# The STORE database; towards FAIR

Paul Schofield, University of Cambridge

# *Open Science*

- Open publication access
- Open scholarly communication
- Open research data

# *Open Science*

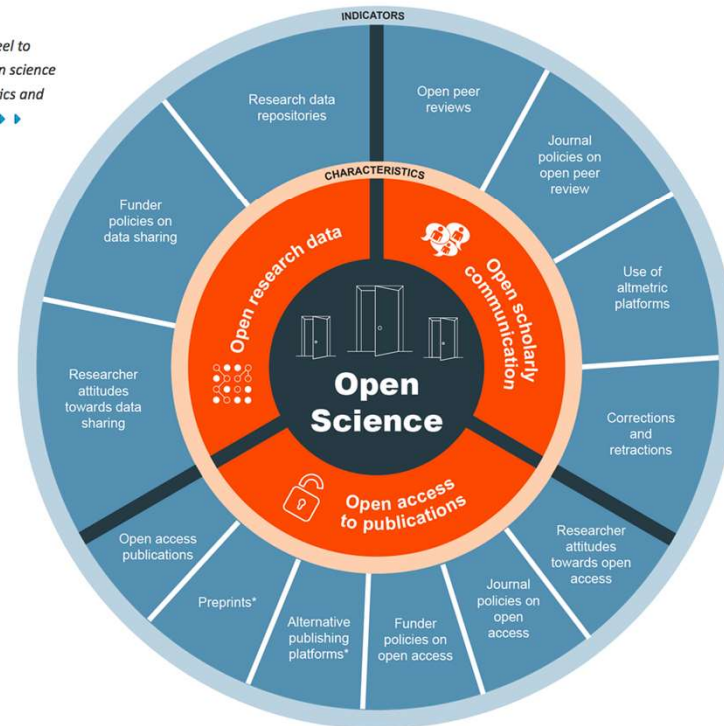"Open science is an approach based on open cooperative work and systematic sharing of knowledge and tools as early and widely as possible in the process"

European Commission

| Open Access to publications | Responsible management of data (FAIR principles) | Open access to data 'as open as possible, as closed as necessary' | Information about outputs / tools / instruments to validate/re-use results and data | Digital /physical access to results to validate the conclusions |
|---|---|---|---|---|

# EC Requirements

- Must manage the digital research data in line with the **FAIR principles** (Findable, Accessible, Interoperable, Reusable)

- **Data Management Plan** (DMP) is required by M6; updated mid-project and at end of project

- **Deposit (meta)data as soon as possible** after production/generation or after processing and quality controls

- Deposit data in a **trusted repository** and make them **open as soon as possible** (deadlines set in DMP), following the "as open as possible, as closed as necessary" (open by default) principles

- Data closed if necessary, but **metadata must be FAIR and under CCO** (trusted repositories will automatically share metadata in CCO)

- Open licence, preferentially CC-BY or CC0 licence

- Detailed information about research outputs or tools/instruments needed to re-use or validate the data (e.g. data, software, algorithms, protocols, models, workflows, electronic notebooks)

**Examples of metadata**
author(s) name, author(s) ORCID, DOI, licence, language, journal, title, etc.

# Why do we need to share data?

Evidence support for publications
    Reproducibility
    Accountability
Data and resource reuse and reanalysis
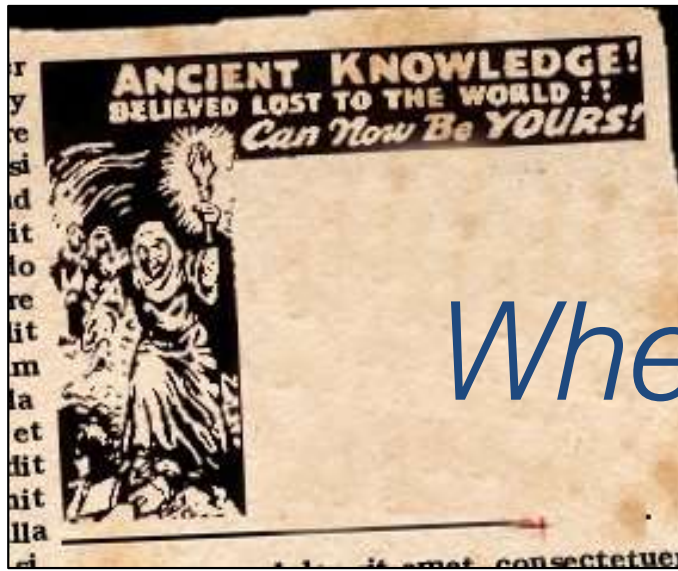    Resilience
    New discovery
    Prevention of duplication
    Supports RRRs
    Value for funding
Coordination and integration of projects
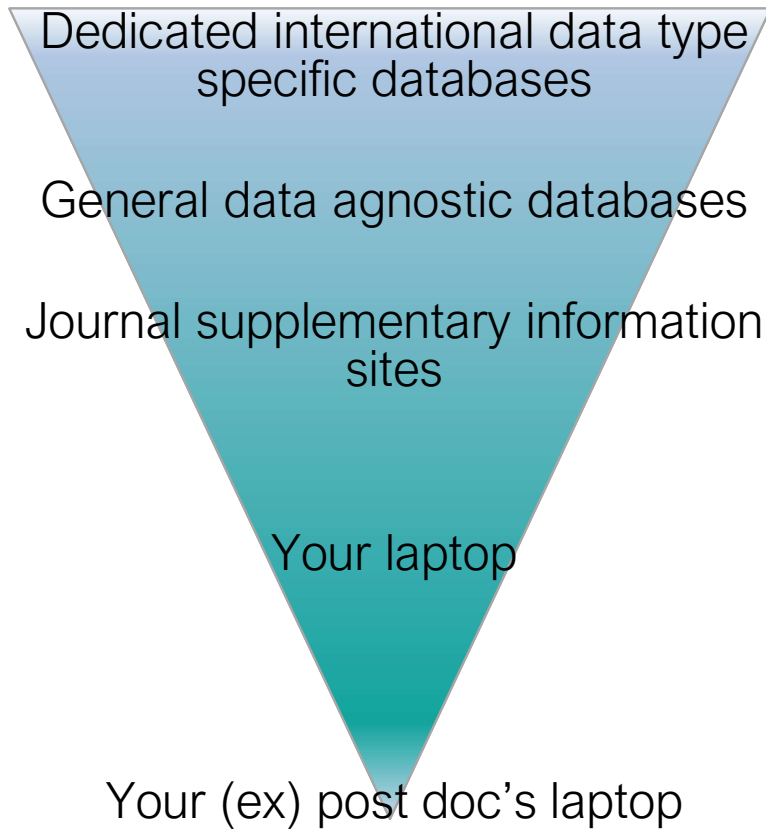Training and education resources

*Where's the data?*

# Data resources

Dedicated international data type specific databases

General data agnostic databases

Journal supplementary information sites

Your laptop

Your (ex) post doc's laptop

# Open Science

"Open science is an approach based on open cooperative work and systematic sharing of knowledge and tools as early and widely as possible in the process"

European Commission

| Open Access to publications | Responsible management of data (FAIR principles) | Open access to data 'as open as possible, as closed as necessary' | Information about outputs / tools / instruments to validate/re-use results and data | Digital /physical access to results to validate the conclusions |
|---|---|---|---|---|

# The STORE database

Funded by the European Commission EURATOM Programme and the Bundesamt fuer Strahlenschutz since 2009

European Commission | Horizon 2020 European Union funding for Research & Innovation

**RadoNorm**
Managing risks from radon and NORM

elixir

re3data.org
REGISTRY OF RESEARCH DATA REPOSITORIES

i a © pi ◎ §

http://doi.org/10.17616/R3732R

**STOREDB**

Free to the entire community and sustained by the BfS Self curation, upload and controlled access.

FAIRsharing.org
standards, databases, policies

DataCite
FIND, ACCESS, AND REUSE DATA

Findable Accessible Interoperable Reusable

**Environmental data**

**Primary experimental data**

**Social sciences and policy documents**

**Epidemiological, health and genetic data**

**STOREDB**

**Legacy datasets**

**Protocols and software**

http://www.storedb.org

CC BY NC SA

Funded by Euratom research and training programme 2014-2018 under grant agreement No 900009.

# *FAIR sharing*

## Use of data at scale by humans and machines

# *OPEN* AND *FAIR* ARE NOT THE SAME THING

"Open" is about data rights and licensing.

"FAIR" is about mechanics.

Anne Raugh

# *FAIR* IS PRIMARILY CONCERNED WITH PROGRAMMATIC PROCESSING

*FAIR* recognizes that data is diverse and scattered across cyberspace. Software processing levels the playing field.

Anne Raugh

The FAIR Principles at *go-fair.org*

# *FAIR is not a binary state*

FAIRness is a spectrum. There are various ways to increase FAIRness and they can be applied incrementally.

BUT

- Use of generic data repositories is not very FAIR

- Lack of granularity in metadata

- Lack of metadata standards relevant to theme or community

- Non-grounded

The FAIR Principles at *go-fair.org*

# Open Science Projects

Open Science Projects primary goals aim to increase collaborative scientific data sharing, analysis and more rapid scientific advancement.

## GeneLab

GeneLab, an open science multi-omics repository, covering transcriptomics, metagenomics, epigenomics, proteomics, and metabolomics. Studies comprise of data from model organisms including microbes, plants, fruit flies, rodents and humans.

Learn more GeneLab

## BSP

The NASA Space Biology Biospecimen Sharing Program (BSP) collects biospecimens to maximize the scientific return from biological spaceflight and associated ground investigations and to encourage and broaden participation from the scientific community in space biology-related research.

Learn more about BSP

## ALSDA

Ames Life Sciences Data Archive (ALSDA) collects, curates, and makes available space-relevant higher-order phenotypic datasets. Datasets that enable scientists to perform retrospective analysis across missions, experiments, life science disciplines, research subjects, and species.

Learn more about ALSDA

## NBISC

NASA Biological Institutional Scientific Collection (NBISC) is a biorepository of non-human samples collected from NASA-funded spaceflight investigations and correlative ground studies. The purpose of NBISC is to receive, store, document, preserve, and make the collection available to the scientific community.

Learn more about NBISC

# GeneLab Data System



1. NASA GeneLab - the first comprehensive space-related omics database (https://genelab.nasa.gov/)
2. Investigators can search, download, submit, privately share, and/or analyze omics data from spaceflight and corresponding ground-analog experiments.
3. GeneLab Data Systems users can explore GeneLab datasets in the Data Repository, submit omics data through the Submission Portal, analyze data using GeneLab Analysis Platform tools, and Visualize study results through Visualization Portal.
4. GeneLab also offers biospecimen processing services and NASA's Institutional Scientific Collection is a space-research biobank that offers potential investigators hundreds of biospecimens for further analysis.

# NASA Life Sciences Data Archive (LSDA)



The LSDA is a collection of NASA's life sciences research data and information, consisting of human, animal, microbe, and plant studies conducted from 1958 to present.

# Standardisation of metadata

## Ontologies

Comprised of standardised hierarchical concepts and relationships

Capture semantic data and metadata

# Radiation Biology Ontology (RBO)

- Facilitates data retrieval and query expansion from NASA's GeneLab 'omics database and STORE



- Contains over 200 annotated classes and instances specific to the study of radiation on biological systems, as well as imports of more than 3500 additional classes from 13 other OBO Foundry ontologies

- Published through the OBO Foundry

- Available through NCBI Bioportal web site and application programming interface at https://bioportal.bioontology.org/ontologies/RBO

### Basic Formal Ontology (BFO)

- Study types
- Environment
- Radiation
- Experimental modality
- Radiation Source
- Disease and anatomy
- Taxon

**radiobiology study type**
- adaptive radiation response study
- anatomical study
- biokinetics study
- Cancer study
- carcinogenesis study
- Clinical study
- dna damage and repair study
- Environmental studies
  - Ecological population modelling study
  - Ecotoxicology study
  - Environmental radiation monitoring study
  - Environmental radionuclide transfer study
  - Environmental Radon study
  - Environmental study_ Abiotic_anthropogenic
    - Building materials radiological safety study
    - Building radiological safety study
  - Natural environment studies
  - Naturally occurring radioactive materials study
  - Non ionising electromagnetic radiation study
  - Radionuclide dispersal modelling study
- environmental study
- epidemiological study
- external exposure study
- fractionated radiation exposure
- gene expression study
- ground analog study
- ground control study
- internal contamination study
- Laboratory study
- Legal and governance study
- lifespan study
- marker discovery study
- Mass media study
- metabolomics study
- mixed exposure route study
- nuclear accident study
- Nuclear industry study
- offspring study
- physiological study
- Preparedness study
  - Civil protection study
  - Disaster planning study
  - Nuclear accident study
  - Situation awareness and decision support study
- proteomics study
- Security and law enforcement study
  - Bioterrorism study
  - Military Defence study
- Social and psychosocial studies
  - Attitudinal study
  - Behavioural study
  - Communication study
  - Community study
  - Holistic approaches to governance study

# *Radiation Biology Ontology (RBO)*

- Facilitates data retrieval and query expansion from NASA's GeneLab 'omics database and STORE

- Contains over 200 annotated classes and instances specific to the study of radiation biological systems, as well as imports more than 3500 additional classes from 13 OBO Foundry ontologies

- Published through the OBO Foundry

- Available through NCBI Bioportal web site and application programming interface at https://bioportal.bioontology.org/ontologies/RBO

# *How FAIR is STORE?*

- ✓ Unique identifier

- ✓ Identifier persistence

- ✓ Resolvable identifier (identifiers.org)

- ✓ Structured metadata

- ✓ Grounded metadata (resolvable metadata IDs)

- • Data identifier explicitly in metadata

- • Metadata identifier explicitly in metadata

- ✓ Programmatic access

- ✓ Licenses

# Work to be done

- Full implementation of the RBO ontology

- Implementation of the REST services

- Implementation of automated versioning

- Refinement of user interface

- Implement ISA-TAB to add licences, metadata and data IDs etc to data header.

- Investigate automatic extraction of candidate metadata classes from data record using NLP.
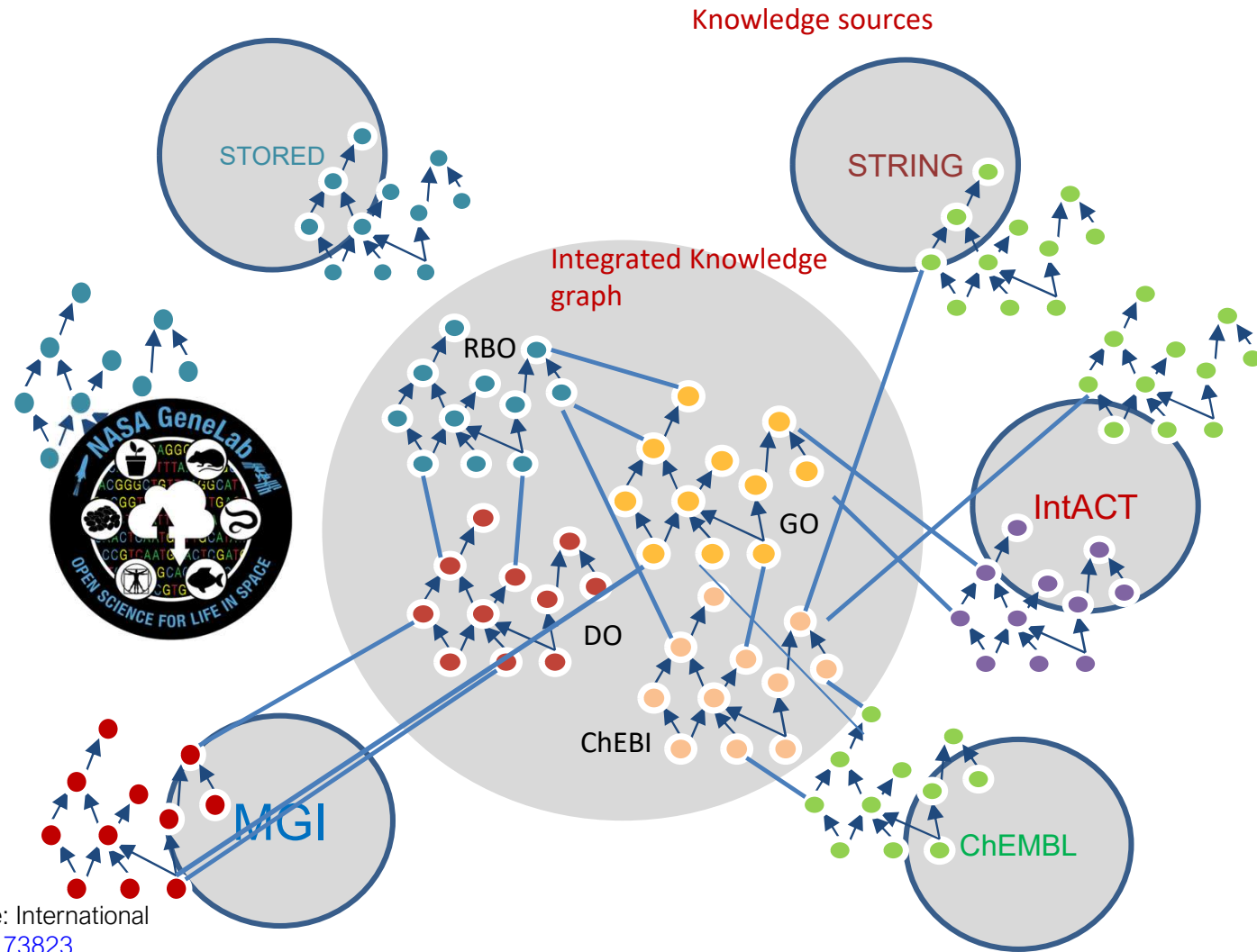
# Using open data to create new knowledge

Semantic data standards and metadata

- Discovering and integrating data between databases
- Federated queries and query expansion
- Semantic integration
- Ontologies permit the assertion of defined relationships between concepts

RBO cross references 13 OBO ontologies directly providing semantic linkage to most relevant databases

- Construction of Knowledge Graphs
- Graph embeddings for data representation
- Classification and similarity
- Inductive inference
- Learning over graph convolutional neural networks

Wilson et. Al. (2023) Machine intelligence for radiation science: International Journal of Radiation Biology, DOI: 10.1080/09553002.2023.2173823

# Acknowledgements

Daniel Berrios

Karin T. Slater

Jack Miller

Kristen Peach

Sylvain V. Costes
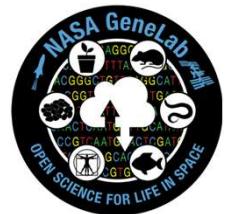
Ulrike Kulka

Michael Gruenberger