

# Metadata from ERB

## ERB 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. In this set, only the publications in Estonian language are used. The dataset has been compiled in digital format since 2002 and aggregates the work of multiple institutions and generations in collecting the publication information.

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: [http://www.elnet.ee/images/pdf/juhendid/vanaraamat\\_MARC21.pdf](http://www.elnet.ee/images/pdf/juhendid/vanaraamat_MARC21.pdf)

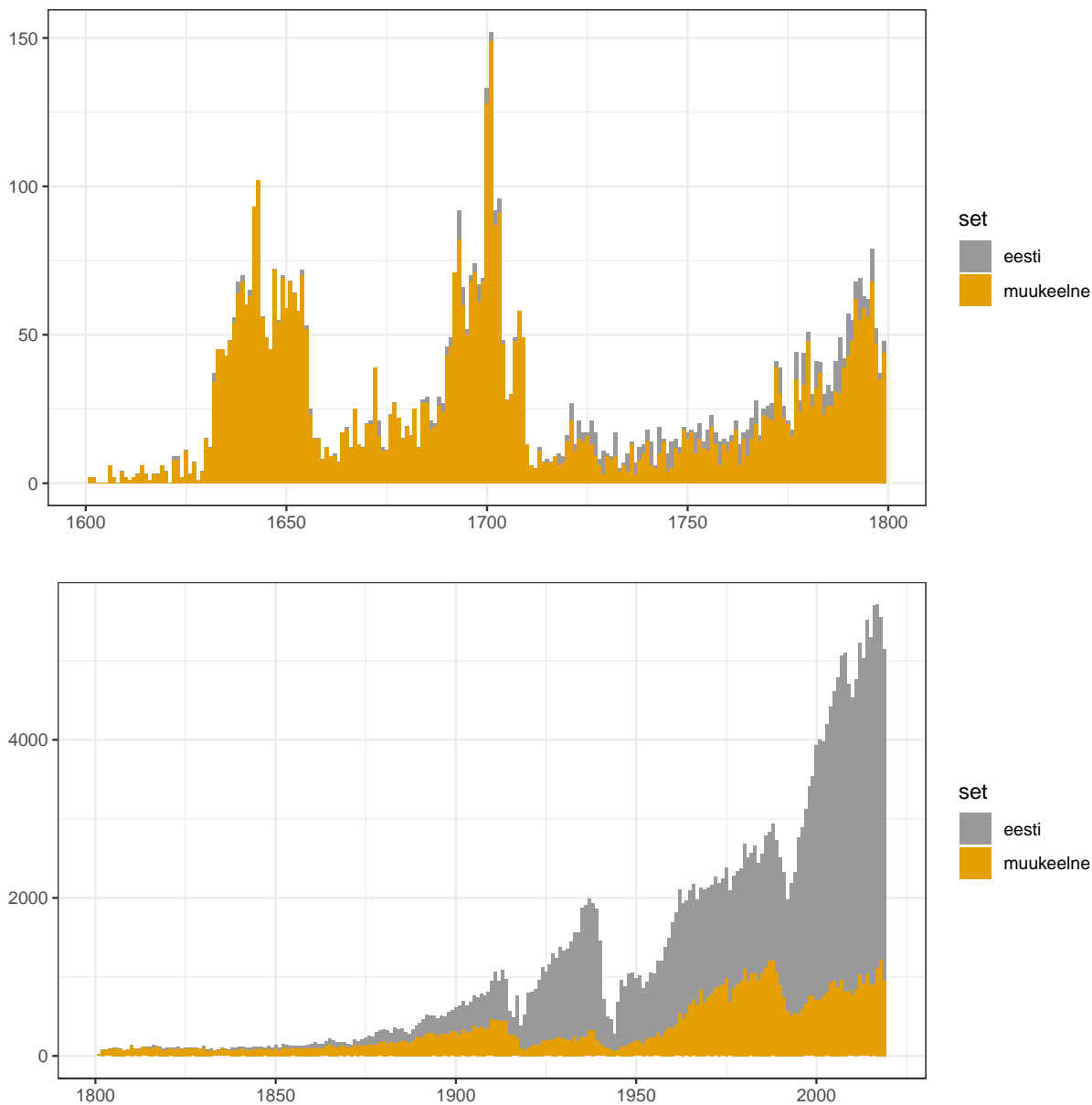
## Intro

The current code organizes the ERB 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

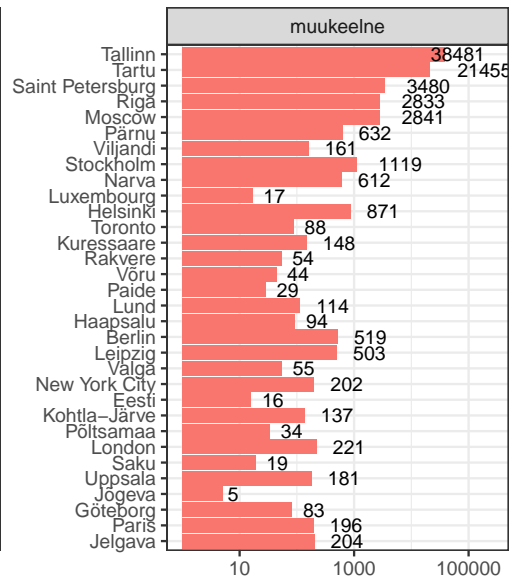
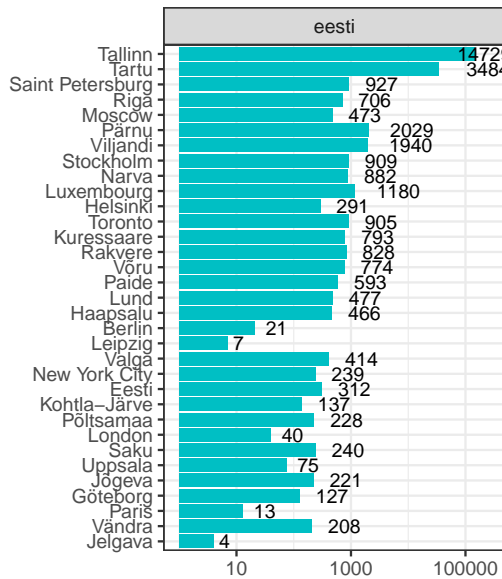
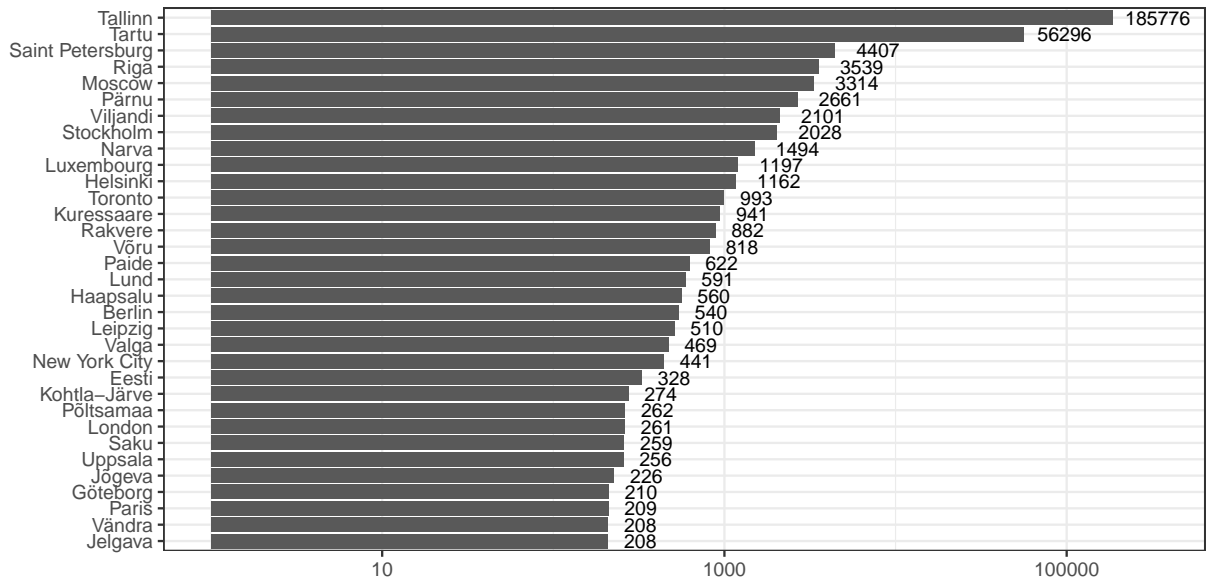
The dataset has altogether information on 291346 printed items. The coverage has been estimated to be better than 95% for all of the relevant published works.

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.



## 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.



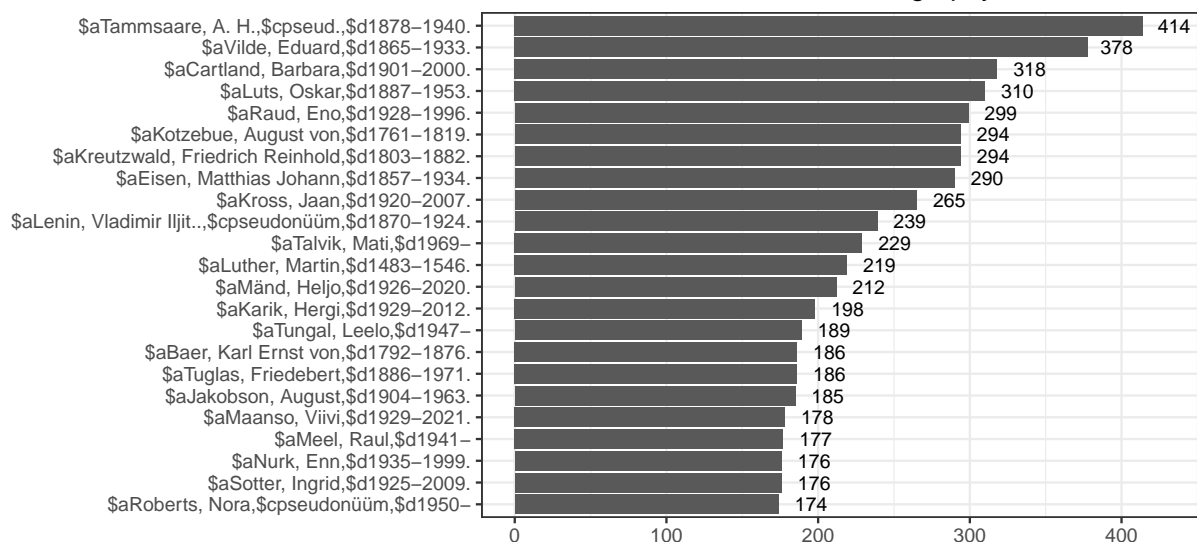
fct\_rev(set)

muukeelne

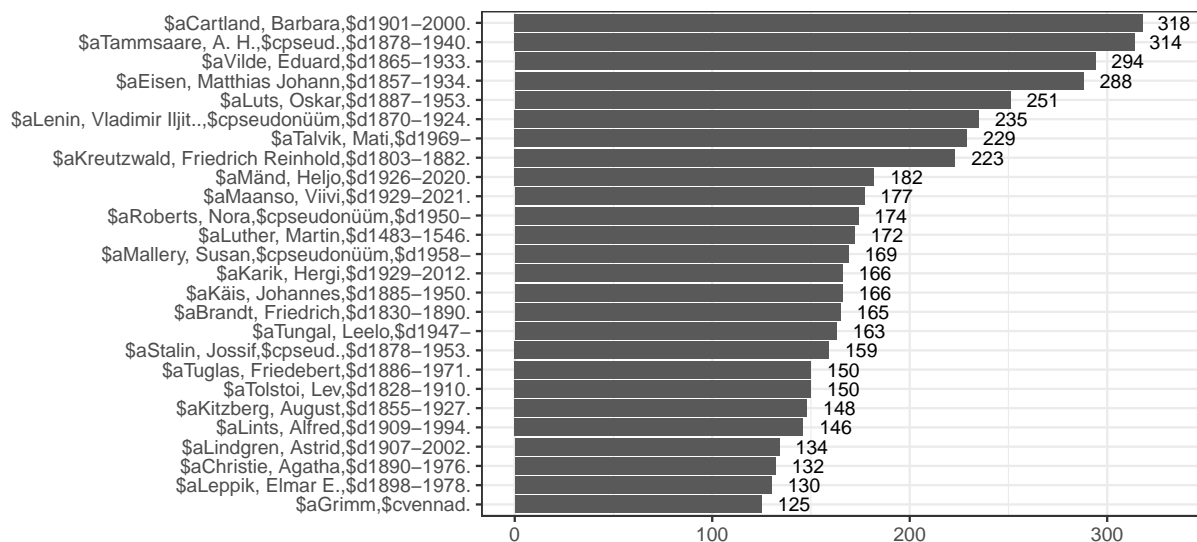
eesti

## Authors

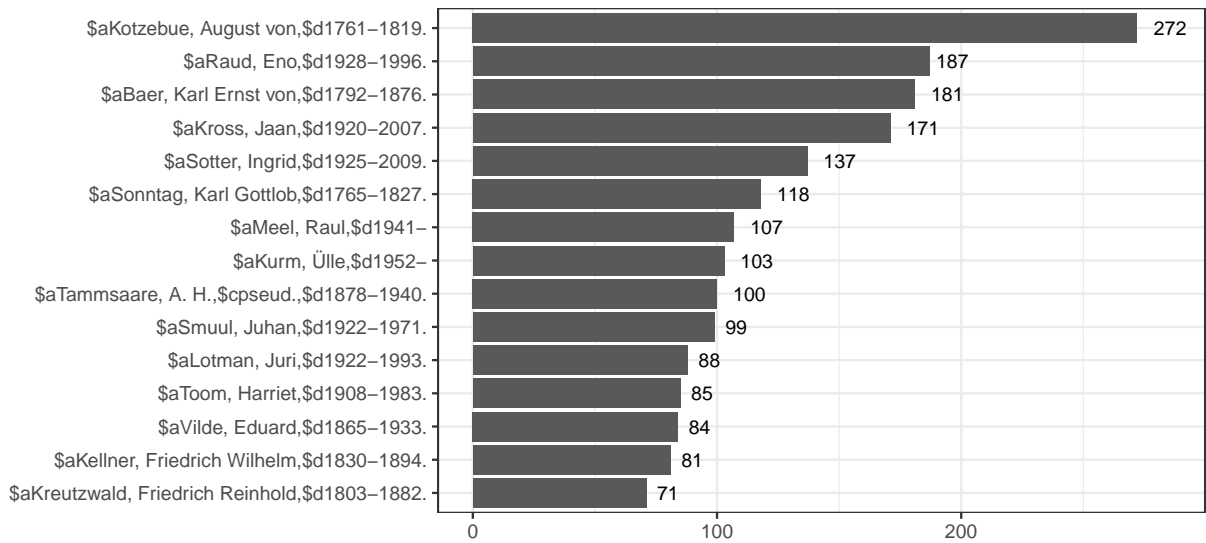
Most common authors in the bibliography



Most common authors in the Estonian language set



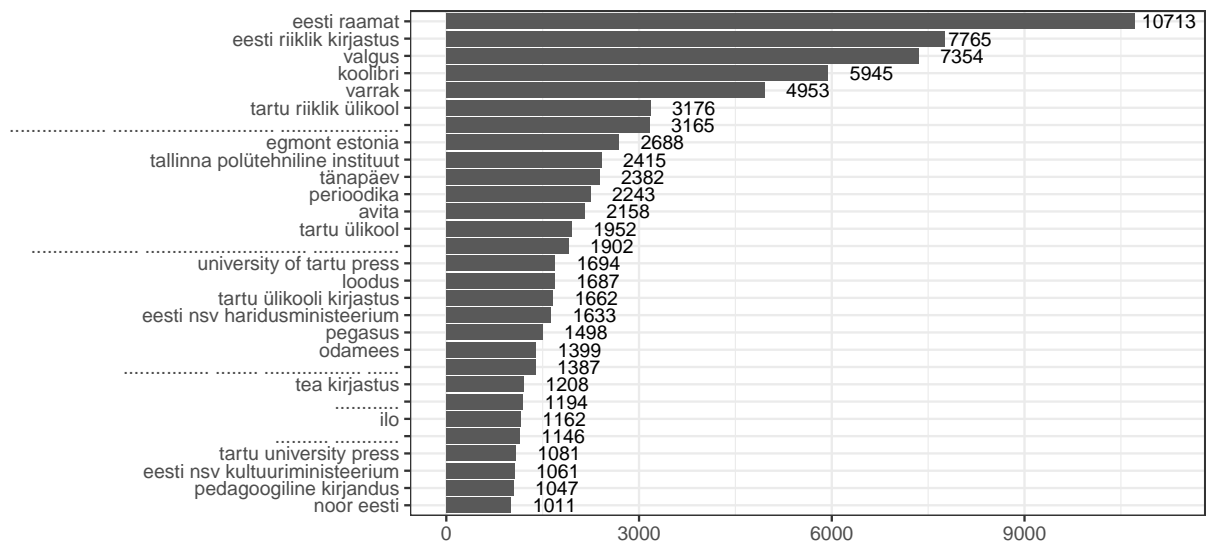
### Most common authors among the publications not in Estonian

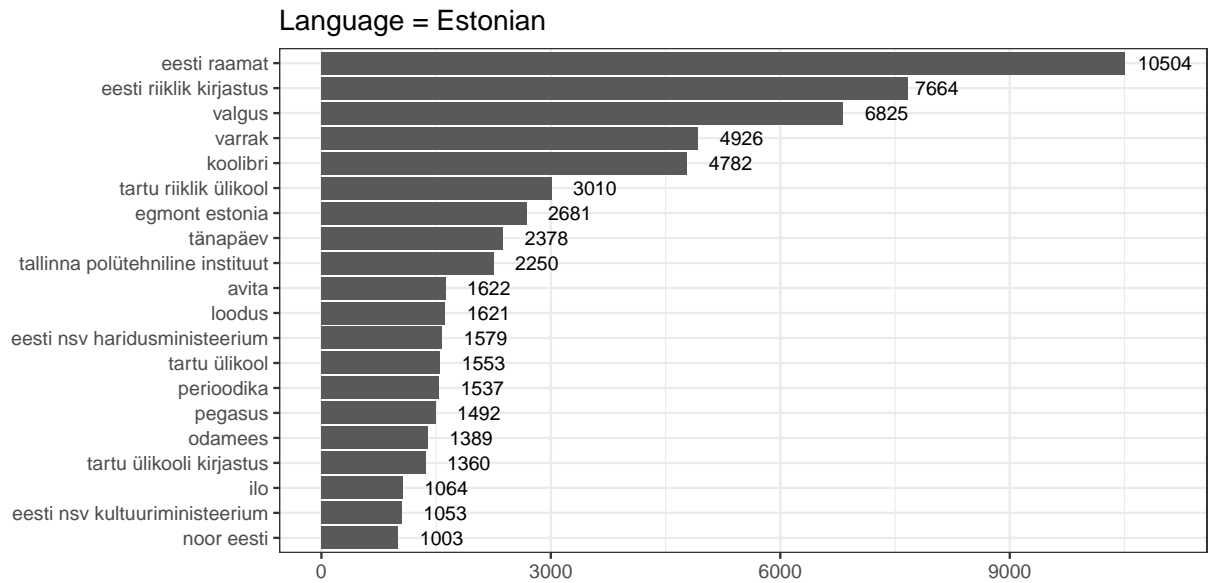


### Publishers

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

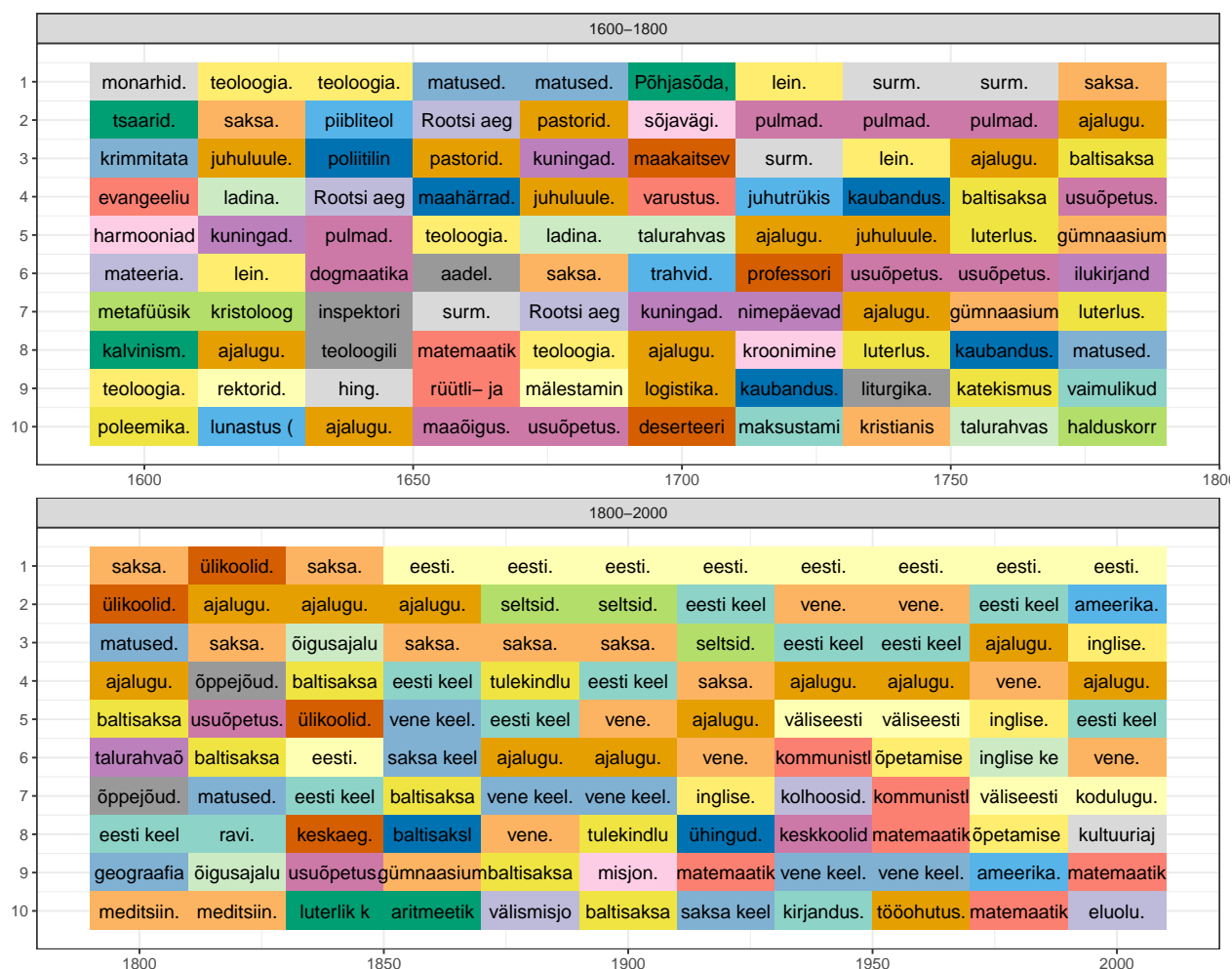
### Language = All





## 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.



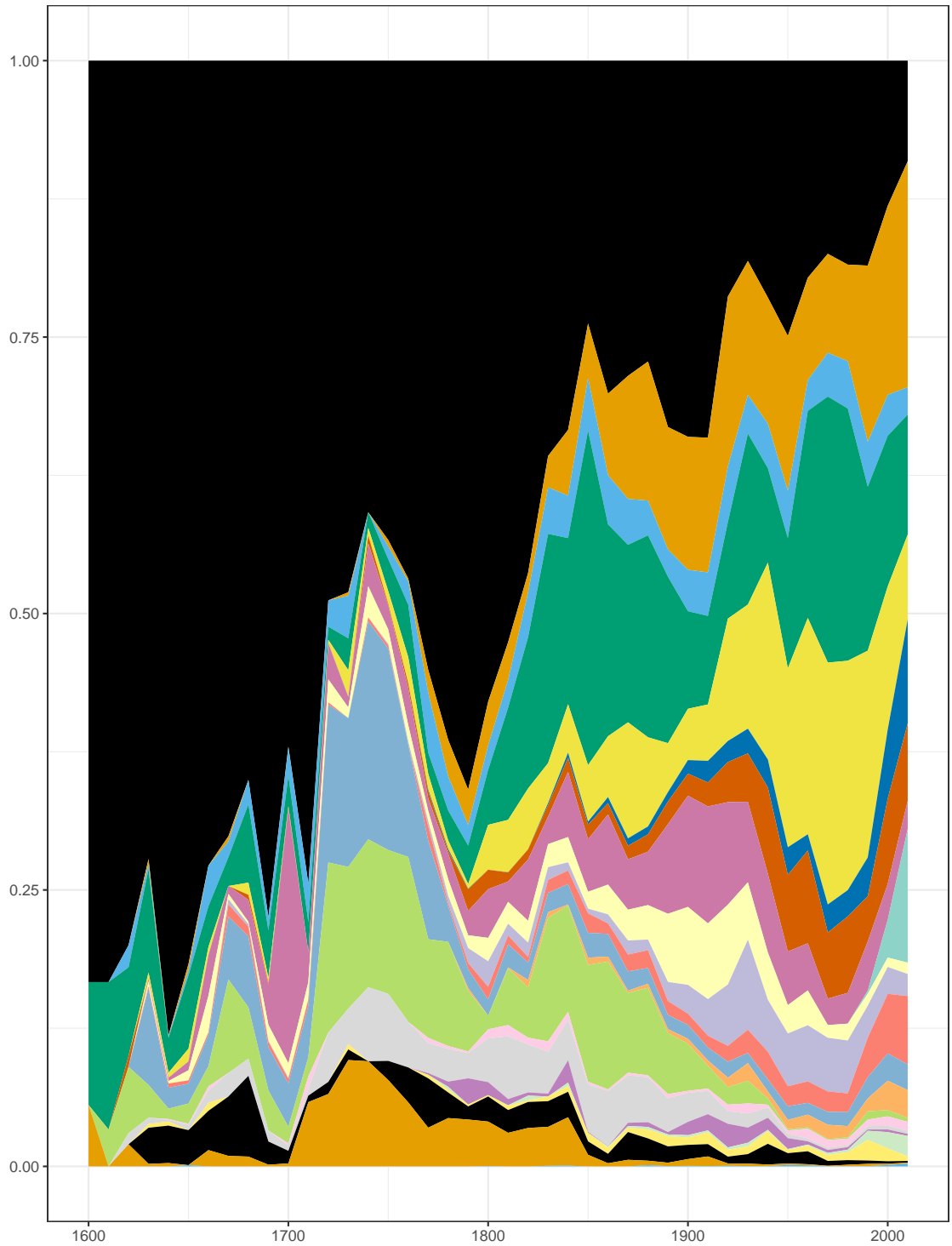
## Genres

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.

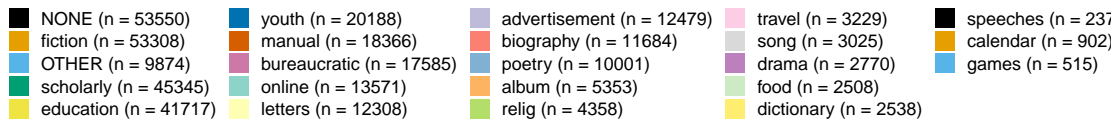
Genres have been grouped into 21 major categories. Many remain uncategorized. A total of 474 unique keywords are currently classified as OTHER. Top genres there are given in the table below. 1 works in the dataset had no keyword marking.

genres	N
annotatsioonid.	342
preprindid.	311
esseed.	306
harjutused.	282
õpetajaraamatud.	259
tabelid.	214
ehitusnormid.	197
anekdoodid.	188
kunstnikuraamatud.	185
aforismid.	184

Language = All

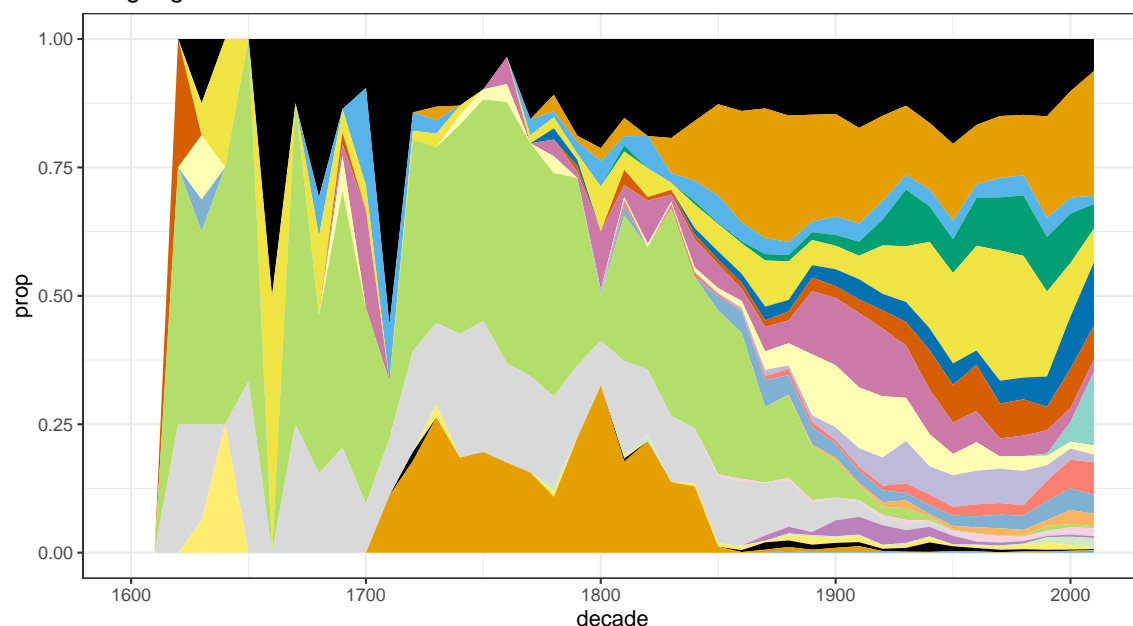


All (n = 279438)

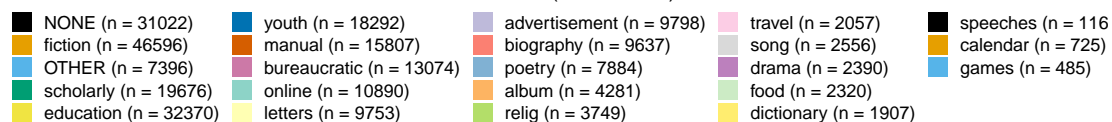




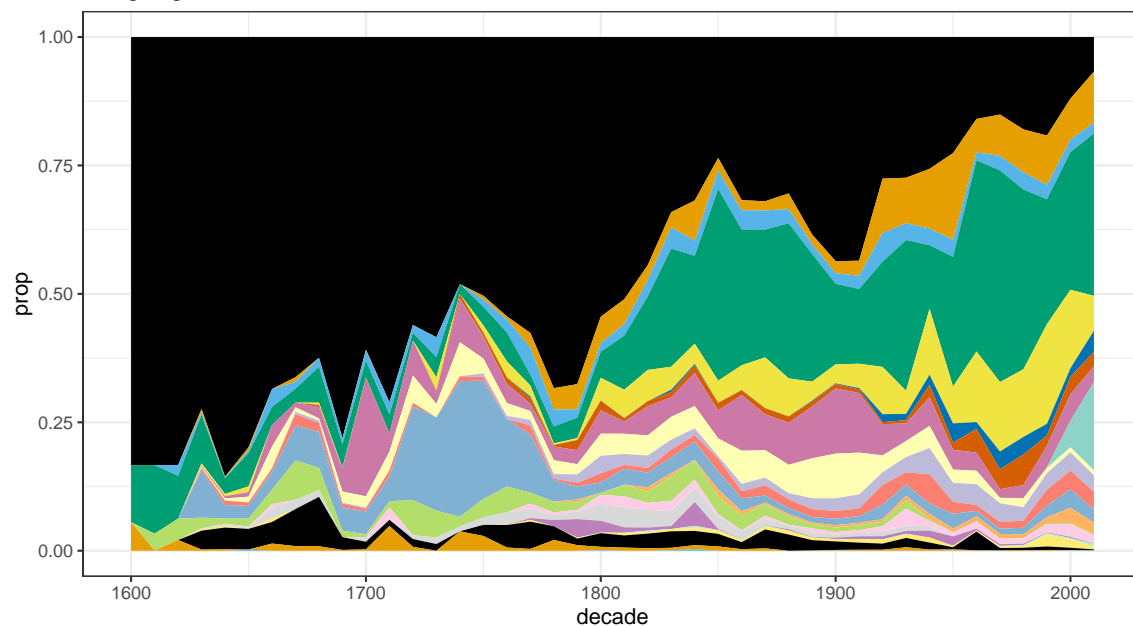
Language = Estonian



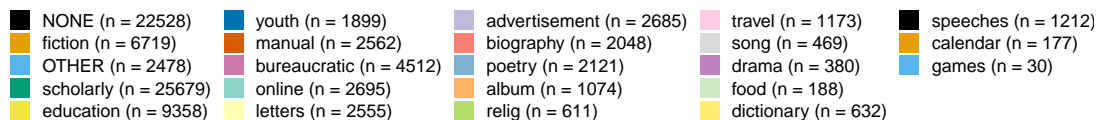
Estonian (n = 198309)



Language = Non-estonian



Other (n = 81174)



## 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](#).

