# Automatic Generation of Wiktionary Entries for Finno-Ugric Minority Languages

Zsanett Ferenczi, Iván Mittelholcz, Eszter Simon
Helsinki, 8 January 2018

Research Institute for Linguistics, Hungarian Academy of Sciences

**General objective:**

to provide linguistically based support for several small FU digital communities to generate online content and help to revitalize the digital functions of some FU minority languages

**Practical objective:**

to create bilingual dictionaries for six FU minority languages paired with four major languages and to deploy the enriched lexical material on the web in the framework of Wiktionary

- generating the translation units
    - creating proto-dictionaries → raw dictionaries
    - manual validation & preprocessing & consistency check → valid dictionaries = translation units
- generating the Wiktionary entries
    - providing POS tags
    - disambiguating the POS tags
    - adding IPA transcription
    - generating the final entries
- uploading the entries
- evaluation