# Open research data for use and re-use

Ana Alice Baptista
https://www.cienciavitae.pt/201D-B2FC-E126

Braga, July / 2023

Universidade do Minho
Escola de Engenharia

# Contents

- Open Science

- Open Data

- FAIR Principles

- Re-use cases

# Open Science

"An approach to the scientific process that focuses on spreading knowledge as soon as it is available using digital and collaborative technology"

"Open Science." *Research-And-Innovation.ec.europa.eu*, research-and-innovation.ec.europa.eu/strategy/strategy-2020-2024/our-digital-future/open-science_en.

"Knowledge is of all and for all"

"SOBRE CIÊNCIA ABERTA." *Ciencia Aberta*, Ministério da Ciência, Tecnologia e Ensino Superior, 2016, www.ciencia-aberta.pt/sobre-ciencia-aberta. Accessed 16 June 2023.

# Ambitions of the EU's Open Science Policy

- **Open Data**

- European Open Science Cloud

- New generation metrics

- Future of scholarly communication

- Rewards

- Integrity & Reproducibility

- Education & Skills

- Citizen Science

"Open Science." *Research-And-Innovation.ec.europa.eu*, research-and-innovation.ec.europa.eu/strategy/strategy-2020-2024/our-digital-future/open-science_en.

# Open Definition

"Open data and content can be **freely used, modified, and shared** by **anyone** for **any purpose**"

Open Knowledge Foundation. "The Open Definition - Open Definition - Defining Open in Open Data, Open Content and Open Knowledge." *Opendefinition.org*, opendefinition.org/.

"Open research data refers to the data underpinning scientific research results that has no restrictions on its access, enabling anyone to access it"

"Facts and Figures for Open Research Data." *Research-And-Innovation.ec.europa.eu*, Directorate-General for Research and Innovation, research-and-innovation.ec.europa.eu/strategy/strategy-2020-2024/our-digital-future/open-science/open-science-monitor/facts-and-figures-open-research-data_en. Accessed 16 June 2023.
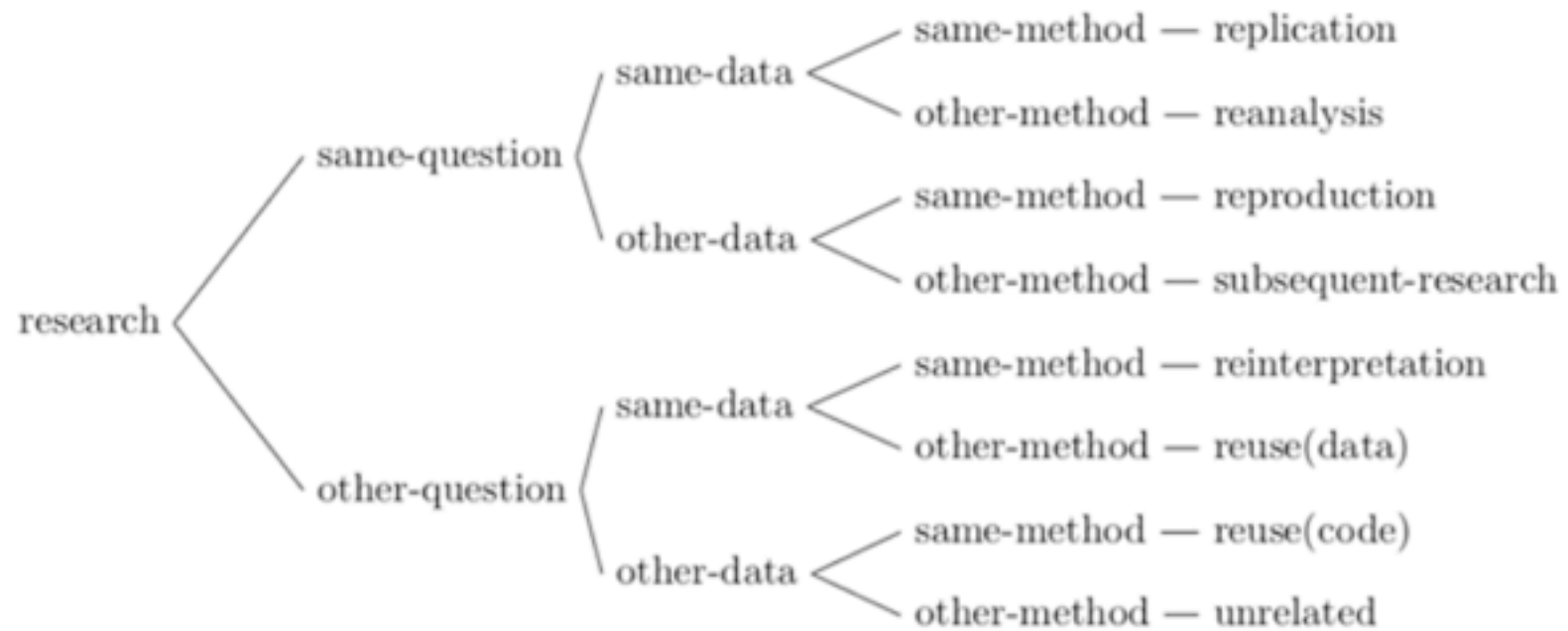
# Why open research data?

- Because it is compulsory (mandates)

- For transparency & reproducibility

- For public scrutiny

- For trustability

- For use & re-use

# What may open research data be used for?

Van de Sandt, S., Dallmeier-Tiessen, S., Lavasa, A. and Petras, V., 2019. The Definition of Reuse. Data Science Journal, 18(1), p.22. DOI: http://doi.org/10.5334/dsj-2019-022

# Why not use Open Research Data?

- Operational:
  - Language
  - Format
  - up-to-dateness
  - Meaning (variable, units of measure,…)
  - License

- Cultural & Political:
  - Replication/reinterpretation/reanalysis studies not so sexy
  - Trust / reliability
  - Culture
  - No explicit rewards

For the data to be used & re-used it has to be made available the right way

# FAIR principles

**F**indable

**A**ccessible

**I**nteroperable

**R**e-usable

# FAIR

**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

# FAIR

**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
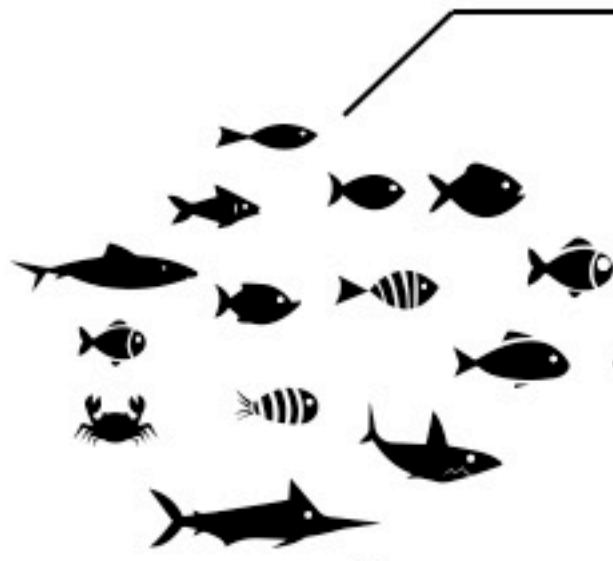
# REUSE cases

# Citizen science in the surveillance and monitoring of mosquito-borne diseases
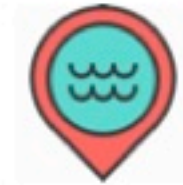
- Cost reduction for mosquito monitoring systems

- Early warning system for mosquito-borne diseases

- Detections further from the known invasion area

European Commission, Directorate-General for Research and Innovation, Switters, J., Osimo, D., *Citizen science in the surveillance and monitoring of mosquito-borne diseases – Open science monitor case study*, Publications Office, 2019, https://data.europa.eu/doi/10.2777/431775

**FishBase**
Raw data released with a cc by-nc licence
(consortium 1)

Institution two processes the data

A third consortium produces AquaMaps

2275 bibliographical citation

Report on fishing stocks discussed at the European Parliament

" FishBase is a digital catalogue of fishes, collecting a variety of data on 34,300 fish species, such as geographical distribution, biometrics, habitats, population dynamics as well as reproductive, metabolic and genetic data."

to know more

# Open research data for use and re-use

Ana Alice Baptista
https://www.cienciavitae.pt/201D-B2FC-E126

Braga, July / 2023



Universidade do Minho
Escola de Engenharia