# Connecting Collections: Using Linked Data in Libraries, Museums and Archives
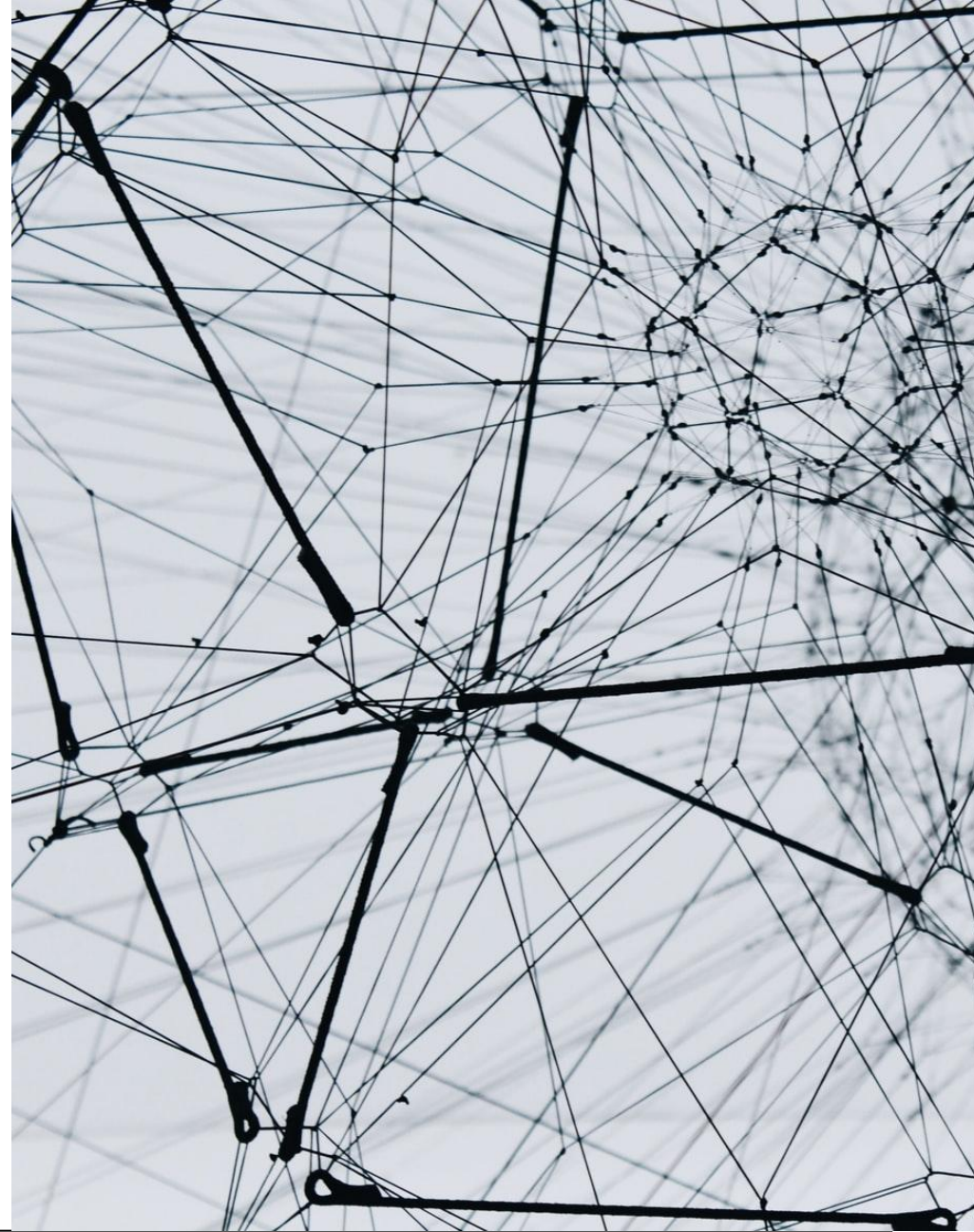
**Unlock the Libraries : VÖB event, 24 May 2022**
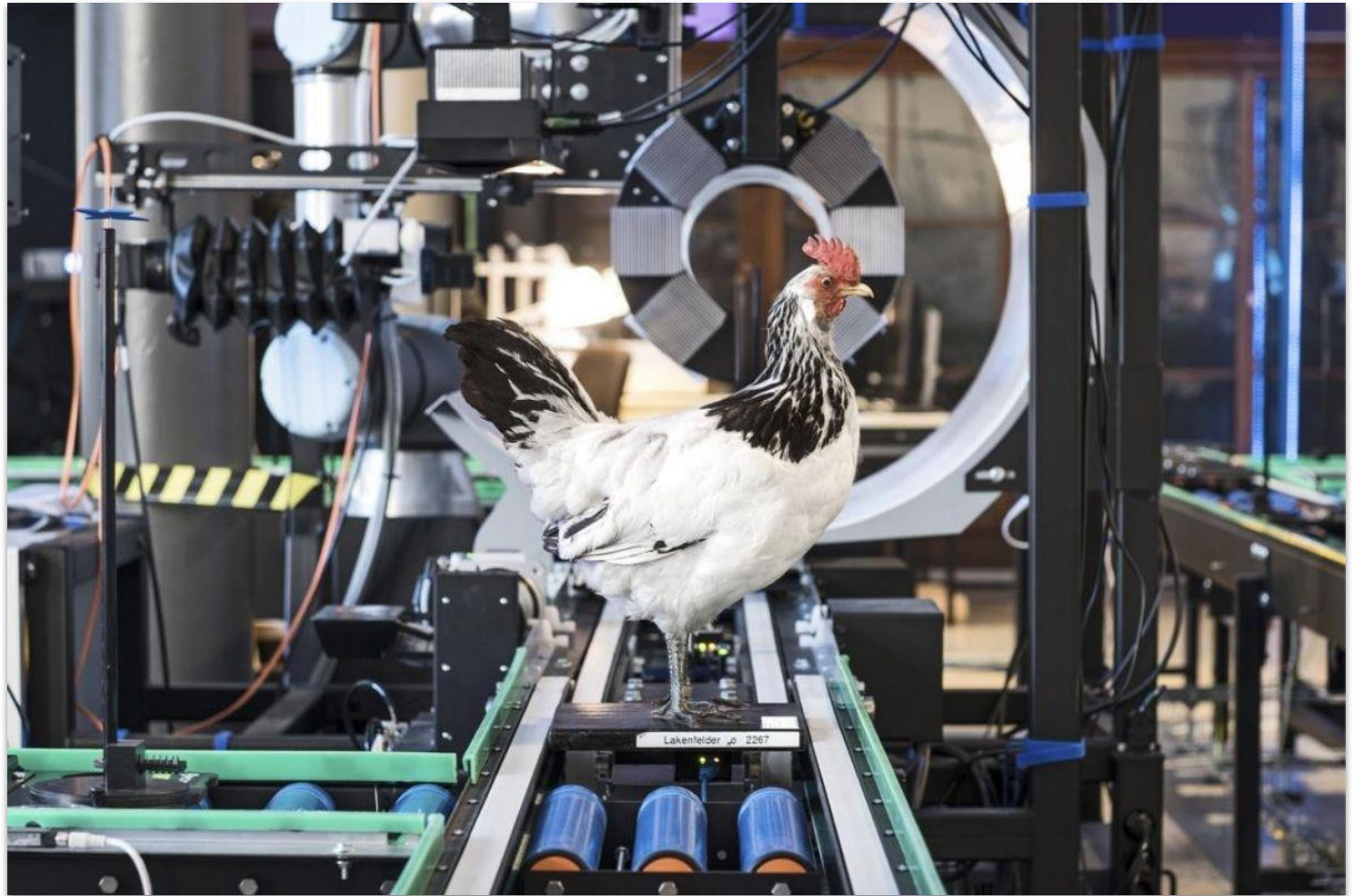
# Overview

- Part 1:  What is Linked Data - and how big is "Big"

- Part 2:  What does Linked Data for GLAM look like?

- Part 3: Collections Data vs Collections as Data

- Data Part 4: Linked Data in the wild: the good, the bad & the ugly
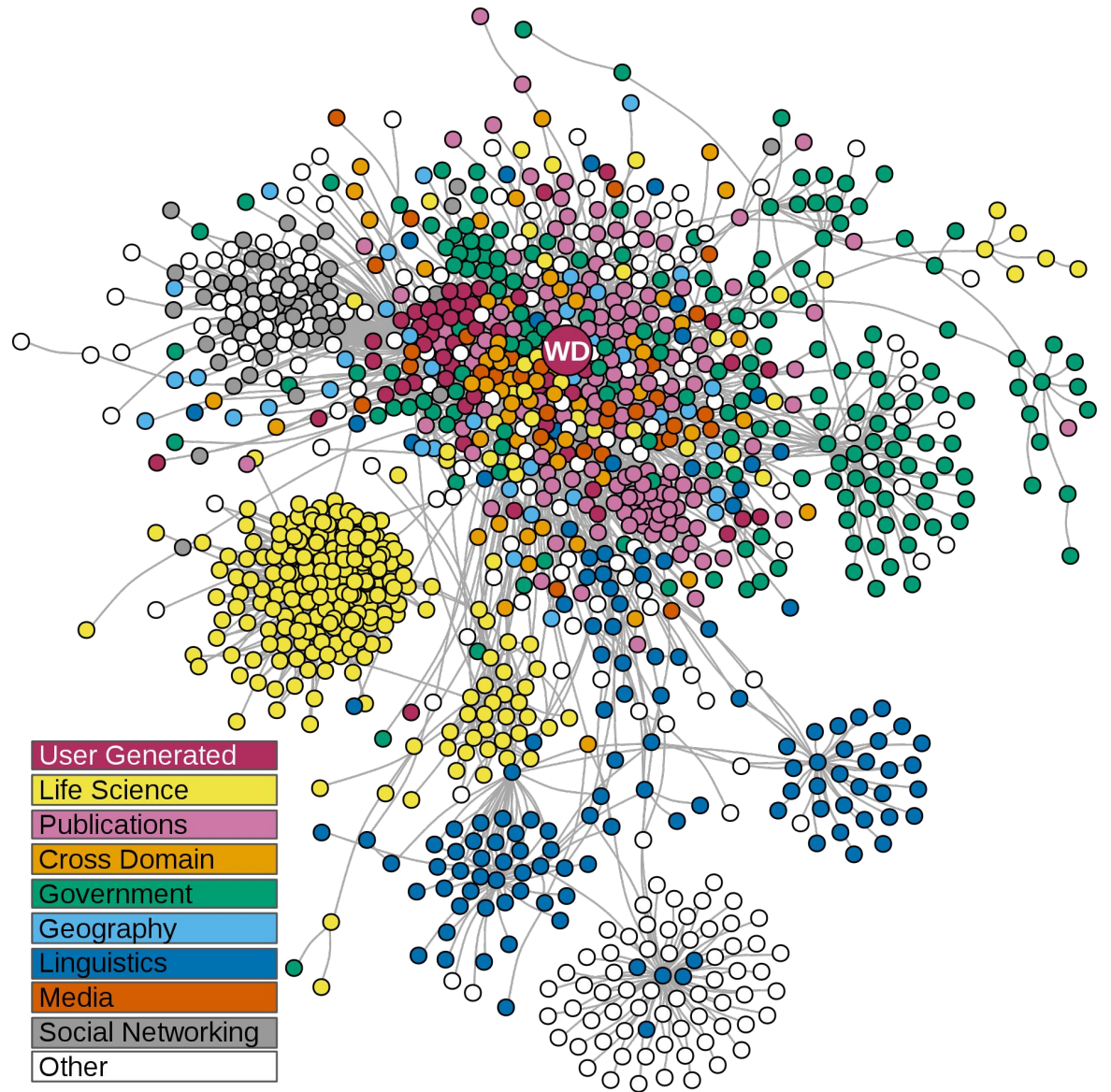
- Part 5: Looking forward

# Heritage digitisation has been going on for a while, with a variety of objectives.
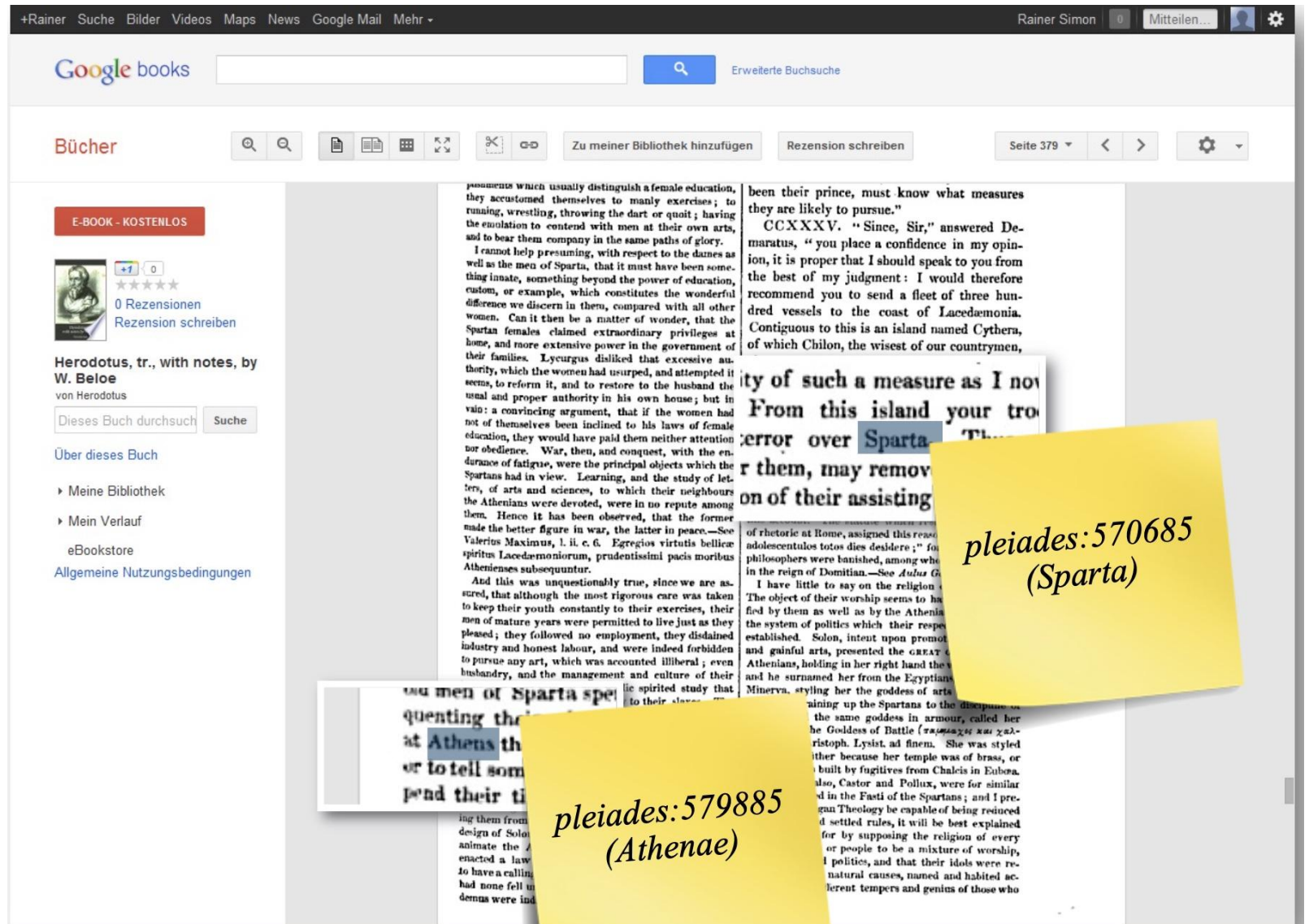


Chicken on a scanner: Museum für Naturkunde, Berlin.

Linked Data is a set of principles for connecting human and machine readable structured data



| Legend |
|---|
| User Generated |
| Life Science |
| Publications |
| Cross Domain |
| Government |
| Geography |
| Linguistics |
| Media |
| Social Networking |
| Other |

WD

# Semantic Recognition
(not quite yet a reality)…

From 5 star LOD to…

**From 5 star LOD to LOUD Data:**

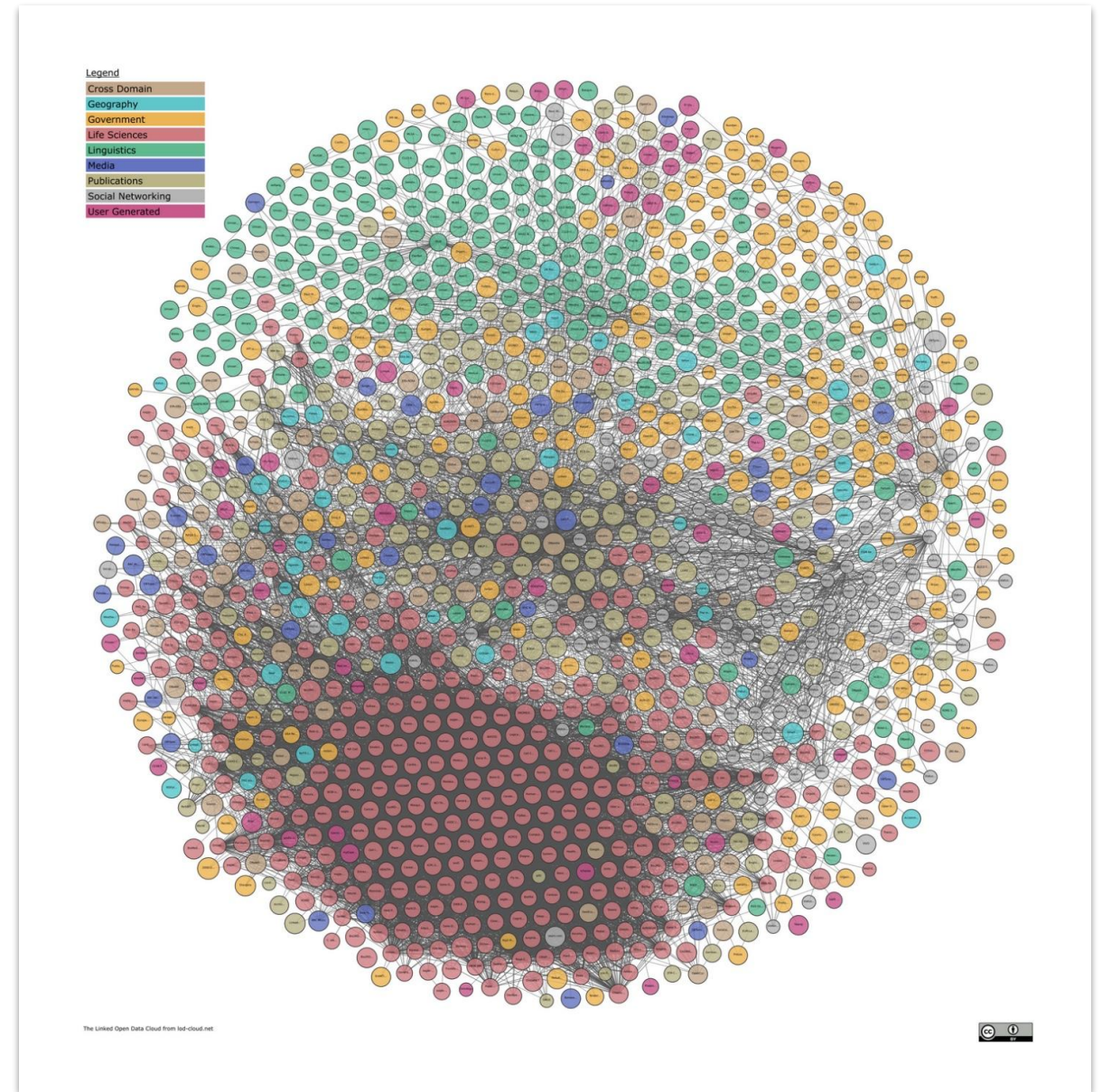**L**inked
**O**pen
**U**sable
**D**ata

→

**A**bstracted appropriately
**B**arriers to entry low
**C**omprehensible
**D**ocumented with working examples
**E**xceptions few, patterns are many

# "Big" is relative…

- Compared to other sectors, the heritage contribution is relatively small:
  - British museum: 4 million objects
  - Europeana: 50 million
  - Wikidata 98 million
- Data arrives in many different formats.
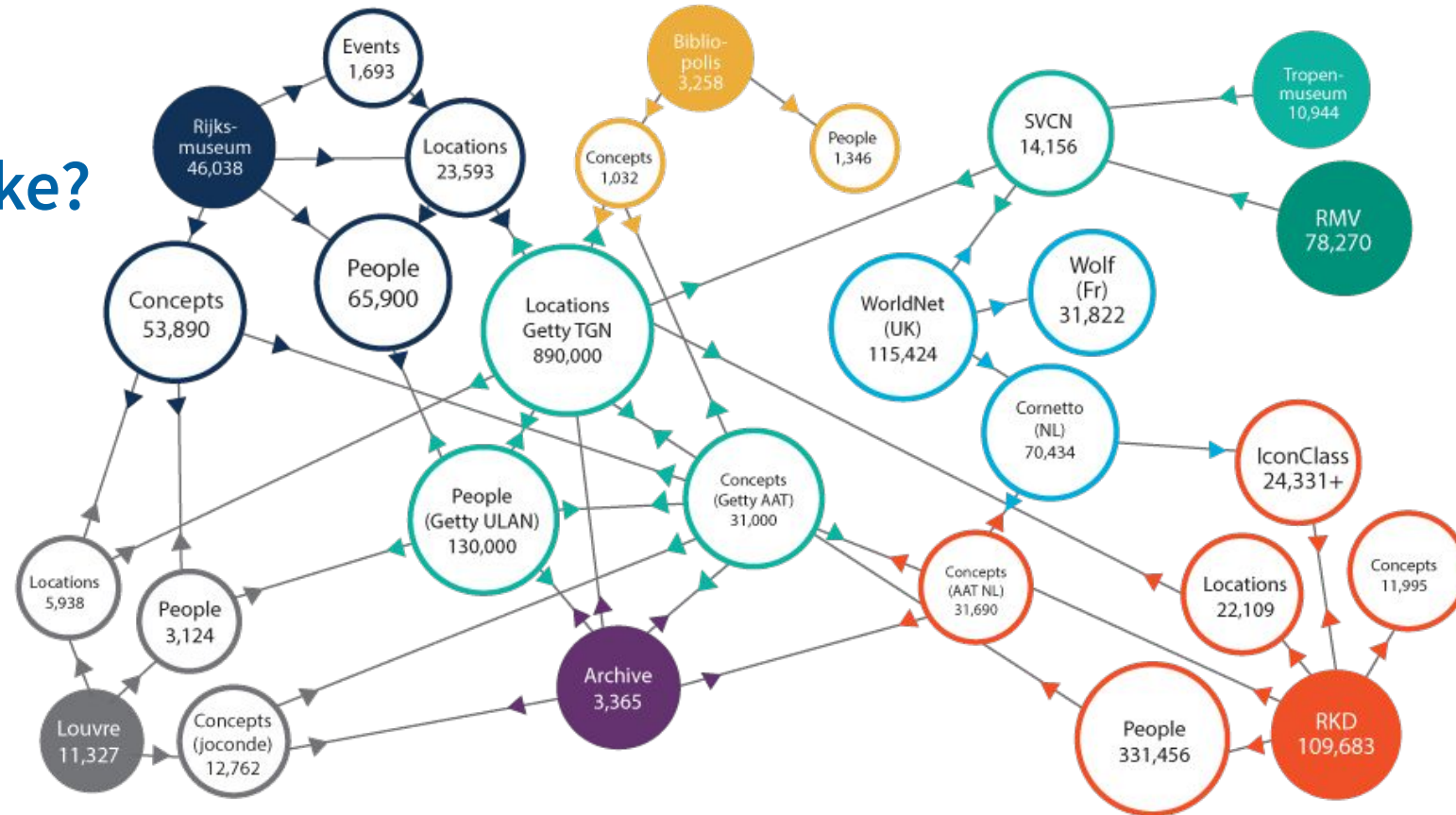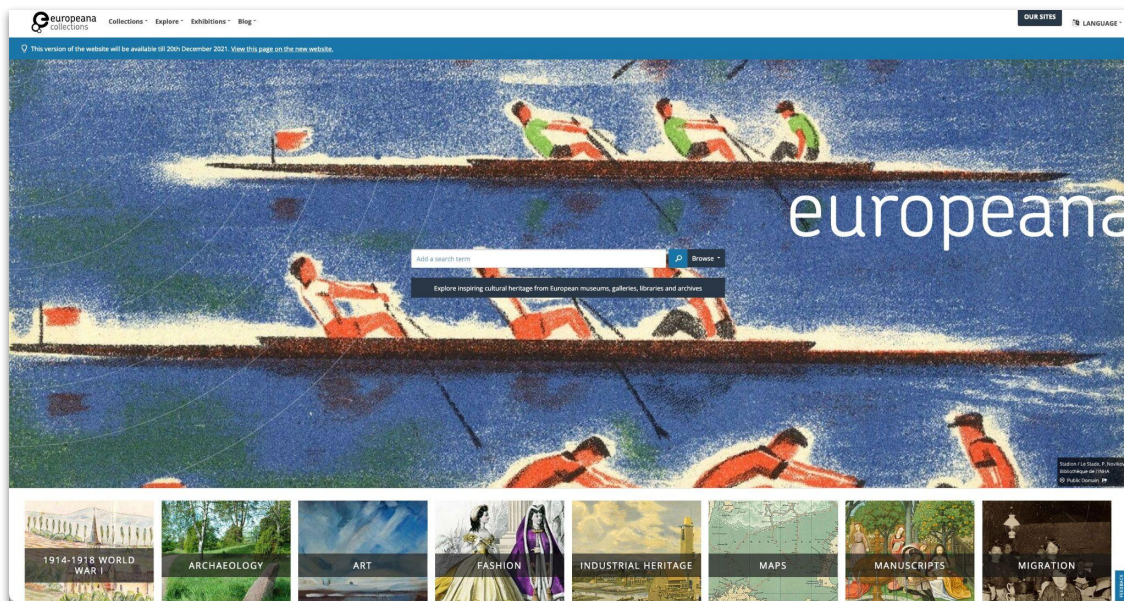- Interoperable systems are hard to build.

https://lod-cloud.net/

# Part 2:
# What Does LD for GLAM look like?

- Heritage institutions were early adopters
- Saw the value in linking collections
- Data is heterogenous
- Too many standards!
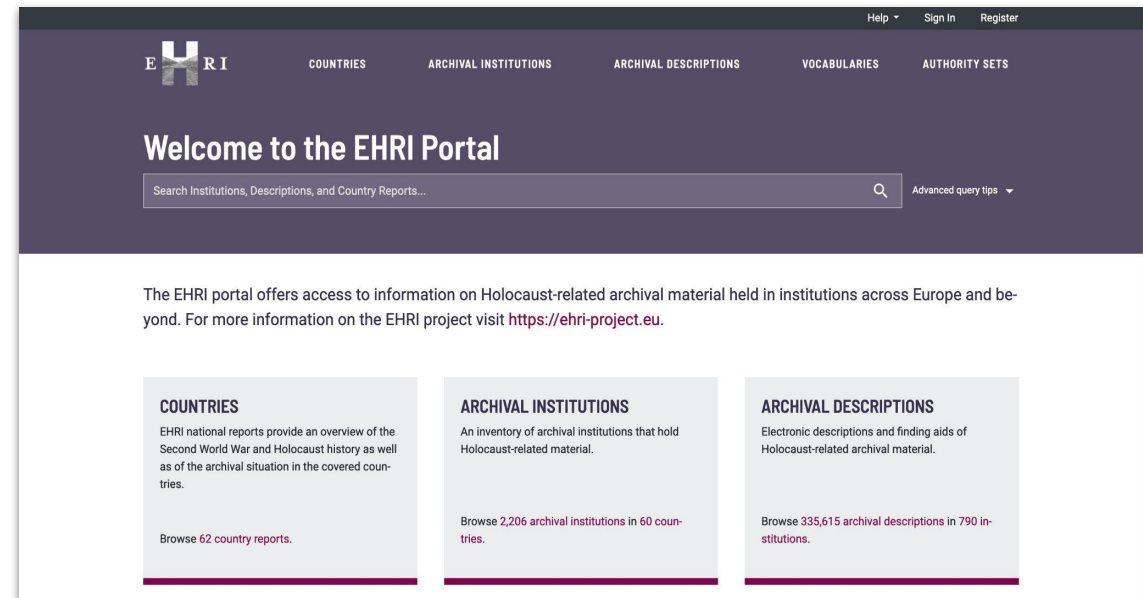
# Large Scale Infrastructures

## Europeana

## European Holocaust Research Infrastructure



https://www.europeana.eu

https://www.ehri-project.eu/

# Wikidata

# National Library of Wales in Wikidata

# Reassembling the Republic of Letters

# Challenges: Technical…



Figure 13.1 Sample original register page



Figure 13.2 An E03 input form, with information from the Oceania register for 1962



Figure 13.3 The equivalent form on a terminal screen, showing the first entry from the form in Figure 13.2

McCutcheon, D. (1986) The British Museum. In Light, R. B., Roberts, D., Stewart, J. D. (Eds.), *Museum documentation systems: Developments and applications*

# Challenges: Often Political



Home     About CENL     Member Libraries     News & Events     Reports & Publications     #Star

## TEL services to be discontinued from 31 December 2016

The European Library (TEL) was launched by the Conference of European National Librarians (CENL) in 2004 as the union catalogue of European national libraries and has since become a web portal and open data hub for national library data in Europe. Its success led to the Commission asking CENL to set up what became Europeana. The European Library has disseminated library data in a variety of ways to promote its wider use. TEL has been the

CENL News, December 13, 2016.
https://www.cenl.org/tel-services-to-be-discontinued-from-31-december-2016/

# Part 3: Collections data vs…

- Heritage institutions shape knowledge,
- Heritage professionals know how selection happens,
- What <u>is</u> included has been valourised

## Digital cultural colonialism: measuring bias in aggregated digitized content held in Google Arts and Culture

Inna Kizhner, Melissa Terras

College of Arts, Humanities and Social Sciences, University of Edinburgh, UK

Maxim Rumyantsev
Department of IT in Creative and Cultural Industries, Siberian Federal University, Russia

Valentina Khokhlova
Department of IT in Creative and Cultural Industries, Siberian Federal University, Russia

Elisaveta Demeshkova
Department of IT in Creative and Cultural Industries, Siberian Federal University, Russia

Ivan Rudov
Department of IT in Creative and Cultural Industries, Siberian Federal University, Russia

Julia Afanasieva
Department of IT in Creative and Cultural Industries, Siberian Federal University, Russia

Correspondence:
Inna Kizhner, Department of IT in Creative and Cultural Industries, School of Humanities, Krasnoyarsk 660049, Russia.
E-mail: ikizhner@sfu-kras.ru

# Collections as Data:
# Santa Barbara Statement on Collections as Data

- Aims to encourage computational use of digitized and born digital collections.
- Guided by ongoing ethical commitments.
- Aim to lower barriers to use.
- Designed for everyone, serve no one.
- Shared documentation helps others find a path to doing the work.
- Default open, except in cases where ethical or legal obligations preclude it.
- Data development values interoperability.
- Data stewards work transparently in order to develop trustworthy, long-lived collections.
- Data, as well as the data that describe those data are considered in scope.
- The development of collections as data is an ongoing process and does not necessarily conclude with a final version.
- **https://collectionsasdata.github.io/statement/**

# Lessons for AI research

○ Pre-curated sources of data are attractive…

○ But also complex

　■ Jo & Gebru, *Lessons from Archives* (2020):
　**https://arxiv.org/pdf/1912.10389.pdf**

○ "Haphazardly categorizing people in the data used to train ML models can harm vulnerable groups and propagate societal biases"

# Part 4: GLAM data projects: The good, the bad and the ugly

# Linked Data projects: the good

## Looted Art Detector

### About

**Objective:** Identify high priority artworks for provenance research

**Description:** Online Free Digital Tool

**Approach:** Automatic text analysis using frequency counts

Note: The frequency counts target textual indicators of UNCERTAINTY, UNRELIABILITY, or ANONYMITY, as well as the possible presence of RED FLAG names related to NAZI-looted art, forced sales and duress sales. The resulting calculations do not signify that an artwork is looted. They simply quantify observations concerning the text for further analysis.

### How it works

**The user uploads a CSV file that contains provenance texts**

Note: The uploaded CSV can contain other information as well - urls, titles, artists, etc. The only requirement is that the CSV also contain one column with the provenance texts.

The program will ask the user to enter the name of the column that contains the provenance text.

**The Provenance Text Analyser calculates the number of times key words appear in each provenance text and downloads a CSV named "results.csv"**

Note: The results.csv file contains all the original information uploaded by the user PLUS additional columns with word counts.

**The user uses his/her own tools to analyse the results.csv.**

### Recommendations: How to analyse "results.csv"

The Text Analysis provides quantitative indicators for the user to integrate in analysis. Which artworks are most likely to have problematic provenances?

1) Look for HIGH UNCERTAINTY

**https://artdata.pythonanywhere.com/about/**

---

Help | Sitemap | Disclaimer | Privacy statement | Contact          Deutsch | English | русский

### German Lost Art Foundation

**Search**

search ...   [Search]

▸ Start ▸ Lost Art-Database

Lost Art-Database
- Basics
- Search
- Advanced Search
- Search Requests
- Found-Object Reports
- Reporting objects

Module "Provenance Research"

#### Introduction

The Lost Art Database contains data on cultural objects which as a result of Nazi persecution or the direct consequences of the Second World War were removed and relocated, stored or seized from their owners, particularly Jews, or on cultural objects where, because of gaps in their provenance, such a story of loss cannot be ruled out as a possibility. In matter of found-objects a distinction in lost cultural assets as a result of the Second World War or National Socialist persecution isn't possible, because of the difficulty to separate in both categories. The owner normally didn't know anything about the provenance of the objects they have often inherited resp. they hardly can class them with historical processes.

The database is divided into two areas:

#### 1. Search Requests

It is possible here to register cultural objects lost by public institutions or private individuals and institutions as a result of National Socialist rule and the Second World War, requesting a world-wide search via the Lost Art Internet Database. Owners or custodians of cultural objects with an uncertain or incomplete provenance can search here whether these objects have been sought elsewhere.

#### 2. Found-Object Reports

It is possible here to register cultural objects where it is known that they were taken illegally from their owners or relocated to another place as a result of the war. The section also contains reports on cultural items with an uncertain or incomplete provenance, suggesting the possibility of illegal dispossession or a

**Basics**

During the years of Nazi domination between 1933 and 1945 a relocation of cultural items took place whose full scope has still not been completely explored and investigated.

▸ More

**General principles**

▸ 📄 General principles for the registration and deletion of reports in the Lost Art Database

**Reporting objects**

The German Lost Art Foundation accepts search reports and found reports from private individuals and institutions for cultural assets that were removed, relocated or confiscated, particularly from Jewish owners, as a result of the National Socialist regime or the Second World War, or for which a history of loss cannot be ruled out due to gaps in provenance.
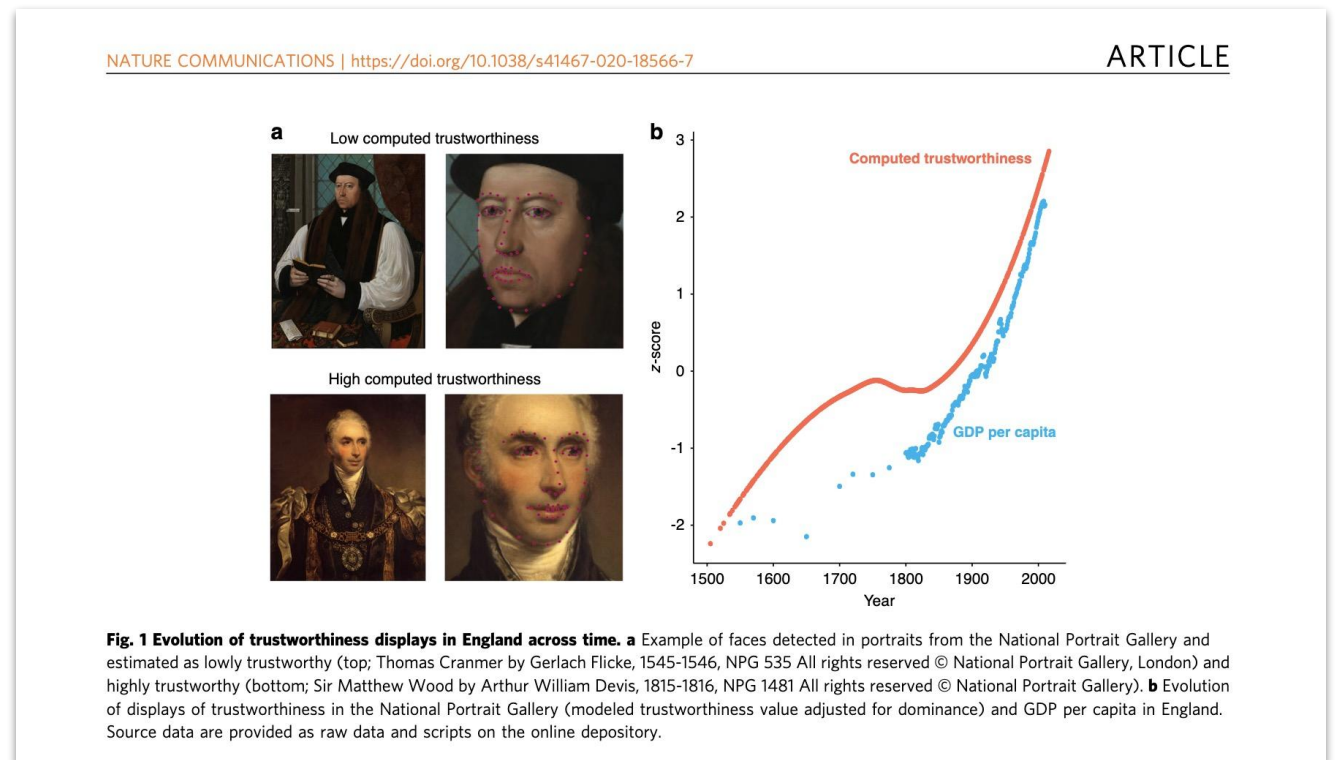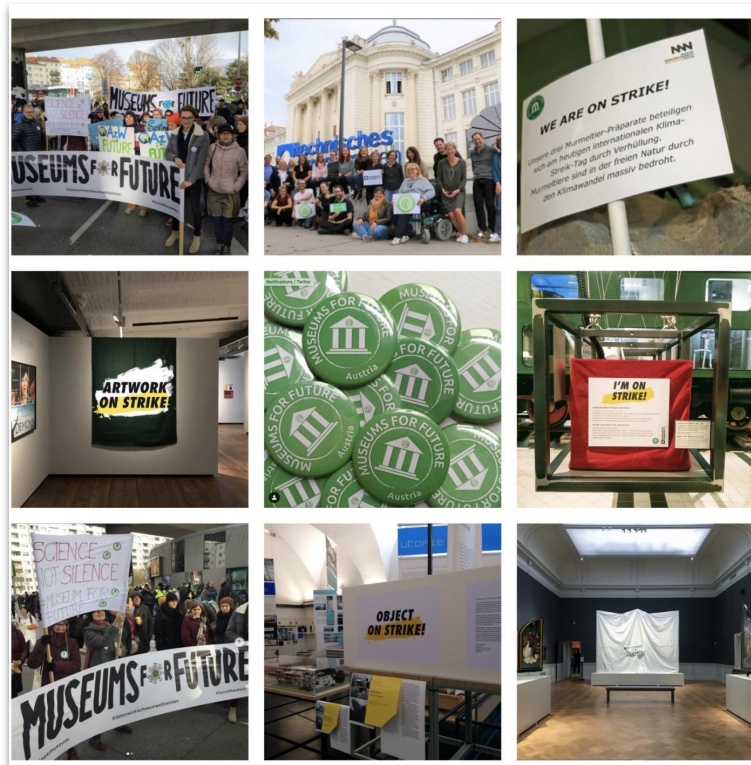
▸ More

---

# Linked Data Projects: the bad





**Screenshot of cover page and figures from the trustworthiness study "Tracking historical changes in trustworthiness using machine learning analyses of facial cues in paintings" conducted by Baumard, et al on European portraits (courtesy and via Nature.com open access) Full paper at: https://rdcu.be/b8PAF**

# Linked Data projects: the ugly little secrets





**Pinned Tweet**

**Is the British Museum's endpoint working?**
@bm_lod_status

Automated

The endpoint is down and it seems unlikely that it's ever coming back. A once leading platform for open heritage data from one of the world's major museum collections, gone. Over two million "persistent" URIs – relied upon as authoritative by many other projects – dead.

🪦 🥓 💀

7:01 PM · Feb 11, 2022 · Twitterrific for iOS

**98** Retweets   **34** Quote Tweets   **217** Likes

# Why this matters:

"Museums are seen as a beacon of trust... The level of criticality museums have when considering collections is the same level they need to have when it comes to developing digital applications."

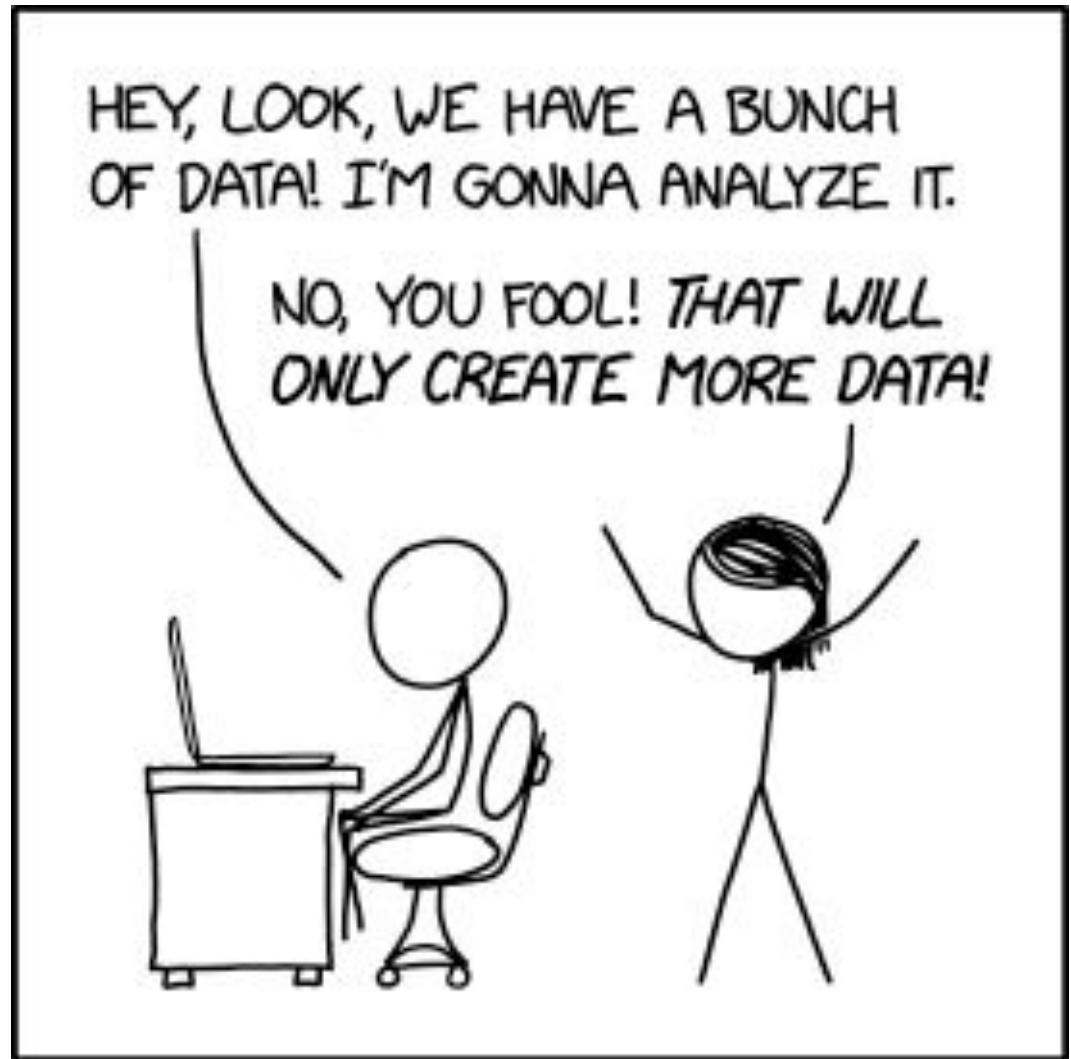Dr Oonah Murphy, Goldsmiths University, Museums + AI Network.



**www.themuseumsai.network/toolkit/**

**Thank you!**

rebecca.kahn@univie.ac.at

@rebamex



Data Trap: xkcd, CC BY-SA https://xkcd.com/2582/