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?)

F

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. <code>dusich</code> vs Suprasliensis. <code>Aoyctxz</code> <*duxěxs; Polish <code>wszak</code> vs Supr. <code>Bcakz</code>, Ru. <code>bcakz</code>, Ru. <code>bcak[uŭ]</code> <*vsx-aks), 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. <code>gwiazda</code>, Cz. <code>květ</code> <*gvězda, *květs, 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 свинеце <*svinькь and полотенеца <*poltsnbka), and those of *g are not unknown: Zaliznjak (2004: 47) admits that palatalised forms of the Germanic loan къназ- <*kъnęg- 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 оусьразн '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 (> Каz. қамыс).)



Table 2: LCS Vowels after the monophthongisation of diphthongs

	Fr	ont	Back			
High	i				у	u
	ŕьĺ			у	ŗъļ	
Mid	e	ę			Q	o
	ě	ð				
Low		Æ				
			a			

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

Labial		Dental		Palatal		Velar	
m		n		ń			
b	р	t	d	ħ	ħ	k	g
		s	Z	š ž		х	
				č			
		1		ĺ			
		r		ŕ			
7	I				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: /k g x f ü/, e.g. in respectively หำาว <*kitъ, เสดแดง <*igemonъ, xนางหว <*xitonъ, นงเน�ร <*ijosifъ³, and мvрง <*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 sourcelanguage (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': afprovs ១% ឬការ vs ១% មាន vs មាន vs មាន vs ១% មាន vs ១ ១៥ មាន vs ១៥ មាន vs ១៥ មាន vs ១៥ មាន vs ១% មាន vs ១% មាន vs ១៥ មា letter for $\frac{g}{g}$ (and the writing of $\frac{g}{g}$ with the palatalisation-diacritic) could be linked to the inadmissibility in the native systems of soft $[k^j, g^i]$ sounds, and whether their replacement with regular $\langle 9, 1 \rangle$ or $\langle r, \kappa \rangle$ was more likely in systems with some level of native $[k^j, g^j]$ (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. pory < *orb-q vs Uspenskij Sbornik nt. acc. sg. paro < *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 $/\xi/$ 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.



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

participle of certain verb-classes whose present-stem ends on a hard-consonant, which in South Slavic remains high and backed, e.g. Supr. zory, Psalterium Sinaiticum a吻a%經歷 <*stergy-jь, Codex Marianus дажя <*jĀdy-jь (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 <.e> for it), but which in most of North Slavic lowered to /a/: Old Polish (Kazania Świętokrzyskie) has both $rec\underline{a}$ and $rec\phi$ (with the special Old Polish letter for the merged reflex of *e and *o) <*reky, and in other texts also bior \$\phi\$ <*bery; Russkaja Pravda река, Uspenskij Sbornik донда <*dojьdy, Ru.Ch.Sl. (Vita Methodii) вьсемог<u>а</u>и <*vьхетоду-jь. Kortlandt's positing of a CS *y (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 *y is a convenient placeholder which allows all the relevant evidence to be retrieved. /ĕ/ is responsible for the NSl. /ĕ/ vs SSl. /e/ shapes of jo-stem masc. acc. pl. and the ja-stem nom./acc. pl. and gen. sg. endings, which are reflected in respectively the post- and prerevolutionary spellings of the Russian nom./acc. pl. long-adjective endings -ые < *vjě < *vjě vs -ыя < *yja < *yję < *yj**ę**.⁵

The need for the retention of the $/\bar{\mathbb{A}}/$ archiphoneme, which represents merged Early Common Slavonic * $\bar{\mathbb{A}}$ 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 $/\bar{\mathbb{A}}/$) was explicitly posited by Kortlandt as far back as 1979 (p.266) as part of his ECS system. In short, a combination of:

- 1.) the lack of any device in the Glagolitic alphabet to render /ja na ra la/ sequences (for which Glagolitic texts must use the jat' <a>> letter whose base-value is /ĕ/);
- 2.) overwhelming spellings of palatal-letter (<சயல்v>) + jat' in the Kiev Folia (the oldest and therefore least distant ms. from the 'original' OCS, as first codified by Cyrill and Methodius and for which the Glagolitic alphabet was devised) for the reflexes of LCS *č/š/ž/ħ + *Ā (e.g. அடிvard <*ob-věħĀlъ, கூயம் **e** dušĀmi), as well as occasional traces of such spellings in later Glagolitic OCS (e.g. Psal. சமய்ச் <*čĀšę); and
- 3.) the evidence of certain modern Bulgarian dialects, which have reflexes of LCS *ě in words like $\varkappa'e\delta a <*\check{z}\bar{\mathcal{A}}$ ba 'toad' (Stojkov 1954: 74–78),

all together point very strongly towards there having occurred a split on the Southeastern periphery of Slavic between LCS dialects which have $/\check{e}/<*\bar{\mathcal{A}}$ and the majority of the rest which got /a/, and that original OCS ('Urkirchenslavisch') was an $*\bar{\mathcal{A}}>/\check{e}/$ dialect. I posit that $*\bar{\mathcal{A}}$ remained until the opposition /a/: $/\check{e}/$ after palatal consonants was reintroduced when PV2 and PV3 brought new softconsonants /c ś dź/ into the system, which could be followed by both /a/ and $/\check{e}/$: fem. nom. sg. $/st_{\Delta}d\acute{z}a/<*st_{\Delta}g$ (PV3) vs fem. loc. sg. $/no\underline{d\acute{z}e}/<*nog\check{e}$ (PV2) (Winslow 2022: 304-305). Thus since my LCS system is based on a point just /a/ and PV3, I must also retain the /a/ archiphoneme.

⁶ It's possible to argue that the short *Æ counterpart to *Æ persisted in East Slavic until after the Fall of the Jers, and that the ESl. so-called e > o shift before hard-consonants / back-vowelled syllables is actually just the resolution of this archiphoneme as /o/ (where palatalisation of the preceding consonant remained, in e.g. Ukr. δ∂жοлα <*pьčÆla, or was newly phonemicised, in e.g. Ru. θёсла <*v'Æsla <*vesla), and that there was never a stage when these words had /e/ (based among other things on <o> spellings regardless of stress after palatal-letters in very early texts, and even after the letters for secondarily-soft LCS plain consonants in the Birchbark documents (Le Feuvre 1993, Nakonečnyj 1962), but there isn't space to elaborate on the issue here (see Winslow 2022: 304 fn.16). Unlike the situation with long *Æ, OCS shows no sign of anything but an /e/ reflex of short *Æ (and indeed the fact that



Other annoying pre-LCS morphological isoglosses reflected in the texts include the masc./nt. instr. sg. *o- and *jo-stem endings *-ьть (N.Sl.) and *-оть (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 *vьčera 'yesterday', and the *-tь (N.Sl.) vs *-tь (S.Sl.) verbal endings of 3rd sg. and pl. present (plus its extension to 2nd and 3rd sgl. aorists like OCS начать, OR (Uspenskij Sbornik) высть, начать). 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 $\frac{f}{h} \int_{0}^{h} \frac{1}{h} \frac{1}{h} dx$ are included as unitary vocalic phonemes, following Schenker (1995: 94), rather than as combinations of \sqrt{b} \sqrt{b} + \sqrt{f} \hat{l} r l, 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 /l/ f/ even in dialects with secondarily-palatalised consonants before fallen weak LCS /ь/). The need for both front and back $*\acute{r}$, *r is unambiguously shown by the East Slavic reflexes /er/ and /or/ (Ru. смерть, морковь), but *ĺ vs *l is more complicated: PIE *plnos, *wlk os > Lithuanian pilnas 'full', wilkas 'wolf' (LCS *pĺnъ, *vĺkъ) vs Lith. stulpas (LCS *stlpъ 'pillar') suggests that Balto-Slavic had differentiated front/back variants of the PIE syllabic *1 (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 (<*mĺčĀti), but the Polish reflexes are complicated and likely have more to do with the surrounding consonants: *pĺnъ by contrast gives pełny with hardened /l/ and the Polish non-palatalising-/e/ reflex of *ъ, and the differing reflexes in wi<u>erz</u>ch <*vŕхъ, śmi<u>er</u>ć <*sъmŕtь and m<u>ar</u>twy <*mŕtуъјь rule out any explanation based on the nature of the LCS syllabic-liquid alone (for more discussion see Bethin op. cit.: 73-75). 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 and <na>, 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 *ŕ *r and *ĺ as one would expect: ഉംഗം - ത്യംകം- ഉംകം - കംകം-<*r॔, ฉะเลย <*r॔, and ฐานาล สาวาร <*lí (Winslow 2022: 313), and even Zographensis spells all 5 occurences of *vĺk- 'wolf' with ษณษะ-/ษณอง- and all 15 instances of its *-mĺč- root with -รรมอง-(e.g. ஊகைச்சுயடி). 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 /br br bl /8.

Consonants

Dejotation

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 ebulanz of the Kipchak word /jovşan/ 'wormwood' in the Hypatian Codex, whose modern cognates (Turkmen /yowşan/jowşan/, Kazakh /zuwsan/, Azeri /zuwsan/, Azeri /zuwsan/, 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 */E.

- In the database I will have to use the single Unicode characters <ṛ ṭ ḷ ḹ>, rather than what's shown in my table, since the latter cannot actually be rendered without using the letters for /r ŕ l ĺ/ 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 /ē y/, which I will have to replace with <ē ȳ>.
- To my mind the only evidence in support of a genuine jer + liquid stage comes from the paradigms of verbs like OCS cathatu < *sutíti, where the syllabic /f/ in the stem alternates with /br/ depending on the vocality of the following morpheme: the e.g. 3sg. pres. *sutrety (Zogr., Supr. cathatus) or (one possibility of the) 3rd sg. aorist *sutre (Supr. cathatus) and the other possibility for the 3rd sg. aorist *sutí (Psal. Sin. 240rd), or with a different prefix Mar. 20rd *stif, being word-final or preconsonantal, must be syllabic /f/. The same alternation occurs in the zero-grade forms of verbs like *umerti, as is clear from the Polish reflexes umarł <*umr/lb vs umre <*umr/lo vs umre <*umr/lo vs umre = vumard = vu

Reflexes of the so-called jot-palatalisation are all written either as unitary palatal phonemes, or in the case of jot-palatalised labials as /vĺ mĺ bĺ pĺ/, rather than as sequences of consonant + /j/, hence /ń ĺ ŕ/ for *nj *lj *rj. The 'dejotated' reflexes of *tj (and *kt+front-vowel) and *dj are denoted using the modern Serbian Cyrillic letters /ħ/ and /ħ/ respectively, because the commonly used alternatives, i.e. /ť d/ (as used in e.g. Olander 2015) or /ḱ g/ (as used by me in Winslow 2022), or variations thereof, are visually too close to symbols used elsewhere in the system. /ḱ, g/ are anyway already used in my system for foreign /k, g/ before front-vowels, and /ť d/ 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 *ħ,ħ were close enough to foreign /g k/ before front-vowels that the original Glagolitic system used <¾ w> for both sets (i.e. alongside attested &ஃ೨೩೪೨-೨ < ἡγεμών would have been **೨೩೪೬-೩೨-೪ <*osoħeni, and alongside attested & ೪-೮-೬೪-೨ × «ἦνσος), 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.

Pre-dejotation *stj and *zdj are differentiated from the PV1 reflexes of *sk and *zg by writing the former as *šħ and *žħ and the latter as *šč and *žǯ, even though their modern reflexes do not differ from each other anywhere and so must've fallen together in the CS period, because they often alternate with their respective un-palatalised counterparts morphologically and derivationally, e.g. očistiti:očišħenьje vs. jьskati:jьščǫ, jĀzditi:jĀžħǫ vs jьzgъnati:jьžǯenǫ.

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 /ħħńĺŕ/ but not the PV2/3 reflexes /c ś dź/, is ahistorical. However it should be reemphasised 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 convicingly 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 /ħħńĺŕ/ 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.

<u>/j/</u>

Word-initial *j\(\bar{E}\)-/*a-

The tendency for ECS * \bar{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 * \bar{j} \bar{E} - in the absence of

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]ħe[tъ]; *ѓон* < Tk. *gön* 'leather', *меѓу* <*meђu.

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 ΔΦΦ- or IAM-, 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. In particular the 1sg. nom. pronoun *azъ/jĀzъ is especially problematic: I follow ESSJA I p.100 which ultimately plumps for *azъ, but Derksen doesn't discuss it at all. (A lengthy discussion of the evidence can be found in Teneva's (2012) article on the subject.)

Forms with insecure etymologies can't under any methodology be used as good evidence in phonological investigations, so in difficult cases like the above I simply mark the lemma in the database and provide some short discussion, so that eventually the web-interface can flag such forms in some way and inform users of the specific difficulties.

Like Derksen, I assume that roots going back to PIE jot-less long *ē or dipthongal *oi-, e.g. the root for 'to eat', PIE *h₁ēd, all took prothetic *j and merged with *jÆ- 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 κατμ etc.). Isolated nominal forms like Ru. *язва* (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. *examь*, Pol. *jeść* are later generalisations from prefixed forms like OR (ΣΝΝΈς ΤΗ, where no jot-prothesis could take place (cf. Schenker 1995: 88, Winslow 2022: 302 fn.14).

Jers before *i

As explored more fully in Winslow (2022: 313-315), OCS spellings seem to suggest that freevariation between

y> and <

yy> was a feature of the pre-Jer Shift Urkirchenslavisch orthographic system for conveying the reflexes of the sound-groups *ьj and *ъj (so-called 'tense jers'), regardless of whether they were in strong or weak position. The examples given were: Zogr. \mathscr{O} рьжэрэ Δ vs \mathscr{O} рьжэр Δ <*znamenь Δ , эаньэ Δ <*udaŕь- Δ vs эхэн Δ <*omočь- Δ ; Mar. эржает ту эржает ту «*osodet» јь; КБ шваэрт у с*milostьјо, ৬৯১৯৯৯৯ vs ৬৯১৯৯৯৯৯ < *vьхотоду-jь). Of importance here is the fact that the same orthographic system characterises both pre- (i.e. KF) and post-Jer Shift texts; that even in strong position in a text like Zographensis, which shows pretty clear signs of having undergone the Jer Shift, spellings like ஆக்கூர், மூல் செர் what in the live dialect underlying Zogr. must surely have been /udaŕij/, /boĺij/, /veštij/, are not infrequent¹⁰. The fact that the same alternation occurs in the pre-Jer Shift KF (i-stem gen. plurals & +гэчач + < *zapovědbjb vs аравя < *ludbjb) suggests that it is a common inheritance from the Urkirchenslavisch spelling system, and thus that in pre-Jer Shift Slavic the difference between /ь ъ/ and /i y/ was neutralised before /j/, and we should perhaps posit archiphonemes (which I call \hat{I} and \hat{Y}) in this position. These archiphonemes are, in slightly different terms, effectively posited by Trubetzkoy (1954: 70) in his analysis of the Urkirchenslavisch phoneme-system¹¹.

However, for simplicity and accessibility's sake it's better to avoid overburdening the indexing-system with unfamiliar and controversial archiphoneme-symbols, so I keep *ьj/*ъj as the denotations for these groups.

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 /Îj/: words like μαριία < Μαρία, «ΤΑΔΙΙΙΙ < στάδιον, 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 *marij $\bar{\mathcal{E}}$, *stadijь,

¹⁰ Marianus and Psalterium Sinaiticum, on the other hand, frequently show a Russian-style /ej/ reflex of strong tense *ьj: Psal. மூடியூர <*vorbьjь, ஈகுறூர <*plътьjь, ஜூருர் (அரும்கு) இது பார்க்கும் குறும் கூறும் கூறும்

¹¹ Though Trubetzkoy, like me, believes Urkirchenslavisch to have been based on a /j/-less dialect, so in that particular system the archiphonemes would be conditioned by the position before *vowels*, rather than before /j/.

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 примти <*prijeti or long-form adjectives like masc. nom. pl. дроузии <*drugi-ji, but within roots it does seem restricted to these post-LCS loanwords.

or as nativised *matjæ, *stadeje, but there are no occurrences of jer-spellings in these words in the OCS texts in TOROT. Other similarly-Greek words like μμαβολία (< διάβολος), however, do show up in OCS with jer-spellings: Supr. μειάβολα, Zogr. Luke 8 and Psal. Psalm 108 δ. ½ Δ. Ψοθέθ, which (alongside the modern Macedonian *fabon* with the reflex of *ħ 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 /Îj/. Old Russian texts even show spellings of марию suggestive of full nativisation: Laurentian Primary Chronicle ເພື່ອ με μα, (μομείθη, (jostem βαίη η μαριά, μαριά, αναριά, αναριά, αναριά, αναριά (jostem βαίη η μαριά (βοσίλειος).

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 *bj occurred: *dij\(\bar{E}\times\to l_b\), *vasilijb, *marij\(\bar{E}\) etc.

Word-initial *jь-/*ji-/*i-

With word-initial *ji-/*jь-, I follow Derksen's (2009: 16) practice of writing *jь-, even though Derksen himself (2003) has argued for a split between *ji- and *jь- 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 *jьdǫtь just like any other and retained the /j/, and Ukr. *съкати* <*jьskati (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 *ьsk groups (cf. *польський*). I make an exception for certain forms of the personal-pronoun *jь, however, and write *jimь, *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 μάκοβς (borrowed via Gk. Ἰακώβ), or ιπέμονς (< ήγεμών), 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-native (cf. my treatment of foreign initial *e- below). An exception is made for μέογες < Gk. Ἰησοῦς, because of the greater likelihood that Slavs will have heard of Jesus even before the first biblical translations, and because spellings like Zogr. ΚΊΙ ΙΚΎ suggest the same $/\hat{Y}/$ archiphoneme reflex of * \bar{z} before * \bar{z} as you get in native Mar. ΒΊΙ ΗΣΤΗΝΆ < *z υς \bar{z} jesting (see above).

Prefixed forms like *do-jьti '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 ∂iйmu (<*dojьti with compensatorily-lengthened /o/ > /i/) vs ∂oïmu. 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' (воинъ, as opposed to *vojьпъ, the gen. pl. of *vojьпа), or the loc. sg/pl. desinences of any jo-stem noun whose stem ends on /j/, e.g. Psal. ‰ьаш⊕т <*žerbъji.

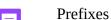
Word-initial *je-/*e-

No Glagolitic text makes any effort to distinguish /je/ (after vowels or word-initially) from post-consonantal /e/, writing both with <3>, unlike the situation with the reflexes of *je vs *e, where in Zogr. and Mar. and partially in Assem. (Velcheva 1981: p.168) the full front-nasal digraph <3e> is

¹³ Spellings like Zogr. Mark 13:3 "петух. и ъковъ. и одинъ." "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 *jъjækovъ/ > /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 *je, while just the second 'nasalising component' <€> is used for post-consontal *e, e.g. Mar. 3rd pl. aorist э€2€ <*<u>je</u>se, as opposed to KF ๛ฅจะ๛ <*prijeti vs ษองจะมฅ <*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, but Suprasliensis, with its jotated <> letter, does 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, evangelьje, egüpьtъ, elisavetь, elinъ, evangelista, egüpataska, elinaska, episkupastvo, evrejaska, eliseja, emamausa, etijopaska), the only spellings with <ค> are คงหนะห, เริกทร, เร็งหหม, and คงหม, 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 $<\epsilon>$, vs 3070 with $<\epsilon>$ ¹⁵. Thus 97.2% of native word-initial *je- in Suprasliensis is spelt with <€>, while 97.5% of the occurrences of the clearly post-LCS Greek-mediated foreign borrowings listed above instead use plain <e>, 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-jotated *e- for such foreign borrowings, and the extent to which they take prothetic *j- and fall together with the native vocabularly is left as something for investigators to determine based on the evidence of each manuscript.

孠



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 ¬ыбыжты <*gyb-noti but оусыжты <*usъp-noti, cf. 3sg. aor. оусыс). Geminate consonants were banned and either simplified (истышты <*jьs-sekti) or dissimilated (процвисты <*pre>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 (αετρα, ατρογια, ατρογια, ατρογια) all concern root-internal *sr where insertion of *t is common also to the Germanic and sometimes Baltic cognates: the examples give by Meillet (1965: 136) include: (for ατρογια) Lith. dial. srauja next to Latvian strauja, then Germanic *straum- (> Eng. stream, Old Norse straumr etc.); (for ατρα) Lith. aštrus, Gk. ἄκρος (here the *s is from PIE *k). 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. Δεσωμανών <*sorm-οπω spelling cited by Diels (p. 122); otherwise new /sr/ from metathesised *sErC groups is tolerated unchanged.

New occurences 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-, *jьz- (e.g. Supr. 3sg. aor. βτ3Δρογ 'roared', from *vъz-ruti), but because of the clitic prepositions *jьz and *bez, which form one phonological word with whatever follows them and thus cause OCS spellings like Mar. Luke 1 გრ გან *jъz _ro.kъ. Meillet (p. 136) also cites the Old Polish adverb zdręki <*jьz _ro.ky, which proves that the phenomenon is not limited to SSI. 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 <e>: ஙகணe, மூலிக்க்க், ஊக்க், ஊக்க், ஊக்க், ஊக்க், ஊக்க், ஊக்க்க், ஊக்க்க், ஊக்க்க், ஊக்க்க்க்க் digitisation, five of which I've confirmed with Altbauer's (1971) facsimile of the manuscript. ஈங்க்களை 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. и́мавъ which Eckhoff's corpus wrongly lemmatises as *jьmati instead of *jьměti, and which thus get reconstructed as *jemlevě instead of *jьmavě. 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 <βαρακ> and so seems to be tolerated, even though Diels cites prepositional forms like Supr. σεζαραζογμα, σεβαραλα 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 paspewumb, and there are traces of this even in Psalterium Sinaiticum: Psalm 48 ষ্ট্রেইলেস্ ২৯৯৩৯েস (Diels 1963: 122). The Old Rus. Uspenskij Sbornik is pretty consistent in keeping prefixed verb-forms like pasapywhtk <*orzrušitb, but by the time of the Laurentian Codex we get forms like kaspaayem and neh3peyennoe.

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

jьskěliti (Meillet 1965: 133) обновити/оходити ошьлъ



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 TATA from the 1229 Treaty between Smolensk, Riga and Gotland (Version A). LCS *tatь 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 (...тง кажетъ владънкы чловъколювые нако пръданника радбойника 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. τατι /tat'/) of this noun now ends on the same class of "soft" consonants as original jo-stem nouns like *pastyŕь > /pastyr'/, where the original LCS palatal *f has fallen together with secondarily-palatalised /r'/ from plain LCS *r before LCS front-vowels, in e.g. the original i-stem *zvěrь > /zvěr'/. This system thus no longer distinguishes between descendants of the original LCS palatals and the newly secondarilypalatalised 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 **TATL** has begun to take jo-stem endings, including the Old Russian /a/ reflex of LCS *\bar{E} in the genitive/accusative singular. LCS $/\bar{E}$, though, by definition can only occur after LCS palatal consonants (see above), so a reconstruction *tat $\bar{\mathcal{A}}$ 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 *mamю*), 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 **TATA**, 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 $/\dot{r}\bar{A}$ / (palatal $/\dot{r}/<*rj$ + the archiphoneme $/\bar{A}$ /) to occur, then a morphological change which replaces the sequence /ri/ with $/\dot{r}\bar{A}$ / 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 pagaphath <*orzaf\bar{E}ti from the prefixed pagophath *orzoriti is irrelevant, because LCS /orzaf\bar{E}ti/ does not violate the rule of LCS phonotactics: palatal /f/ can be followed by /\bar{E}/ because such a combination exists in the paradigms of wholly securely reconstructable jo-stem nouns, e.g. nt. gen. sg. *mof\bar{E} (> Pol. morza, OCS mopha, Ru. mopha, etc.)

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 *žь́nǫ, bastard Suprasliensis has PPP заклать (a noisome foulness), etc. -Could talk about the impossibility of dealing with sъmęšę deviances of S-aorist sъmęsę (vs. sъmętǫ sъmętošę), since unlike with nasal-stems, the deviance-slot here is taken up with the -oxaorists, leaving no room for deviantly-RUKI'd S-aorists
- -Could use the овьсяный OR adjective as another example of derivational-morphology made possible only be the rise of soft /s'/ (if it can be confirmed that the *ěnъ adjectival suffix in e.g. OCS оловънъ 'leaden' is LCS)
- -възлакати would be a good one to use to talk about the unmetathesised groups like *old-, *olketc., 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 /le/ desinence (i.e. with the consontant-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 *voĺьпъ should have a palatal /l/, or whether the *volÆ is specifically differentiated from the root *vol- by a *-jÆ nounforming suffix. Spellings are similarly suggestive of *vol-
- -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 *-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.