/ > /zdr/ rule in OCS, /zr/ from the

14 Psalterium Sinaiticum contains just six occurrences of non-digraph <ϣ>: **ⲡⲃⲁⲙⲓⲛⲁⲓⲧⲓⲕⲓⲙ**, **ⲙⲉⲧⲉⲛⲁⲓⲧⲓⲕⲓⲙ**, **ⲛⲁⲓⲧⲓⲕⲓⲙ**, **ⲛⲁⲓⲧⲓⲕⲓⲙ**, **ⲛⲁⲓⲧⲓⲕⲓⲙ**, and **ⲛⲁⲓⲧⲓⲕⲓⲙ**, according to Eckhoff's digitisation, five of which I've confirmed with Altbauer's (1971) facsimile of the manuscript. **ⲛⲁⲓⲧⲓⲕⲓⲙ** 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. *ймавѣ* which Eckhoff's corpus wrongly lemmatises as **jъmati* instead of **jъmĕti*, and which thus get reconstructed as **jemlĕvĕ* instead of **jъmavĕ*. 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.

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 *разрушить*, and there are traces of this even in Psalterium Sinaiticum: Psalm 48 𐌹𐌺𐌹𐌸𐌹𐌸𐌹𐌸𐌹𐌸 𐌹𐌺𐌹𐌸𐌹𐌸𐌹𐌸𐌹𐌸 (Diels 1963: 122). In Old Russian, the Uspenskij Sbornik is pretty consistent in keeping prefixed verb-forms like 𐌹𐌺𐌹𐌸𐌹𐌸𐌹𐌸𐌹𐌸 *<**orzrušitŭ, but by the time of the Laurentian Codex we get forms like 𐌹𐌺𐌹𐌸𐌹𐌸𐌹𐌸𐌹𐌸 and 𐌹𐌺𐌹𐌸𐌹𐌸𐌹𐌸𐌹𐌸.

For this reason I don't include /zdr/ <*zr at prefix or preposition-boundaries in my LCS system, so that investigators can see for themselves the extent of each text's adherence to the expected phonological development vs restoration of /zr/ under morphological pressure.

jъskěliti (Meillet 1965: 133)
обновити/оходити ошълъ

16 Conversely, sequences of *sk, *zg at prefix-boundaries which show PV1 reflexes, like Mar., Zogr. $\text{b}+\underline{\text{u}}\text{sa}\text{g}\text{a}\text{g}\text{a}\text{g}\text{a}\text{g}$ <*orš-čytetŕ (ECS *skit- > *ščit-), Psal. $\text{b}+\underline{\text{u}}\text{r}\text{g}\text{a}\text{g}\text{a}\text{g}\text{a}\text{g}$ <*orž-žigajetŕ (ECS *zgig- > *žžig-) are kept as *šč, *žž. Such forms may well not go all the way back to the time of PV1, and instead be just the result of a synchronic rule prohibiting /zž/ and /sč/ (> /žž/ and /šč/) that remained active until much more recently, especially given prepositional-phrase forms like Psal. $\text{p}+\underline{\text{u}}\text{sa}\text{g}\text{a}\text{g}\text{a}\text{g}$ <*js *červa, so this is arguably inconsistent with my treatment of *ss, *sš etc. My justification is that *sk, *zg > *šč, *žž are *conspicuously* PV1-changes, which we *know* originated well before our target LCS point, whereas de-gemination or simplification of *sš are less clear-cut.

case those of the original jo-stems). Consequently, a word like **ТАТЪ** has begun to take jo-stem endings, including the Old Russian /a/ reflex of LCS $*\bar{A}$ in the genitive/accusative singular. LCS $/\bar{A}/$, though, by definition can only occur after LCS palatal consonants (see above), so a reconstruction $*tat\bar{A}$ 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 $/\acute{r}\bar{A}/$ (palatal $/\acute{r}/ < *rj$ + the archiphoneme $/\bar{A}/$) to occur, then a morphological change which replaces the sequence $/ri/$ with $/\acute{r}\bar{A}/$ 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** $< *orza\acute{r}\bar{A}ti$ from the prefixed **разорити** $*orzoriti$ is irrelevant, because LCS $/orza\acute{r}\bar{A}ti/$ does not violate the rule of LCS phonotactics: palatal $/\acute{r}/$ can be followed by $/\bar{A}/$ because such a combination exists in the paradigms of wholly securely reconstructable jo-stem nouns, e.g. nt. gen. sg. $*mor\bar{A}$ ($>$ Pol. *morza*, OCS **морѣ**, Ru. *морѣ*, etc.)

In the case of Supr. Gsg. masc. **звѣръ**, for an original i-stem (**звѣри** $< *zv\acute{e}ri$), a direct LCS ancestor for the attested form can still be given ($*zv\acute{e}r\bar{A}$), because palatal $/\acute{r}/$ 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 $/\acute{r}/$, so that the nom. sg. $*zv\acute{e}r\bar{A}$ became $/zv\acute{e}r'/$, and its stem now ended on the same consonant $/r'/$ as original ja- and jo-stems ending on LCS $*\acute{r}$ 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 $*\acute{z}h\acute{u}q$, bastard Suprasliensis has PPP **заклатъ** (a noisome foulness), etc.
 -Could talk about the impossibility of dealing with $s\acute{m}\acute{e}\acute{s}\acute{e}$ deviances of S-aorist $s\acute{m}\acute{e}\acute{s}\acute{e}$ (vs. $s\acute{m}\acute{e}\acute{t}q$ $s\acute{m}\acute{e}\acute{t}o\acute{s}\acute{e}$), 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 **овсяный** OR adjective as another example of derivational-morphology made possible only by the rise of soft $/s'/$ (if it can be confirmed that the $*\acute{e}n\bar{A}$ adjectival suffix in e.g. OCS **оловѣнь** 'leadен' 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 *teí- 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 /ʎe/ 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.

- related is derivation-morphology difficulties such as whether adjective *volъnъ should have a palatal /j/, or whether the *volĀ is specifically differentiated from the root *vol- by a *-jĀ noun-forming suffix. Spellings are similarly suggestive of *volĀ-

-Talk about the pres. forms of *telhi and link back to the discussion about the difficulties with syllabic *ŕ, saying that Derksen and the two Czech dictionaries cite *tŕk 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 ⲡⲁⲣⲱⲛⲉⲛⲟⲓⲛ <*protolčenojo in Psal. Psalm 138

-The seeming impossibility of reconstructing aberrations like Supr. жласти, жладьба, from what can only be an original *geld- root and likely a Germanic loan (cf. 1X жлѣдетъ in the same text, or OR желести) are real barriers to

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 *-jъ 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.