# The Alan Turing Institute

*The Turing Way* **Workshop:**
Reproducible, Open and
FAIR Research

## Session 1

Emma Karoune & Malvika Sharan
Pronouns: she/her/hers

# Emma Karoune


Historic England


The Alan Turing Institute


The Turing Way

- **Archaeobotanist/Palaeoecologist**

  - FAIR Phytoliths project
  - Open reference collections


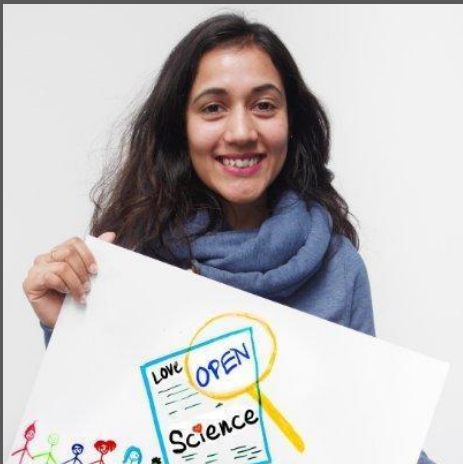International Phytolith Society

- **Open Researcher & Community Manager**
  - Tools, practices and systems programme
    - Turing Way
    - DECOVID/ Turing-RSS Lab
  - SSI Fellow/ UK-Elixir FAIR data Fellow


Software Sustainability Institute

Link to Emma's SSI Fellows page

- PhD in Bioinformatics

- Open Access & Open Source research publications

- Computational and Open Science skill training (2015-)

- Community Building in Open Science (2016-)

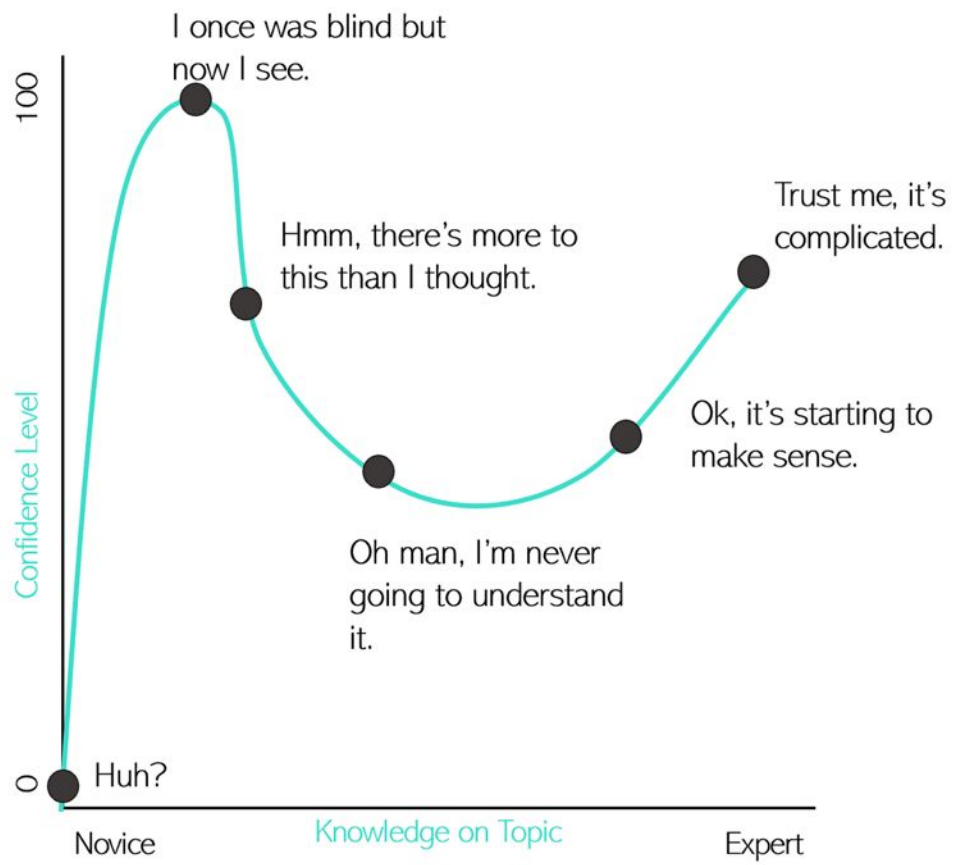- Co-lead of *The Turing Way* & Open Life Science (2019-)

Senior Researcher

Tools, Practices and Systems

The Alan Turing Institute, UK

# Learning Objectives

– **<u>Problem</u>:** Scientific errors have real world effect

– Define what reproducible research is

– Understand what open science practices are

– Identify FAIR principles for your research

– Differentiate between FAIR and open data/research

– Learn how to implement a reproducible workflow

# Disclaimer:

You probably already know all about it!

Adapted from: Squad. (2018, December 13). Dunning-Kruger Effect: Definition, Test, Examples & Quiz. Science Terms. https://scienceterms.net/psychology/dunning-kruger-effect/

# Scientific errors have real world effects

*Researchers have an obligation to consider the ethical standards (right actions) and their impact on society.*

# Scientific errors have real world effects

# Scientific errors have real world effects

"We will redouble our efforts to avoid such errors in the future …" - Authors

"mistaken way to examine (country specific) data provides an intellectual rationalisation for things that affect how people think about the world."
- Daniel Hamermesh

https://www.gapminder.org/, Rosling, H., Rönnlund, A. R., & Rosling, O. (2018). Factfulness: Ten Reasons We're Wrong About the World--and Why Things Are Better Than You Think. https://www.bbc.co.uk/news/magazine-22223190

# Price of Popularity

Researchers in Australia used unreliable COVID-19 data and misreported that an anti-parasite drug could stop the virus from replicating in cells.



*In a desperate attempt to save dying patients, doctors began justifying the drug's use against COVID-19 as the virus spread aggressively throughout Latin America.*

*Reardon, S. (2021). Flawed ivermectin preprint highlights challenges of COVID drug studies. Nature, 596, 173–174. doi: 10.1038/d41586-021-02081-w. Mega, E. R. (2020). Latin America's embrace of unproven COVID treatment hinders drug trials. Nature, 586, 481–482. doi: 10.1038/d41586-020-02958-2,* Slides under DOI: 10.5281/zenodo.6337939

# Post-publication Peer Review is Important, But …

*To consult the [experts] after an experiment is* finished is often merely to ask to conduct a *post mortem examination.* […] can perhaps say what the experiment died of. - Ronald Fisher

# Avoid Errors Before the Harm Occurs

1.  Mistaken research design or analysis processes

    ○ **Project design for open, FAIR & reproducible research**

# Avoid Errors Before the Harm Occurs

1. Mistaken research design or analysis processes

   ○ **Project design for open, FAIR & reproducible research**

2. Wrong choices of tools and methods

   ○ **Data handling, data management, collaboration process**

# Avoid Errors Before the Harm Occurs

1. Mistaken research design or analysis processes

   ○ **Project design for open, FAIR & reproducible research**

2. Wrong choices of tools and methods

   ○ **Data handling, data management, collaboration process**

3. Inappropriate baseline comparison

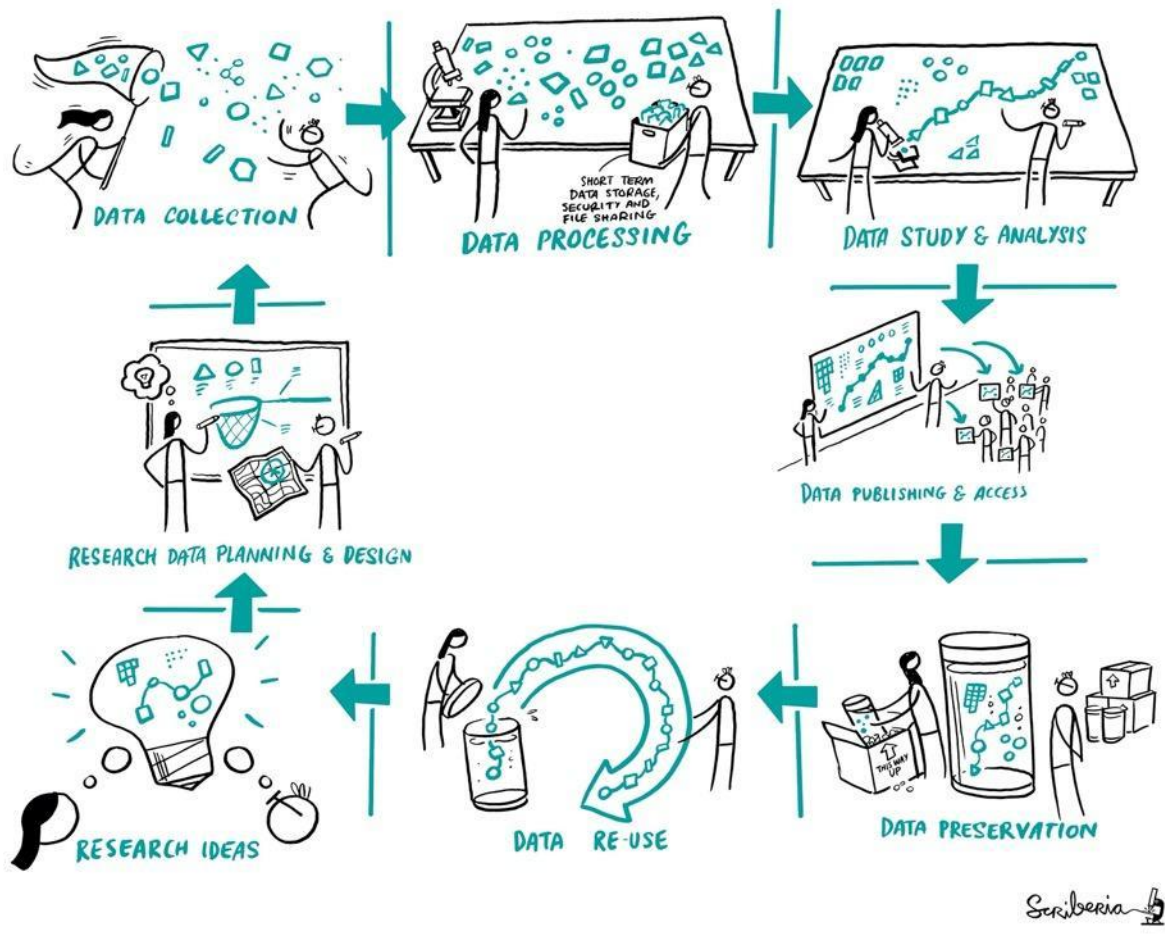   ○ **Lack of technical understanding (we won't discuss this!)**

# Session 1 - Reproducible, Open and FAIR research

- What is reproducible research

- What is open research

- Concerns about opening up research

- Understanding FAIR data/research

- Making your research FAIR

- Setting up repositories and working collaboratively

# Research Reproducibility

*Reproducible research save valuable time in verifying and building upon existing solutions.*

|  |  | Data | |
|---|---|---|---|
|  |  | Same | Different |
| **Analysis** | Same | Reproducible | Replicable |
|  | Different | Robust | Generalisable |

# Reproducible Research Workflows

| | | Data | |
|---|---|---|---|
| | | Same | Different |
| **Analysis** | Same | Reproducible | Replicable |
| | Different | Robust | Generalisable |

# Why do reproducible research?

Applying best practices in our research requires intention, resources, time and collaboration, which can be overwhelming.

# slido

What are some barriers to reproducibility?

ⓘ Start presenting to display the poll results on this slide.

Barriers to reproducible research

Held to higher standards than others

Is not considered for promotion

Publication bias towards novel findings

Requires additional skills

Plead the 5th

Support additional users

Takes time

# *The Turing Way*

An Open Source project that involves and supports its

**diverse community** to make data science **reproducible,**

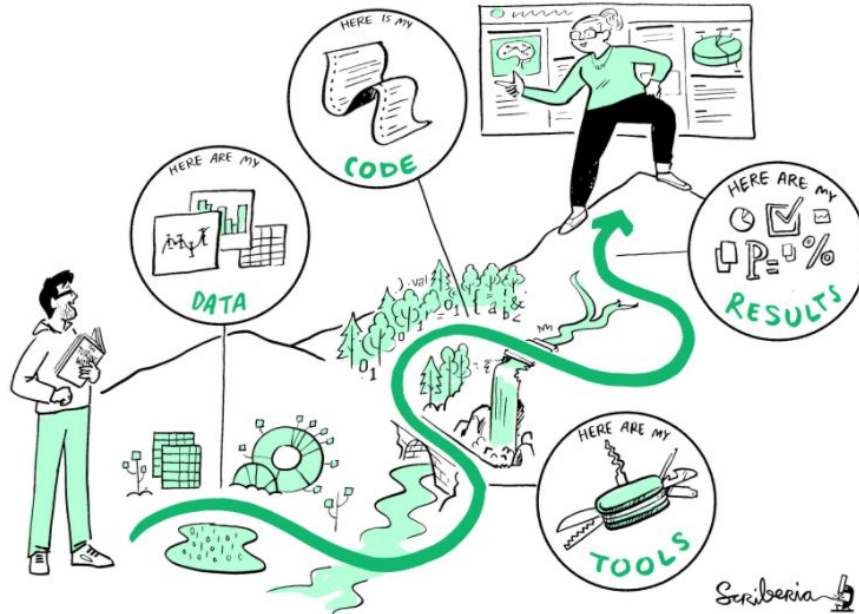**ethical, collaborative and inclusive** for you.

# Guide for Reproducible Research

*This guide covers topics related to skills, tools and best practices for research reproducibility.*

*The Turing Way* defines reproducibility in data research as data and code being available to fully rerun the analysis.

There are several definitions of reproducibility in use, and we discuss these in more detail in the Definitions of Reproducibility section of this chapter. While it it absolutely fine for us each to use different words, it will be useful for you to know how *The Turing Way* defines *reproducibility* to avoid misunderstandings when reading the rest of the handbook.

# A book: "Work in Progress"

# Reproducible Research Workflows



## Inner Source vs Open Source

# Privacy and Sensitivity concerns

# Reflection Exercise 🤔

– What are your motivations to share your work?
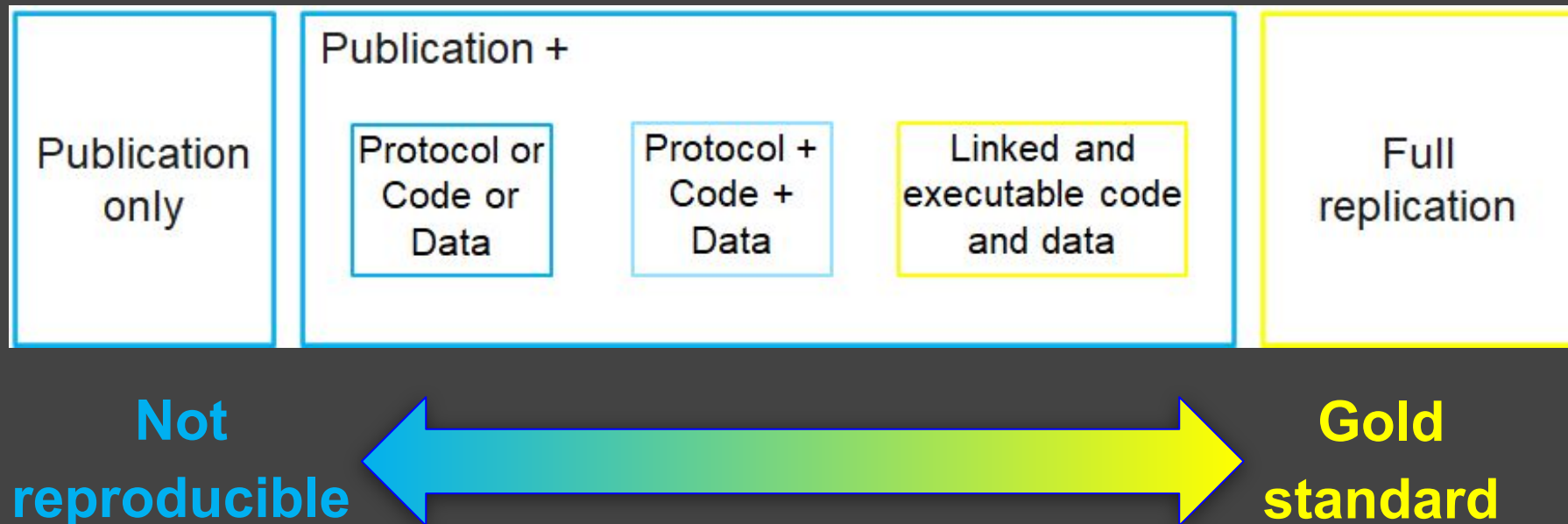
– What are your concerns about sharing your work?

# Breakout Room: 15 Minutes

What do you need to consider to ensure that your work is reproducible?

– Research object: (1) Data, (2) code, (3) workflow, (4) documentation

– Practices: Version control, licensing, data management, communication, collaboration, reusability, long-term archiving

– Instructions: (i) Each room is assigned a research object. (ii) Nominate a notetaker in your room. (iii) Select 1 or 2 practices. (iv) Discuss at what stages of research cycle are these practices applied. (v) Report!

# Reproducible Research Spectrum



*Adapted from Peng 2011*

# Reproducibility: Where should we start?
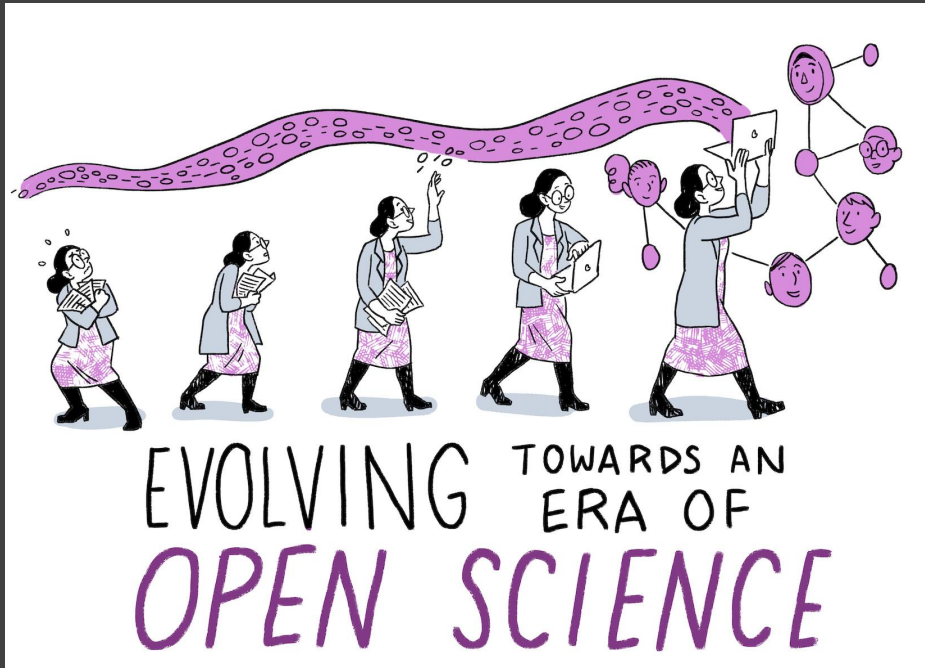
Reproducible research workflow:

1. Collaboration

- Using collaborative, version controlled and open ways to work with others.

2. Transparency

- Clear documentation of methods, data, code - openly shared research compendium.

# **Open Science** to enable Collaboration and Transparency



- Open Science practices remove barriers from sharing and using scientific resources at all stages of research.

- Make research findings accessible to all rather than keeping them locked away (for example, behind a paywall).

Open Science aims to transform research by making it more transparent → accessible → reliable → reproducible → reusable → collaborative → beneficial to society.
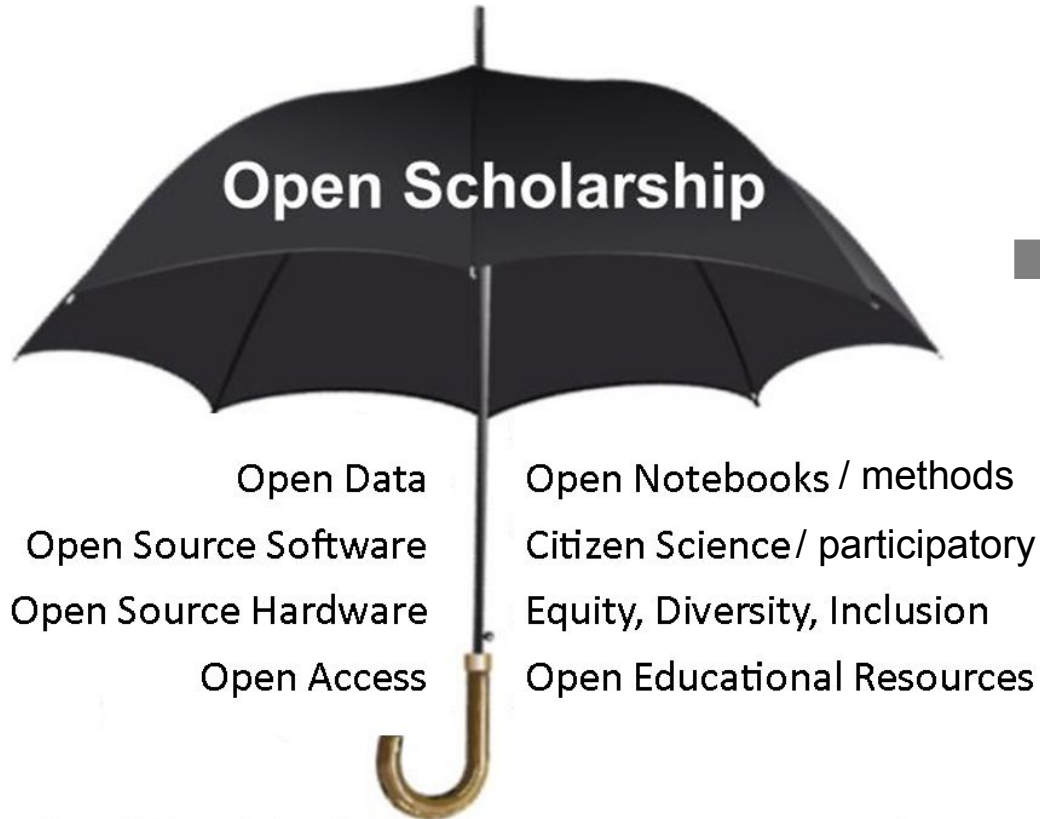
To achieve this openness in research, we need to make each element of the research process ...

– *publicly available,*
– *with permission to view, use, modify and distribute, and*
– *description for how one can collaborate.*

# Open Science is an umbrella term for open research practices

# What are some barriers we should work to remove?



Open Scholarship

Open Data          Open Notebooks / methods
Open Source Software   Citizen Science / participatory
Open Source Hardware   Equity, Diversity, Inclusion
Open Access         Open Educational Resources

Transparency →
Reproducibility →
Research Quality →
Sustainability

**+**

Collaboration →
Inclusive Research →
Equity and Diversity →
Global Accessibility

.meetup.com/Berlin-Open-Science-Meetup/          Robin Champieux and Danielle Robinson

Truegif.com

TOO MANY CHOICES

Open Science

Open Research

Open Scholarship

We are going to use these terms interchangeably today to cover important considerations for you!
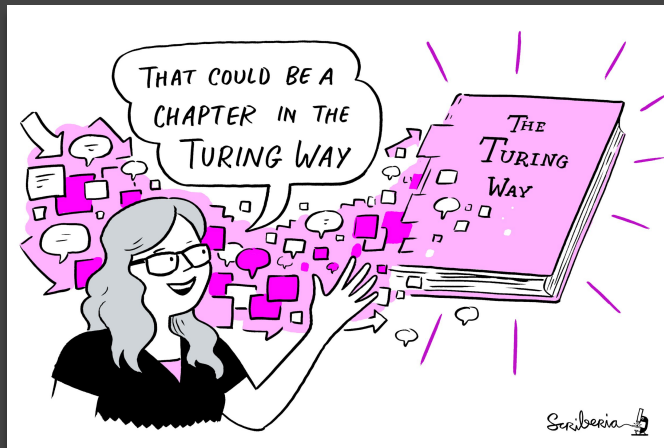
**Book:**
the- turing-way.netlify.app/

**GitHub:**
github.com/alan-turing-
institute/the- turing-way

**Twitter:**
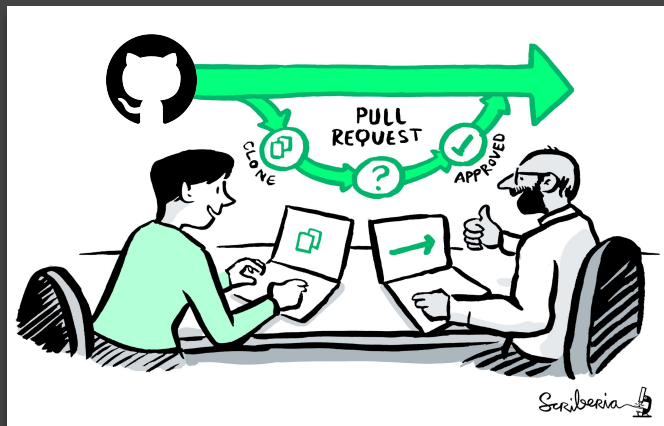twitter.com/turingway

**Email:**
theturingway@gmail.com

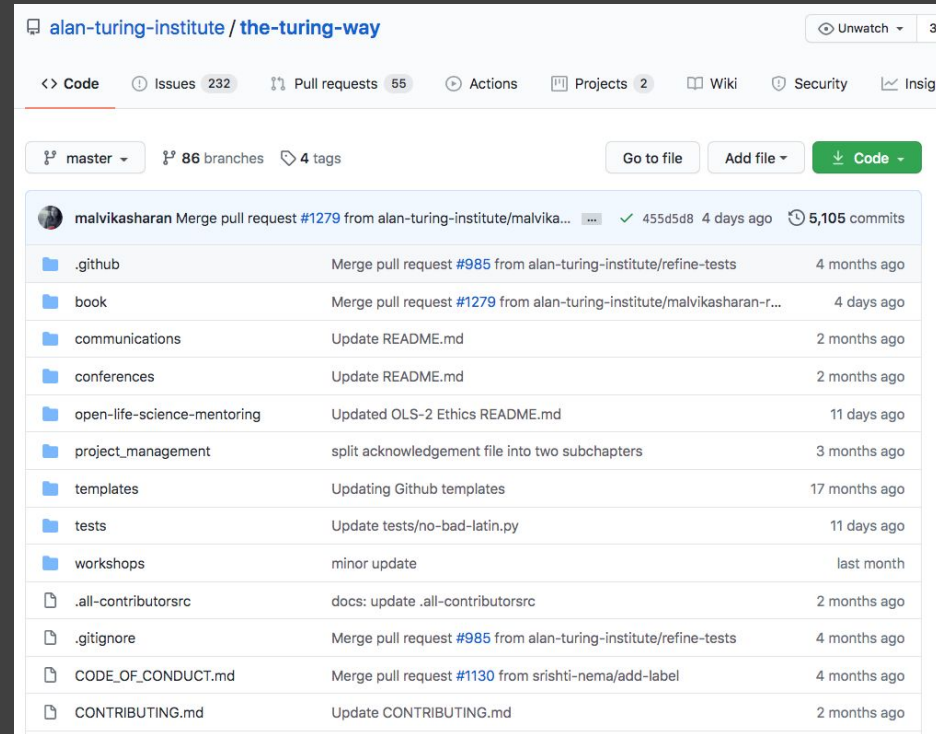A Book

A Community

An Open Source Project

A Culture of Collaboration

# A Collaborative, Version Controlled and Open Project

– Shared online

– Designed for open collaboration

– Hosted on GitHub with history and versions

– Described with open license

– Community oriented

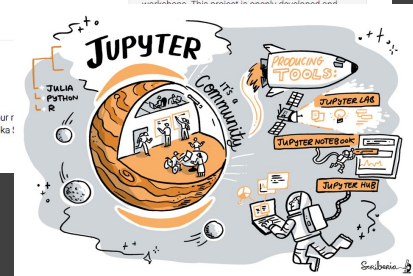| 📁 .github | Remove prettier configuration | |
| 📁 book | minor update | |
| 📁 communications | Fix typos | |
| 📁 conferences | Add KW formatting pedantry | |
| 📁 project_management | Update online-collaboration-cafe.md | |
| 📁 templates | Updating Github templates | |
| 📁 tests | Add "et cetera" as a deprecated Latinism | |
| 📁 workshops | Remove mis-pasted text | |
| 📄 .all-contributorsrc | Merge pull request #991 from alan-turing-institute/all-contributors/a... | 5 days ago |
| 📄 .gitignore | ignore pptx in workshop folder | 9 months ago |
| 📄 .travis.yml | add html-proof file again | last month |
| 📄 CODE_OF_CONDUCT.md | her -> their | 6 months ago |
| 📄 CONTRIBUTING.md | Update CONTRIBUTING.md | 2 months ago |
| 📄 GOVERNANCE.md | Read through months later | 5 months ago |
| 📄 LICENSE.md | Fix typo in licence | 2 months ago |
| 📄 README.md | Merge pull request #991 from alan-turing-institute/all-contributors/a... | 5 days ago |
| 📄 book_skeleton.md | Update book_skeleton.md | 13 months ago |
| 📄 contributors.md | Add myself to contributors.md | 11 months ago |
| 📄 tips_and_tricks_survey.md | Update tips_and_tricks_survey.md | 14 months ago |
| 📄 ways_of_working.md | Adjust team contact section | 5 months ago |


WELCOME!

# An Open Science project

– everyone can freely read, reuse, distribute, modify and help develop

– the project belongs to the research community (CC-BY license)

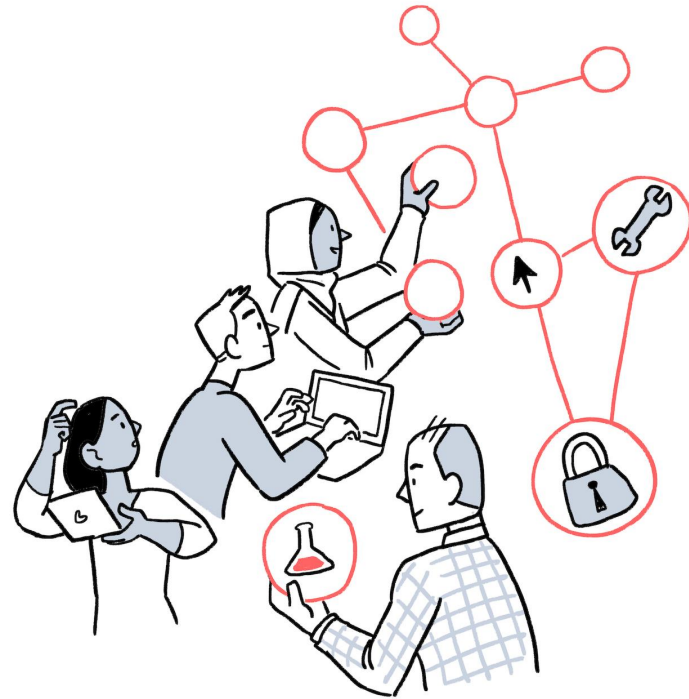– Builds in collaboration with other projects

🌻 Open Science:
Where should we start?

# FAIR principles

- **F**indable,
- **A**ccessible
- **I**nteroperable
- **R**eusable

FAIR principles from **Wilkinson *et al*. (2016)**
**DOI: 10.1038/sdata.2016.18**

# FAIR data analogy



Annotation makes it easier to find important things

You would not buy food with no labels!

Labels make different foods easier to find and access in stores, combine with other foods (interoperable) and use in different ways.

Adapted from talk by Philippe Rocca-Serra (2020)

# What is the meaning of  F A I R  data ?

**F = Findable**

*in an **online** data catalogue / archive / portal*

findable by **humans** and by **machines**

- **ENA** for DNA sequences
- **GBif** and **OBIS** for biodiversity data
- **BioImage Archive** for images of biological material
- **Zenodo** as a a general-purpose open-access repository

**Standardised** and **rich** discovery **Metadata** explaining:

☑ **Who**:  is the **author** / **contact person** for questions

☑ **How**:  were the data created  --> **procedures** / **protocols**

☑ **How**:  to **access** the data, consider **licenses**

☑ **What**: **keywords** describe the data

☑ **What**: **parameters** were measured, **species** & **geography** covered

☑ **When**:  were the **data** and **updates** created

# What is the meaning of F A I R data ?



**A = Accessible**   **Data & Metadata**

*from catalogue/archive/portal*

*via **machine to machine** and **human interfaces***

- **Web interfaces** for human searches & downloads
- **APIs** for searching & accessing
- Clear **instructions** for access (download, request access,.)
- **Keeping metadata** when data is deleted
- **Metadata update** when updating data / information
- **All data levels** should be archived: raw data is the most important and at a minimum must be provided

# What is the meaning of F A I R data ?



HOW DO YOU OPEN A .XZQ FILE?

INTEROPERABLE

## I = Interoperable
*readable & understandable by **humans / code** :*

❑ **Community-accepted** data formats & file types
  ○ open (non-proprietary)
  ○ sustainable (think in 10 years from now)

❑ **Clear**, **controlled vocabulary** for data & metadata
  ○ describing all relevant terms/values/units
  ○ specific → data/metadata "dictionary"

❑ Your data should be **standalone**, packaged up with
  ○ all necessary information and files to allow the data to
    be understood by anyone at any time

❑ **Readable** by code:
  ○ machine readable descriptions of data: files and format

# What is the meaning of F A I R data ?



**R = Re-usable**

Know *how I can **trust, repeat, re-analyse, re-use*** the data.
Necessary to provide:



time         year 20XX

❑ Data **usage licence** --> full terms & conditions

❑ Data **provenance** --> metadata and information on:

  ○ every data life-cycle stage
  ○ documentation / protocols / references
  ○ link to accompanying data and publications
  ○ instruments & software used

❑ **Relationship** between the different levels of data you provide is documented:

**raw**--> **quality controlled** -->**processed**-->**published**

# What is a data life-cycle?
# Data provenance



## Data Life-Cycle:

- *covers the entire period of time over which data exists*
- *encompasses all the stages: **first Capture** → **data re-use***

1. **Sample acquisition** → raw data:
   sample preparation, experimental settings / parameters, raw data acquisition

2. Data **quality control**: checking and updating of collected data
   documentation of QC procedures

3. Data **processing** & **analysis**: guided by scientific question
   documentation of processing steps, analysis methodology

4. **Archiving** & **publication**: data are placed in an online catalogue
   discovery metadata, provenance metadata, provenance files, references and links

5. Data **dissemination** / **integration**: adding data to well known portals, brokers

6. Data **reuse**: only possible with sufficient provenance information!

# Breakout Room: FAIR Reflection

*What part of your research lifecycle/outputs can be made open?*

*What you will have to do in your project to maintain "this" aspect of FAIR.*

- Room 1: Findable

- Room 2: Accessible

- Room 3: Interoperable

- Room 4: Reusable

# Reproducible research doesn't always mean open

- Reproducibility can be facilitated by open, but open is a choice

- Reproducibility needs to be considered at all stages

- As open as possible, as closed as necessary

  - Open principles should be applied when you can

  - NEVER for private, confidential or sensitive data

- Always apply FAIR (**F**indable, **A**ccessible, **I**nteroperable, **R**eusable)

# Reproducibility

- Is my code **correct**?

- Can others **read and test** it?

- Is my **workflow robust**?

- Have I provided **guidance**?

- Is my work **citable**?

*Reproducibility should ensure higher scientific standards allowing others to test and reuse your work ...*

# Open Source

- Is my code **freely available**?

- Can others **modify and share** it?

- Is my **workflow reusable**?

- Have I provided **permission**?

- Is my work **open for collaboration**?

*... and Open Source should allow anyone to reuse, report errors, fix issues, build on and collaborate*

# FAIR doesn't need to be open

- FAIR does not require data to be open

- FAIR requires open metadata

- Detailed information about research/data should be open

- FAIR applies open standards for interoperability

**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

Metadata: information about the "data descriptors" that facilitate cataloguing data and data discovery

# Is this Reproducible, Open or FAIR?



- TARO is a PhD researcher in a health data science team

- Collaborates with people in their team of engineers and data scientists

- They have written Python code for a commonly used dataset

- Their supervisor suggested them to publish their work online

- Created a public repository to share their code, data and documentation

- Sent out an email to their team members to use their code

The question is not "Should I share my work?", but "How can my work benefit other collaborators?"
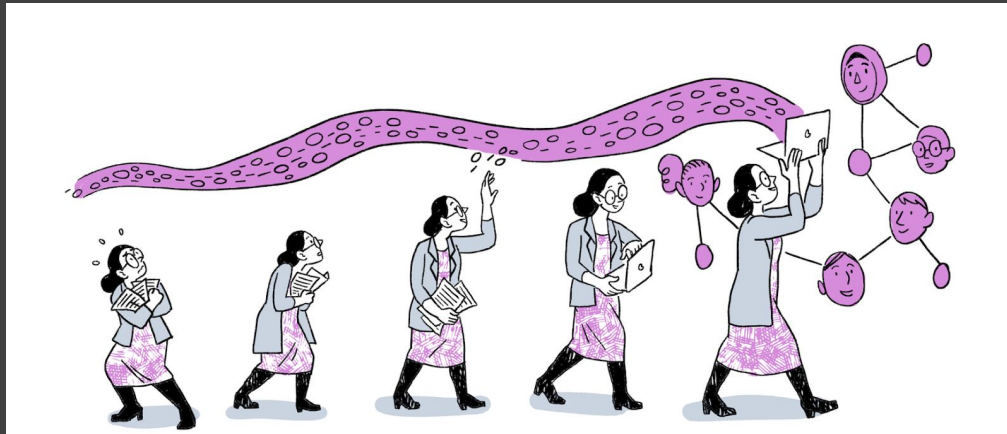
You are your number one collaborator!

Good practices benefit (future) you!

# Open Science path is self defined!



- Reproducibility is essential.
- Open Science involves many concepts - take the steps that are possible.
- FAIR provides guidance for actionable steps that make your work reproducible.
- *Ask for feedback and help whenever you can.*

# Reproducible & Open Science Projects

Online Repository

License/Permission

Documentation
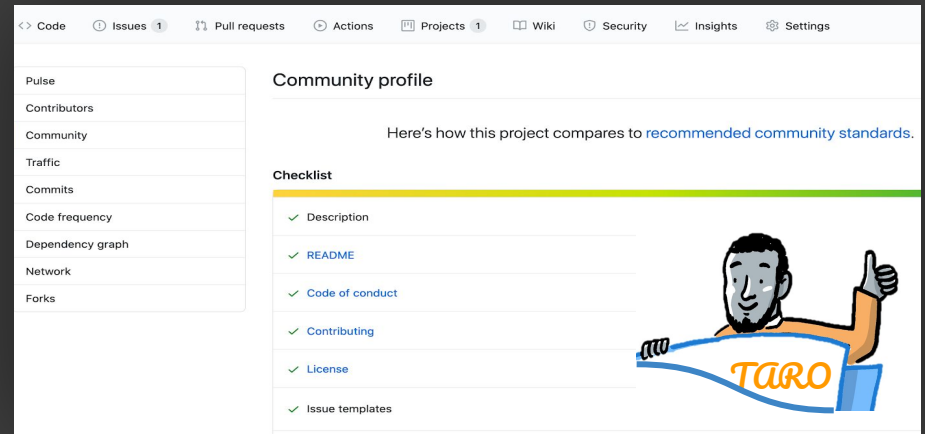
Version Control

Dependencies and data

Review & Verify

Report/Fix Bugs

Release & Cite

# Ensure Use and Reuse



- Create a project repository

- Create a README file with information on their scripts

- Add an open license in their repository for reuse

- Make it easy to test

- Add a minimum guideline for reporting errors

- Release citable versions of documentation/code/data

# READMEs for Open and Collaborative Projects

**Motivation:**

Learn how to communicate your project effectively.

**Method**:

Write clear description of the project in README file.

Why do READMEs matter?

# What is a README?

- Found in the root directory of your repository

- In ALL CAPS, a request for all to "read me!"

- First stop for your collaborators

*Could also be: website landing page, list of bullet points*

# In your README, show:

- what you're doing, what's your motivation

- what makes your project special and exciting

- who your collaborators are, how do they contributors

- where are the key resources, where to find them

# A closer look at a README



**STEMM Role Models App**

*Inspire future generations by providing the most exciting and diverse speakers for your conference.*

chat on gitter

## Welcome!

First and foremost, Welcome! 🎉 Willkommen! 🎊 Bienvenue! 🎈🎈🎈

Thank you for visiting the STEMM Role Models app project repository.

This document (the README file) is a hub to give you some information about the project. Jump straight sections below, or just scroll down to find out more.

- What are we doing? (And why?)
- Who are we?
- What do we need?
- How can you get involved?
- Get in touch
- Find out more
- Understand the jargon

README example: <u>STEMM Role Models App</u>

- ● Welcome message!

- ● Project description & vision

- ● How to:
  - ○ Test and verify
  - ○ Fix errors
  - ○ use (license)
  - ○ Get involved
  - ○ Report issues

# README = Project Documentation



- project description and features

- installation instructions

- tutorials and requirements

- how to run associated tests

- list of authors/contributors

- contact information

- links to related material

# Assignment: Create a project repository

Add top-level files: README and LICENSE

- README to communicate about your work
- Also try to add a License - to allow others to use, modify, build upon your work
    - We will cover licenses in the next session in more detail.

Use The Turing Way chapter for README to guide your assignment

https://the-turing-way.netlify.app/project-design/project-repo/project-repo-readme.html

# Create a GitHub Repository with README file

## Create a new repository

A repository contains all the files for your project, including the revision history.

**Initialize this repository with:**
Skip this step if you're importing an existing repository.

☐ **Add a README file**
This is where you can write a long description for your project. Learn more.

☐ **Add .gitignore**
Choose which files not to track from a list of templates. Learn more.

> ℹ️ **Note**
>
> **Three lessons about README**
>
> - Know your users and what they need
> - Get users doing powerful things quickly
> - Watch out for jargon!
>
> Source: Hao Ye. (2021, March). Collaborations Workshop 2021 Mini-Workshop: README tips to make your project more approachable (Version v1.0.0). Zenodo. http://doi.org/10.5281/zenodo.4647391

# End of Part 1: Further Reading and Examples

- *The Turing Way* README

- Open Life Science README

- Purple Booth's README Template

- Thoughtbot's Blog on How to Write a Good README

- Matias Singer's curated List of Awesome READMES

# Assignment: Create a project repository

Add top-level files: README and LICENSE

- – README to communicate about your work
- – Also try to add a License - to allow others to use, modify, build upon your work
    - We will cover licenses in the next session in more detail.

Use The Turing Way chapter for README to guide your assignment

https://the-turing-way.netlify.app/project-design/project-repo/project-repo-readme.html

# The
## Alan Turing
# Institute

---

**End of Session 1**

Several slides in this presentation uses *The Turing Way* and Open Life Science contents shared under CC-BY 4.0 License.