

# Open Science & FAIR

*Introduction to Data Management Practices course*

NBIS DM Team

data@nbis.se

<https://nbisweden.github.io/module-open-science-dm-practices/index.html>



---

Make scientific research and its dissemination  
accessible to all levels of society.

- Open methodology
- **Open source**
- **Open data**
- Open access
- Open peer review
- Open educational resources

**What do you think are reasons for Open Data?**

- Democracy and transparency
  - Publicly funded research data should be accessible to all
  - Published results and conclusions should be possible to check by others
- Research
  - Enables others to combine data, address new questions, and develop new analytical methods
  - Reduce duplication and waste
- Innovation and utilization outside research
  - Public authorities, companies, and private persons outside research can make use of the data
- Citation
  - Citation of data will be a merit for the researcher that produced it



---

*Doing “sloppy” science & not being open and transparent*

Waste of resources

Contributing to the current research credibility crisis

Contributing to the current reproducibility crisis

*Harming the profession*

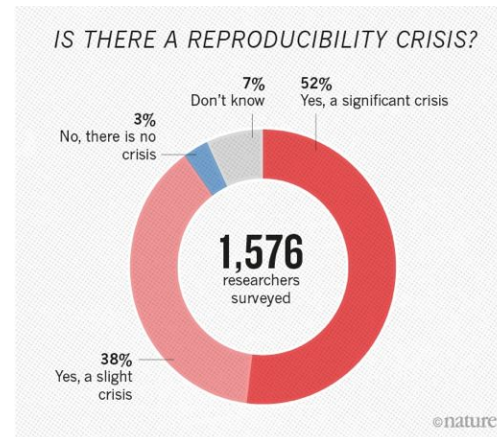
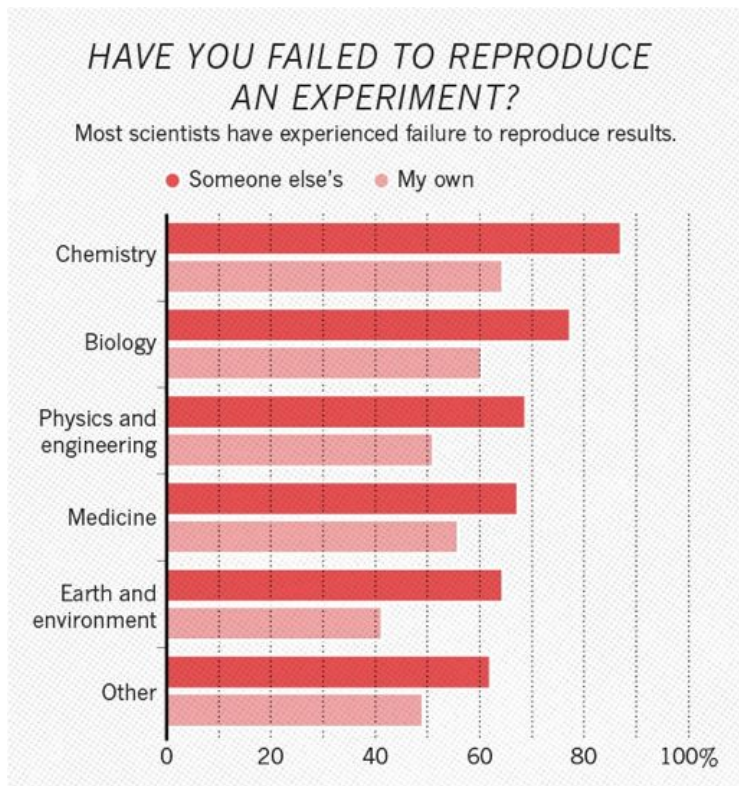
*Harming public trust in research*

My take of material by Rochelle Tractenberg “[Unexpected Ethical Challenges in Bioinformatics and Genomics.](#)”

---

Do you think we have a **credibility** and/or  
**reproducibility** crisis?

If so, what are some of its causes?

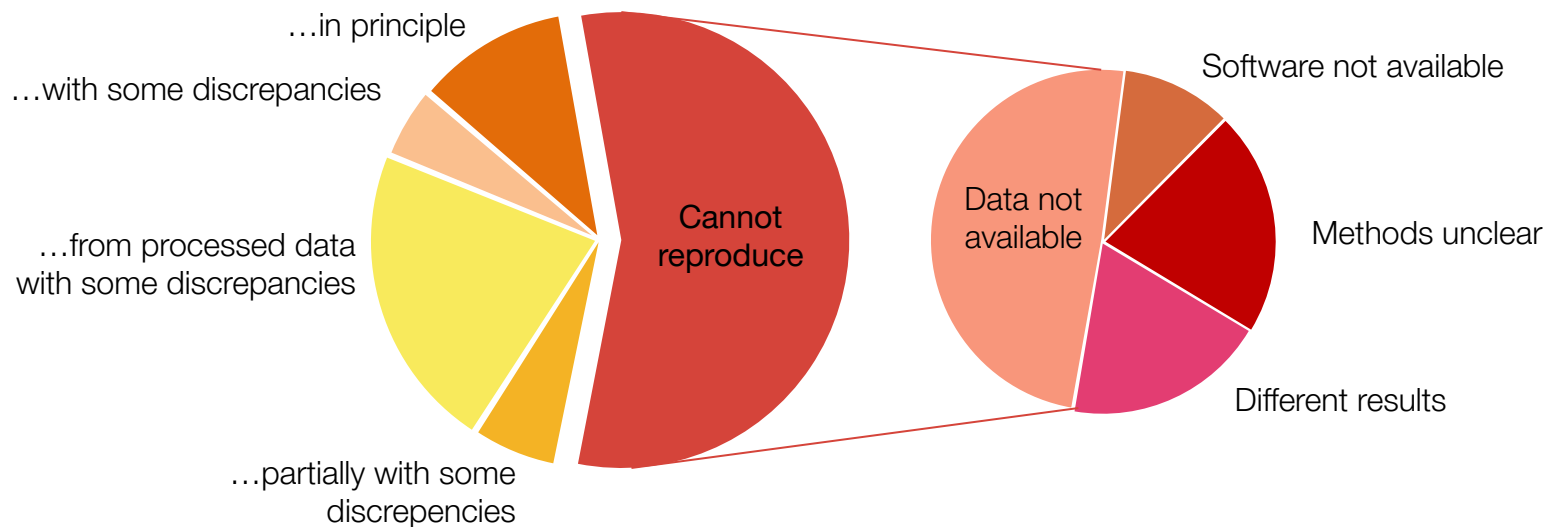


[1] "1,500 scientists lift the lid on reproducibility". *Nature*. 533: 452–454

[2] Begley, C. G.; Ellis, L. M. (2012). "Drug development: Raise standards for preclinical cancer research". *Nature*. 483 (7391): 531–533.

Reproduction of data analyses in 18 articles on microarray-based gene expression profiling published in Nature Genetics in 2005–2006:

Can reproduce...



Summary of the efforts to replicate the published analyses.

Adopted from: Ioannidis et al. Repeatability of published microarray gene expression analyses.

*Nature Genetics* 41 (2009) doi:10.1038/ng.295

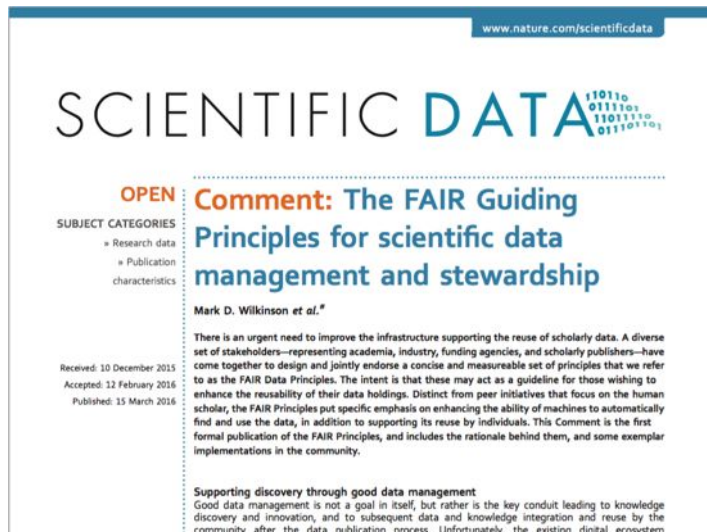




<https://www.youtube.com/watch?v=N2zK3sAtr-4>

- To be useful for others data should be
  - **FAIR** - Findable, Accessible, Interoperable, and Reusable  
*... for both Machines and Humans*

Wilkinson, Mark et al. “The FAIR Guiding Principles for scientific data management and stewardship”. Scientific Data 3, Article number: 160018 (2016) <http://dx.doi.org/10.1038/sdata.2016.18>



## Box 2 | The FAIR Guiding Principles

### To be Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata clearly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

### To be 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

### To be 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

### To be 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

- Data have a **globally unique persistent identifier**
  - *e.g. a DOI, database accession number, etc*
- Data are described by **metadata**
  - *Information that explains the data*
- Data and metadata are findable in a **search resource**
  - *There must be ways of searching for the data*

- Data is retrievable through a **standardised communication protocol** (open, free, allowing authentication & authorisation where necessary)
  - *e.g. http, sftp, etc*
- Metadata are accessible, **even if data is no longer available**
  - *Information about the data can be found even if data is no longer available*

- Metadata use a formal, accessible, shared **language for knowledge representation**
  - *Metadata is available in a form that even a computer can make use of*
- Metadata use **vocabularies** that follow the FAIR principles
  - *Standardised ways of capturing information about the data (that are in themselves FAIR)*
- Metadata include qualified **references** to other metadata
  - *If the data relies on other data, there must be links to those*

- Data have a clear **data usage license**
  - *It is obvious under what conditions the data can be reused*
- Metadata are associated with **detailed provenance**
  - *The metadata is detailed enough to understand for what research questions it is relevant to reuse*
- Metadata meet domain-relevant community **standards**
  - *Metadata is described according to existing standards in the research field*

- 
- Both humans and machines are intended users of data
  - The principles are not necessarily about *open* data
    - “As open as possible, as closed as necessary”
  - FAIRness is not something absolute
    - Different levels of FAIR maturity
  - FAIR does not enforce any particular technical standards

## FAIR at source?



Retroactively?



- **Data Management Plans**, to do your thinking ahead of time
- **Using standard metadata descriptions**, to clearly define your data
- **Organising your analysis**, so you and others can understand what you have done
- **Use versioning control** to keep track of changes you do
- **Clean up metadata and data** to be consistent with the standards you have chosen
- **Submit your data to international public repositories**, so others can find and reuse your data
- **Use scripted analysis of your data**, that can be understood by others

- Strong international movement towards Open Science
- European Commission recommended the member states to establish national guidelines for Open Access
  - Swedish Research Council (VR) submitted proposal to the government Jan 2015
- Research bill 2017–2020 – 28 Nov 2016
  - “*The aim of the government is that all scientific publications that are the result of publicly funded research should be openly accessible as soon as they are published. Likewise, **research data** underlying scientific publications should be **openly accessible** at the time of publication.*” [my translation]
- 2018 – VR assigned by the government to coordinate national efforts to implement open access to research data



## G20 HANGZHOU SUMMIT

· 杭州 2016年9月4-5日

HANGZHOU, CHINA 4-5 SEPTEMBER

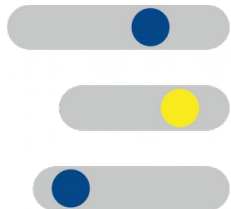
**'We support appropriate efforts to promote open science and facilitate appropriate access to publicly funded research results on findable, accessible, interoperable and reusable (FAIR)'**





## EUROPEAN OPEN SCIENCE CLOUD

The EOSC will offer 1.7 million European **researchers** and 70 million professionals in science, technology, the humanities and social sciences a virtual environment with **open and seamless services for storage, management, analysis and re-use of research data**, across borders and scientific disciplines by federating existing scientific data infrastructures, currently dispersed across disciplines and the EU Member States.



**EOSC**FAIR  
Executive Board Working Group

- [Directive \(EU\) 2019/1024](#) of the European Parliament and of the Council of 20 June 2019 on open data and the re-use of public sector information
- To be implemented into national member state laws by 16 July 2021

*"EU countries must adopt policies and take action to make **publicly funded research data openly available**, following the principle of ‘**open by default**’ and support the dissemination of research data that are findable, accessible, interoperable and reusable (the ‘**FAIR**’ principles)"*

**Funders**  
Data Management Plans  
Open Data

Vetenskapsrådet, FORMAS, Riksbankens Jubileumsfond

[illegible]



- **Data Management Plans**, to do your thinking ahead of time
- **Using standard metadata descriptions**, to clearly define your data
- **Organising your analysis**, so you and others can understand what you have done
- **Use versioning control** to keep track of changes you do
- **Clean up metadata and data** to be consistent with the standards you have chosen
- **Submit your data to international public repositories**, so others can find and reuse your data
- **Use scripted analysis of your data**, that can be understood by others