Krzysztof Nowak, Dorota Mika, Wojciech Łukasik
**TEI Modelling of Lexicographic Data in the DARIAH-PL Project**

The main goal of the "DARIAH-PL Digital Research Infrastructure for the Arts and Humanities" project is building the Dariah.lab infrastructure, which would allow for sharing and integrated access to digital resources and data from various fields of the humanities and arts. Among numerous tasks that the Institute of Polish Language, Polish Academy of Sciences coordinates, we are working towards the integration of our lexicographic data with the LLOD resources (Chiarcos et al. 2012). The essential step of this task is to convert the raw text of a dozen of paper-born dictionaries into TEI-compliant XML format (TEI Consortium).

In this paper we would like to outline the main issues involved in TEI XML modelling of these heterogeneous lexicographic data.

In the first part, we give a brief overview of the formal and content features of the dictionaries. For the most part, they are multi-volume works developed between the 1950s and 2010s with the research community in mind, and as such they are rich in information and structurally complex. They cover diachronic development (from medieval Polish and Latin to present-day Polish) and the functional variation of Polish (general language vs. dialects, proper names).



*Dictionary of Old Polish*
*Personal Names*



*Dictionary of Old Polish*

On a practical level, this means that, first, substantial effort had to be put into optimizing the quality of the OCR output. Since, except for `grobid-dictionaries` (Khemakhem et al. 2018), there are no tools at the moment that would enable easy conversion of lexicographic data, the subsequent phase of structuring of dictionary text had to be applied on a *per resource* basis.

TEI XML annotation has three main goals. First, it is a means of preserving the textuality of dictionaries which make heavy use of formatting conventions to convey information and employ a complex system of text-based internal cross-references. Second, TEI modelling aims at a better

understanding of each resource and its explicit description. The analysis is performed by lexicographers who may, however, come from a lexicographic tradition different from the one embodied in a particular dictionary, and thus need to make their interpretation of the dictionary text explicit. Regardless, this way we may also detect and correct editorial inconsistencies, which are natural for collective works developed over many years. Third, the annotated text is meant to be the input used in the alignment and linking tasks, it is therefore crucial that functionally equivalent structures are annotated in a systematic and coherent way. As we plan to provide an integrated access to the dictionaries, the TEI XML representation is also where the first phase of data reconciliation takes place. It does not only concern the structural units of a typical dictionary entry, such as `<sense/>` or `<form/>`, but also mapping between units of analytical language the dictionaries employ, such as labels, bibliographic reference system etc.

**References**

Christian Chiarcos, Sebastian Hellmann and Sebastian Nordhoff. 2012. Linking linguistic resources: Examples from the Open Linguistics Working Group, In: Christian Chiarcos, Sebastian Nordhoff and Sebastian Hellmann (eds.), Linked Data in Linguistics. Representing Language Data and Metadata, Springer, Heidelberg, p. 201-216.

Mohamed Khemakhem, Axel Herold, Laurent Romary. 2018. Enhancing Usability for Automatically Structuring Digitised Dictionaries. In: GLOBALEX workshop at LREC 2018, May 2018, Miyazaki, Japan. 2018.

TEI Consortium, eds. "9 Dictionaries." *TEI P5: Guidelines for Electronic Text Encoding and Interchange.* [Version number]. [Last modified date]. TEI Consortium. https://tei-c.org/release/doc/tei-p5-doc/en/html/DI.html (10 June 2022).