# FIN-CLARIAH MEETING

December 1, Linna Building, Kalevantie 5, Tampere

### PRELIMINARY PROGRAM

11.00-11.10 Welcoming Words by Sanna Kumpulainen, Associate Professor in Information Studies, Tampere University

11.10-12.00 Keynote I on Studying SSH Research Needs: Elina Late, Senior Research Fellow in Information Studies, Tampere University

#### 12.00-13.00 Lunch

13.00-13.45 Keynote II on Language Models: Sampo Pyysalo, Associate Professor at the Department of Computing, University of Turku

#### 13.45-14.30 Work Package Presentations I

13.45-14.00 WP1.3 Veronika Laippala: Noise-Tolerant NLP

14.00-14.20 WP 1.1, 1.2, 2.1, 2.2, 2.3 Mietta Lennes: Kielipankki – The Language Bank of Finland

14.20-14.25 WP2.4 Harri Kettunen: Helsinki Term Bank for the Arts and Sciences

14.25-14.30 WP2.5 Jenny Tarvainen: Automated Text Tools for Learner Language

#### 14.30-15.00 Coffee

#### 15.00-16.00 Work Package Presentations II

15.00-15.05 WP3.1 Martin Matthiesen: Pipeline from the National Library to CSC 15.05-15.10 WP3.2 Tanja Välisalo: Named Entity Recognition for NARC Data 15.10-15.15 WP4.3 Eetu Mäkelä: Evaluation and Subsetting

15.15-15.20 W4.1 Julia Matveeva: Metadata Harmonization

15.20-15.25 WP4.4 Mikko Laitinen: Twitter

15.25-15.30 WP4.2 Eero Hyvönen: LOD

15.30-15.35 WP3.4 Raine Koskimaa: Game Streams

15.35-15.40 WP3.3 Maria Valaste: Qualitative Surveys

15.40-15.45 WP3.5 Kimmo Elo (Risto Turunen replacing): Text Networks

15.45-15.50 WP5 Sanna Kumpulainen: Evidence-based RI Development + Education & Resources

16.00-17.00 Free Chilling & Refreshments / Parallel session: Executive Board Meeting (with Zoom option)



### MATERIALS

The presentation slides will be made available after the event at <a href="https://www.kielipankki.fi/organization/fin-clariah/fin-clariah-2023-12-01/">https://www.kielipankki.fi/organization/fin-clariah/fin-clariah-2023-12-01/</a>

### PROJECT DELIVERABLES

https://www.kielipankki.fi/organization/fin-clariah/deliverables/

### WP1.3 NOISE-TOLERANT NLP, 15 MIN





## FIN-CLARIAH 2022-23:

Overview of the work completed in Modules 1 and 2

University of Tampere, 1.12.2023

# https://www.kielipankki.fi

#### LANGUAGE BANK ACCESS CORPORA TOOLS ORGANIZATION SUPPORT

0

SUOMEKSI PÅ SVENSKA



KIELIPANKKI

B

Apply for rights to use our language resources.



The Language Bank of Finland

Browse our corpora.



Try our tools.



Who are the Language Bank?





Researcher of the Month: Aleksi Sahala

#### News

 New resource (beta): Christmas Gospel text-tospeech in four Uralic languages, Korp (29,11,2023)



Help and instructions.



Books and manuscript s

Newspapers and journals

Online discussions and other data from the Internet (e.g., Suomi 24, Ylilauta) Sessions of the Parliament (videos + transcripts)

Parallel corpora and other multilingual resources

> Lexicons and terminologie s

Learner language corpora

Sign languages

**Dialect corpora** 

Lahjoita puhetta – Donate Speech (puhelahjat)

See the complete list of corpora at https://www.kielipankki.fi/corpora/

## https://www.kielipankki.fi/corpora

					Etsi:		
Abbreviation	Name and metadata	License	Apply	Location 🕈	Service level 🗢	Help 🗘	Cite ≑
acquis-ftb3	The Finnish Sub-corpus of the JRC-Acquis Multilingual Parallel Corpus	PUB		Korp	В	0	99
acquis-ftb3-dl	Finnish Sub-corpus of the JRC-Acquis Multilingual Parallel Corpus, Downloadable Version	PUB		Download	В	0	<b>55</b>
agricola-v1-1-korp	The Morpho-Syntactic Database of Mikael Agricola's Works version 1.1, Korp	PUB		Korp	В	0	<b>55</b>
ai2d-rst-v1-1	AI2D-RST: A multimodal corpus of 1000 primary school science diagrams version 1.1	PUB		Download	В	0	<b>99</b>
aku-egg-dl	Speech and EGG (Electroglottography) Simultaneous Recordings, downloadable version	ACA		Download	В	0	<b>77</b>
amph	amph-Corpus	ACA	<b>+</b> )	Download	В	0	<b>77</b>
ArkiSyn-korp	ArkiSyn Database of Finnish Conversational Discourse, Helsinki Korp Version	PUB		Korp	В	0	99
AVOID	Corpus of Age-related Voice Disguise (AVOID)	RES	*)	Download	В	0	<b>77</b>
BeserCorp	The Corpus of Beserman Udmurt	PUB		Korp	В	0	<b>99</b>
ccmh-src	Corpus Cyrillo-Methodianum Helsingiense: Corpus of Old Church Slavonic Texts, source	PUB		Download	В	0	<b>77</b>
ceal-dl	The Downloadable Version of Classics of English and American Literature in Finnish	RES		Download	А	0	<b>9</b> 9
ceal-o	Classics of English and American Literature in Finnish, Sentences and Paragraphs in the Original Order	RES	۲	Korp	А	0	<b>55</b>
ceal-par-korp	Classics of English and American Literature as translated by Kersti Juva, English-Finnish parallel corpus, Korp	RES		Korp	А	0	<b>99</b>
coal par c dl	The Downloadable Version of Classics of English and American Literature as translated by Kersti Juva, English-Finnish parallel	ACA		Download		0	

## https://www.kielipankki.fi/corpora

					Etsi:		
Abbreviation	Name and metadata	License	Apply :	• Location •	Service level 🗢	Help 🗘	Cite 🗢
acquis-ftb3	The Finnish Sub-corpus of the JRC-Acquis Multilingual Parallel Corpus	PUB		Korp	В	0	99
acquis-ftb3-dl	Finnish Sub-corrector the Acquis Multilingual Parallel Corpus, Downloadable Version	PUB		Download	В	0	99
agricola-v1-1-korp	The Morpho-Syntactic Database of Mikael Agricola's Works version 1.1, Korp	PUB		Korp	В	0	99
ai2d-rst-v1-1	AI2D-RST: A multimodal corpus of 1000 primary school science diagrams version 1.1	PUB		Download	В	0	<b>99</b>
aku-egg-dl	Speech and EGG (Electroglottography) Simultaneous Recordings, downloadable version	ACA		Download	в	0	99
amph	amph-Corpus link to access location	ACA	+)	Download	В	0	<b>99</b>
ArkiSyn-korp	ArkiSyn Database of Finnish Conversational Discourse, Helsinki Korp Version	PUB		Korp		0	99
AVOID	Corpus of Age-related Voice Disguise (AVOID)	RES	•)	Download	в	0	99
BeserCorp	The Corpus of Beserman Udmurt	PUB		Korp	В	0	<b>99</b>
ccmh-src	Corpus Cyrillo-Methodianum Helsingiense: Corpus of Old Church Slavonic Texts, source	PUB		Download	В		<b>9</b> 9
ceal-dl	The Downloadable Version of Classics of English and American Literature in Finnish	RES		Download	А	0	<b>99</b>
ceal-o	Classics of English and American Literature in Finnish, Sentences and Paragraphs in the Original Order	RES	•	Korp	А	0	<b>99</b>
ceal-par-korp	Classics of English and American Literature as translated by Kersti Juva, English-Finnish parallel corpus, Korp	RES		Korp	А	0	<b>99</b>
coal par c dl	The Downloadable Version of Classics of English and American Literature as translated by Kersti Juva, English-Finnish parallel	ACA		Download		0	

# **Citation instructions**

amph	amph-Corpus	;	ACA	•)	Download	В	0	<b>99</b>
ArkiSyn-korp	ArkiSyn Data Korp Version	PUB		Korp	В	0	<b>77</b>	
AVOID	Corpus of Age-related Voice Disguise (AVOID)			*)	Download	В	9	"
BeserCorp	The Corpus of Beserman Udmurt				Korp		0	<b>9</b> 9
ccmh-src	Corpus Cyril Slavonic Text Reference instructions: AVOID Please cite the language resource as follows: Kinnunen, T., Hautamäki, R. G., Sahidullah, M., Ha Disguise (AVOID) [speech corpus]. Kielipankki. Ref Show: [Bibtex] [Zotero]			Werner, http://ur olar	S., & Bentz, M ( m.fi/urn:nbn:fi:ll	[suomek Corpus of Age 5-20180606	si] [in Englis e-related Voie 21	h]

# Do researchers follow the citation instructions for corpora & resources?

 Google Scholar 2023-11-30:
 571 references cite the URN prefix of the Language Bank of Finland (*urn:nbn:fi:lb-*)

=	Google Scholar	um:nbn:fi:lb-
٠	Artikkelit	Noin 571 tulosta (0,02 sekuntia)
	Mikä tahansa päiväys Vuodesta 2023 Vuodesta 2022 Vuodesta 2019 Oma ajanjakso	[SITAATTI] Karjalan suomen sanomalehtikorpus         J Mäkisalo, H Kemppanen - 2017 - erepo.uef.fi         The corpus contains issues of the Karjalan Sanomat newspaper published in 2012-2014.         The corpus is available in Kielipankki - the Language Bank of Finland (http://urn.fi/urn.fi/urn.fi/lb         ☆ Tallenna 切 Viittaa Aiheeseen liittyviä artikkeleita Kaikki 2 versiota ≫         Istraattu The Karelian Finnish Newspaper Corpus

### META SHARE O:

curated metadata catalogues for resources

### Platforms and services provided by CSC

### **Download service:**

downloadable versions of resources



concordances and statistics from corpora; links to external resources

HeLI-OTS: automatic language identification from text **Demo** tools: examples of software for text and speech analysis & processing

Aalto-ASR, Tekstiks: automatic speech recognition (ASR) tools

>> Language models...

See the complete list of tools at https://www.kielipankki.fi/tools/

### The Language Bank of Finland Researchers of the Month

### Suomeksi

Do you know researchers who use the Language Bank of Finland and who might be good candidates for Researcher of the Month? Would you be one of them? Inform us!

		Etsi:	
Published	Researchers	Corpora/Tools	Publications +
2023.06	Mikael Varjo	ArkiSyn	2022, 2020, 2019, 2018
2023.05	Rosa González Hautamäki	AVOID	2019, 2018, 2017, 2016
2023.04	Johanna Vaattovaara	Suomi24	2022a, 2022b, 2019, 2012, 2011, 2009
2023.03	Noora Hoffrén	CFINSL, snowfrog-dl	2019
2023.02	Maria Sarhemaa	Suomi24	2022, 2021





WP1.1:

# TEXT PROCESSING AND ANNOTATION ENVIRONMENTS

FIN-CLARIN

### The Finnish Sub-corpus of the Newspaper and Periodical Corpus of the National Library of Finland version 2, Korp (*klk-fi-v2*) <u>http://urn.fi/urn:nbn:fi:lb-202009152</u>

Finnish   Swedish   Other languages   Parallel			Log in Suomi   Sve	nska   English 🝂 🔁 Cite Korp MENU 🚍
	Select corpora: 0 of 1246 corpora selected - 0 of 37.30G tokens			KIELIPANKKI The Language Bank of Fieldard
Simple Extended Advanced Compare 3	✓ Select all     ✓ Select non     Tetokonevälitteistä viestintää (33)	e		
	Akateemisia tekstejä (23)		Size	
KWIC:         hits per page: 25 v         sort within corpora: not sorted v	<ul> <li>Historiallisia aineistoja (115)</li> <li>Oppijansuomen aineistoja (4)</li> <li>Kirjallisuusaineistoja (26)</li> </ul>	Show word picture	2,040,320,047 Se	entences
	<ul> <li>Käsin annotoituja aineistoja (1)</li> </ul>		22,452,829,704 1	Tokens
	<ul> <li>Lehti- ja uutisaineistoja (880)</li> <li>Kansalliskirjaston lehtikokoelman (KLK) suomenkieliset lehdet (179)</li> </ul>		Time Coverage	
	<ul> <li>Kansalliskirjaston lehtikokoelman (KLK) suomenkieliset lehdet, versio 2 (test) (225)</li> <li>1990- ja 2000-luvun suomalaisia aikakaus- ja</li> </ul>	Kansalliskirjaston lehtikokoelman (KLK) suoi (test)	1771 - 2021	
	sanomalehtiä (369) ▶ ☐ Ylen suomenkielinen uutisarkisto (22) ▶ ▲ Suomen kielen tekstikokoelma (SKTP/FTC): lehdet [RES	Kansalliskirjaston sanoma- ja aikakauslehtikokoelma versio 2, Korp (klk-fi-v2-korp)	n suomenkielinen osakorpus	
	<ul> <li>(49)</li> <li>▲ Kotuksen aikakauslehtikorpus (8)</li> <li>▶ STT:n uutisarkisto 1992–2018 (beta) (27)</li> </ul>	Aineistossa olevat linkit sivun kuviin ja PDF-tiedostoil digitaalisten aineistojen verkkosivuille. Useimmille vu sivun kuvien ja PDF-tiedostojen käyttö on sallittu vai	hin vievät Kansalliskirjaston otta 1939 uudemmille lehdille n tutkimustarkoituksiin ja mistä Kansalliskisiastan	
	🖨 Karjalansuomi [ACA]	palvelussa.	amista nansalliskifjästön	
	<ul> <li>Parlamenttiaineistoja (2)</li> <li>Puheaineistoja (138)</li> </ul>	Huomaa, että korpus on vielä testivaiheessa, joten s muutoksia ilman ilmoitusta.	iihen voi tulla merkittäviäkin	

nnish   Swedish   Other languages   Parallel	Log in Suomi   Svenska   English A111 Cite Korp
Simple Extended Advanced Compare	16 corpora selected - 900.28M of 37.30G tokens
Simple     Extended     Advanced     Compare       "ilter:     Add publication       word     is <any word="">     Aa       or     Aa       <any word="">     Aa       <any word="">     Aa       Search     within sentence</any></any></any>	In 2020 and 2021, get statistics by the name of publication (search for 'any word')
WIC: hits per page: 25 v sort within corpora: not sorted v Sta	atistics: compile based on: publication    Show statistics  Show word picture  date (ISO)
Results:	digitization date     identified languages in text     Style Cate
	<ul> <li>issue name</li> <li>issue number</li> <li>label</li> </ul>
	metadata filename     original filename
	<ul> <li>□ page number</li> <li>✓ publication</li> </ul>
	<ul> <li>sentence language (identified)</li> <li>sentence language identification</li> </ul>

KWIC: hits per page: 25 v sort within co	rpora: not sorted ~ Statistics: compile based on: p	ublication 🚽 🗹 Show statistics 🗌 Show	word picture
KWIC Statistics Word picture	St	atistics by the nam	ne of
<sup>⊮A</sup> ∖ Show Trend Diagram Show map	. pu	iblication (search f	or 'any word')
Number of rows: 301	in	2020-2021 (klk-fi-	v2)
publication	Total	KLK suomi v2: 2020	KLK suomi v2: 2021
🔽 Σ	1 000 000 (900 280 994)	1 000 000 (537 488 468)	1 000 000 (362 792 526)
Helsingin Sanomat	60 091 (54 098 749)	53 738,8 (28 883 982)	69 501,9 (25 214 767)
🔲 Kaleva	30 475,7 (27 436 662)	27 746,6 (14 913 470)	34 518,9 (12 523 192)
Aamulehti	30 133 (27 128 185)	31 520,6 (16 941 960)	28 077,3 (10 186 225)
□ Ilta-Sanomat	27 548,9 (24 801 789)	24 395,1 (13 112 080)	32 221,5 (11 689 709)
🔲 Satakunnan Kansa	27 221,2 (24 506 771)	24 712,4 (13 282 606)	30 938,2 (11 224 165)
Turun Sanomat	26 823,8 (24 148 960)	26 775,4 (14 391 468)	26 895,5 (9 757 492)
Etelä-Suomen Sanomat	25 511,2 (22 967 280)	23 027,7 (12 377 137)	29 190,6 (10 590 143)
🔲 Lapin Kansa	24 648,8 (22 190 848)	23 039,7 (12 383 552)	27 032,8 (9 807 296)
Ilkka-Pohjalainen	24 100,2 (21 696 991)	21 497,8 (11 554 840)	27 955,8 (10 142 151)
🔲 Keskipohjanmaa	23 916,5 (21 531 589)	21 454,5 (11 531 548)	27 564,1 (10 000 041)
Savon Sanomat	22 006 (19 811 571)	19 544 (10 504 666)	25 653,5 (9 306 905)
Keskisuomalainen	20 734,3 (18 666 739)	21 396,5 (11 500 380)	19 753,3 (7 166 359)
🔲 Karjalainen	20 191,2 (18 177 773)	20 279 (10 899 721)	20 061,2 (7 278 052)
Hämeen Sanomat	19 833,2 (17 855 486)	18 257,4 (9 813 119)	22 168 (8 042 367)
Iltalehti	17 134 (15 425 423)	15 278,6 (8 212 076)	19 882,8 (7 213 347)
🔲 Länsi-Savo	16 873,6 (15 190 967)	14 872,8 (7 993 951)	19 837,8 (7 197 016)
Salon Seudun Sanomat	15 075.9 (13 572 560)	14 015.8 (7 533 315)	16 646.6 (6 039 245)





## Language identification from text: HeLI-OTS 1.5

Available on Zenodo: https://doi.org/10.5281/zenodo.10071264







WP1.2:

# **SPEECH PROCESSING AND ANNOTATION**

FIN-CLARIN

### Speech recognition: speech to text

Automated speech transcription service for Estonian speech and a user interface for transcription editing.

GET STARTED VIEW

VIEW DEMO

# **Tools:** tekstiks.ee

#### How does it work?

Tekstiks.eeis a public speech recognition service of TalTech's Laboratory of Language Technology. The system demonstrates technologies and models developed in this lab. These currently achieve state-of-theart Estonian speech recognition results, even compared to commercial alternatives. The system is fully automated and can process multiple files in parallel. There can be a queue and delays, especially during business hours. The average processing time is about half of the recording's length. Coming soon: real-time progress and expected time of completion information

	• • • Tekstik	5	HOME	FILES	DEMO	🍀 EN	•	ς a
	ЭC						Ł DO	WNLOAD
	<b>Reene Leas</b>	00:11 Ta on 16. mai kell sai kuus nin uudised. Mina olen Reene Lea usaldusküsimusega.	g Päevakaja s. Valitsus si	võtab kok dus lisaee	ku tänase elarve vasi	päeva tä tuvõtmise	htsan e riigil	nad kogus
200:22 Survet on jaanipäeval peaks mingi asja vastu mõtlema, mis toob suuri rah püsikulusid ilma selle seda eelarve kontekstis vaatamata, et noh, tegeliku ole ühtegi põhiust.						rahal likult	isi seal ei	

#### 1. Upload a speech recording in Estonian

Most of the popular audio and video formats are supported. Max size limit is 500MB

#### 2. Wait for the speech recognition to complete.

The system trained using machine learning methods will search for Estonian speech segments and tries to differentiate multiple speakers. Then it will transcribe the speech segments into text and finally will add punctuation. Many Estonain celebrities and radio personalities can be identified by name as well.

#### 3. Correct speech recognition mistakes

The editing of the transcription is interactive. The integrated audio player and the text are in sync. The currently playing word in coloured to simplify the manual editing.

### 4. Download the result

Download the transcription, currently the DOCX format is supported.





# **Tools:** tekstiks.ee



	pohjantuuli_ja_aurinko.wav Kesto: 00:00:36	/ Ladattu: 1. joulukuuta 23	Muokkaustila: Tavallinen (korjaa tunnistusvirheet) Merkinnät (teksti ja puhe)					
	\$		5 C	🕹 LATAA				
	<u>8</u> S1	<sup>00:01</sup> pohjantuuli ja aurinko pohjant	uuli ja aurinko väittelivät kummalla	olis enemmän voimaa				
	<u> </u>	<sup>00:06</sup> kun he samalla näkivät kulkija voimakkaampi joka nopeamm	<sup>00:06</sup> kun he samalla näkivät kulkijan jolla oli yllään lämmin takki silloin he sopivat että se on voimakkaampi joka nopeammin saa kulkijan riisumaan takkinsa					
	<u></u> ద్ది క1	<sup>00:16</sup> voihan tuuli alkoi puhaltaa niir kääri <mark>miestäkin</mark> ympärilleen	n että viuhuu mutta mitä kovempaa	se puhalsi sitä tarkemmin				
	<u>S</u> S1	00:23						
	1.0x •			S1				
10	15	20	25	30	35			



-11	le Tekstiks
-----	-------------

	pohjantuuli_ja_aurinko.wav <sub>Kesto:</sub> 00:00:36	Ladattu: 1. joulukuuta 23	Muokkaustila: Tavallinen (korjaa tunnistusvirheet) Merkinnät (teksti ja puhe)		
	\$		5 C		
	<u>S</u> S1	<sup>00:01</sup> pohjantuuli ja aurinko pohjantuu	uli ja aurinko väittelivät kummalla	olis enemmän voimaa	
	<u>ද</u> s1	<sup>00:06</sup> kun he samalla näkivät kulkijan voimakkaampi joka nopeammir	jolla oli yllään lämmin takki silloir saa kulkijan riisumaan takkinsa	n he sopivat että se on	
	<u>2</u> S1	<sup>00:16</sup> voihan tuuli alkoi puhaltaa niin e kääri <mark>mies takin</mark> ympärilleen	että viuhuu mutta mitä kovempaa	a se puhalsi sitä tarkemmin	
	<u>ද</u> s1	00:23			
	1.0x -				
			S1	S1	
10	15	20	25	30	35









### WP2.1:

# SOCIAL DATA SCIENCE



# What kind of data can be deposited?

- Text or speech in any natural language
  - text documents, audio and video recordings
  - annotations and transcripts may be included
  - must be in useful and well-described formats
- Make sure you have the rights to distribute the data, at least for research purposes
  - Copyright and other Intellectual Property Rights
  - Personal data

Contact FIN-CLARIN for details fin-clarin@helsinki.fi

# Submit the basic details about your dataset <u>http://urn.fi/urn:nbn:fi:lb-2021121422</u>

With this form you can ask <u>FIN-CLARIN</u> to publish on-line the essential metadata of the corpus or tool that you wish to deposit with <u>Kielipankki (The Language Bank of Finland)</u> for distribution. The corpus or tool can be completed or still in progress.

- Please fill in all relevant parts of the form, even if the information provided is still preliminary.
- If necessary, the information you provide can be edited and completed together with FIN-CLARIN.
- Completing the form does not oblige you to conclude the deposition agreement, but the information may be of great help if you need further advice on your resource later.
- FIN-CLARIN may also contact you, or the responsible person you have indicated, to agree on follow-up measures concerning the resource.
- Once you have requested that we add the metadata of the language resource in the language resource catalogue, your resource can immediately gain more visibility, even if it is not yet ready for publication.
- FIN-CLARIN is happy to help you with any questions related to the deposition and distribution of the resource. You can reach us by sending email to fin-clarin (ATT) helsinki.fi

Other language resources to be published by Kielipankki (The Language Bank of Finland) of FIN-CLARIN

Contact details

(\*) Name of the information provider \*

(\*) Email address of the information provider \*

ORCID identifier of the information provider (instructions)

What is ORCID?

The home organization of

# Are safety measures required for sharing?

- The Language Bank of Finland can offer various means for protecting personal data and other types of restricted content:
  - 1. Access management, if needed: university login (ACA); access granted on an individual basis upon application (RES)
  - Resource-specific data protection terms and conditions (must be accepted by the end-users)
  - 3. Data encryption (first resource: <u>findarc</u>)
  - 4. Sensitive Data (SD) services at CSC



# Two types of agreements are usually needed:

- 1. Deposition (distribution) license agreement (DELA) <u>https://www.kielipankki.fi/support/dela/</u>
  - The Language Bank (University of Helsinki) obtains the right to distribute the resource to end-users under specific terms and conditions.
  - a DELA is usually required, unless the University of Helsinki is the rightholder or the resource is already available under a public license.
  - For resources containing personal data, the agreement also includes the details of the data controller and the resource-specific data protection terms and conditions regarding the redistribution of the data via the Language Bank.
- 2. End-User License Agreement (EULA)
  - The End-User agrees to use the data under specific conditions which the rightholder has approved.

# **CLARIN** license categories



Publicly available



Available for academic, logged in users



Access is based on an individual application Language Bank Rights: lbr.csc.fi

# More detailed license conditions

- +BY the author(s) must be cited
- +NC non-commercial use only
  - +ID login is required
- +PLAN a research plan is required
- +PRIV contains personal data
- +NORED redistribution is not allowed

+DEP modified versions can be redistributed via CLARIN

and other resource-specific conditions, if required (e.g., data protection terms and conditions)

# Are you using a resource that includes personal data (license condition **+PRIV**)?

- As a user, you are required to submit the title of your project and a link to the **public privacy notice** regarding your project, to be <u>published by Kielipankki</u>.
  - The privacy notice must be submitted when applying for access to a restricted resource that contains personal data.
  - If the resource is available without an application, submit your privacy notice via the <u>e-form</u>.




WP2.2:

# LEARNERS' ASSESSMENT ENVIRONMENTS



## Donate Speech Corpus v1.0 http://urn.fi/urn:nbn:fi:lb-2020090321

Tweet #lb\_puhelahjat

### Donate Speech datasets (puhelahjat) for research use

Suomeksi

Donate Speech datasets for commercial use: see further details on another page

#### Important information for all users of this resource: Removal requests

Versions of this resource:	
Donate Speech Corpus, version 1.0	Apply for access rights, academic research use only
<ul><li>License (for researchers)</li><li>Attribution instructions</li></ul>	<b>+PRIV:</b> This resource contains personal data. Submit public information about personal data processing
	Download the resource
Donate Speech Corpus: Sample	Download the resource





### Pulistaan meille parempia palveluja

Nyt kerätään kaikenlaista puhuttua suomea! Lahjoittamasi puheen puheen avulla esimerkiksi ääniohjatut laitteet voivat oppia ymmärtämään erilaisia murteita ja puhetapoja. Valitse alta aihe, josta haluat puhua.



### KIELIPANKKI The Language Bank of Finland

#### LATAUKSET DOWNLOADS

🏠 НОМЕ 👚 UP

Logged in as:

#### Location: /download/puhelahjat/complete/

Name	Size	Description
puhelahjat_v1_2020_part_01_alignments_manual_2022-04-22.zip	11M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_01_annotations_manual_2021-10-28.zip	3.3M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_01_audio.zip	11G	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_01_meta.zip	252K	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_02_alignments_manual_2022-04-22.zip	11M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_02_annotations_manual_2021-10-28.zip	3.1M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_02_audio.zip	11G	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_02_meta.zip	239K	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_03_alignments_manual_2022-04-22.zip	10M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_03_annotations_manual_2021-10-28.zip	3.1M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_03_audio.zip	11G	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_03_meta.zip	268K	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_04_alignments_manual_2022-04-22.zip	11M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_04_annotations_manual_2021-10-28.zip	3.2M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_04_audio.zip	11G	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_04_meta.zip	253K	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_05_alignments_manual_2022-04-22.zip	10M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_05_annotations_manual_2021-10-28.zip	3.1M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_05_audio.zip	11G	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_05_meta.zip	250K	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_06_alignments_manual_2022-04-22.zip	11M	Donate Speech Corpus 1.0
puhelahjat_v1_2020_part_06_annotations_manual_2021-10-28.zip	3.1M	Donate Speech Corpus 1.0
puhelahjat v1 2020 part 06 audio.zip	11G	Donate Speech Corpus 1.0



#### > audio

#### Name

clt0000001\_ses01\_rec0001.flac clt000001\_ses01\_rec0002.flac clt0000001\_ses01\_rec0003.flac clt0000002 ses01 rec0001.flac clt0000002\_ses01\_rec0002.flac clt000002\_ses01\_rec0003.flac clt0000002 ses01 rec0004.flac clt0000002 ses01 rec0005.flac clt0000002 ses01 rec0006.flac clt0000003 ses01 rec0001.flac clt0000003 ses01 rec0002.flac clt0000003 ses01 rec0003.flac clt0000003 ses01 rec0004.flac clt0000003 ses01 rec0005.flac clt000003\_ses01\_rec0006.flac clt000003\_ses01\_rec0007.flac clt000004\_ses01\_rec0001.flac clt0000004 ses01 rec0002.flac clt0000004\_ses01\_rec0003.flac clt0000004 ses01 rec0004.flac clt0000004 ses01 rec0005.flac clt0000004 ses01 rec0006 flac

#### 2020\_part\_01

alignments

- 🖊 🚞 manual\_2022-04-22
  - clt0000001\_ses01\_rec0002.TextGrid
  - clt000003\_ses01\_rec0002.TextGrid
  - clt0000003\_ses01\_rec0004.TextGrid
  - clt0000003\_ses01\_rec0006.TextGrid
  - clt0000004\_ses01\_rec0002.TextGrid
  - clt0000004\_ses01\_rec0003.TextGrid
  - clt0000004\_ses01\_rec0006.TextGrid
  - clt0000004\_ses01\_rec0007.TextGrid
  - clt0000004\_ses01\_rec0008.TextGrid
  - clt0000005\_ses01\_rec0002.TextGrid
  - clt0000005\_ses01\_rec0004.TextGrid
  - clt0000005\_ses01\_rec0006.TextGrid
  - clt0000005\_ses02\_rec0005.TextGrid
  - clt0000005\_ses02\_rec0007.TextGrid
  - clt0000006\_ses01\_rec0002.TextGrid
  - clt000006\_ses01\_rec0004.TextGrid
  - clt0000006\_ses01\_rec0005.TextGrid
  - clt000006\_ses01\_rec0006.TextGrid
  - clt0000007\_ses01\_rec0002.TextGrid



# The Donate Speech Corpus has helped in training automatic speech recognizers

- Automatic speech recognition systems, trained by Aalto Speech Research team on the Donate Speech Corpus, are available on Zenodo.
- More on the content of the Donate Speech Corpus:

Moisio, A., Porjazovski, D., Rouhe, A. et al. Lahjoita puhetta: a large-scale corpus of spoken Finnish with some benchmarks. *Lang Resources* & *Evaluation* (2022). <u>https://doi.org/10.1007/s10579-022-09606-3</u>



# L2 speech corpora (soon available!)

- Speech recorded from second-language learners in the DigiTala project (UHEL, Aalto, JYU), including oral skills assessment data:
  - DigiTala: L2 Finnish data from upper secondary schools and university, autumn 2021
  - DigiTala: L2 Finnish data from upper secondary schools, spring 2021
  - DigiTala: L2 Swedish data from adult language learners, spring 2023
  - DigiTala's YKI data (Yleiset kielitutkinnot, general language tests)





# ASR systems for L2 Finnish

- Models trained by Aalto Speech Research team for performing automatic speech recognition (ASR) and pronunciation rating for children's L2 Finnish.
  The systems are available on HuggingFace Hub.
- Tools for the automatic speaking assessment of spontaneous L2 Finnish and Swedish (DigiTala project)







WP2.3:

# **TRANSLATION AND INTERPRETATION**





#### Heijän nähtyy paimoit ruvettih sanelemah, midä heile oli sanottu lapseh näh.

Сэсся висьталісны кага йылысь кывломторсо.

base form: lapsi baseform (compound boundaries):



Download hit page as...

юс да лясниын куйлысь кагаёс .

### PEECH (KRL) (BETA)

n da pani žiivatoin soimeh, ku heile ei olluh sijua mat pšen ta pani Hänet šoimeh , šentäh kun heilä ei löyty neh, ku heile ei olluh sijua matkuniekoin talois. meh , šentäh kun heilä ei löytyn tilua matkuštajien tal viruu soimes. »

kumpani makasi šoimešša .

WORD ATTRIBUTES base form: lapsi baseform (compound boundaries): lapsi part of speech: noun msd: Case=GenlNumber=Sing dependency relation: apposition Show Dependency Tree Listen utterance in Korp (OLO) Listen utterance in Korp (kpv) Listen utterance in Korp (krl) Listen utterance in Korp (myv)

>> try this query in Korp <<

## KELIPANKKI The Language Bank of Finland



The Language Bank of Finland		
		🏫 номе 會 UP
Location: /download/xmas-gospel-tts/		Logged in as:
Name	Size	Description
LICENSE.txt	19К	Licence information
README.txt	2.0K	More information
xmas-gospel-tts-src.zip	38M	CC BY NC

CLARIN PUB- tai ACA-lisensoidut valikoidut versiot ladattavista kieliain hakemistosta /appl/data/kielipankki/. Palvelimien käyttö edellyttää CSC-

Selected versions of the downloadable CLARIN PUB or ACA licensed cor /appl/data/kielipankki/. To use the servers you need a CSC account. Tech Parallel corpus with linked audio: Christmas Gospel text-to-speech in four Uralic languages, source (= downloadable version)





Kiitos! Tack! Thank you!

## www.kielipankki.fi

**General support** fin-clarin@helsinki.fi

**Technical support** *kielipankki@csc.fi* 

# WP 2.4 The Helsinki Term Bank for the Arts and Sciences

#### FIN-CLARIAH Meeting in Tampere, December 1st, 2023

Harri Kettunen (harri.kettunen@helsinki.fi) & Tiina Onikki-Rantajääskö (tiina.onikki@helsinki.fi) tieteentermipankki-info@helsinki.fi



### TIETEEN TERMIPANKKI

VETENSKAPSTERMBANKEN I FINLAND THE HELSINKI TERM BANK FOR THE ARTS AND SCIENCES



# WP2.5 Automated Text Tools for Learner Language



- trained with 4 CEFR-annotated data sets (CEFR = Common European Framework of Reference, from A1 to C2).
  - ICLFI: International Learner Corpus of Finnish
  - LAS2: The Advanced Finnish Learners' Corpus
  - CEFLING: Linguistic Basis of the Common European Framework for L2 English and L2 Finnish
  - TOPLING Paths in Second Language Acquisition
    - = 7700 texts, 1.6 million tokens
- trained on FinBERT (by Turku NLP group)

Contact:

Jenny Tarvainen: jenny.h.tarvainen@jyu.fi Ida Toivanen: ida.m.toivanen@jyu.fi Ari Huhta: ari.huhta@jyu.fi

### PRELIMINARY PROGRAM

11.00-11.10 Welcoming Words by Sanna Kumpulainen, Associate Professor in Information Studies, Tampere University

11.10-12.00 Keynote I on Studying SSH Research Needs: Elina Late, Senior Research Fellow in Information Studies, Tampere University

12.00-13.00 Lunch

13.00-13.45 Keynote II on Language Models: Sampo Pyysalo, Associate Professor at the Department of Computing, University of Turku

13.45-14.30 Work Package Presentations I

13.45-14.00 WP1.3 Veronika Laippala: Noise-Tolerant NLP

14.00-14.20 WP 1.1, 1.2, 2.1, 2.2, 2.3 Mietta Lennes: Kielipankki – The Language Bank of Finland

14.20-14.25 WP2.4 Harri Kettunen: Helsinki Term Bank for the Arts and Sciences

14.25-14.30 WP2.5 Jenny Tarvainen: Automated Text Tools for Learner Language

#### 14.30-15.00 Coffee

#### 15.00-16.00 Work Package Presentations II

15.00-15.05 WP3.1 Martin Matthiesen: Pipeline from the National Library to CSC 15.05-15.10 WP3.2 Tanja Välisalo: Named Entity Recognition for NARC Data 15.10-15.15 WP4.3 Eetu Mäkelä: Evaluation and Subsetting 15.15-15.20 W4.1 Julia Matveeva: Metadata Harmonization 15.20-15.25 WP4.4 Mikko Laitinen: Twitter 15.25-15.30 WP4.2 Eero Hyvönen: LOD 15.30-15.35 WP3.4 Raine Koskimaa: Game Streams 15.35-15.40 WP3.3 Maria Valaste: Qualitative Surveys 15.40-15.45 WP3.5 Kimmo Elo (Risto Turunen replacing): Text Networks

16.00-17.00 Free Chilling & Refreshments / Parallel session: Executive Board Meeting (with Zoom option)



# WP 3.1

## Data of the National Library

- Contains of out-of-copyright digitized publications
- Newspapers (1771-1918), journals (1771-1912), books
  - news, articles, advertisements, discussion, weather, ... etc
- Page text in ALTO XML, Access images (JPG)
- 9 TB (zip compressed)

### Added value for researchers

- Accessibility: Dataset is closer to CSC processing environments (Puhti, Mahti, LUMI). We handle:
  - catalogue discovery
  - metadata parsing
  - API access
  - downloading tens of millions of files.
- Interoperability: Different research groups can be sure they use the same data
- **R**eproducibility: Datasets can be referenced, older versions can be stored.
- See http://urn.fi/urn:nbn:fi:lb-202311261

### Dataflow

Metadata description Digitization Formats



### **Download control**



### Lessons learned

- Do not underestimate scaling
  - "Harmless" access errors
  - Airflow memory requirements
- Well defined APIs have value
- Less is more
  - o squashfs in containers vs zip files

# Future work / Open questions

- Defining the prerequisites for offering copyrighted material in the future
- Making the reference URN machine actionable
  - E.g.: http://urn.fi/urn:nbn:fi:lb-202311261#a0c1b2
- Long-term availability
  - Ongoing research needed to keep data available
  - Lost snapshots cannot be recovered, only the last version of the data
  - Download of full dataset takes several weeks



# WP3.2 Named entity recognition (NER) for state authority archival data

Antero Holmila, Venla Poso, Ida Toivanen, Tanja Välisalo In co-operation with the National Archives of Finland / DALAI project The usability of mass digitised state authority data in the National Archives of Finland was improved by developing an AI based model for Named Entity Recognition

# → First prototype of the model published by NARC and can be tested at arkkiivi.fi and <u>https://huggingface.co/Kansallisarkisto/finbert-ner</u>

Final model to be published in Spring 2024

## Current unpublished NER model

- obtains comparable results with non-OCR'd data, while
- significantly improving results with OCR'd state authority archival data

# **FUTURE**

Automated **per-unit** metadata identification with NARC data:

- **Document date:** utilizing current work with NER of journal number (diaarinumero) detection
- **Document content type** (*e.g. correspondence, report*) using language tools and computer vision

# WP4.3 EVALUATION AND SUBSETTING

### WP4.3 EVALUATION AND SUBSETTING - ACHIEVEMENTS

### Developed a prototype for "universal" data analysis and subsetting pipeline:

- Researched and adopted a sustainable and customizable software package (Elasticsearch & Kibana stack)
- The software is usable on multiple levels of expertise: API and UI, 'clickable' data browsing interface and more serious programmatic digital humanities application.
- Developed practices for easily integrating new datasets into the solution for relatively rapidly adding new datasets for future projects.
- Developed (and still developing) documentation for easy adaptation and duplication of the process: Deploying the software on CSC or similar hosting service, adding and extracting datasets.
- Tested the solution pipeline in a practical project: <u>https://receptionreader.com/latin</u>

### WP4.3 EVALUATION AND SUBSETTING - CHALLENGES

### **Challenges (and lessons learned)**

- Developing complex plugins from scratch quite a demanding process this work is still in progress.
- Even "ready to use" harmonized data usually needs some further fine refinement before being ready for integration into a database.
- Adopting into use in any hosting service always seems to need some amount of jumping through hoops. But as a positive surprise: here, relatively little.
- Positive lessons learned: The system is way more usable almost out of the box than we expected: Serves as a very smooth and responsive backend for the test use case (Latin reuse reader).
- The standard data analysis tools are excellent for getting a quick overview of a dataset.

### WP4.3 EVALUATION AND SUBSETTING - FUTURE

### Future use & steps for this solution:

- Adopting Elasticsearch & Kibana in place of our current custom data API for use by various research groups in Helsinki.
- Developing language analyzer plugins to better serve scientific use.
- Smoothing out the workflow leading from data indexing to analysis and extraction towards a standardized model based on our test use cases.

### WP4.1 METADATA HARMONIZATION



### WP4.4 TWITTER (X)

- Five massive datasets collected
  - NTS: 799 million words (c. 74m messages from 888,098 accounts)
    - A user interface for easy data access + data shared through Kielipankki
  - Network datasets
    - 19,345 ego networks with 759,495 nodes
    - C. 11 billion words
    - Nordic region, UK, US and Australia (broad geographic coverage)
- The biggest challenge was obviously closing down the API > unique datasets from 2006-May2023;
- [We cannot solely rely on eccentric billionaires and the Big Tech for data access]
- Digital Single Markets directive and the subsequent changes in the Finnish copyright law
- Next up: data enrichment for background variables (age, social class in addition to networks)

## WP 4.2 National Linked Open Data Infrastructure: Sampo Data Services & Semantic Portals

Eero Hyvönen Professor, Department of Computer Science, Aalto University Director, Helsinki Centre for Digital Humanities (HELDIG), University of Helsinki <u>https://seco.cs.aalto.fi/u/eahyvone/</u>
## We Aim at a Paradigm Change: 4 Generations of Publishing Data for Humanities

- 1. Texts (Engravings, Hand-written, and Printed)
- 2. Online Systems for Searching and Exploring
- 3. Publishing Content as Linked Data with Tools for DH
- 4. Automatic Knowledge Discovery and Artificial Intelligence

\succ Our Focus

E. Hyvönen, : Using the Semantic Web in Digital Humanities: Shift from Data Publishing to Data-analysis and Serendipitous Knowledge Discovery, *Semantic Web*, 1 (11), 2020 <u>pdf</u>

## Problem of generations 1 and 2: Content is available as texts for humans to read but not as data for computers!

#### Using Sampo Model for Cultural Heritage Data by 1) human users and 2) machines



Jupyter

## **Sampo LOD Service**



## Linked Data Finland Service: https://ldf.fi

# <section-header>

## Digital Humanities Research Tools





#### Elements of a National Semantic Web Infrastructure for Open Cultural Heritage FAIR Data



Eero Hyvönen: How to Create a National Cross-domain Ontology and Linked Data Infrastructure and Use It on the Semantic Web, forth-coming, 2023.

## **Sampo Series of LOD Services & Portals**





- 2. CultureSampo Finnish Culture on the Semantic Web (2008) [107 000 users]
- 3. TravelSampo Mobile Contextualized Services of Cultural Tourism (2011)
- 4. BookSampo Finnish Fiction Literature on the Semantic Web (2011) [1.6 million users in 2022]
- 5. WW1LOD World War I Linked Open Data (2014)
- 6. WarSampo Finnish World War 2 on the Semantic Web (2015-19) [1 100 000 users]
- 7. Norssi Alumni on the Semantic Web Historical person registry using LOD (2017)
- 8. U.S. Congress Prosopographer U.S. Congress Legislators 1789-2018
- 9. BiographySampo Finnish Biographies on the Semantic Web (2018-20) [381 000 users]
- D. NameSampo Linked Data Workbench for Toponomastic Research (2019) [55 000 users]
- I. WarVictimSampo 1914-1922 National War History [80 000 users]
- Mapping Manuscript Migrations (MMM) medieval manuscripts (2020) [9100 users]
- AcademySampo Finnish Academic People 1640 1899 (2021) [26 000 users]
- FindSampo Archaeological Finds on the Semantic Web (2021) [7 000 users]
- WarMemoirSampo Memoirs of Finnish WW2 veterans (2021) [3 800 users]
- LetterSampo Early Modern Letters on the Semantic Web (2022)
- **ParliamentSampo** Parliament of Finland on the Semantic Web (2023)
- LawSampo Finnish Legislation and Case Law on the Semantic Web (2023)
- **BookSampo II** Semantic Search, Browsing, and Data-analyses (2023)
- **OperaSampo** Opera and music theater performances in Finland 1830-1960 (2023)
- 21. ConfermentSampo Conferment Ceremonies of the University of Helsinki 1640-1899 (2023)

Year of 5 Sampos 2023

Sampo = Mythical artifact of the Finnish Epic Kalevala that gives to its owner riches and good fortune. A metaphore of amazing ancient technology. Links to all Sampos: <u>https://seco.cs.aalto.fi/applications/sampo/</u>

Eero Hyvönen: Digital Humanities on the Semantic Web: Sampo Model and Portal Series. Semantic Web, vol. 14, no. 4, pp. 729-744,, 2023. pdf link



A. Gallen-Kallela: The Defense of the Sampo, 1896, Turku Art Museum,

## PARLAMENTTISAMPO AVAA EDUSKUNNAN MILJOONA PUHETTA JA KANSANEDUSTAJIEN VERKOSTOT KAIKKIEN TUTKITTAVIKSI



#### Project homepage: https://seco.cs.aalto.fi/projects/semparl/

Eduskunnan täysistunnoissa on vuosina 1907–2022 pidetty lähes miljoona puheenvuoroa, ja yhteensä puhujia on ollut noin 2 800. Uusi Parlamenttisampo-palvelu tarjoaa kaiken tämän aineiston yhtenäisenä datana, mikä mahdollistaa kansanedustajien sekä poliittisen kielen, kulttuurin ja verkostojen tutkimisen entistä helpommalla tavalla.

ähes miljoona eduskunnan täysistunnoissa pidettyä puheenvuoroa on Parlamenttisampo-hankkeessa ensimmäistä kertaa muunnettu linkitetyksi avoimeksi dataksi (englanniksi *Linked Open Data*). Avoimen datapalvelun päälle on kehitetty semanttinen Parlamenttisampo-portaali, jota voivat tutkijoiden lisäksi käyttää media, poliitikot ja suuri yleisö.

Parlamenttisampo on uusi jäsen Sampoportaalien sarjassa, joilla on ollut jopa miljoonia käyttäjiä semanttisessa webissä. Muita Sampo-portaaleja ovat esimerkiksi lainsäädäntöä ja oikeustapauksia kokoava *Lakisampo,f*i ja *Biografiasampo,f*i, joka sisältää Suomalaisen Kirjallisuuden Seuran Kansallisbiografiassa julkaistujen suomalaisten merkkihenkilöiden elämäkerrat ja heidän verkostoja.

#### EDUSKUNNAN AINEISTOT AVOIMEKSI FAIR-DATAKSI JA SOVELLUKSIKSI

Suomalaisen demokratian yksi perusta on, että eduskunnan päätöksenteko ja lainsäädäntötyö ovat avoimia ja niihin liittyvät aineistot saatavilla. Keskeinen tähän liittyvä

Pitkäaikainen kansanedustaja Veikko Vennamo (1913–1997) käytti eduskunnan täysistunnoissa enemmän puheenvuoroja kuin kukaan muu. Vuonna 1974 hänet jopa kannettiin ulos istuntosalista, kun hän ei suostunut poistumaan täysistunnosta puhemiehen kehotuksesta huolimatta.

#### And More Sampos are Being Forged at the Smithy

- **CoCoSampo** Finnish letters on the Semantic Web 1809–1917
  - <u>https://seco.cs.aalto.fi/projects/coco/</u>
- **FinEst-LawSampo** Cross-border Multilingual Legislation Search
  - <u>https://seco.cs.aalto.fi/projects/finestsampo/</u>
- **PASampo** British Museum Portable Antiquity Scheme case study
  - https://seco.cs.aalto.fi/projects/diginuma/
- CoinSampo Coin finds of the Nartional Museum of Finland
  - https://seco.cs.aalto.fi/projects/diginuma/
- NomismaSampo Numismatics based on Nomisma.org model and data
  - <u>https://seco.cs.aalto.fi/projects/diginuma/</u>
- ArtSampo Finnish Art on the Semantic Web
  - <u>https://seco.cs.aalto.fi/projects/taidesampo/</u>
- **SampoSampo** Finnish Linked Open Data Cloud of Cultural Heritage



A. Gallen-Kallela: The Forging of the Sampo, 1893, Ateneum Museum,



## https://seco.cs.aalto.fi/projects/fin-clariah/

## WP3.4 GAME STREAMS

Twitch Chat Collector Tool Ready

both live streams and recorded videos

Analysis Tool (Discussion Dynamics, Uniqueness Classification)

Functional, some polishing needed

Publication of code in Github (under which organization?)
Publication of collected data, permission from Twitch required

https://drive.google.com/file/d/1hS4t4XSE6KQNvmDFOx8RQHqN-Xd629bN/view?usp=sharing

## WP3.3 QUALITATIVE SURVEYS

WP members: Adeline Clarke (University of Helsinki), Ida Toivanen (University of Jyväskylä), Jani-Matti Tirkkonen (University of Eastern Finland), Jaakko Peltonen (Tampere University) and Maria Valaste (University of Helsinki)

- Github:
  - R package: <u>https://github.com/DARIAH-FI-Survey-Concept-Network/finnsurveytext</u>
    - Website: <u>https://dariah-fi-survey-concept-network.github.io/finnsurveytext/</u>
  - Other materials: <u>https://github.com/DARIAH-FI-Survey-Concept-Network/</u>
- Our video:
  - <u>finnsurveytext video</u>
  - If you want to comment, please contact us! (<u>adeline.clarke@helsinki.fi</u>)





## WP 3.5: Text network analysis of political texts

WP presentation 01.12.2023

Adjunct professor Kimmo Elo

Senior Research Fellow, Centre for Parliamentary Studies

E-mail: kimmo.elo@utu.fi





## Summary

- . Core data: FinParl (plenary debates of the Finnish eduskunta since 1905)
- . Deliverables:
  - KWIC tool for FinParl corpus: This tool provides a user interface to query word embeddings with KWIC method. The tool offers a simple, yet intuitive user interface built with R Shiny. Output: KWIC table, n-grams table, collocation network
  - TNA tool for the analysis of speeches of Finnish MPs: This tool provides functionalities for vocabulary based content analysis of political speeches. The user selects an MP and can then study 1) a timeline of the MP's plenary speeches, 2) a wordcloud of max. 500 most used words by the MP, as well 3) a speaker-to-concept network consisting of the 50 most frequently used concepts of the selected MP and and of his/her most similar colleagues by vocabulary
- **Release:** release beta versions will be opened for public use/testing on December 19, 2023



#### Funded by the European Union NextGenerationEU

KWIC-tool (1/2)

UNIVERSITY OF TURKU

#### Text networks

ilmasto#muutos	
KWIC / n-gram context size	
3	\$
Case insensitive?	
Choose years	
2.000	2,009 2,011 2,01
2,000 2,002 2,004 2,006 2,00	18 2,010 2,012 2,0

KWIC n-gram						
docname	from	to	pre	keyword	post	pattern
FI2009_9_97.1	428	428	, finanssikriisi ja	ilmastonmuutoksen	aiheuttama lähtö entisiltä	ilmastonmuutoksen
FI2009_6_4.2	215	215	. Mainitsen esimerkkeinä	ilmastonmuutoksen	, energian saatavuuden	ilmastonmuutoksen
FI2009_6_7.1	79	79	Kansalliset rajat ylittävää	ilmastonmuutosta	on mahdotonta estää	ilmastonmuutosta
FI2009_6_7.1	184	184	. Globaalit haasteet	ilmastonmuutoksesta	talouskriisiin ja edelleen	ilmastonmuutoksesta
FI2009_6_8.1	217	217	syvenevä talouslama ja	ilmastonmuutos	. Näitä uhkia	ilmastonmuutos
FI2009_6_9.2	44	44	. Globaalisti katsottuna	ilmastonmuutos	, ympäristötuhot ,	ilmastonmuutos
FI2009_6_9.2	72	72	huomioitu, että	ilmastonmuutos	ja ympäristökysymykset voivat	ilmastonmuutos
FI2009_6_13.2	119	119	, luonto ,	ilmastonmuutos	, influenssaepidemiat ,	ilmastonmuutos
FI2009_33_7.1	414	414	mikā auttaa myös	ilmastonmuutoksen	vastaisessa taistelussa .	ilmastonmuutoksen
FI2009_33_7.1	447	447	on erittäin perusteltu	ilmastonmuutoksen	nākōkulmasta . Aivan	ilmastonmuutoksen
FI2009_5_5.1	1382	1382	ja ilmastopolitiikkaa ja	ilmastonmuutokseen	sopeutumista . Maailmantalouden	ilmastonmuutokseen
FI2009_5_10.3	103	103	tāllā hetkellā ,	ilmastonmuutosta	. Voidaksemme jatkossa	ilmastonmuutosta
FI2009_6_81.1	42	42	, ympäristöhaasteet ,	ilmastonmuutos	, luonnonkatastrofit —	ilmastonmuutos
FI2009_6_97.1	111	111	bruttokansantuotteesta . Myös	ilmastonmuutoksen	torjunta edistāā rauhantyötā	ilmastonmuutoksen
FI2009_6_97.1	134	134	. Lisäksi kiihtyvä	ilmastonmuutos	aiheuttaa ankaraa kuivuutta	ilmastonmuutos
FI2009_6_122.1	475	475	tarttuvista taudeista,	ilmastonmuutoksesta	, suuronnettomuuksista ja	ilmastonmuutoksesta
FI2009_6_123.1	97	97	mainittiin tällaisina uhkina	ilmastonmuutos	, energian ja	ilmastonmuutos
Fl2009_6_175.1	136	136	tuleekin huomioida niin	ilmastonmuutos	, luonnonkatastrofit ,	ilmastonmuutos
FI2009_6_192.1	377	377	. Ympäristö ja	ilmastonmuutos	on mainittu tässä	ilmastonmuutos
FI2009_6_192.1	394	394	. Ympäristöhaasteet ja	ilmastonmuutos	ovat omana kappaleena	ilmastonmuutos
Fl2009_7_138.1	328	328	kehitettävä . Myös	ilmastonmuutoksen	hillitseminen lisää julkisen	ilmastonmuutoksen
FI2009_7_148.1	124	124	on huomioitava myös	ilmastonmuutoksen	torjunnan työmarkkinoille tuoma	ilmastonmuutoksen
Fl2009_7_152.1	160	160	näkökulmasta kestäviä ja	ilmastonmuutoksen	hillinnän kannalta tehokkaita	ilmastonmuutoksen
FI2009_7_167.1	172	172	myös verotus torjuu	ilmastonmuutosta	, suojelee luonnonvaroja	ilmastonmuutosta
Fl2009_14_82.1	450	450	uudistaa verotusta jo	ilmastonmuutoksen	vuoksi . Tässä	ilmastonmuutoksen
FI2009_17_112.1	237	237	aivan nykyisin puhuttavasta	ilmastonmuutoksesta	huolimatta niin kauan	ilmastonmuutoksesta
Fl2009_17_114.1	77	77	jo viittasi ,	ilmastonmuutoksen	kysymykseen, Sahelin	ilmastonmuutoksen
FI2009_17_119.1	361	361	joita muun muassa	ilmastonmuutos	aiheuttaa jo tänä	ilmastonmuutos
Fl2009_17_122.1	53	53	tukemaan kehittyviä maita	ilmastonmuutoksen	vastaisessa kamppailussa .	ilmastonmuutoksen
FI2009_17_122.1	64	64	puheenjohtaja Salolainen otti	ilmastonmuutoksen	vastaiset toimet esille	ilmastonmuutoksen
=12009_19_212.1	85	85	ja ennen muuta	ilmastonmuutoksen	kannalta kestävällä ja	ilmastonmuutoksen
Fl2009_21_39.1	87	87	riippumatta meistä —	ilmastonmuutoksesta	, joka täällä	ilmastonmuutoksesta
FI2009_22_153.1	109	109	kestävällä tavalla,	ilmastonmuutosta	hillitsevällä tavalla ,	ilmastonmuutosta
Fl2009_22_153.1	151	151	. Se on	ilmastonmuutoksen	ja ympäristöteknologian kehittäminen	ilmastonmuutoksen
FI2009_26_65.1	463	463	joita teemme taistelussa	ilmastonmuutosta	vastaan , ja	ilmastonmuutosta



#### Funded by the European Union NextGenerationEU

## KWIC-tool (2/2)



#### Text networks

ilmasto#muutos					
KWIC / n-gra	im context	size			
3					3
Case insensitive?					
Choose years			2.0	09 2,0	1] 2
2,000					

KWIC	n-gram	
x		Freq
Arvoisa h	erra puhemies	181
on se etta	1	158
ja se on		85
että se or	1	65
että meill	ā on	62
ja sitā kar	utta	54
Arvoisa n	ouva puhemies	51
on tärkeä	ā ettā	44
Arvoisa p	uhemies On	43
se ei ole		36
sitä mielta	ā ettā	35
On tärkeä	hā ettā	33
että Suor	ni on	31
myös se	että	31
on hyvā e	että	30
On selvä	a että	30
on tällä h	etkellä	30
tarkoittaa	sitā ettā	30
Arvoisa p	28	
ei ole var	aa	28
Arvoisa p	uhemies Hallituksen	27
meillä ei (	ole	27
prosenttii	n vuoteen mennessä	27
sillä taval	la että	27
Arvoisa p	uhemies Suomen	26
ja ennen	kaikkea	26
ja tāmā o	n	26
aikavälin	ilmasto- ja	25
ei ole ollu	it	25
Se ei ole		25
niin että s	ie	24
se että m	e	24
ei ole mit	ään	23
pitkān aik	avälin ilmasto-	23
Suomi on	sitoutunut	23



Mikko Pesälä (KESK) -Gunnar Jansson (RKP) -Sauli Hautala (KD) -

Huomioitavien sanojen määrä:

200

0.00

300

0.50

190 320 350 380 410 440 470 500

100





#### Timeline Wordcloud Concept network Valitse puhuja: Puhuja: Kimmo Sasi (KOK) Kimmo Sasi (KOK) -Sanapilven sanojen määrä: 180 200 220 240 260 280 Eniten samankaltaiset puhujat Pentti Mäki-Hakola (KOK) -Sirkka-Liisa Anttila (KESK) Raimo Vistbacka (SMP) -0.787 Mikko Elo (SDP) lkm Anneli Jäätteenmäki (KESK) liro Viinanen (KOK) ojen Juha Korkeaoja (KESK) lörn Donner (RKP) -Erkki Pulliainen (VIHR) -Paavo Nikula (VIHR) -Puh 0.50 0.00 Vähiten samankaltaiset puhujat Eva Biaudet (RKP) -Matti Viljanen (KOK) -0.267 Pirkko Ikonen (KESK) -Håkan Malm (RKP) -1990 1992 1994 1996 1998 2000 Elsi Hetemäki-Olander (KOK) -Valtiopäivät Saara-Maria Paakkinen (SDP) -Tuula Paavilainen (SDP) -

#### Kansanedustajat ja eduskuntapuhe (1990-2021)





Timeline Wordcloud Concept network



#### Kansanedustajat ja eduskuntapuhe (1990-2021)

















## WP 5.1 and 5.2

Sanna Kumpulainen Anna Sendra Toset Farid Alijani Jaakko Peltonen



## Deliverables

- D5.1.3: Protocol for collecting workshop data
- How to ensure that information flows from the users to the developers of the infrastructure?
- <u>https://doi.org/10.5281/zenodo.10217404</u>

#### D5.2.2: Educational material

- · Collection of resources related to the tools and materials shared in the DARIAH-FI
- Dynamic document published on DARIAH-FI website
- D5.1.2: Log Data Collection and Analysis
- · Proof of concept of user-based recommendation
- https://github.com/mrgransky/DARIAH-FI



## Publications

- Sendra Toset, A., Kumpulainen, S., & Late, E. (2023). Putting the User in the Loop: Developing a Research Infrastructure for Social Sciences and Humanities Research. 1119-1121. Poster session presented at 86th Annual Meeting of the Association for Information Science and Technology, London, United Kingdom. https://doi.org/10.1002/pra2.964
- Sendra Toset, A., Late, E. & Kumpulainen, S. (2023) More than data repositories: Perceived information needs for the development of social sciences and humanities research infrastructures. Information research. [Accepted for publication]



## Future

- Tiny steps towards an iterative process
- Mature infrastructure, i.e., organization and support services