# A DigiHum approach to HistSocLing of Estonian

Peeter Tinits

25.07.2019

National Library of Latvia
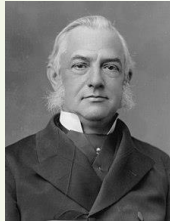
BSSD2019

# A case study in DH

- Project:
  - What mechanisms are responsible for spelling standardization at the end of 19th century Estonian?

- Data sources:
  - Anything relevant

- Timeline:
  - PhD dissertation

# Some context

- Language change: intentional or not



Müller (1861), Saussure (1916), Lass (1999): „Intentional language change is peripheral to language"

Jesperson (1925), Ferguson (1989), Thomason (2007) „Intentional language change is quite central to language phenomenon"

# Written grammars as part of language

- Language institutions and written norms have become normal for us.

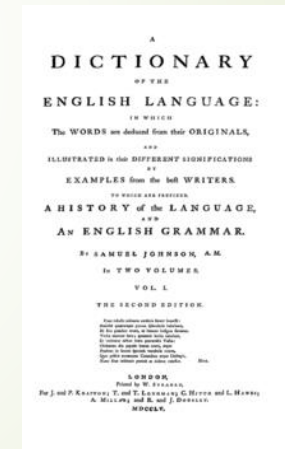Members of *Académie Française (2019)*

# Languages without these grammars

- Sociolinguistic research has shown that this is not always the case

- Of ~7,000 languages, ~3,000 still are not written



Living Languages of the World

# Even much writing is without standards

*Shakespere,
Shackspeare,
Shakespear,
Shakspere,
Shaxspere,
Shaxper,
Shakspeare,
Shackespeare,
Shackspere,
Shackespere*

- E.g. English became standardised from ~18c onwards

Samuel Johnson (1750): *Dictionary of English Language*

- *Shakespeare* standardized from 1860 onwards only

# Standardization research

- Spelling practices vary by community.

- Mechanisms and developments may depend on particular circumstances.

- How did this take place for (some of) Estonian?

# Estonian in the late 19c

- From 1800-1900:
  - Peasant population becomes independent
  - Quick developments in education
    - Literacy from 20% to 95%
  - Intense travel and movement
    - Urbanization, railroads
  - Social and cultural modernization
    - Cultural production and identity
    - Strive to become a nation

# Language standardization

- One form instead of many (colour & color -> color)

- How can this happen?
  - Talk like your neighbour
  - Follow a book
  - Survive longer

- For mechanisms, we should follow the process

# Finding the data

What is needed? What is available?

# Texts

- To study writing, we need texts:
  - Ideally from a variety of people
    - Education, dialect, profession

- Sources:
  - Linguistic corpora
  - Digitized works

# Linguistic corpora

- Linguistic corpora are balanced collections of texts to facilitate language research.

- For the period, there's two:
  - Written Estonian Corpus
    - some 300 book snippets in 1890-1930
  - Old Written Estonian Corpus
    - all texts 1500-1700, some 50 texts from later

- Small-ish and not much metadata



**Eesti Kirjakeele Korpus: 1890ndad**

**Statistika ja bibliograafia**

- Korpuse koostisosad valdkondade kaupa.
- Algallikad paberkandjatel failide kaupa.

**Tekstid**

**Märgendamata tekst, iga lause eraldi real**

- Ajakirjandustekstid (zip-fail 586 Kb)
- Ilukirjandustekstid (zip-fail 438 Kb)



**Vana kirjakeele korpus**

| Avaleht | Otsing | Juhend | Tekstid | Väljaanded |

Eesti vana kirjakeele korpus sisaldab 15. kuni 19. sajandi tekste. Vanemad tekstid on morfoloogiliselt märgendatud, st neist tekstidest saab infot otsida sõnade tänapäevases kirjaviisis algvormide ning vormiinfo järgi.

- **15. ja 16. sajandist** on korpusesse lisatud kõik teadaolevad ja säilinud eestikeelsed tekstid (v.a nimeloendid), nii käsikirjad kui ka trükised.
- **17. sajandist** on korpusesse lisatud enamik säilinud trükitekste.
- **18. ja 19. sajandist** on lisatud valik trükitekste. Märgendatud on osa tekstidest.
- **19. sajandi II poolest** on lisatud vallakohtute protokollid, mis on automaatselt märgendatud; vt täpsemalt M.-L. Pilvik, K. Muischnek, G. Jaanimäe, L. Lindström, K. Lust, S. Orasmaa, T. Türna "Mõistus sai kuulotedu: 19. sajandi vallakohtuprotokollide tekstidest digitaalse ressursi loomine".

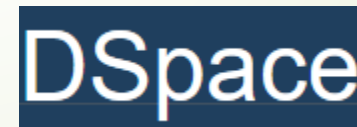Esimene Eestis trükitud eestikeelne raamat (Stahl, 1637)

# Linguistic corpora

- Texts have high quality (though spelling was not in focus)

- Relatively small though.

# Digitized texts

- For the past 20-30 years, a huge number of texts have been digitized.
  - In various collections and formats
  - With varying editing practices

Eesti Kirjandusmuuseum
Estonian Literary Museum

Kreutzwaldi sajand
EESTI KULTUURILOOLINE VEEB

DIGAR
DIGITAALARHIIV

DSpace

KULTUURIMINISTEERIUM

ETERA

WIKISOURCE

# Examples of text collections

- A variety of interfaces and access points
  - Usually text mining is not an option they consider
  - But don't mind either

# Getting the text

- Files of different format into text

# Getting the text

- Finally, will have files that contain raw text
  - Quality will vary

# Metadata

- Again in various formats
- Ought to be collected

```
pro0004___Bornhöhe, Eduard___"Tallinna narrid ja narrikesed"___1892___lk 3-12
pro0005___Bornhöhe, Eduard___"Willu wõitlused"___1890___lk 5-43
pro0006___Bornhöhe, Eduard___"Würst Gabriel"___1893___lk 229-240
pro0007___Eisen, M. J.___"Hiiu köster ja Saare kirikhärra"___1893___lk 43-52
pro0008___Eisen, M. J.___"Tartu saladused"___1891___lk 7-18
pro0009___[???]___"Hella"___1890___lk 3-12
pro0010___Hermann, K. A.___"Lapse mälestus"___1896___lk 1-8
pro0011___Hermann, K. A.___Rikka ja waese pulmad. Külajutukene.___1899___Lk. 29-37
pro0012___Hermann, K. A.___Uudisjutud Eesti rahwa elust.___1895___Lk. 37-45
pro0013___[???]___"Hugo ja Tekla". Jutustus Toolse (ranna) lossist ja ritterist, mis seal 15. aast
pro0014___Tüll, He                    ema imelikud juhtumised wõera mere-saarte peal.___1891
pro0015___Jaanus,                     olm juttu noorerahwa elust."___1893___lk 3-11
pro0016___Jakobson                    399___lk 34-39
pro0017___Järw, J.                    ustus Eesti minewikust."___1892___Lk. 127-136
```

Ein einfeltige weise zu Beten, fur einen guten freund Mart. Luther.

View/Open
- r_iii_i_252i_10.pdf (17.68Mb)
- Tekstituvastusega (1.431Mb)

Date
1535

Author
Luther, Martin

Metadata
Show full item record

sordiaretaja

1.IX 1883 Harjumaal Kloostri valla Karilepa-Tõnul. Taluperemees Siim A., Liisa Vrager. Ae 1917 Anna Maria Volmer. Ants (1918, vt), Valve Jaagus (1920, vt), Ilmar (1927--92, vt). Vasalemma vallakool 1893--97, Paldiski algkool 1897--1900, Haapsalu linnakool 1902--03, Peterburi linnukasvatuskursus 1911--12, põllumajandusdoktor 1947. Sõjaväes Peterburis, vangistatud 1905--06, I maailmasõjas suurtükiväeametnik Tallinnas 1914--16, Eesti Sõjaväelaste Keskbüroo juhatuse liige, Harju maakonnanõukogu liige ja sekretär 1917--20, 1918 Eesti diviisi staabi mobilisatsiooniosakonna asjaajaja, Vabadussõjas ülemjuhataja staabi majandusülem, Eesti Sordiparanduse Seltsi Jõgeva sordikasvatuse osakonnajuhataja 1920--50, "Väikelooma-kasvataja" toimetaja 1919--21, ENSV TA korrespondentliige 1946, ÜN 1947—50 (II ks); ENSV teeneline teadlane 1945, NE preemia 1947, Stalini preemia 1948. Surnud 19.I 1950 Jõgeval, maetud Tartu Raadi kalmistule. EAT, EE, EBLi, EAT2, ENE1, ENE2, EABL, VjV, ENE2(14), ETeadBL – ISOTAMM 2

# Metadata

- General indexes (e.g. Estonian National Bibliography)

Estonian National Bibliography

Publication info

Links to digitized text

Person entry ids

RAHVUSRAAMATUKOGU AVAANDMED

data.digar.ee

Search

DIGAR    DEA    ERB    ISIKUD/KOLLEKTIIVID    Andmekaevurile

VIAF Näidiskirje

ISIKUD MARC21XML formadis
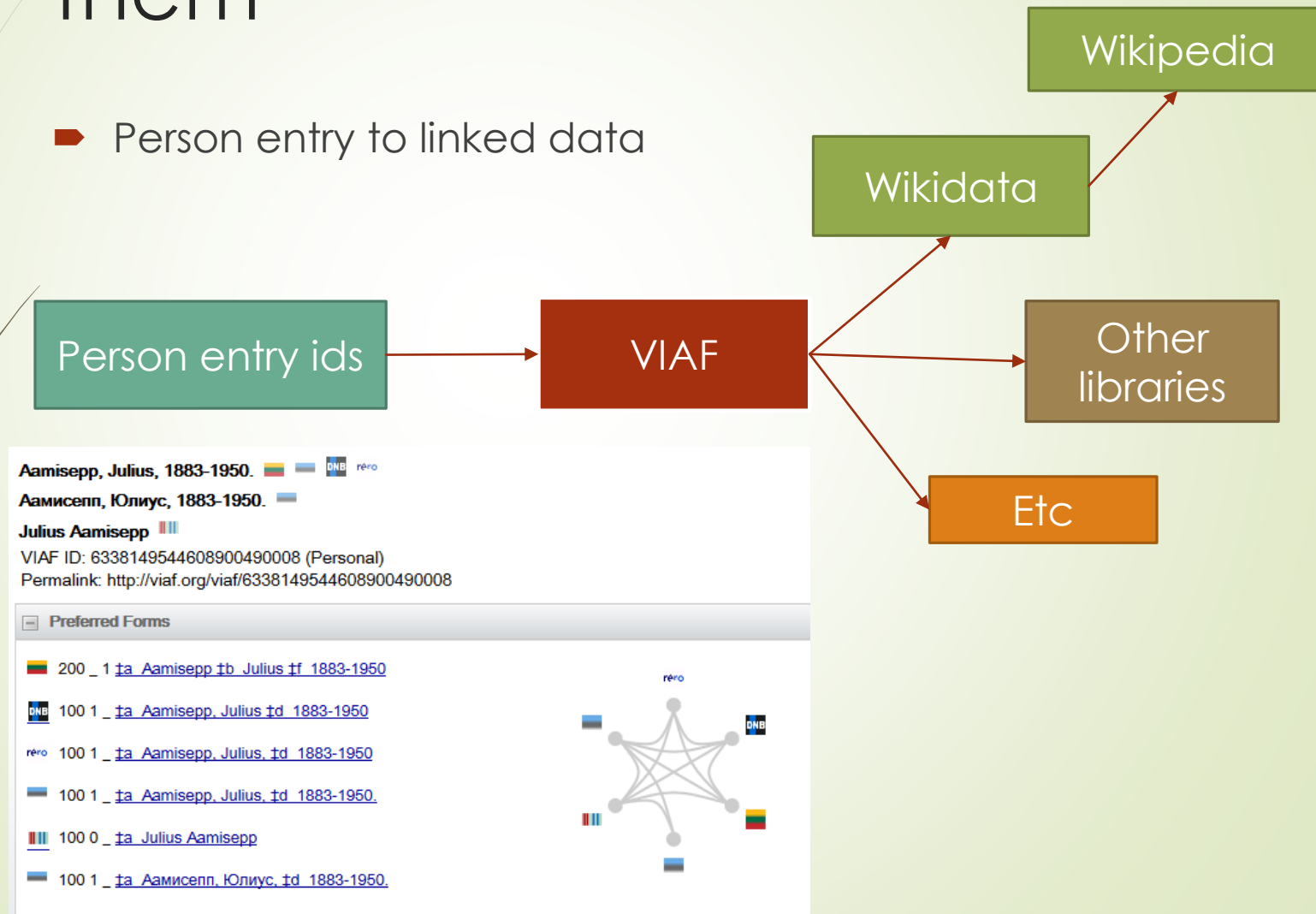
VIAF Näidiskirje

ISIKUD RDF formaadis
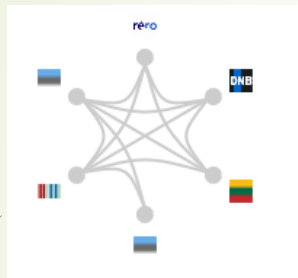
MARC21XML OAI-PMH

Eestikeelne raamat

MARC21XML OAI-PMH

Muukeelne raamat

# Collecting them all, linking them

■ Person entry to linked data

Wikipedia

Wikidata

Person entry ids → VIAF

Other libraries

Etc

Aamisepp, Julius, 1883-1950.
Аамисепп, Юлиус, 1883-1950.
Julius Aamisepp
VIAF ID: 6338149544608900490008 (Personal)
Permalink: http://viaf.org/viaf/6338149544608900490008

Preferred Forms

200 _ 1 ‡a Aamisepp ‡b Julius ‡f 1883-1950

100 1 _ ‡a Aamisepp, Julius ‡d 1883-1950

100 1 _ ‡a Aamisepp, Julius, ‡d 1883-1950

100 1 _ ‡a Aamisepp, Julius, ‡d 1883-1950.

100 0 _ ‡a Julius Aamisepp

100 1 _ ‡a Аамисепп, Юлиус, ‡d 1883-1950.
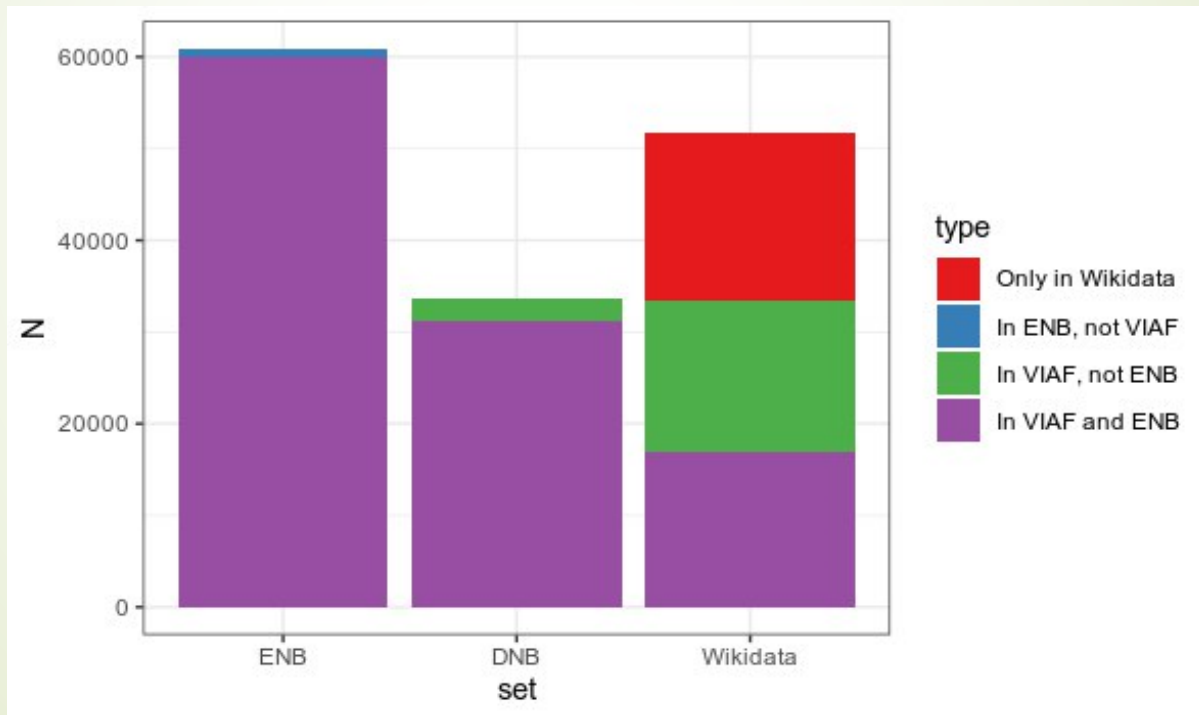
# Aggregating the data

# Manual linking



Manual linking

# Harmonization

- For archival puproses, accuracy is key.
- For analysis, data harmonization is often needed.

```
183  works[kirjastus=="w.bormi pär.",kir
184  works[kirjastus=="f.w.borm",kirjastu
185  works[str_detect(kirjastus,"noor ee
186  works[str_detect(kirjastus,"mutsu")
187  works[str_detect(kirjastus,"ploompu
188  works[str_detect(kirjastus,"postime
189  works[str_detect(kirjastus,"postime
190  works[str_detect(kirjastus,"h\\.laa
191  works[str_detect(kirjastus,"eesti k
192  works[kirjastus=="eesti kirjanduse
193  works[kirjastus=="hermann",kirjastu
194  works[kirjastus=="hermann'i raamatu
195  works[kirjastus=="pealadu hermanni rmtkpl.",kirjastus:="k.a.hermann"]
196  works[kirjastus=="hermanni raamatukauplus",kirjastus:="k.a.hermann"]
197  works[kirjastus=="hermann'i rmtkpl.",kirjastus:="k.a.hermann"]
198  works[kirjastus=="pealadu hermanni kaupluses",kirjastus:="k.a.hermann"]
199  works[kirjastus=="leoke'se kirjastus",kirjastus:="h.leoke"]
200  works[kirjastus=="leoke",kirjastus:="h.leoke"]
201  works[kirjastus=="leoke'se raamatuäri antikvariaat",kirjastus:="h.leoke"]
```

```
29  #works <- works#[!is.na(koht),.N,by=koht]
30  works <- works[,koht:=str_replace_all(koht,"s$","")]
31  works <- works[,koht:=str_replace_all(koht,"l$","")]
32  #works <- works[!duplicated(works)]
33  #places[places!=""&places!="S.l."&places!="S.l.,"&places!="S. l."&places!="S. l.,"]
34  works <- works[koht!=""&koht!="S.l."&koht!="S.l.,"&koht!="S. l."&koht!="S. l.,"]
35  works[str_detect(koht,"Paide"),koht:="Paide"]
36  works[str_detect(koht,"Weissenstein"),koht:="Paide"]
37  works[str_detect(koht,"Вейсенштейн"),koht:="Paide"]
38  works[str_detect(koht,"Вейссенштейн"),koht:="Paide"]
39  works[str_detect(koht,"Haapsalu"),koht:="Haapsalu"]
40  works[str_detect(koht,"Hapsa"),koht:="Haapsalu"]
41  works[str_detect(koht,"Гапсаль"),koht:="Haapsalu"]
42
43  works[str_detect(koht,"Keila"),koht:="Keila"]
44  works[str_detect(koht,"Rakvere"),koht:="Rakvere"]
45  works[str_detect(koht,"Wesenberg"),koht:="Rakvere"]
46  works[str_detect(koht,"Везенберг"),koht:="Rakvere"]
```
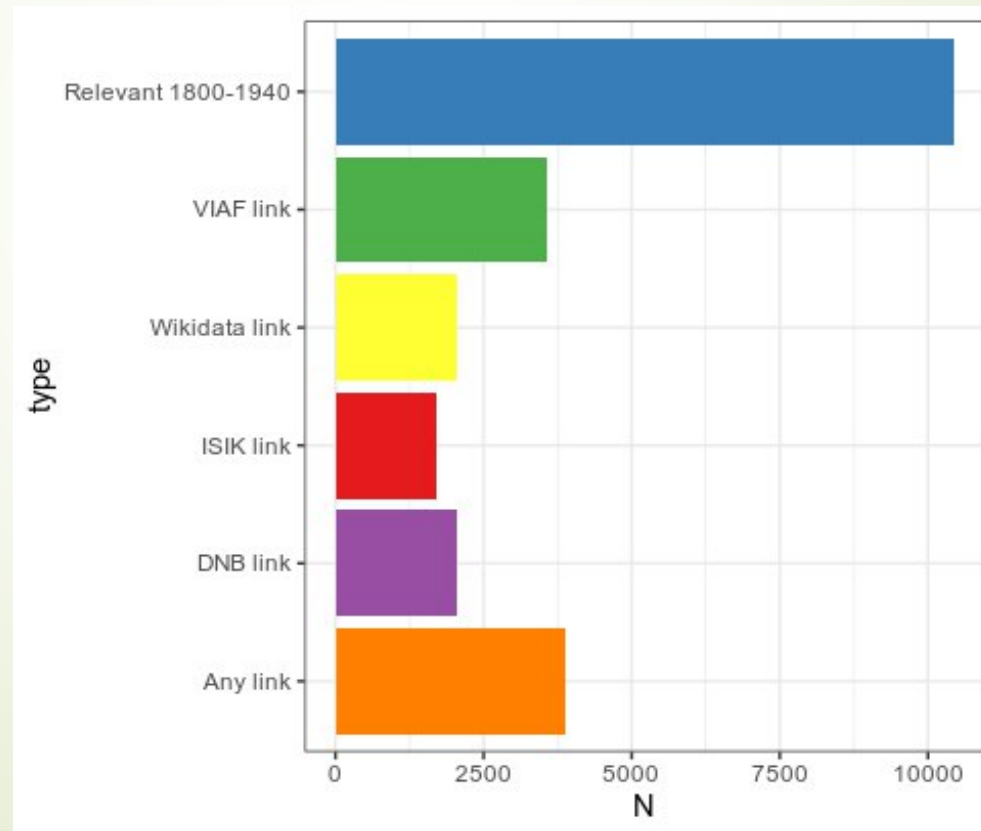
# Results of linking

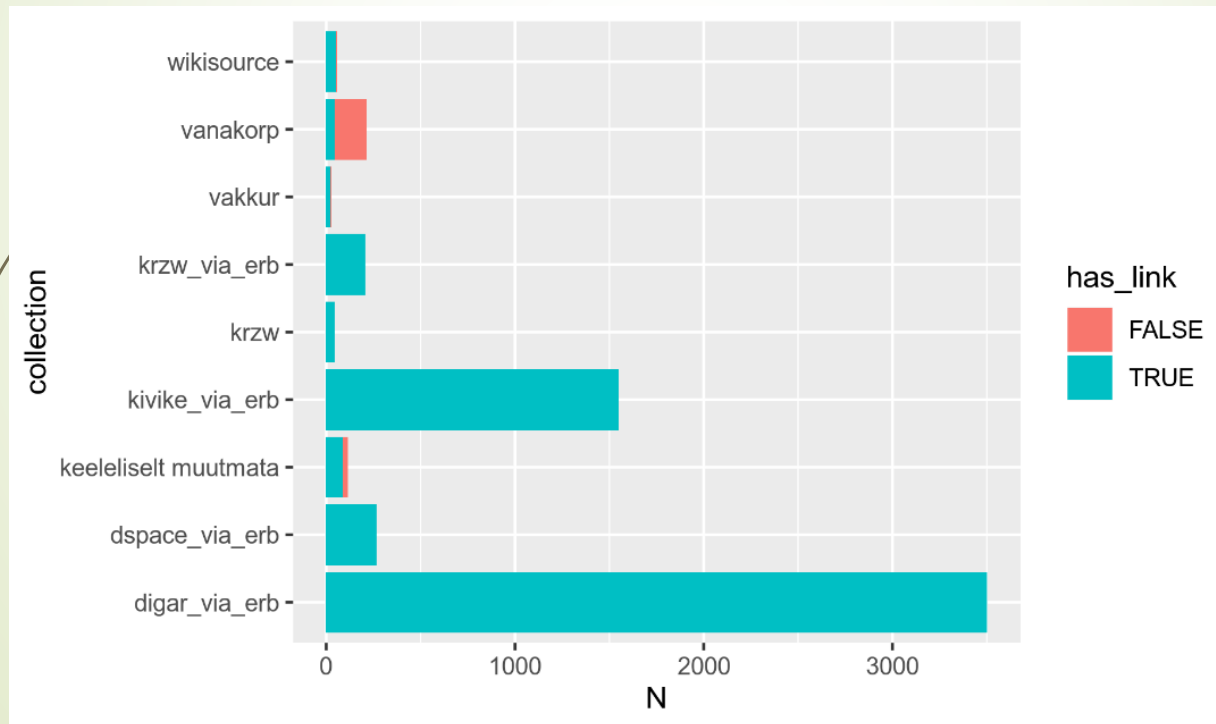- Linkability of archives

# Results of linking

- Linked author metainformation

# Results of text collection

- 5132 texts total, of 4186 unique publications

# Processing the corpus

- Raw texts can be processed
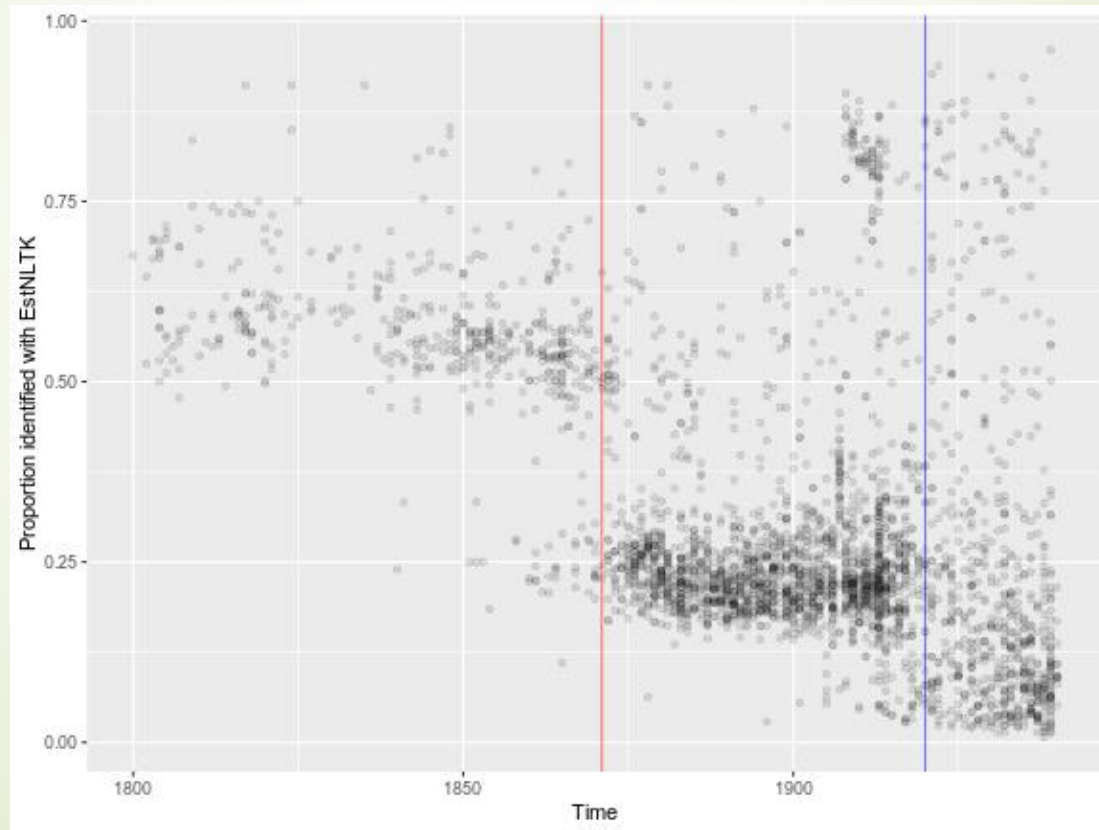- Lemmatize, POS-tag etc (used Python EstNLTK)

```python
from estnltk import Text
text = Text('mõeldud')
text.tag_analysis()
```

```
[('Usjas', 'A', 'omadussõna algvõrre'),
 ('kaslane', 'S', 'nimisõna'),
 ('ründas', 'V', 'tegusõna'),
 ('künklikul', 'A', 'omadussõna algvõrre'),
 ('maastikul', 'S', 'nimisõna'),
 ('tünjat', 'A', 'omadussõna algvõrre'),
 ('Tallinnfilmi', 'H', 'pärisnimi'),
 ('režissööri', 'S', 'nimisõna')]
```
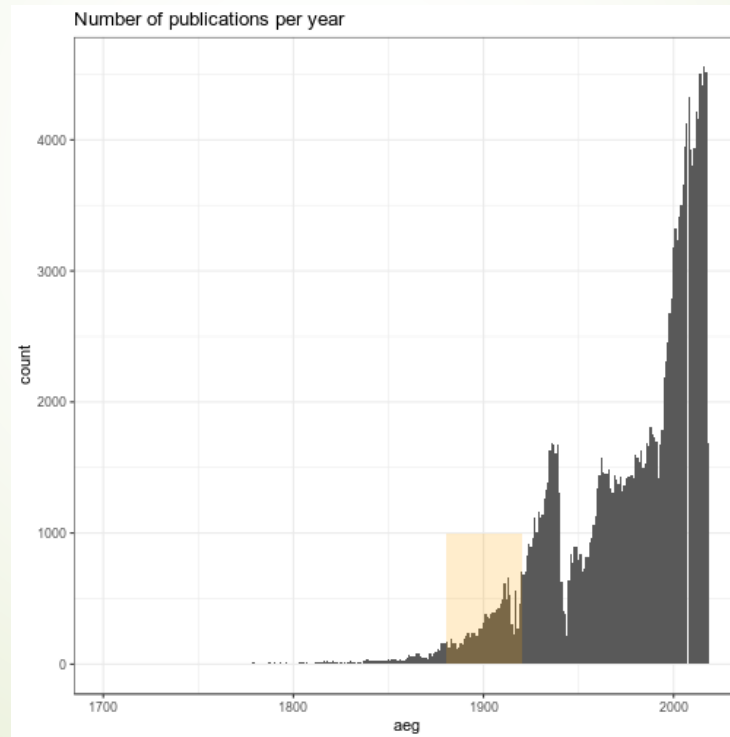
```
{'paragraphs': [{'end': 7, 'start': 0}],
 'sentences': [{'end': 7, 'start': 0}],
 'text': 'mõeldud',
 'words': [{'analysis': [{'clitic': '',
     'ending': '0',
     'form': '',
     'lemma': 'mõeldud',
     'partofspeech': 'A',
     'root': 'mõel=dud',
     'root_tokens': ['mõeldud']},
    {'clitic': '',
     'ending': '0',
     'form': 'sg n',
     'lemma': 'mõeldud',
     'partofspeech': 'A',
     'root': 'mõel=dud',
     'root_tokens': ['mõeldud']},
    {'clitic': '',
     'ending': 'd',
     'form': 'pl n',
     'lemma': 'mõeldud',
```

# Results of the corpus

- NLP success rate over time

# The published texts



Number of publications per year

# Representativeness

# Publishing industry

➮ More authors and open genres

# A few cities dominate

# In corpus

- s

# Publishers

- Top publishers (harmonized)

# Publishers

- Top publishers (harmonized) within corpus

# Authors

- Top authors (harmonized) within corpus

# Other types of data

- There is a lot of relevant scholarship.
- Just thinking about it in terms of „data" is new.

- Needs creativity in finding and utilizing the data

# Biographic data



- ➤ Assemble life histories
  - ➤ from various sources

### Eesti biograafiline andmebaas ISIK
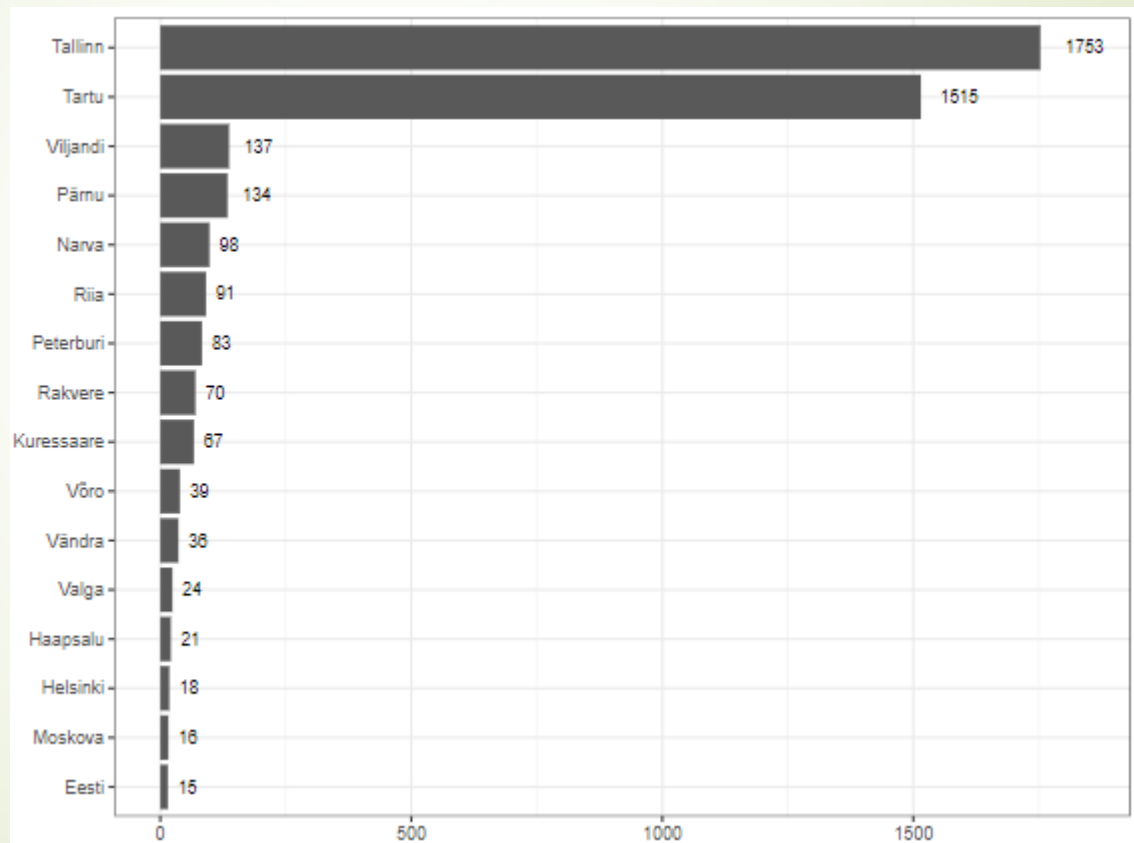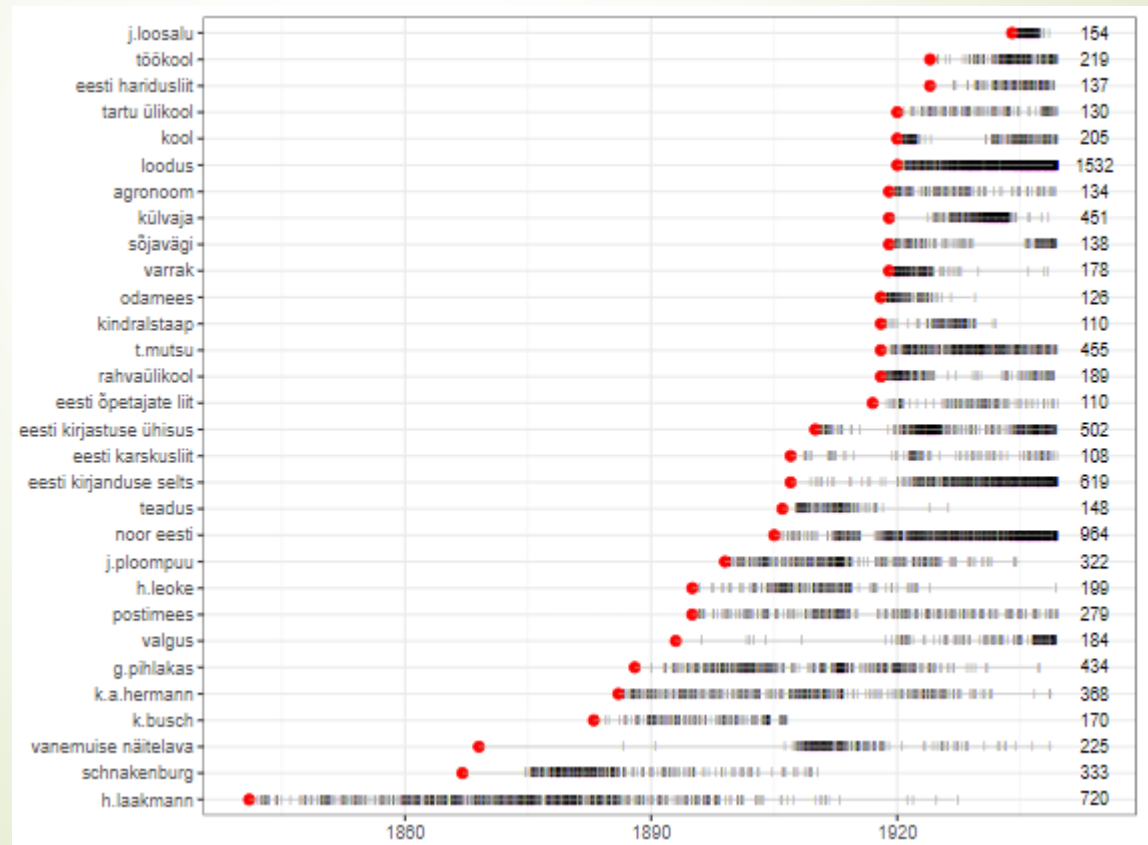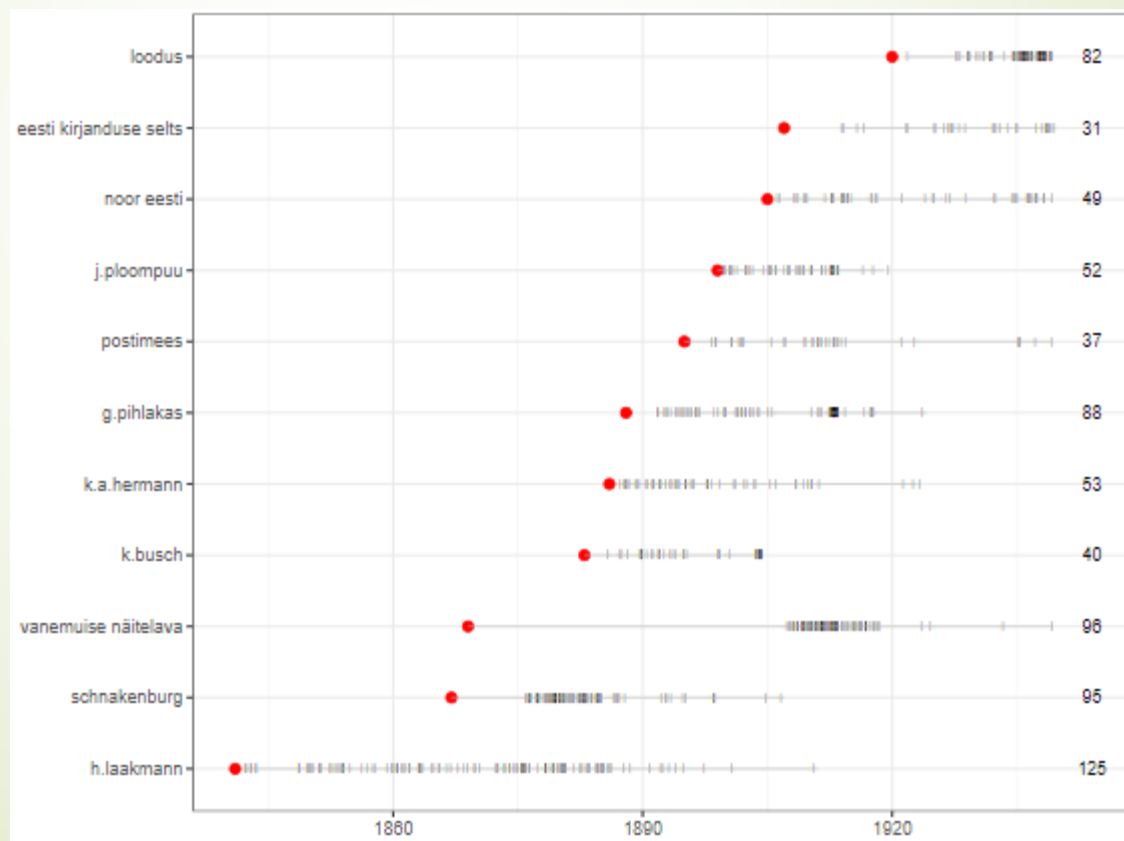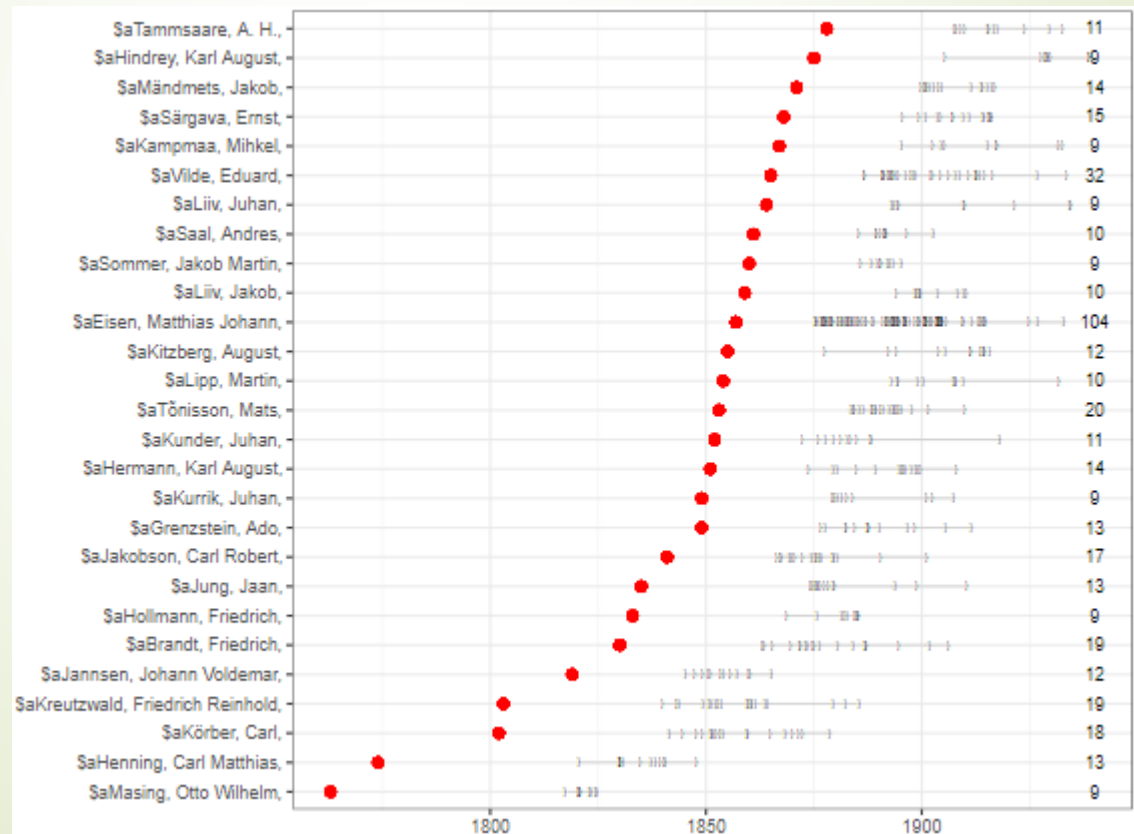
avalehekülg · nimekiri ja otsing · ametikohad [108] · asutused [468] · isikud [4445] · kaastöö väljaannetele [1131] · nime erikujud [925] · pseudonüümid [9179] · sugulased [309] · sünnikohad [624] · surmakohad [186] · tegevusala [599] · lühendid ja allikad · admin

**SAAL, ANDRES**

**Nimi**

| Nime normitud kuju | Saal, Andres |
| Allikad | Postimees 1857-1907, 235-236; EAE-40, lk. 200, EE (2000), XIV, lk. 447. |

**Eludaatumid**

| Sünniaeg | 1861-05-21 |
| Sünnikoht | Selja k. Tori v. Pärnumaa |
| Surmaaeg | 1931-06-23 |
| Surmakoht | Los Angeles |

**Üldandmed**

| Tegevusala | kirjanik, õpetaja |
| Biograafilised andmed | SAAL, Andres. 21.V 1861 Tori khk Tori v Laasil. Taluom (ms puusepp) Jaak S. (1817--1902), Ano Lindebaum (ms toatüdr, 1818--92). Ae Emilie Rosalie Moks (1871--1954). Rosa Regina Boiley (1902--30, maalikunstn), Leo Henry Wladimir (1904--65, ins). Selja v-k, Tori khkk, vallak.õp eksam 1880, TU vabakuulaja 1886--89 Cronenbergi fototehnikak Baieris, Viini paljundustehnikainst. Selja v-k õp 1880--84. Oleviku toim.l 1884--90 ja tsinkograaf 1893--97, Trt Jaani kir.k õp 1884--90, toiduainete kaupmees Trt-s 1890--93, fotograaf Frankfurdis 1897--98, topograaf Jaaval 1898--1920, Los Angeleses 1920--31. jäi pimedaks 1928. Kirjanik. Srn 23.VI 1931 Hollywoodis, tuhastati (urn 1932 Trt Maarja kir-s, 1944 kirj.muus-s). PmA, EBL, EE, EBLI, ENE1, EKBL, EKrL, ENE2, EKoBL, ENE2(14) |

| Person | Aamisepp, Julius |
|---|---|
| Geschlecht | männlich |
| Quelle | Pressearchiv des Herder-Instituts Marburg Wikipedia (Stand: 25.09.2018): https://en.wikipedia.org /wiki/Julius_Aamisepp |
| Zeit | Lebensdaten: 1883-1950 |
| Land | Estland (XA-EE) |
| Geografischer Bezug | Geburtsort: Karilepa (Gem. Lääne-Harju) Sterbeort: Jõgeva |
| Beruf(e) | Agrarwissenschaftler Revolutionär |

# Placename locations



Eesti Keele Instituudi kohanimeandmebaas (KNAB)
Place Names Database (KNAB)

# Birthplaces of writers

- Some birthplaces of writers in corpus

# Birthplaces by decade

# Education by decade



Joon. 1. Eesti rahvakoolide võrk 1851. aastal

Public elementary school network in 1851

# Education by decade



Public elementary

Public elementary school network in 1881

# Life histories of writers

# Demographic data

- Censuses and analyses

# Dialect geography

# Dialect predicted values

# Dialect compositions in cities

# Migration and dialects

# Getting to research

Operationalizing the question

# Making the study

- Specify the measured variables

- Measure in suitable texts

- Combine with metadata

- Check the mechanisms

# Variants and variation

- A set of variables:
    - w / v – wabariik, vabariik
    - ää / ea – hää, hea; sääl, seal
    - üi / üü – nüid, nüüd
    - herra / härra
    - naine / naene
    - om / on
- Fairly common in frequency
- Have interesting variation in 1880-1920

# Variants and variation



õนน/õน

# Study on particular variants

# Average trends

# Language prescription

- Historian's overviews



Overview of proposed norms in 1909 (Raag 2008)

# Language prescription

- Prescriptive events on the timeline

- Blue (suggested norm towards 1)

- Red (suggested norm towards 1)

# Language prescription

- Prescriptive events align very well with trends in usage.

- => at least some degree of prescriptive influence.

# Language prescription

# Mechanisms of standardization

- Influences on language usage
  - Home dialect
  - Everyday language use
  - Education
  - Prescription
  - Etc
- What made a suggestion successful?
- Generational change vs individual change?

# Phases of change

- Looking at the aftermath of prescriptive events

- Blue (suggested norm towards 1)

- Red (suggested norm towards 1)

# Simple tendencies in data

- Significant predictors (logistic regression)
  - Model does not include dialect info

# Simple tendencies in data

- Significant predictors (logistic regression)
  - Model does not include dialect info

# Simple tendencies in data

- Significant predictors (logistic regression)
  - Model includes dialect info

# Simple tendencies in data

- Significant predictors (logistic regression)
  - Model includes dialect info

**Ww**

Phase 1 (n = 91) — Trend, Intended, tradition, popsize, edu, youth
Phase 2 (n = 56) — Trend, Intended, tradition, popsize, edu, youth

**ste**

Phase 1 (n = 170) — Trend ↗, Intended ↗, **tradition** ↘, popsize, edu, **youth** ↗
Phase 2 (n = 86) — Trend →, Intended ↘, tradition, **popsize** ↘, edu, youth ↗

**te**

Phase 1 (n = 265) — Trend ↗, Intended ↗, **tradition** ↘, popsize, edu, **youth** ↗

**naene**

Phase 1 (n = 199) — Trend →, Intended ↘, **tradition** ↗, **popsize** ↗, **edu** ↗, **youth** ↘
Phase 2 (n = 87) — Trend ↗, Intended ↗, tradition ↗, **popsize** ↗, **edu** ↗, youth

**ial**

Phase 1 (n = 210) — Trend →, Intended ↗, tradition, **popsize** ↘, edu, youth
Phase 2 (n = 71) — Trend ↗, Intended ↗, tradition, popsize, edu, youth

**eila**

Phase 1 (n = 101) — Trend, Intended, tradition, popsize, edu, youth

**kuda**

Phase 1 (n = 117) — Trend, Intended, tradition, popsize, edu, youth
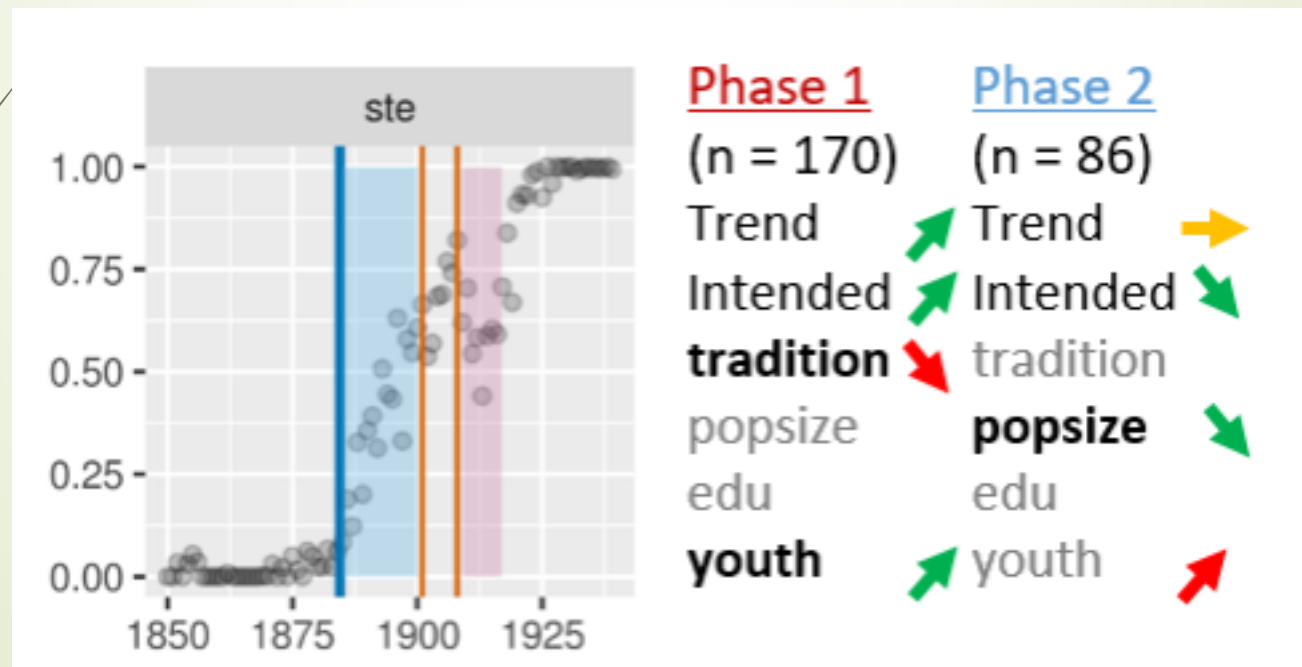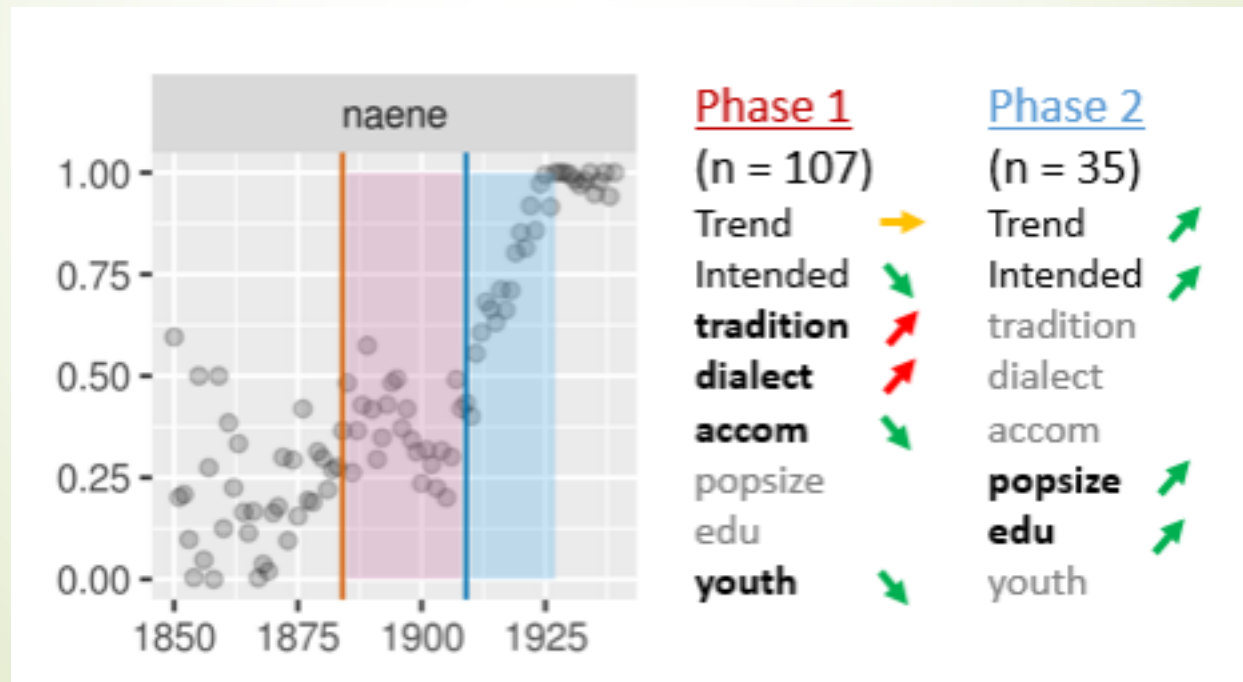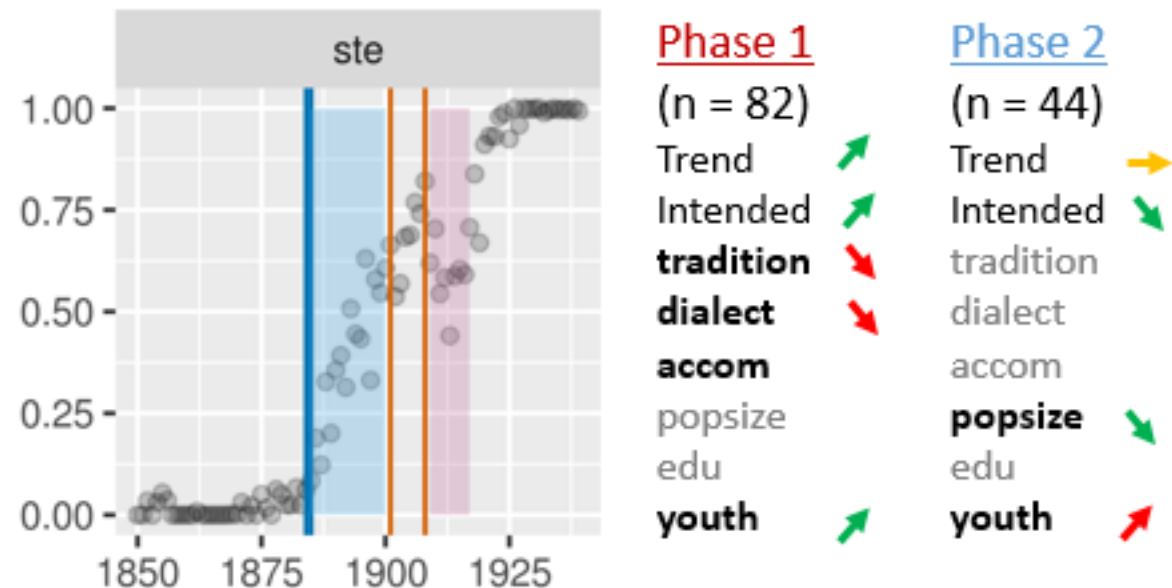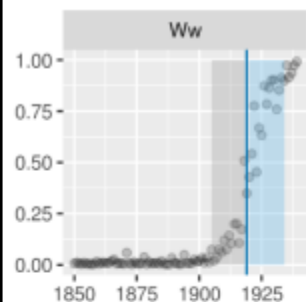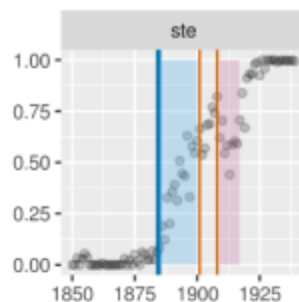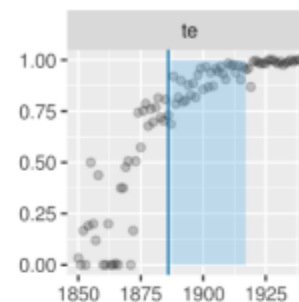Phase 2 (n = 323) — Trend ↗, Intended, tradition, **popsize** ↘, edu, youth

**auu**

Phase 1 (n = 206) — Trend →, Intended ↗, tradition, popsize, edu, **youth** ↘
Phase 2 (n = 71) — Trend ↗, Intended ↗, tradition ↘, **popsize** ↗, **edu** ↗, youth

**õuu**

Phase 1 (n = 200) — Trend →, Intended ↗, tradition, **popsize** ↗, edu, **youth** ↘
Phase 2 (n = 77) — Trend ↗, Intended ↗, tradition, **popsize** ↗, edu, youth

**ää**

Phase 1 (n = 158) — Trend ↗, Intended ↘, **tradition** ↗, popsize, edu, **youth** ↘
Phase 2 (n = 57) — Trend ↗, Intended ↘, **tradition** ↘, **popsize** ↗, edu, **youth** ↘
Phase 3 (n = 93) — Trend ↘, Intended ↘, tradition, **popsize** ↗, edu, youth

**kåsa**

Phase 1 (n = 166) — Trend, Intended, tradition, popsize, edu, youth
Phase 2 (n = 59) — Trend, Intended, tradition, popsize, edu, youth

**üi**

Phase 1 (n = 110) — Trend ↗, Intended, tradition, popsize, edu, youth
Phase 2 (n = 261) — Trend ↗, Intended ↗, tradition, **popsize** ↘, edu, youth

Panel "ste":

Phase 1 (n = 82) | Phase 2 (n = 44)
- Trend ↗ | Trend →
- Intended ↗ | Intended ↘
- **tradition** ↘ | tradition
- **dialect** ↘ | dialect
- **accom** | accom
- popsize | **popsize** ↘
- edu | edu
- **youth** ↗ | **youth** ↘

Panel "naene":

Phase 1 (n = 107) | Phase 2 (n = 35)
- Trend → | Trend ↗
- Intended ↘ | Intended ↗
- **tradition** ↗ | tradition
- **dialect** ↗ | dialect
- **accom** ↘ | accom
- popsize | **popsize** ↗
- edu | **edu** ↗
- **youth** ↘ | youth

Panel "kuda":

Phase 1 (n = 59) | Phase 2 (n = 161)
- Trend → | Trend ↗
- Intended | Intended ↗
- tradition | tradition
- dialect | dialect
- accom | accom
- popsize | **popsize** ↘
- edu | edu
- youth | youth

Panel "ää":

Phase 1 (n = 77) | Phase 2 (n = 28) | Phase 3 (n = 45)
- Trend ↘ | Trend ↗ | Trend ↘
- Intended ↘ | Intended ↗ | Intended ↗
- **tradition** ↗ | tradition | tradition
- dialect | dialect | dialect
- accom | accom | accom
- popsize | popsize | **popsize** ↗
- edu | edu | edu
- **youth** ↘ | youth | **youth** ↗

Panel "üi":

Phase 1 (n = 49) | Phase 2 (n = 128)
- Trend → | Trend ↗
- Intended | Intended ↗
- tradition | tradition
- dialect | **dialect** ↗
- accom | accom
- popsize | **popsize** ↘
- edu | edu
- youth | youth

# A basic overview

- Significant factors across all phases
  - Top (no dialect info), bottom (with dialect info)

# Study conclusions

- Main hypothesis (late 19c Estonian):
  - Spontaneous changes align with an influx of new authors
  - Prescriptive events line up clearly with trends
  - Some generational influence apparent: successful changes are led by young and large cities.
- Study generally:
  - Data is everywhere and has not been collected
  - Already „first look" at it can give insights
  - Thinking in terms of data also adds to theory

# General discussion

- 80% of DH research is finding the data and getting it ready
    - Manual + technological work
    - Good-enough solutions for current tasks
- Open collections can have a variety of uses
    - Difficult to predict
- Also by-products of research can be useful
    - Databases & processing
    - Visualizations & understanding

# Some advice/thoughts

- Start & stay open, document everything!
- Think how your current work fits to the whole project
  - Work with manageable chunks and fit them in right away
- Don't worry about lack of skills beforehand, but consult specialists!
  - Usually your problems are not new and have solutions.
  - However working with „data" is new to humanities so you may need to step outside your institute/department.

# Interactive plot links

- [https://peetertinits.github.io/slides/DH-labs-2019/prese_DH-lab_2019.html#14](https://peetertinits.github.io/slides/DH-labs-2019/prese_DH-lab_2019.html#14)

- [https://peetertinits.github.io/slides/EDHC2018/prese_EDHC_2018.html#45](https://peetertinits.github.io/slides/EDHC2018/prese_EDHC_2018.html#45)

- [https://peetertinits.github.io/slides/EDHC2018/prese_EDHC_2018.html#50](https://peetertinits.github.io/slides/EDHC2018/prese_EDHC_2018.html#50)

- [https://peetertinits.github.io/slides/plots/authorbirths.html](https://peetertinits.github.io/slides/plots/authorbirths.html)

- [https://peetertinits.github.io/slides/plots/samamaakond.html](https://peetertinits.github.io/slides/plots/samamaakond.html)