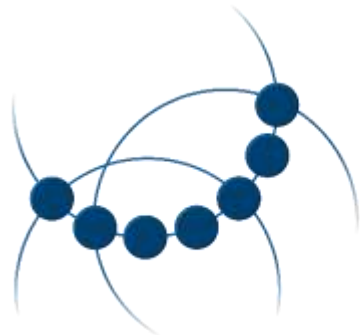


**CLARIN**  
Common Language Resources and  
Technology Infrastructure



# FIN-CLARIAH

Krister Lindén

Research Director of Language Technology

Director of the Language Bank of Finland (Kielipankki)

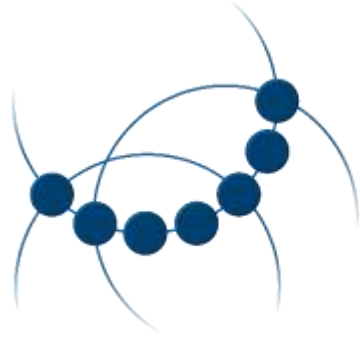
National Coordinator of FIN-CLARIN

Principal Investigator of FIN-CLARIAH



UNIVERSITY OF HELSINKI

**CLARIN**  
Common Language Resources and  
Technology Infrastructure



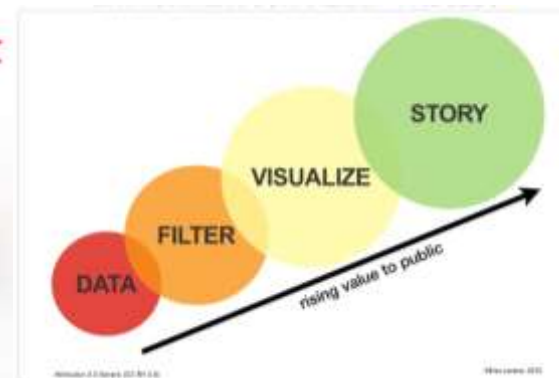
**KIELIPANKKI**  
The Language Bank of Finland

A blue logo graphic for DARIAH-EU, featuring a stylized five-petaled flower or star shape with a white center.  
**DARIAH-EU**

Digitizing, analysing and visualising data for research in  
Social Sciences and Humanities

– What is a Research Infrastructure in SSH?

# What do SSH vs. ITC need for science?



# <https://vlo.clarin.eu>

## CLARIN Virtual Language Observatory

Welcome to the VLO!

Use the **search bar** below to start searching through hundreds of thousands of language resources, or [continue](#) to browse everything and use **facets** to narrow down to your area of interest or discover new resources.

[See all records](#) [Take a quick tour](#)

about  
1.1 million  
data sets

Search

Showing #  Results per page: 10

- Use the categories below to limit the search results to those matching the selected value(s).
- Language
  - Collection
  - Resource type
  - Modality
  - Format
  - Keyword

### EXMARaLDA Demo corpus

(Part of Hamburger Zentrum für Sprachkorpora (HZSK))

A selection of short audio and video recordings in various languages to be used for instruction or demonstration of the EXMARaLDA system.; HIAT (simplified); HIAT; free comment; suprasegmental information; accentuation/stress; English translation; Standard German translation; German translation; English translation; code-switch



### The Hamburg MapTask Corpus (HAMATAC)

(Part of Hamburger Zentrum für Sprachkorpora (HZSK))

Audio and two video recordings of map tasks with adult L2 users of German and one L1 speaker. The speakers' L1 and their L2 proficiencies vary. The maps used for the tasks are available.; orthographic transcription/simplified HIAT; Fine-grained part of speech tagging using TreeTagger and the STTS tagset.; superordinate...



# CLARIN Virtual Language Observatory

Welcome to the VLO!

Use the **search bar** below to browse everything and use filters to narrow down your search.

[See all records](#) [Take a tour](#)

Showing all 1640603 records

Use the categories below to limit the search results to those matching the selected value(s).

- Language
- Collection
- Resource type
- Modality
- Format
- Keyword



- Computer-mediated communication corpora
- Historical corpora
- L2 learner corpora
- Newspaper corpora
- Parallel corpora
- Parliamentary corpora
- ...

speakers' L1 and their L2 proficiencies vary. The maps used for the tasks are available, orthographic transcription/simplified HIAT; Fine-grained part of speech tagging using TreeTagger and the STTS tagset; superordinate...

# Resource families & WGs

## Corpora

- [Computer-mediated communication corpora](#)
- [Corpora of academic texts](#)
- [Historical corpora](#)
- [L2 learner corpora](#)
- [Literary corpora](#)
- [Manually annotated corpora](#)
- [Multimodal corpora](#)
- [Newspaper corpora](#)
- [Parallel corpora](#)
- [Parliamentary corpora](#)
- [Reference corpora](#)
- [Spoken corpora](#) , ...

## Lexical Resources

- [Lexica](#)
- [Dictionaries](#)
- [Conceptual Resources](#)
- [Glossaries](#)
- [Wordlists](#) , ...

## Tools

- [Normalization](#)
- [Part-of-speech tagging and lemmatization](#)
- [Named entity recognition](#)
- [Tools for sentiment analysis](#) , ...

**FAIR** - Findable, Accessible, Interoperable, Reusable

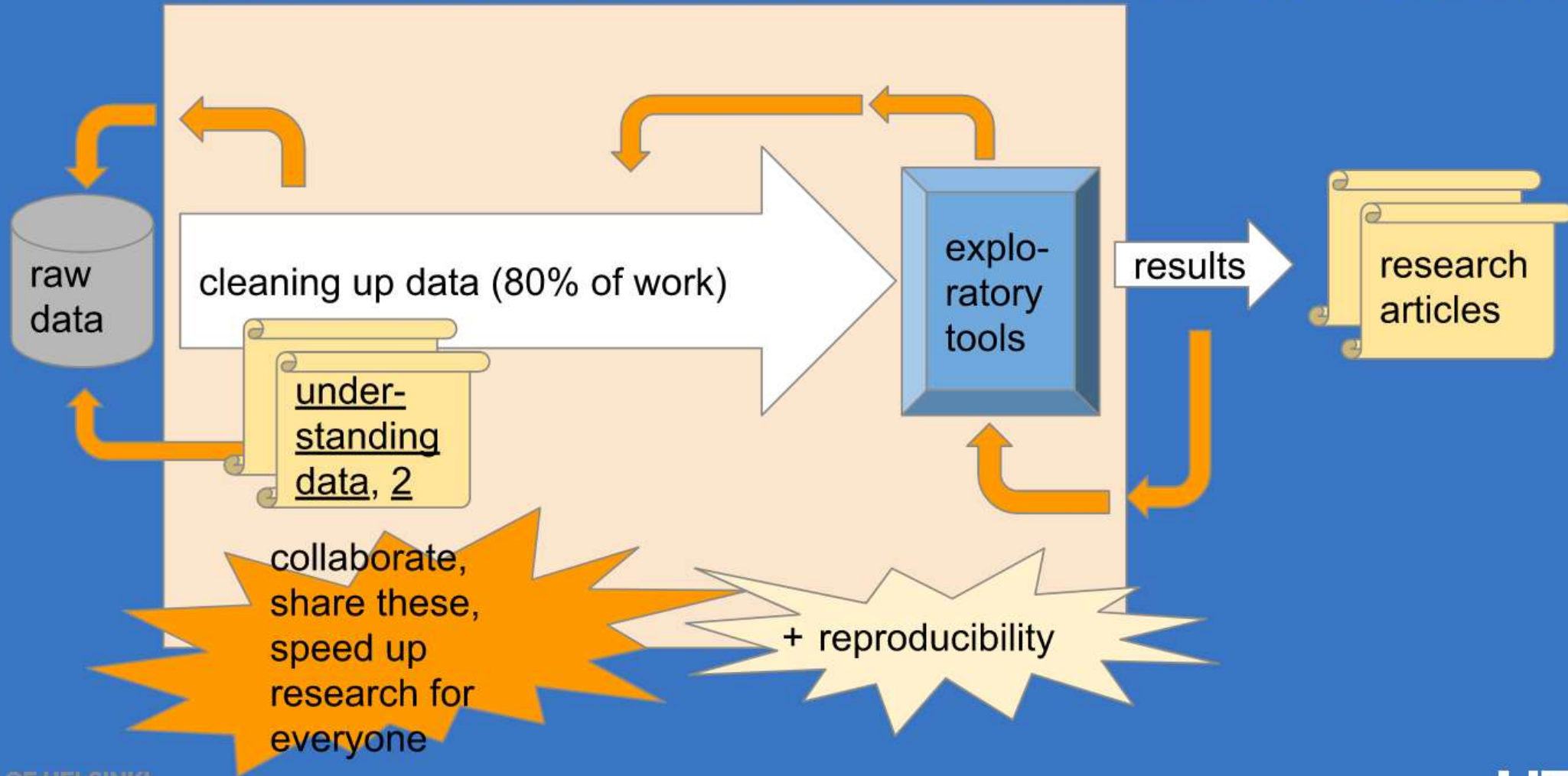
Resources	2022
<b>Text</b>	
Magazines and newspapers 1770- (NLF and Web publ.)	15 Gw
Social media and similar sources 2000- (Suomi24, Ylilauta, ...)	5 Gw
Literature and manuscripts (Gutenberg, Fennica, archives)	70 Mw
<b>Speech</b>	
Video sessions from the Finnish Parliament 2008-2020	3000 h
Dialect and everyday speech (Kotus, Turku, Donate Speech)	4500 h
Sign language resources (Aalto, Kuurojen liitto)	100 h
<b>Multilingual and Other Resources</b>	
Multilingual Resources (EuroParl, laws, Bible, Opus, ...)	5 Gw
Learner's resources (Oulu, Jyväskylä, Kotus, Aalto)	5 Mw
Open source lexicons and terminologies (Helsinki, Tromssa)	400 Kw

Kielipankki has approx. 25 GW in >1400 databases

Resource families in Kielipankki



# Leverage collaboration, open science workflows to reduce individual workload





# Cleaning and Exploration

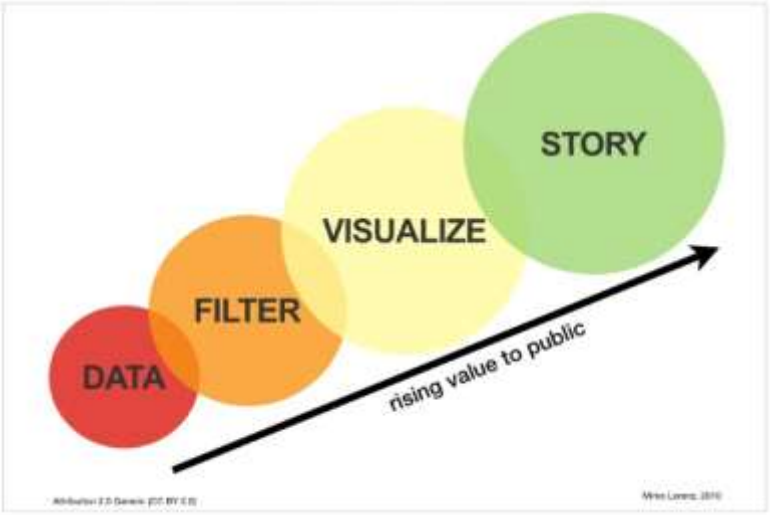
## Cleaning tools

- **Ingestion**
  - Text: Suomi24, web publications, blogs, ...
  - Images: Clay tablets, magazines , manuscripts, ...
  - Audio: Radio programs, ...
  - Video: Parliamentary and news videos, ...
- **Conversion**
  - html2txt
  - OCR & HTR
  - ASR
  - Video2txt
  - ...
- **Structuring**
  - Boundary detection: speech act, sentence, paragraph, ...

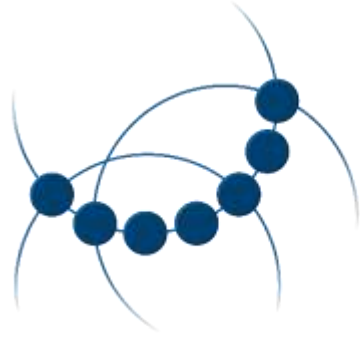
## Exploratory tools

- **Annotation**
  - Picture element annotation
  - Language identification
  - Base form annotation
  - Named-entity annotation
  - Linguistic annotation
  - Sentiment & topic annotation
- **Feature aggregation**
  - Feature counting
  - Feature vectors
  - Clustering & linking
- **Overview & visualisation**
  - Concordance
  - Trend diagrams
  - Location maps
  - Topic modelling
  - Network analysis
  - Summarization
- **Reading & understanding**
  - Subset extraction
  - Full-text, picture & video viewers

# Infrastructure for digitizing, analysing and visualising data for research in SSH



**CLARIN**  
Common Language Resources and  
Technology Infrastructure



**KIELIPANKKI**  
The Language Bank of Finland

 **DARIAH-EU**

**FIN-CLARIAH**

**International Background**



# CLARIN ERIC

European Research Infrastructure Consortium  
founded on February 29, 2012

<https://www.clarin.eu> / NL



SWE-CLARIN



INSTITUT FÜR  
DEUTSCHE SPRACHE



CLARIN-D

[ ... ]



*International cooperation  
and sharing of resources*

## Member countries:

- 22 members
- 2 observers

## The Netherlands

- Austria
- Belgium
- Bulgaria
- Croatia
- Cyprus
- Czech Republic
- Denmark
- Estonia
- Finland
- Germany
- Greece
- Hungary
- Iceland
- Italy
- Latvia
- Lithuania
- Norway
- Poland
- Portugal
- Slovenia
- Sweden

## Observers:

- South Africa
- UK
- USA / CMU



# DARIAH ERIC

European Research Infrastructure Consortium

founded on August 6, 2014

<https://www.dariah.eu/>

*WG: Theatralia*

*WG: Digital Numismatics*

*WG: Women writers  
in history*

*[ ... ] WG: Visual media  
and interactivity*

*Empower research communities  
with digital methods*

## Member countries:

- 19 members
- 8 coop. partners

### France

Austria  
Belgium  
Bulgaria  
Croatia  
Cyprus  
Czech Republic  
Denmark  
Germany  
Greece  
Ireland  
Italy  
Luxembourg  
Malta  
Netherlands  
Poland  
Portugal  
Serbia  
Slovenia

### Coop. partners

Finland, Hungary  
Norway, Romania  
Slovakia, Sweden  
Switzerland  
United Kingdom

# Resource families & WGs



## Corpora

- [Computer-mediated communication corpora](#)
- [Corpora of academic texts](#)
- [Historical corpora](#)
- [L2 learner corpora](#)
- [Literary corpora](#)
- [Manually annotated corpora](#)
- [Multimodal corpora](#)
- [Newspaper corpora](#)
- [Parallel corpora](#)
- [Parliamentary corpora](#)
- [Reference corpora](#)
- [Spoken corpora](#)

## Lexical Resources

- [Lexica](#)
- [Dictionaries](#)
- [Conceptual Resources](#)
- [Glossaries](#)
- [Wordlists](#)

## Tools

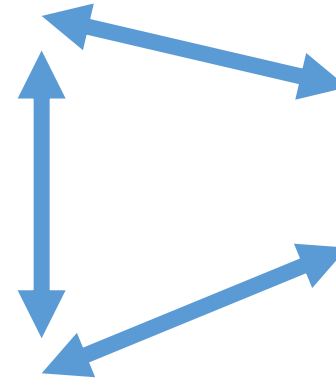
- [Normalization](#)
- [Named entity recognition](#)
- [Part-of-speech tagging and lemmatization](#)
- [Tools for sentiment analysis](#)

**FAIR** - Findable, Accessible, Interoperable, Reusable

# FIN-CLARIAH – Partners

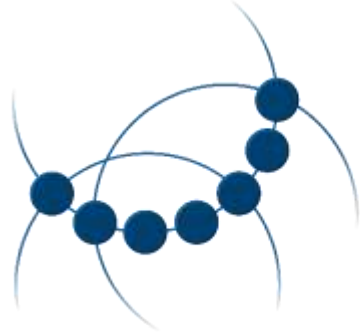
- University of Helsinki
- CSC – IT Center for Science
  
- Aalto University
- Tampere University
- University of Eastern Finland
- University of Jyväskylä
- University of Oulu
- University of Turku
- University of Vaasa
- Institute for the Languages of Finland
- National Archives
- National Library

Coordination, access to and long-term storage of large **centrally provided resources and tools**



Provide access to **resources and tools developed locally** by individual researchers or research groups

**CLARIN**  
Common Language Resources and  
Technology Infrastructure



**KIELIPANKKI**  
The Language Bank of Finland

A blue graphic logo for DARIAH-EU, featuring a stylized flower or star shape with five petals, each containing a white arrow pointing outwards.  
**DARIAH-EU**

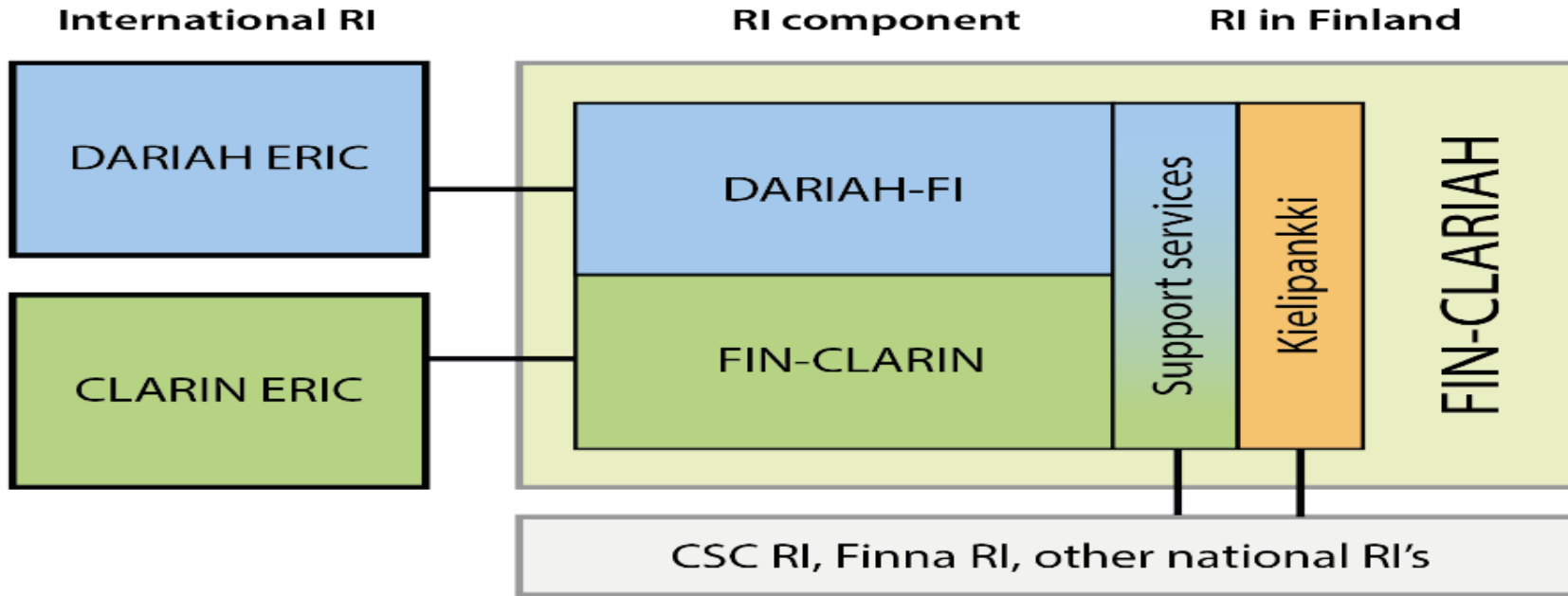
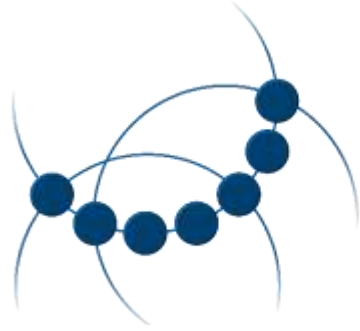
# Roadmap 2021-2030

In 2020, FIN-CLARIAH was accepted on the national roadmap comprising the components: FIN-CLARIN and DARIAH-FI



**CLARIN**

Common Language Resources and  
Technology Infrastructure



Research Infrastructure for the Humanities and Social Sciences

**FIN-CLARIN**

Krister Lindén, chair, PI  
Mikko Tolonen, vice chair



UNIVERSITY OF HELSINKI

# FIN-CLARIAH Roadmap for 2022-2030

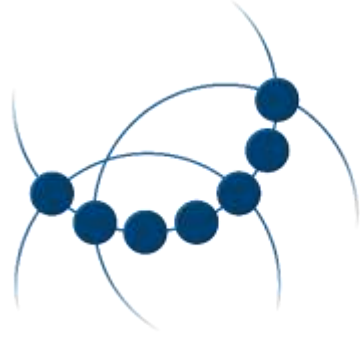
Module 5	Information Interaction (IIA)	UTA
Module 4	Analysing structured data	UHEL
Module 3	Structuring data	UHEL
Module 2	Language Research Infrastructure (LRI)	UHEL
Module 1	Natural Language Processing (NLP)	UHEL

**DARIAH-FI**

**FIN-CLARIN**

**FIN-CLARIAH**

**CLARIN**  
Common Language Resources and  
Technology Infrastructure



**KIELIPANKKI**  
The Language Bank of Finland

 **DARIAH-EU**

# Project 2022-2023

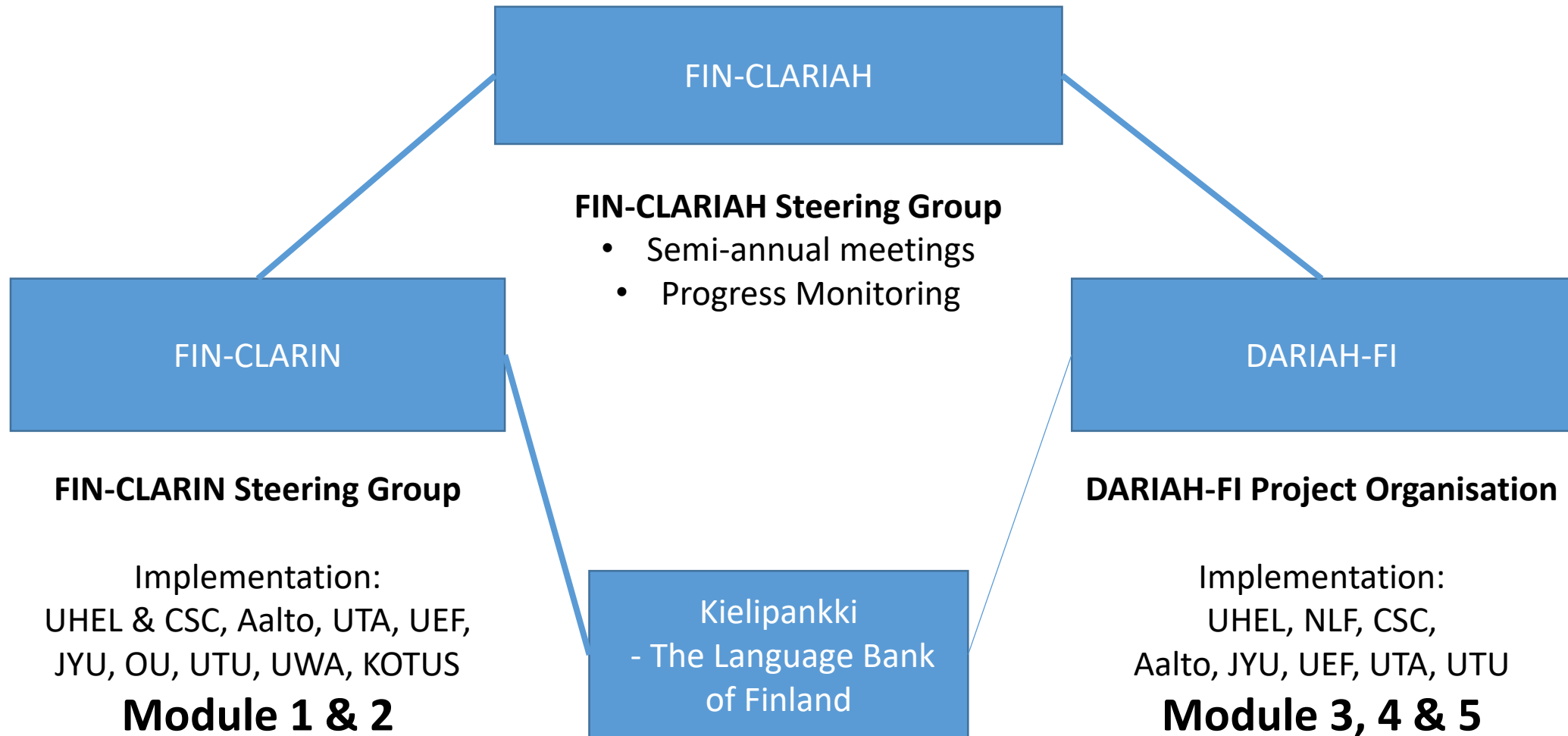
Applied for Academy of Finland infrastructure funding in 2021  
for FIN-CLARIAH consolidation

# Focus areas

1. Reach beyond processing of **spoken standard Finnish into colloquial speech** (Goal 1) [*FIN-CLARIN*]
2. Cater to a broad range of SSH research needs for processing **unstructured text** (Goal 2) [*FIN-CLARIN/DARIAH-FI*]
3. Facilitate research based on **metadata** (Goal 3) [*DARIAH-FI*]

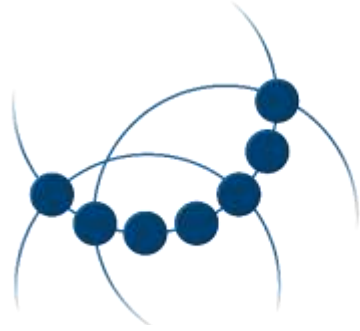
<b>Module 1: Natural Language Processing</b>	Goal 1	Goal 2	Goal 3
<b>Module 2: Language Research Infrastructure</b>	Goal 1	Goal 2	Goal 3
<b>Module 3: Structuring data</b>		Goal 2	Goal 3
<b>Module 4: Analysing structured data</b>		Goal 2	Goal 3
<b>Module 5: Information Interaction</b>	Goal 1	Goal 2	Goal 3

# RI organisation



**CLARIN**

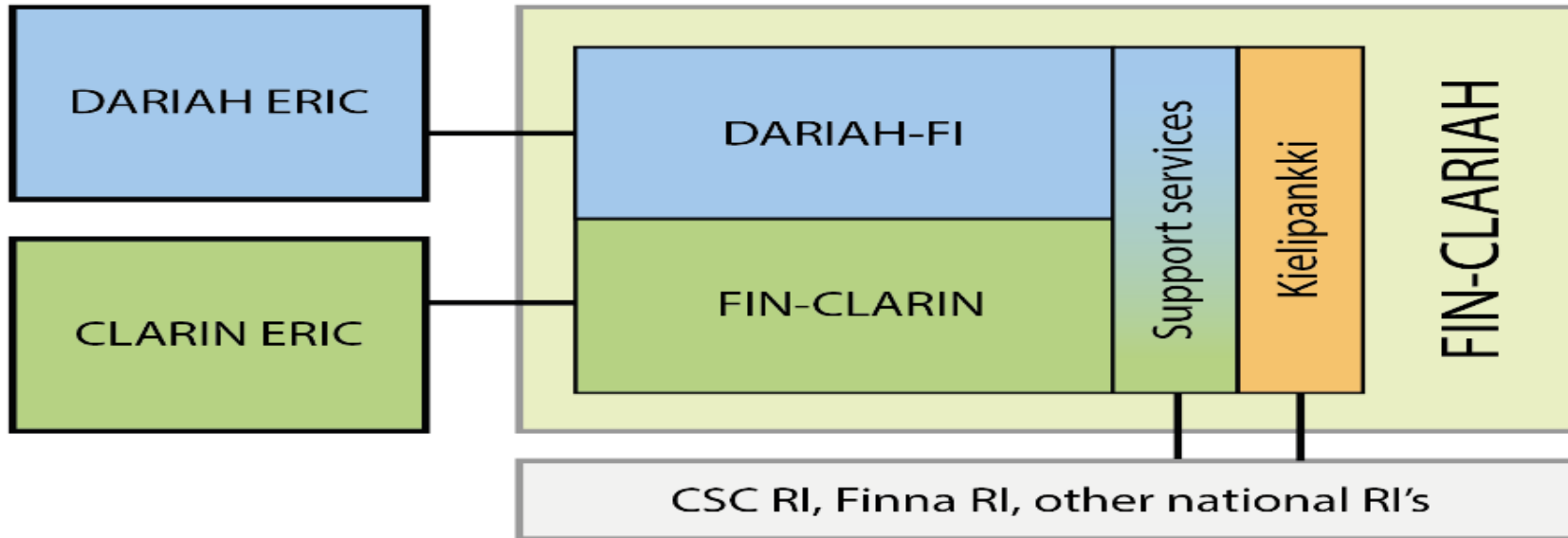
Common Language Resources and  
Technology Infrastructure



**International RI**

**RI component**

**RI in Finland**



Research Infrastructure for the Humanities and Social Sciences

**FIN-CLARIN**

Krister Lindén, chair, PI  
Mikko Tolonen, vice chair



UNIVERSITY OF HELSINKI