

# Finno-Ugric Digital Natives – Linguistic Support for Finno-Ugric Digital Communities



Kristiina Jokinen<sup>1,2</sup>, Graham Wilcock<sup>1,3</sup>, Tamás Varadi<sup>4</sup>, Eszter Simon<sup>4</sup>

<sup>1</sup>University of Helsinki, <sup>2</sup>University of Tartu, <sup>3</sup>CDM Interact Oy

<sup>4</sup>Research Institute of Linguistics, Hungarian Academy of Sciences

## Collaboration project between Academy of Finland and Hungarian Academy of Sciences

- Support digital visibility and viability of minor Finno-Ugric languages in the current globalising world
- Develop tools and resources for automatic language processing
- Generate online content on the target languages
- Experiment with interactive applications

### Sami languages

- The nine Sami languages belong to the Finno-Ugric language family
- Spoken in the northern parts of Finland, Norway and Sweden and in the Kola Peninsula in Russia
- North Sami (*davvisámegiella*) is the biggest with ca 20 000 speakers
- Mainly morphophonetic variation, syntactic differences fairly small
- e.g. in Eastern dialects:
  - the vowel system contains three low vowels /ä:/a:/â/ Western dialects have only two /ä:/a/
  - diphthongs become open:
    - miessi /miessii/ 'fawn' : geassi /kiässii/ : /kiessii/ 'summer'
    - luohhti /luohhtii/ 'a yoik' : goahti /kuähtii/ : kuohtii/ 'teepee, trad.dwelling'
  - /k/ in between vowels disappears and /p/ changes into /v/
    - johka /joohka/ : joga /joo.a/ 'river'
    - lohpi /lohpii/ : lobi /loovii/ 'permission'



### Multimodal Interaction – SamiTalk

- Human-Robot Application for navigating information from Wikipedia
- Laughter and body posture recognition

Robot: *Man fáttá birra don háliidat gullat? Soames mu favorihta fáttát leat Sámeigiella, Roavvenjárga, Mohtorgielká, Tiina Sanila-Aikio, Manga.*  
 What topic do you want to hear about? Some of my favourites are: Sami language, Rovaniemi, Snowmobile, Tiina Sanila-Aikio, Manga.  
 User: *Sámeigiella* (Sami language)  
 Robot: *Sámeigiella dili Norgga, Ruota, Suoma ja Ruošša bealde lea earálgán.*  
 The Sami language situation is different in Norway, Sweden, Finland and Russia  
 User: *Dilli Suomas* (The situation in Finland.)  
 Robot: *Bealli Suoma 9350 sámis máhtet sámeigiela. Suomas hállujuw dawwi-sámeigiella, anárašgiella ja nuortalašgiella, main anárašgiella dušše Suomas. ...*

### Dictionary Creation

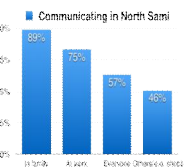
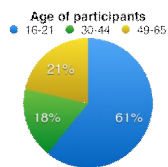
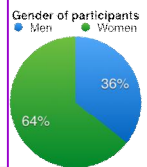
Create bilingual dictionaries for six small Finno-Ugric languages (Udmurt, Komi-Permyak, Komi-Zyrian, Hill Mari, Meadow Mari and Northern Saami) paired with four major languages which are important for these communities (English, Finnish, Hungarian, Russian)

method	RESULTS FOR THE METHODS		RESULTS FOR THE MERGED DICTIONARIES	
	all (#)	useful (%)	lang pair	useful (%)
WikiTitle	2,989	99.33	sme-eng	6,042 92.29
W2D ext	921	98.91	sme-fin	7,100 91.44
W2D tri	11,714	93.51	sme-hun	4,969 90.72
KDE4	8,401	89.03	sme-rus	4,373 95.95

### DigiSami Data Collection

- North Sami speech corpus: formal reading and free conversation
- Collected in 3 villages in Finland: Utsjoki, Inari and Ivalo, and in 2 villages in Norway: Kautokeino and Karasjok
- Participants bilingual (Finnish or Norwegian), most had lived their life in the Sami area
- Data transcribed and used to study spoken interaction (speaker identification, multimodal conversation analysis, laughing)
- Available: <https://goo.gl/SW9Goz>

Dialects	Read speech		Conversational speech	
	#Spk	Duration	#Spk	Duration
Kautokeino	4	1.03	-	-
Karasjoki	6	0.72	6(1)	1.5
Ivalo	6	0.72	7(1)	0.72
Utsjoki	5	1.07	6(1)	1.03
Inari	4	0.73	-	-
<b>Total</b>	<b>25</b>	<b>3.26</b>	<b>19</b>	<b>4.28</b>



- Jokinen, K., Hiovain, K., Laxström, N., Rauhala, I., Wilcock, G. DigiSami and Digital Natives: Interaction Technology for the North Sami language. *IWSDS*, 2016
- Jokinen, K. Open-domain Interaction and Online Content in the Sami Language. *LREC*, 2014
- Jokinen, K., Wilcock, G. Community-based Resource Building and Data Collection. *SLTU*, 2014
- Jokinen, K., NgoTrong, T., Hautamäki, V. Variation in Spoken North Sami Language, *Interspeech* 2016

- Benyeda, I., Koczka, P., Váradi, T. Creating seed lexicons for under-resourced languages. *GLOBALEX 2016*
- Simon, E., Benyeda, I., Koczka, P., Ludányi, Zs. Automatic creation of bilingual dictionaries for Finno-Ugric languages. *The 1st International Workshop on Computational Linguistics for Uralic Languages*, 2015
- Wilcock, G., Laxström, N., Leinonen, J., Smit, P., Kurimo, M., Jokinen, K. Towards SamiTalk: a Sami-speaking robot linked to Sami Wikipedia. *IWSDS*, 2016