



REPRODUCIBILITY AND REPLICABILITY: OPEN SCIENCE ON TRUSTED DATA

Christine Choirat | Swiss Data Science Center

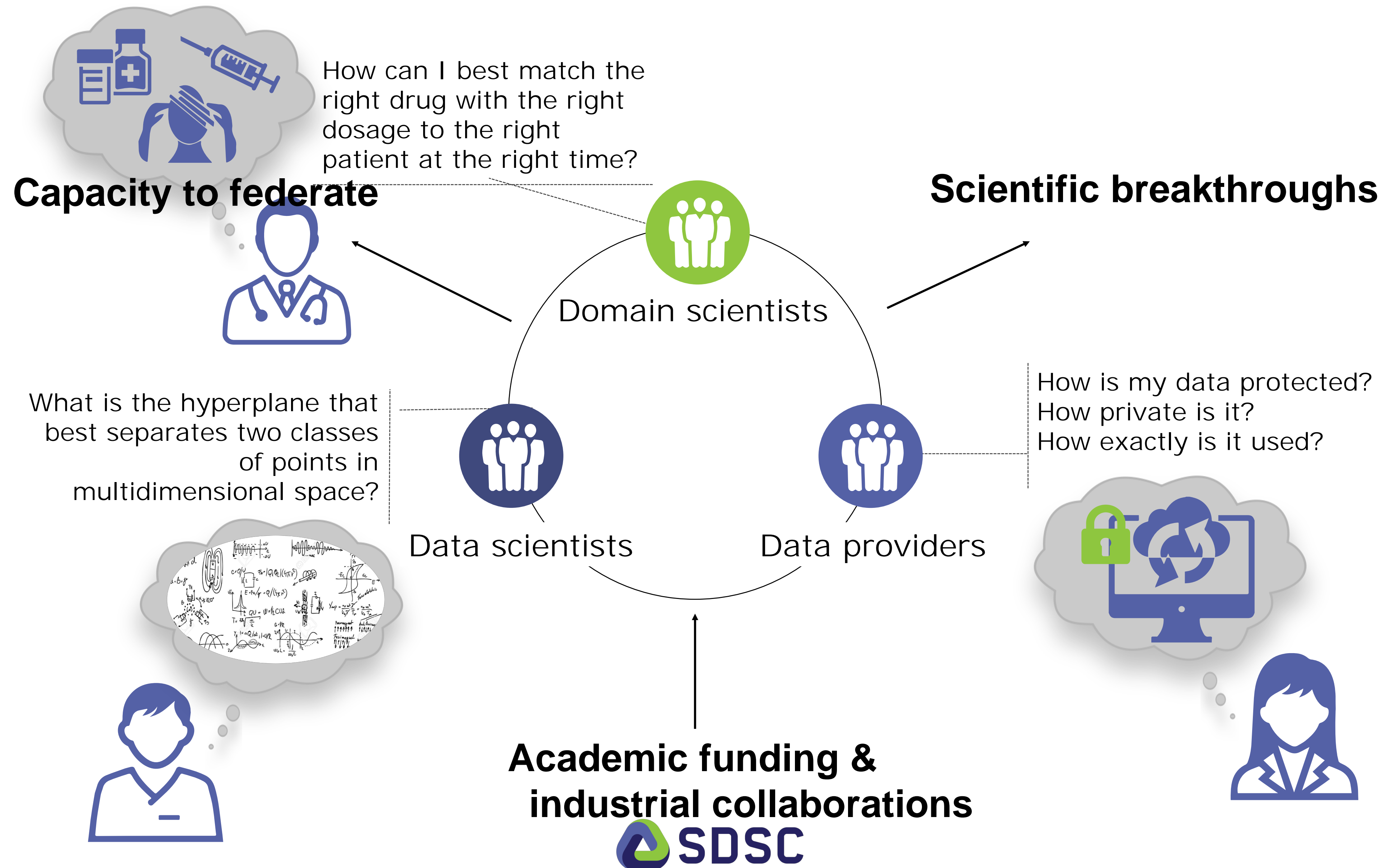


<https://datascience.ch/>
[@SDSCdatascience](#)



FOREWORD: SWISS DATA SCIENCE CENTER

Closing the gaps in the data science journey



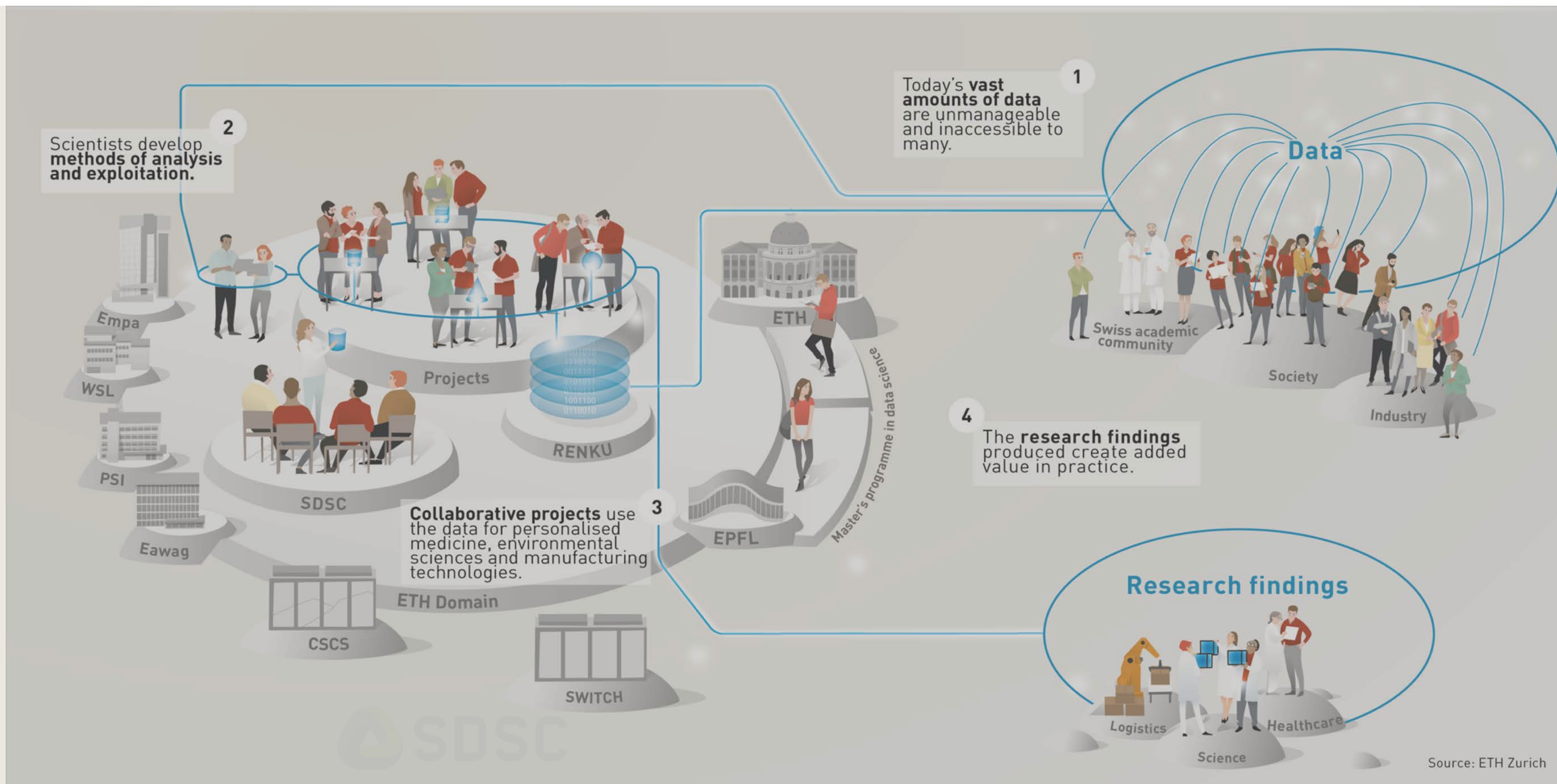
A multinational team of professionals

Joint venture between EPFL and ETH Zurich, with offices in Lausanne and Zurich
Fully operational since January 2017, growing to 50+ full-time professionals
Mission: accelerate the adoption of data science in academia and industry

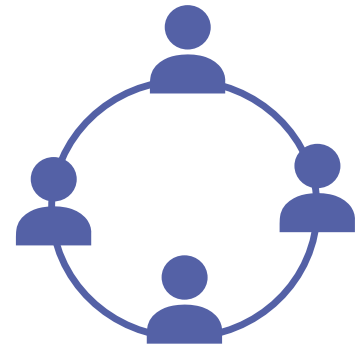


Uniquely positioned in a broad landscape

- Very effective collaborations between two world-renowned engineering schools
 - Dedicated team of scientists to bridge the various gaps in the data science journey
- ... enabling adoption of data science **at a scale only possible when pooling resources**



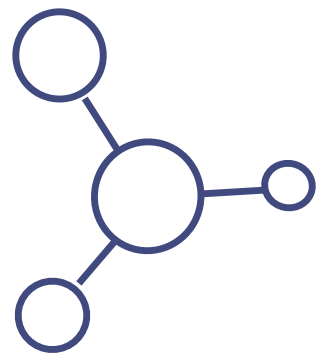
Main focuses



Embedded R&D
collaboration

Academic and industry research collaborations

- 30 academic projects in personalized health and environmental science
- Several signed industry partnerships
- Initiating discussions at the international level



Insights as a Service

RENKU, the SDSC analytics platform (Open Source)

- Facilitates multi-disciplinary collaborations in data science & AI
- Promotes reusability / reproducibility of science (respects FAIR principles)
- Extremely positive response from academia and industry



Data Custodian, a multisided platform to establish trust and transparency in data usage

- A data vault and secure multiparty compute ecosystems (reference architecture lead)
- Enables cooperation between mutually non-trusting parties (trusted intermediary)
- Promotes a shift from data ownership to the ownership of its use (preserve data sovereignty)



Educational Services

Contribution to the education in Data Science

- Involved in EPFL and ETH Zurich Master's Degrees in Data Science
- Contributes to continuous education through EPFL / HEC Lausanne CAS and ETH Zurich DAS programs
- Coaching programs for industry partners





REPRODUCIBILITY “CRISIS”

PR situation to avoid...

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

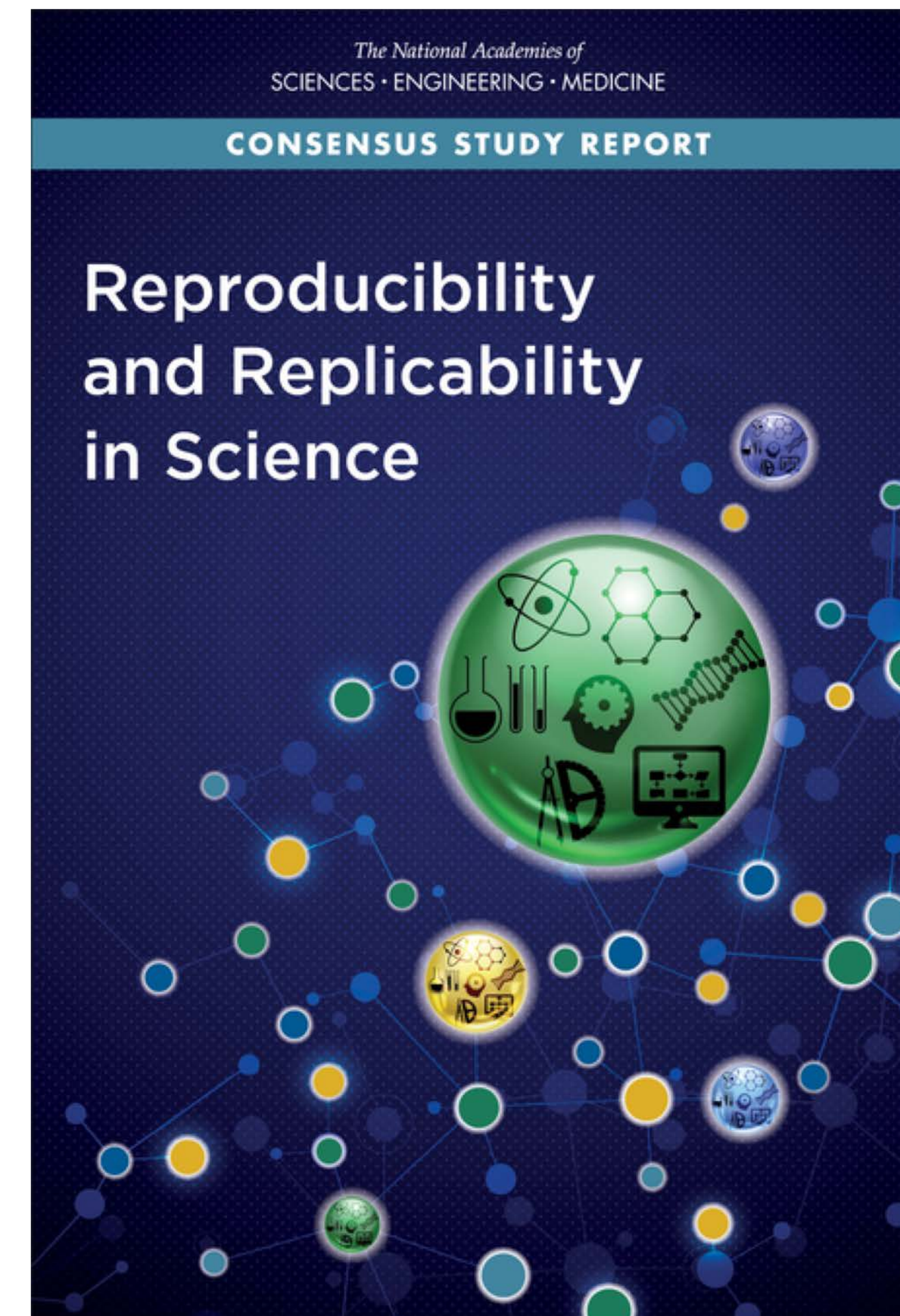
WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

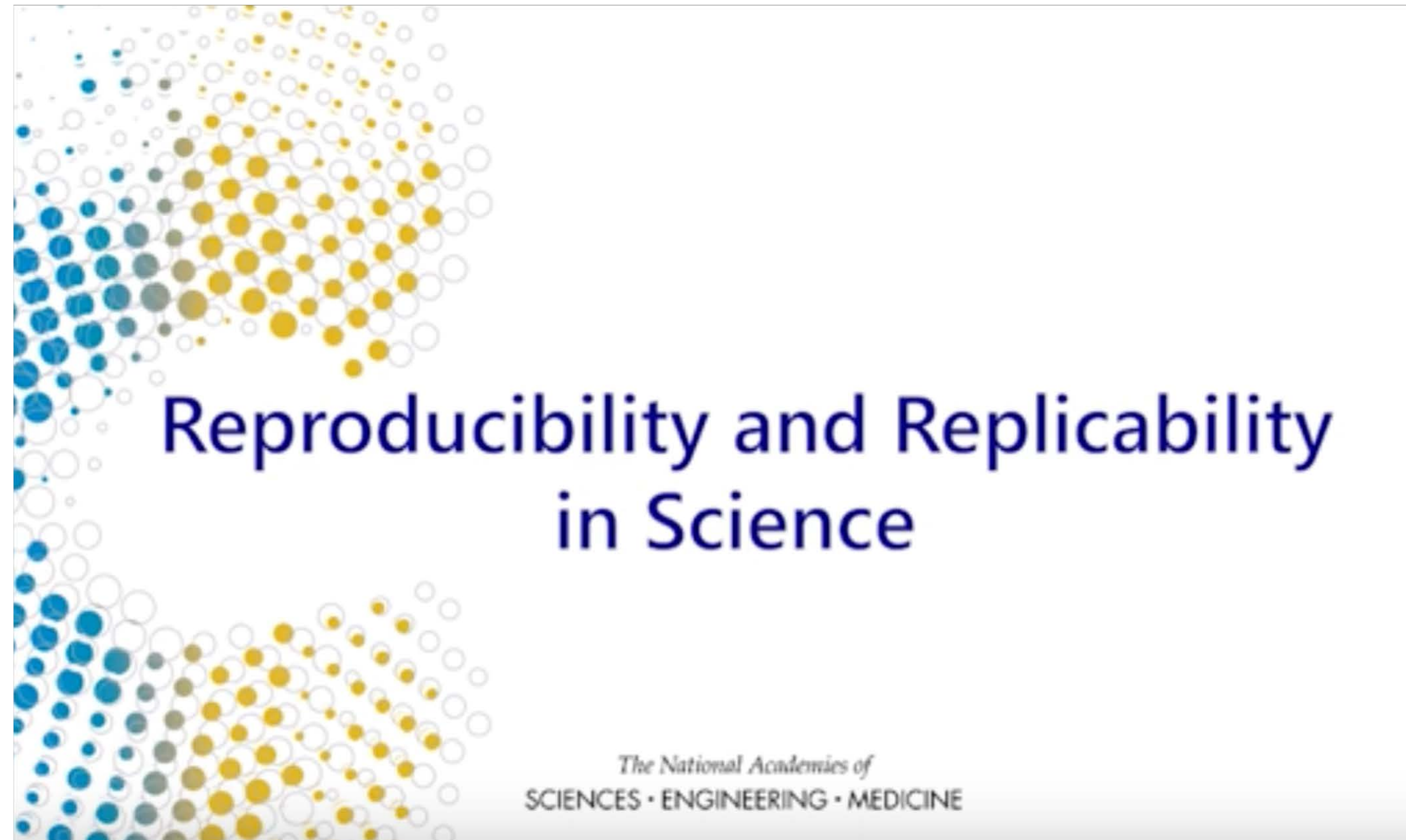


Guidelines

- Researchers: include a clear, specific, and complete description of how the reported results were reached
- Funding agencies and organizations: should consider investing in research and development of open-source, usable tools and infrastructure that support reproducibility
- Journals: should consider ways to ensure computational reproducibility
- The NSF: should take steps to facilitate the transparent sharing and availability of digital artifacts, such as data and code



Guidelines



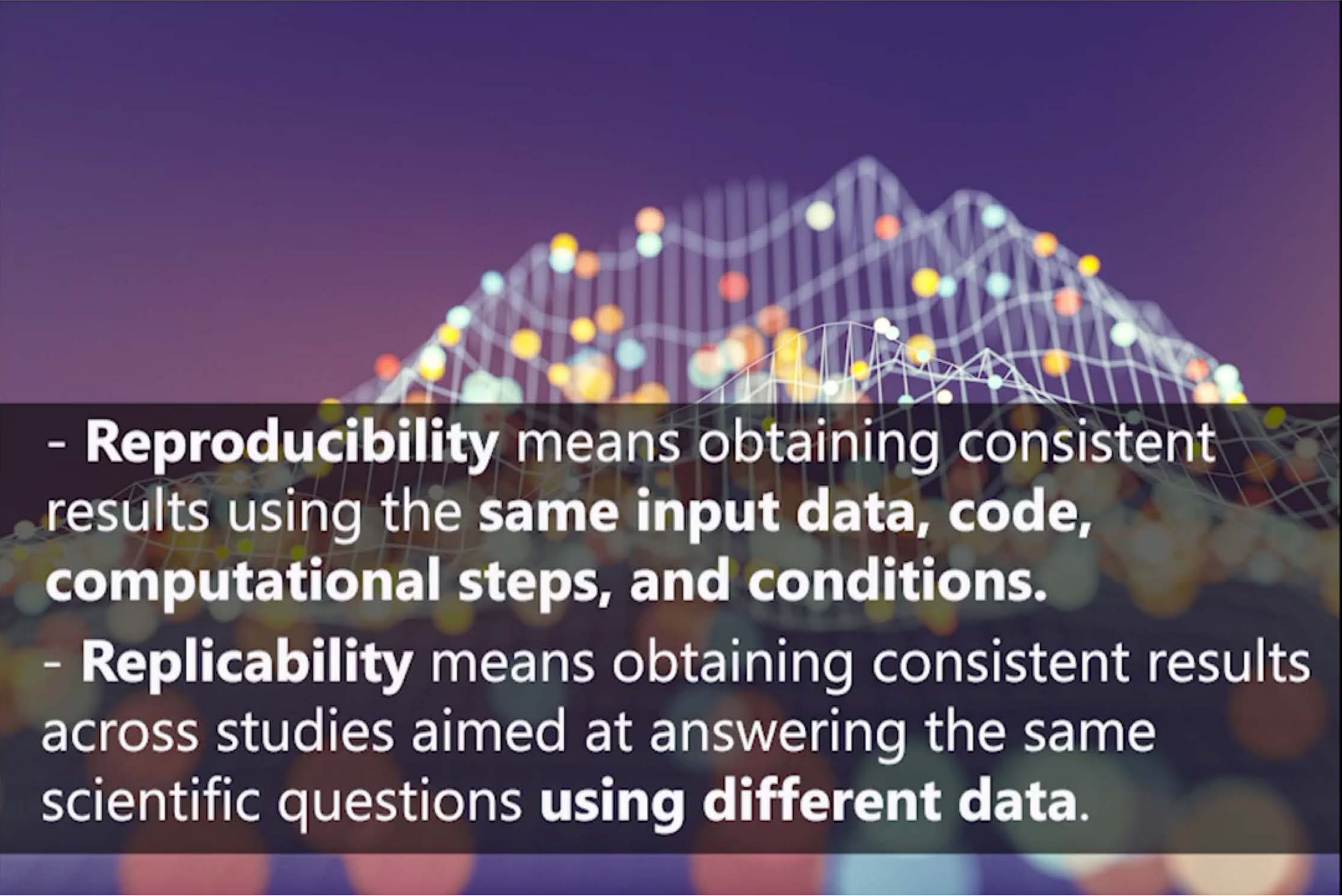
Source: [NASEM](#)

Reproducibility and Replicability in Data Science



One of the ways that the scientific community confirms the validity of a new scientific discovery is by repeating the research that produced it.

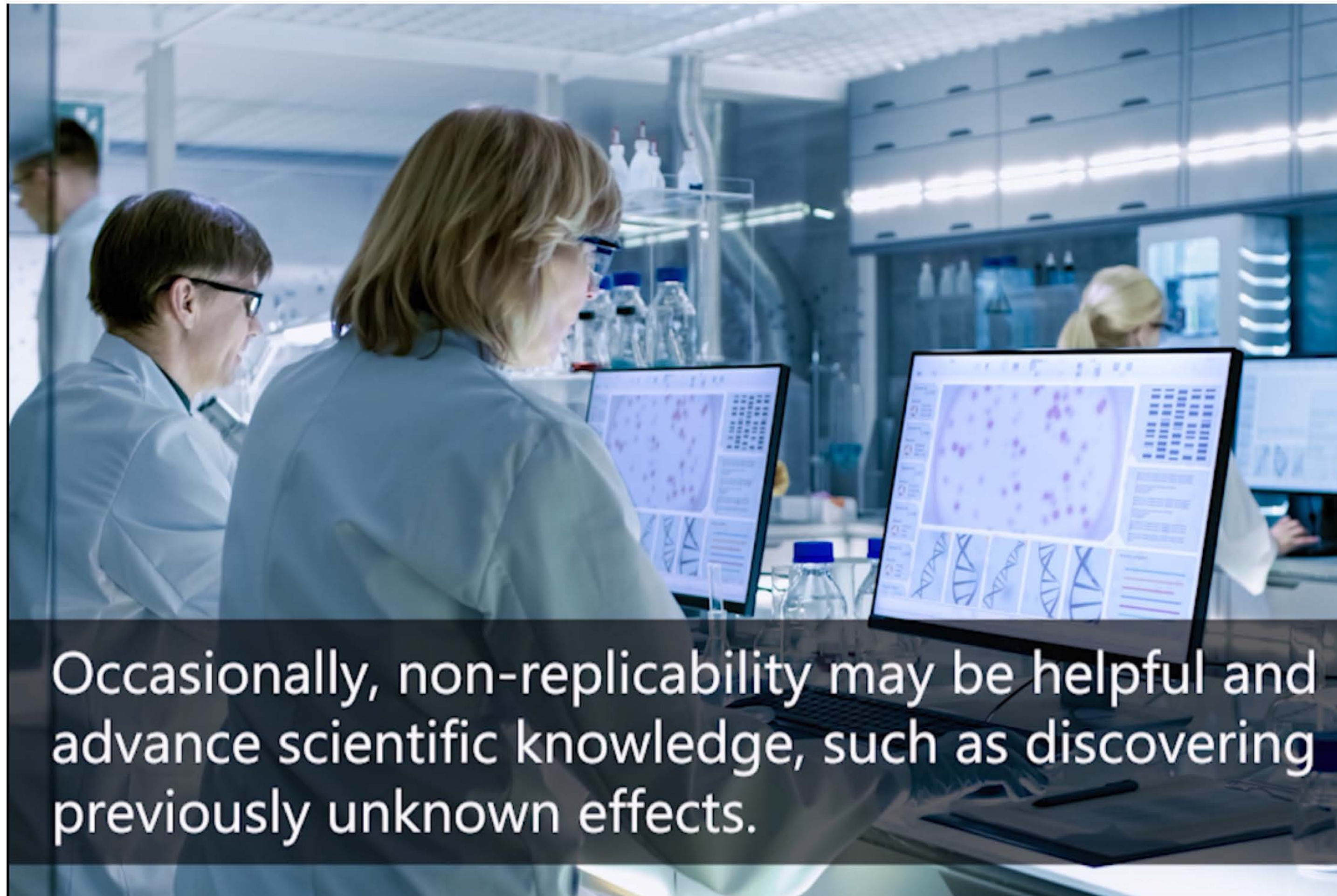
Reproducibility and Replicability in Data Science

- 
- **Reproducibility** means obtaining consistent results using the **same input data, code, computational steps, and conditions.**
 - **Replicability** means obtaining consistent results across studies aimed at answering the same scientific questions **using different data.**

Reproducibility and Replicability in Data Science



Reproducibility and Replicability in Data Science



Occasionally, non-replicability may be helpful and advance scientific knowledge, such as discovering previously unknown effects.

Reproducibility and Replicability in Data Science

At other times, a study cannot be replicated due to reasons ranging from simple mistakes to bias and, rarely, fraud.

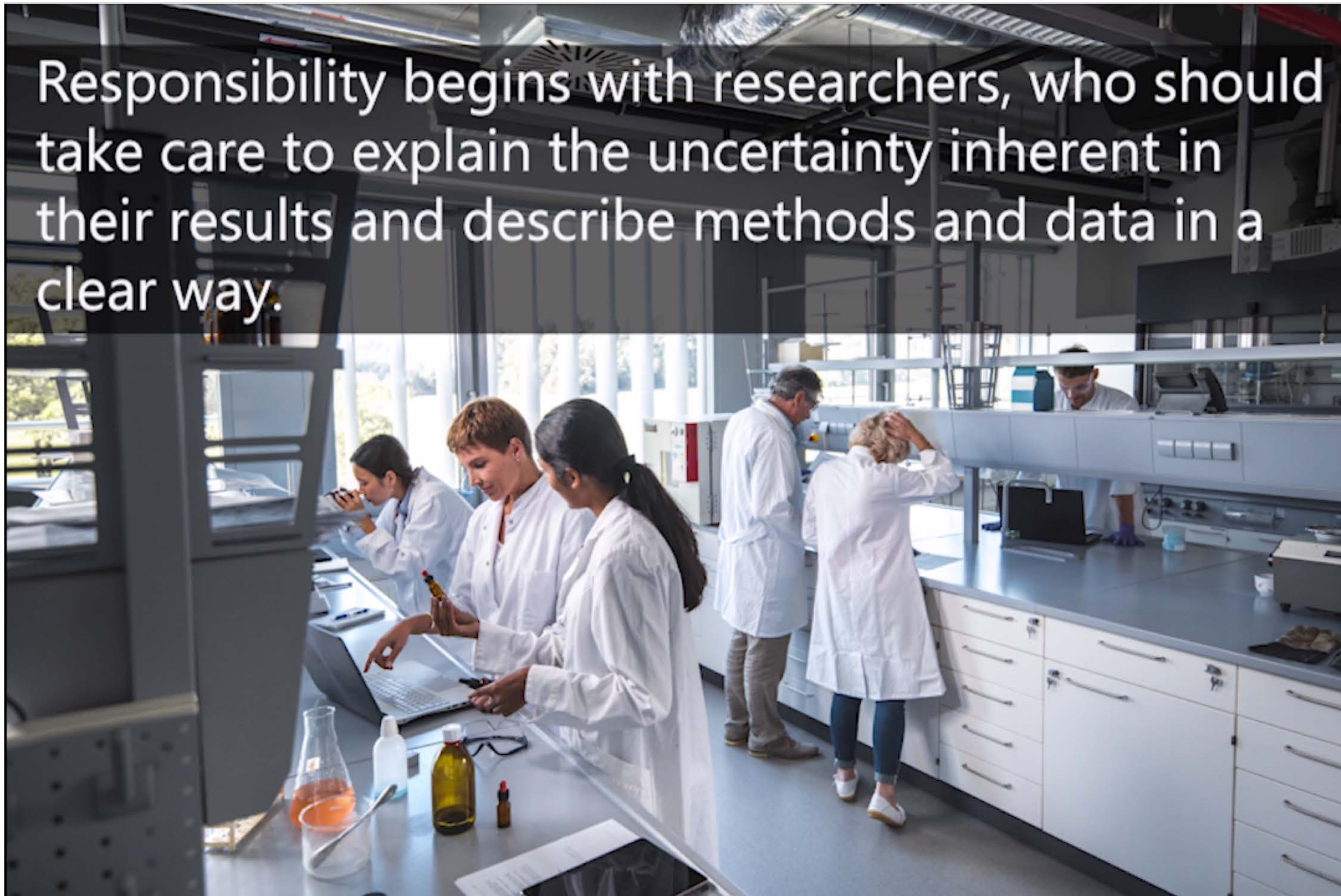


Reproducibility and Replicability in Data Science



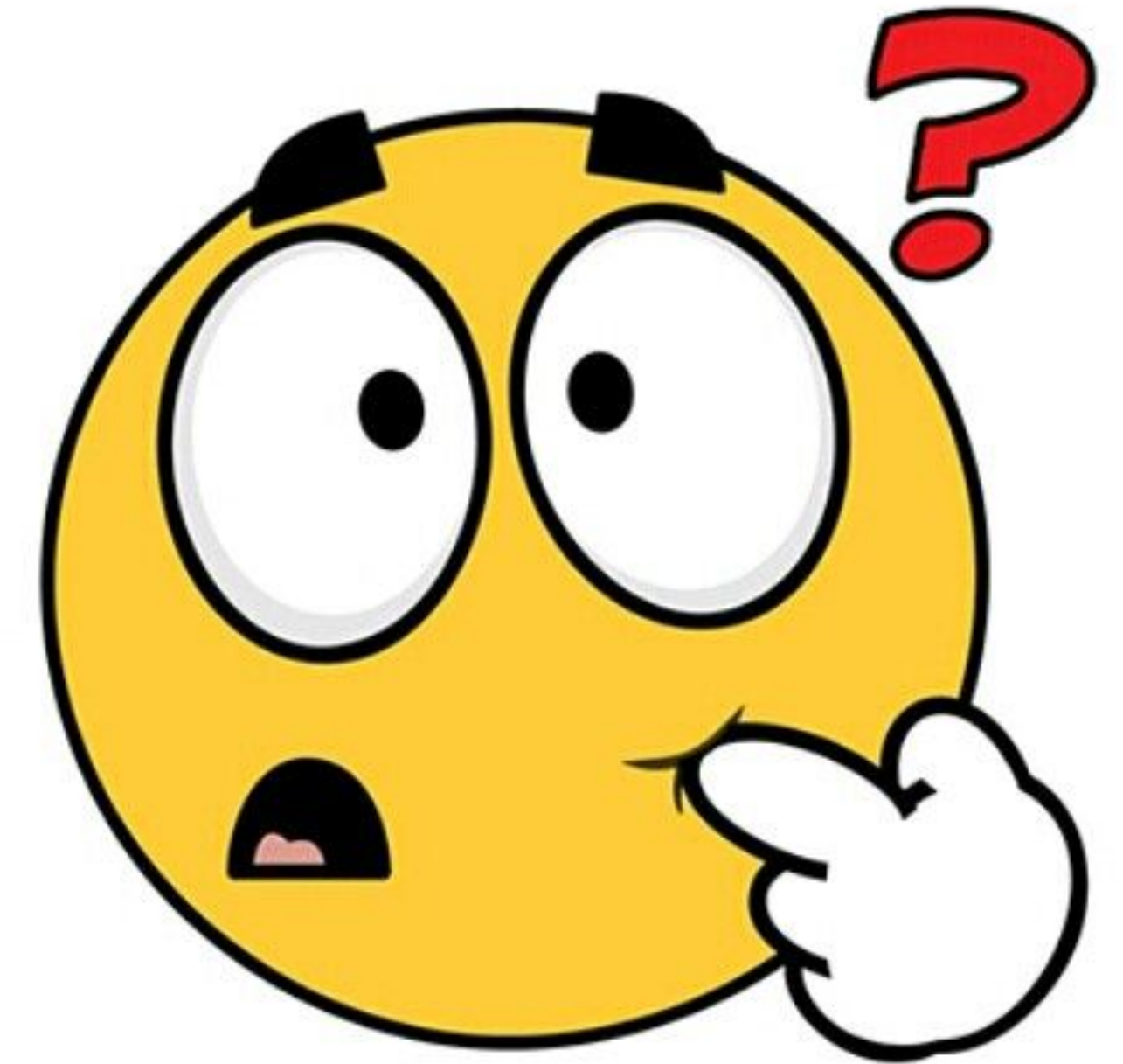
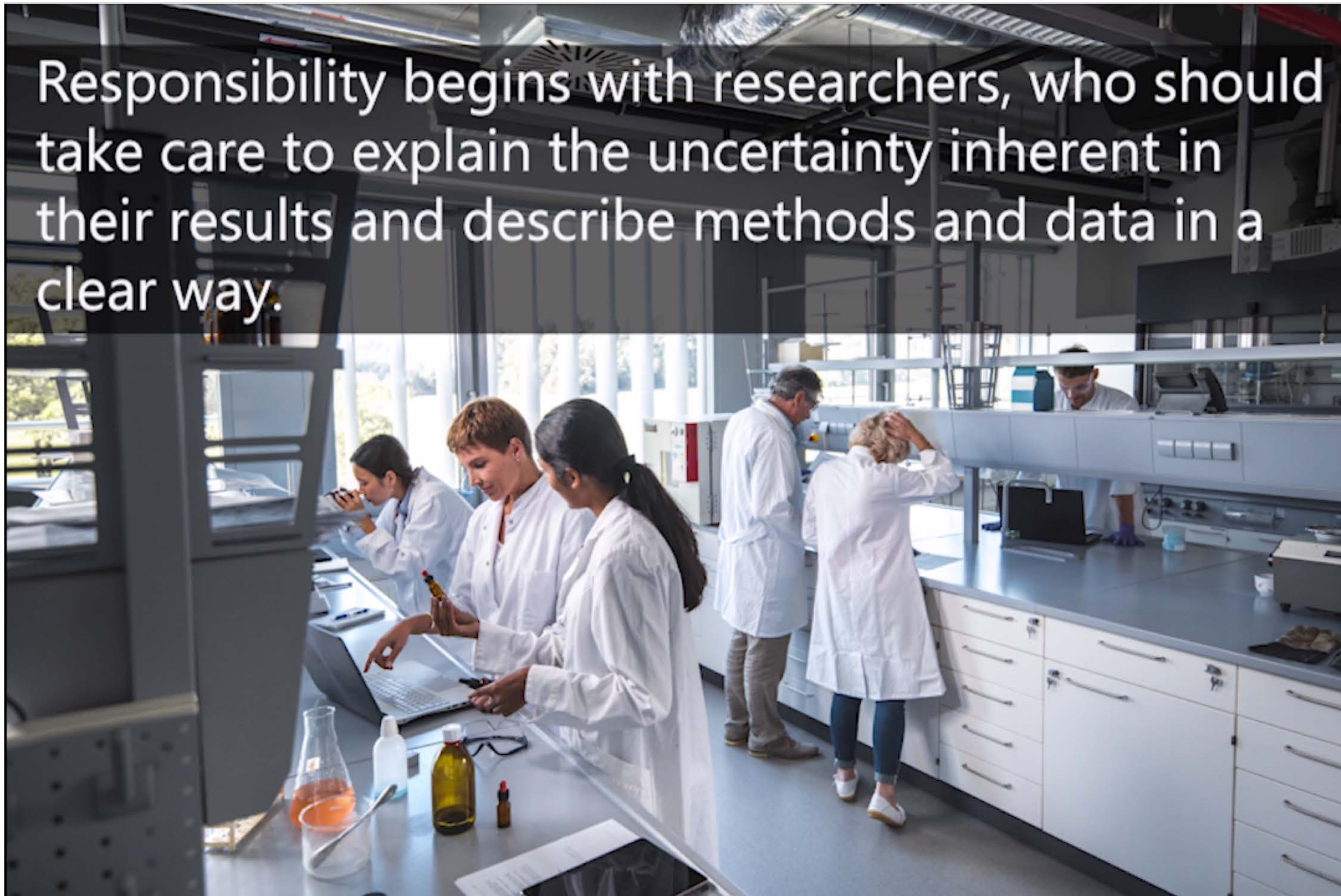
Reproducibility and Replicability in Data Science

Responsibility begins with researchers, who should take care to explain the uncertainty inherent in their results and describe methods and data in a clear way.



Reproducibility and Replicability in Data Science

Responsibility begins with researchers, who should take care to explain the uncertainty inherent in their results and describe methods and data in a clear way.



Five FAQs in Data-Driven Research

1. How did I compute this result?

2. How does new data change this result?

3. How did you compute *your* result?

Can I use your data to reproduce it?

With your code?

On your infrastructure?

4. Has anyone ever used an <XYZ-algorithm>
on this data?

Did it work?

5. Who is using the data and algorithms?

In which context?

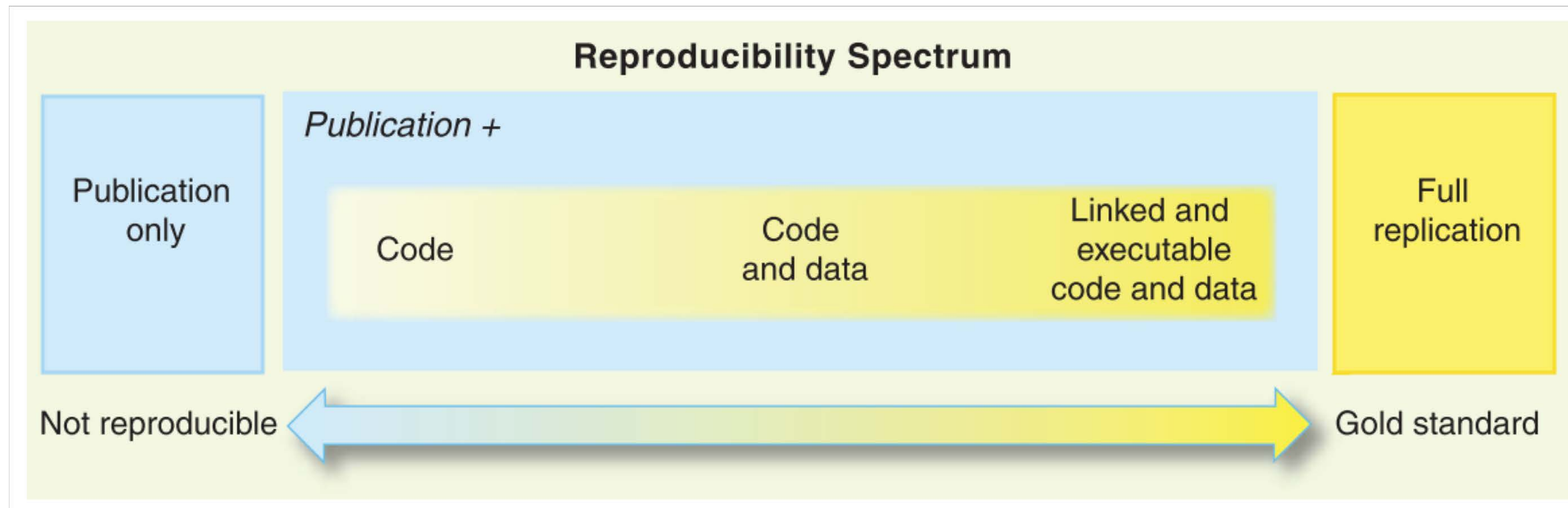
Five Questions —> Three Words

Reproducibility

Reusability

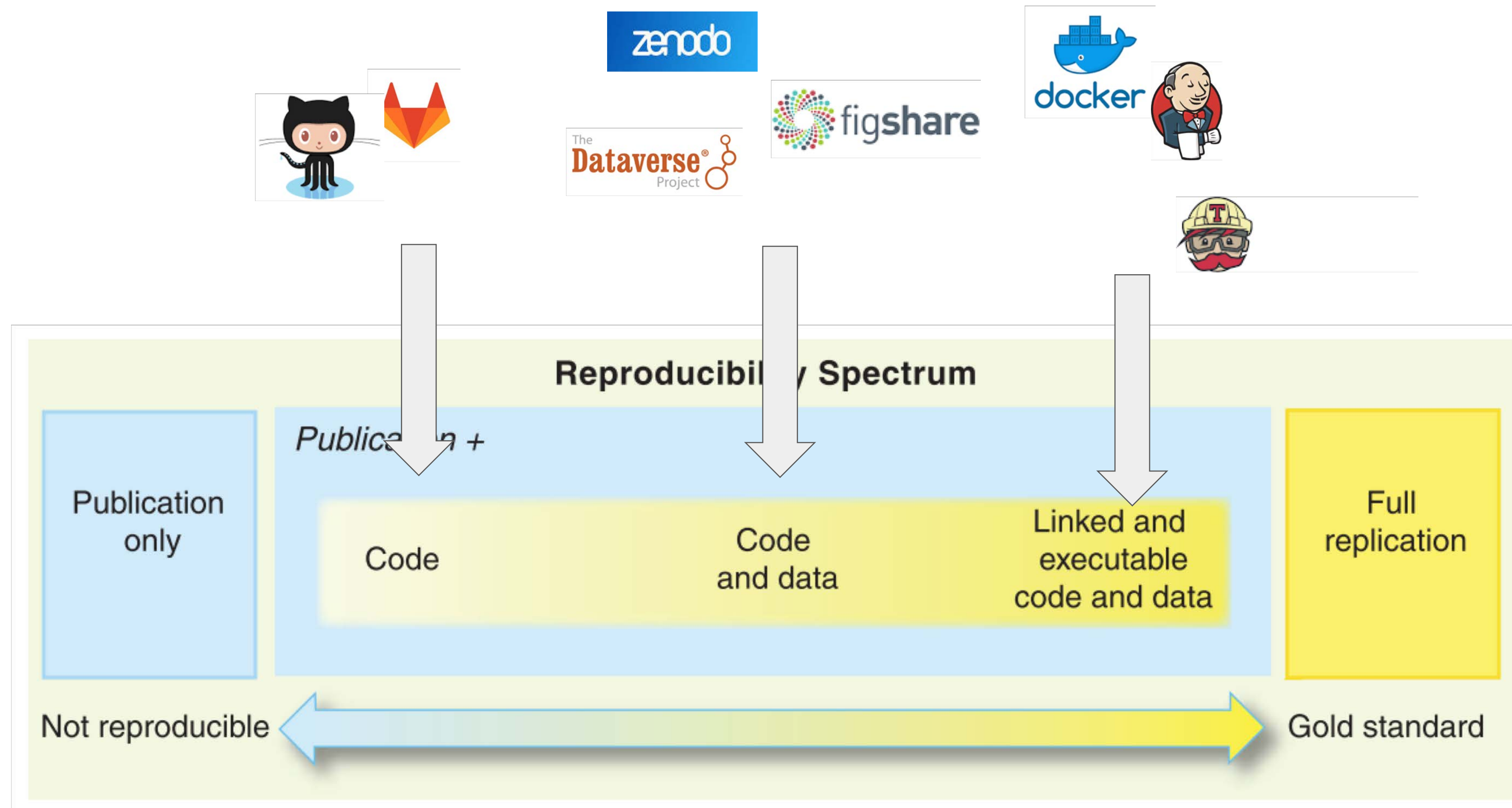
Collaboration

Reproducibility spectrum



Source: [Peng, 2011](#)

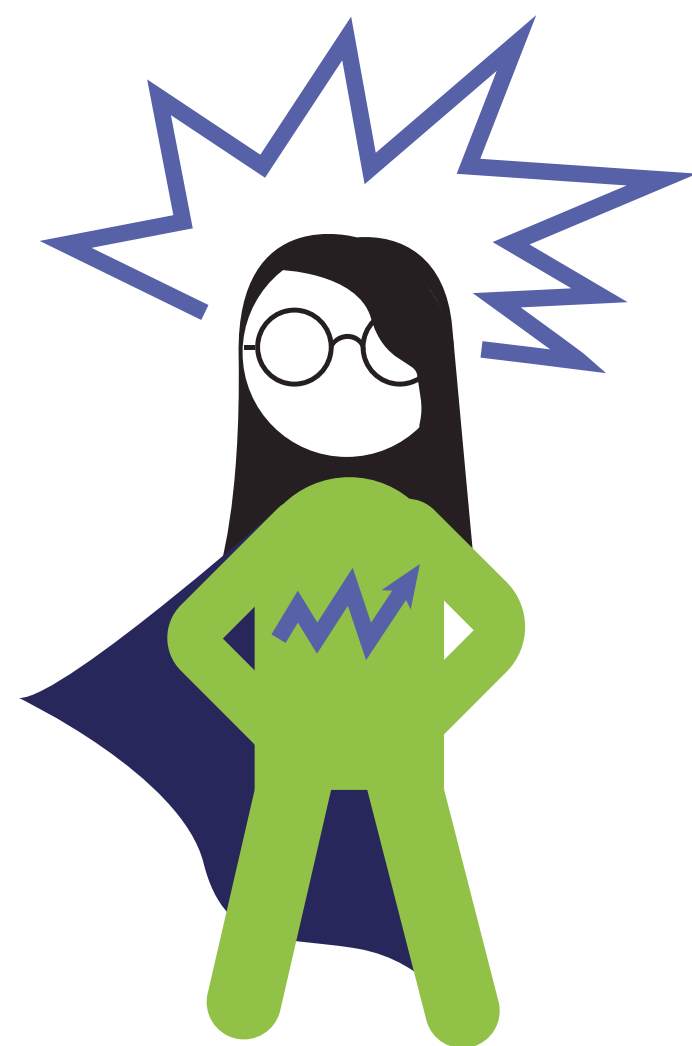
Reproducibility spectrum



Source: [Peng, 2011](#)

A Myriad of Tools

Hard **work** to make science reproducible, accessible, open etc...

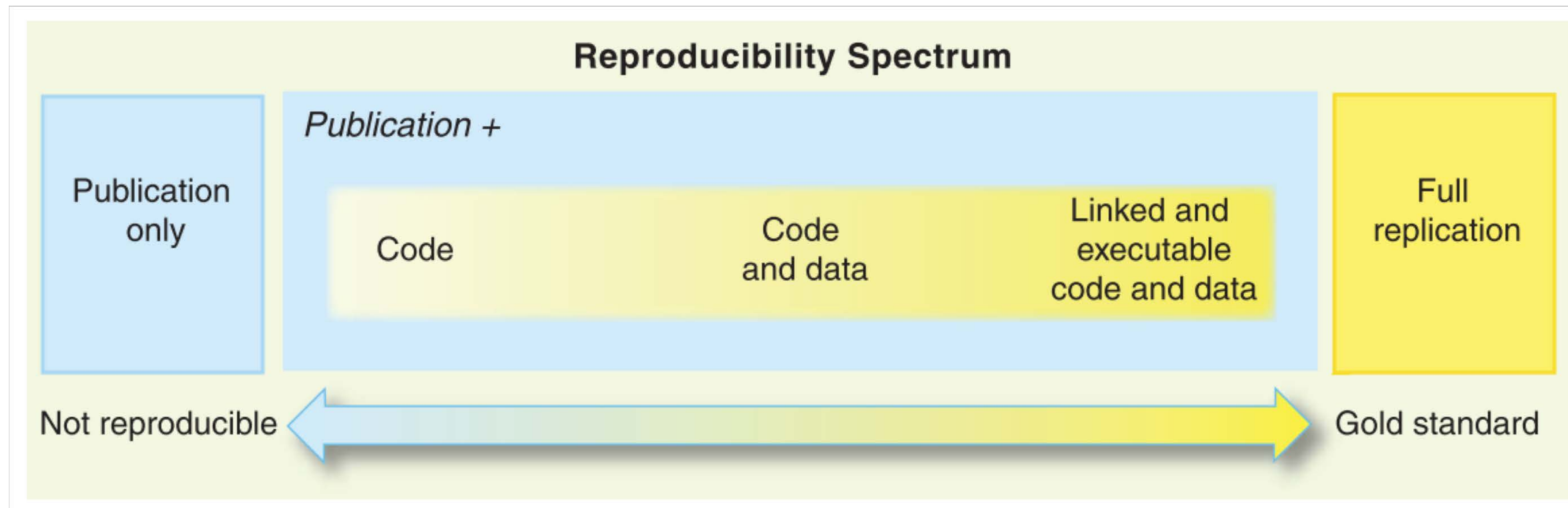


- Version control (git)
- Code sharing (GitHub, GitLab, BitBucket)
- Data sharing (Zenodo, Figshare, institutional digital archives)
- Presentation and communication (Jupyter, RStudio)
- Correctness – testing (CI, e.g. travis)
- Packaging, containerization (docker, singularity)

Difficult to stay **productive** and worry about all of the above!

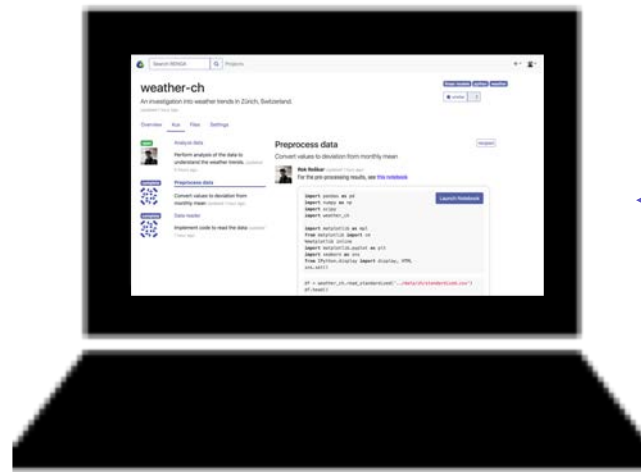
Reproducibility spectrum

RENKU 連句



Source: [Peng, 2011](#)

RENKU - 連句 - A platform for reproducible (data) science



local



Jupyter
Notebook

RStudio

MATLAB

...

RENKU

- Reproducibility
- Reuse
- Collaboration

Gitlab

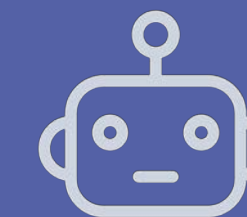
RDF

Docker

Keycloak

JupyterLab

Kubernetes



On premise,
Or on the cloud

Terminology

- We borrow the **RENKU** name from the Japanese word for *linked-verse poetry*
- A “**ku**” is a verse in a renku poem
- We use “**ku**” to mean a piece of the data analysis process – includes discussion, code, and results

1.

Provide the means to create **reproducible** (data) science

2.

Facilitate the **sharing** and **reuse** of research artefacts

3.

Foster a **collaborative environment** for interactive prototyping

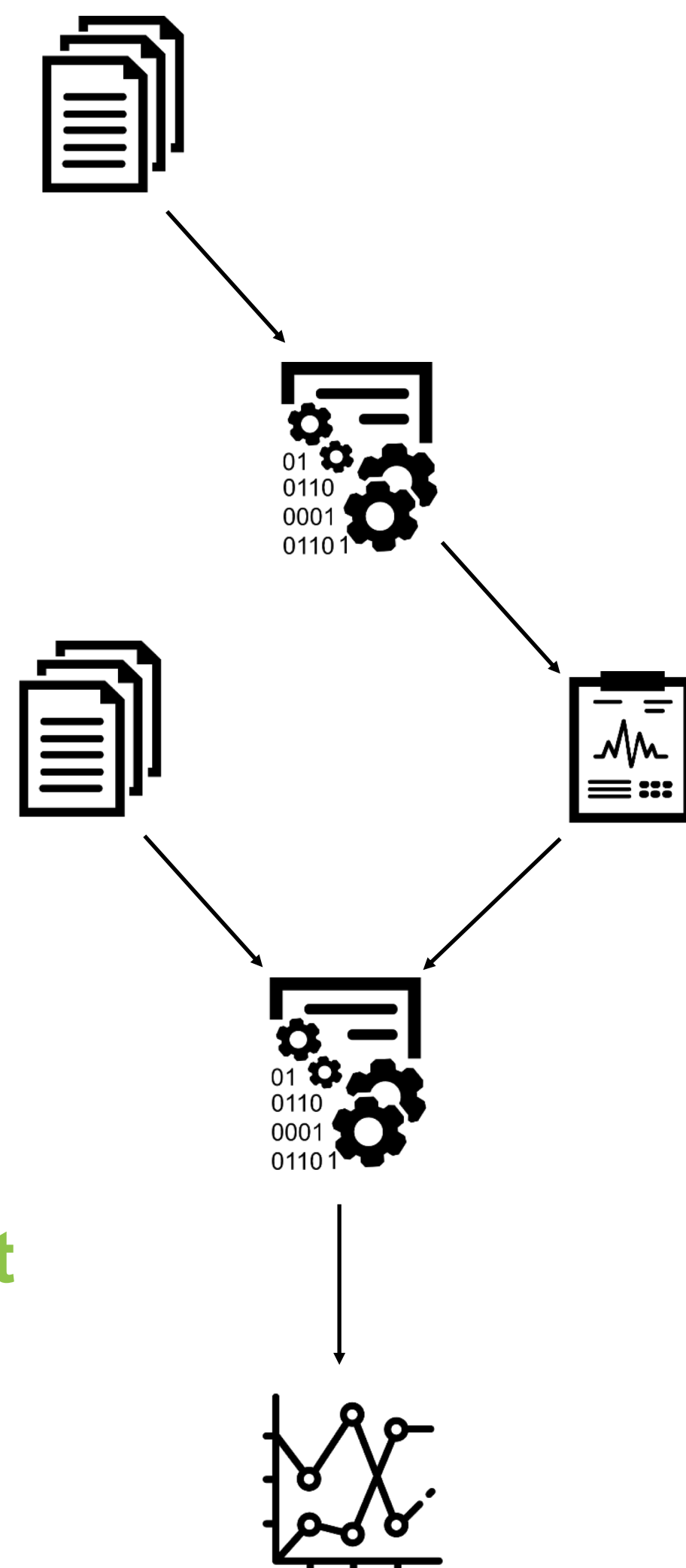
4.

Enable the **discovery** of relevant data and methods

5.

Allow **federated access** across institutions giving each the freedom to impose its own access controls over resources

Capture the scientific process



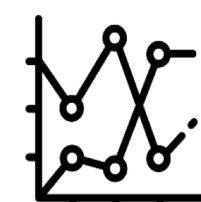
Raw data



Intermediate result



Software



Results

1. Inputs and outputs of analysis steps are recorded into a **knowledge graph** *while the work is being done*
2. Steps can be **repeated** or integrated into more complex **workflows**
3. **Provenance** of all data products is always accessible via simple tools
4. **Version control** is built-in for data, code, and workflows

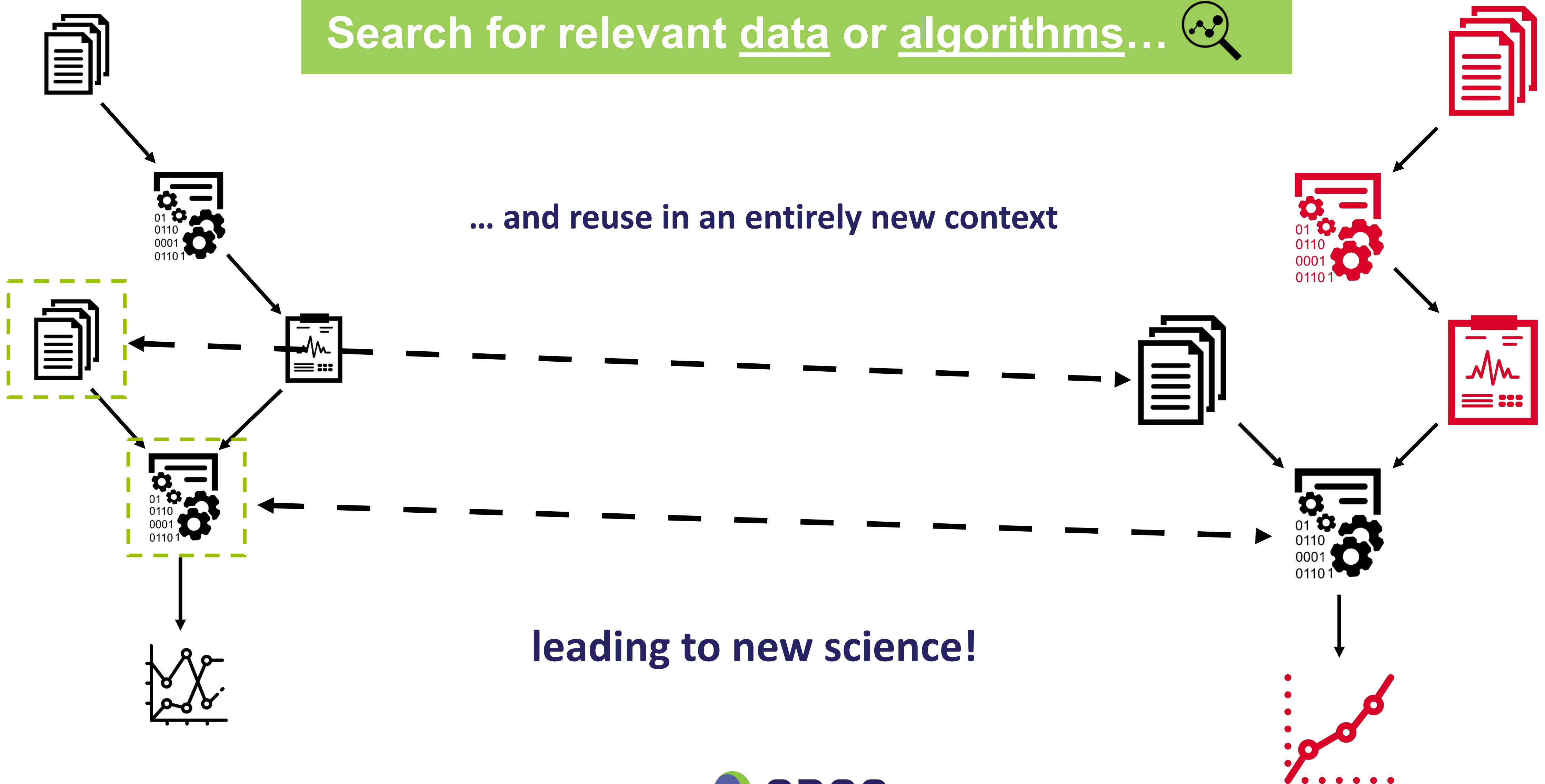
Reuse and repeat

Search for relevant data or algorithms...

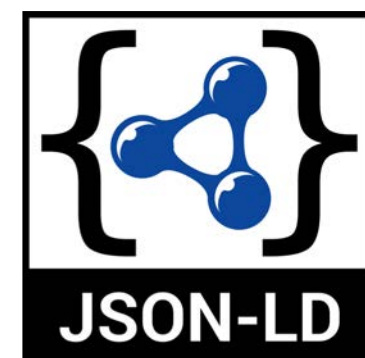


... and reuse in an entirely new context

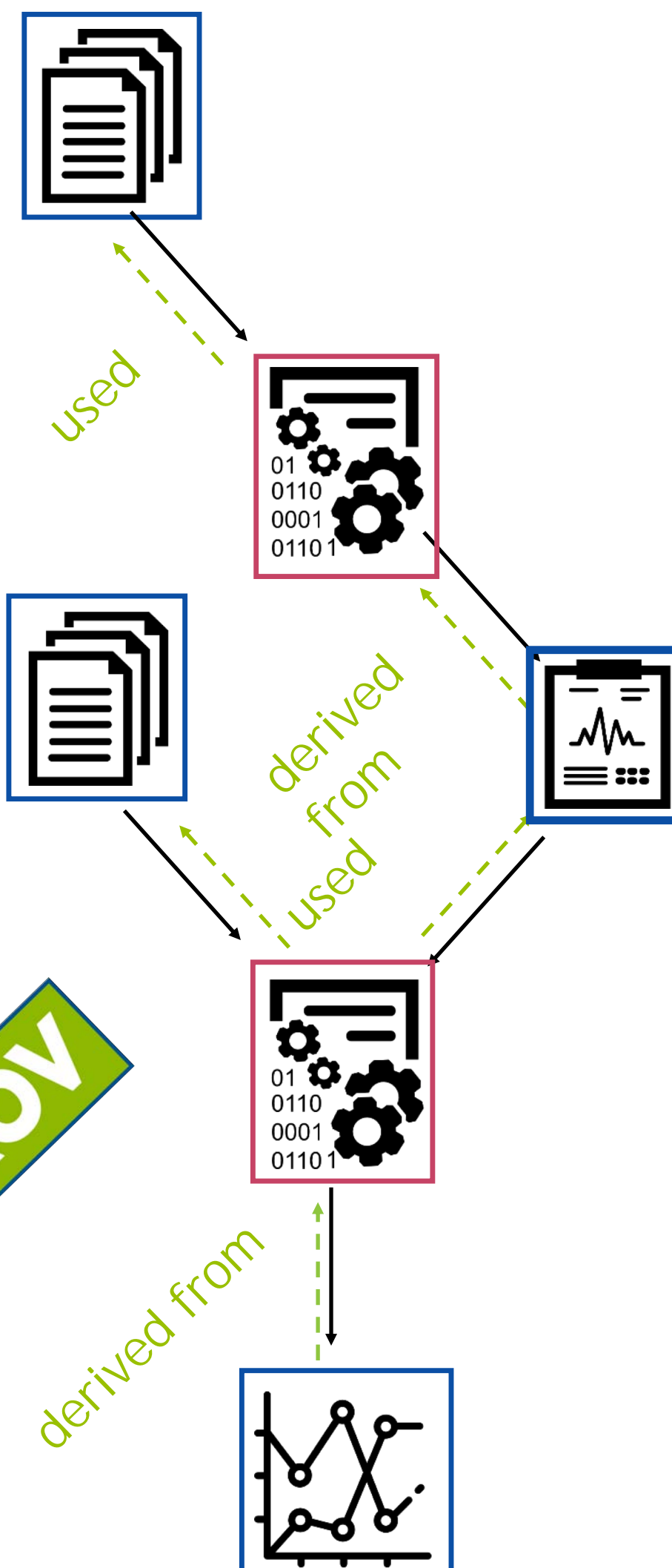
leading to new science!



Encapsulate with rich metadata



- Metadata use Dublin Core, FOAF, and Schema.org
- Provenance graph is based on PROV-O W3C recommendation



- CWL for representing all computational steps
- Capture individual steps from user input
- Tools for constructing workflows from basic pieces
- Rely on container technologies to ensure reproducibility

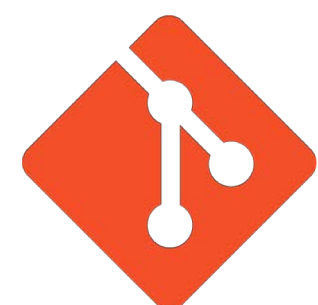
Reusability and collaboration

- Who is using the data and how?
- Which algorithms are used to answer which questions?
- How to regenerate results if new data becomes available? If old data is now off-limits?
- Who to credit?
- How popular is my work/the work of my lab/my unit?

TRUST

What is Renku, really?

A platform integrating:
git, Jupyter/RStudio, docker, analysis workflows linked with a **knowledge graph**



Jump to or search... Projects Notebooks

vom_natt
rokroskar/vom_natt

Overview Kus **Files** Pending Changes Notebook Servers Settings

File View

- .renku
- data
- img
- notebooks
- .gitkeep
- run_vom.ip...
- src_py
- src_sh
- work
- .gitattributes
- .gitignore
- .gitlab-ci.yml
- .gitmodules
- .renku.lock
- Dockerfile
- README.md

notebooks/run_vom.ipynb File view

Run VOM with RENKU

This notebook as meant as documentation in addition to the knowledge graphs generated with renku. It contains all renku commands run to achieve the final results, with explanations and argumentation.

Data preparation

The first step involves running a python script to generate the input files for the VOM. Weather data is taken from the Australian SILO: <https://legacy.longpaddock.qld.gov.au/silo/datadrill/format.php>

CO2 data is taken from the Mauna Loa records: http://scrippsco2.ucsd.edu/data/atmospheric_co2/mlo

Interpolation between missing values is done linearly. Eventually the file dailyweather.prn is produced, which contains meteorological data in the format for the VOM.

View in GitLab

Jump to or search... Projects Notebooks

vom_natt
rokroskar/vom_natt

Overview Kus Files Pending Changes **Notebook Servers** Settings

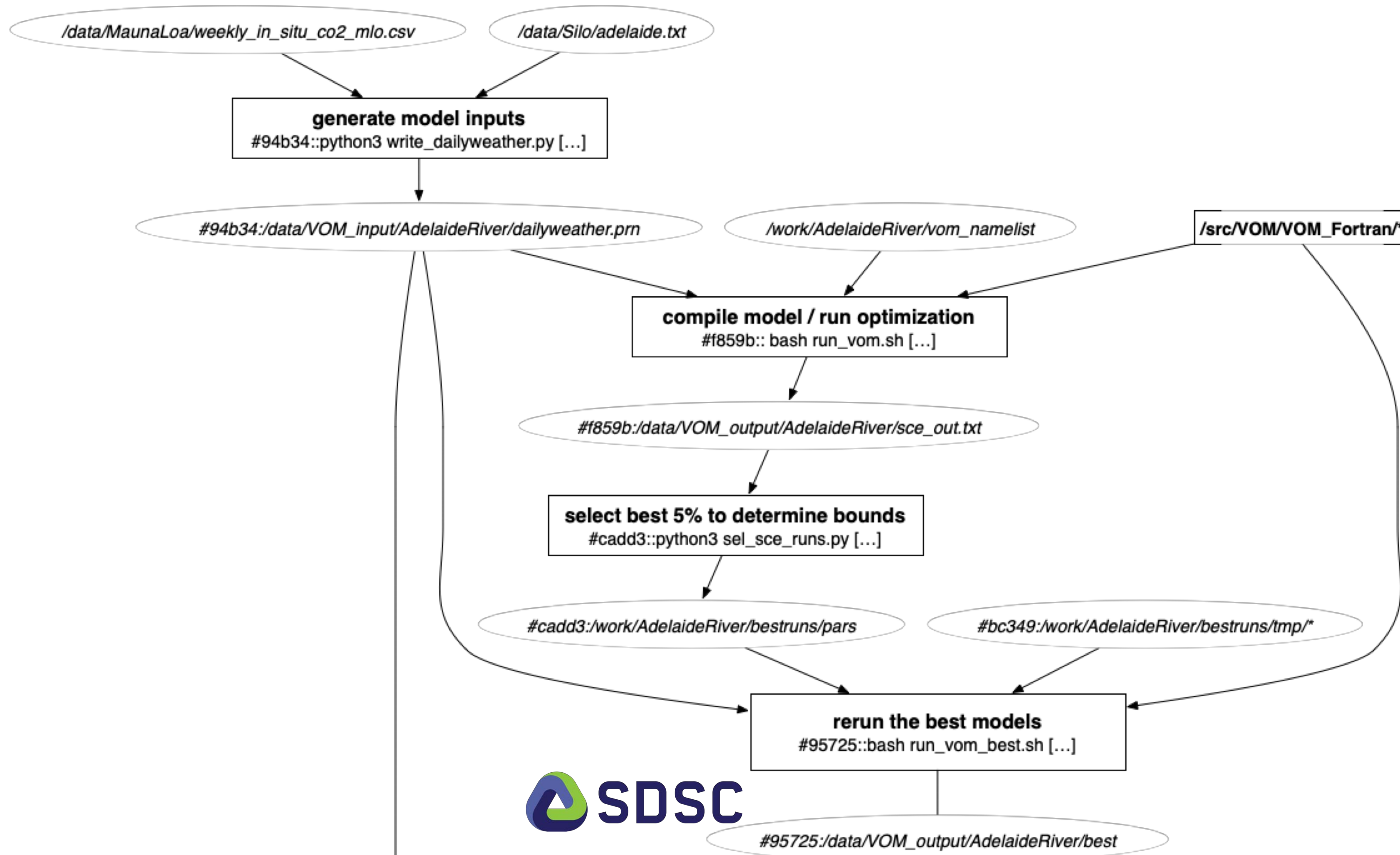
Branch	Commit	Action
master	c495b1ca	Running

Start new server

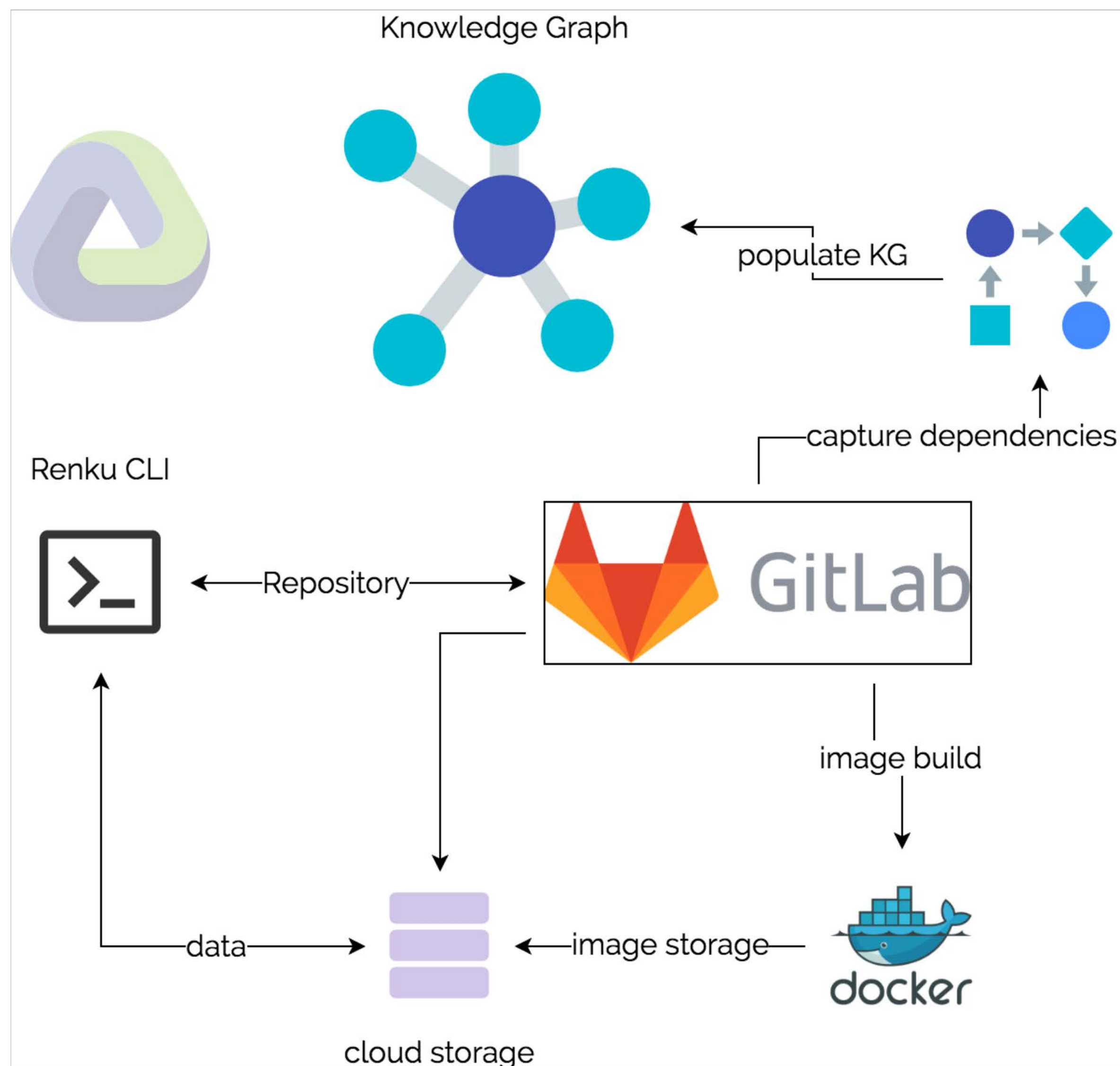
View in GitLab

What is Renku, really?

A platform integrating:
git, Jupyter/RStudio, docker, analysis workflows linked with a **knowledge graph**

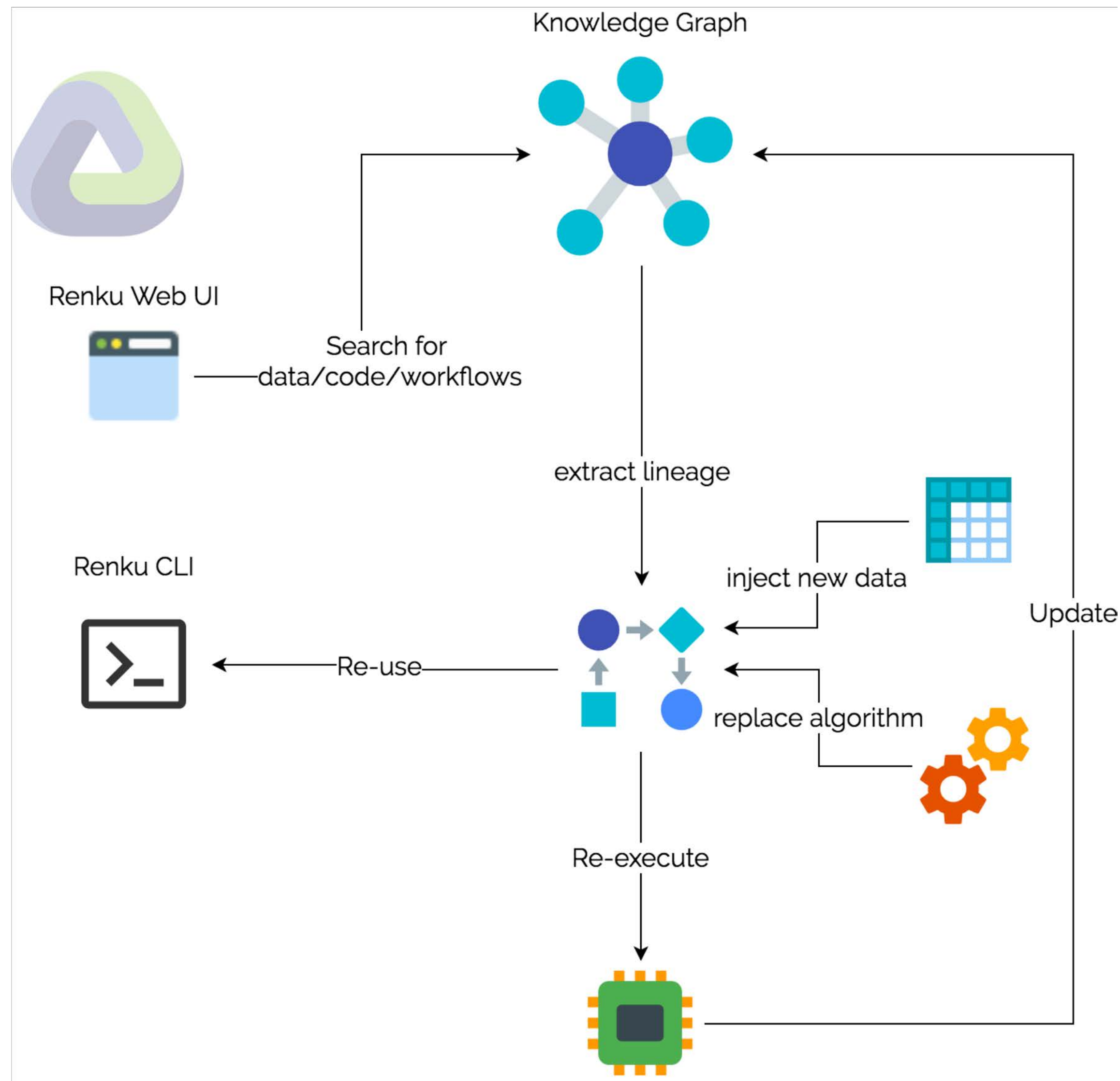


Reproducibility



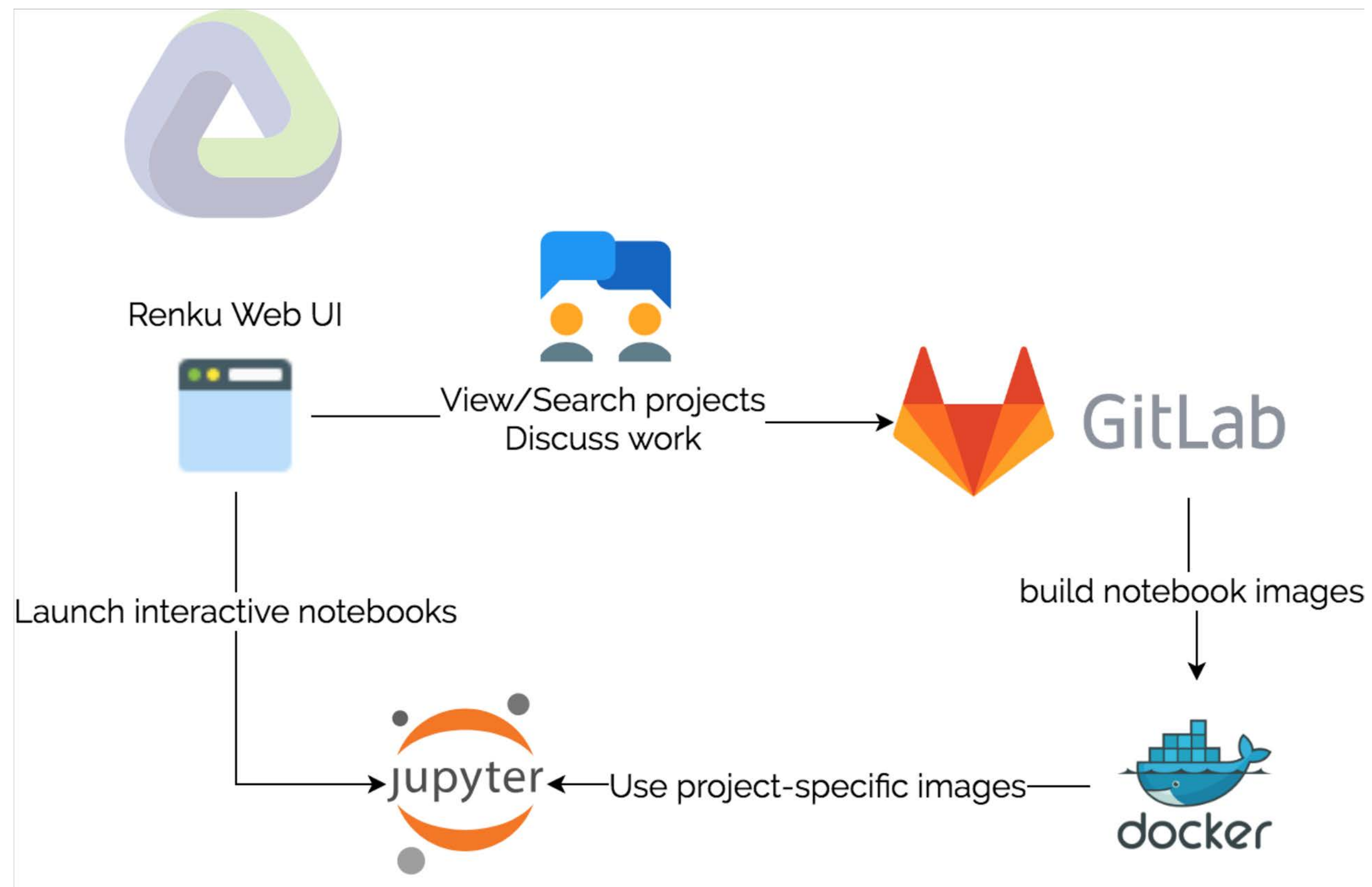
- CLI developed from ground up taking advantage of Git to capture lineage
- Designed with interoperability in mind, using linked data standards to express dependencies
- Using version control for data and code – overcoming usability challenges

Reusability



- Understanding lineage means we can always re-execute and update results as methods or data change
- Workflow construction and re-execution possible with the local client
- Developing the means to execute workflows in the cloud and on HPC (REANA, Cromwell, Toil, ...)
- Search for graph artifacts being developed

Collaboration



- Web UI serves as primary point of contact for users
- Allows creation of projects, discussions with media embedding
- Creation of hosted interactive sessions with version-controlled environments (Docker images)
- Various supported languages, e.g. Python, MATLAB (with GUI), R, RStudio

Where is Renku used?

- We currently run 6 deployments, 4 are “production”
- Several SDSC academic projects are using RENKU from the start
- Large use-case from EPFL life sciences
 - reproducibility/collaboration
 - end-to-end lineage (from wet lab through computation)

Status and Roadmap Q3/Q4 2019

- Knowledge graph is live and first use by UI coming to renkulab.io
 - Many UI/UX improvements, more coming
 - Integrations for data import/export with external repositories e.g. Zenodo
-
- UI features for easier adoption e.g. environment configurations
 - Tighter integration of UI with the knowledge graph, improving metadata
 - Support for more environments, e.g. HPC
 - Workflow execution in the cloud/HPC
 - Solidifying our operations and deployment practices

Open Science ≠ Open Data

What is FAIR DATA?



Data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier.

FINDABLE



Metadata and data are understandable to humans and machines. Data is deposited in a trusted repository.

ACCESSIBLE



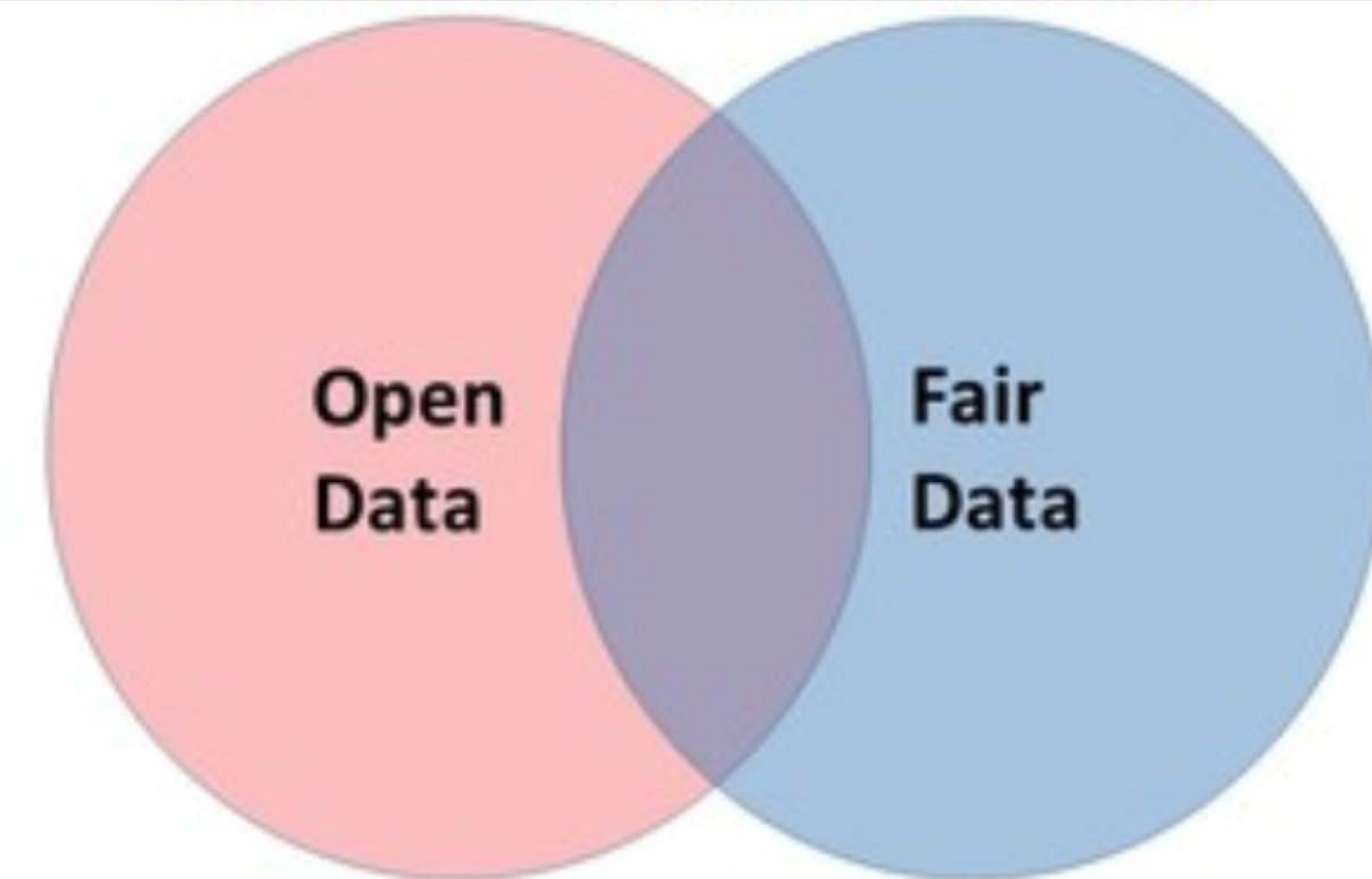
Metadata use a formal, accessible, shared, and broadly applicable language for knowledge representation.

INTEROPERABLE



Data and collections have a clear usage licenses and provide accurate information on provenance.

REUSABLE



Source: [LIBER](#)

Swiss Data Custodian



- A multisided service to maximize economic opportunities and meet the societal challenges brought about by the digital revolution
- It is a data vault + secure multiparty compute ecosystems
 - Assign data ownership to the rightful person
 - Preserve data sovereignty
 - Establish trust and transparency in data sharing
 - Promote economic and societal incentives for sharing data
 - Enable cooperation between mutually non-trusting parties
- Governed by a trusted entity
 - Monitor and maintain compliance with DPA obligations

Thank you!

Renku is under very active development

Preview with tutorial at: <https://renkulab.io>

Open-source:
<https://github.com/SwissDataScienceCenter>

All feedback welcome!

<https://datascience.ch/>
[@SDSCdatascience](#)

