# Foundational Data Stewardship Workshop

# Outline for this week

- Day 1
  - Intro to RDM, open research, DMPs and FAIR (S. Venkataraman)
- Day 2
  - Open and responsible research (Louise Bezuidenhout)
- Day 3
  - Practical implementation (Joy Davidson)

# Introduction to Research Data Management and Open Research

S. Venkataraman PhD, Research Data Specialist, Digital Curation Centre

s.venkataraman@ed.ac.uk

*Data Stewardship Workshop, 17th May 2021, University of Botswana (Virtual)*

# About the DCC

o Established in 2004.

o Based in Edinburgh and Glasgow.

o Works at national and international levels.

o One of leading organisations in the world specialising in training, consultancy, policy making and advocacy in digital data management best practice and services provision.

o Involved in many international consortia and schools.

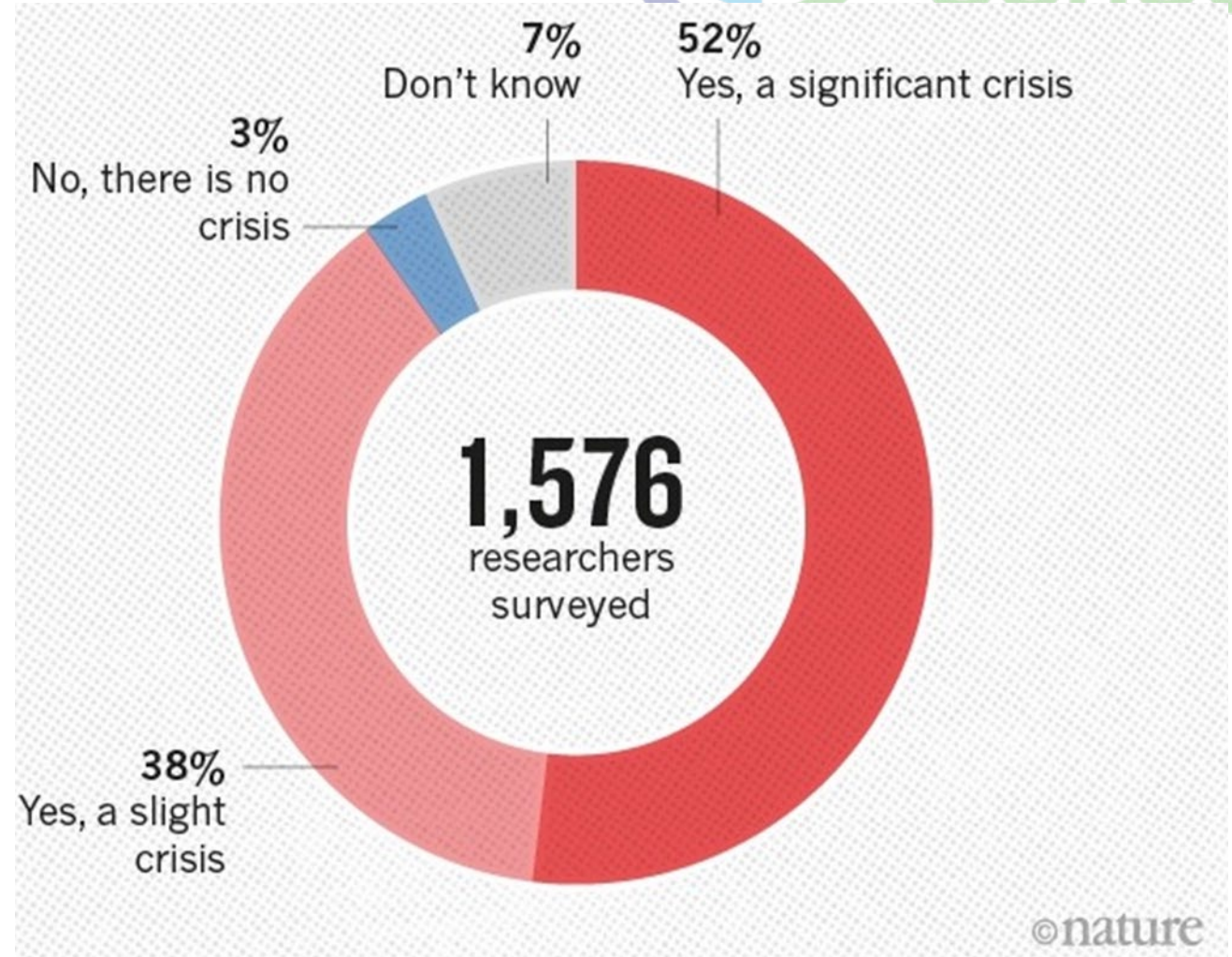o (We do not curate any data ourselves!)

# Learning outcomes

o Be familiar with the curation lifecycle.

o Understand the standardisation methods and principles available to add value to your data.

o Learn about resources to aid your workflows.

o Increase/encourage your level of openness.

o Learn about data management plans and the value in implementing them.

o Understand how data stewards integrate this knowledge

# Is there a reproducibility crisis?

Baker, M. "1,500 scientists lift the lid on reproducibility" *Nature* 533: 452-454 (2016).

http://www.nature.com/news/1-500-scientists-lift-the-lid-on-reproducibility-1.19970

doi:10.1038/533452a

# Where do data stewards fit in existing landscape?

- Increase in RDM policies, DMPs, and awareness of best practices.

- Not enough people with knowledge in data stewardship to meet demand.

- Formal training even less.

- The data steward is at the boundary between researchers and support community.



Research community     Support community

Data scientist     Data steward

Knowledge of research discipline

Research Software Engineer     Data librarian

e-Research

Knowledge of data management and curation

Digital preservation

# RDM & the Data Lifecycle

Create

# What is Research Data Management?

"the active management and appraisal of data over the lifecycle of scholarly and scientific interest"

**Data management is part of good research practice.**

# Data creation tips

- Ensure consent forms, licences and agreements don't restrict opportunities to share data.

- Choose appropriate formats.

- Adopt a file naming convention.

- Create metadata and documentation as you go.

# Ask for consent for data sharing

If not, data centres won't be able to accept the data – regardless of any conditions on the original grant.

# Choose appropriate file formats

o Different formats are good for different things.

o *open*, *lossless* formats are more sustainable e.g. rtf, xml, tif, wav.

o proprietary and/or compressed formats are less preservable but are often in widespread use e.g. doc, jpg, mp3.

o One format for analysis then convert to a standard format.

o Data centres may suggest preferred formats for deposit.

| Type of data | Recommended formats | Acceptable formats |
|---|---|---|
| Tabular data with extensive metadata variable labels, code labels, and defined missing values | SPSS portable format (.por) delimited text and command ('setup') file (SPSS, Stata, SAS, etc.) structured text or mark-up file of metadata information, e.g. DDI XML file | proprietary formats of statistical packages: SPSS (.sav), Stata (.dta), MS Access (.mdb/.accdb) |
| Tabular data with minimal metadata column headings, variable names | comma-separated values (.csv) tab-delimited file (.tab) delimited text with SQL data definition statements | delimited text (.txt) with characters not present in data used as delimiters widely-used formats: MS Excel (.xls/.xlsx), MS Access (.mdb/.accdb), dBase (.dbf), OpenDocument Spreadsheet (.ods) |
| Geospatial data vector and raster data | ESRI Shapefile (.shp, .shx, .dbf, .prj, .sbx, .sbn optional) geo-referenced TIFF (.tif, .tfw) CAD data (.dwg) tabular GIS attribute data Geography Markup Language (.gml) | ESRI Geodatabase format (.mdb) MapInfo Interchange Format (.mif) for vector data Keyhole Mark-up Language (.kml) Adobe Illustrator (.ai), CAD data (.dxf or .svg) binary formats of GIS and CAD packages |
| Textual data | Rich Text Format (.rtf) plain text, ASCII (.txt) eXtensible Mark-up Language (.xml) text according to an appropriate Document Type Definition (DTD) or schema | Hypertext Mark-up Language (.html) widely-used formats: MS Word (.doc/.docx) some software-specific formats: NUD*IST, NVivo and ATLAS.ti |
| Image data | TIFF 6.0 uncompressed (.tif) | JPEG (.jpeg, .jpg, .jp2) if original created in this format GIF (.gif) TIFF other versions (.tif, .tiff) RAW image format (.raw) Photoshop files (.psd) BMP (.bmp) PNG (.png) Adobe Portable Document Format (PDF/A, PDF) (.pdf) |
| Audio data | Free Lossless Audio Codec (FLAC) (.flac) | MPEG-1 Audio Layer 3 (.mp3) if original created in this format Audio Interchange File Format (.aif) Waveform Audio Format (.wav) |
| Video data | MPEG-4 (.mp4) OGG video (.ogv, .ogg) motion JPEG 2000 (.mj2) | AVCHD video (.avchd) |
| Documentation and scripts | Rich Text Format (.rtf) PDF/UA, PDF/A or PDF (.pdf) XHTML or HTML (.xhtml, .htm) OpenDocument Text (.odt) | plain text (.txt) widely-used formats: MS Word (.doc/.docx), MS Excel (.xls/.xlsx) XML marked-up text (.xml) according to an appropriate DTD or schema, e.g. XHMTL 1.0 |

# Documentation

Think about what is needed in order to evaluate, understand, and reuse the data.

- Why was the data created?

- Have you documented what you did and how?

- Did you develop code to run analyses? If so, this should be kept and shared too.

- Important to provide wider context for trust.

# What are metadata?

- Metadata
  - Standardised
  - Structured
  - Machine and human readable

- Metadata helps to cite and disambiguate data.

- Documentation aids reuse.

# Metadata standards

These can be general – such as Dublin Core

Or discipline specific:

o Data Documentation Initiative (DDI) – social science

o Ecological Metadata Language (EML) - ecology

o Flexible Image Transport System (FITS) – astronomy

Search for standards in catalogues like:

o http://rd-alliance.github.io/metadata-directory/

o https://rdamsc.dcc.ac.uk/

# Controlled vocabularies

*"MTBLS1: A metabolomic study of urinary changes in type 2 diabetes in……"*



Legend:
- H. sapiens
- Homo Sapien
- Homo sapiens
- homo sapiens
- Homo sapiens (L.)
- Human
- Human
- humans
- sapiens

# ...and ontologies?

o e.g. SNOMED CT (clinical terms) or MeSH

• Defined terms + taxonomy.

o Useful for selecting keywords to tag datasets.

o You can find many ontologies in the BARTOC catalogue and elsewhere.

| ➢ **Organism A** | ► **Organism B** | ❖ **Organism *n*** |
|---|---|---|
| ➢ Term A1 | ► Term A1 | ❖ Term A1 |
| ➢ Term A2 | ► Term A2 | ❖ Term A2 |
| ➢ Term A3 | ► Term A3 | ❖ Term A3 |
|   ➢ Term B1 |   ► Term B1 |   ❖ Term B1 |
|   ➢ Term B2 |   ► Term B2 |   ❖ Term B2 |
| ➢ Term C4 | ► Term C4 | ❖ Term C4 |
| ➢ . | ► . | ❖ . |
| ➢ . | ► . | ❖ . |
| ➢ . | ► . | ❖ . |
| ➢ Term *n* | ► Term *n* | ❖ Term *n* |

# Where will you store the data?

- Your own device (laptop, flash drive, server etc.)
  - And if you lose it? Or it breaks?
- Departmental drives or university servers.
- "Cloud" storage.
- Do they care as much about your data?

**The decision will be based on how sensitive your data are, how robust you need the storage to be, and who needs access to the data and when.**

# Collaborative platforms and third-party tools

- OSF - open platform for sharing data in active phase with fellow researchers and others in secure environment.

- Third-party - commercial (e.g. Dropbox, G Drive, OneDrive) or open source (e.g. ownCloud)



https://osf.io

https://owncloud.org

# Backup vs. preservation

**Backups**

o Used to take periodic snapshots of data in case the current version is destroyed or lost.

o Backups are copies of files stored for short or near-long-term.

o Often performed on a somewhat frequent schedule.

**Archiving**

o Used to preserve data for historical reference or potentially during disasters.

o Archives are usually the final version, stored for long-term, and generally not copied over.

o Often performed at the end of a project or during major milestones.

# How will you allow others to use your data?

Apply licences to disambiguate reuse restrictions.

# Secondary vs primary data

# License research data openly

Try the EUDAT online licence wizard:

https://ufal.github.io/public-license-selector/

CREATIVE COMMONS LICENSES

| | COPY & PUBLISH | ATTRIBUTION REQUIRED | COMMERCIAL USE | MODIFY & ADAPT | CHANGE LICENSE |
|---|---|---|---|---|---|
| PUBLIC DOMAIN | ✓ | ✗ | ✓ | ✓ | ✓ |
| CC BY | ✓ | ✓ | ✓ | ✓ | ✓ |
| CC BY-SA | ✓ | ✓ | ✓ | ✓ | ✗ |
| CC BY-ND | ✓ | ✓ | ✓ | ✗ | ✗ |
| CC BY-NC | ✓ | ✓ | ✗ | ✓ | ✓ |
| CC BY-NC-SA | ✓ | ✓ | ✗ | ✓ | ✗ |
| CC BY-NC-ND | ✓ | ✓ | ✗ | ✗ | ✗ |

✓ You can redistribute (copy, publish, display, communicate, etc.)

✓ You have to attribute the original work

✓ You can use the work commercially

✓ You can modify and adapt the original work

✓ You can choose license type for your adaptations of the work.

# Deposit in a data repository

**Long-term preservation of data.**

# Deposit in a data repository

- The Re3data catalogue can be searched to find a home for data.

- www.fosteropenscience.eu/content/re3data-demo

- Better to use a domain specific repository if available.

- Check they match particular data needs e.g. formats accepted, mixture of Open and Restricted Access.

- Do they assign a persistent and globally unique identifier for sustainable citations and to links back to particular researchers and grants?

- Look for certification as a '*Trustworthy Digital Repository*' with an explicit ambition to keep the data available in long term.

www.re3data.org

# What is a Persistent Identifier (PID)?

*a long-lasting reference to a document, file or other object*

- PIDs come in various forms e.g. ORCID, DOI, ISBN...

- Typically they're actionable i.e. type it into web browser to access.

- Many repositories will assign them on deposit.

- Important for **provenance**.





www.re3data.org

# European perspective...



Final Report and Action Plan from the European Commission Expert Group on FAIR Data

**TURNING FAIR INTO REALITY**

2018

# What FAIR means: 15 principles

## Findable:

**F1.** (meta)data are assigned a globally unique and persistent identifier;

**F2.** data are described with rich metadata;

**F3.** metadata clearly and explicitly include the identifier of the data it describes;

**F4.** (meta)data are registered or indexed in a searchable resource;

## Interoperable:

**I1.** (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.

**I2.** (meta)data use vocabularies that follow FAIR principles;

**I3.** (meta)data include qualified references to other (meta)data;

## Accessible:

**A1.** (meta)data are retrievable by their identifier using a standardized communications protocol;

**A1.1** the protocol is open, free, and universally implementable;

**A1.2.** the protocol allows for an authentication and authorization procedure, where necessary;

**A2.** metadata are accessible, even when the data are no longer available;

## Reusable:

**R1.** meta(data) are richly described with a plurality of accurate and relevant attributes;

**R1.1.** (meta)data are released with a clear and accessible data usage license;

**R1.2.** (meta)data are associated with detailed provenance;

**R1.3.** (meta)data meet domain-relevant community standards;

doi: 10.1038/sdata.2016.18

Comprehensive descriptions can be found at
https://www.go-fair.org/fair-principles/

# Common misconceptions

o FAIR data does not have to be open.

o The principles do not specify particular technologies or implementations e.g. semantic web.

o FAIR is not a standard to be followed or strict criteria – it's a spectrum/continuum.

o It doesn't only apply to the life sciences.

# Increasing that which is FAIR & open

# FAIR ≠ Open

as open as possible, as closed as necessary



Image: 'Balancing rocks' by Viewminder CC-BY-SA-ND www.flickr.com/photos/light_seeker/7780857224

# Bringing together what you've learnt

- Make informed decisions to anticipate and avoid problems.

- Avoid duplication, data loss and security breaches.

- Develop procedures early on for consistency.

- Ensure data are accurate, complete, reliable and secure.

- Save time and effort to make your life easier!

- Useful both to researchers and institutions

Schiermeier, Q. "Data management made simple" *Nature* **555**, 403-405 (2018).
https://www.nature.com/articles/d41586-018-03071-1
doi: 10.1038/d41586-018-03071-1

# DCC Checklist for a DMP

o The DCC assessed existing funder requirements, DMP templates and other best practice to see what should be included in plans. This was synthesised down into common themes and questions.

o 13 questions on what's asked across the board.

o Prompts/pointers to help researchers get started.
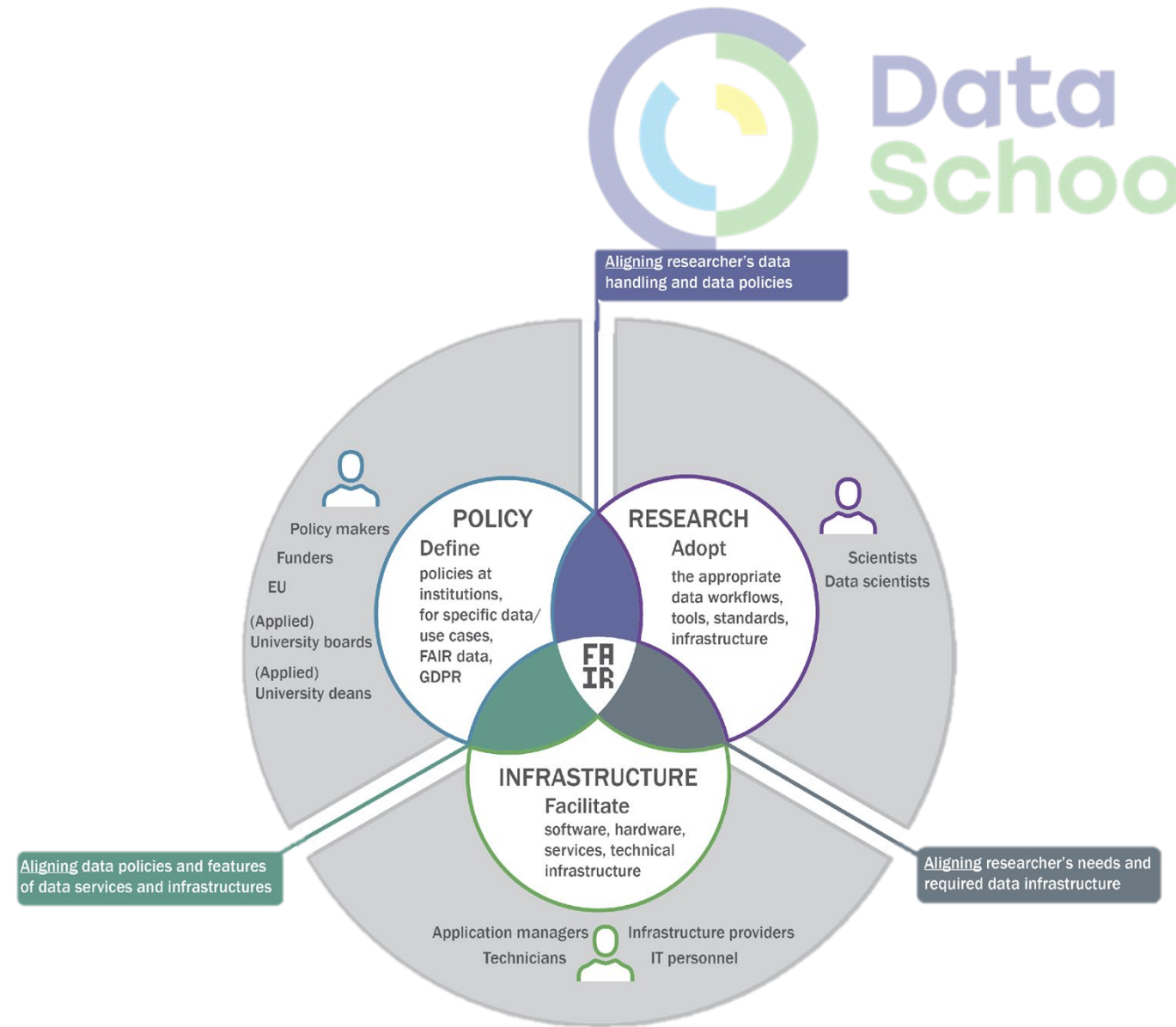
o Guidance on how to answer.

# Example plans

o Plans from several funders and disciplines via DCC www.dcc.ac.uk/resources/data-management-plans/guidance-examples

o Scientific DMPs submitted to the NSF (USA) provided by DataOne https://www.dataone.org/data-management-planning

o DMPs published in RIO journal http://riojournal.com/browse_user_collection_documents.php?collection_id=3&journal_id=17

o Share yours! - www.dcc.ac.uk/share-DMPs

# The different roles of a data steward

- Unlike most other roles, data stewards traverse the researcher-service provider barrier.

- Need knowledge from both perspectives.

- Three key areas requiring data steward training: **Policy**, **Research** and **Infrastructure**.

- **At the heart of these are the FAIR principles.**

- NB. For **Infrastructure**, please refer to e.g. the RISE tool.



Staiger et al. *Data stewards function landscape and its stakeholders*. (2019) Zenodo. http://doi.org/10.5281/zenodo.3460552

# FOSTER Open Science



| What is Open Science? | Best Practice in Open Research | Open Access Publishing | Open Peer Review | Sharing Preprints |
|---|---|---|---|---|
| | | | | |
| Data Protection & Ethics | Open Source Software & Workflows | Managing & Sharing Research Data | Open Science & Innovation | Open Licensing |
| | | | | |

# OpenAIRE

# Research Data Alliance

# Acknowledgements

# Homework exercise

Imagine you are a biologist who is doing microscopy experiments imaging tissue specimens. The data captured by the imaging is 100s of GB in size and is then cleaned and analysed to produce derivatives of the original captured data. Some of these derivatives may eventually be published. In preparation for publication, the data will also be segmented and annotated using standard ontologies. Documentation will also include metadata standards that will sufficiently describe the experimental procedure to allow reproducibility. Publication of the data is mandatory due to funder policy and must be deposited in a repository within 3 years of data production and must use an open licence without restrictions on reuse.

Now…please split into groups and see if you can answer the following questions using the tools and guidelines that have been described:

o What **file format(s)** should data be captured/preserved in?

o Which **metadata standard(s)** should be used?

o What **ontology(ies)** should be used?

o Which **licence(s)** should be used?

o Which **repository** would be the best fit for these data?

o Do you foresee any problems with the data?

o (Hint: not all the questions can be answered definitively! – but why not?)

o **Please use the FAIR Data Forum to post answers and discuss!**