

Integrating Open Science in the Humanities: the Case of Computational History

Better Science Through Better Data 2019; #scidata19
Wellcome Collection, London
6 November 2019

Mikko Tolonen (University of Helsinki)

`mikko.tolonen@helsinki.fi`

Outline of the talk

1. (Open) Science in the Humanities?
2. Challenges of Humanities Data
3. Integrating Open Science to Humanities in Computational History
4. (Experiments to go around problems of noisy data)
5. Conclusion

Helsinki Computational History Group

“Computational history” refers to an integrated mixed methods approach to study large digitized historical sources. “Integrated” means that data science is combined to specialized subject knowledge; in the case of COMHIS, intellectual history and book history.

<http://helsinki.fi/computational-history>

Better Science Through Better Data

(Open) Science in the Humanities?

Science and hermeneutics

Tangible objects



Subjective experience



Need for
mixed
methods!

photo: Time Machine project

Three aspects of OA in the Humanities

Raw/primary data

- seldom access to full data
- special arrangements when access (no possibility to share further)
- data providers reluctant to share openly

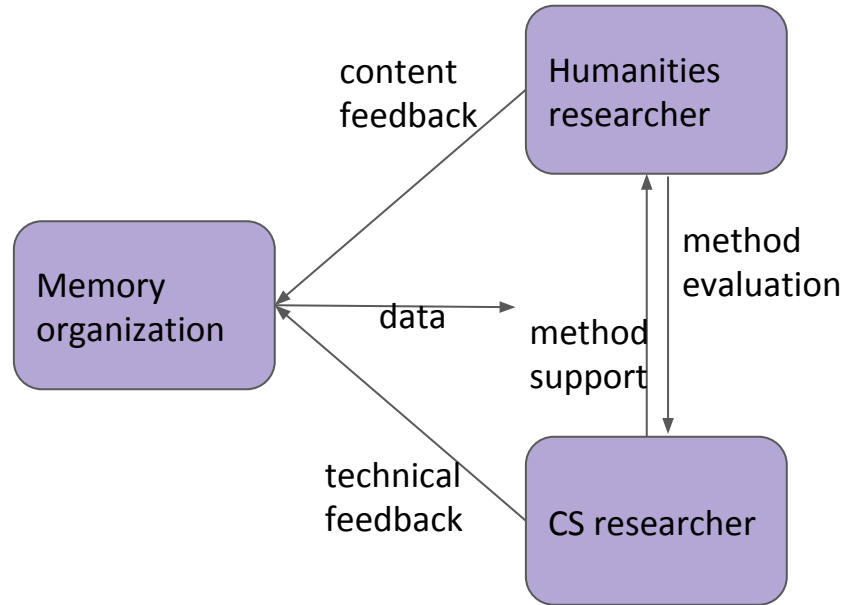
Research data

- Scholars often think they have no research data
- idea of reproducibility vague/non-existent
- little/no credit in opening research data

Open publications

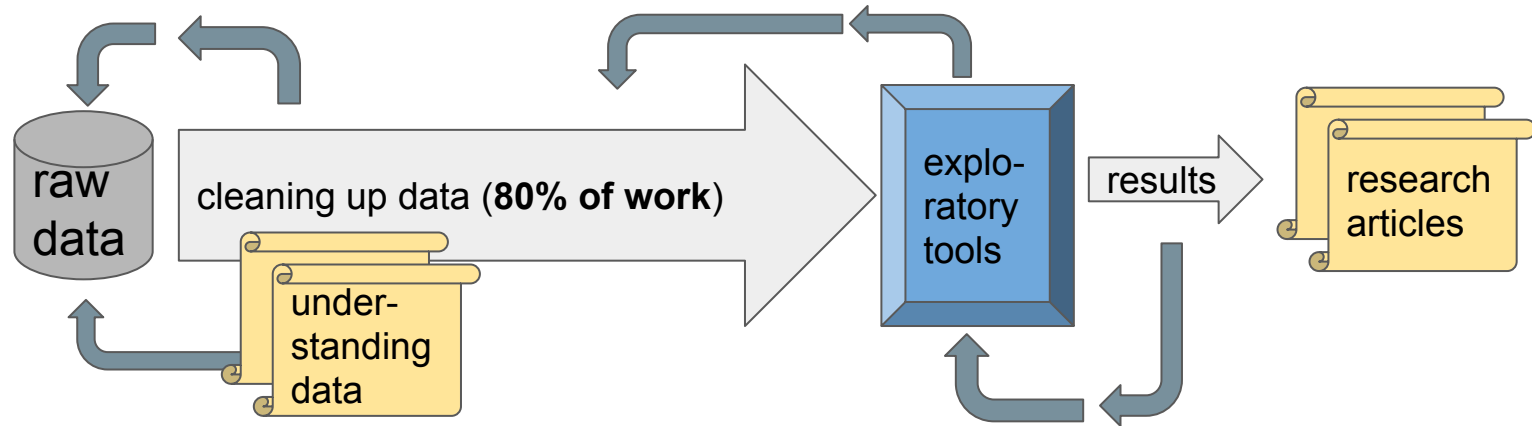
- Idea that humanities cannot afford OA
- Some think Plan-S is hurting humanities
- Scholars don't grasp the value or relevance of OA even with respect to publications

Better data in Humanities calls for well-functioning ecosystems



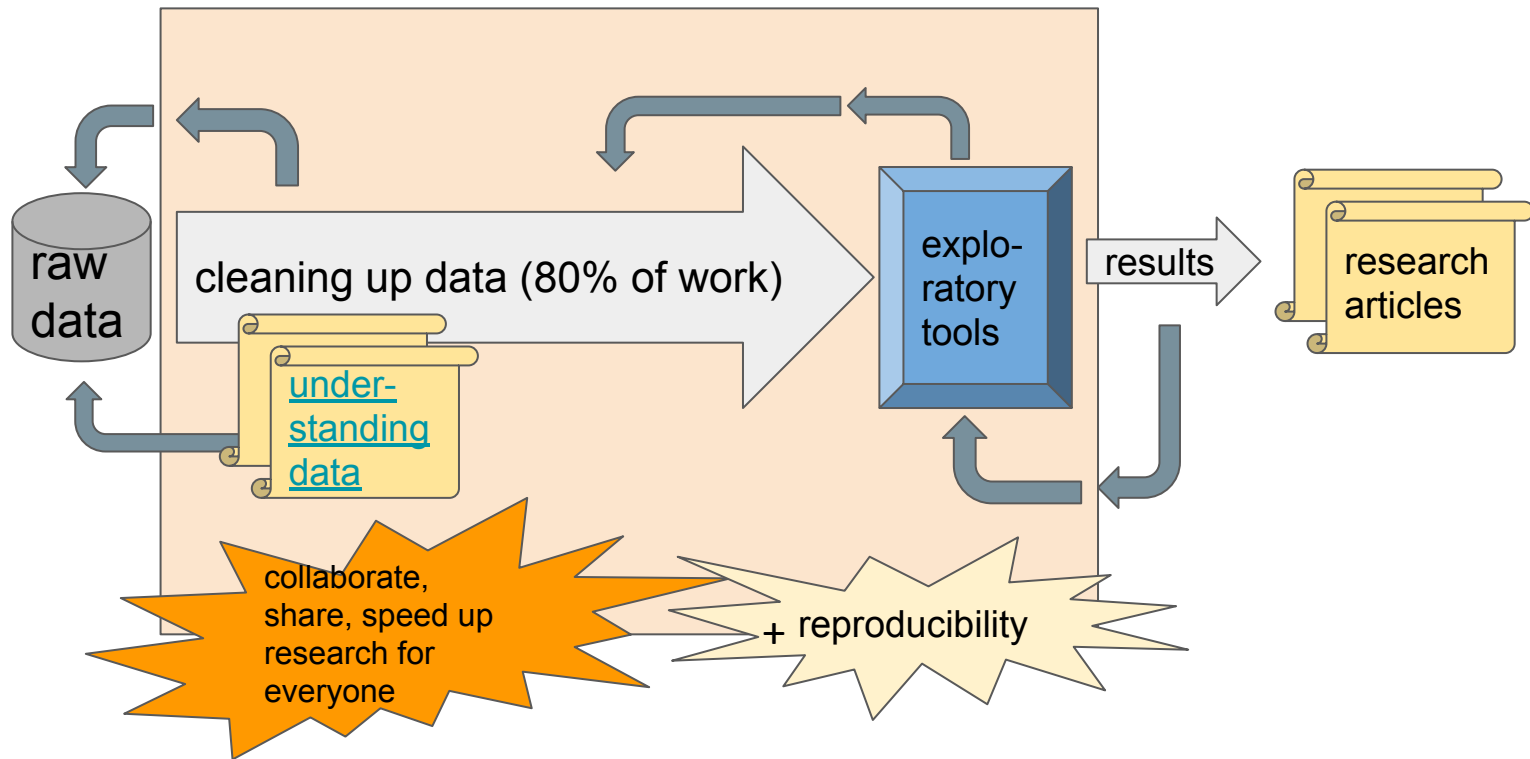
UNIVERSITY OF HELSINKI

Computational history research process



80% of your time for data cleanup, another
80% for algorithms, ...

In an ideal world: collaboration & open science workflows to reduce individual workload



credit: Eetu Mäkelä & Mikko Tolonen

Better Science Through Better Data

Challenges of Humanities Data

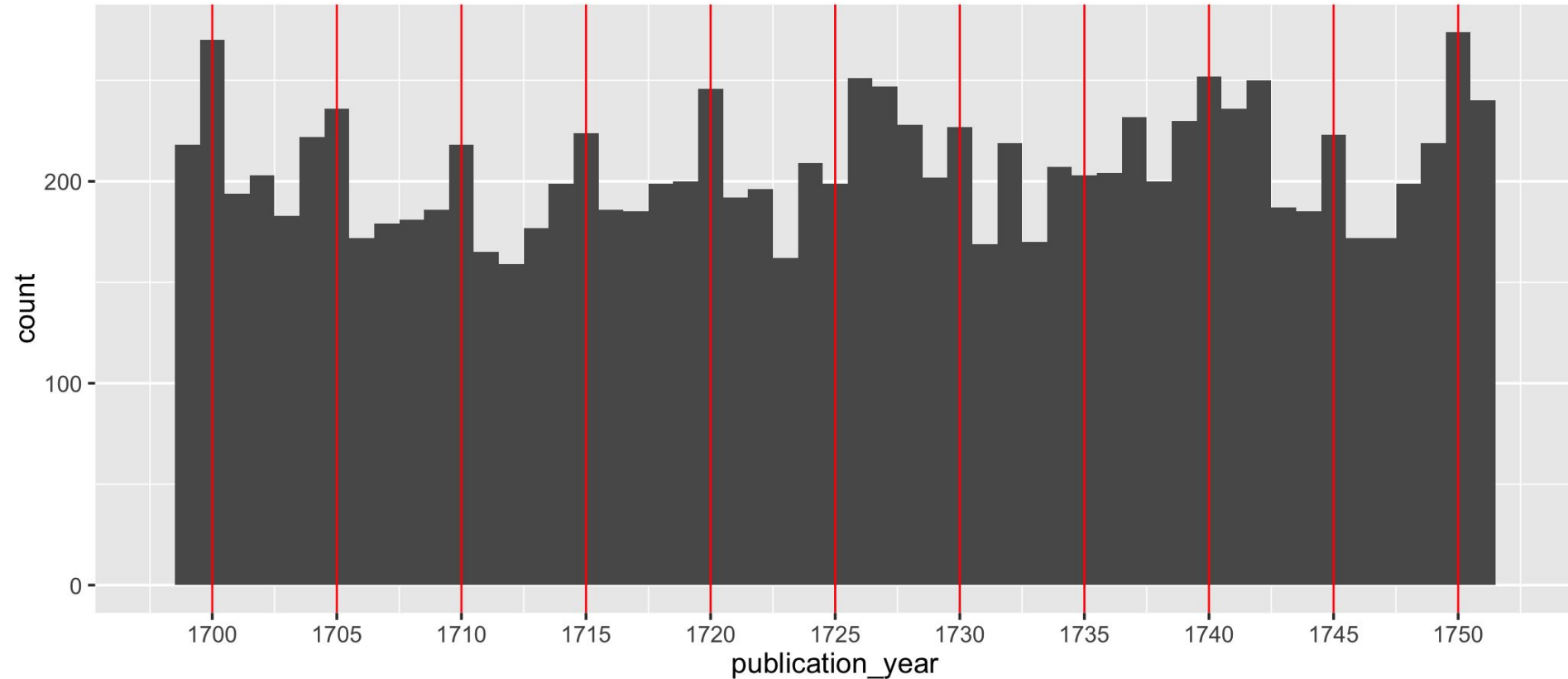
Nature of digitized humanities data

- Collections of collections
- Relevance of source criticism
- Stacking more sources does not often solve problems of bias!



Wellcome
Collection

Example of bias that is particular to humanities data:
The 5-year theory with respect to ESTC data



Burden of empty promises

**3.5 Million Books
1800-2015**

**Internet Archive +
HathiTrust**

Challenges with open data

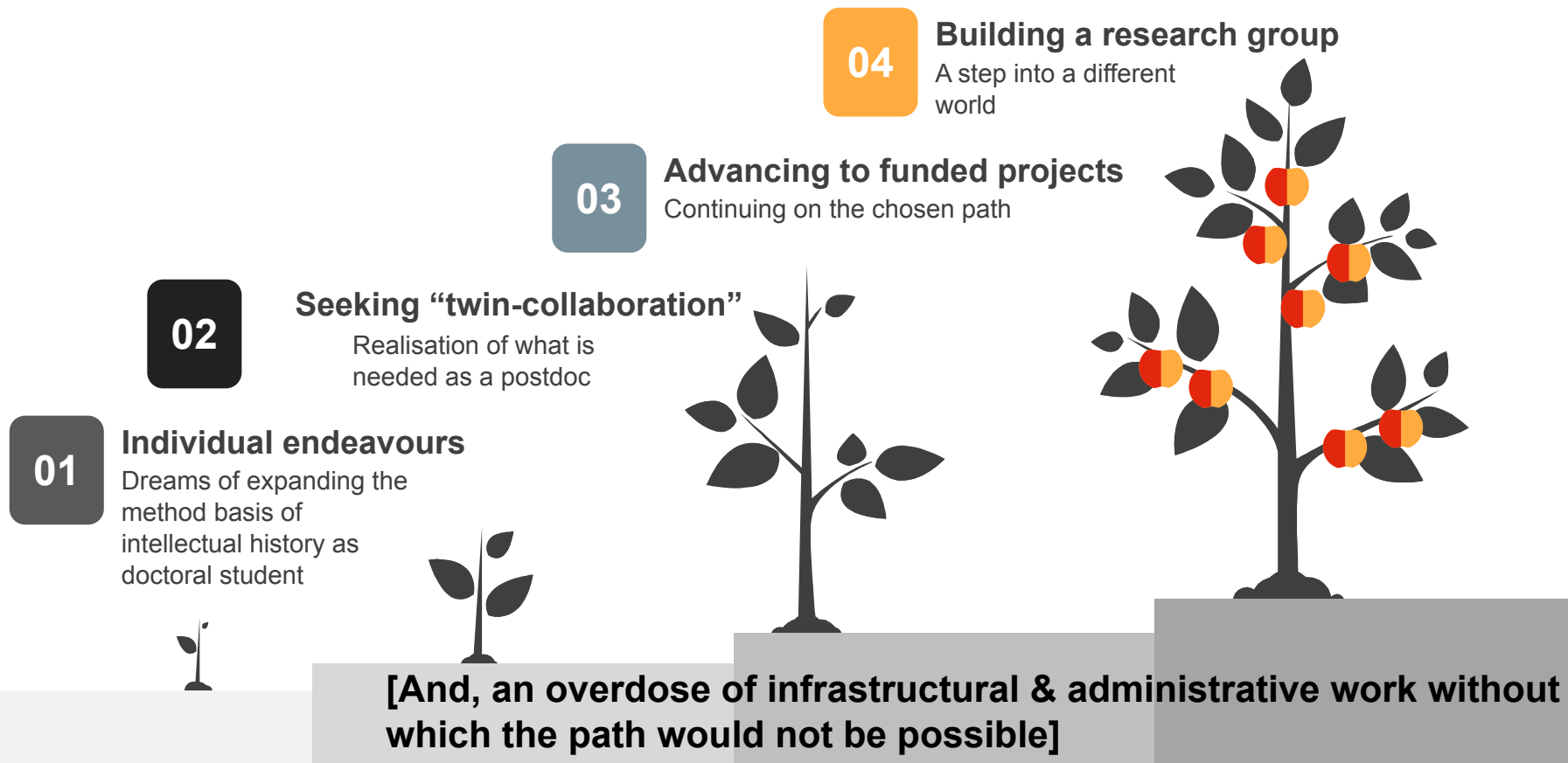
- **Institutions** reluctant to give full access to data. Why?
- **Research process** is not opened and research data is not shared in the Humanities. Transparency, reproduction, collaboration, new initiatives are missing. Why?

Short answer: **Cultural change takes time**. We need concrete examples in the core field of the Humanities that actually prove OPEN DATA PRINCIPLES as useful.

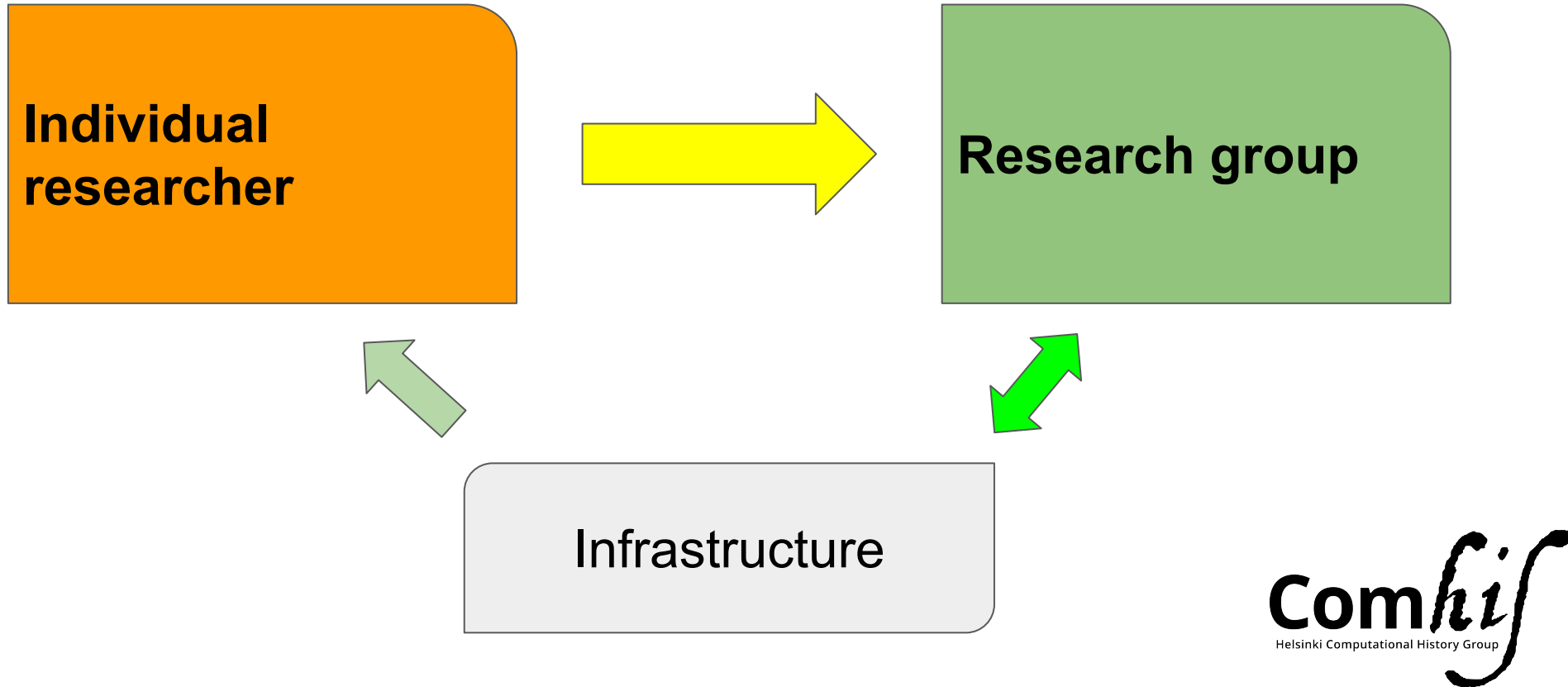
Better Science Through Better Data

Integrating Open Science to Humanities in Computational History

Tradition & research culture



Key factor for the humanities of the future



Helsinki Computational History Group's aim:

Understanding public communication in early modern Europe

Movement of ideas

- Metadata work based on several different library catalogues
- genres (poetry, pamphleteering); intellectual traditions (natural law tradition, ancient texts)
- text reuse: genres (historical works, quoting practices)

Research data releases

- ESTC; Fennica; Kungliga; CERL; ECCO text reuse (+ EEBO text reuse); Finnish Newspapers

Conceptual change

- concepts are crucial, but not directly jumping into this for various reasons
- Theoretical underpinning (historians + linguists)
- Concepts as linguistic objects (linguists + historians + CS)

Tools for others

- UIs, APIs, shiny apps etc.

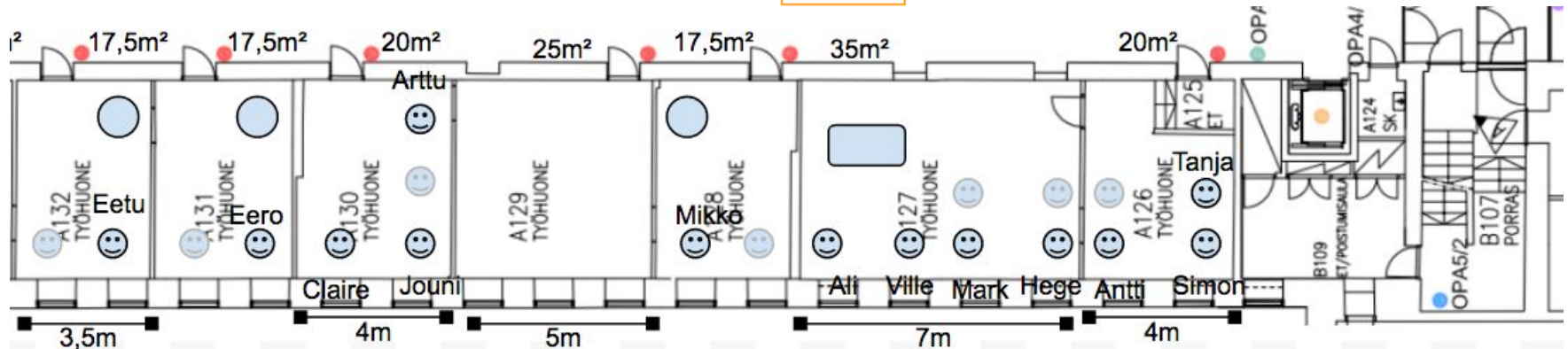




2013



2017



😊 = pysyvä paikka

😊 = fleksipaikka (Jani, Viivi, Ylva, Erik, Mikko, Erkki, Esko, Petri, ...)

Data mining in computational history

Text mining of large corpora ↔ Metadata as a quantitative tool

- **Objective:** understanding conceptual change, uses of language
 - **Sources:** full-text databases (ECCO, EEBO, Finnish Newspapers etc.)
 - **Potential:** Theoretically great, the future?
 - **In practice:** raw data almost never openly available; if it is, tied to limited interfaces
 - **Scalability with open research data:** data-driven approach
 - **Methodological perspective:** Messy to study historical sources, intellectual input not guaranteed.
- **Objective:** Quantitative study of material objects
 - **Sources:** World is full of different metadata collections
 - **Potential:** Greatly underestimated (even by librarians)
 - **In practice:** difficulties with open access to raw data and supporting data sources, but not impossible.
 - **Scalability with open research data:** fantastic
 - **Methodological perspective:** excellent for borrowing best practices from other scientific fields. Quality of catalogues varies.

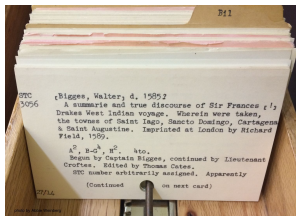
➡ It is the combination of harmonized (and better) metadata and full-text sources that will lead to better science in early modern intellectual history (ECCO and ESTC), for example!

Build it yourself! Humanities questions guiding our tool building

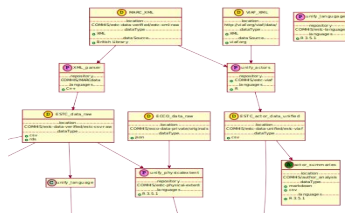
- **Bibliographic sources** as basis of networks and time series to inform the intellectual endeavours.
- **Text-reuse detection** to study influence (using BLAST to deal with OCR-mistakes).
- **Materiality explorations** of printed items based on information derived from layout, font etc.
- **Stylometry** to study particular questions of authorship.

- **Topic modelling and word embeddings** (etc.) to explore conceptual change
- **Detecting argumentative structures** based on syntax.

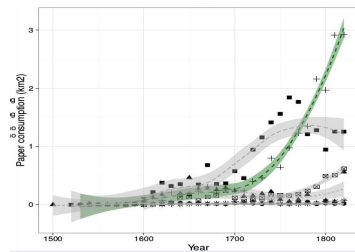
Research potential



Open data science ecosystem



Research cases



Bibliographic Data Science and the History of the Book (c. 1500–1800)

Leo Lahti , Jani Marjanen , Hege Roivainen & Mikko Tolonen

Our virtuous cycle of better data & better science

- Combining **harmonized metadata to full-text sources (ESTC & ECCO)** -> Enables text mining in a new way, upcoming this academic year.
- Using **full-texts to enrich metadata (ESTC & ECCO)** -> Feeding back to the loop, better quality data, detecting subject/topics for example.
- Combining **text reuse information to metadata (ESTC & ECCO)** -> feeds back to edition information.
- **Re-OCRing (ECCO)** -> Feeds back to all processes that combine ECCO and ESTC.

Experiments to go around problems of noisy data

Text reuse detection to study influence

- BLAST -bioinformatics software
- computational analysis on the whole ECCO, that is +200,000 texts
- ~130,000 reuse fragments in Hume's *History* alone
(Whole ECCO has millions of interlinked reuse fragments)
- ~150 - 3000 characters / reuse fragment

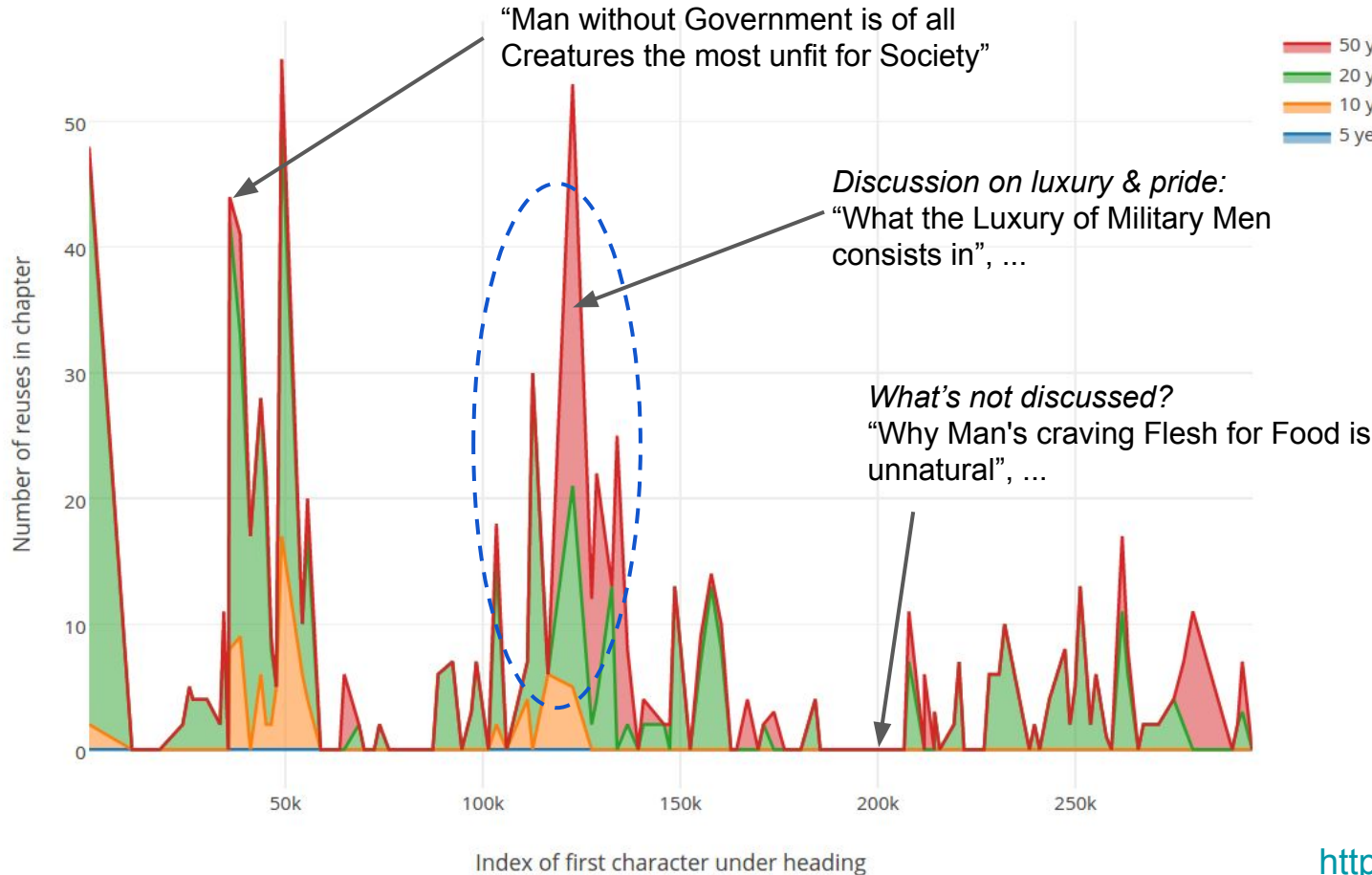
“Boscobel: or, the compleat history of His Sacred Majesty's most miraculous preservation after the battle of Worcester, which was fought Sept. 3, ...” (1660), Blount, Thomas, 1618-1679

“broad pieces to the king, judging they would be necessary to him in his present condition; for he durf carry no money about him in his mean garb and short cut hair, except about ten or twelve Lhillings in silver. Windham hereupon went to Lime, and spoke with Elef- don about hiring a lhip, which he undertook; but not till he was told, it was for His Ma- jefty's transportation. During the four or five dayv\” which the King this first time staid at Windham's, where he was was known by most of the family, e heard the bells ring, and feeing a company got to- gether in the church-yard, which wa4 very near the” [...]

A general history of England. Containing an Account of the first Inhabitants of the Country, and the Transactions in it, from the earliest ... (1754) Carte, Thomas, 1686-1754

" 300 broad pieces to the king, judging they would be necefhry for him in his prefelit condition; for he durst carry no money about him in his mean garb and his short cut hair, ex- cept about ten or twelve shillings in silver. Windham hereupon went to Lyme, and 1poke to Elefi'on about hiring a lhip, which he undertook: but not till he was told, it was for his majetfy's transportation. During the four or five days, which the king, this sirss time, flaid at Windam's (where he was known to most of the fa- mily) he heard the bells ring, and feeing a company got together in the church- yard, which was very near the” [...]

Fable (1714) reuses by chapter heading



CONTENTS.

<i>Man pretends to have for his Species,</i>	153
<i>Why Man’s craving Flesh for Food is unnatural,</i>	ibid.
<i>We ought not to judge of Nature’s design, but from the effects she shews,</i>	155
<i>Man never acknowledges Superiority without Power,</i>	156
<i>The feeling of Brutes proved from several concurring Symptoms,</i>	157
<i>A Definition of Frugality,</i>	158
<i>What the Lavishness or Frugality of Nations depend upon,</i>	159
<i>Maxims to make a People great and flourishing.</i>	162
<i>To make a Society good and honest,</i>	162
<i>The present Grandeur of the Dutch is not owing to the Virtue and Frugality of their Anc-stors,</i>	164
<i>The Hardships and Calamities they have suffered</i>	164
<i>Their natural Wants,</i>	165
<i>The Dutch not frugal by Principle</i>	168
<i>’Tis Policy and not Virtue that makes the Dutch encourage Frugality,</i>	169
<i>How they promote Lavishness when it suits with their Interest,</i>	170
<i>What</i>	

Better Science Through Better Data in the Humanities

Conclusion

What to learn from other fields?

- **Research support** databases
- **Collaboration & reuse** through open data
- Importance of **standards**



Nature Reviews Genetics **14**, 89-99 (February 2013) | doi:10.1038/nrg3394

Reuse of public genome-wide gene expression data

Johan Rung¹ & Alvis Brazma¹ [About the authors](#)

Data Science in Humanities

Potential

- New methods, classical research questions
- New scale of quantitative analysis
- Quality through collaboration

Pitfalls

- Data quality easily overlooked
- Existing tools drive research
- Expertise lacking

Science and hermeneutics

Tangible historical objects

Subjective historical experience



**Need for
new
methods
and clear
principles
for data
sharing!**

Humanities collaboration for better data

- Crowdsourcing experts
- Collaboration with different field of science, national libraries, infrastructures and projects
- Collaboration with companies that do digitization
- Interoperability & dealing with noise and bias

→ **We need right kind of infrastructures for specific purposes that enable collaboration between researchers, companies and libraries.**

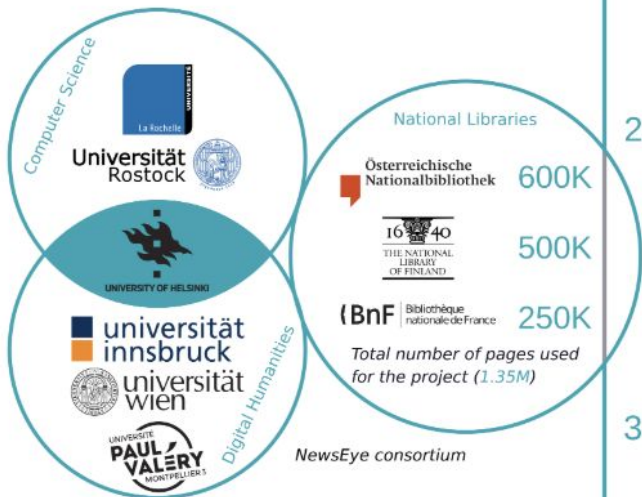


NewsEye is a research project advancing the state of the art and introducing new concepts, methods and tools for digital humanities by providing enhanced access to historical newspapers for a wide range of users. With the tools and methods created by NewsEye, crucial user groups will be able to investigate views and perspectives on historical events and development and, as a consequence, the project will change the way European digital heritage data is (re)searched, accessed, used and analysed.

Workflow

The core concept of NewsEye is a seamlessly integrated armoury of tools and methods that will improve the users' capability to access, analyse and use the content in the digital libraries of historical newspapers.

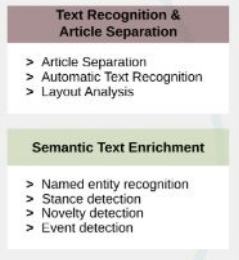
Four Case Studies: the inner test material used by the project's humanities research groups will be from European newspaper datasets from the three partner libraries focussing on the period 1850-1950.



4 NewsEye case study topics



2. Gender



3. Migration

From projects to large-scale sharing?

- Dariah-FI: Researcher-driven ecosystem of services for data-intensive social sciences and humanities (SSH) research

Themes	Modules	Focus areas
Data access and documentation	Social Sciences and Humanities Big Data	Digitised and born-digital data
Research methods and tool development	Analytica	Computational techniques and environments
Dissemination of best practices	Information Interaction	Researcher support

extras

Initial data

ECCO

ESTC

Evolving set of analysis and processing tools

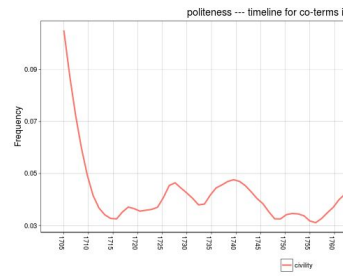
APIs

```
{
  "docFreq" : 6322
  "totalTermFreq" :
  "collocations" :
  "conversation" :
  "politeness" :
  "style" : 0.00
  "eloquence" :
  "obliging" : 0
  "genteel" : 0.1
```

Applications

le n-gram viewer

ms:



Scripts

```
21 return(publications_yearly)
22 }
23
24
25 get_relative_hits_yearly <- function(hits_sub
26 averages_yearly <- hits_subset_yearly
27 averages_yearly["total_titles"] <- hits_all
28 averages_yearly["frequency"] <- averages_yea
29 return(averages_yearly)
30
31
32
33 get_query_set_with_freqs <- function(query_fil
34 pubs_year
35 dataset)
36
37 filter_list <- list(id = get_query_ids_df_fr
38 idfiltered_dataset <- get_filtered_dataset(f
39 query_titles_yearly <- get_publications_year
40 get_relative_hits_yearly(query_titles_year
41 query_filename <- basename(query_file)
42
```

Etc...

...

Project goals

Research publications

Keywords: [history publishing](#), [short-title catalogue](#)

How to Cite: Lahti, L., Ilomäki, N. & Tolonen, M., (2015). A Quantitative History in the English Short-Title Catalogue (ESTC), 1470–1800. *L* 25(2), pp.87–116. DOI: <http://doi.org/10.18352/lq.10112>

596

62

1

Public tools, APIs, code



This repository

Search

Pull requests

Issues



rOpenGov / bibliographica

Code

Issues 2

Pull requests 0

Projects 0

Wiki

Tools for bibliographic data analysis — Frit

Refined data

publisher	latitude	longitude	publication_1
printed for J. Cooke; and sold by the booksellers of B...	51.50853	-0.125740	London
re-printed at the Old Post Office in Fishamble street	53.33306	-6.248890	Dublin
Sold at the Bible and Heart, in Cornhill	51.50853	-0.125740	London
printed for J. Bew	51.50853	-0.125740	London
N/A	37.56596	14.282404	N/A
Printed by J. Evans, Long-lane, West-smithfield, Lon...	51.50853	-0.125740	London
printed for Thomas Ewing, and William Smith	53.33306	-6.248890	Dublin
sold by J. Williams, and all the booksellers in town a...	51.50853	-0.125740	London

Challenges in reuse & sustainability

Data

- Often siloed for various different reasons
- Most humanities data not digitized currently (< 5% of relevant cultural heritage)
- data providers reluctant to share openly

Methods

- Often scattered
- Borrowed from other fields of science; do not necessarily suit the research questions and start guiding the work
- Reinventing the wheel

Expertise

- Often fragmented
- Cross-disciplinary collaboration often difficult because of the gaps in research cultures
- Ineffective collaboration between partners