# Linked Data Scopes
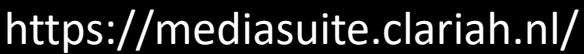
Victor de Boer     Ivette Bonestroo     Marijn Koolen     and   Rik Hoekstra

# (Digital) Humanities research has increasingly a data component



CLARIAH MEDIA SUITE

https://mediasuite.clariah.nl/

Humanities Data in R
Exploring Networks, Geospatial Data, Images, and Text

Springer

THE SHAPE OF DATA IN DIGITAL HUMANITIES
MODELING TEXTS AND TEXT-BASED RESOURCES
Edited by
Julia Flanders and Fotis Jannidis

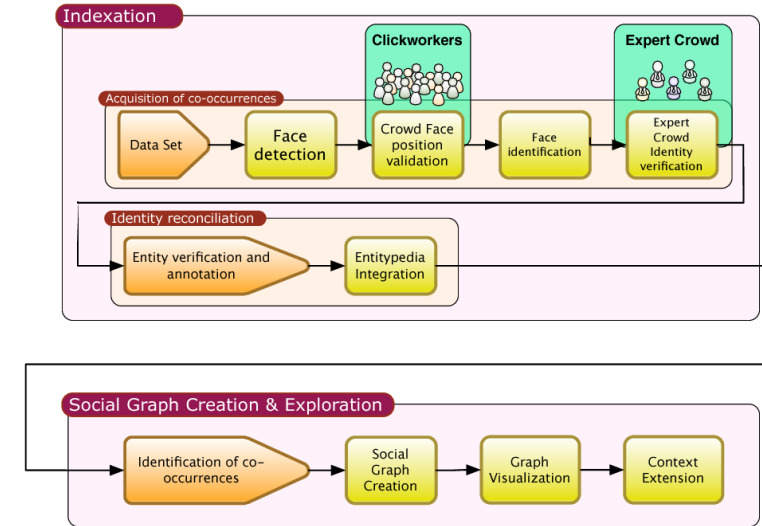# Variety of processing and data manipulation workflows/pipelines



tions.



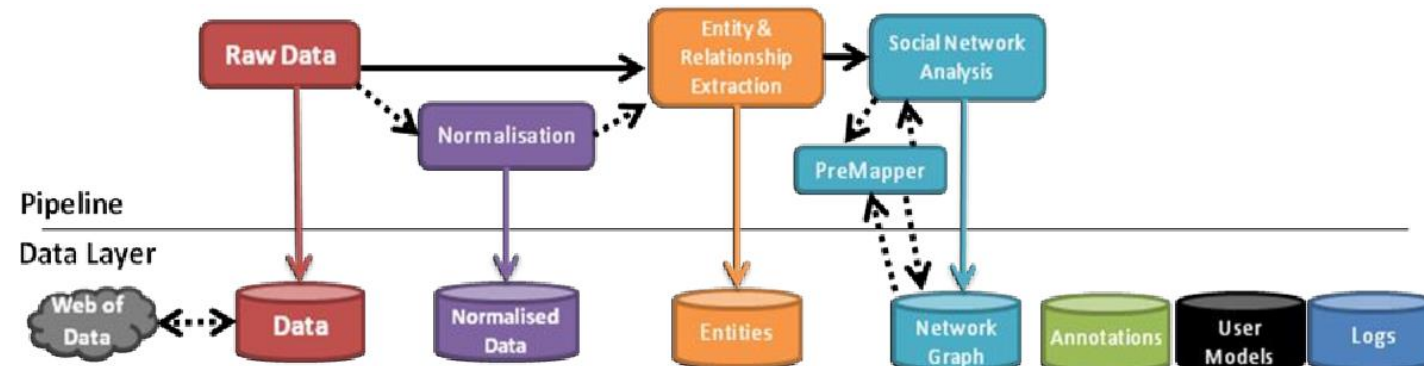**Fig. 3.** High-level view of the HoE indexation pipeline.



Figure 1. The CULTURA Pipeline and Data Layer

# Are these steps recognizable in resulting publications?

Interactions in this pipeline **change** the data and are essential in understanding any subsequent analysis. It makes them **part of historical research methodology**, but there is little consensus on how these steps can or should be performed.
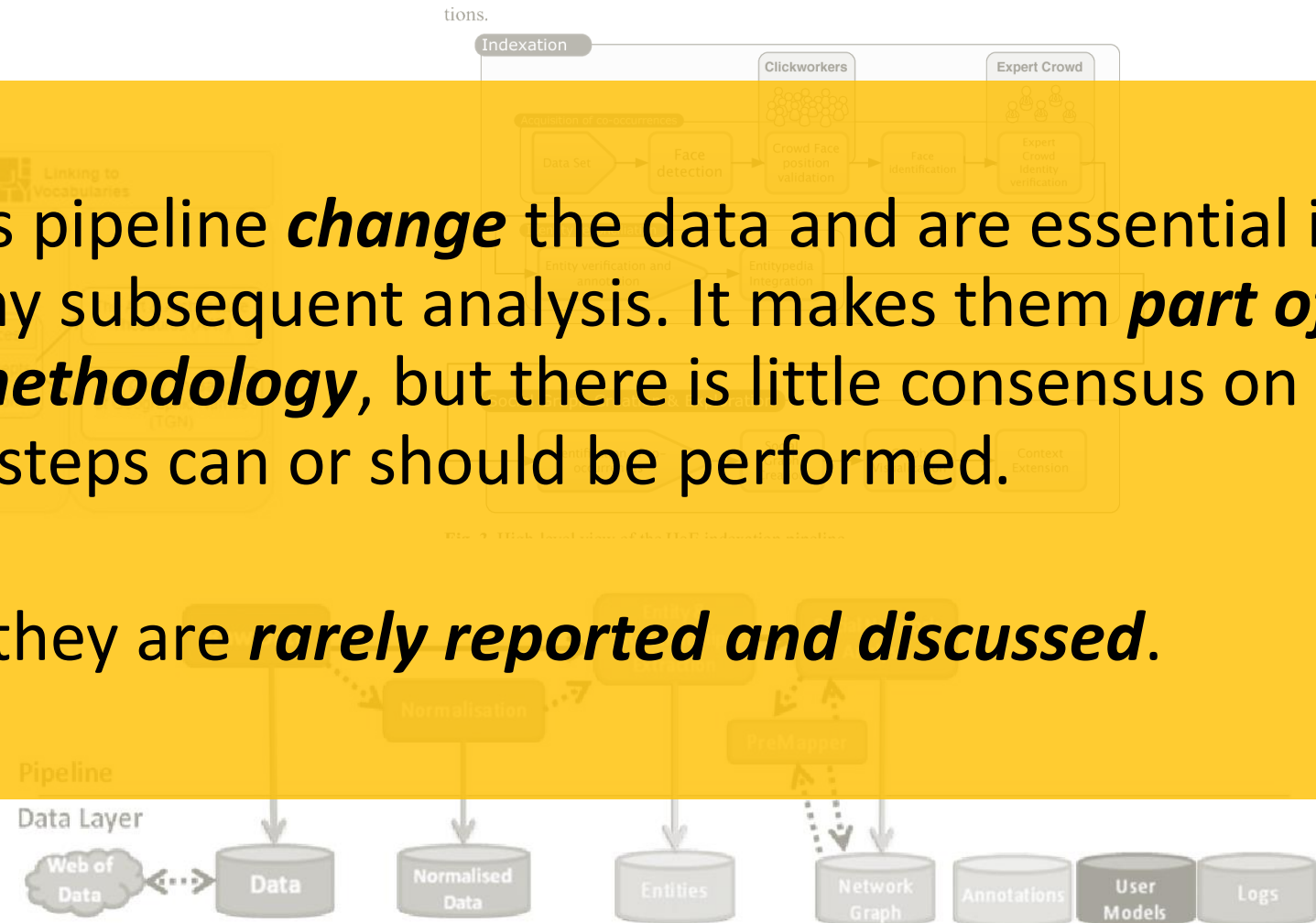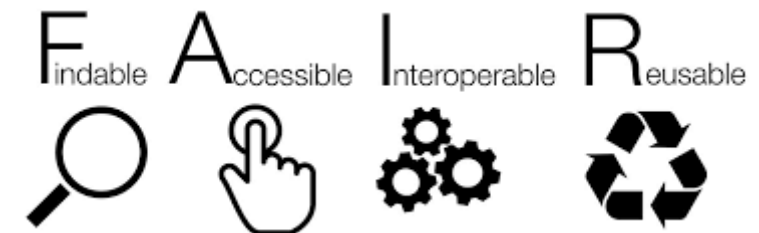
Moreover, they are **rarely reported and discussed**.

Rik Hoekstra & Marijn Koolen (2019) Data scopes for digital history research, I

# Goal: increase reusability and transparency in humanities research output

Through explicit descriptions of data transformations users of the datasets can assess "the context in which the data was created, its quality and validity, and the appropriate conditions for use." (Groth et al. 2012)

Describe and share
- Methodology
- Intermediate results
- Datasets
- Data enrichments

Findable   Accessible   Interoperable   Reusable

Groth, P., Y. Gil, J. Cheney, and S. Miles. 2012. Requirements for provenance on the web.

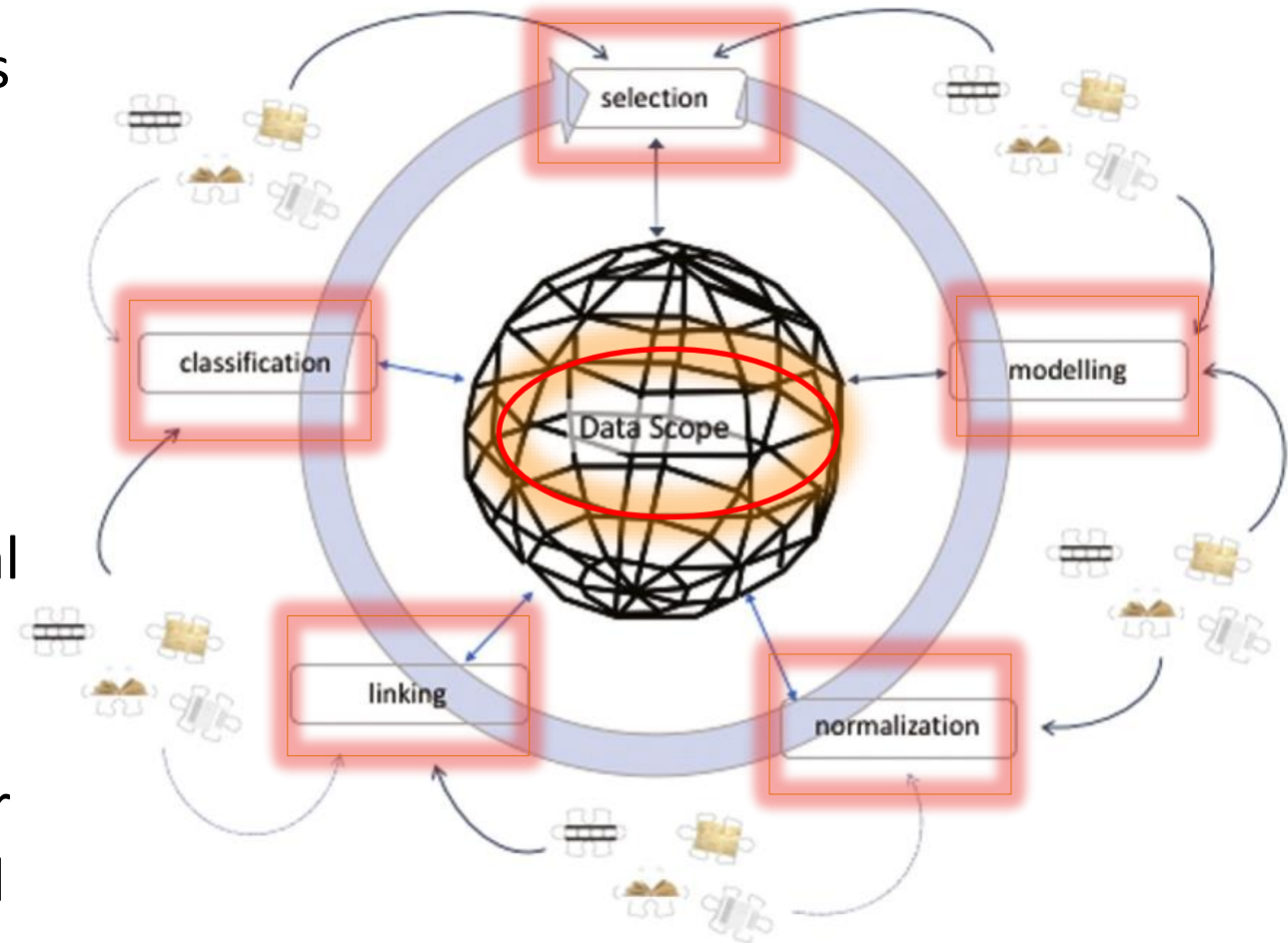# Preliminary: Data scopes for digital history research

Data scopes are a method to 'characterize the interaction between researchers and their data and the transformation of a cluster of data into a research instrument.' (Hoekstra & Koolen,2019)

Rik Hoekstra & Marijn Koolen (2019) Data scopes for digital history research, Historical Methods: A Journal of Quantitative and Interdisciplinary History, 52:2, 79-94, DOI: 10.1080/01615440.2018.1484676
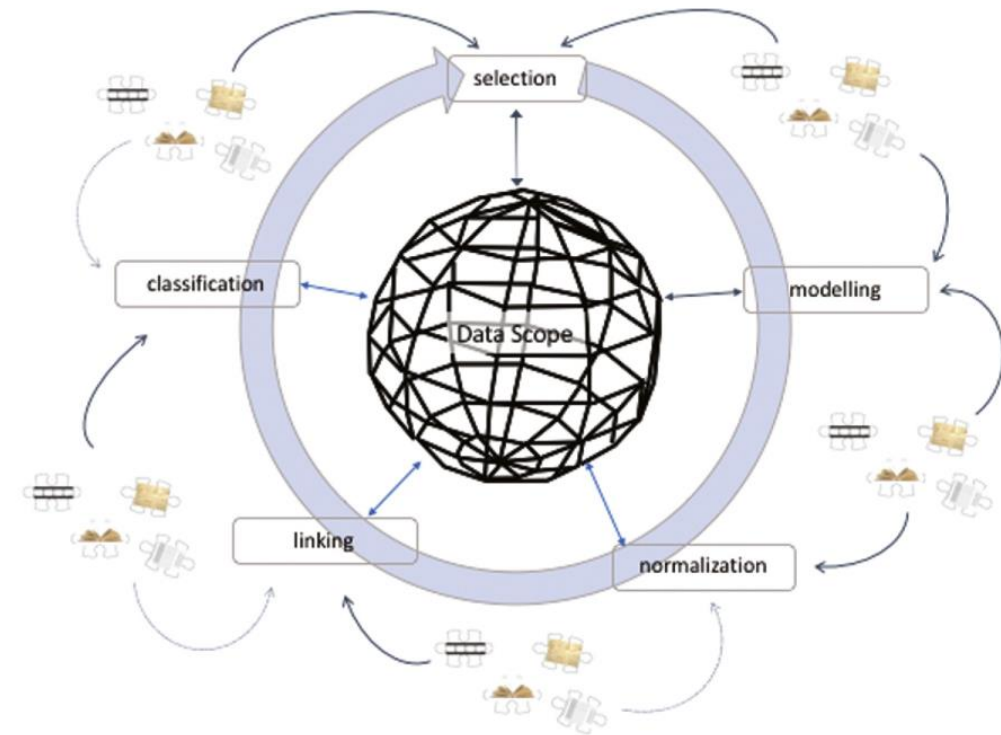
# Data scope denotes 5 data manipulation steps

1. **Selection**: which data and sources are selected? (corpus forming)

2. **Modelling**: how are the relevant elements in sources represented? (using implicit or explicit models)

3. **Normalization**: how are surface forms mapped to a normalized form? (e.g. "Firstname, Lastname")

4. **Linking**: what explicit internal and external connections are established? (Includes deduplication, NE resolution etc.)

5. **Classification**: how are objects grouped or categorized? (includes internal or external schemes or theories)
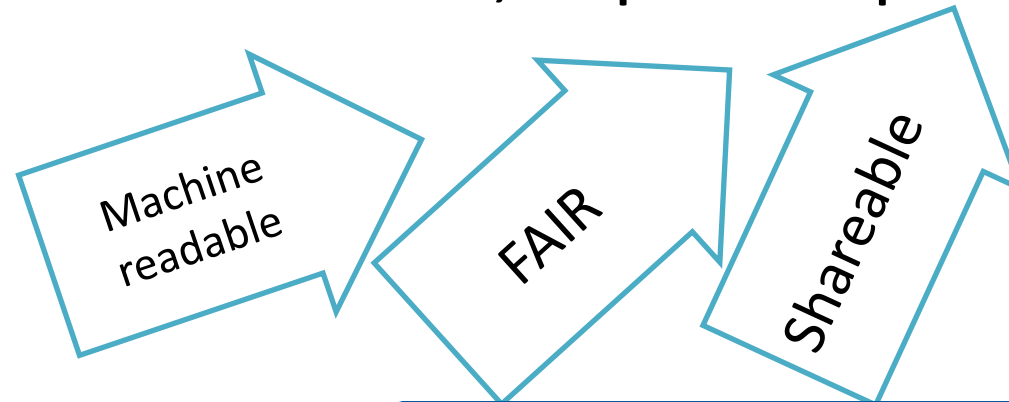


Rik Hoekstra & Marijn Koolen (2019) Data scopes for digital history research,

# This paper: Contribution 1

From qualitative conceptualization …



…to reusable, explicit representation

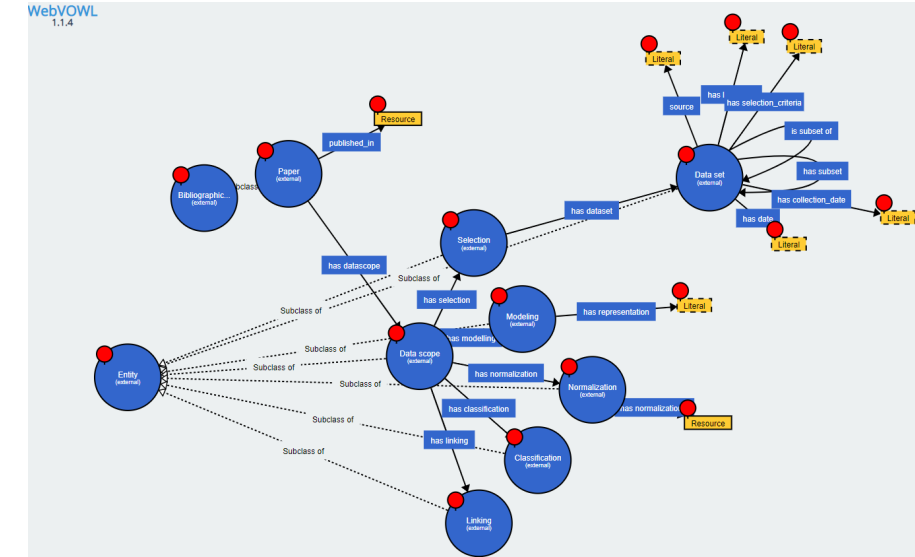# The Data scopes ontology



<http://w3id.org/datascope#>
permanent identifiers
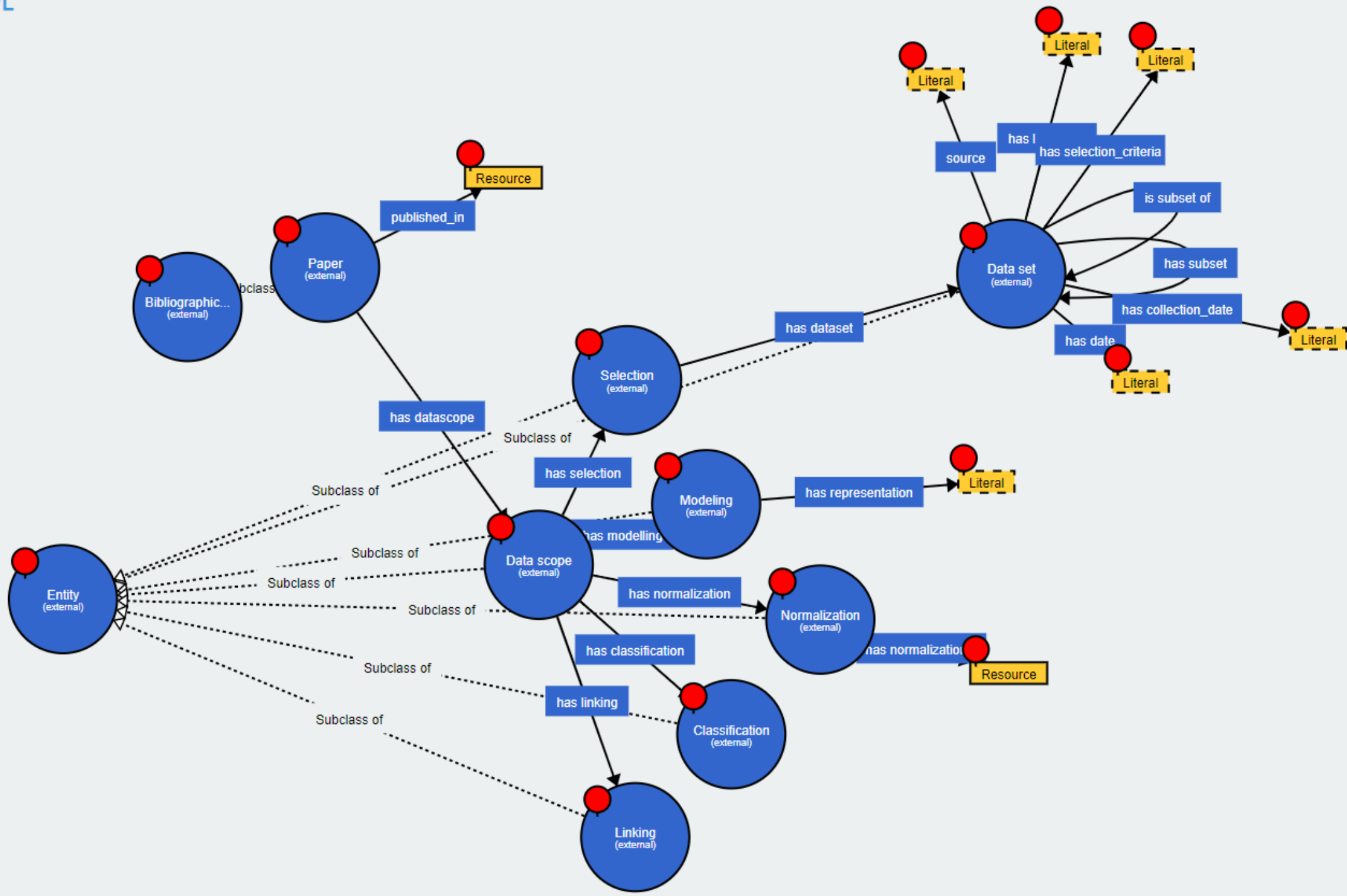
Central class is **dsont:DataScope**

Classes for the five main manipulation activities

Additional classes and properties for Datasets,
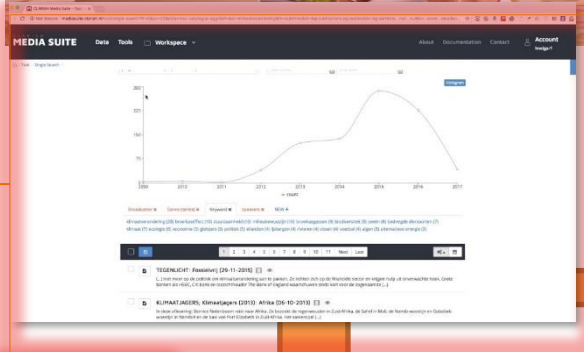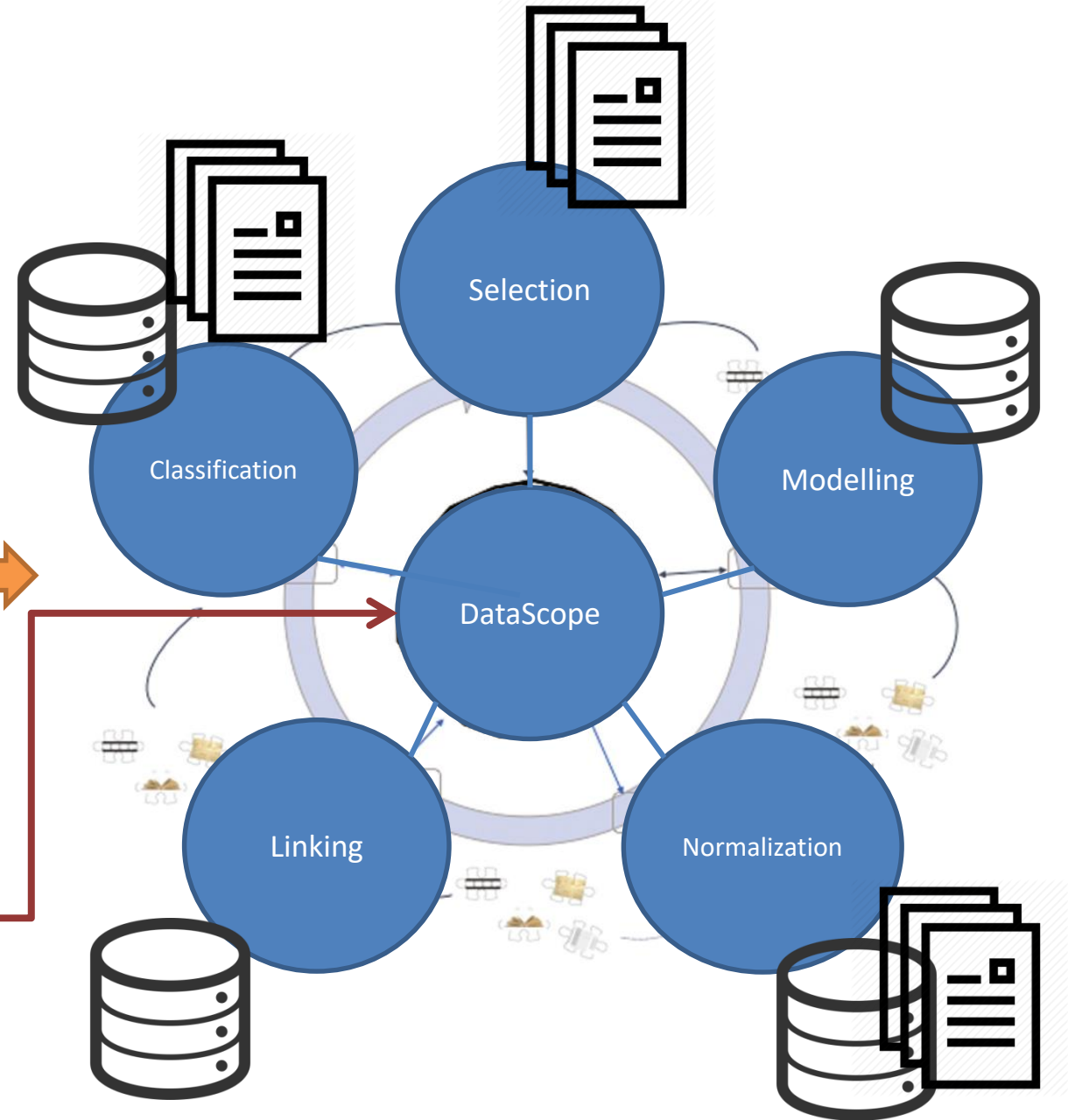selection steps, linksets, Concept Schemes etc.

Aligned with Dcterms, PROV

http://biktorrr.github.io/datascope/
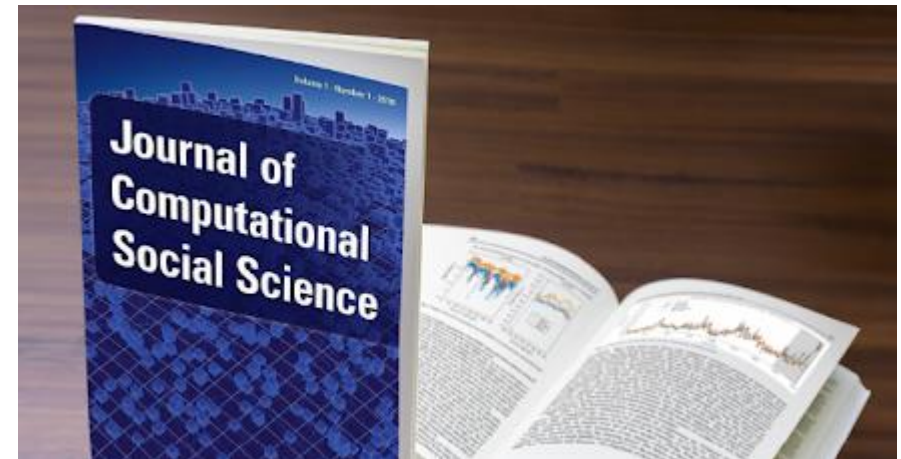
# Use of The Data scopes ontology

# This paper: Contribution 2

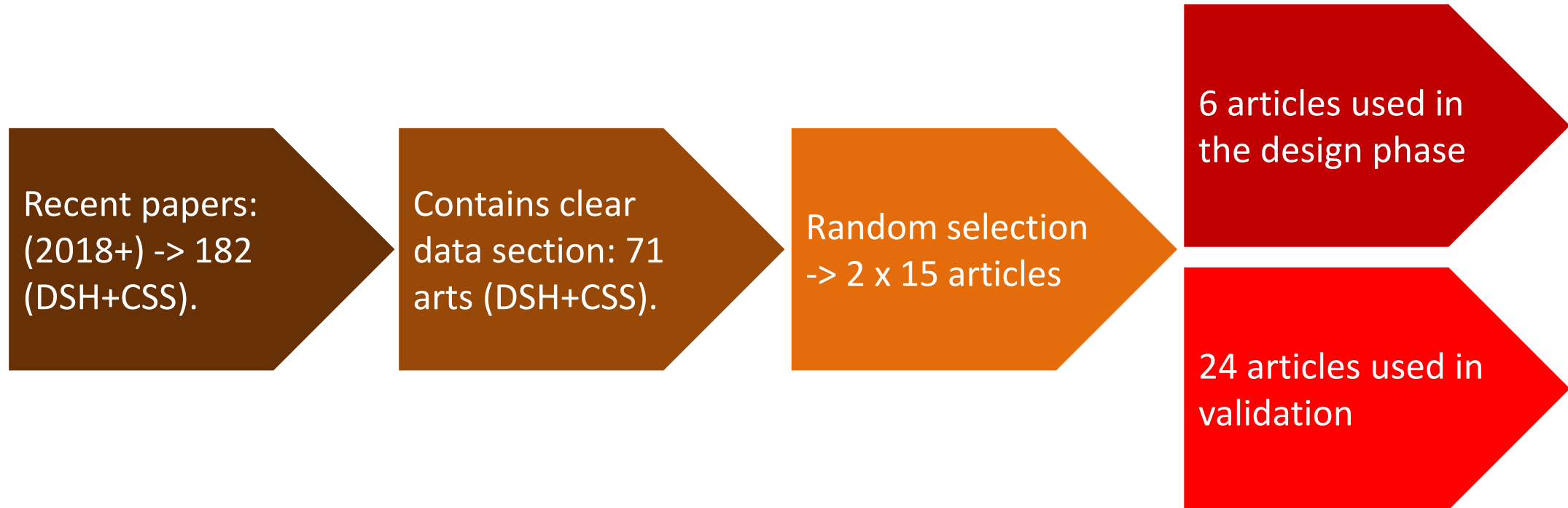Initial validation of the model in two research domains: Digital Humanities and Computational Social Sciences



Digital Scholarship for the Humanities (DSH)



Computational Social Science (CSS)

# Method

Recent papers: (2018+) -> 182 (DSH+CSS).

Contains clear data section: 71 arts (DSH+CSS).

Random selection -> 2 x 15 articles

6 articles used in the design phase

24 articles used in validation
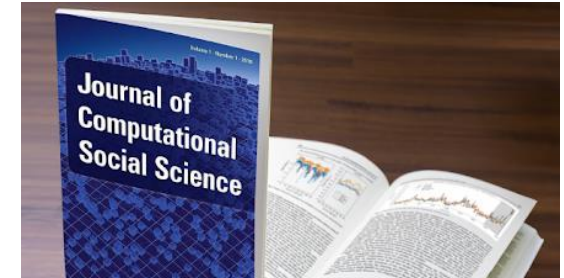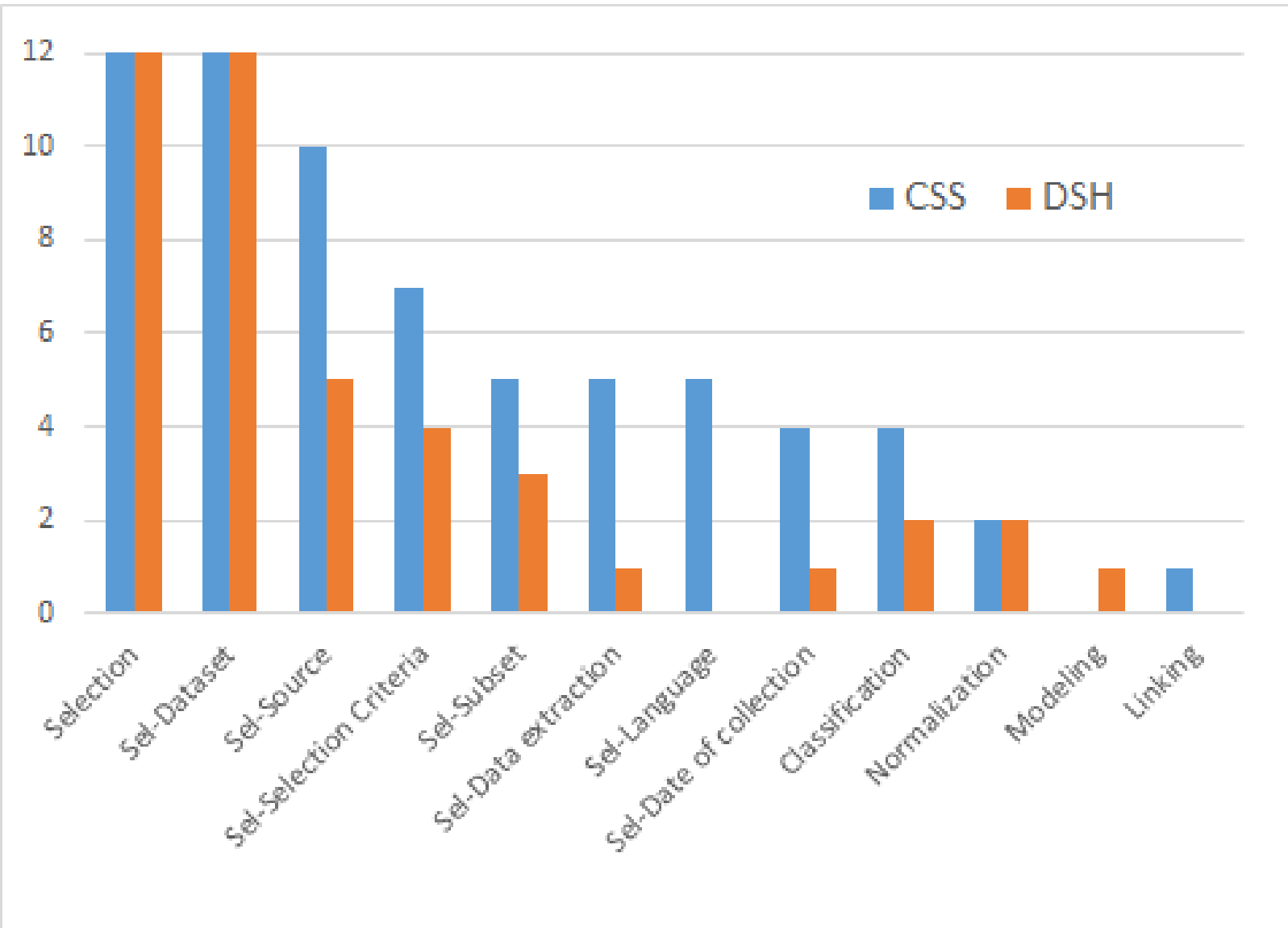
Set of coding guidelines
Two independent coders

# Results



Each element appears at least once

But most papers do not have complete descriptions

Selection most used (as hypothesized by Hoekstra & Koolen)

Can be used to compare fields

CSS papers have more complete descriptions
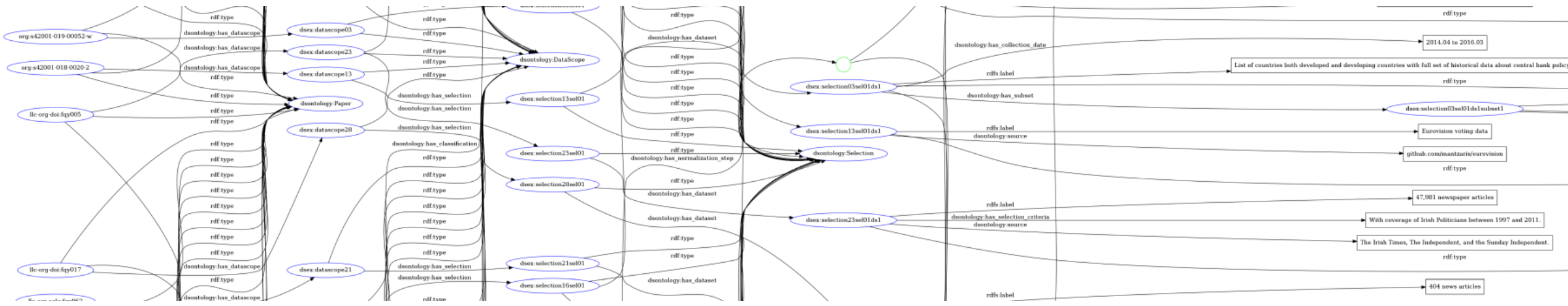
# Results represented as RDF data

511 RDF triples describing the datascopes of
24 papers, relating back to the original papers

SPARQLable at

https://semanticweb.cs.vu.nl/test/query

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX dsont: <https://w3id.org/datascope#>
SELECT ?s ?norm (COUNT(?ds) as ?dscount) WHERE {
    ?s rdf:type dsont:DataScope .
    ?s dsont:has_Selection ?sel .
    ?sel dsont:has_Dataset ?ds .
    ?s dsont:has_Normalization ?norm }
GROUP BY ?s ?norm
```

*"Show me the # of datasets for which a normalization step is registered "*

# Next steps

Further refinement and validation

Integrate these with existing tools
   CLARIAH mediasuite
   OpenRefine...

Allow direct and FAIR publishing
   Using Nanopublications/nanobench

Linking to other ontologies
   Workflow ontologies...

# Thank you!

http://biktorrr.github.io/datascope/

http://mediasuite.clariah.nl

@victordeboer

v.de.boer@vu.nl