

Niko Partanen

University of Helsinki

ORCID: 0000-0001-8584-3880

Jack Rueter

University of Helsinki

ORCID: 0000-0002-3076-7929

Rogier Blokland

Uppsala University

ORCID: 0000-0003-4927-7185

Old Permic Universal Dependencies Treebank

Old Permic, also known as Old Komi, is an extinct variety of Komi that was spoken in the late Middle Ages in the lower Vychegda river basin in northeastern European Russia, in an area that currently is not Komi-speaking. This language variety is attested in fragmentary records from the 14th to 17th century written both in the Old Permic alphabet and in Cyrillic. These records are of significant importance for research on the history of the Komi language. Here we introduce our attempt towards a new Universal Dependencies treebank that will contain the existing corpus of Old Permic in a structured and CoNLL-U annotated format. This is the first time this material is being made available in digital format, and our contribution describes the current state of the art and remaining challenges.

Keywords: Uralic, Permic languages, Old Permic, Old Komi, CoNLL-U, Universal Dependencies

Introduction

Komi is a pluricentric Uralic language spoken in north-eastern European Russia. There are two main varieties with separate written standards, Permian Komi and Zyrian Komi¹, and a severely endangered less-documented variety called Yazva Komi. The largest variety, Zyrian Komi, has approximately 180,000 speakers. This paper focuses on the Zyrian Komi variety, which has existing written records from the 14th century onwards. The variety used in these early written records is usually known as Old Permic, and it consists of about 830 words (about 240 in the Old Permic script and nearly 600 in Cyrillic) of running text of religious nature; they are all translations of Church Slavonic texts (Baker 1983; Baker 1985: 24-29; Rédei 1993), and

¹The terms 'Komi-Permyak' and 'Komi-Zyrian' are also used; speakers of either variant will, however, usually just say they speak Komi. Hence, our terminology provides a clearer indication of the pluricentricity of the language.



Figure 1. The 14th century Troitsa icon with Komi text (Стефан Пермский, Public domain, via Wikimedia Commons)

The texts are all of a religious nature, and have existing parallels among the Russian recensions of Old Church Slavonic religious material. This is also important for the interpretation of the texts. The texts themselves are not always easy to read, and many characters are close to one another and in clear variation between different spellings of the individual words, but comparison with other variants is usually helpful for arriving at a specific reading.⁷

Writing system

The Old Permic script is read from left to right, and it consists of 38 characters. Five of the letters are combining. The texts, as is common for Russian writing from the period, are all in *scriptio continua*, i.e. word or sentence boundaries are not marked; only some wider spaces occur occasionally to denote a longer pause. The script was used between the 14th and 16th century, and its creation is usually attributed to Stephen of Perm (c. 1340–1396). There are two texts in Old Permic that can be directly dated; one to 1486 (cf. Grinščenko & Ponarjadov 2021: 12) and one to 1510 (cf. Lytkin 1952: 32). The exact time period when the script was used still remains a topic of further investigation, but can with certainty be said to have been from the end of the 14th century till the middle of the 16th century. This dating is supported by the linguistic features of Old Permic: as the divergence from modern Komi dialects is relatively minor, the script cannot be much older.

The Old Permic Unicode block was proposed to be included in Unicode in 2011, which was accepted in 2014, and currently there are several fonts that support these characters. An Android mobile keyboard has also been published⁸. The script is not actively used in modern communication, although small clusters of enthusiasts do exist at least in Russia and Finland.

The script has been compared to Cyrillic and Greek characters (Penttilä 1924; Stipa 1960), as well as, less convincingly, to e.g. the Old Turkic script (Ponarjadov 1996) and the Phoenician alphabet (Turkin 1996), and the possible influence of Komi traditional written signs, known as *pas* '(property) sign', has also been extensively discussed (cf. Korolev & Savel'eva 1996), but to the authors' knowledge as yet no thorough comparative study of how the different scripts and signs exactly relate to each other at the character level, and which influence is visible where, has been conducted.

⁷Penttilä (1924: 37) points out that copies of the same text often differ significantly in terms of graphemic legibility.

⁸<https://play.google.com/store/apps/details?id=com.majbyr.keyboard>

The Old Permic script was used to write the first texts in Old Komi, but later, after knowledge of the script declined, Old Komi texts were transliterated into Cyrillic; there are two extant texts in Cyrillic totalling almost 600 words, and there is therefore much more material in Old Komi in Cyrillic than in the Old Permic script. These texts in Cyrillic are also included as they clearly represent the same language form.⁹ Inclusion of texts in the Old Permic corpus is essentially done based on language, not the script, although the Old Permic script is very distinctive for this language variety.

To illustrate this, the text known as the inscription of Vasyuka Kyldashev in Nomokanon, dated to 1510, can be represented as follows in Unicode:

```
VIPI 7LL ɔLP ɓVVVPI ɣ7 VILVLRIL VIIVPI<br/> ɔVPLZ VIVVIL ɔZRPIL ʒVIV<br/>
ɔLV<br/>VL<br/>RV<br/>P
```

To our knowledge no one has yet published all Old Permic texts in the original script using the currently available Unicode encoding. We aim to complete this by the next Universal Dependencies release. At the same time, due to the difficulties in reading and interpreting the characters, it is clear this will not be the final word on the matter.

Corpus

Texts are included in Old Permic script or Cyrillic alphabet, depending on the text; for texts in Old Permic also Lytkin's 1952 Cyrillic transliteration will be included. All texts found before the 1950s have been exhaustively described in Lytkin 1952; since then there have been two major discoveries: a text with 18 words published in Sidorov in 1962, and a number of texts with a total of 24 words published in Grinščenko & Ponarjadov in 2021.

The word form corresponds to the string that is present in the original sources. The word tokenization follows Lytkin's conventions. We agree that it could be possible to also represent each text without word boundaries, following the *scriptio continua* of the original texts, but as we operate with a treebank structure, divisions into words and sentences are essential, and in any case we are working with an interpretation of the text as we are publishing it in a highly complex annotated structure. The original line boundaries are marked as `</>` in the word form and lemma columns, and in the MISC-column the feature value is represented by `<NewLine=Yes>`. This follows the conventions of other treebanks of historical texts, e.g. as in Old East Slavic Birchbark.

The goal of these conventions is that the corpus user should be able to retrieve a digital representation of the entire original text from the material in a condition where they can read in the original sources, and represented in the characters that can be considered as correct or plausible interpretations. Combined with the information about missing or destroyed characters this gives a good starting point also for situations where the text becomes more readable. Lately

⁹The cryptographic texts in Russian in the Old Permic script are obviously not included.

there have been great advances in image restoration using e.g. image processing algorithms (Knox et al. 2008) or x-ray imaging and AI (Parsons et al. 2023), so one wonders if unreadable parts of Old Permic texts could also eventually be retrieved with modern technology. Even in parts where the text on icons has been removed, it is worth pondering if the pigments used have left some detectable traces in the wood layer beneath. In a wider context we believe our annotation format serves as a good example in republishing electronically annotated versions of ancient texts without loss of information. The problem remains that we have only few texts in high resolution (the abovementioned Troitsa [Figure 1] and the texts in Grinščenko & Ponarjadov 2021), and no easy way to obtain high resolution images of the other ones.

In lemmatization the forms are given in the contemporary Zyrian Komi orthography. When the corresponding word does not exist in the standard language, but does occur in dialects, the form follows the correspondences deductible as presented in the recent Zyrian Komi dialect dictionary (Beznosikova et al. 2012). The dialect underlying the Old Permic corpus seems to be close to the contemporary Udora and Lower Vychegda dialects to which it has been geographically closest as well. As the current Komi dialect treebanks also contain lemmas in the standard Zyrian Komi (Partanen et al. 2018), this choice seems to be well founded for our purposes here.

Part-of-speech tagging follows the tagset in Universal Dependencies project, as used in the Permyak Komi and Zyrian Komi treebanks (Partanen et al. 2018), included in the UD release 2.5. (Zeman et al. 2019). In addition, the morphological analysis follows the UD conventions in the Komi and other Uralic language treebanks. The recent efforts to systematize the annotation conventions within Uralic languages in general (Partanen & Rueter 2019), and in more detailed questions such as numerals (Rueter et al. 2021), have been taken into account as closely as possible. Old Permic itself does not contain morphosyntactic features entirely absent from modern Komi dialects, or our current level of annotation at least has not encountered any. One morphological difference with regard to current Komi treebanks would be the seeming lack of some tense oppositions that would correspond to the contemporary present and future tenses. The same pattern found in Old Permic is also present in the modern Komi dialect of Udora (Partanen & Kellner 2021: 178). The current annotation conventions, however, are Tense=Pres Tense=Fut, although this may not be the optimal solution.

One difference between the Old Permic treebank and the other Zyrian Komi treebanks is that the morphological analysis produced by an FST is not included in the MISC column. The main reason is that there is no analyser of this type for Old Permic, and as the number of texts is very small and do not represent a currently spoken dialect, it does not seem necessary to integrate it either in the GiellaLT infrastructure that is mainly targeted toward language maintenance efforts.

The annotation of dependencies is currently in progress, and we have not yet encountered any annotation problems not known from other existing treebanks. Texts are, in the end, relatively straightforward Komi, although it has to be emphasized that if some words are to be read differently, this would also often result in a new syntactic structure.

As sentence boundaries are not indicated in the Old Permic texts, they have to be deduced by the context. This allows multiple interpretations, which in turn results in different tree structures in the annotations; we have mainly followed the interpretation of Lytkin 1952 here. The treebank is characterized by a fairly large number of lists and repetitions, which seems to follow from the religious nature of the material. There are often annotations with APPOS relation. In some cases the word order is not natural for Komi, and follows the original Church Slavonic source. The possessive constructions annotated with NMOD regularly have an opposite word order from the one expected based on other Komi treebanks.

Conclusion

We hope this work leads the way to further research on Old Permic. There is ample space for future critical editions and new analyses of these texts, especially as our knowledge of Komi dialects and other Permic varieties increases. The language resource presented in this paper concentrates primarily on making these materials available in a new digital form, with basic linguistic analysis that is complete, though not exhaustive.

The complete treebank will be made available in the Universal Dependencies release (v2.14), which is scheduled for May 15, 2024 (data freeze on May 1).¹⁰

Further work that could complement the current landscape of Komi treebanks would include different early Komi texts spanning newer periods of the written Komi tradition (from the 18th century onwards). Having more openly available materials on Komi dialects would be important. There is a discontinuity between the Old Permic script and later orthographies used for Komi, but through the contemporary dialectal variation these registers are still interwoven, despite the interruption in the written record.

References

Baker, R. W. (1983). Slavonic influence upon the language of the Old Permian texts. *Finnisch-Ugrische Forschungen* 45: 82-106.

Baker, R. W. (1985). The development of the Komi case system. *Suomalais-Ugrilaisen Seuran toimituksia* 189. Helsinki: Suomalais-Ugrilainen Seura.

Baraksanov, G. G. (1992). *Važ komi gižöd"äs (Важ коми гижӧдъяс)*. – Syktyvkar.

Beznosikova, L. M., Zaboeva, N. K., Ajbabina, E. A. & Kosnyreva, R. I. (2012). *Коми сёрнисикас кывчукӧр I-II. Словарь диалектов коми языка I-II. Сыктывкар: ООО Издательство «Кола».*

Derin, M. O., & Harada, T. (2021). Universal Dependencies for Old Turkish. In *Proceedings of*

¹⁰https://github.com/UniversalDependencies/UD_Komi-OldPermic/tree/dev

the Fifth Workshop on Universal Dependencies (UDW, SyntaxFest 2021) (pp. 129-141).

Grinščenko, A. I. & Ponarjadov, V. V. (2021). Новые находки памятников древнепермского языка и письма. Урало-алтайские исследования 4(43) (pp. 7-34).

Knox, K. & Easton, R. & Christens-Barry, W. (2008). Image restoration of damaged or erased manuscripts. 16th European Signal Processing Conference (EUSIPCO 2008), Lausanne, Switzerland, August 25-29.

Korolev, K.S. & Savel'eva, É.A. (1996). К проблеме происхождения коми письменности. Бараксанов, Г.Г., Тираспольский, Г. И. & Напалков, А.Д. (отв. ред.). Стефан Пермский и современность. Сыктывкар: Российская академия наук. Уральское отделение. Коми научный центр. (pp. 24-29)

Ljašev, V.A. (1980). Диалектное членение древнекоми языка. Серия препринтов. "Научные доклады". Выпуск 60. Сыктывкар: Академия Наук СССР, Коми филиал.

Lytvynenko, V. V. & Griščenko, A.I. (2022). Unlocking two marginalia in Old Permic script in a fifteenth-century Slavonic manuscript (Russian State Library, Volok. 437) with Athanasius' *Orations Against the Arians*. *Byzantinoslavica, Revue internationale des études byzantines* 80/1-2: 146-162.

Parsons, S., Parker, C. S., Chapman, C., Hayashida, M., & Seales, W. B. (2023). EduceLab-Scrolls: Verifiable Recovery of Text from Herculaneum Papyri using X-ray CT. arXiv preprint arXiv:2304.02084.

Partanen, N., Blokland, R., Lim, K., Poibeau, T., & Rießler, M. (2018). The first Komi-Zyrian universal dependencies treebanks. In *Second Workshop on Universal Dependencies (UDW 2018)*, November 2018, Brussels, Belgium (pp. 126-132).

Partanen, N., & Kellner, A. (2021). On the interplay between tense marking, aspect and temporal continuity in Udora Komi. *Finnisch-Ugrische Forschungen*.

Partanen, N., & Rueter, J. (2019). Survey of Uralic Universal Dependencies development. In *Proceedings of the Third Workshop on Universal Dependencies (UDW, SyntaxFest 2019)* (pp. 78-86).

Penttilä, A. (1924). Muutamia paleografisia huomioita Pyhän Tapanin syrjäniläisestä kirjaimistosta. *Turun historiallinen arkisto* 1 (pp. 32-45).

Ponarjadov, V.V. (1996). О возможной связи стефановской письменности с древнетюркской руникой. Савельева, Е.А. (отв.ред.). Христианизация коми края и ее роль в развитии государственной культуры. Том II. Филология. Этнология. Сыктывкар: Коми научный центр УрО Российской академии наук. (pp. 238-241).

Rédei, K. (1993). Óegyházi szláv szemantikai és szintaktikai hatás az ózürjén nyelvben. Budapest: Akadémiai Kiadó.

Rueter, J., Partanen, N., & Pirinen, T. A. (2021). Numerals and what counts. In Fifth Workshop on Universal Dependencies. The Association for Computational Linguistics.

Sidorov, A. S. (1962). Новые памятники древнекоми письменности (С комментариями, подстрочными примечаниями и заключением В. И. Лыткина). Вопросы финно-угорского языкознания. Вып. 1. Moskva. (pp. 178-211)

Stipa, G. 1960. Der Ursprung der permischen Schrift. Zeitschrift der Deutschen Morgenländischen Gesellschaft. 110 (Neue Folge 35) (pp. 342-364)

Terent'ev, Semën (compl) (2011a). ЗЫТЛНЛТТНУ ҺГҮҮН-ҮТЛНЦДЛН ҺЛЪОПГ / Antologiâ finno-ugorskoj poëzii. – Syktyvkar.

Terent'ev, Semën (2011b). ЛМЗН ???Ү ҮҮѠЗГ. Omar Khajam. *Rubaj*. Syktyvkar. [Note: The optimal representation of this source in the Old Permic script is still being revised by the authors.]

Turkin, A. 1996. Некоторые размышления о происхождении древнепермской письменности. Савельева, Е.А. (отв.ред.). Христианизация коми края и ее роль в развитии государственной культуры. Том II. Филология. Этнология. Сыктывкар: Коми научный центр УрО Российской академии наук. (pp. 278-281)

Zeman, D.; et al. (2019) Universal Dependencies 2.5, LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University, <http://hdl.handle.net/11234/1-3105>.