

Data Science Tools zur Analyse des globalen Repository-Netzwerks

Dirk Pieper – Friedrich Summann
Bielefeld University Library

Kolloquium Wissensinfrastruktur, 31.5.2019

Abstract

- Der Vortrag skizziert technischen Möglichkeiten, das globale Publikationsnetzwerk und seine Entwicklung zu analysieren und zu beschreiben.
Seit 2004 fallen zahlreiche Metadaten im BASE-Umfeld an: insbesondere Publikationsmetadaten, Metadaten zu den Quellen, Indexdaten und damit verbunden Informationen, die zu den Quellen die Verteilung bei Publikationsstyp, Publikationsdatum, Sprache des Dokuments, OA Status und Lizenzinformationen liefern.
Skizziert wird das Konzept eines Monitoring-Systems inklusive des Datenflusses und der zu implementierenden Schnittstelle.
An einigen konkreten Beispielen wird prototypisch gezeigt, wie mit geeigneten Tools eine visuelle Umsetzung implementiert wurde und erweitert werden kann.

Overview

- Introduction
- The Data Basis and its Pre-Processing
- The API Layer
- Concrete Examples
- Summary

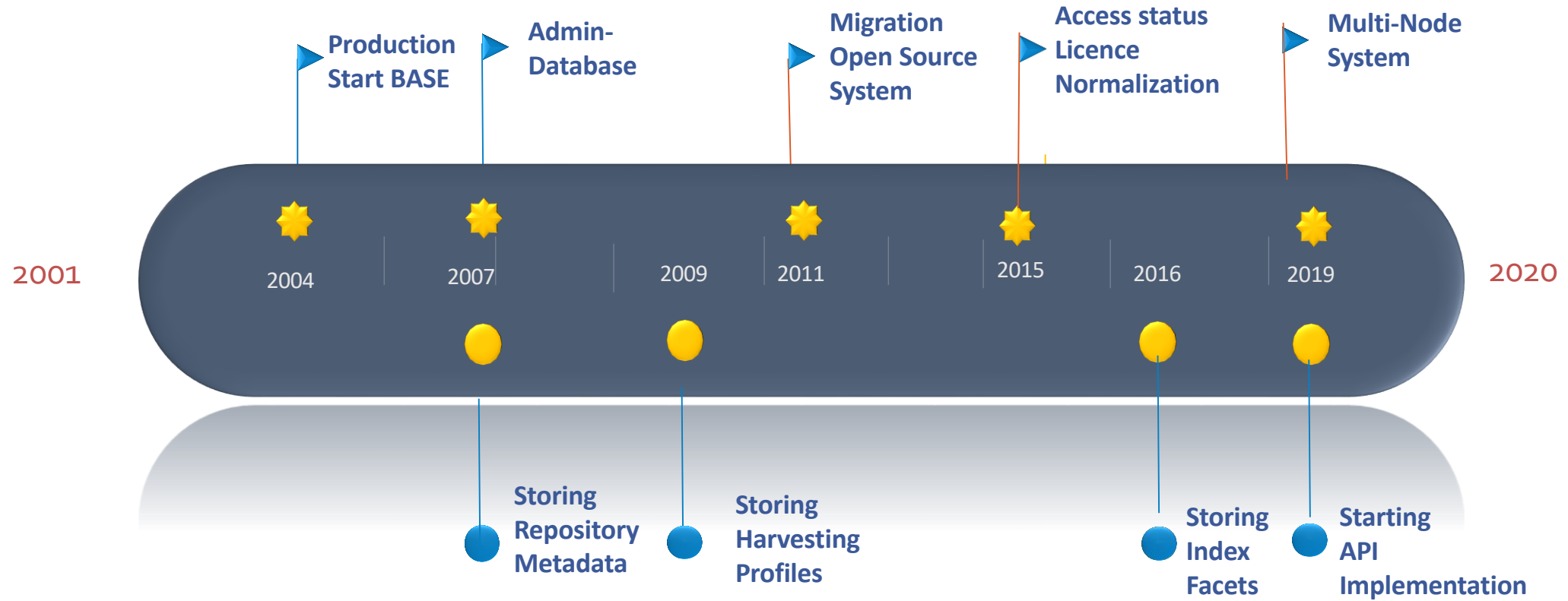


- 7115 Repositories included (mostly via OAI-PMH, some via Crossref, few via proprietary ...)
- From 124 Countries world-wide
- Ca. 146 Mill. Documents/Objects
- Ca. 70 % Open Accessible
- Internal Scheme: extended Dublin Core Format
- Inst. Repositories, Journal Platforms, Research Data, Digital Collections, Multimedia

Overview

- Introduction
- **The Data Basis and its Pre-Processing**
- The API Layer
- Concrete Examples
- Summary

- Repository Metadata –
 - Common resource description (country, platform, repository type, OA status, geocodes, name)
- Harvesting Metadata
 - Number of Documents, Number of OA Documents
 - Technical figures (Metadata Formats, Deleting Strategy)
- Index Data
 - - Number of repositories, number of documents
 - - the Index Facets (date, language, publication type, Access status, licences)



Einstellungsdatum*:	20.03.2006	Datum auswählen
Bearbeiter:	Friedrich Summann	▼
Quellenname:	PUB - Publications at Bielefeld University	
Harvestername (OAI)*:	UBBielefeld-PUB	
System:	LibreCat	▼
	Sonstiges (hat Vorrang!):	
Klassifikation:	DDC	
Herkunftsland (Domain / TLD):	de	▼
Bundesstaat (z.B. denw):	denw	
DINI-Zertifikat:	2016	▼
URL (Applikation-Startseite):	http://pub.uni-bielefeld.de/	
OAI Basic URL:	http://pub.uni-bielefeld.de/oai	
Repository Id:	od_2294	
OA:	Unbekannt	▼
OASet:	open_access	
Breitengrad (Finder):	52.037600	
Längengrad:	8.493600	
Adresse:	----- From OpenDOAR: Universität Bielefeld 25, D-33615 Bielefeld Germany -----	
Kontakt:	----- From OpenDOAR: Friedrich Summann Administrator friedrich.summann@uni-bielefeld.de -----	
Typ:	Hochschulpublikationen	▼



Basic search Advance

Entire Document

country:eu

Verbatim search Additional word frequency

Boost open access documents

5,726,792 hits in 146,724,642 documents

Refine Search Result

Author

Subject

Dewey Decimal Classification

Year of Publication

Content Provider

Language

Document Type

Access

Terms of Re-use

Document Type

(2,784,193) Article contribution

(712,020) Unknown

(675,438) Book

(489,873) Still image

(211,473) Text

(167,435) Doctoral and postdoctoral thesis

(164,609) Bachelor thesis

(137,507) Dataset

(94,286) Conference object

(75,503) Book part

(66,202) Master thesis

(49,808) Report

(25,904) Review

(20,340) Other non-article

(14,970) Manuscript

(14,055) Course material

(11,349) Lecture

(11,066) Thesis

(6,880) Audio

Terms of Re-use

(936,467) PDM

(760,128) CC-BY-NC-ND

(395,597) CC-BY

(151,575) CC-BY-NC-SA

(97,643) CC0

(34,933) CC-BY-NC

(5,849) CC-BY-ND

(5,387) CC-BY-SA

(53) GPL

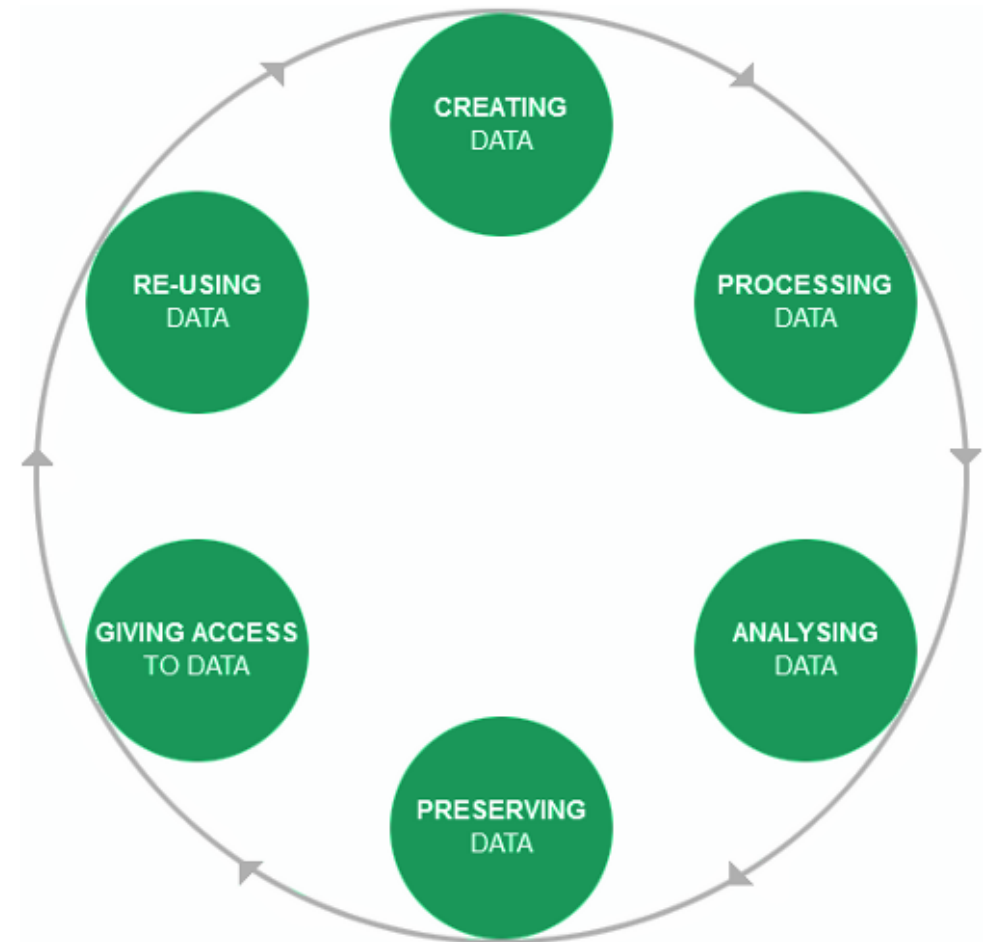
(22) MIT

(1) BSD

Bielefeld Center for Data Science

Konzeptioneller Rahmen	Einführung in Daten
Datensammlung	Datenschließung und -sammlung
	Evaluierung und Sicherstellung der Qualität von Datenquellen
Datenmanagement	Datenorganisation
	Datenmanipulation
	Datenkonvertierung
	Erzeugung und Verwendung von Metadaten
	Datenheilung, -sicherheit und -wiederverwendung
	Datenaufbewahrung
Datenevaluation	Datenwerkzeuge
	Grundlegende Datenanalyse
	Dateninterpretation (Datenverständnis)
	Nutzung von Daten zur Identifizierung von Problemen
	Datenvisualisierung
	Datenpräsentation (verbal)
	Datengetriebene Entscheidungsfindung
Datenanwendung	Kritisches Denken
	Datenkultur
	Datenethik
	Datenzitierung
	Datenteilung
	Evaluieren von Entscheidungen basierend auf Daten

UK Data Archive research data lifecycle model



Collection Data

Data Management

Data Provision

Data Re-Use

BASE-Metadatas-Store

Protocol Data

BASE-Admin-Data

BASE-Indexdata

Normalized/Processed
Data Store

API

Client 1

Client 2

Client x

Client -

External data (to compute correlations)

- Number of Scientists, Number of Students per country
- Size of Population per Country
- Publication Output per Country
- Funding Budgets
- Bibliometric figures

Level of Aggregation

- Repository
- *Institutional (more than 1 resource per institution)*
- Country Aggregation
- Repository Type Aggregation
- Continent Aggregation
- Global Aggregation

Overview

- Introduction
- The Data Basis and its Pre-Processing
- **The API Layer**
- Concrete Examples
- Summary

API (to be Developed)

- Restful API (via http)
- Query Language Definition
- Format Definition Responses
- Export Formats (xml,json,csv)
- Documentation

Flexible search query-like with Multi-Level - Support

Query Examples:

- List number of publications with publication type ‚thesis‘ per country
- List percentage of OA documents per repository/country/continent
- Showing the Increase of Research Data Objects per country per year
- List languages of documents with type books

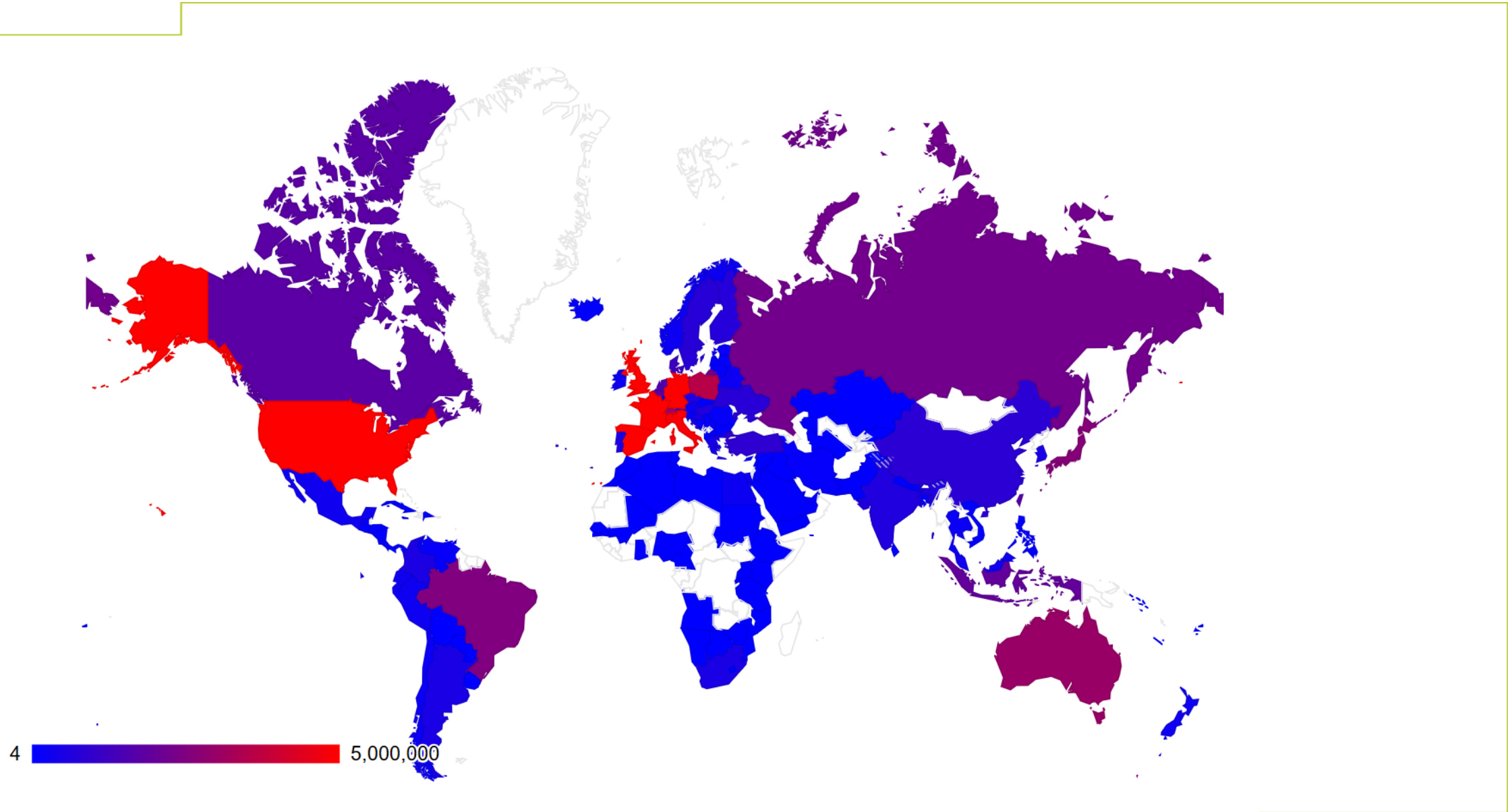
Timeline support

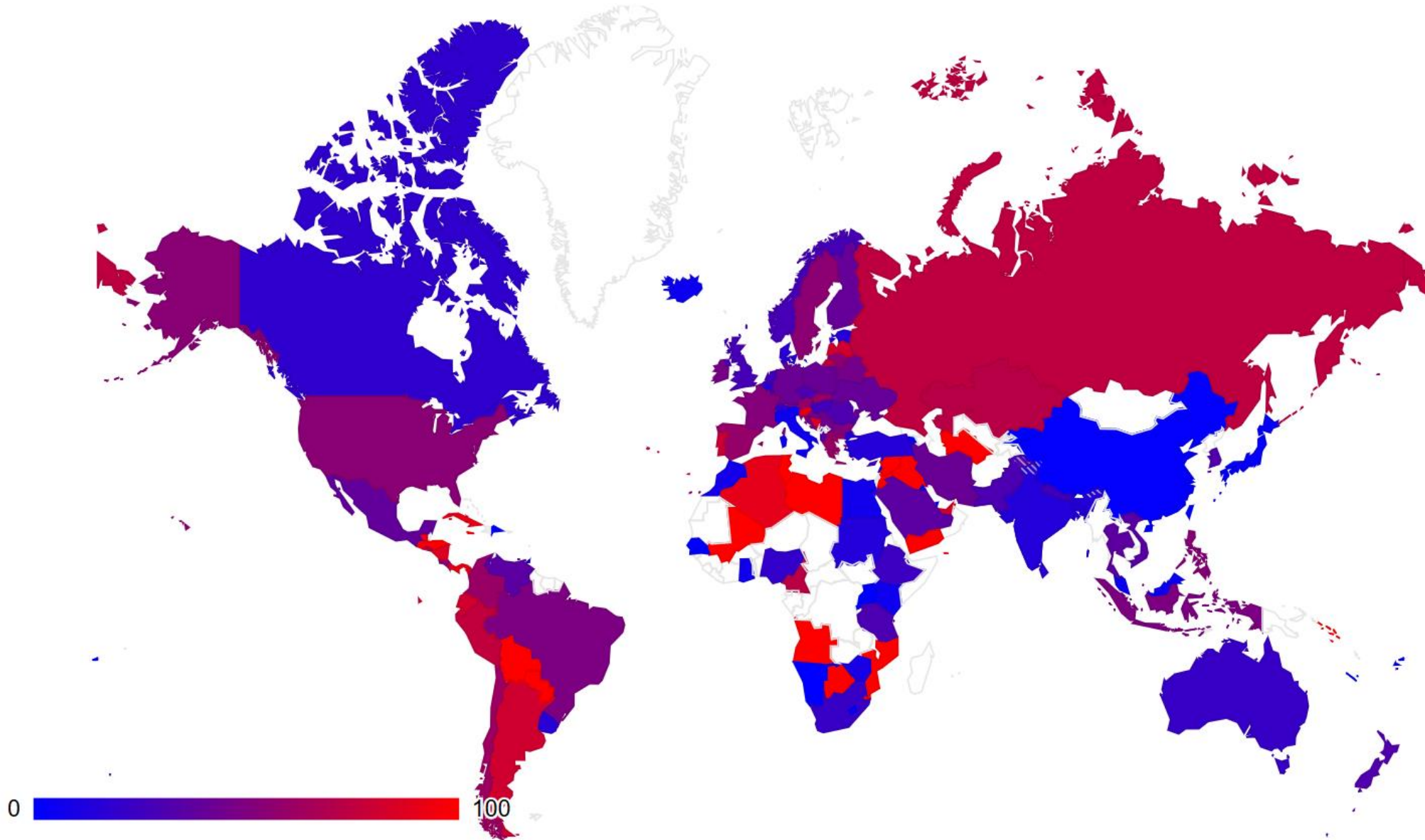
Overview

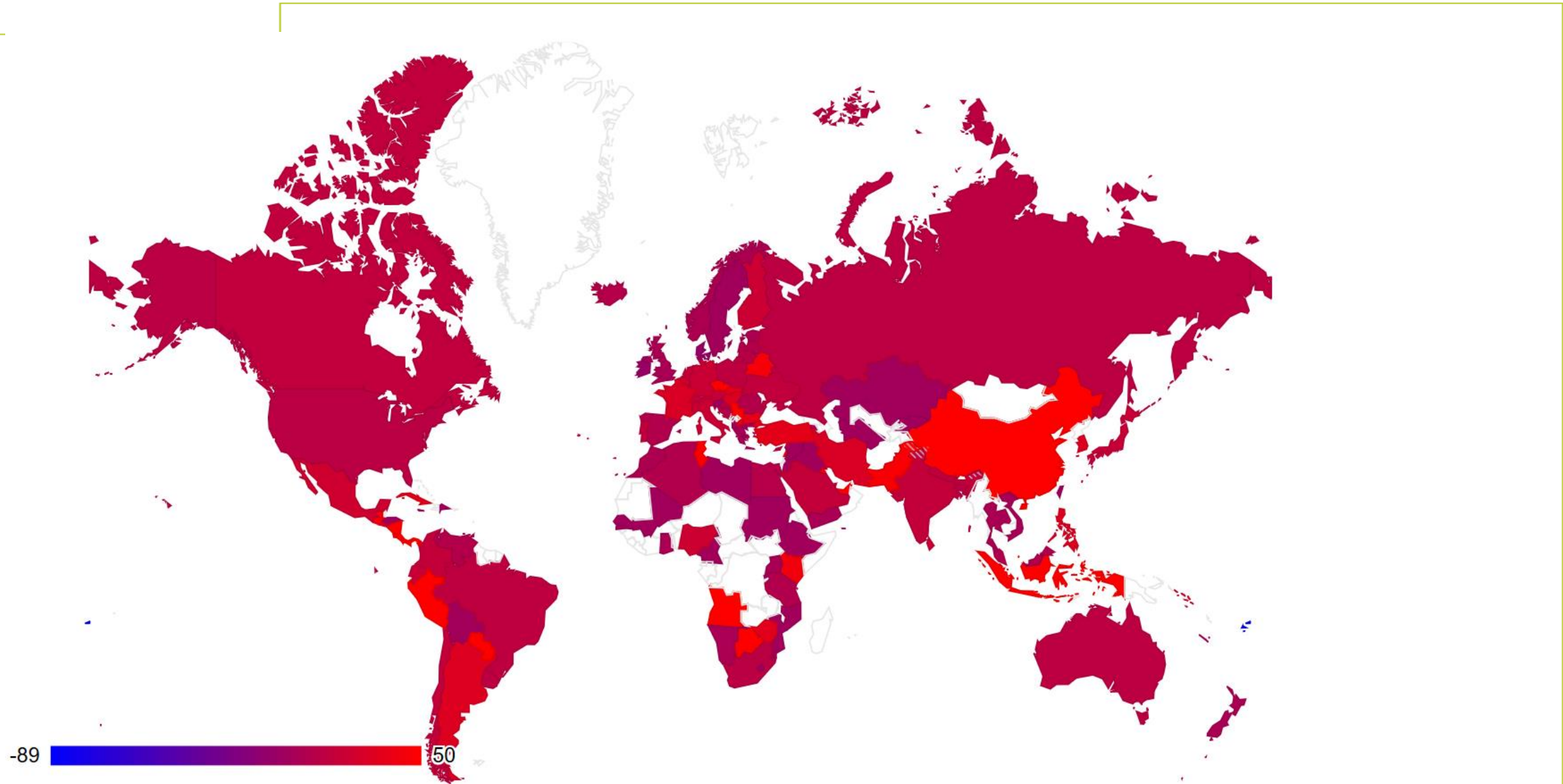
- Introduction
- The Data Basis and its Pre-Processing
- The API Layer
- **Concrete Examples**
- Summary

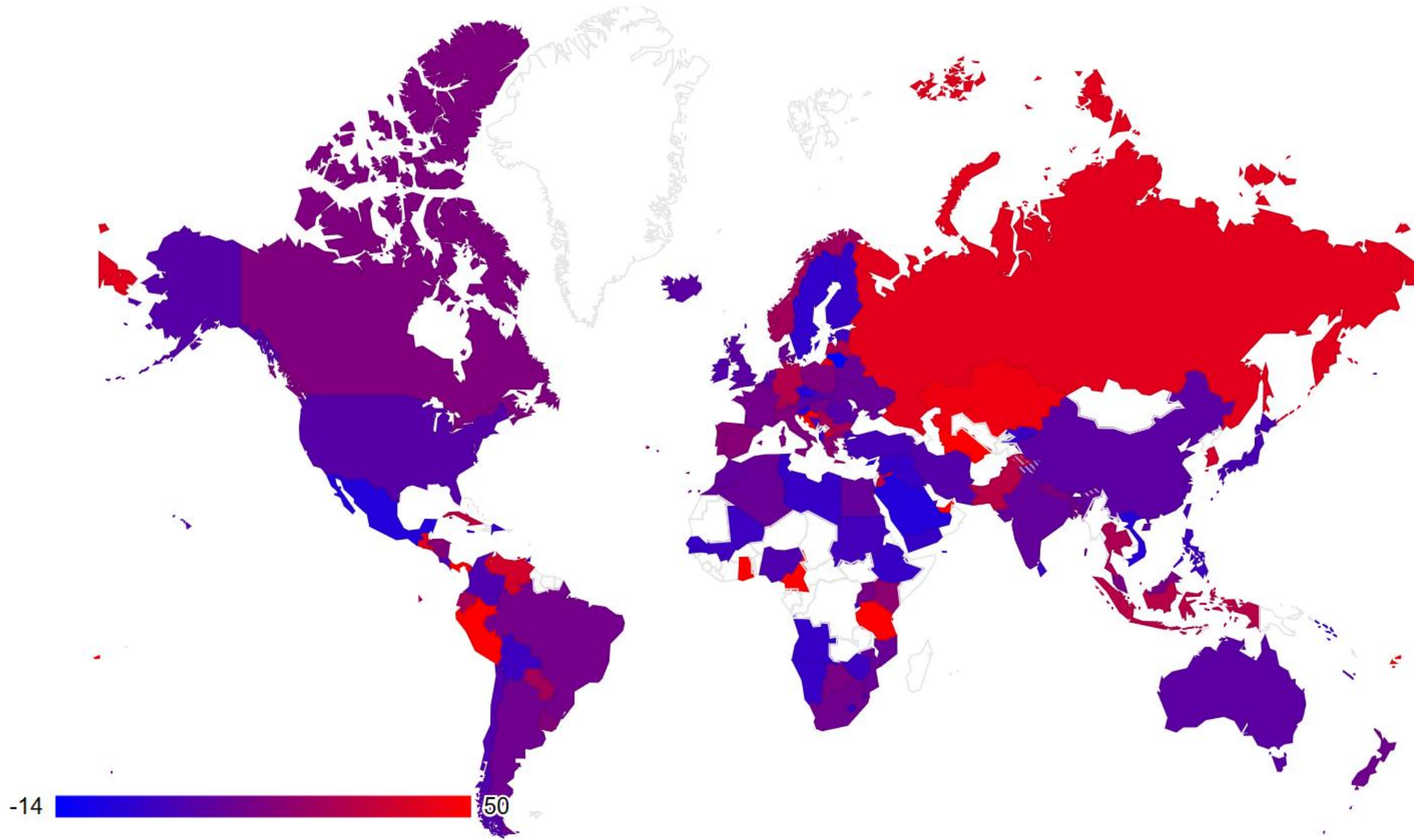
Existing Tools for Visualisation

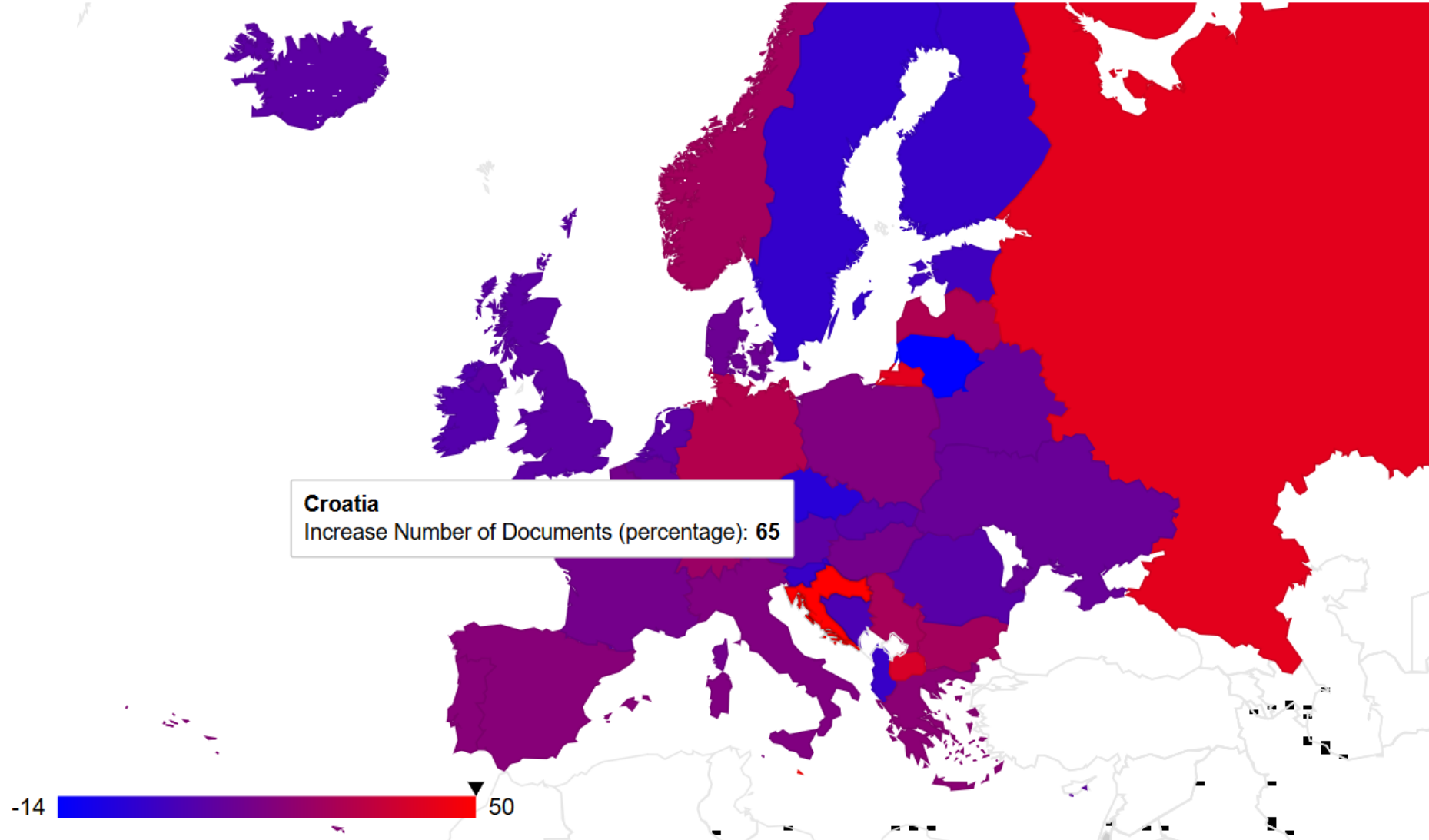
- World Maps
 - Google Charts
- Repository/Country Evaluation
 - D3js (Javascript Data Visualisation Framework)

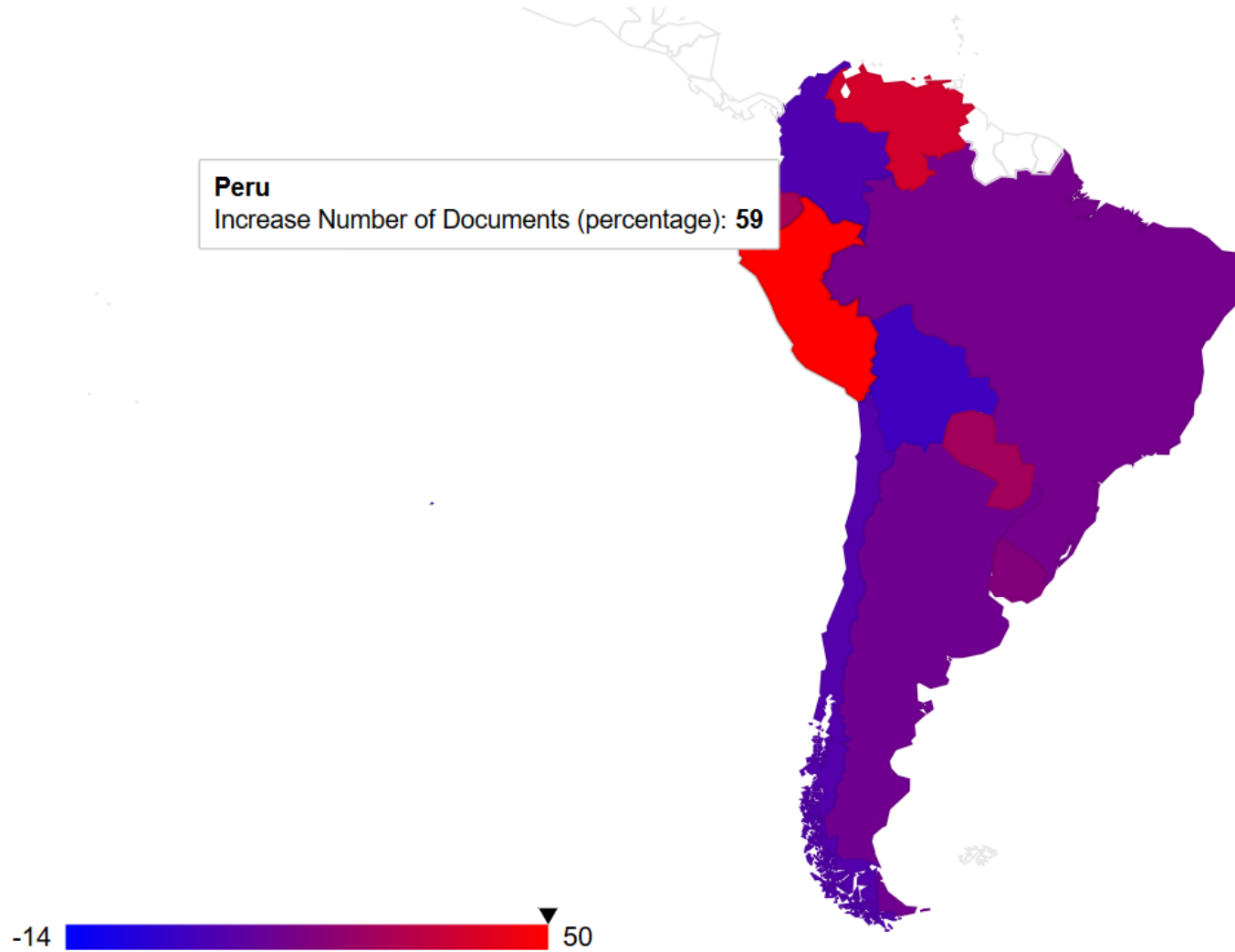




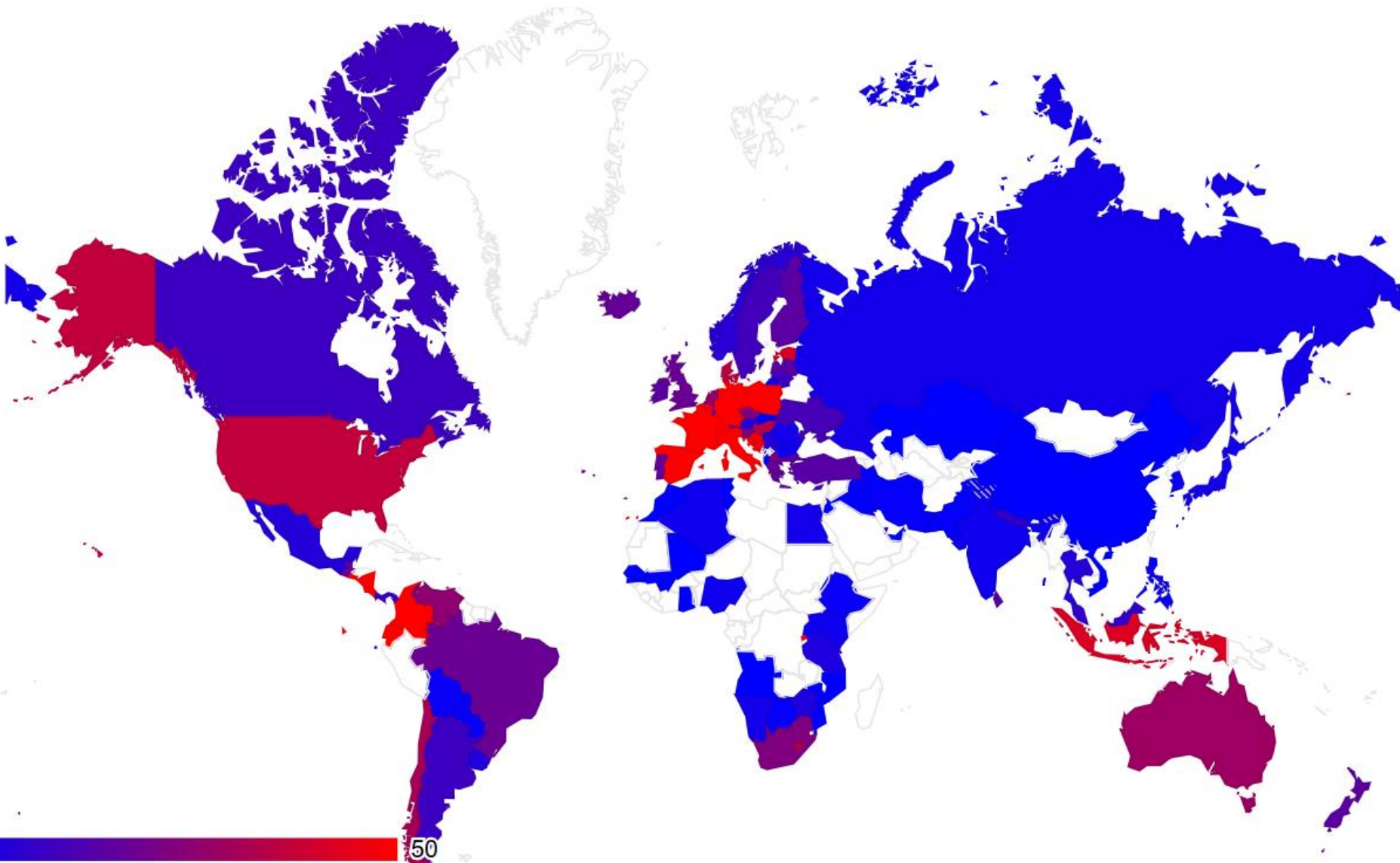








Number of Documents per Researcher (May 2019)



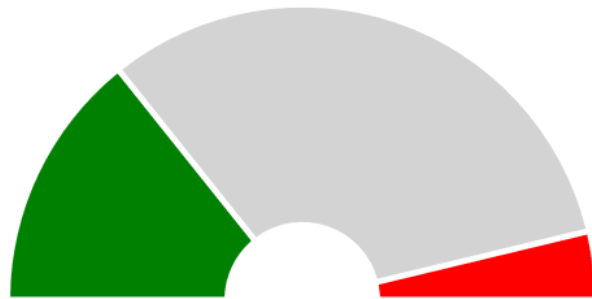
BASE Statistics: - Content Sources



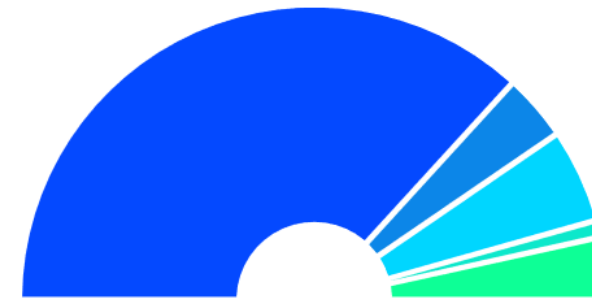
United Kingdom

5284718 Documents from 338 Repositories

4 % of Documents world-wide, 5 % of Repositories world-wide



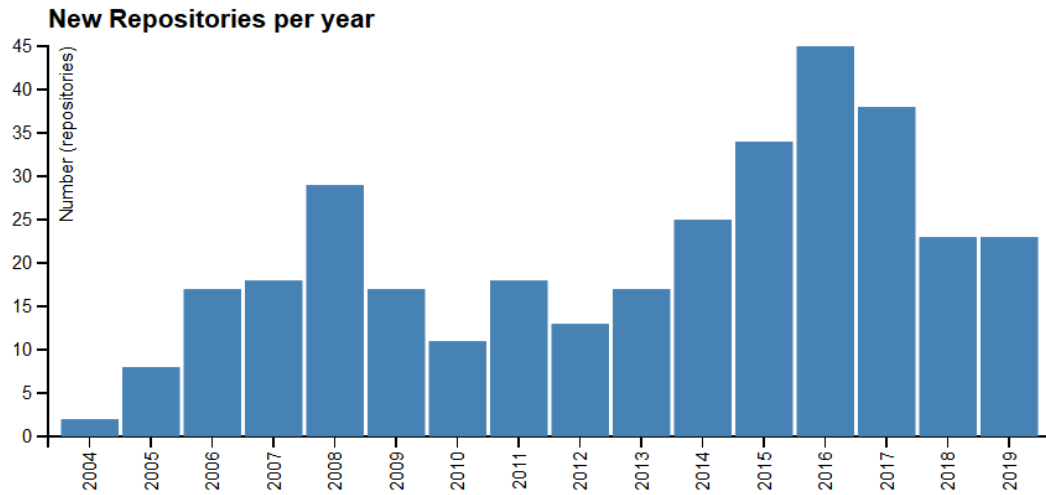
OA Status Distribution (Number of Documents)



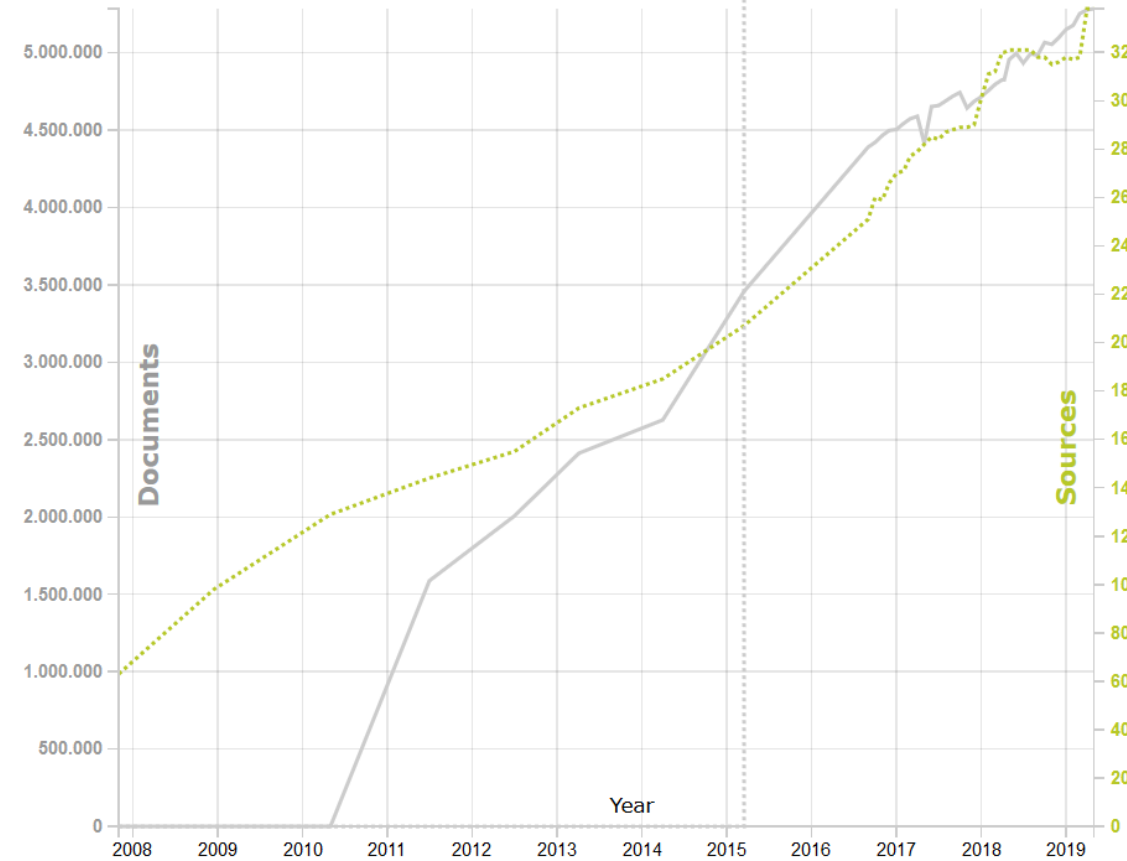
Repository Type Distribution (Number Documents)



Repository Development



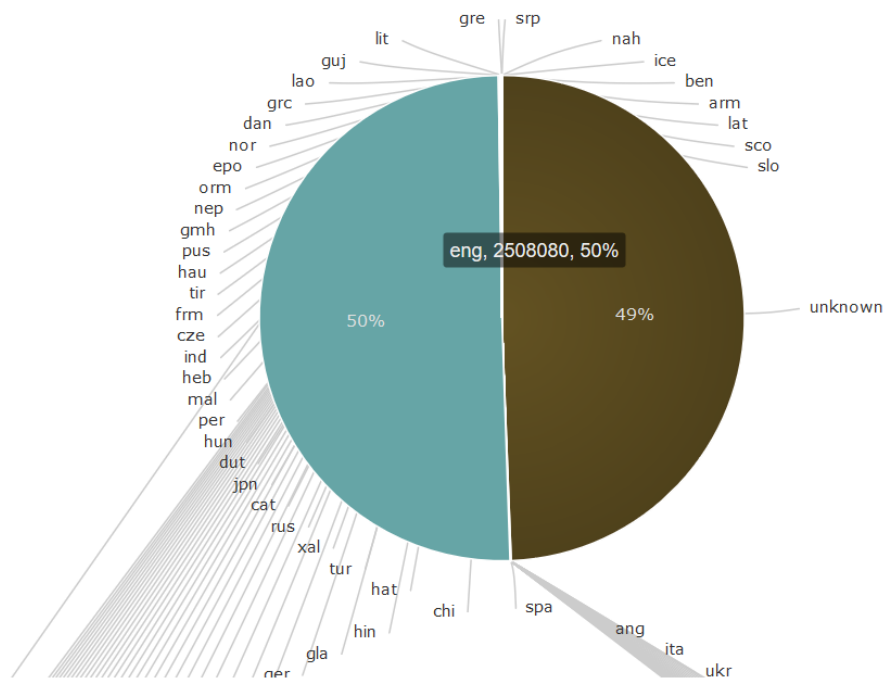
Number of indexed Documente und Sources in BASE for United Kingdom



Repository Contents Analysis

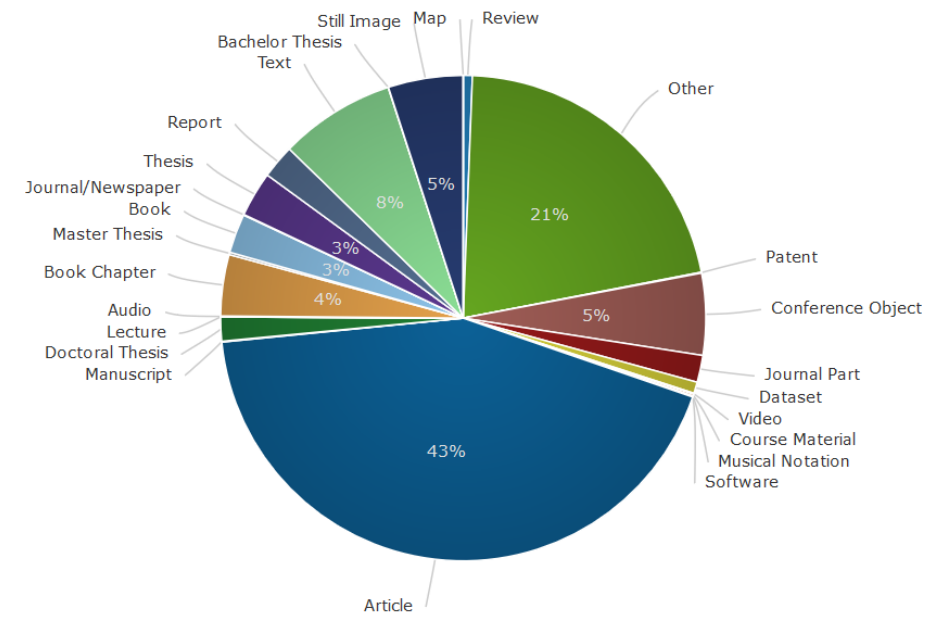
Language Distribution

(number of documents)



Publication Type Distribution

(number of documents)



- Presentation Slides
- Data Provision: Thesis OA Development
- DINI List of Repositories
- Bibliometrics
- ORCID2-DE

Danke!