

Scientific Report of Short Term Scientific M

COST STSM Reference Number: COST-STSM-IS1305- 26615

Period: 20-04-2015 to 24-04-2015

Duration: 5 working days.

COST Action: IS1305

STSM type: Regular (from FINLAND to ESTONIA)

STSM Title: Development of Sketch Grammar and GDEX (Good Dictionary E

Guest/STSM applicant: Tarja Heinonen, Institute for the Languages of Finla

Host: Jelena Kallas, Institute of the Estonian Language, jelena.kallas@eki.ee

1. Purpose of the STSM

The aim of the STSM was to learn about the corpus query tool Sketch Engine and to write grammatical rules that operate on a morphologically preanalyzed corpus. The full potential of the Sketch Engine can be put to use. The goal was to produce sketches, which are one-page, corpus-derived summaries of a word's grammatical behaviour (ibid.) for a representative selection of Finnish headwords and compare them with the current entries in the Dictionary of Contemporary Finnish. The ultimate goal was to find new and more efficient ways to update the dictionary. An additional purpose was to outline preliminary good example sentence specifications for the automatic generation of example sentences (Kilgarriff et al., 2008; Kosem et al., 2011, Kosem et al., 2013). The idea behind this was to make a lexicographer's work more efficient by suggesting a small number of example sentences for the lexicographer to choose from, instead of full concordances. At the end of the STSM there was an opportunity to participate in the EAAL (Estonian Association of Applied Ling

Finnish Word Sketch Grammar is geared towards the specification of the sketch engine (available at <https://the.sketchengine.co.uk/auth/corpora/>).

Day 3: Grammar writing, Word Sketches and other functionalities of the Sketch Engine for Finnish and Estonian Word Sketches.

Day 4: Preliminary introduction to GDEX with a quick overview on how the sketch engine works for the Estonian Collocation Dictionary (Kallas & Tuulik, 2014). We also discussed the dictionary with Iztok Kosem. On top of that, I learned about a forthcoming dictionary of Estonian from chief editor Margit Langemets. The Estonian dictionary will be corpus-based while ours on Finnish is going through only minor updates based on new sources.

Day 5: Defining parameters for Finnish GDEX and writing a preliminary GDEX configuration on the basis of the one used for the Estonian Collocation Dictionary. Changes were made to the conditions on word frequencies and the lists of the words to be excluded. The next step would be the next step. I also continued writing rules on grammatical features.

During my stay in Tallinn I had a chance to participate in the EAAL conference (invited speakers Simon Krek and Iztok Kosem). In this context, on day 4, I gave a presentation to the conference audience about the updating process of the Dictionary of Collocations and what kind of methods and sources we use.

As a result of the STSM, Finnish Sketch Grammar with about 70 rules was developed and the parameters for GDEX were set as well. The results are presented in the appendix of the Finnish Sketch Grammar; Appendix 2: GDEX configuration for Finnish).

3. References

Kallas, J. (2013). Eesti keele sisusõnade süntagmaatilised suhted korpuskeeleloomingus. [Syntagmatic relationships of Estonian content words in corpus lexicography.] Tallinn: Tallinn University. Dissertations on Humanities. http://e-ait.tlulib.ee/303/1/kallas_jelena.pdf.

toc/euralex_2008/.

Kilgarriff, A. (2013). Using Corpora and the Web as Data Sources for Dictionaries. In *The Bloomsbury Companion to Lexicography*. Bloomsbury, London. Chapter 10.

Kosem, I.; Husák, M. & McCarthy, D. (2011). GDEX for Slovene. In *Proceedings of the eLex 2011 conference, 17–19 October 2011, Tallinn*, 159. Available at: http://elex2011.trojina.si/Vsebine/proceedings/eLex2011_03_Kosem_Husak_McCarthy.pdf

Kosem, I., Gantar, P. & Krek, S. (2013). Automation of lexicographic work: from lexicographers and crowd-sourcing. In: I. Kosem, J. Kallas, P. Gantar, M. Tuulik (eds.) *Electronic lexicography in the 21st century: third proceedings of the eLex 2013 conference, 17–19 October 2013, Tallinn*. Available at: http://eki.ee/elex2013/proceedings/eLex2013_03_Kosem_Gantar_Krek.pdf