# What Linked Data can tell about Geographical Trends in Finnish Fiction Literature

## − *Using the BookSampo Knowledge Graph in Digital Humanities*

Telma Peura[1,2], Petri Leskinen[1] & Eero Hyvönen[1,2]

[1] Semantic Computing Research Group (SeCo),
Aalto University, Finland;
[2] Helsinki Centre for Digital Humanities (HELDIG),
University of Helsinki, Finland

# Introducing BookSampo (1)

**Part of the Sampo series**

- building **a national ontology** (FinnONTO) and infrastructure for **cultural heritage** (CH) data that can be accessed through a semantic portal (Hyvönen, 2020)

- open source

- based on semantic web, residing on a SPARQL endpoint

- unifying and connecting heterogenous data silos of CH data

- CultureSampo (2009), BookSampo (2011), WarSampo (2015), BiographySampo (2018).

- Three components of model building

  1) LD creation and publishing model

  2) an interface with multiple perspectives

  3) two-step filtering cycle

Table 1

Sampo Model Principles P1–P6

| P1 | Support collaborative data creation and publishing |
| --- | --- |
| P2 | Use a shared open ontology infrastructure |
| P3 | Make clear distinction between the LOD service and the user interface (UI) |
| P4 | Provide multiple perspectives to the same data |
| P5 | Standardize portal usage by a simple filter-analyze two-step cycle |
| P6 | Support data analysis and knowledge discovery in addition to data exploration |

Hyvönen et al. (2022)
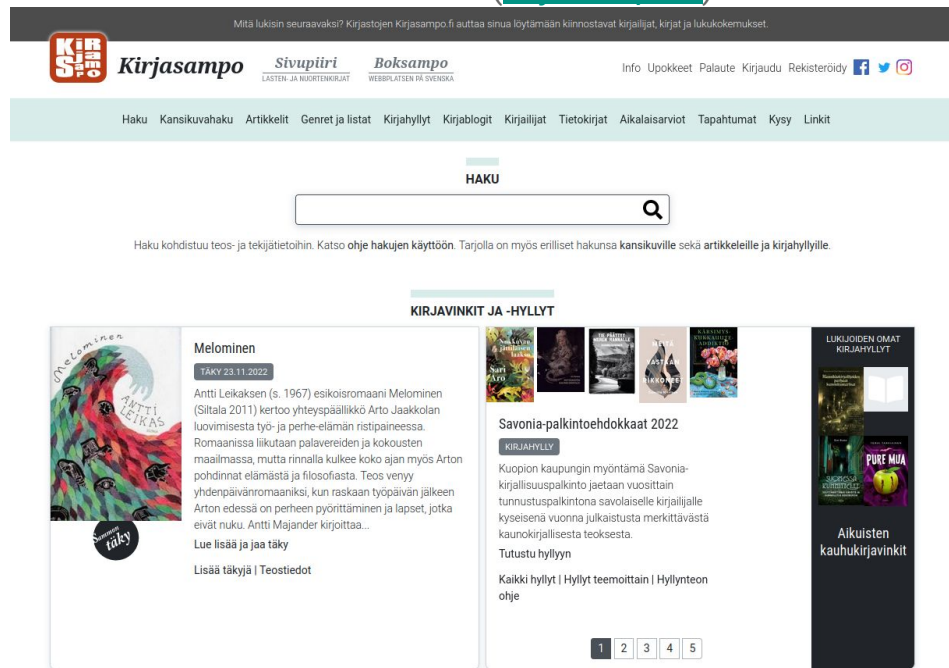
# What does BookSampo contain?

| Type | Number of Instances |
|---|---|
| Literary Works | 93,000 |
| Editions | 127,000 |
| Book Covers | 27,000 |
| Fictional Characters | 19,000 |
| Contemporary Reviews | 15,000 |
| Weblinks | 10,000 |
| Literary Series | 2,900 |
| Literary Awards | 2,700 |
| Literary Award Series | 200 |
| Movies | 1,100 |
| People (e.g. Authors) | 29,000 |
| Author's Pictures | 2,600 |
| Publishers | 2,600 |

| Type | Nr. of Instances |
|---|---|
| Literary Works | 207,771 |
| Editions | 213,161 |
| Edition Parts | 79,447 |
| All Editions | 285,518 |
| Book Covers | 112,971 |
| Main Characters | 45,397 |
| Contemporary Reviews | 14,644 |
| Weblinks | 25,017 |
| Literary Series | 8,374 |
| Literary Awards | 6,310 |
| Literary Award Series | 290 |
| Movies | 2,010 |
| People (e.g. Authors) | 62,207 |
| Author Pictures | 4,140 |
| Publishers | 5,455 |

Mäkelä et al. (2013)

August 2022

# Introducing BookSampo (2)

- a collaboration between FInnish public libraries and SW researchers
  - started in 2007
  - published 2011
- original data came from a dump from the Helsinki metropolitan area libraries
- new material automatically converted every night from BTJ Finland Ltd (Mäkelä et al. 2013)
- Kaunokki ontology for fiction literature
- Maintained by a team of librarians
  - improves the quality of the metadata
- contains practically all literature of Finnish public libraries (Mäkelä et al. 2011)
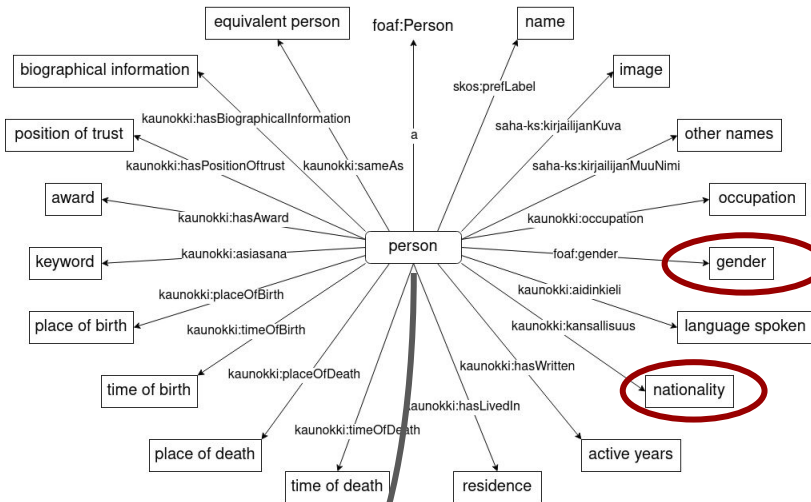  - a good data source for DH research

Current interface ([kirjasampo.fi](kirjasampo.fi)):

# BookSampo Data Model

**Author**

foaf:Person

- equivalent person
- biographical information — kaunokki:hasBiographicalInformation
- position of trust — kaunokki:hasPositionOftrust
- award — kaunokki:hasAward
- keyword — kaunokki:asiasana
- place of birth — kaunokki:placeOfBirth
- time of birth — kaunokki:timeOfBirth
- place of death — kaunokki:placeOfDeath
- time of death — kaunokki:timeOfDeath
- residence — kaunokki:hasLivedIn
- name — skos:prefLabel
- image — saha-ks:kirjailijanKuva
- other names — saha-ks:kirjailijanMuuNimi
- occupation — kaunokki:occupation
- gender — foaf:gender
- language spoken — kaunokki:aidinkieli
- nationality — kaunokki:kansallisuus
- active years — kaunokki:hasWritten

person — kaunokki:sameAs
a

**Physical publication**

kaunokki:fyysinen_teos

- other author — kaunokki:toimittaja
- title — skos:prefLabel
- illustrator — kaunokki:kuvittaja
- cover — kaunokki:kansikuva
- translator — kaunokki:kaantaja
- page count — kaunokki:sivuLkm
- original version — kaunokki:onEnsimmainenVersio
- language — kaunokki:kieli
- publisher — kaunokki:hasPublisher
- publication year — kaunokki:ilmestymisvuosi

publication
a

**Abstract work**

kaunokki:romaani

- physical work — kaunokki:manifest_in
- description — dce:description
- ISBN — sch:isbn
- time of story — kaunokki:hasTimeOfStory
- concrete place — kaunokki:worldPlace
- setting — kaunokki:paikka
- character — kaunokki:toimija
- main character — kaunokki:paahenkilo
- award — kaunokki:onPalkinto
- original language — kaunokki:alkukieli
- keyword — kaunokki:asiasana
- theme — kaunokki:teema
- genre — kaunokki:genre
- author — kaunokki:tekija
- title — skos:prefLabel

novel
a

(Annastiina Ahola)

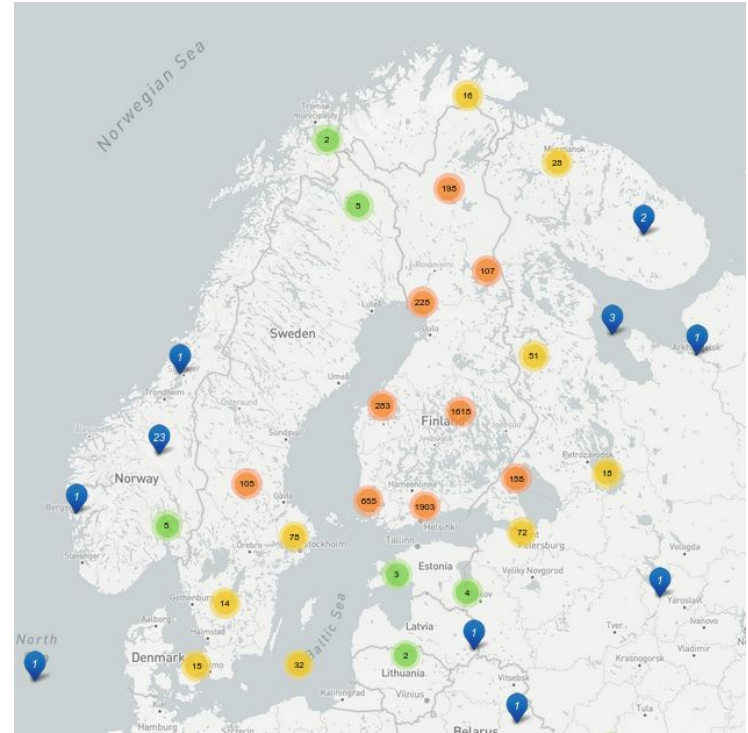# Case Study: Geographical diversity in novels

Distant Reading (Moretti, 2013)

- adopting a quantitative perspective to literature
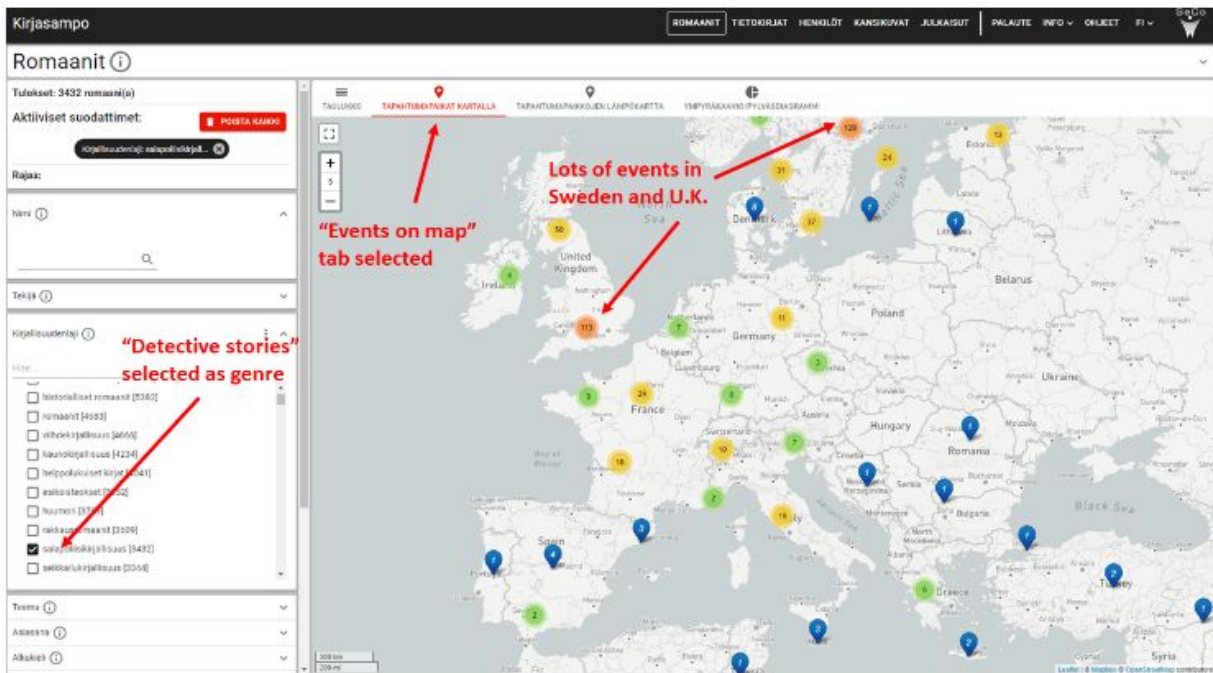- capturing large patterns in the literary production

Literary Geography (e.g. Evans & Wilkens, 2018)

- in terms of author nationalities and languages
  - extra-textual
- in terms of novel settings
  - contents

**How has the geography of novels published in Finnish developed in the past 50 years?**

# BookSampo for research



Example: detective stories on a map (Hyvönen et al., 2022)

**Observations on data quality**

Lack of coherence:

- over 600 URIs for genres, and over 400 themes
    - without a clear hierarchy
- "everything that helps in searching"
- e.g. How to group war literature?
    - 2nd WW ⇒ wars (events) ⇒ <u>societal events</u>
    - hostilities (warfare) ⇒ <u>societal events</u> ⇒ events
    - wartime ⇒ societal periods of time ⇒ time periods
    - the cold war ⇒ mutual action ⇒ action

# Data preprocessing and methods

| | |
|---|---|
| Finnish novels | 16568 |
| Author Gender known | 16367 |
| Author Nationality known | 16503 |
| Both known | 16311 |
| Translated novels | 17539 |
| Author Gender known | 17153 |
| Author Nationality known | 16672 |
| Both known | 16465 |

**SPARQL-querying the BookSampo SW and constructing necessary subgraphs**

- constructing a sub-graph with the wanted attributes
  - dealing with annotation errors and irregularities to simplify the query process
  - e.g. Scottish, English, Welsh ⇒ Brits; Ålandians, Swedish-speaking Finns ⇒ Finns
  - Sami people, Finnish Swedes
  - removing duplicated nationalities (several URIs)
  - unification of the language codes
    - most based on lexvo, some irregularities
- Jupyter Notebooks
  - SPARQLWrapper, RDFLib
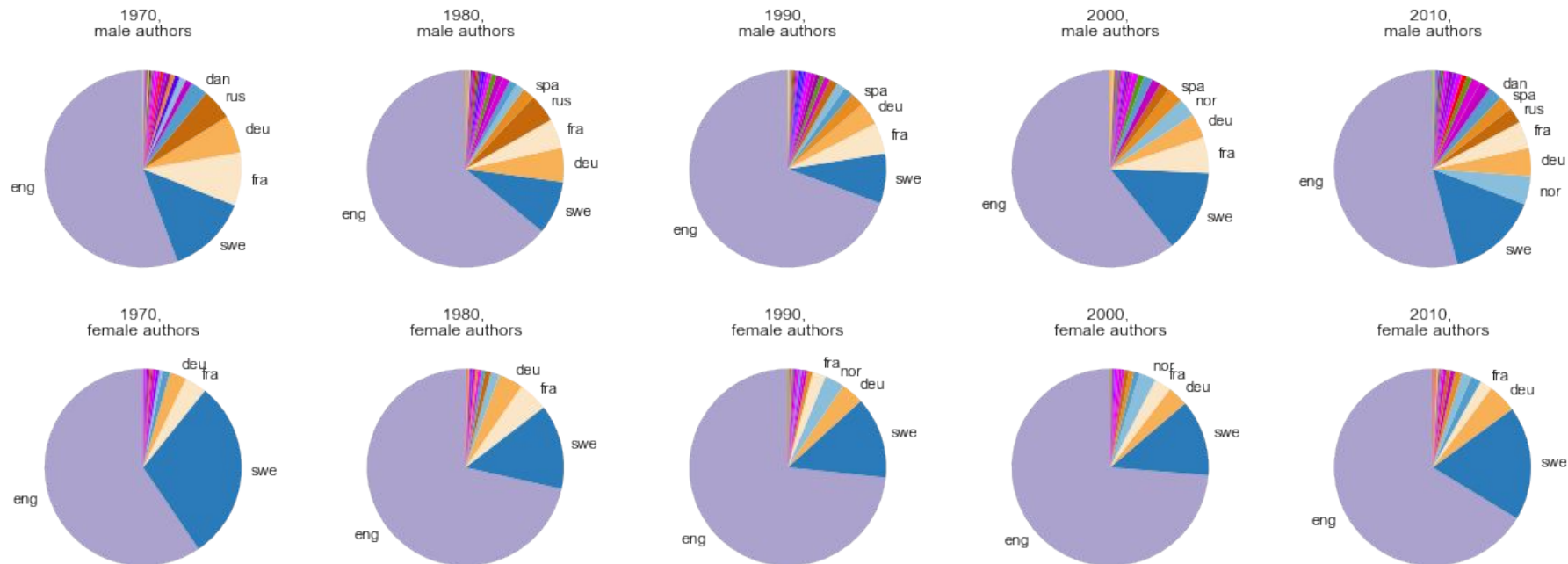- Filtering on publication year (1971-2020)

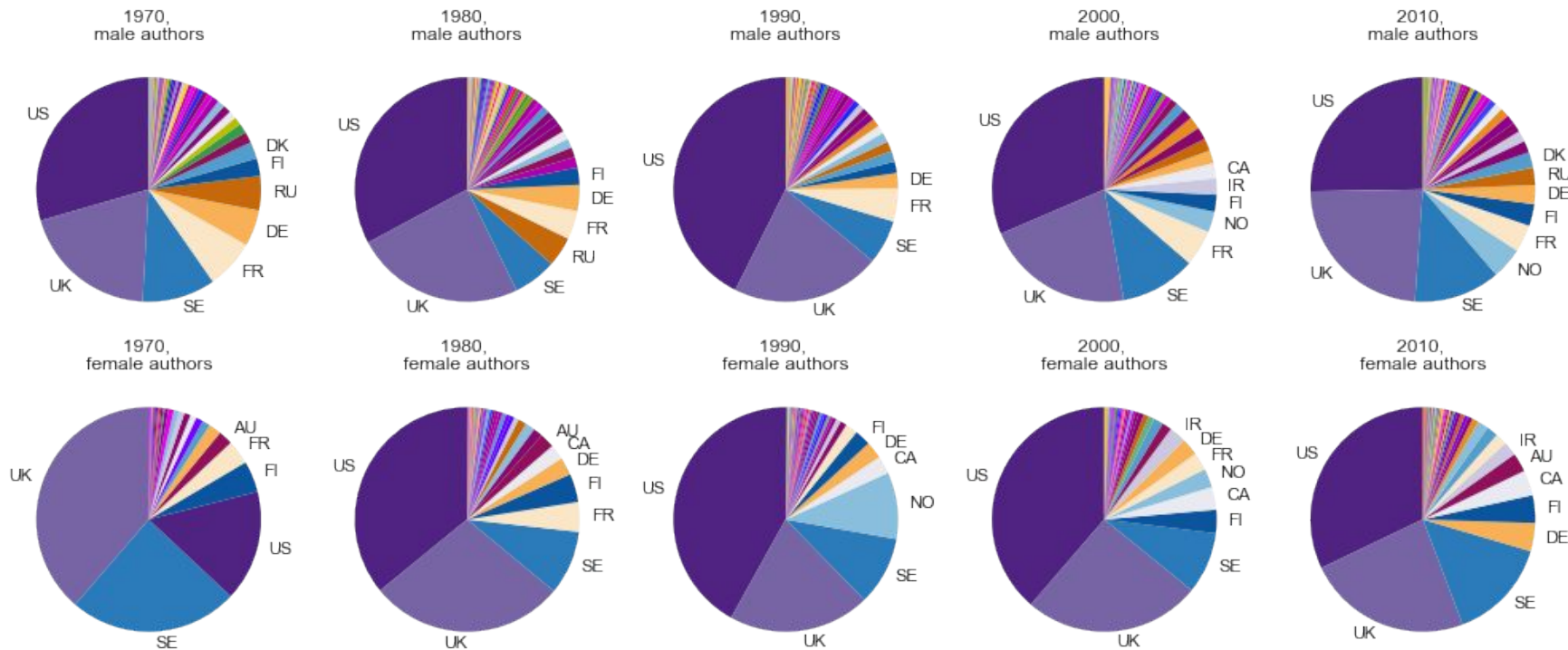# Unbalanced representation of languages and nationalities

# Unbalanced representation of languages and nationalities

# A Gender Bias in Language Diversity
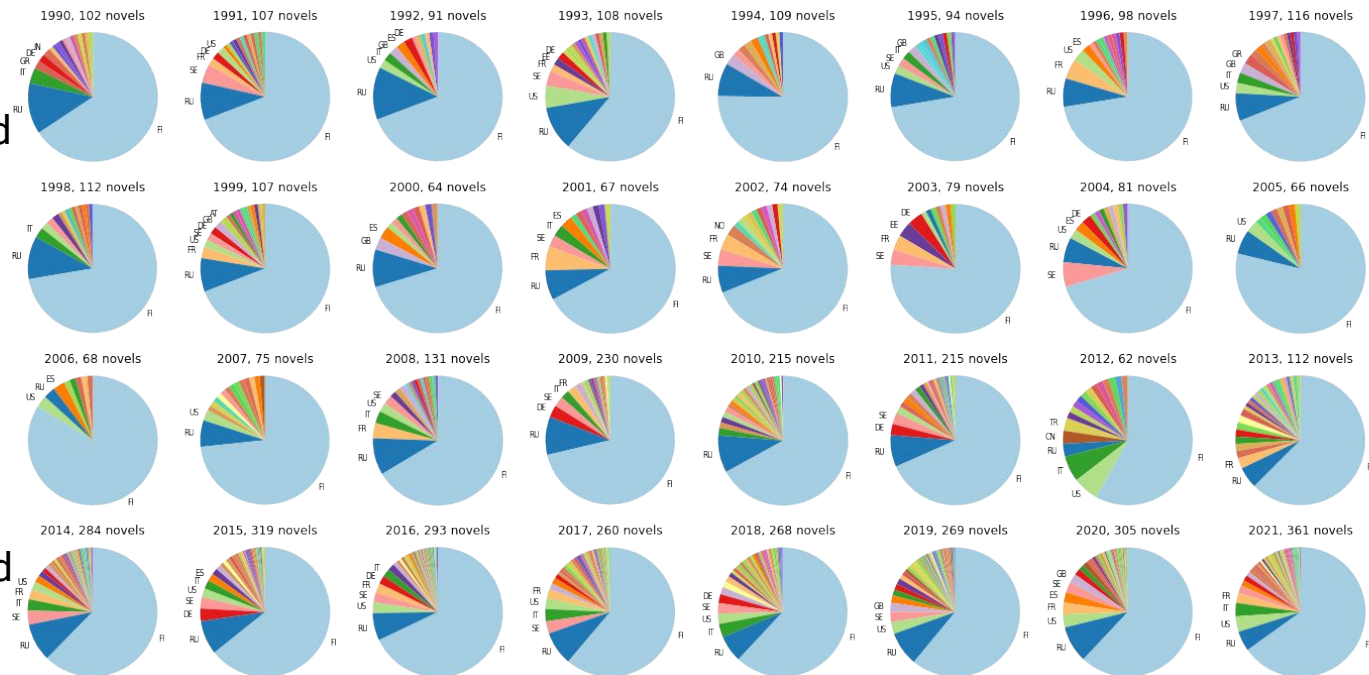
# A similar pattern in nationalities
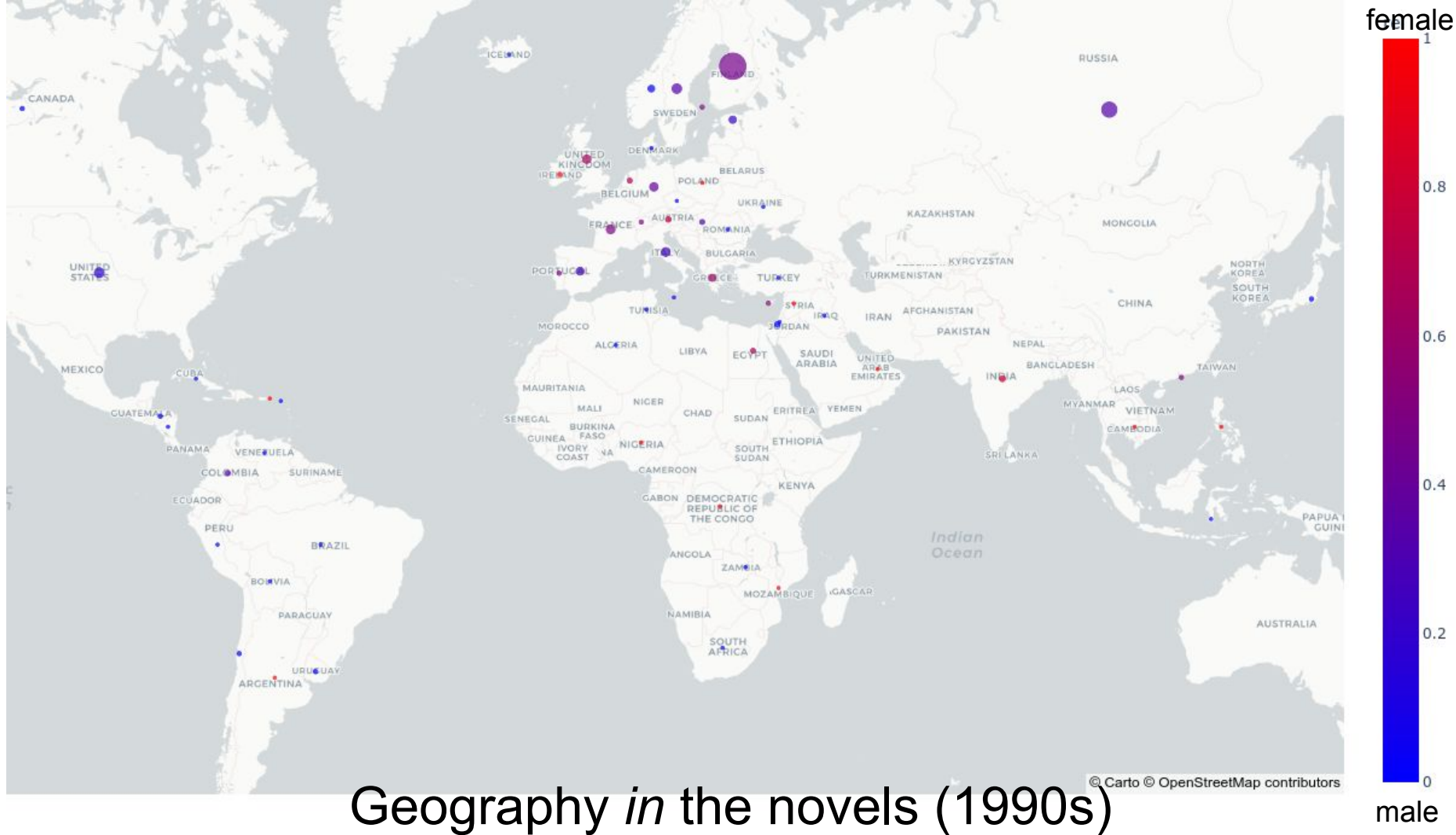
# Geography *in* the novels written by Finnish authors

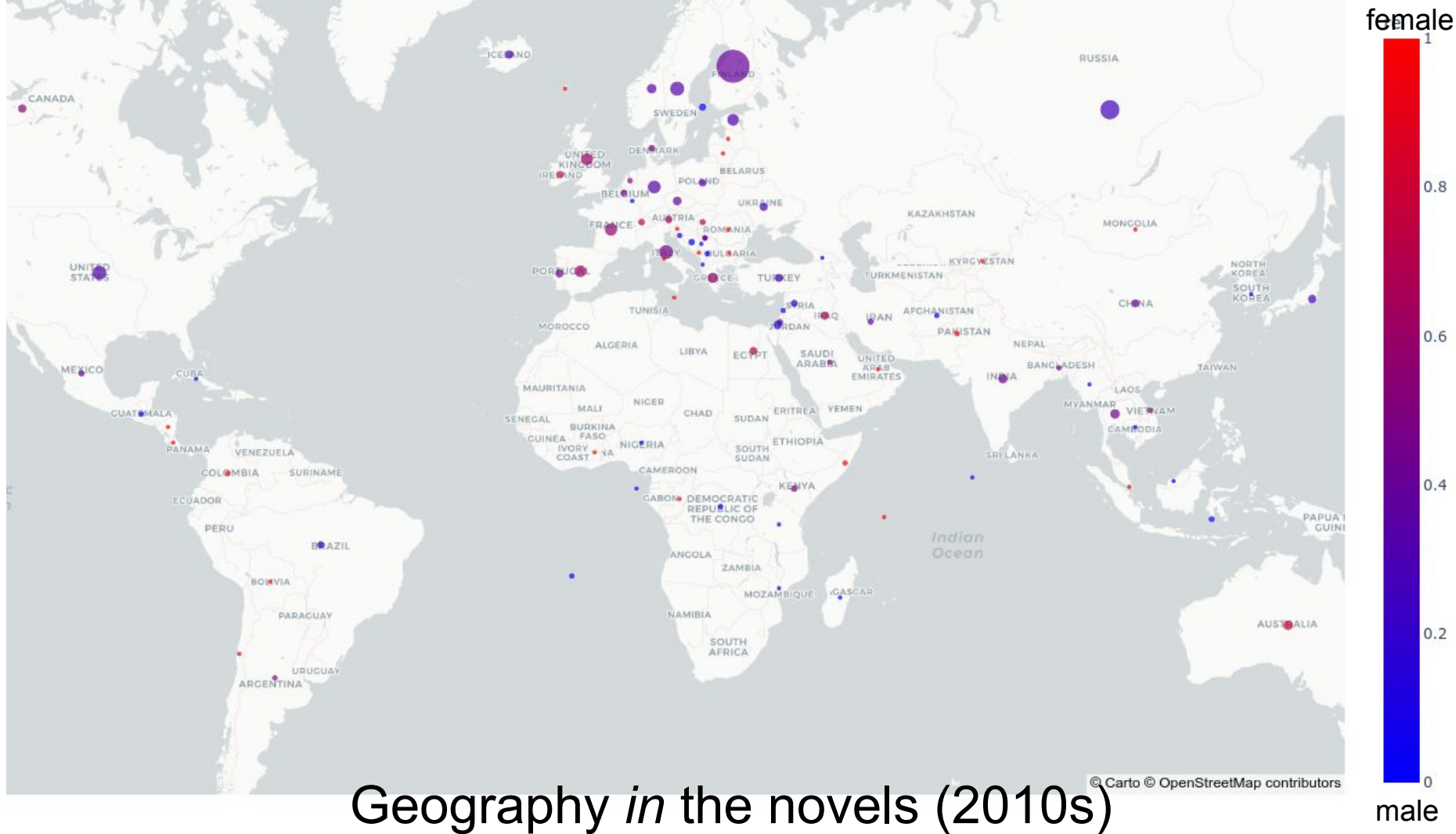Only about ⅓ of yearly publications is annotated for the geographical setting

Grouping place names by country and year

Russia the 2nd most popular country in almost every year

The proportion of Finland slightly decreasing

Geography *in* the novels (1990s)

Geography *in* the novels (2010s)

# What do authors write about when they write about Russia?

Clustering based on genre and theme keywords

1) biographical, historical war novels, memories            48/96/144
2) war literature, societal novels                          5/110/117
3) personal development, relationships, moving, psychological    42/30/72
4) thrillers, crime fiction, espionage                      6/52/61
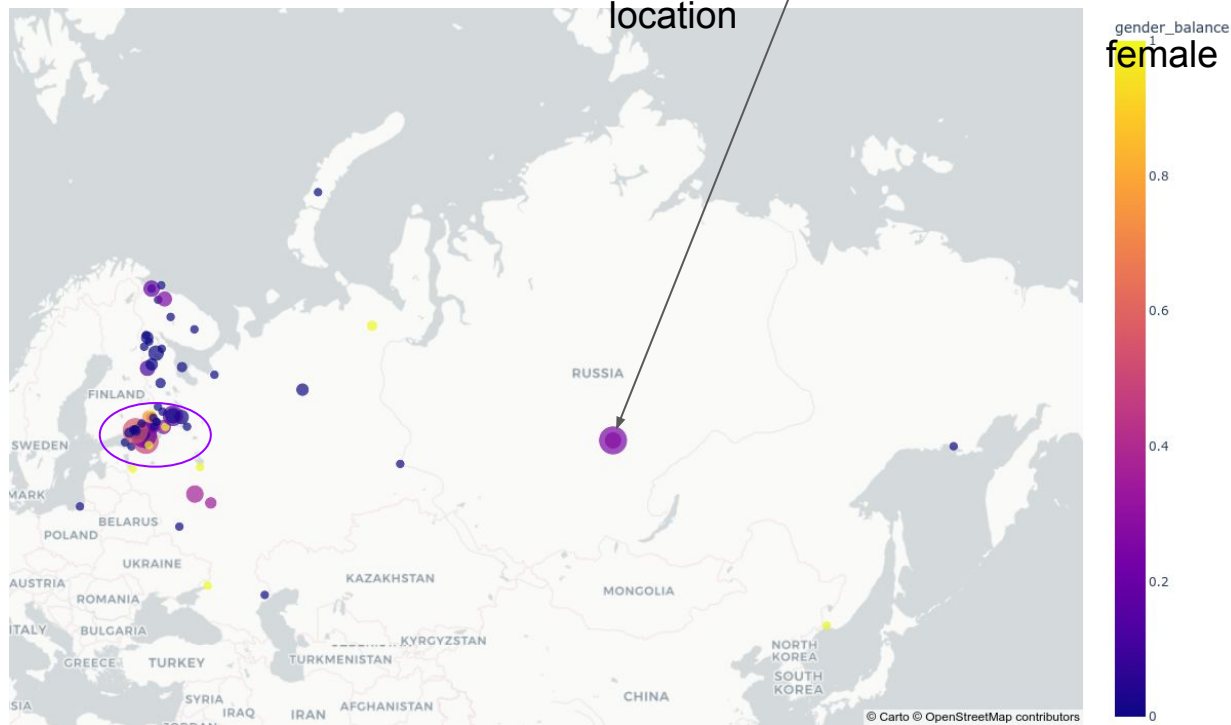5) women, historical romance                                23/5/28

⇒ *Female and male authors have different perspectives on Russia in the context of novels*

# On the map

Novels about Russia actually tell about Finland.

Focus on Karelia.

Confirms the previous claim (Hallila & Hägg, 2007) that Finns are interested in historical, national topics.

The country of Russia as the annotated place, not a specific location

# Conclusions and Future Directions

The BookSampo Knowledge Graph can be used in DH research, but it requires careful preprocessing and compromising on detailed/specific claims.

Novels seem to **become more diverse** both in terms of extra-textual indices and textual settings.

Bias in the annotations?

- Geographical setting annotated only in a third of the novels
- **When is the geographical setting seen important?**
    - E.g. genre literature
    - **When should it be annotated?** How fine-grained?

Insights from this exploratory research can help to improve the annotation practices further.

Collaboration between DH and SW researchers, and librarians is needed.

# References

Evans, E., & Wilkens, M. (2018). Nation, ethnicity, and the geography of british fiction, 1880-1940. *Journal of Cultural Analytics*, *3*(2), 11037.

Hallila, M., & Hägg, S. (2007). history and historiography in contemporary Finnish Novel. *AVAIN-Kirjallisuudentutkimuksen aikakauslehti*, (4), 74-80.

Hyvönen, E. (2020). "Sampo" Model and Semantic Portals for Digital Humanities on the Semantic Web. In *Proceedings of the Digital Humanities in the Nordic Countries 5th Conference (DHN 2020)*. CEUR-WS. org.

Hyvönen, E., Ahola, A., & Ikkala, E. (2022). BookSampo Fiction Literature Knowledge Graph Revisited: Building a Faceted Search Interface with Seamlessly Integrated Data-Analytic Tools. In *International Conference on Theory and Practice of Digital Libraries* (pp. 506-511). Springer, Cham.

Moretti, F. (2013). *Distant reading*. Verso Books.

Mäkelä, E., Hypén, K., & Hyvönen, E. (2013). Fiction literature as linked open data—The BookSampo dataset. *Semantic Web*, *4*(3), 299-306.

# Thank you! Questions?