

Metadata from ENB

ENB metadata

Abstract

Estonian National Bibliography is a metadata set that aims to collect information on all publications written in any language in Estonia and all texts written in Estonian in whichever country. The dataset has been compiled in digital format since 2002 and aggregates the work of multiple institutions and generations in collecting the publication information.

The data here has been extracted on Jan 27 2022.

This dataset presents the Estonian National Bibliography dataset in wide instead of long format used in Marc21, with some of the variables that may be useful for text-mining studies. It includes the following information

- publication ID in National Library
- time of publication (aeg)
- place of publication (koht_raw)
- partly standardized place of publication (koht)
- publisher (kirjastus_raw)
- partly standardized publisher information (kirjastus)
- title (title, subtitle, comptitle)
- author (name, date of birth,id)
- other associated authors (translator, editor etc)
- number of copies printed
- font information
- keyword
- genre
- link to fulltext

Information on coding the variables can be found here: [meta file on github]

Helpful information on the metadata available can be found here: <http://data.digar.ee/#page5>

The rules followed in adding information on older books can be found here: https://www.elnet.ee/wiki/lib/exe/fetch.php?media=kataloogimine:vanaraamatmarc_2019.pdf

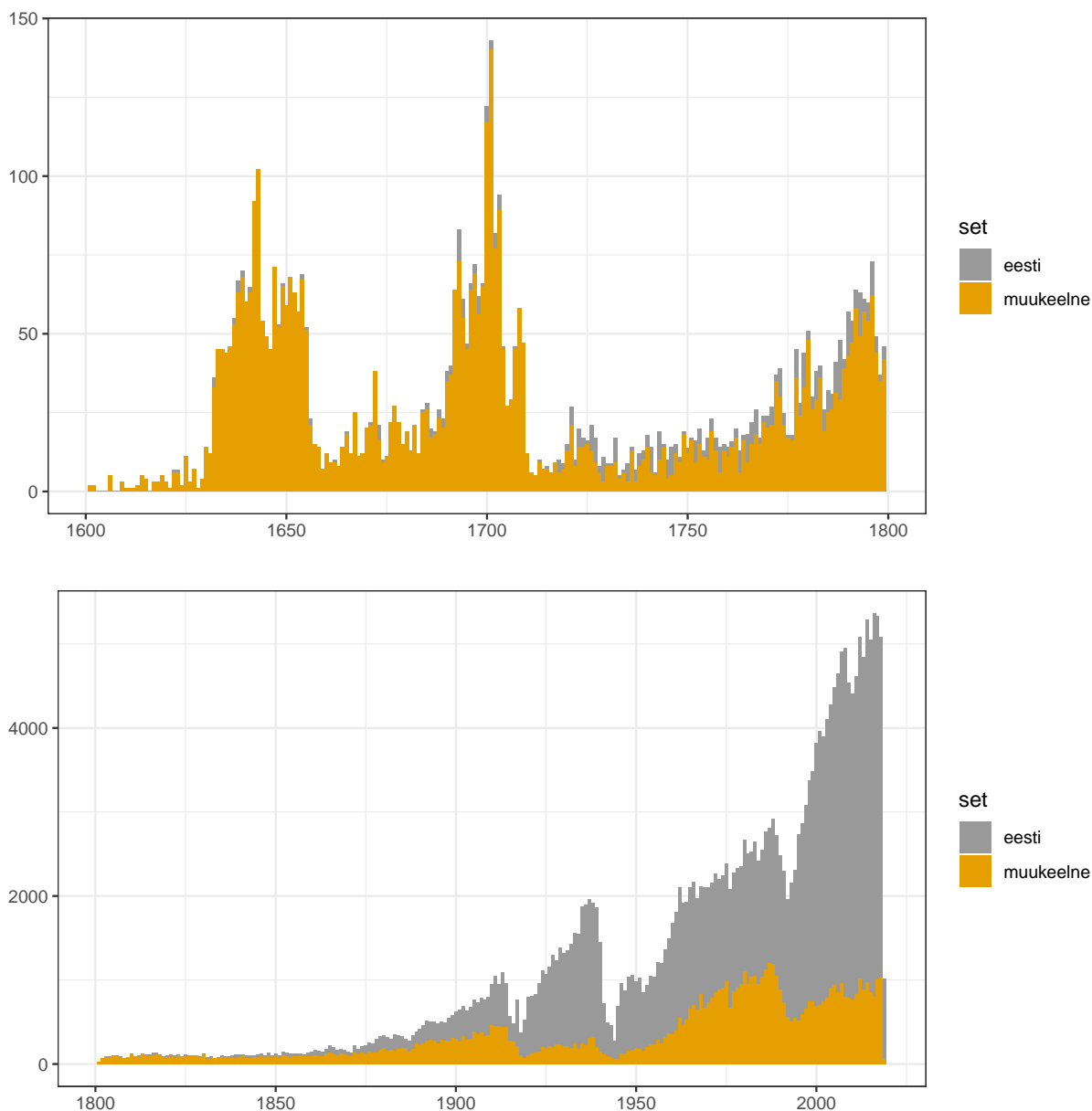
Intro

The current code organizes the ENB metadata into a tidy format: with one publication per line, and with the basic metadata information distributed into few relevant categories. This has been done to facilitate further data analytic steps, particularly when someone may be unfamiliar with Marc21 database structures.

Summary

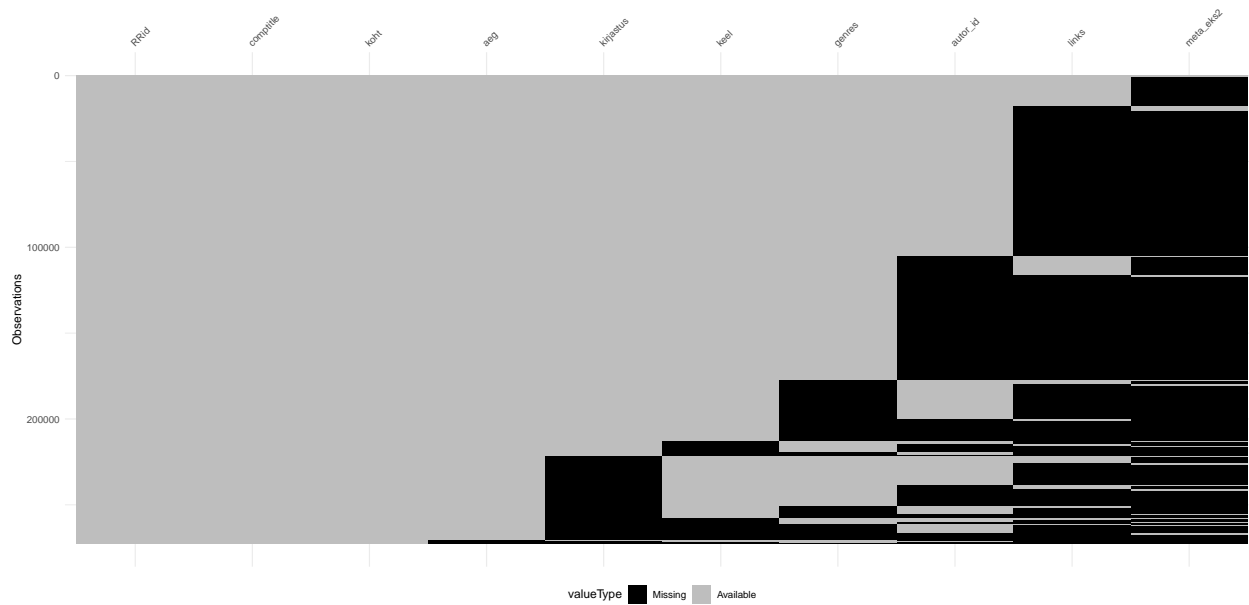
The dataset has altogether information on 272411 printed items. The coverage has been estimated to be better than 95% for all of the relevant published works.

The dataset is divided into two: works in Estonian language, and works in other languages. The two sets are displayed over the year of publication here in different colors.



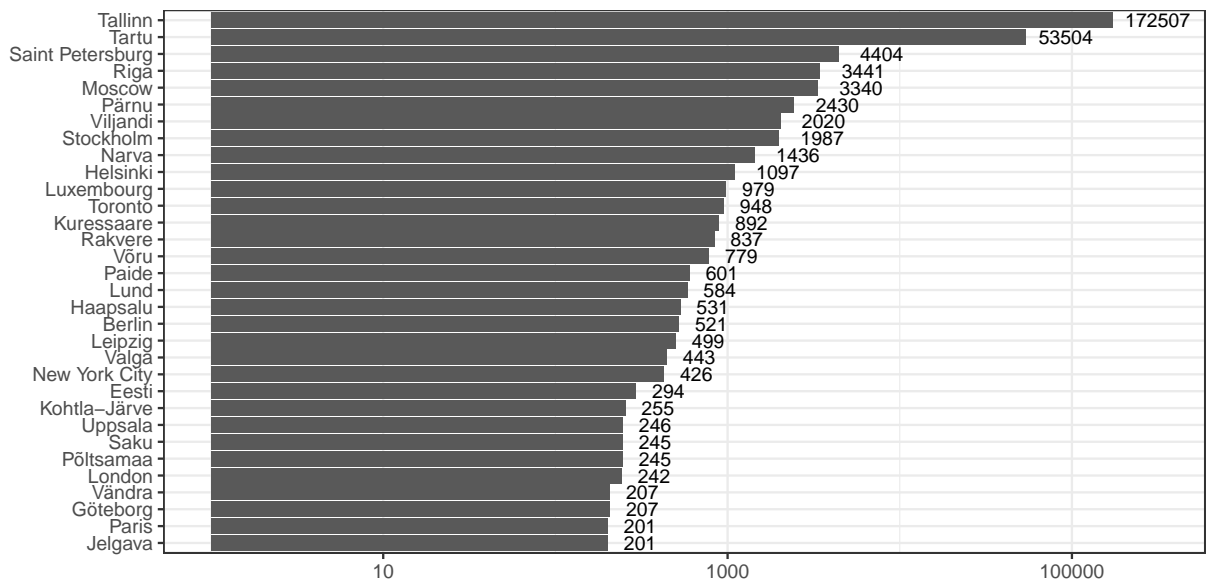
Coverage of metainformation

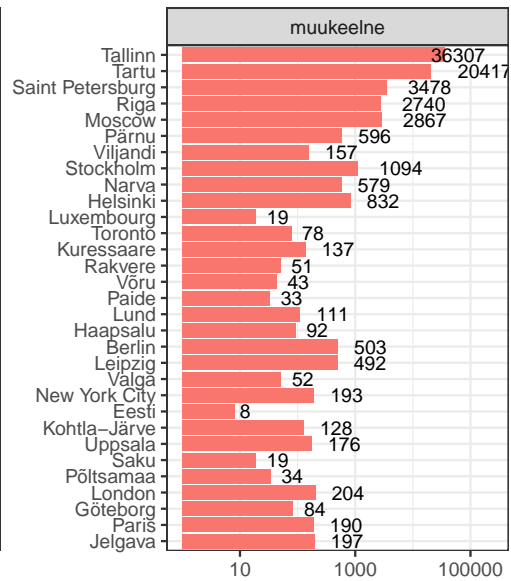
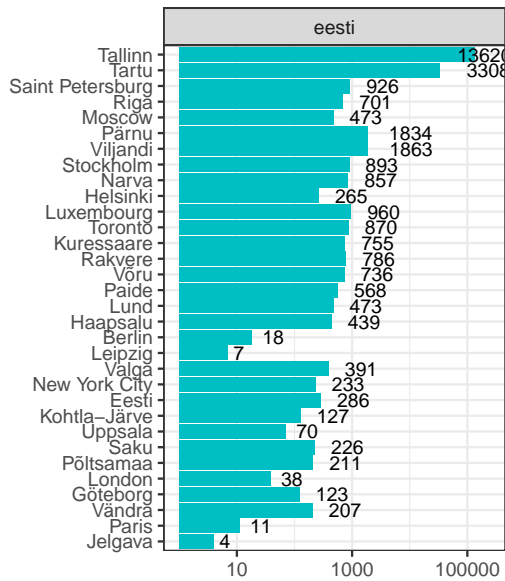
An overview of the coverage of the data is given below. Grey areas indicate datapoints that do have the information of that column, black areas the datapoints that do not have the information of that column. For example 82% of the books have information on the publisher, however just 59% of the books have information on the author. 16% of the books have a link to an online digital copy.



Cities

The city names where the works have been published have been harmonized manually and through a few heuristic algorithms. The tokens that appear more than 40 times should be mostly harmonized, while rarer tokens have been included only through algorithmic processing. Depicted below are the most common cities of publication and the number of publications in them separately for language sets.





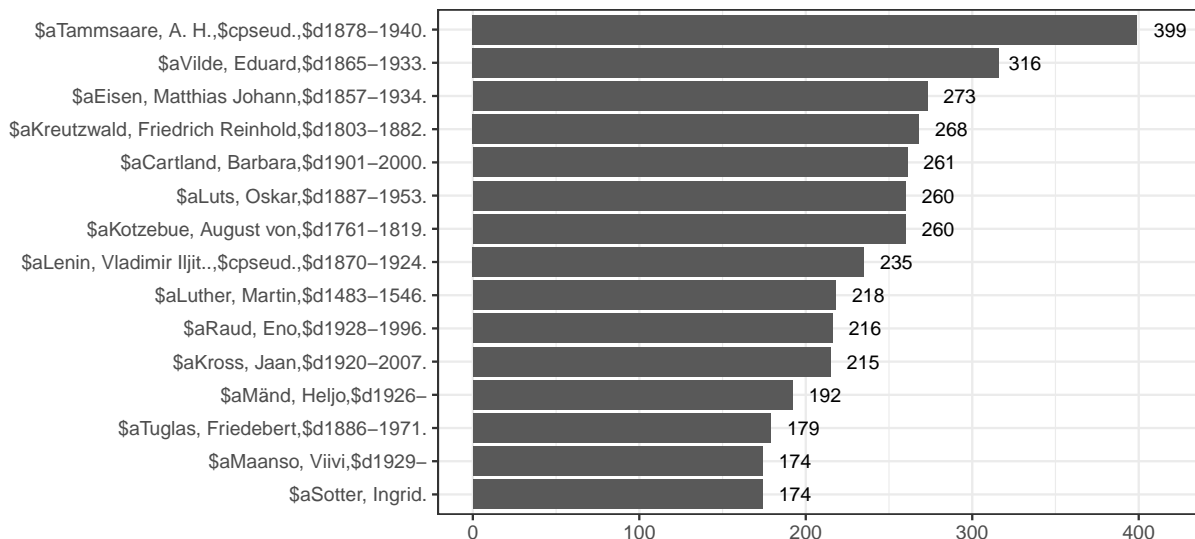
fct_rev(set)

muukeelne

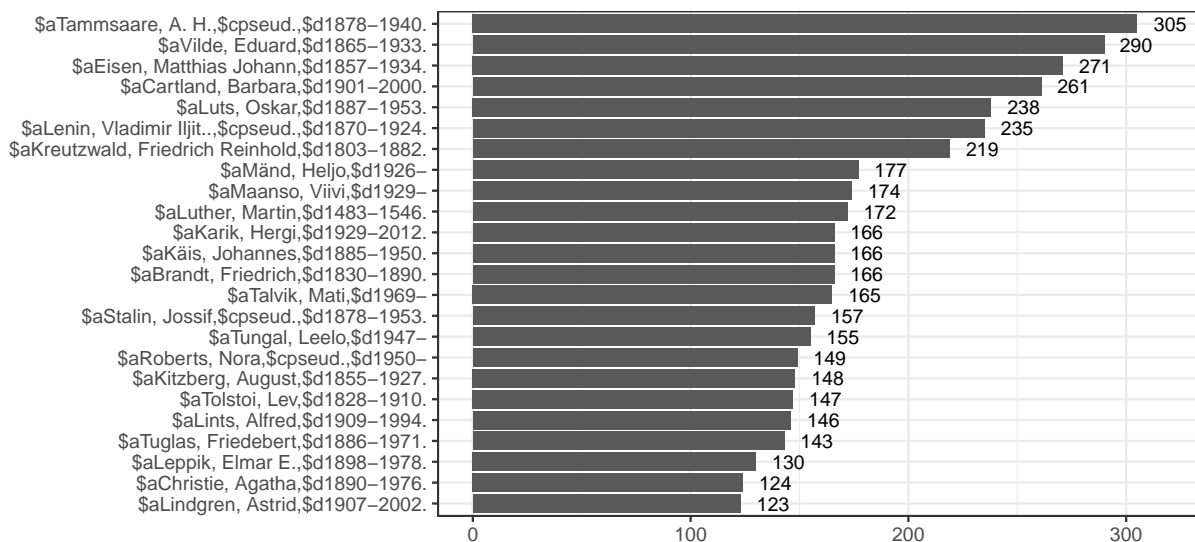
eesti

Authors

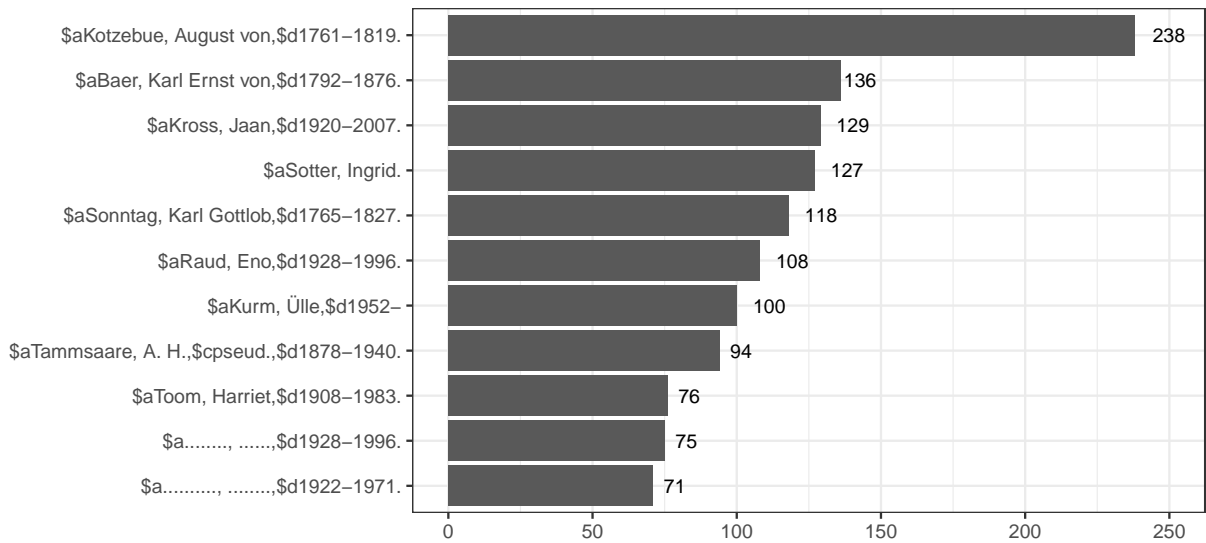
Most common authors in the bibliography



Most common authors in the Estonian language set



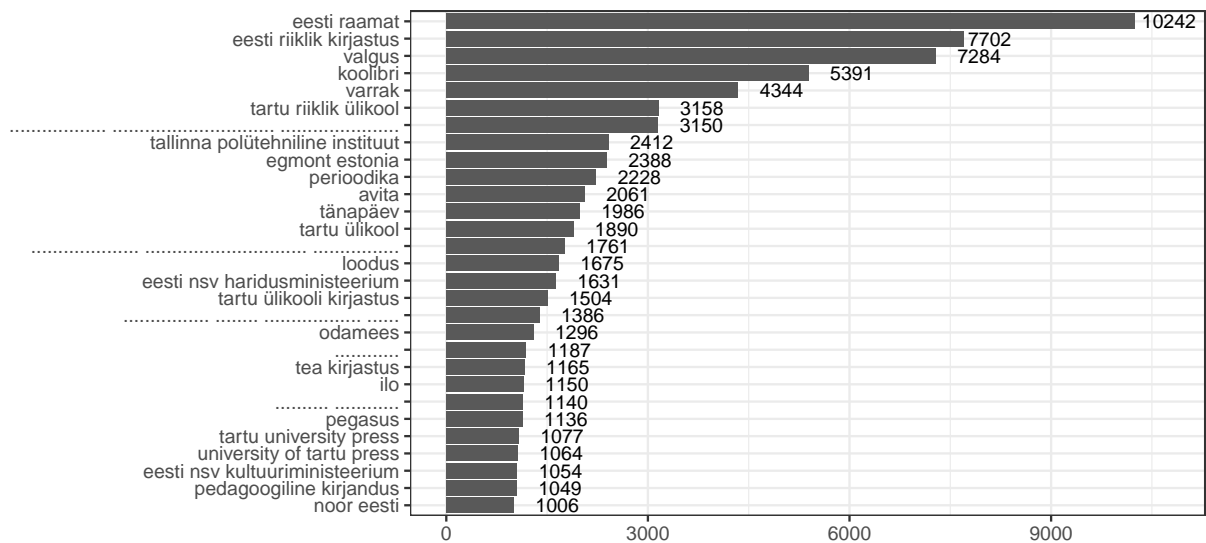
Most common authors among the publications not in Estonian

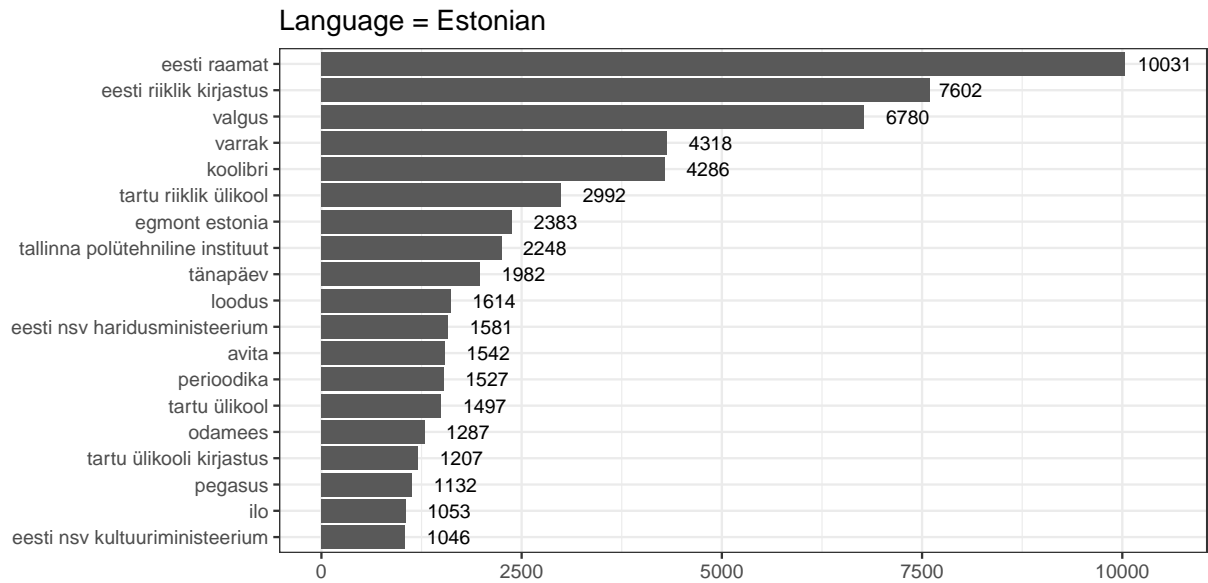


Publishers

Publisher names have also been harmonized manually and algorithmically with more frequent names in focus.

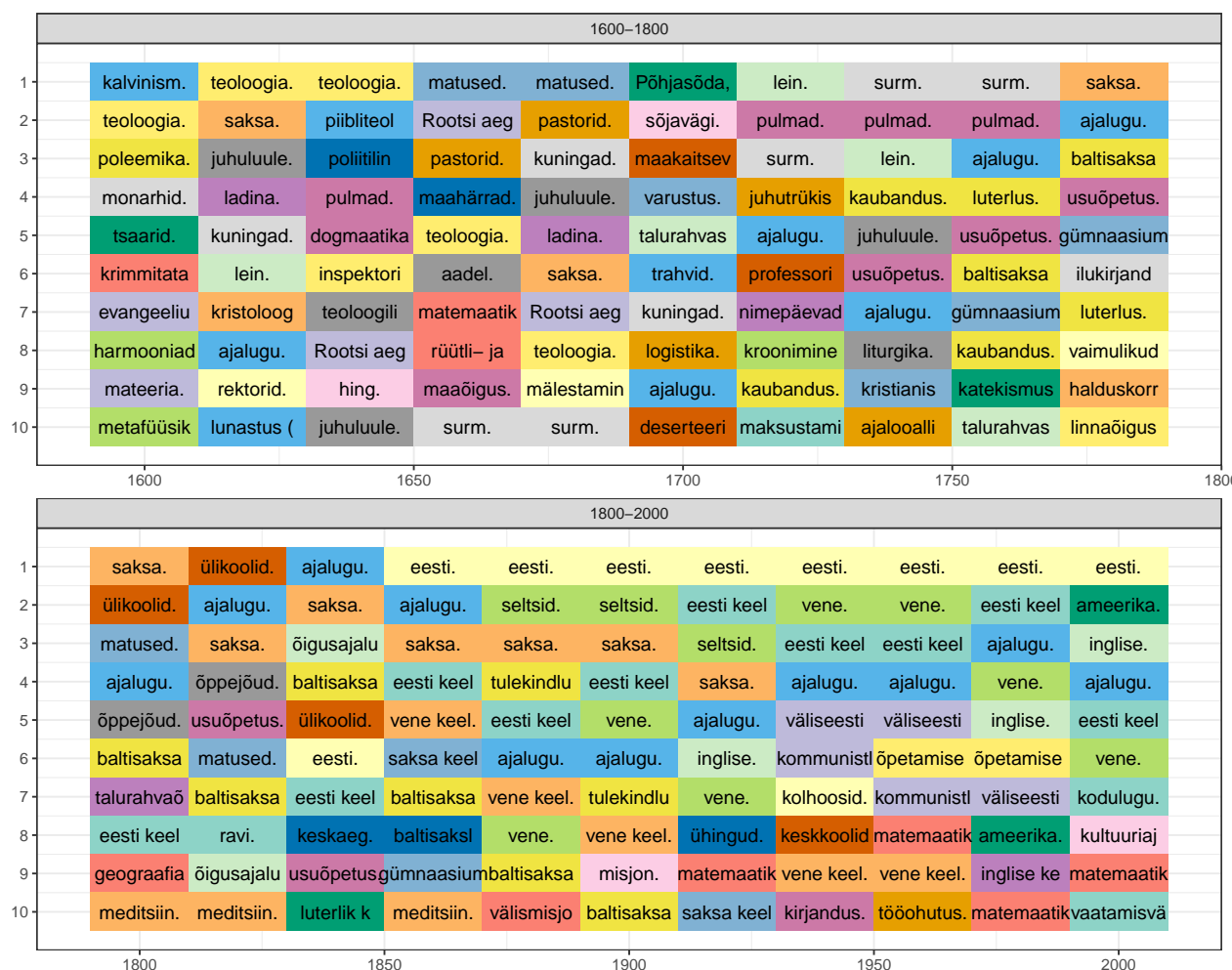
Language = All





Topics

The bibliography has each publication marked by some general topics. Some works have no topic marking, many works have several topic markings. By 20-year intervals, here are the most common topics in the set. For visualization purposes, only the first 10 characters of the topic are shown.



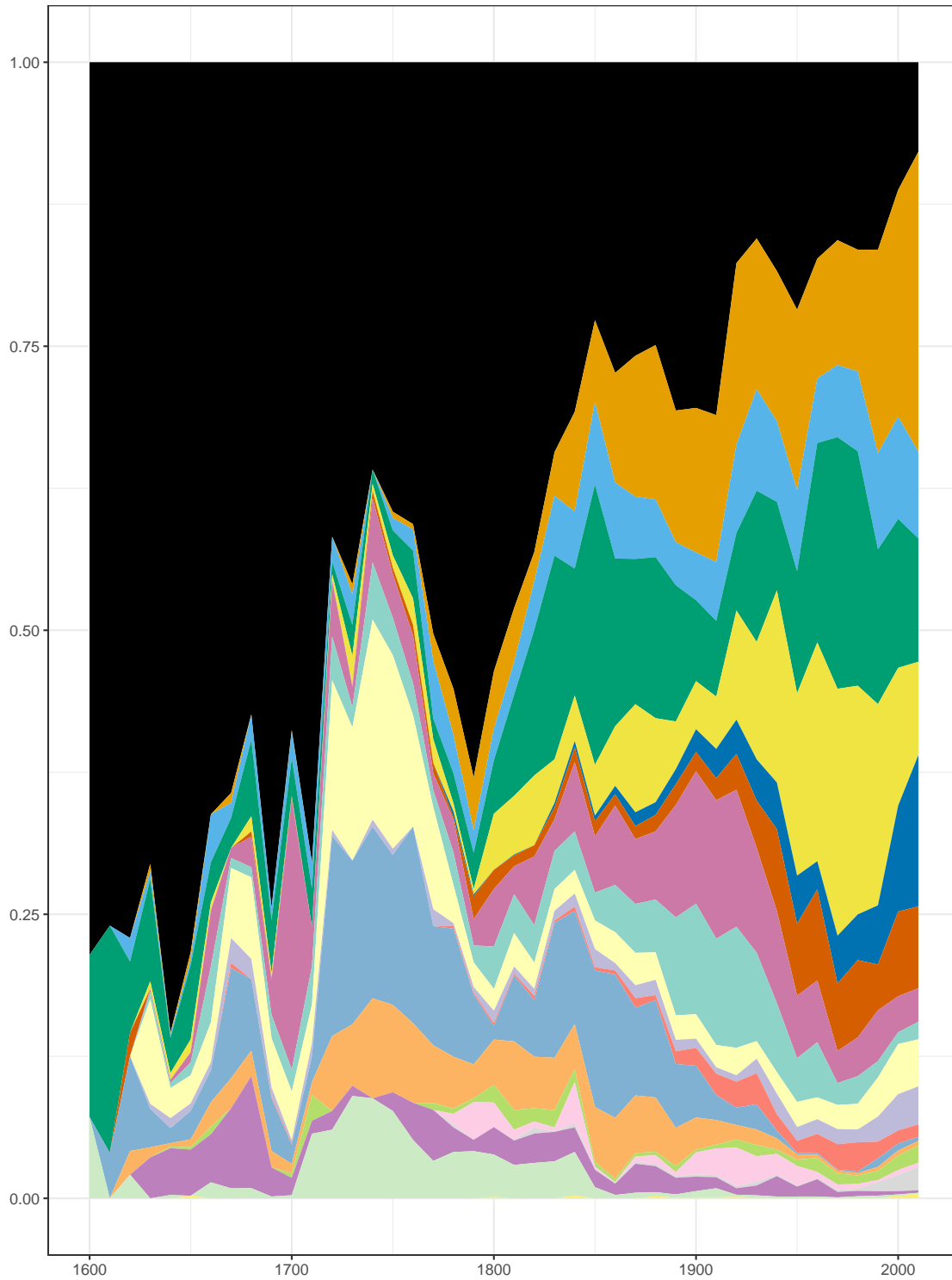
Genres

The bibliography also includes genre markers on many works. Some works have no markers, many have several. Here, we have built a few larger categories based on the dataset and visualized their frequencies over time. NONE means that the works did not have any genre markers.

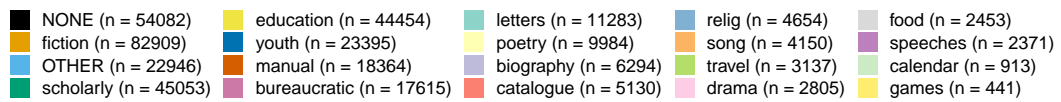
Genres have been grouped into 18 major categories. Many remain uncategorized. A total of 559 unique keywords are currently classified as OTHER. Top genres there are given in the table below. 54647 works in the dataset had no genre marking.

genre_orig	N
mälestused.	2621
e-raamatud.	1258
võrguväljaanded.	948
juubeliväljaanded.	936
fotoalbumid.	854
kavad.	826
statistilised andmed.	813
infotrükised.	749
sõnaraamatud.	643
oskussõnastikud.	465

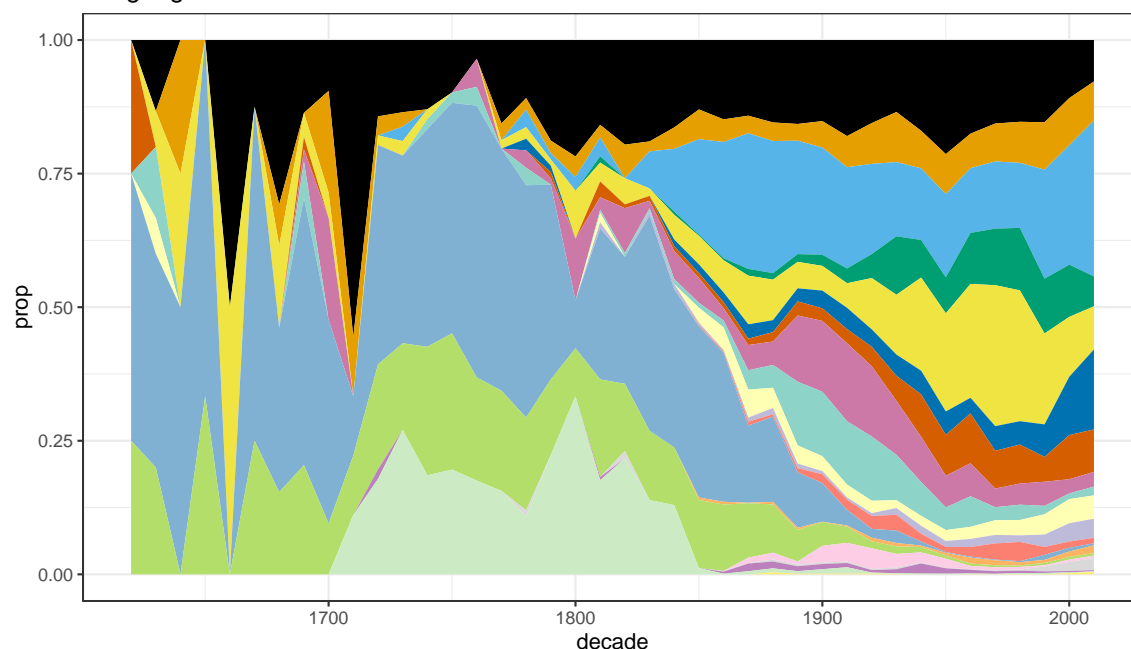
Language = All



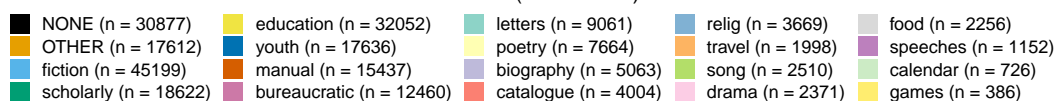
All (n = 270462)



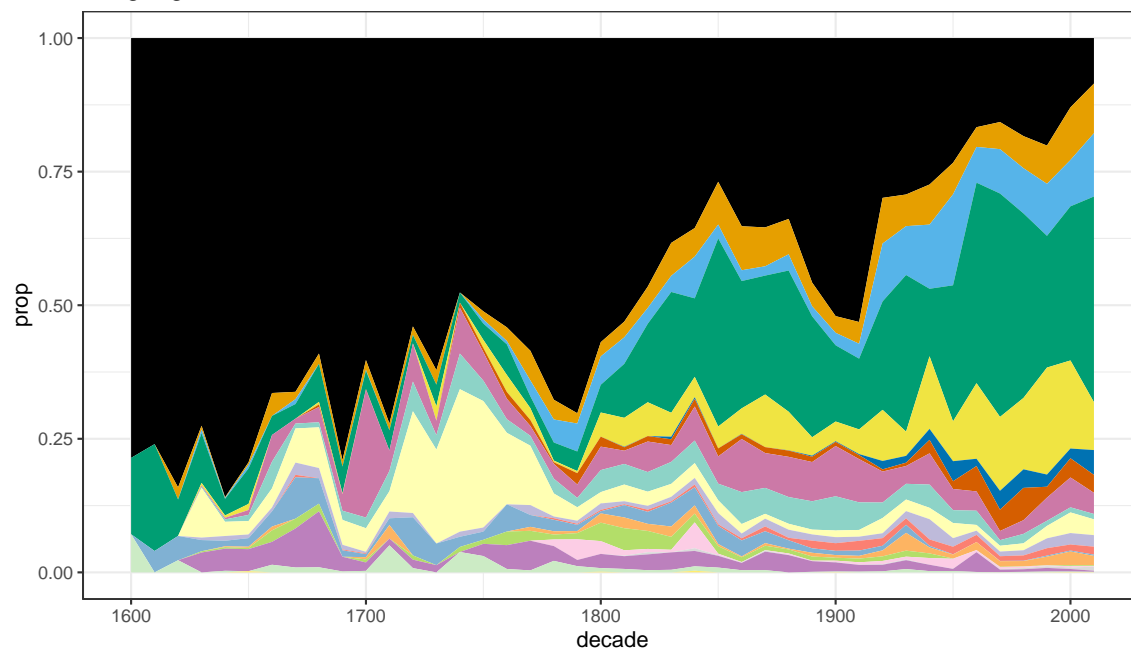
Language = Estonian



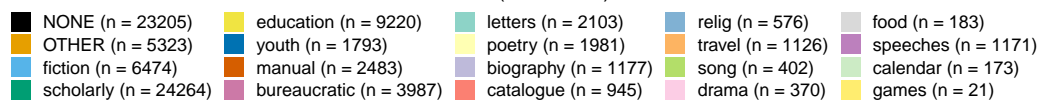
Estonian (n = 191569)



Language = Non-estonian



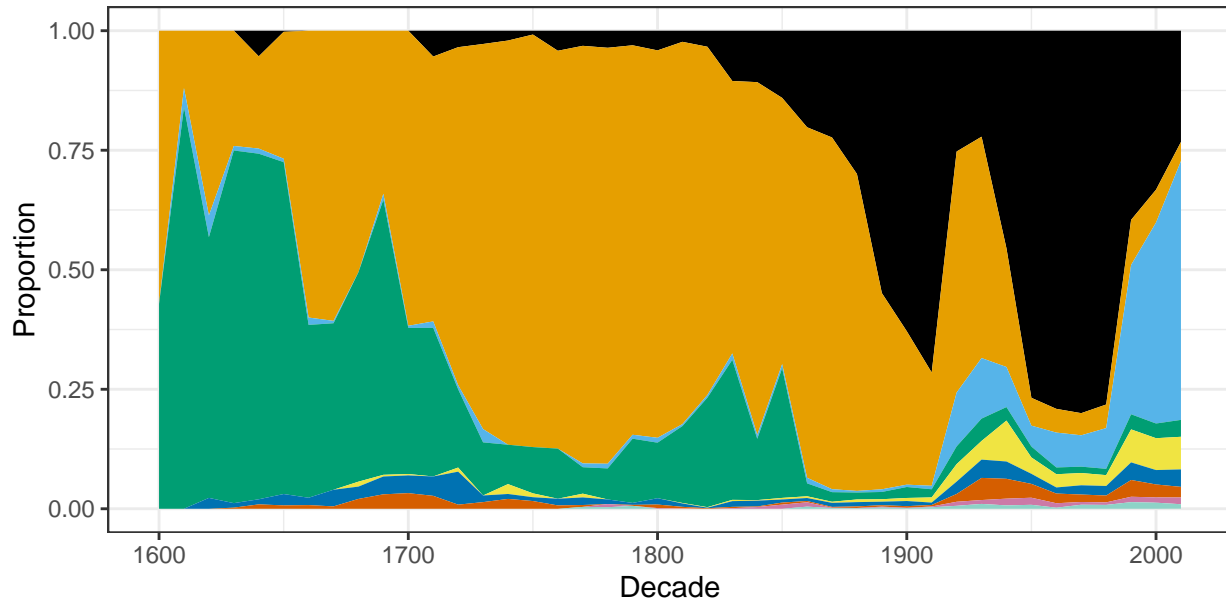
Other (n = 78918)



Languages

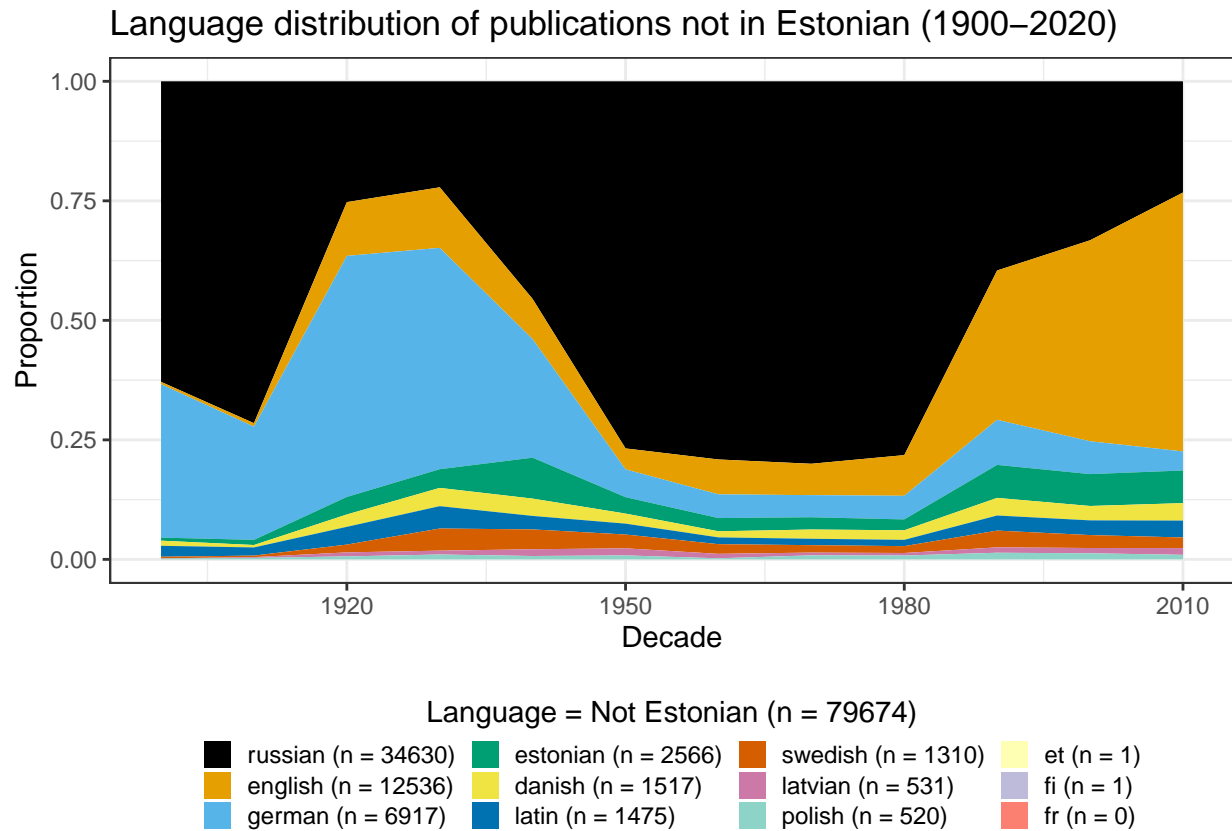
For the dataset that is not in Estonian, language has been marked in combination of existing metainformation and automatic language extraction based on the title of the work. An overview of the languages is given below. Each work is just given one language here and multilinguality is not included in this measure.

Language distribution of publications not in Estonian (1600–2020)



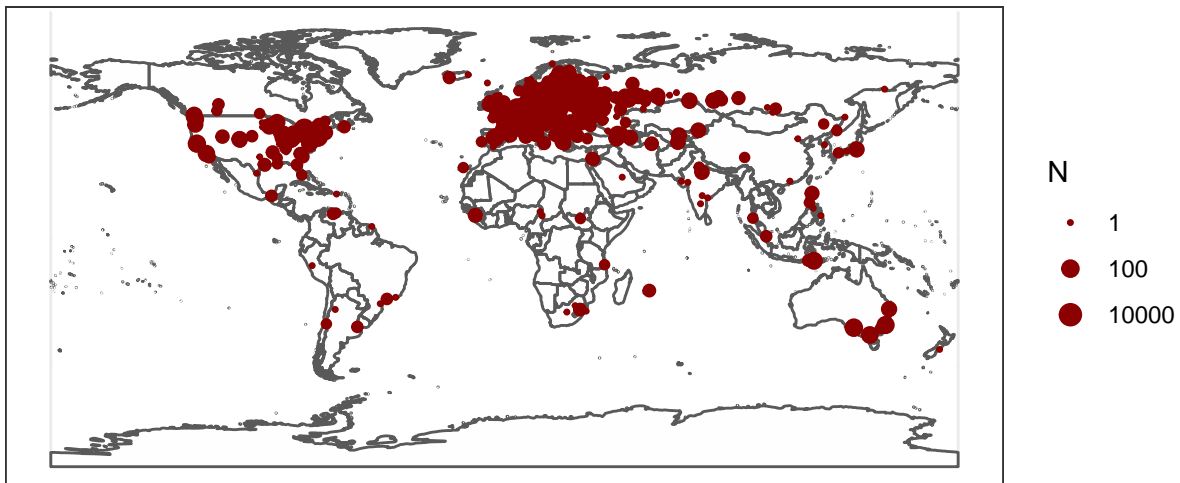
Language = Not Estonian (n = 79674)

■ russian (n = 37631)	■ latin (n = 4663)	■ swedish (n = 1411)	■ et (n = 1)
■ german (n = 17178)	■ estonian (n = 2617)	■ latvian (n = 563)	■ fi (n = 1)
■ english (n = 12662)	■ danish (n = 1720)	■ polish (n = 541)	■ fr (n = 1)

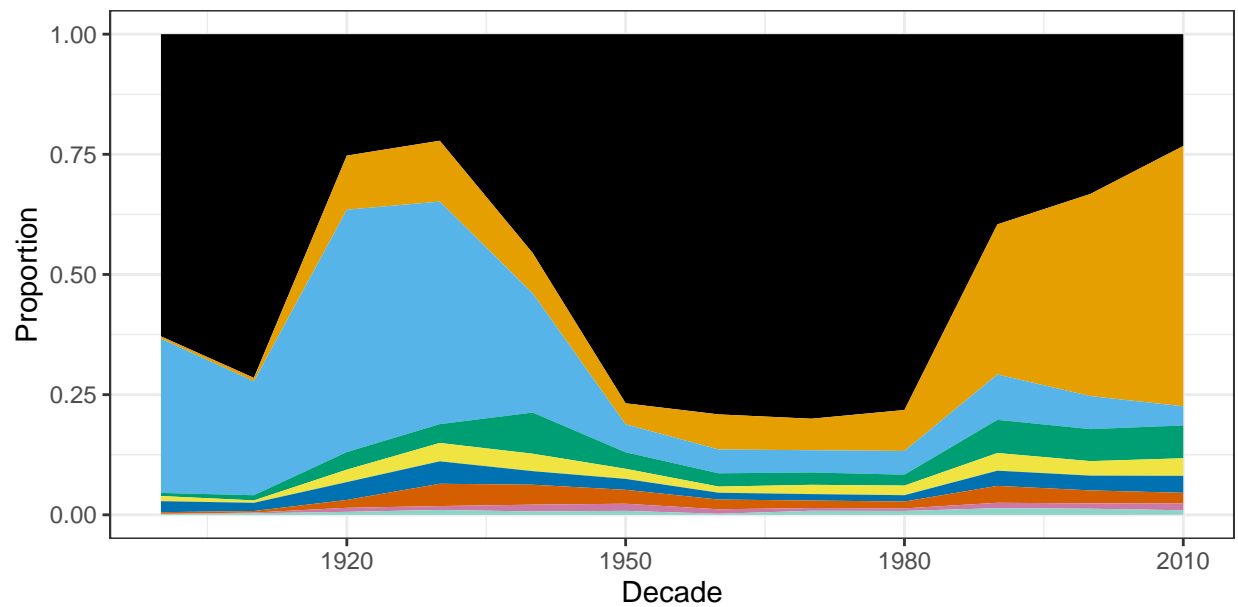


Publication locations

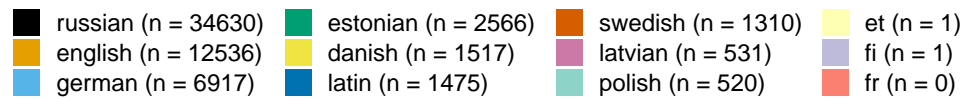
Publication locations have also been joined with geographic information in the GeoNames database. This allows the publication locations also to be displayed on a map. To see the interactive version of the map, use the [html overview](#).



Language distribution of publications not in Estonian (1900–2020)



Language = Not Estonian (n = 79678)





V6

~decade2: 1500

Play

