**Silga, Sviķe,** *Dr. philol.*
**Ventspils University of Applied Sciences**

# Report
# on ELEXIS Transnational Research Visit Grant at the Austrian Centre for Digital Humanities of the Austrian Academy of Sciences (ACDH-OeAW).
# (Vienna, Austria March 16-20 2020)

**Travel Grant: Call 3**

**Project title: "German-Latvian LSP Glossary of Kawall's "Dieva radījumi pasaulē" and its Original Work"**

**Introduction**

My visiting grants project proposal is part of a larger project aiming to research Latvian botanical terminology used in H. Kawall's work "God's Creatures in the World" ("Dieva radījumi pasaulē"). This work (a textbook) is one of the first translations from German into Latvian, in which the author mentions Latvian botanical terms for the first time ever, in addition to the terms of zoology and mineralogy, the book has a separate chapter – Plant Kingdom (Augu valsts). A detailed research of botanical lexis requires a digital corpus of language material on which to compare and study special lexis used in the original language and translation. Therefore, the **goal** of the research stay is to create a bilingual digital LSP corpus based on the original book in German and its translation into Latvian, along with the aim to compile a bilingual LSP glossary that includes a collection of special botanical vocabulary used in H. Kawall's translation and original work.

In this report, I will present the workflow of my research visit, i.e. preparatory work,

used books, support and main work at the ACDH-OeAW, materials studied and tools used to the achieve the goal – create two corpora and a glossary. I will summarize the conclusions and single out some possible solutions for the further research process.

## Workflow and description of steps for performing specific tasks

**Conducting preparatory work before the research visit.**
Preparation of two printed books: H. Kawall "Deewa raddijumi pasaulē" („Dieva radījumi pasaulē" (DRP)) (1860) and "Die Naturgeschichte für Kinder und Elementarschüler, oder erster Unterricht über das Mineralreich, Pflanzenreich und Tierreich, mit über 300 kolorierten Abbildungen" (1855). Scanning and saving both books in PDF format.

**Digitising paper books (OCR scans).**
For solving theoretical issues, finding tools and methods of digitisation of both scanned books the support of the National Library of Latvia (NLL) was sought. The scanned documents had to be prepared in OCR (optical character recognition) format, i.e. in such a format to make it possible to find, edit and process specific fragments of the text using the search function. The documents could not be pictures, such as jpg format. Doc-Works programme environment was used to process texts by working with the scanned material. The texts of both books are written in the old script, and it was the greatest challenge of the project. The new supervised machine learning approach offered by NLL was used to digitise the texts, and these two books were the first ones to be processed like this. Initially the text was recognised by an untrained algorithm. Next, the text was edited.

For the computer to recognise text accurately, a sample of the correct representation of the text must be made first – at least 10 000 perfect, human-edited lines. Each line must be checked by two writers editing individual lines. You can do it and view the process on the NLL website: https://frakturs.lnb.lv/ (see Fig. 1)

**Figure 1. Editing Gothic texts manually**

To enable automatic recognition, the Tesseract text recognition software was used. It works using the LSTM (long short-term memory) neural network model. LSTM operates more accurately than the early neural network models, and it is well-suited for script and speech recognition. LSTM belongs to the deep learning algorithms. The Tesseract software is a completely new solution of the year 2019, and it was used in this project as an experiment. This program used the German text model as the basic model, where the letters have diacritical marks.

**During research visit.**

Participation in the virtual meeting with deputy head of the ACDH-OeAW, Dr. Karlheinz Mörth, who shared expertise regarding best practice and standards in lexicography (see Figure 2). During the virtual meeting there was a possibility follow insights into different projects carried out in the ACDH-OeAW e.g. Vienna Corpus of Arabic varieties (VICAV). It was a good opportunity to read the study "Best practices for lexicography – intermediate report" by C. Tiberius, R. Costa, T. Erjavec, S. Krek, J. McCrae, C. Roche, T. Tasovac, 31 January 2020 (available at: http://ejuz.lv/estpractices).

**Figure 2. Virtual Meeting with Dr. Karlheinz Mörth**

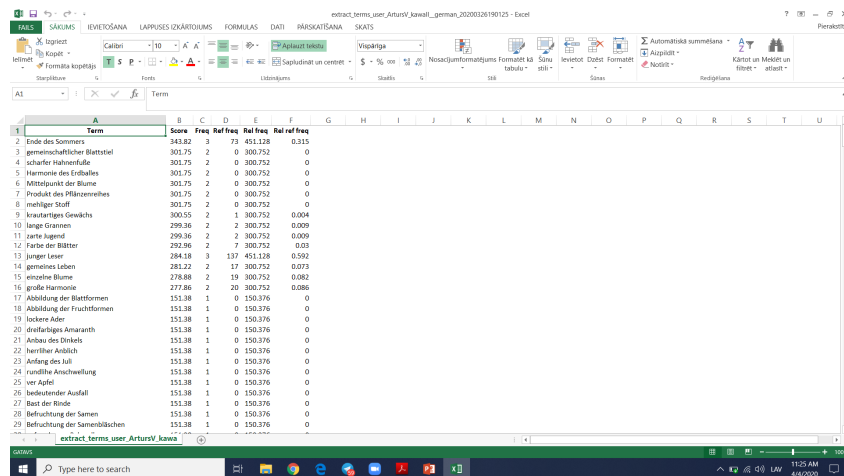**Aligning the original work and its translation.**

Aligning was performed by taking a sample from Chapter 2 of both books and manually copying the original work and translation segments into Excel Sheets (see Figure 3), as Excel is one of the formats required for future work with Sketch Engine. As the German text was relatively erroneous and required a lot of manual editing work, only a small part of the text could be processed in this way during the research visit.



**Figure 3. Aligned texts in Excel Sheets**

**Extracting data for a bilingual glossary.**

To obtain data for the glossary, the Sketch Engine software was used. Texts in Excel tables were first uploaded there, and two text corpora were created from those – German and Latvian language corpus. Using the functionalities of Sketch Engine extractions, keyword lists and term lists with detailed information were created, e.g. score, frequency (see Figure 4 for German). They are intended for future use in creating glossaries and in terminology research.



**Figure 4. Extracted Terms from German corpus**

**Short conclusions and future prospects after research visit.**

To sum up the experience, it must be concluded that the prepared text material plays a significant role. In this case the text was in old script in two languages (German and Latvian). The greatest challenge of this research was to successfully prepare the old script material, and it needs to be further developed. Another objective would be to research the available software for editing German Gothic script, as manual editing work is time-consuming. The result of the research visit are digitised versions of two sections of printed books mentioned above and extracted keyword and term lists for the text samples. The extracted material will be further analysed by frequency data, as well as clarifying the most popular word collocations. In this way it will be possible to identify specific translation models and strategies in terminology translation developed by H. Kawall. The extracted Latvian botanical terms used in H. Kawall's

work will also be analysed in relation to the contemporary botanical terminology, and the results of the research in the form of an article will be included in the multilingual "New Botanical Dictionary" (a mobile app prototype), developed at Ventspils University of Applied Sciences.