# Data citation for the Humanities and Social Sciences: a special case?

Barbara McGillivray (University of Cambridge and The Alan Turing Institute)

Nicolas Larrousse (CNRS / Huma-Num - SSHOC Project WP3)

Daan Broeder (CLARIN ERIC - SSHOC Project WP3)

# Open Science and Open Knowledge

- Our common goal is **Open Science,** meaning we promote openness, integrity, and reproducibility in research;
- and the results to be treated as **Open Knowledge**; knowledge that is free to use, reuse, and redistribute without legal, social or technological restriction
- The **FAIR principles** help us with both
  - *Findability*, Accessibility, Interoperability, *Reusability*
- We (researchers, librarians etc.) like it, public funders like and require it

**Project:**

# SSHOC
social sciences & humanities open cloud

Horizon 2020
European Union Funding
for Research & Innovation

**Type of action & funding:**
**Research and Innovation action**
(INFRAEOSC-04-2018)

**Partners: 48**
(23 beneficiaries + 25 LTPs)

**SSH ESFRI Landmarks and Projects**
& international SSH data infrastructures

**Project budget:**
**€ 14,455,594.08**

**Project website:**
**www.SSHOpenCloud.eu**

**Duration: 40 months**
(January 2019 – 30 April 2022)

Objectives:

- creating the social sciences and humanities **(SSH)** part of European Open Science Cloud **(EOSC)**
- maximising **re-use** through **Open Science** and **FAIR** principles (standards, common catalogue, access control, semantic techniques, training)
- interconnecting existing and new infrastructures (clustered cloud infrastructure)
- establishing appropriate **governance model** for SSH-EOSC
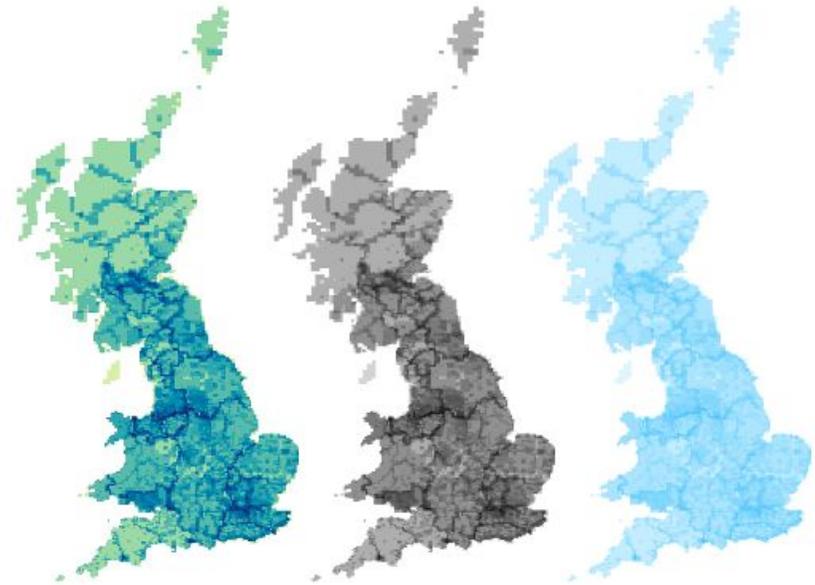
# Specificities of SSH



Growing adoption of data-intensive methods

New types of digitized and born-digital data, and very diverse types of data, including cultural heritage

Growing attention to data, with new challenges and needs

Increased number of data professionals and stewards in libraries

# A variety of data types

# Data traditions and data practices

**Variety**

Different traditions for quantitative/computational methods

**Focus**

Historically, on publications (book) rather than data

# From data to data papers

Links

| Deposit and publication of data | Data documentation | Data paper |
|---|---|---|
| Typically in open repositories assigning DOI | +intellectual organisation | Peer-reviewed publications describing datasets and crediting their creators |

# The Journal of Open Humanities Data

- Launched in 2015; open access and peer-reviewed
- Vision: be a key part of a thriving community of scholars sharing humanities data
- Fastest growing data journal for humanities research
- 25 articles published so far (6 in 2021)
- Two publication types
- Short data papers (1000 words): descriptions of datasets with high reuse potential
- Full-length research papers (3000-5000 words): discussion of methods, challenges, and limitations in the creation, collection, management of humanities data



https://openhumanitiesdata.metajnl.com/

# Data citation in SSH

- The notion of "publishing data" is relatively new
- Data until a recent period was not really "*noble*" … the focus was more on the "final product" (e.g. an article or a book)

- Very diverse and no specific common approach to data citation A lot of initiatives … so far no real standard, but more "communities of practices" (See SSHOC deliverable " Inventory of SSH citation practices …."  https://doi.org/10.5281/zenodo.4436736)

REVUE D'ÉGYPTOLOGIE 67

PUBLIÉE PAR LA
SOCIÉTÉ FRANÇAISE D'ÉGYPTOLOGIE

PARIS – ÉDITIONS PEETERS

2016

-> In order to cite data, you need **infrastructures** to publish them in a proper (FAIR) way but you also must develop mechanism(s) to cite SSH data:
- Build stronger links between data and publication
- Facilitate and develop a "culture of data publication/citation"

etc.

DIGITAL MEDIEVALIST

One goal of SSHOC is to address the problems of the current data landscape with its disciplinary silos

**SSHOC Project and Data citation**

- Recommendations based on "Force11 Citation Principles" adapted to SSH specificities
- Make citations "actionable" as a classic citation isn't machine actionable

-> Proposition to create a prototype of *"FAIR SSH Citation Infrastructure"*

- Process existing citations or create new ones
- Enrich them automatically (other sources) and manually (annotations)
- Standardized citations
- Publish citations into the "Citation Infrastructure" to make them actionable

# The future for data publication and citation

Although the potential for data creation, processing and (re)sharing have increased  enormously

- Data citation practices still mainly based on traditional publication model ie. paper citing paper citing paper
- More complex collaborative and dynamic workflows are possible
- and also needed for using distributed and dynamic data sources e.g. social media and virtual collections of heterogeneous and distributed data
- Key is proper data identification, description and provenance tracking during the whole data life-cycle

# Libraries in the centre?

- Increased need for Openness and  Sharing of data and information at all phases of the SSH DLC
- More data and  new types of data
    - how do we cite millions of data points?
    - how do we cite dynamic data e.g. social media
- Are the current practices wrt data publication and citation still valid?

- Libraries may ask if they can go far beyond caring only for the end-result (papers & data),
- they are  already involved in
    - storing project result data-sets for reproducibility
    - training and facilitating researchers wrt  Open Science and Open Knowledge and FAIR principles
    - Providing catalogues and meta-science information
- However  the physical and conceptual distance between Libraries and research labs poses limitations

Open science

Data citation

Data publication