



HAL
open science

Pilot study on the intercalibration of a categorisation system for FAIRer digital objects related to sensitive data in the life sciences

Romain David, Christian Ohmann,, Jan-Willem Boiten,, Steve Canham,,
Maria Panagiotopoulou,, Theresia Mayrhofer, Michaela, Dario Longo,,
Florence Bietrix,, Luisa Chiusano, Maria, Arnaud Laroquette,, et al.

► To cite this version:

Romain David, Christian Ohmann,, Jan-Willem Boiten,, Steve Canham,, Maria Panagiotopoulou,, et al.. Pilot study on the intercalibration of a categorisation system for FAIRer digital objects related to sensitive data in the life sciences. RDA Plenary 18, Nov 2021, Edimbourg, United Kingdom. , 2021. hal-03663991

HAL Id: hal-03663991

<https://hal.science/hal-03663991>

Submitted on 10 May 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

R. David (ERINHA), C. Ohmann (ECRIN, pilot study coordinator), M. Cano Abadia (BBMRI), F. Bietrix (EATRIS), J.-W. Boiten (EATRIS/Lygateure), S. Canham (ECRIN), M.L. Chiusano (EMBRC), W. Dastru (Euro-BioImaging), A. Laroquette (EMBRC), D. Longo (Euro-BioImaging), M.Th. Mayrhofer (BBMRI), M. Panagiotopoulou (ECRIN), A. Richard (ERINHA), P. Verde Contact: christianohmann@outlook.de

Background: Sharing sensitive data is a specific challenge within EOSC-Life. A toolbox is developed, providing information to researchers who wish to share and use sensitive data and to support the workflow of handling these kinds of data objects.

The **objective** of the pilot study (Ohmann *et al.*, 2021) was to evaluate the first version of a categorisation system for digital objects from life sciences by human experts from different research infrastructures with respect to **applicability, reliability, sustainability and consistency**. The primary objective is to identify the main improvements that are necessary to arrive at an inter-community, well understood and **commonly approved categorisation system**. From the results, an improved version of the categorisation system has been provided and will be used as a basic component for the toolbox under development.

Pilot study method: Each involved infrastructure nominated two experts, willing to perform the individual assessment of a number of typical resources around sensitive data to be included in the EOSC-Life toolbox. The experts from the involved infrastructures **selected up to 25 resources**, spanning a wide range of resource types (e.g. legislation & regulations, position papers, policies & principles, background & explanatory material, tools). *The experts selected the resource types that are relevant for their research infrastructure and then assessed these resources independently of each other using the categorisation system that we developed.*

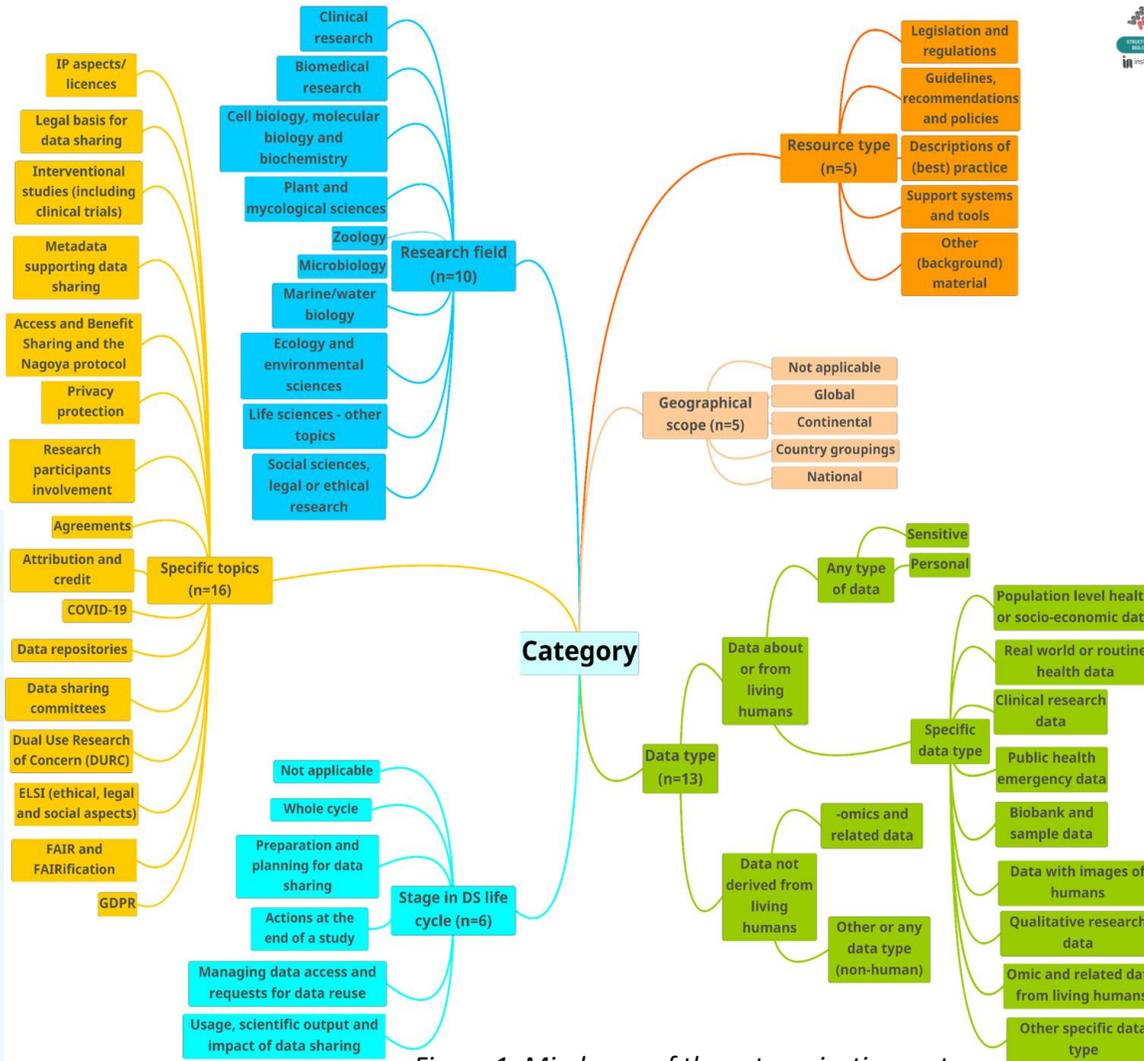


Figure 1: Mind map of the categorisation system



What is EOSC-Life? EOSC-Life is an H2020-funded project bringing together the 13 Life Science 'ESFRI' Research Infrastructures to create an open, digital and collaborative space for life science research.

www.eosc-life.eu

Main results: The aspects to be improved to reach the goal of a consistent, sustainable, cross-disciplinary tagging in life sciences are:

- Standard definitions of tags required
- Elimination of ambiguity between tags
- Guiding the number of tags assigned for a category
- Filling gaps particularly in missing "resource types" and "research fields"
- Elimination of meaningless tags
- Simplifying research stage
- Simplifying data type tags
- Guidance note and training needed for taggers

The result is a categorisation system with 55 tagging values spread over 6 dimensions (Figure 1), as compared to the original 93 tags spread over 8 dimensions.

Next step / take-home message
Avoiding polysemic bias is one of the most difficult issues to resolve and is necessary across disciplines. This categorisation system will **need iterative improvement**. The pilot study on tagging intercalibration could serve as a model for a future cross-disciplinary workflow tagging system in Canonical Workflow Framework of Research (CWFR).