# Mining Texts at the National Library of Estonia

Potential for Learning and Discovery

Peeter Tinits

12.03.2020

Tallinn, Estonia

# Digitizing texts

Large digitization programmes:

*Estonian written word to the world* (2014-2018)
    Accomplished: 9% of books, 26% of periodicals

*Action Plan for Digitization of Cultural heritage* (2018-2023)
    Goal: 28% of printed publications (+700,000 newspaper pages, + 800,000 book pages)

*Legal Deposit Copy Act* – since 2017 digital copies of estonian Publications

Current state: 13% of books, 36% of periodicals in the registry

Digital print archives (DIGAR):
    61,042 prints
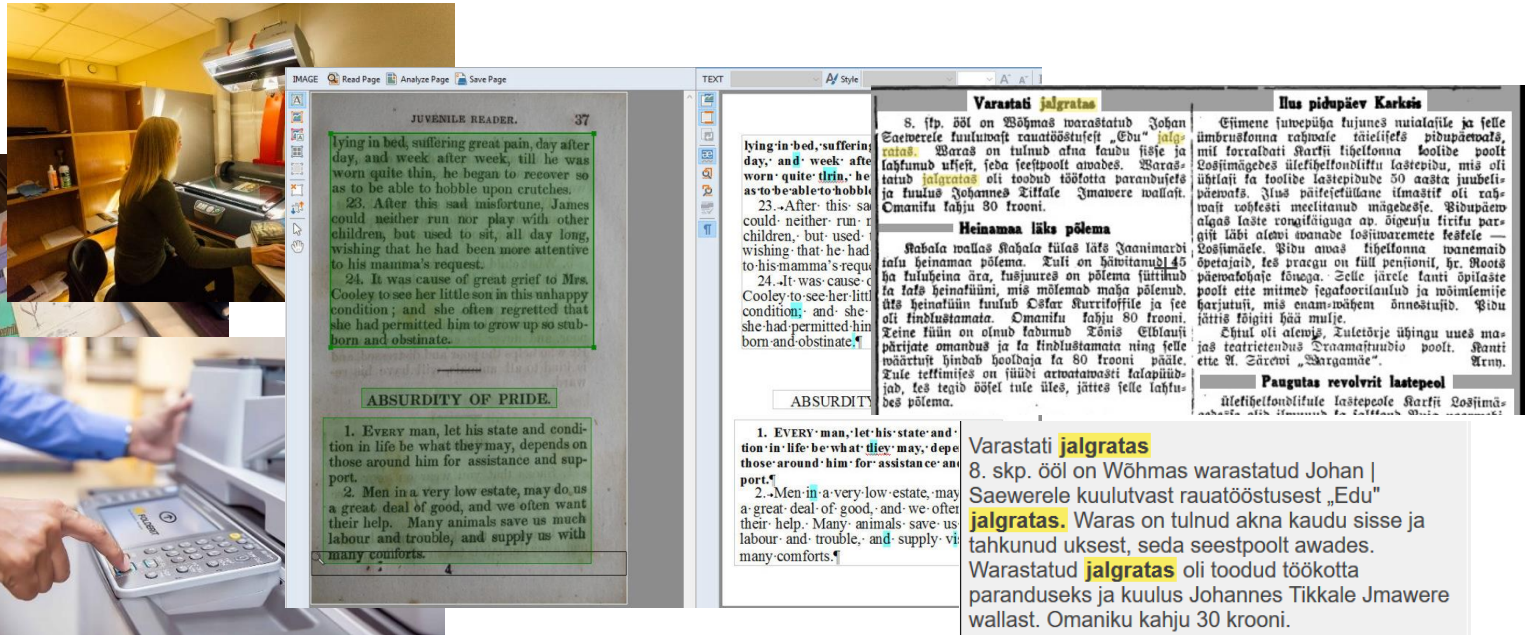
Digitized periodicals (DEA.DIGAR):
    390,186 issues,
    3,725,232 pages
    6,441,167 articles
    2285 titles

# Digitization



DOCUMENT SCAN → SCANNED IMAGE FILE → OCR (Optical Character Recognition) → TEXT DOCUMENT

# DIGAR: search possibilities

# Text and data mining approach

Developed by a # of institutions and researchers

# Text and data mining in NLE

Improving access for text and data mining

# Overviews

Number of issues that can currently be accessed publically

EESTI
RAHVUS-
RAAMATUKOGU

# Overviews

Number of issues that can currently be accessed publically

# Bibliography as data

Data on 351,301 printed publications

Aims for complete coverage (quite close)
Published in Estonia, or abroad but related to Estonia

DIGAR    DEA    ERB    ISIKUD/KOLLEKTIIVID    Andmekaevurile    Näited    Kontakt

Eesti Rahvusbibliograafia

```
<mx:datafield ind1="1" ind2=" " tag="700">
    <mx:subfield code="a">Stendhal,</mx:subfield>
    <mx:subfield code="d">1783-1842</mx:subfield>
    <mx:subfield code="0">(BLBNB)0005643
</mx:datafield>
<mx:datafield ind1="0" ind2=" " tag="70
    <mx:subfield code="a">ذبربمؤ«نذ ٱٻۃ۲۵ۻۮ
    <mx:subfield code="d">1783-1842</mx:
    <mx:subfield code="0">(ISNI)00000001
</mx:datafield>
<mx:datafield ind1=" " ind2="1" tag="70
    <mx:subfield code="a">ذبربمؤ«نذ ٱٻۃ۲۵ۻۮ
    <mx:subfield code="d">1783-1842</mx:
    <mx:subfield code="0">(NLR)RU NLR AU
</mx:datafield>
```
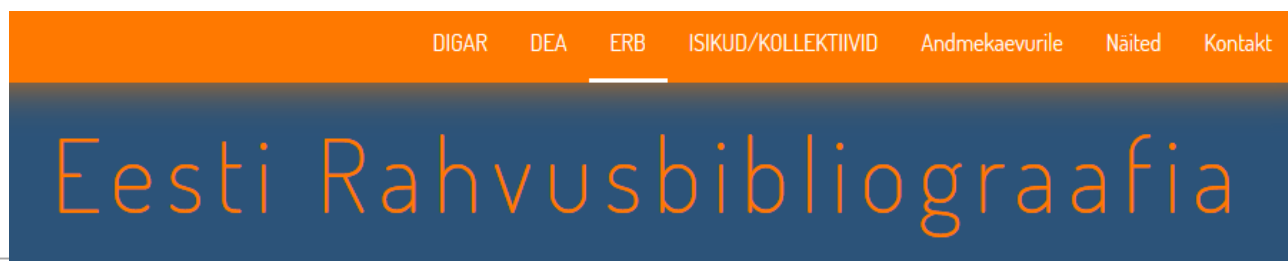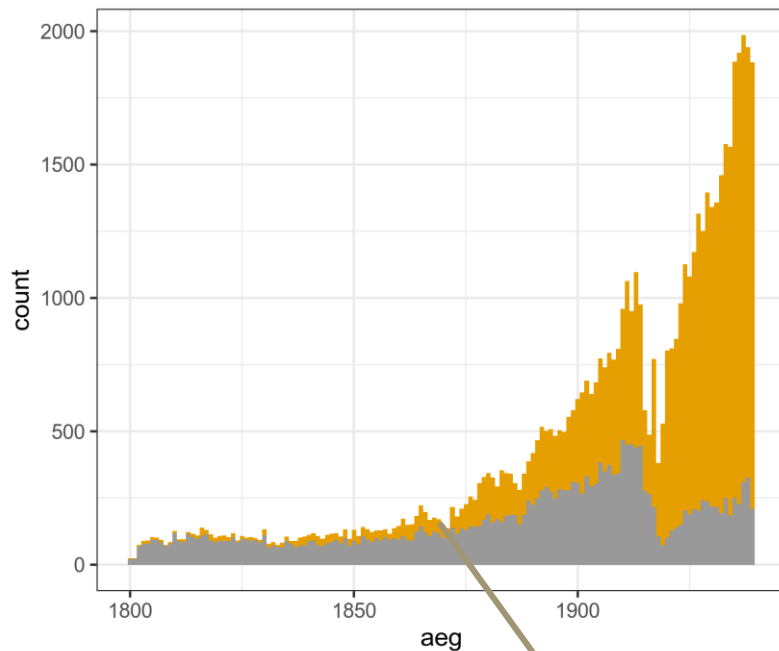
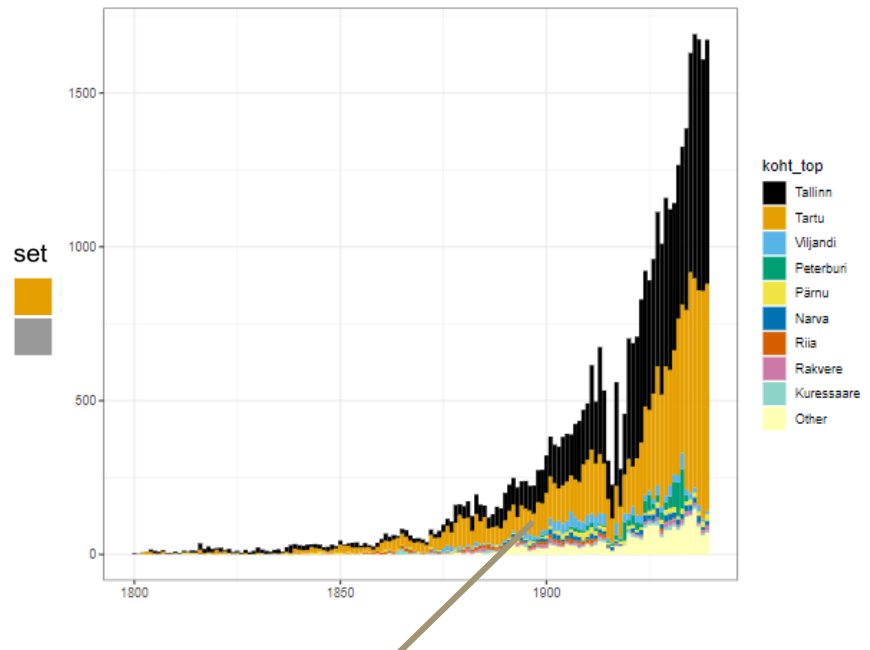| RecordNumber | c | aeg | koht | kirjastus | title |
|---|---|---|---|---|---|
| 115097 | $c[1894] | 1894 | Narva | R. PĆμder, | Tumm armastaja : |
| 115134 | $c[1894] | 1894 | Paide | A. Seidelberg, | Saardami tĆꞮĆꞮmees : |
| 115406 | $c1894 | 1894 | Viljandi | Viljandi Eesti Karskuse Selts """"Vabadus"""", | Kolm esimest elu-aastat : |
| 115679 | $c1894 | 1894 | Tallinn | A. Busch, | WiieteistkĆ¼mne aastane kapten / |
| 115739 | $c1894 | 1894 | Tallinn | A. Brandt, | Willem ja Luise ehk Ćμnnetu ja Ćμnnelik perekon |
| 115891 | $c1894] | 1894 | Tallinn | s.n., | Eeskirjad ĆꞮlemtalurahwa ja walla kohtutele Talli |
| 116002 | $c1894 | 1894 | Jurjev | H. Laakmann, | Weike meelejahutaja : |
| 116089 | $c[tsens. 1894] | 1894 | Riia | Paltimaa Ć•igeusu vennaste-selts, | Ć•petus Ćμigest usust : |

EESTI RAHVUS-RAAMATUKOGU

# Uses

Some examples in collections as data

# For learning

Number of publications by language, by place



„National Awakening"

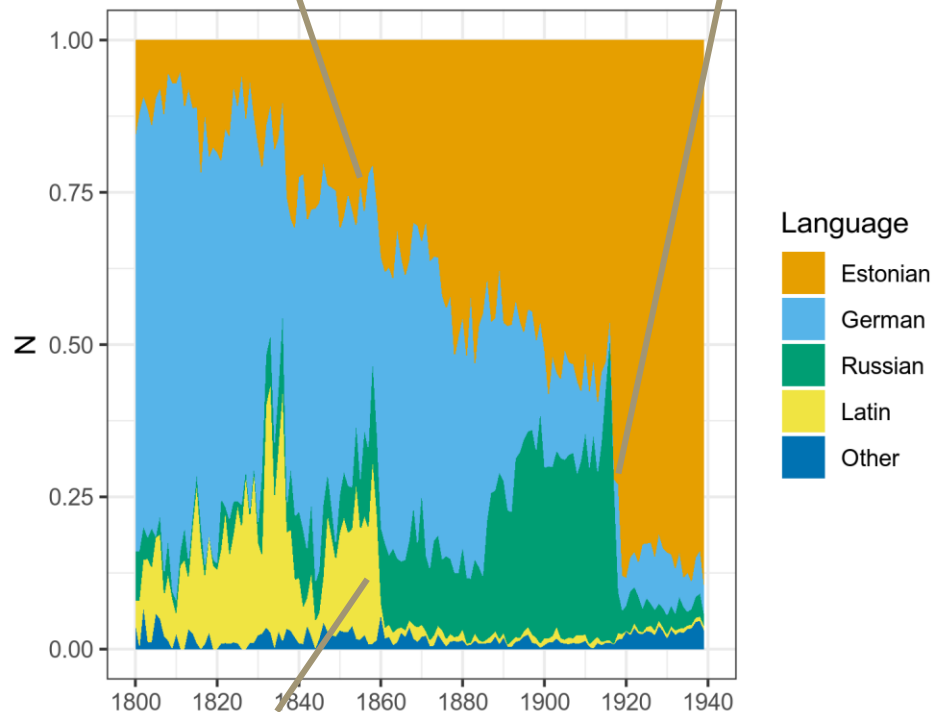Tartu leads publishing

# For learning

Languages of prints 1800-1940

The use of Estonian grows gradually and surely

Independent Estonia



Baltic German Publications throughout

Loss of Latin as the language of science

EESTI
RAHVUS-
RAAMATUKOGU

# For learning

Genres of prints 1700-2000

Religious texts          Fiction

# For discovery

Top keywords 1600-2000 (20-year chunks)



Religious texts

Estonia(n) as keyword

Civic societies

Disappearance of German keyword

# For discovery

# For discovery

Texts featuring „marlene" and „dietrich"

# For discovery

Texts featuring „marlene" and „dietrich"

# For discovery

Texts by topic on Marlene Dietrich

variable
- Cinema
- Concerts
- Fashion
- General
- History

EESTI
RAHVUS-
RAAMATUKOGU

# For discovery

## Words in texts over time

# For discovery

Words by frequency in texts on Marlene Dietrich

# For discovery

Words that were common only for a short time

# For discovery



Co-occurrence network of top words

Thank you!

# All links
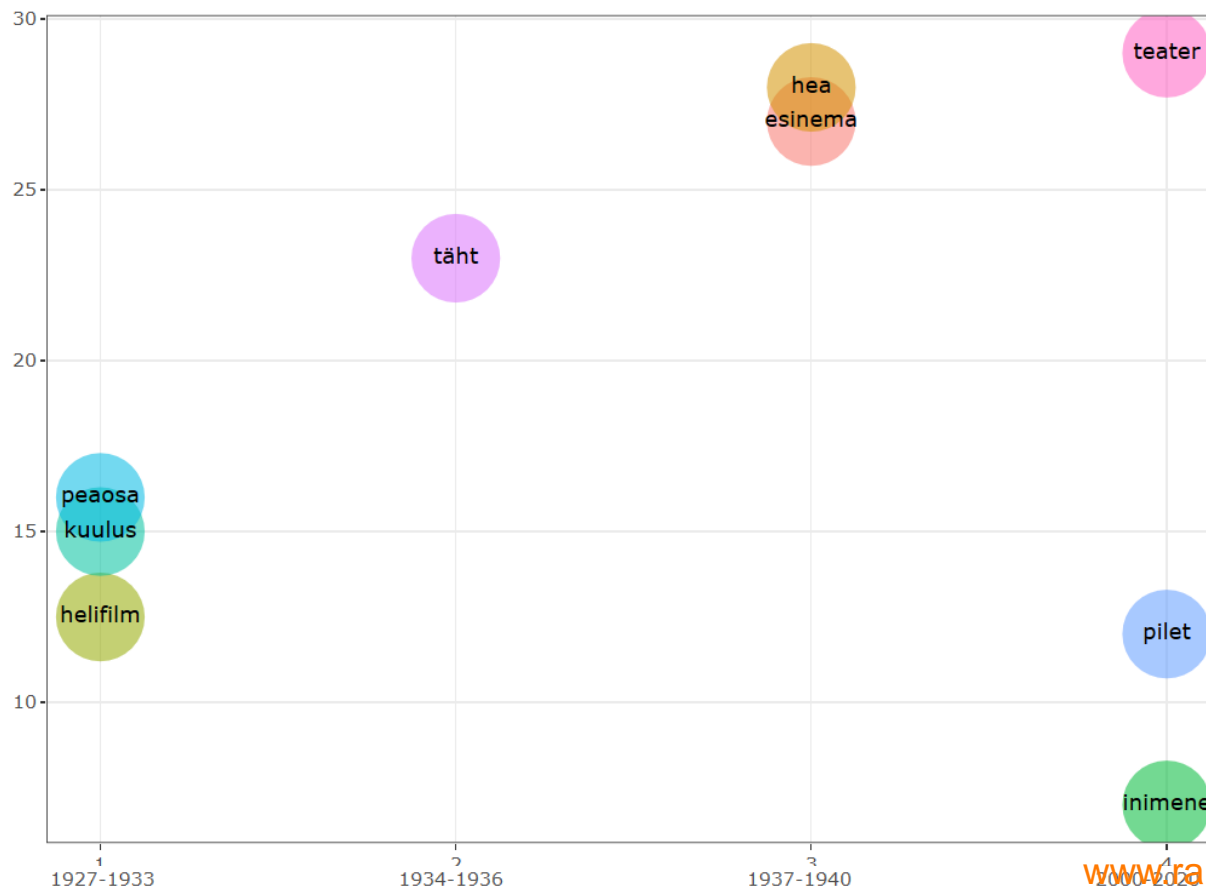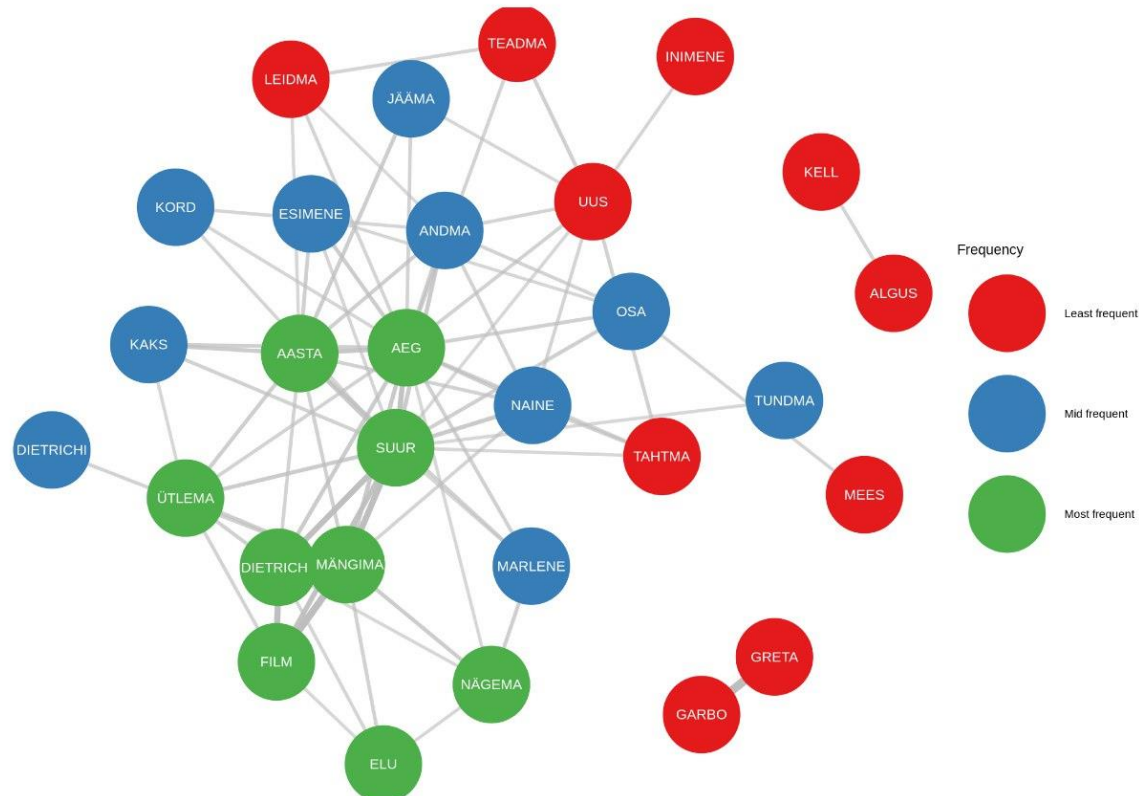
https://peetertinits.github.io/slides/plots/RR2020/current_state_color.html

https://peetertinits.github.io/slides/plots/RR2020/future_state_color.html

https://peetertinits.github.io/slides/plots/RR2020/dietrich_articlecount.html

https://peetertinits.github.io/slides/plots/RR2020/dietrich_articlecount_wtypes.html

https://peetertinits.github.io/slides/plots/RR2020/dietrich1.html

https://peetertinits.github.io/slides/plots/RR2020/dietrich_all_lines.html

https://peetertinits.github.io/slides/plots/RR2020/dietrich_continued.html

https://peetertinits.github.io/slides/plots/RR2020/dietrich_singles.html