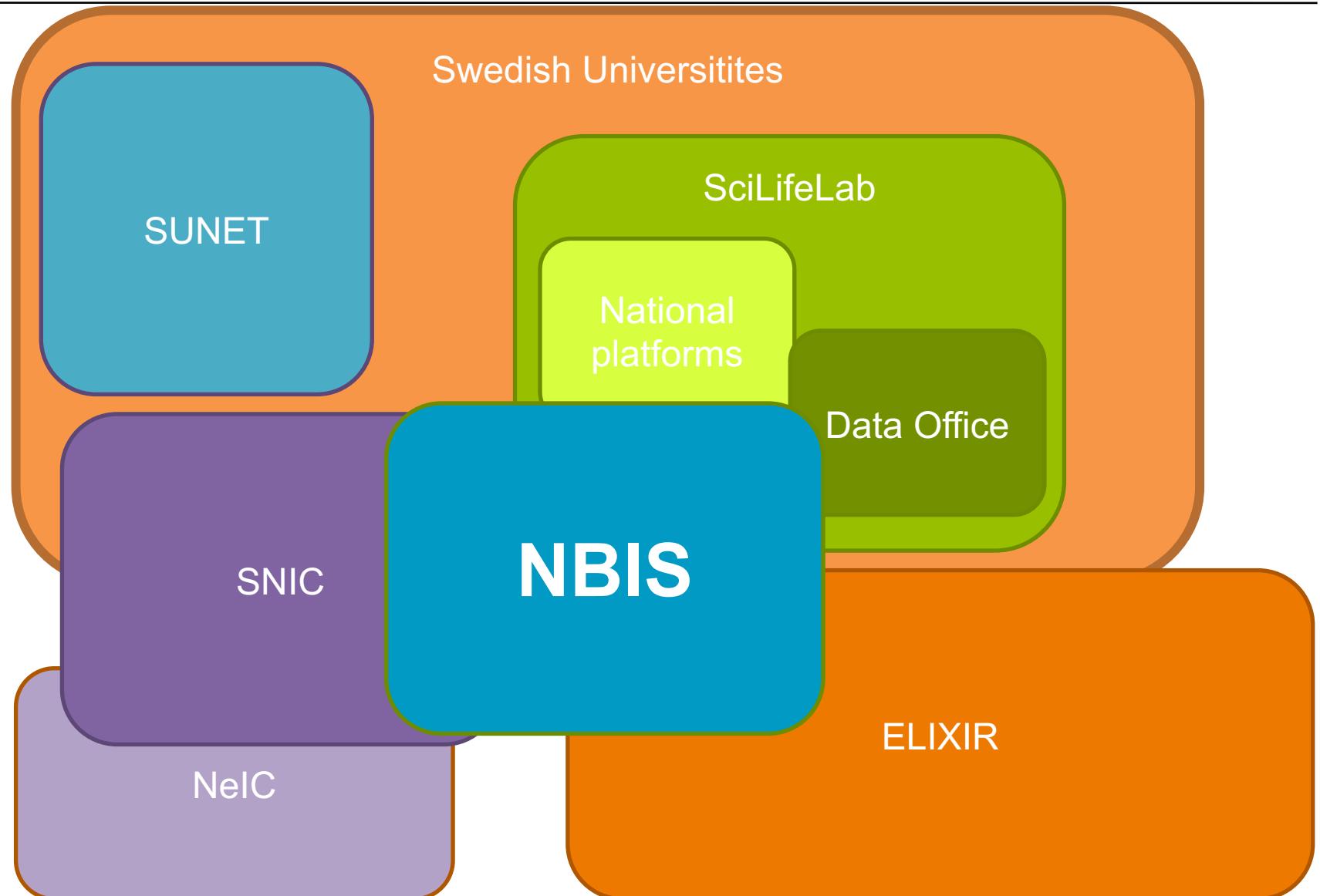

Managing your data

Niclas Jareborg, NBIS
niclas.jareborg@nbis.se

2018-09-03



- To make your research easier!
- To stop yourself drowning in irrelevant stuff
- In case you need the data later
- To avoid accusations of fraud or bad science
- To share your data for others to use and learn from
- To get credit for producing it
- Because funders or your organisation require it



Well-managed data opens up opportunities for
re-use, integration and new science

Science

LETTERS

Cite as: J. Berg., *Science*
10.1126/science.aan5763 (2017).

Editorial Retraction

Jeremy Berg

Editor-in-Chief

After an investigation, the Central Ethical Review Board in Sweden has recommended the retraction of the Report "Environmentally relevant concentrations of microplastic particles influence larval fish ecology," by Oona M. Lönnstedt and Peter Eklöv, published in *Science* on 3 June 2016 (1). *Science* ran an Editorial Expression of Concern regarding the Report on 1 December 2016 (2). The Review Board's report, dated 21 April 2017, cited the following reasons for their recommendation: (i) lack of ethical approval for the experiments; (ii) absence of original data for the experiments reported in the paper; (iii) widespread lack of clarity concerning how the experiments were conducted. Although the authors have told *Science* that they disagree with elements of the Board's report, and although Uppsala University has not yet concluded its own investigation, the weight of evidence is that the paper should now be retracted. In light of the Board's recommendation and a 28 April 2017 request from the authors to retract the paper, *Science* is retracting the paper in full.

REFERENCES

1. O. M. Lönnstedt, P. Eklöv, *Science* **352**, 1213 (2016).
2. J. Berg, *Science* **354**, l242 (2016); published online 1 December 2016.

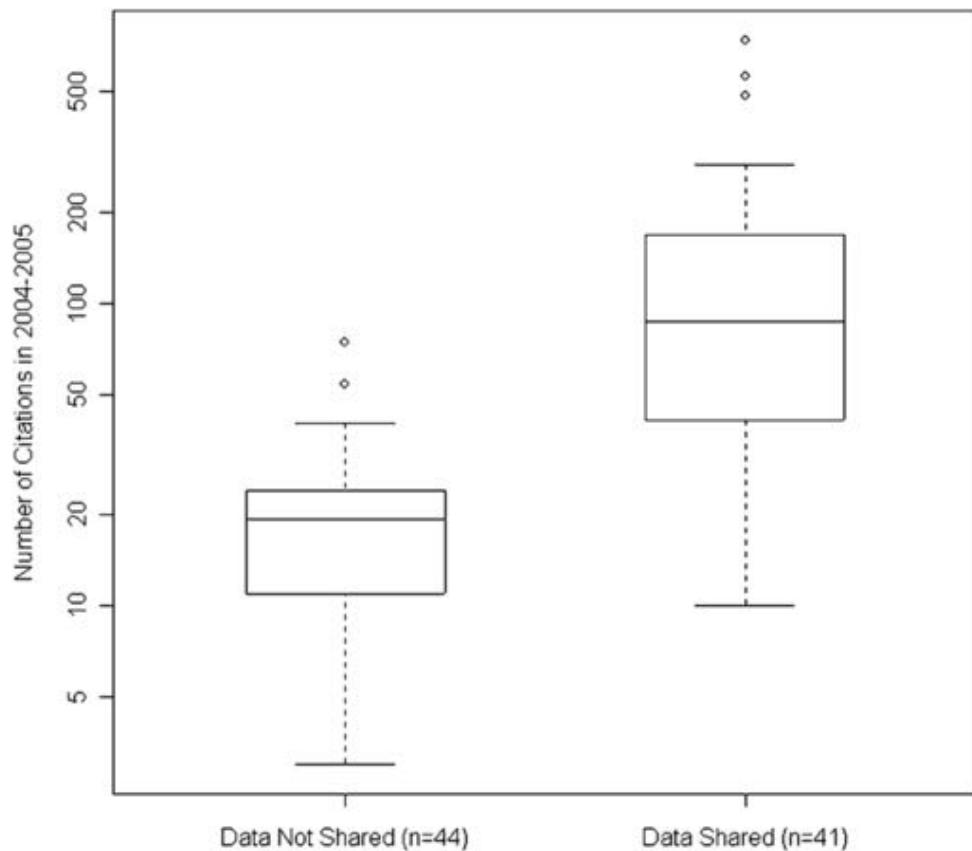
Published online 3 May 2017
10.1126/science.aan5763

- Be able to show that you have done what you say you have done
- Universities want to avoid bad press!

- To make your research easier!
- To stop yourself drowning in irrelevant stuff
- In case you need the data later
- To avoid accusations of fraud or bad science
- To share your data for others to use and learn from
- To get credit for producing it
- Because funders or your organisation require it



Well-managed data opens up opportunities for re-use, integration and new science



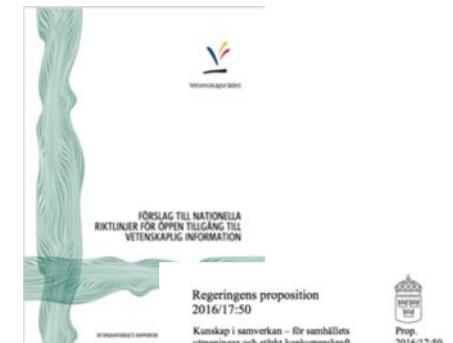
- Sharing Detailed Research Data Is Associated with Increased Citation Rate

- To make your research easier!
- To stop yourself drowning in irrelevant stuff
- In case you need the data later
- To avoid accusations of fraud or bad science
- To share your data for others to use and learn from
- To get credit for producing it
- Because funders or your organisation require it



Well-managed data opens up opportunities for
re-use, integration and new science

- *The practice of providing **on-line access** to scientific information that is **free of charge** to the end-user and that is **re-usable**.*
 - Not necessarily unrestricted access, e.g. for sensitive personal data
 - “As open as possible, as closed as necessary”
- Strong international movement towards Open Access (OA)
- European Commission recommended the member states to establish national guidelines for OA
 - 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



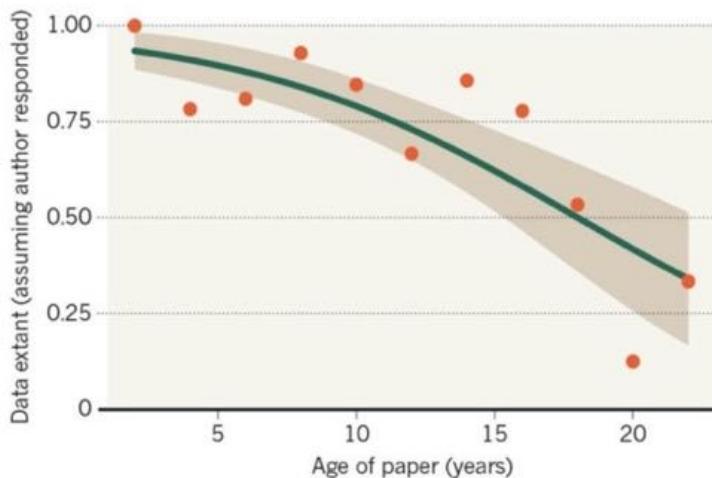
- Democracy and transparency
 - Publicly funded research data should be accessible to all
 - Published results and conclusions should be possible to check by others
- Research
 - Enables others to combine data, address new questions, and develop new analytical methods
 - Reduce duplication and waste
- 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



Data loss is real and significant, while data growth is staggering

MISSING DATA

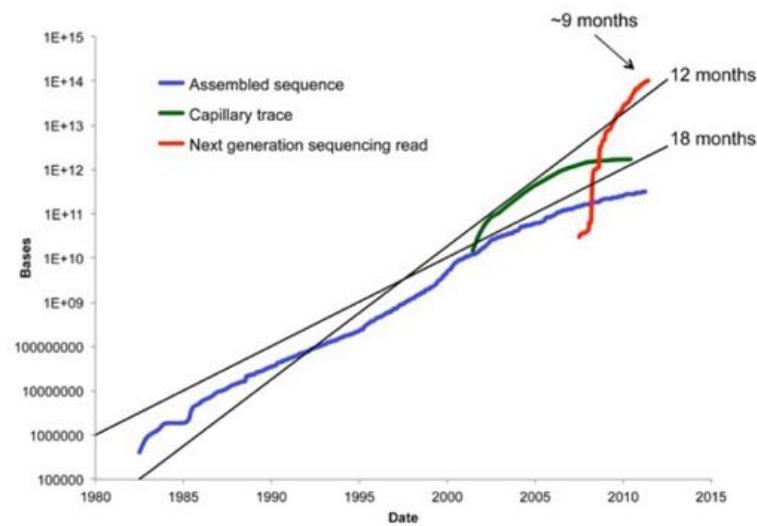
As research articles age, the odds of their raw data being extant drop dramatically.



Nature news, 19 December 2013

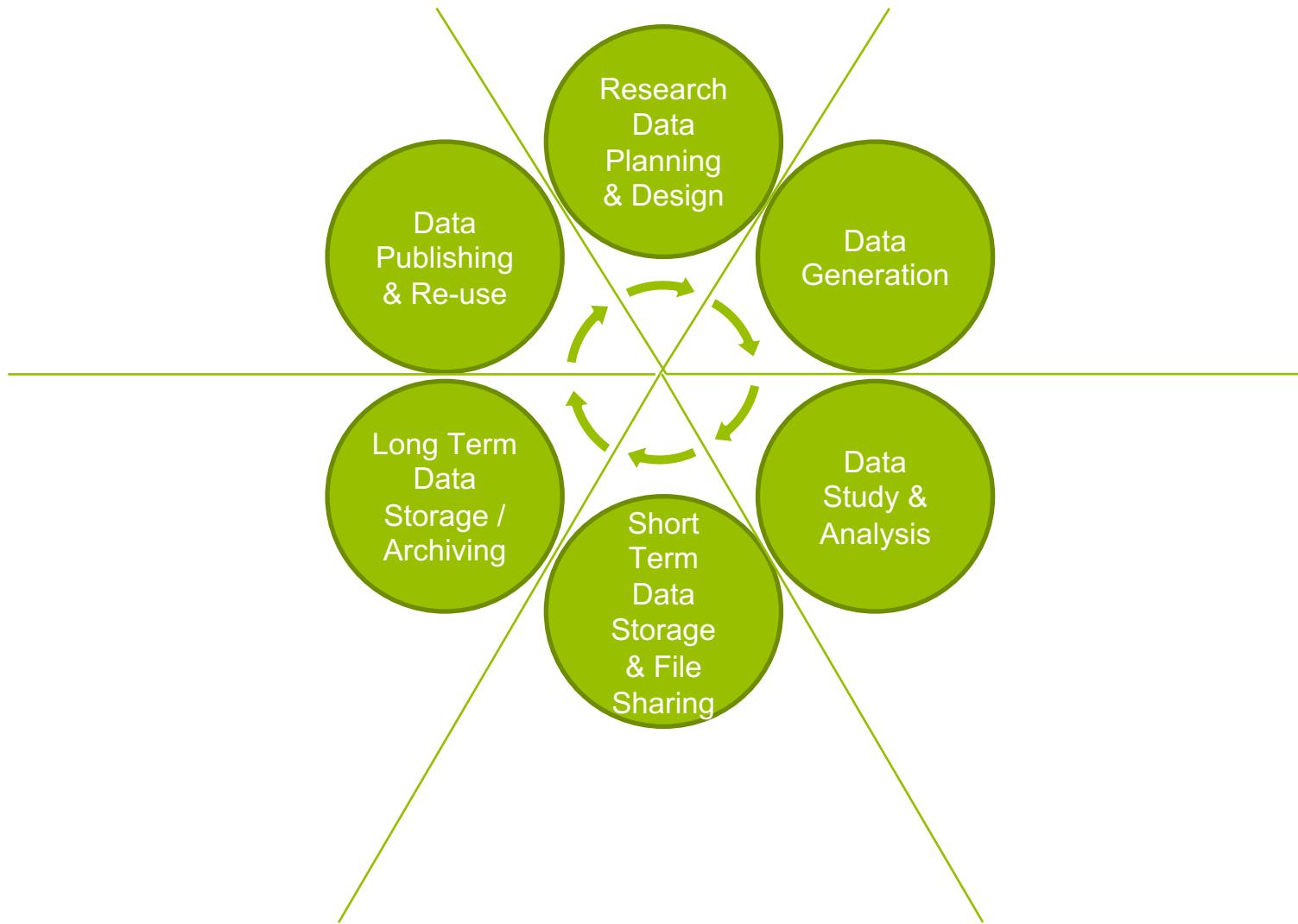


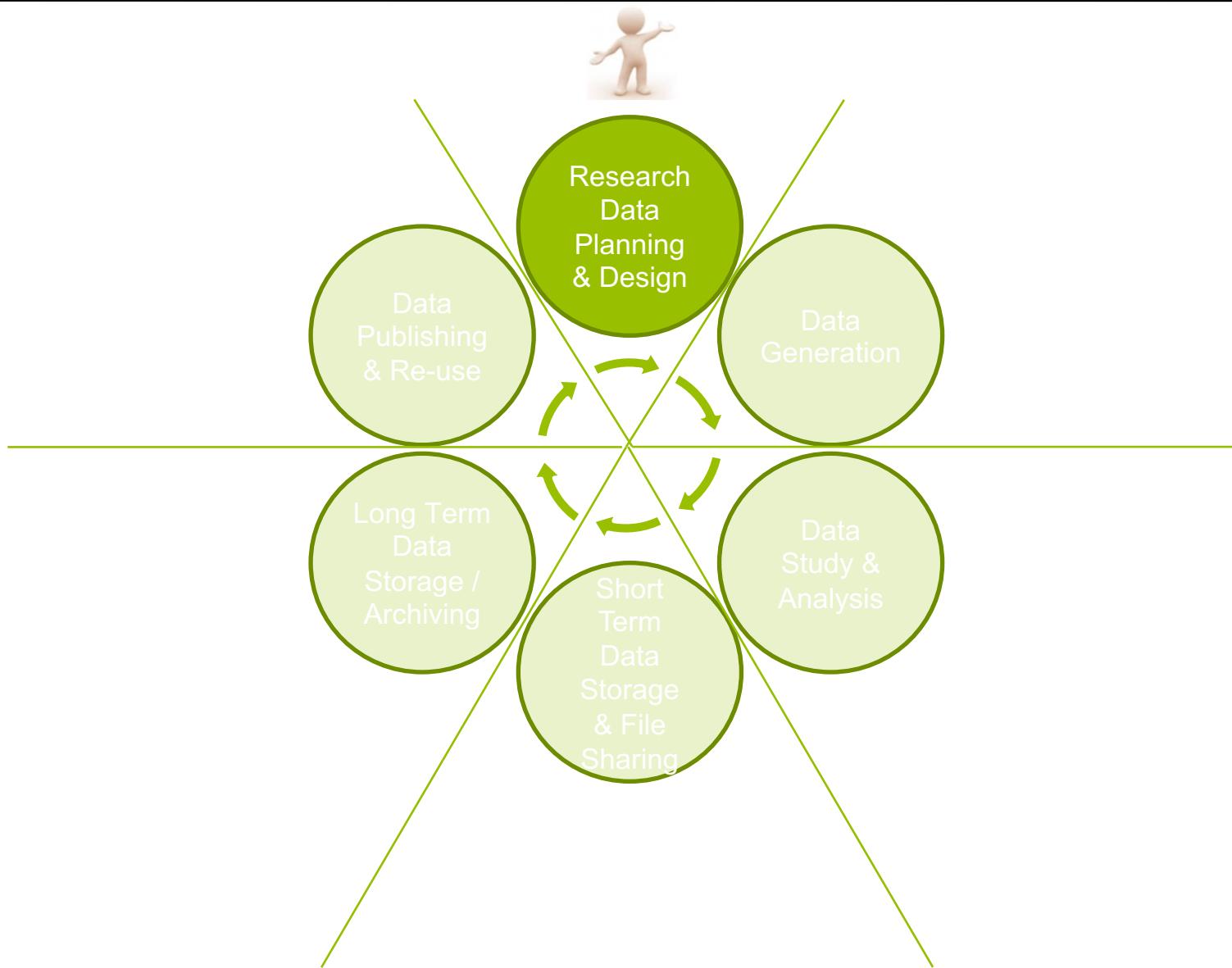
'Oops, that link was the laptop of my PhD student'



- DNA sequence data is **doubling every 6-8 months** and looks to continue for this decade
- Projected to surpass astronomy data in the coming decade

Slide stolen from Barend Mons





- Data Management planning
 - Data types
 - Sizes, where to store, etc
 - **Metadata**
 - Study, Samples, Experiments, etc
 - Use standards!



Will become a standard part of the research funding application process

- **Data collection** - data types and volumes, analysis code
- **Data organization** - folder and file structure, and naming
- **Data documentation** - data and analysis, metadata standards
- **Data storage** - storage/backup/protection & time lines
- **Data policies** - conditions/licences for using data & legal/ethical issues
- **Data sharing** - *When and How* will *What* data (and code) be shared
- **Roles and responsibilities** - who's responsible for what & is competence available
- **Budget** - People & Hardware/Software

The image shows a screenshot of a Nature journal article. At the top left is the Nature logo with the subtitle "International Journal of science". To the right are four buttons: "Search", "E-alert", "Submit", and "Login". Below the header, the text "EDITORIAL • 13 MARCH 2018" is visible. The main title of the article is "Everyone needs a data-management plan". Below the title is a quote: "They sound dull, but data-management plans are essential, and funders must explain why."

A cognitive bias in which relatively unskilled persons suffer illusory superiority, mistakenly assessing their ability to be much higher than it really is.

-Wikipedia

DMP tools

DMPonline

 [Home](#) [Public DMPs](#) [Funder requirements](#) [Help](#) [Language +](#)

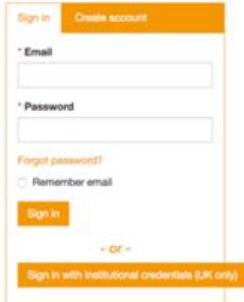
Welcome

DMPonline helps you to create, review, and share data management plans that meet institutional and funder requirements. It is provided by the Digital Curation Centre (DCC).

Join the growing international community that have adopted DMPonline:

	17,622 Users
	203 Organisations
	23,083 Plans
	89 Countries

Some funders mandate the use of DMPonline, while others point to it as a useful option. You can [download funder templates](#) without logging in, but the tool provides tailored guidance and example answers from the DCC and many research organisations. Why not sign up for an account and try it out?





<https://dmponline.dcc.ac.uk/>

ELIXIR Data Stewardship Wizard



Design of experiment

Data design and planning

Data Capture/Measurement

Data processing and curation

Data integration

Data interpretation

Information and insight

Is there any pre-existing data?

Are there any data sets available in the world that are relevant to your planned research?

No

Yes

Will you be using any pre-existing data (including other people's data)?

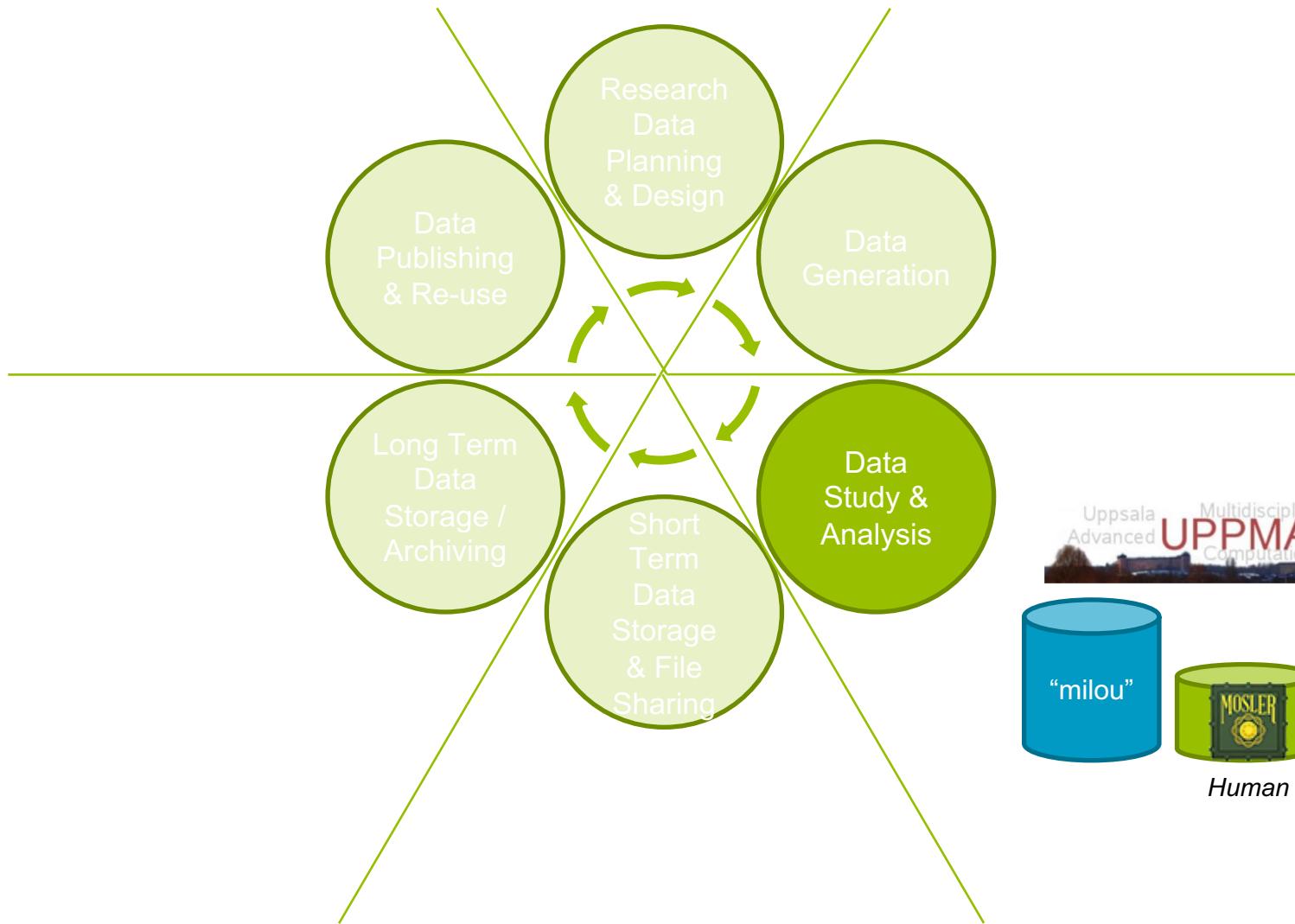
Will you be referring to any earlier measured data, reference data, or data that should be mined from existing literature? Your own data as well as data from others?

No

Yes

What reference data will you use?

<https://dsw.fairdata.solutions/>



- Guiding principle
 - “*Someone unfamiliar with your project should be able to look at your computer files and understand in detail what you did and why.*”
- Research reality
 - “*Everything you do, you will have to do over and over again*”
 - Murphy’s law



Trevor A. Branch
@TrevorABranch

Follow

My rule of thumb: every analysis you do on a dataset will have to be redone 10–15 times before publication. Plan accordingly. #Rstats



- Poor organizational choices lead to significantly slower research progress

“Your primary collaborator is yourself six months from now, and your past self doesn’t answer e-mails.”
- It is critical to make results reproducible

From bloodjournal.hematologylibrary.org by guest on September 2, 2011. For personal use only.
HEMOSTASIS, THROMBOSIS, AND VASCULAR BIOLOGY

Gene-expression patterns predict phenotypes of immune-mediated thrombosis

Anil Potti, Andrea Bild, Holly K. Cressman, Deborah A. Lewis, Joseph R. Nevins, and Thomas L. O’Connor

Antiphospholipid antibody syndrome (APS) is a complex autoimmune thrombotic disorder with defined clinical phenotypes. Although not all patients with APS develop thrombosis, those with aPLA are at high risk for complications, particularly of the cardiovascular system. Mandating aggressive and extended lifelong anti-phospholipid antibody therapy. One hundred forty-four patients (157 patients with APS and 47 with venous thromboembolism [VTE]), 32 patients with VTE without aPLA, 32 patients with aPLA only, and 8 healthy patients) were enrolled. RNA from peripheral blood collections was used for DNA microarray analysis. Patterns of gene expression that characterize APS, as well as antibodies in the presence of aPLA were identified by hierarchical clustering and binary regression methods. Gene-expression profiles were used to predict individuals with APS from patients with VTE without aPLA. Importantly, similar methods identified expression profiles that accurately predicted those patients with aPLA at high risk for thrombotic events. All profiles were tested in independent cohorts of patients. The ability to predict APS, but more importantly, those patients at risk for venous thrombosis, represents a paradigm for a genomic approach that can be applied to other populations of patients with venous thromboembolism, providing for more effective clinical management of disease, while also rejecting the possible underlying biological processes. (Blood, 2006;107:1391-1396)

© 2006 The American Society of Hematology

Retracted

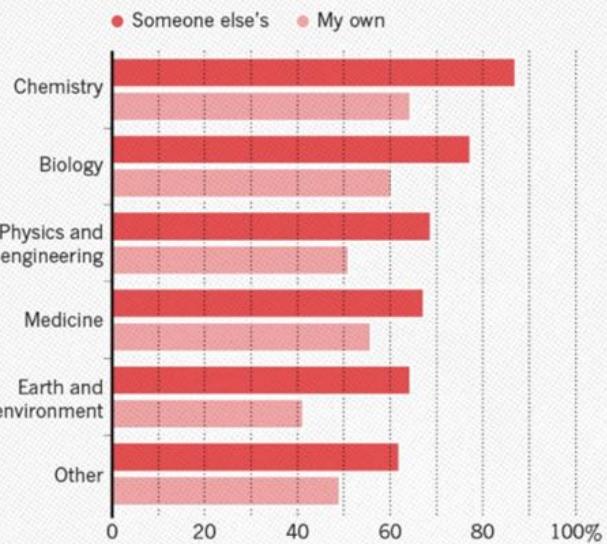


"I think you should be more explicit here in step two."

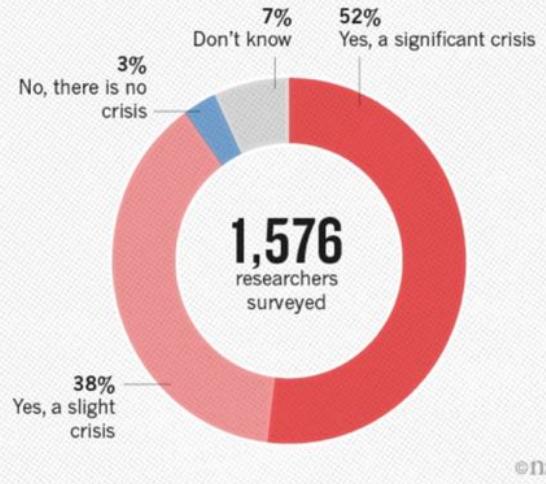
A reproducibility crisis

HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.



IS THERE A REPRODUCIBILITY CRISIS?



A recent survey in Nature revealed that irreproducible experiments are a problem across all domains of science¹.

Medicine is among the most affected research fields. A study in Nature found that 47 out of 53 medical research papers focused on cancer research were irreproducible².

Common features were failure to show all the data and inappropriate use of statistical tests.

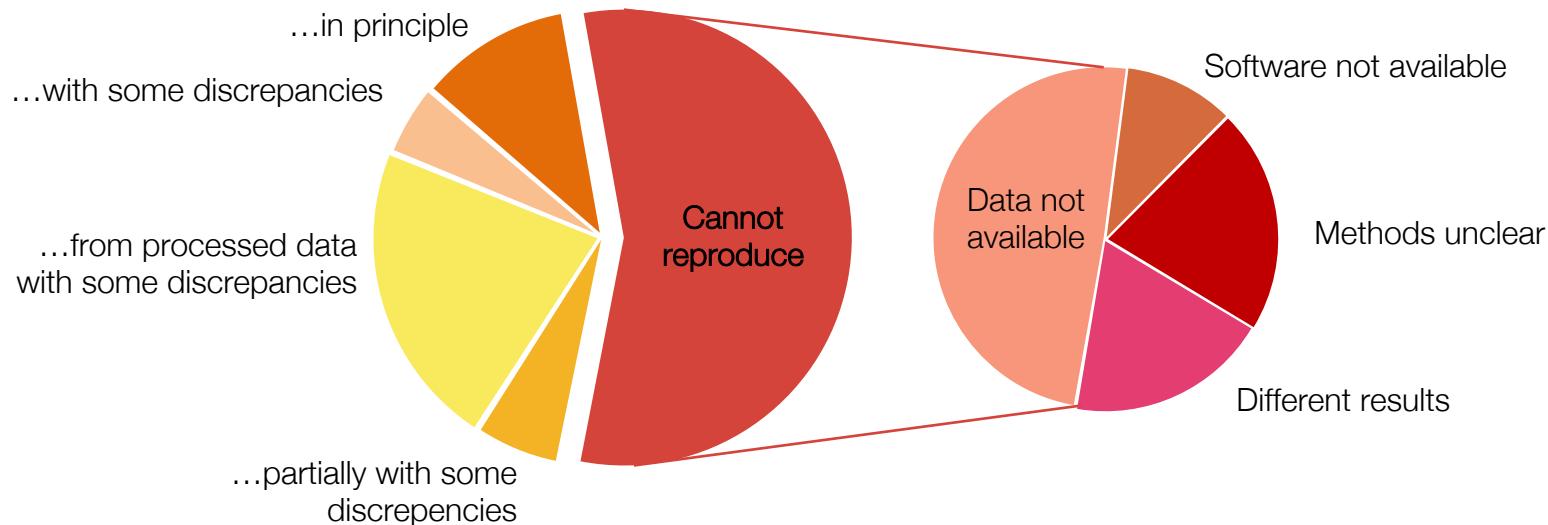
[1] "1,500 scientists lift the lid on reproducibility". Nature. 533: 452–454

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

A reproducibility crisis

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

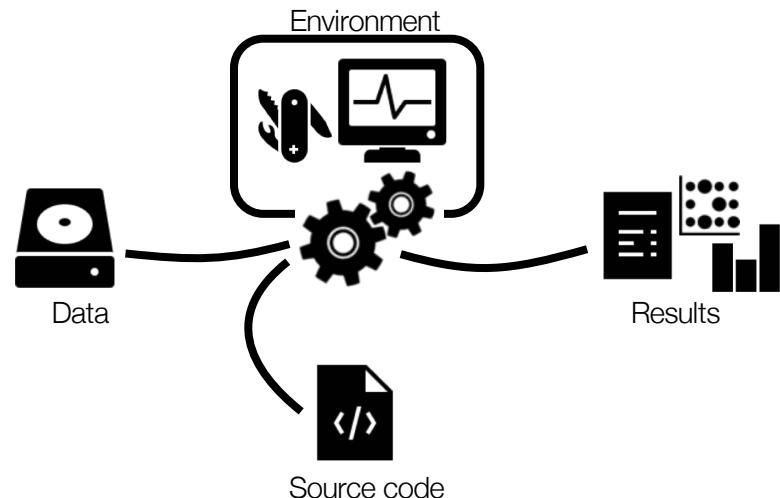
What do we mean by reproducible research?

		Data	
		Same	Different
Code	Same	Reproducible	Replicable
	Different	Robust	Generalizable

Is it really any point doing this?

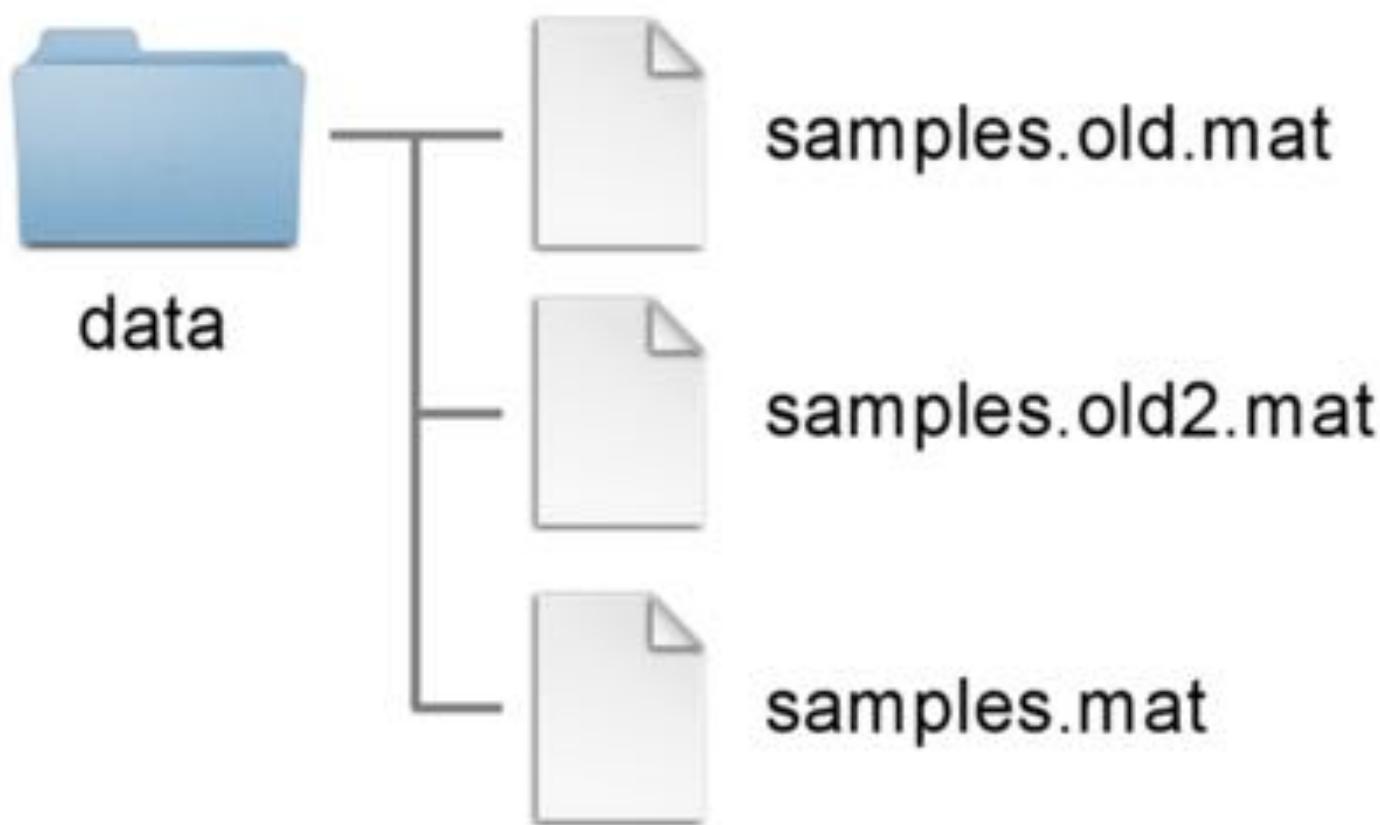
- Primarily for ones own benefit!
Organized, efficient, in control.
Dynamic team members.
- Transparent what has been done
- Some will be interested in parts of the analysis. Make it easy to redo, then adapt to own data.

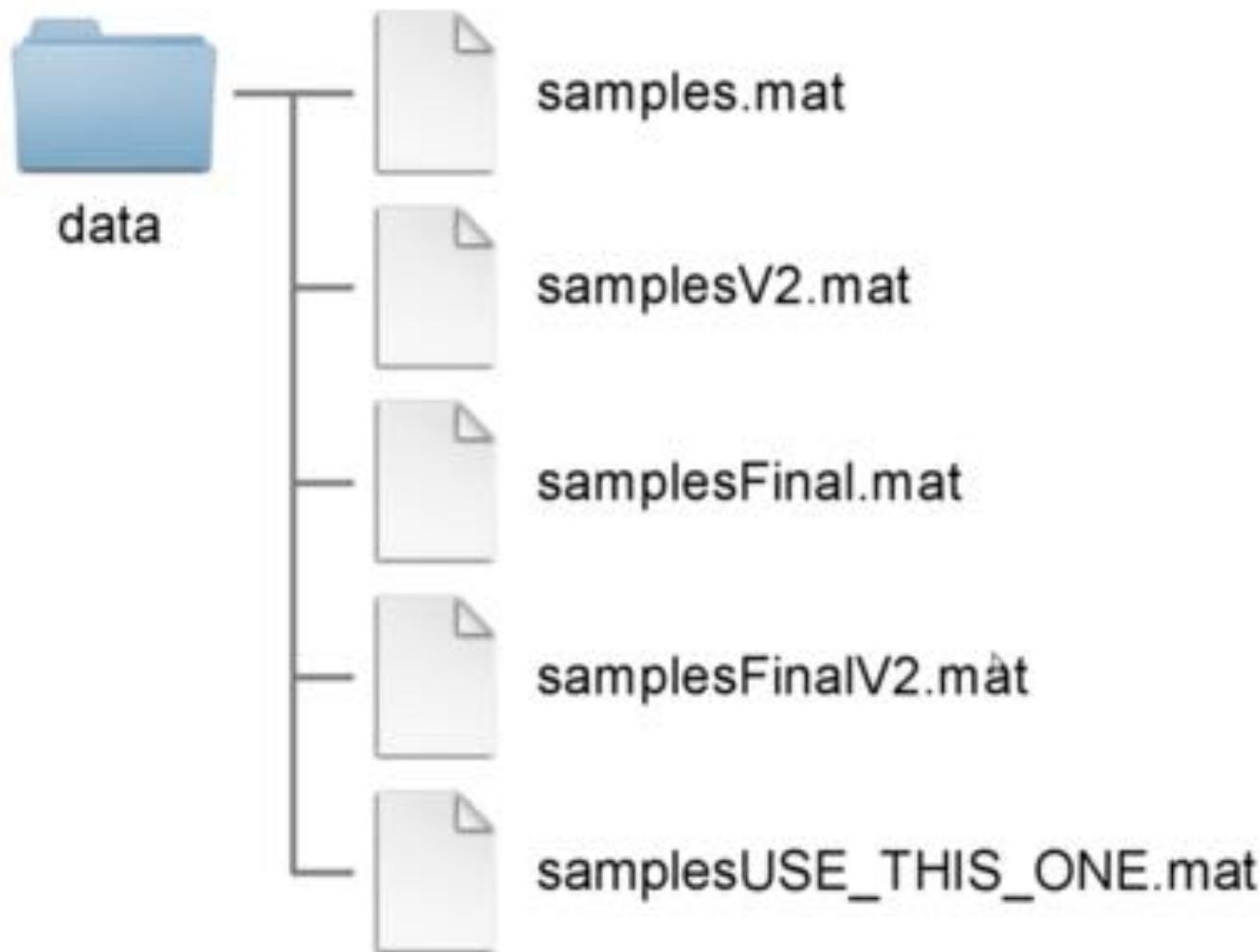
All parts of a bioinformatics analysis have to be reproducible:

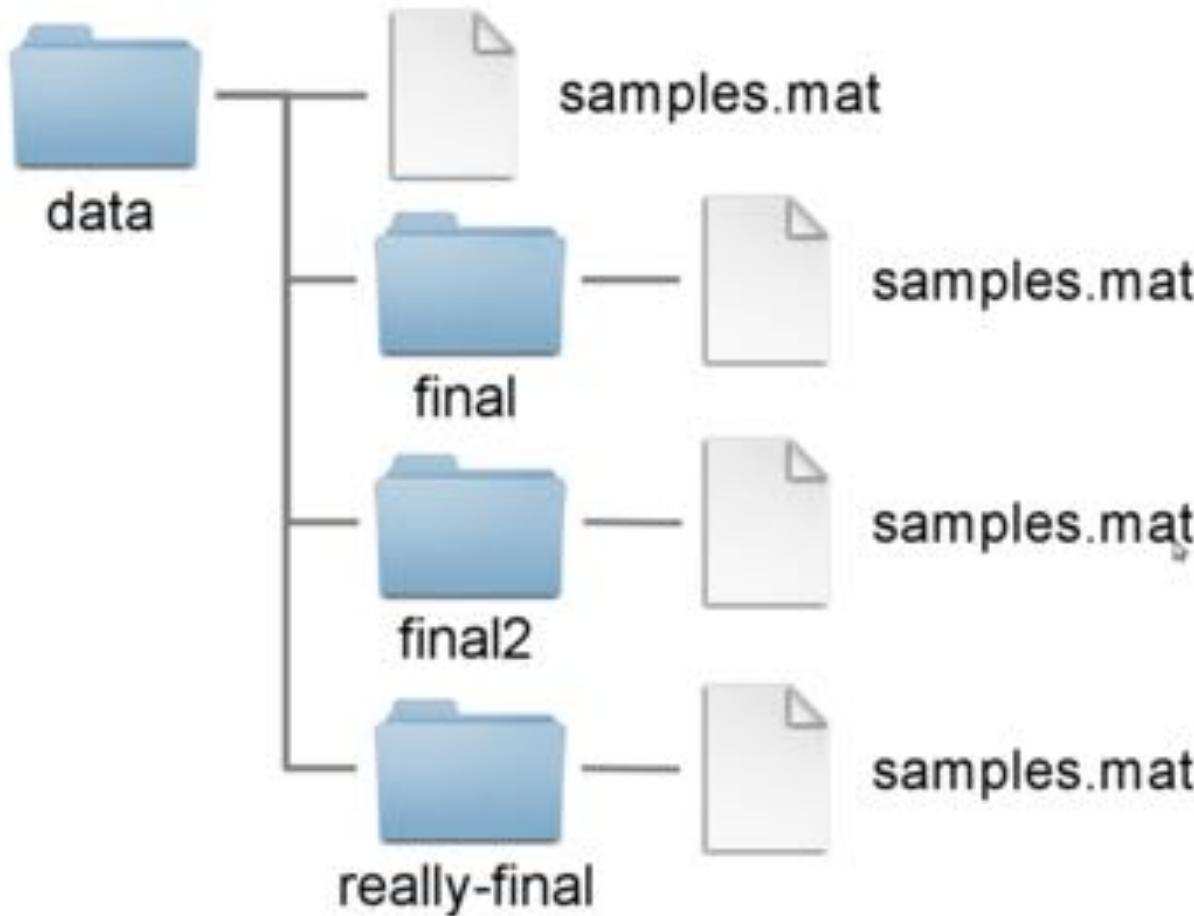










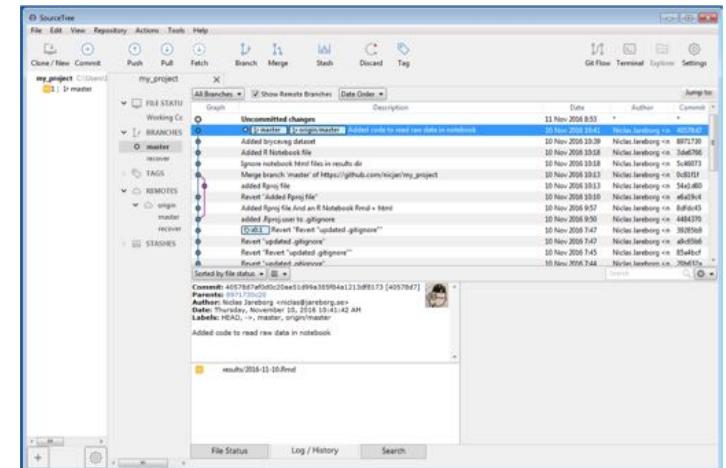


A possible solution



- There is a **folder for the raw data**, which do not get altered, or intermixed with data that is the result of manual or programmatic manipulation. I.e., derived data is kept separate from raw data, and **raw data are not duplicated**.
- **Code is kept separate from data.**
- Use a **version control system** (at least for code) – e.g. **git**
- There is a **scratch directory for experimentation**. Everything in the scratch directory can be deleted at any time without negative impact.
- There should be a **README in every directory**, describing the purpose of the directory and its contents.
- Use **non-proprietary formats** – .csv rather than .xlsx
- Etc...

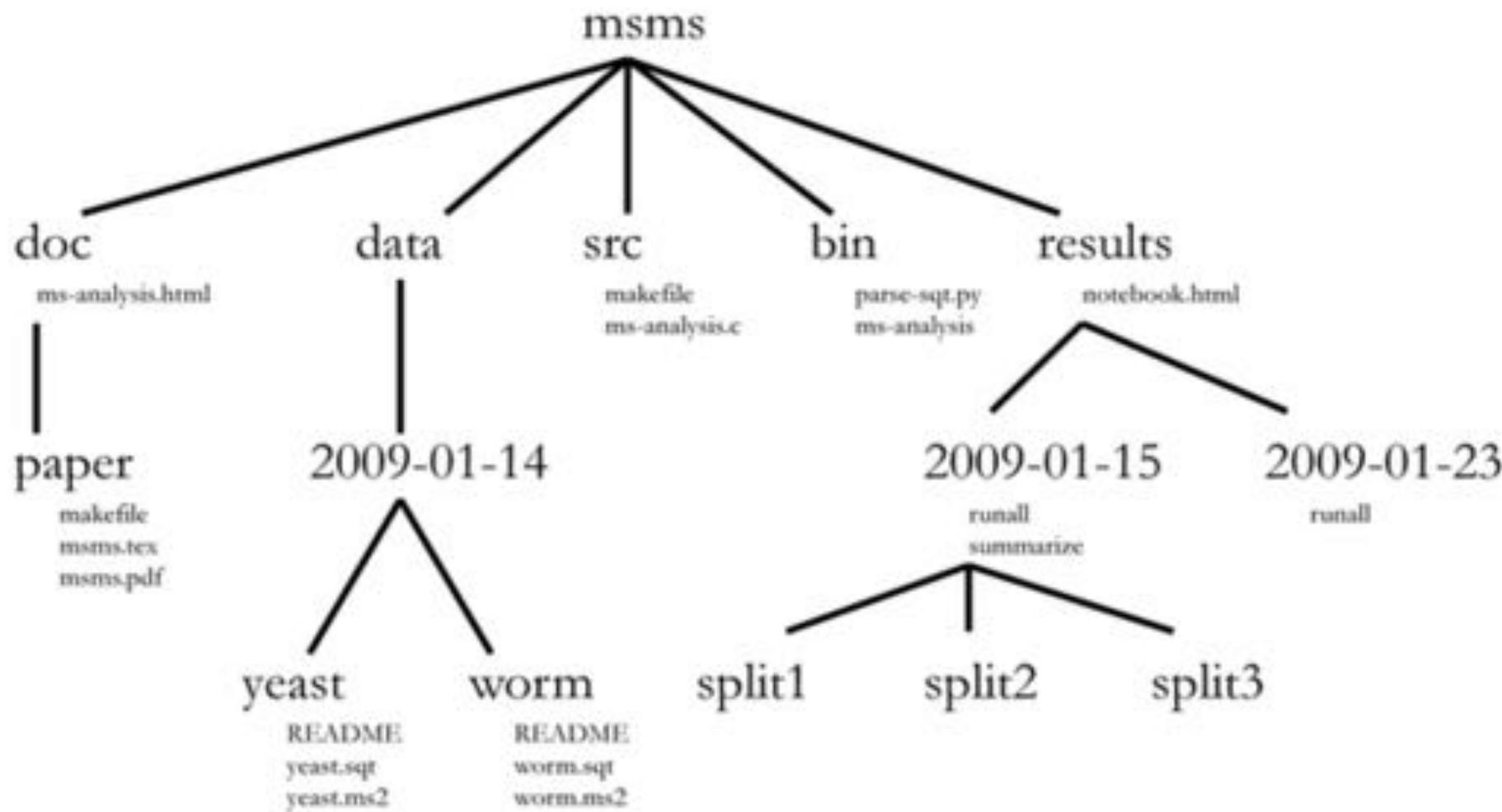
- What is it?
 - A system that keeps records of your changes
 - Allows for collaborative development
 - Allows you to know who made what changes and when
 - Allows you to revert any changes and go back to a previous state
- Several systems available
 - git, RCS, CVS, SVN, Perforce, Mercurial, Bazaar
 - git
 - Command line & GUIs
 - Remote repository hosting
 - GitHub, Bitbucket, etc



- There is a **folder for the raw data**, which do not get altered, or intermixed with data that is the result of manual or programmatic manipulation. I.e., derived data is kept separate from raw data, and **raw data are not duplicated**.
- **Code is kept separate from data.**
- Use a **version control system** (at least for code) – e.g. **git**
- There is a **scratch directory for experimentation**. Everything in the scratch directory can be deleted at any time without negative impact.
- There should be a **README in every directory**, describing the purpose of the directory and its contents.
- Use **non-proprietary formats** – .csv rather than .xlsx
- Etc...

- A text-based format is more future-safe, than a proprietary binary format by a commercial vendor
- ***Markdown*** is a nice way of getting nice output from text.
 - Simple & readable formating
 - Can be converted to lots of different outputs
 - HTML, pdf, MS Word, slides etc
- *Never, never, never use **Excel** for scientific **analysis!***
 - Script your analysis – bash, python, R, ...



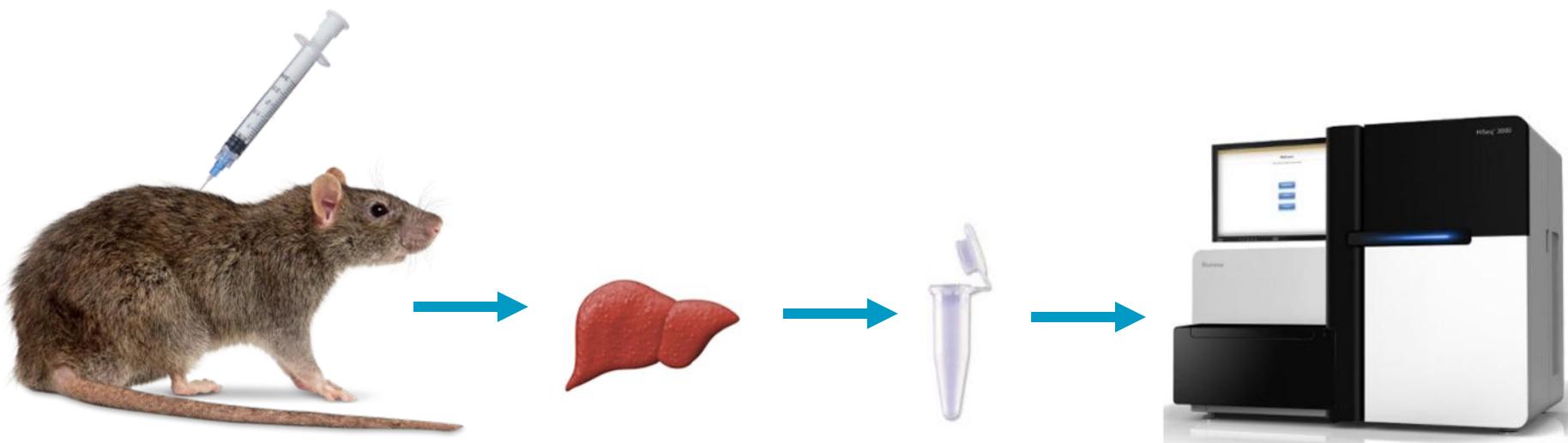


Noble WS (2009) A Quick Guide to Organizing Computational Biology Projects. PLoS Comput Biol 5(7): e1000424.
doi:10.1371/journal.pcbi.1000424

<http://journals.plos.org/ploscompbiol/article?id=info:doi/10.1371/journal.pcbi.1000424>

```
project
|- doc/           documentation for the study
|
|- data/          raw and primary data, essentially all input files, never
edit!
|   |- raw_external/
|   |- raw_internal/
|   |- meta/
|
|- code/          all code needed to go from input files to final results
|- notebooks/
|
|- intermediate/ output files from different analysis steps, can be deleted
|- scratch/       temporary files that can be safely deleted or lost
|- logs/          logs from the different analysis steps
|
|- results/
|   |- figures/
|   |- tables/
|   |- reports/
```

- Need context → document **metadata**
 - From what was the data generated?
 - How do the samples differ?
 - What where the experimental conditions?
 - Etc



- Standards
 - Controlled vocabularies / Ontologies
 - Agreed terms for different phenomena

HOW STANDARDS PROLIFERATE:
(SEE: AAC CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)



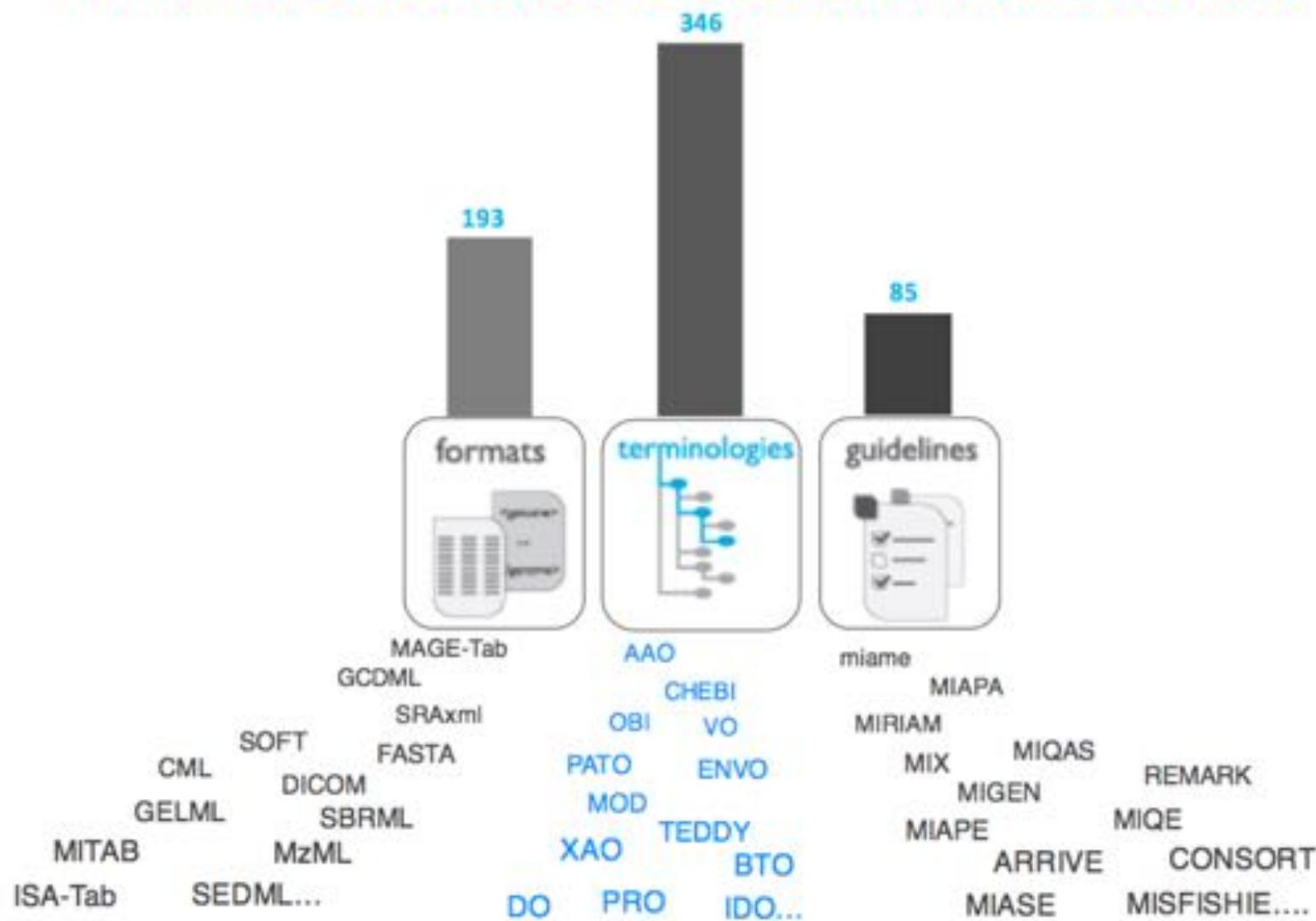
Human Phenotype Ontology

Summary Classes Properties Notes Mappings Widgets

Jump To: All Clinical modifier Abnormality of blood and blood-forming tissues Abnormal bleeding Abnormal thrombosis Abnormality of bone marrow cell morphology Abnormality of coagulation Abnormal platelet parameters Abnormality of thrombocytes Extramedullary hematopoiesis Hematological neoplasm Leukemia Acute leukemia Acute lymphocytic leukemia Acute myelomonocytic leukemia Acute monocytic leukemia **Acute myeloid leukemia** Acute myelomonocytic leukemia - Acute promyelocytic leukemia - B-lymphocyte acute leukemia Chronic leukemia Lymphoid leukemia Myeloid leukemia Myeloproliferative disorder Lymphoma Lymphoproliferative disorder

	Details	Visualization	Notes (0)	Class Mappings (21)
Preferred Name	Acute myeloid leukemia			
Synonyms	Acute myeloblastic leukemia Acute myelogenous leukemia Acute myelocytic leukemia			
Definitions	A form of leukemia characterized by overproduction of an early myeloid cell.			
ID	http://purl.obolibrary.org/obo/HP_0004808			
database_cross_reference	MedGenID:3470 UMLS:C0023467			
definition	A form of leukemia characterized by overproduction of an early myeloid cell.			
has_alternative_id	HP:0004843 HP:0001914 HP:0006728 HP:0006724 HP:0005516			
has_exact_synonym	Acute myeloblastic leukemia Acute myelogenous leukemia Acute myelocytic leukemia			
has_obo_namespace	Human_phenotype			
id	HP:0004808			
label	Acute myeloid leukemia			
notation	HP:0004808			
prefLabel	Acute myeloid leukemia			
treeView	Acute_leukemia			
subClassOf	Acute_leukemia			

In the life sciences there are >600 content standards



FAIRsharing.org
standards, databases, policies

Standards Databases Policies Collections Add/Claim Content Stats Log in or Register

A curated, informative and educational resource on data and metadata *standards*, across all disciplines, inter-related to *databases* and *data policies*.

Find

 **Recommendations**
Standards and/or databases recommended by journal or funder data policies.

Discover

 **Collections**
Standards and/or databases grouped by domain, species or organization.

Learn

 **Educational**
About standards, their use in databases and policies, and how we can help you.

 Search FAIRsharing

Search

Advanced Search

 Fine grained control over your search.

Search Wizard

 FAIRsharing
Let us guide you to your results.



699 Standards

Terminology	Artifact	343
Model/Format		239
Reporting Guideline		117

[View all](#)



974 Databases

Life Science	733
Biomedical Science	181
General Purpose	10

[View all](#)

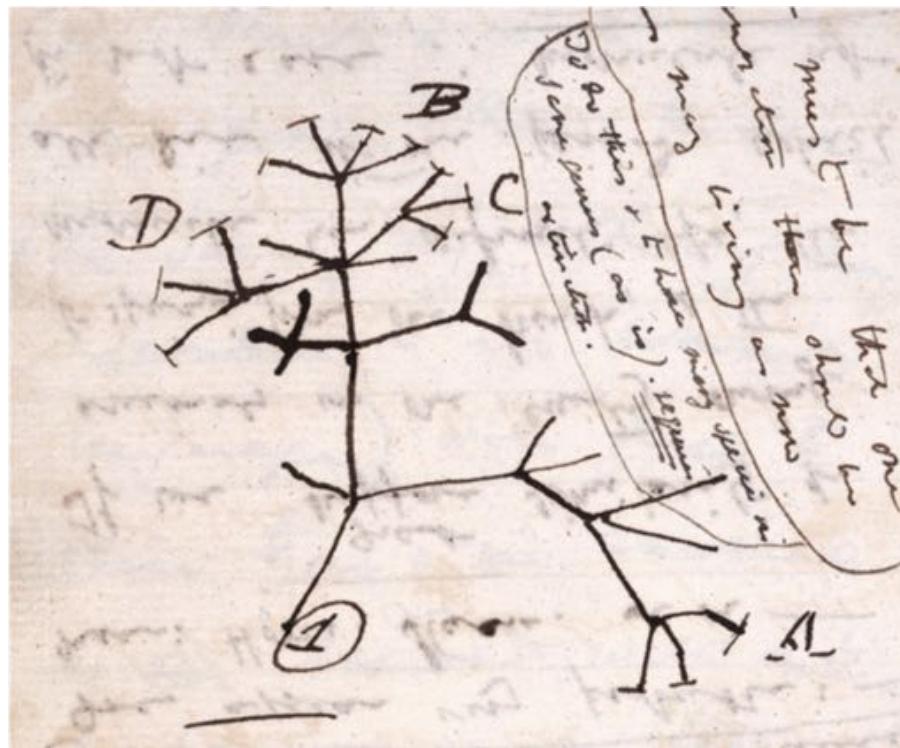


97 Policies

Funder	22
Journal	68
Society	3

[View all](#)

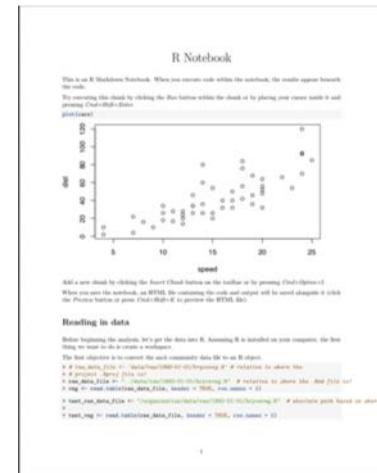
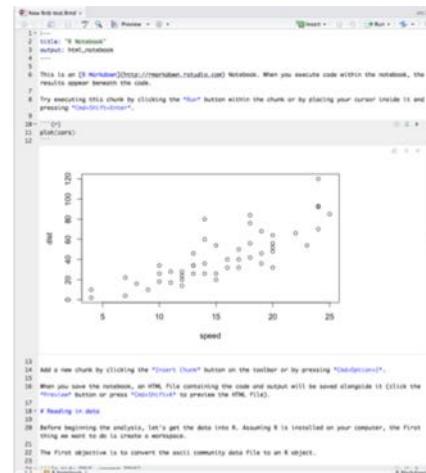
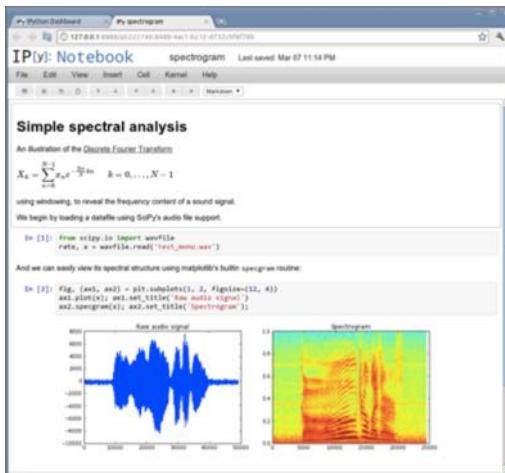
- Why?
 - You have to understand what you have done
 - **Others should be able to reproduce what you have done**



- Put in *results* directory
- *Dated* entries
- Entries relatively verbose
- Link to *data* and *code* (including versions)
 - Point to commands run and results generated
- Embedded images or tables showing results of analysis done
- Observations, Conclusions, and *ideas* for future work
- Also document analysis that *doesn't* work, so that it can be understood why you choose a particular way of doing the analysis in the end

Where to take down notes

- Paper Notebook
- Word processor program / Text files
- Electronic Lab Notebooks
- 'Interactive' Electronic Notebooks
 - e.g. [jupyter](#), [R Notebooks](#) in RStudio
 - Plain text - work well with version control (Markdown)
 - Embed and execute code
 - Convert to other output formats
 - html, pdf, word



- R Markdown makes your analysis more reproducible by connecting your code, figures and descriptive text.
- You can use it to make reproducible reports, rather than e.g. copy-pasting figures into a Word document.
- You can also use it as a notebook, in the same way as lab notebooks are used in a wet lab setting

```

1  ---
2  title: "Differential expression analysis using linear models"
3  output:
4    pdf_document
5  ---
6
7
8  ---{r setup, echo=F, message=FALSE, warning=FALSE, eval=T}
9  library(knitr)
10 opts_knit$set(root.dir=normalizePath("./"))
11 ---
12 ---{r, echo=F, message=FALSE, warning=FALSE, eval=T}
13 source("source/rmarkdown_functions.R")
14 ---
15
16 # Introduction
17 This report and analysis sets out to explore the possibility to use the blood biochemistry
18 factors to predict gene expression.
19
20
21 # Check which blood biochemistry factors can be used together in linear regression
22
23 First load the biochemistry data (this file has been cleaned up and preprocessed by the `bb_pretreatment.Rmd` script):
24 ---{r, echo=TRUE, message=FALSE, warning=FALSE}
25 bioch <- read.delim("intermediate/bb_pretreated.tsv", row.names=1, check.names=F)
26 bioch <- na.omit(bioch)
27 sampleAnno <- read.delim("intermediate/sample_annotation.tsv", stringsAsFactors=F)
28 rownames(sampleAnno) <- sampleAnno$Vcode
29 # Remove subjects that are outliers, has iron deficiency, or without blood biochemistry
30 # data (also remove secondary anemia and macrocytic anemia subjects when this info is available...)
31 samplesKeep <- sampleAnno$Vcode[with(sampleAnno, !hasIronDeficiency & !isOutlier &
32                               !hasSecondaryAnemia & !hasMacrocyticAnemia & !is.na(Date))]
33 bb <- bioch[rownames(bioch)%in%samplesKeep,]
34 ...
35
36

```

Differential expression analysis using linear models

Contents

Introduction	1
Check which blood biochemistry factors can be used together in linear regression	1
VIF scores	1
Pairwise correlation	2
Linear regression using limma	5
Correlate blood biochemistry to PCI of gene expression	5
Separate linear regression for each factor (Alt. 1)	6
Combining many factors in the same regression (Alt. 2)	16
Differential expression between Hb extreme quantiles	22
Session info	24

Introduction

This report and analysis sets out to explore the possibility to use the blood biochemistry factors to predict gene expression.

Check which blood biochemistry factors can be used together in linear regression

First load the biochemistry data (this file has been cleaned up and preprocessed by the `bb_pretreatment.Rmd` script):

```

bioch <- read.delim("intermediate/bb_pretreated.tsv", row.names=1, check.names=F)
bioch <- na.omit(bioch)
sampleAnno <- read.delim("intermediate/sample_annotation.tsv", stringsAsFactors=F)
rownames(sampleAnno) <- sampleAnno$Vcode
# Remove subjects that are outliers, has iron deficiency, or without blood biochemistry
# data (also remove secondary anemia and macrocytic anemia subjects when this info is available...)
samplesKeep <- sampleAnno$Vcode[with(sampleAnno, !hasIronDeficiency & !isOutlier &
    !hasSecondaryAnemia & !hasMacrocyticAnemia & !is.na(Date))]
bb <- bioch[rownames(bioch)%in%samplesKeep,]

```

VIF scores

The Variance Inflation Factor (VIF) score is a way to check if there is a problem of multicollinearity in a regression analysis. This is done by regressing each factor (or predictor) on the remaining ones and calculating the R^2 value. The VIF scores is then calculated as $VIF = 1/(1 - R^2)$. A VIF score of 1 means that there is no correlation among a given factor/predictor and the remaining predictor variables. A general rule of thumb is that $VIF > 4$ need further investigation, while $VIF > 10$ indicate serious multicollinearity problems.

Below we calculate VIF scores for each blood biochemistry factor, and also including age and gender since these two factors will likely be used in a linear regression model to control for such biases:



Study of velocity/energy relationship (1722-03-17, W. Gravesande)

If I drop brass balls from various heights and measure penetration depth in a block of clay, can I settle the dispute regarding conservation of energy?

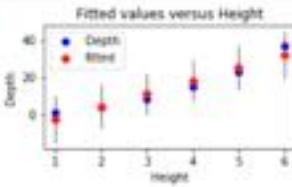
```
In [146]: import statsmodels as sm; import matplotlib.pyplot as plt; import numpy
plt.rcParams["figure.figsize"] = (4,2)
data = {"Depth": [1.1, 4.3, 6.7, 15.5, 23.2, 37.0], "Height": [1, 2, 3, 4, 5, 6]}
```

Test Bessler's hypothesis that kinetic energy, and thereby penetration depth, is linear with respect to height.

```
In [147]: res = sm.formula.api.ols(formula = 'Depth ~ Height', data = data).fit()
sm.graphics.api.plot_fit(res, 1)
```

Out[147]:

Fitted values versus Height



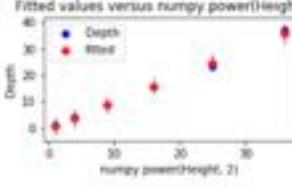
Height	Depth	Fitted
1	1.1	~1.1
2	4.3	~4.3
3	6.7	~6.7
4	15.5	~15.5
5	23.2	~23.2
6	37.0	~37.0

Okayish, but maybe Bernoulli's quadratic model fits better?

```
In [148]: res = smf.ols(formula = 'Depth ~ numpy.power(Height,2)', data = data).fit()
sm.graphics.api.plot_fit(res, 1)
```

Out[148]:

Fitted values versus numpy power(Height, 2)



Height	Depth	Fitted
1	1.1	~1.1
2	4.3	~4.3
3	6.7	~6.7
4	15.5	~15.5
5	23.2	~23.2
6	37.0	~37.0

Near!

- In-browser editing for code, with automatic syntax highlighting, indentation, and tab completion/introspection.
- The ability to execute code from the browser, with the results of computations attached to the code which generated them.

- *There's no perfect set-up*
 - Decide on **a** strategy
 - Example starting points/templates
 - <https://github.com/chendaniely/computational-project-cookie-cutter>
 - <https://github.com/Reproducible-Science-Curriculum/rr-init>
 - <https://github.com/nylander/pTemplate>
- Communicate and discuss structure and ways of working with collaborators
- Document as you go
- Done well it might reduce post-project explaining



Reproducible research for bioinformatics projects

Leif Wigge (leif.wigge@scilifelab.se)
Rasmus Ågren (rasmus.agren@scilifelab.se)
Bioinformatics long-term support (WABI)

Everything can be a project

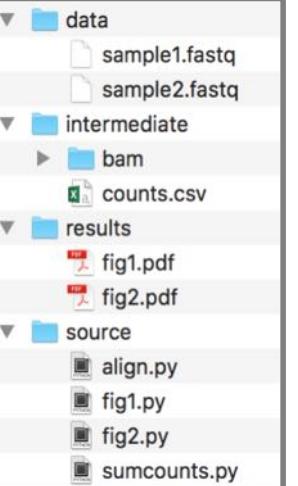
Divide your work into distinct projects and keep all files needed to go from raw data to final results in a dedicated directory with relevant subdirectories (see example).

Many software support the “project way of working”, e.g. Rstudio and the text editors Sublime Text and Atom.

Tip! Learn how to use git, a widely used system (both in academia and industry) for version controlling and collaborating on code.



<https://git-scm.com/>

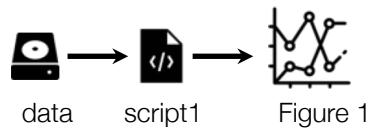


Treasure your data

- Consider your input data static. Keep it readonly!
- Don't make *different* versions. If you need to preprocess it in any way, script it so you can recreate the steps (see box below).
- Backup! Keep redundant copies in different physical locations.
- Strive towards uploading it to its final destination already at the beginning of a project (e.g. specific repositories such as ENA, or GeneExpress, or general repositories such as Dryad or Figshare).

Organize your coding

- Write scripts/functions/notebooks for specific tasks (connect raw data to final results)
- Keep parameters separate (e.g. top of file, or input arguments)



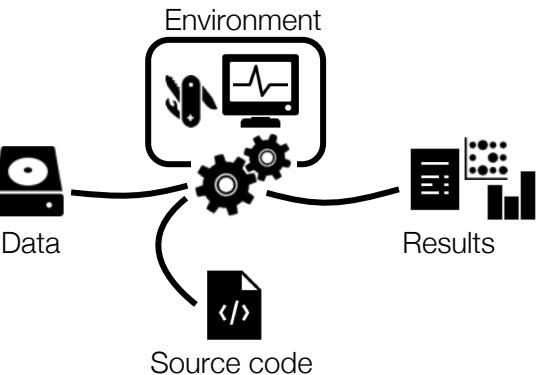
Avoid generating files interactively on the fly or doing things by hand (no way to track how they were made).



Take control of your research by making it reproducible!

By moving towards a reproducible way of working you will quickly realize that you at the same time make your own life a lot easier! The added effort pays off by gain in control, organization and efficiency.

Below are all the components of a bioinformatics project that have to be reproducible.



For the advanced

As projects grow, it becomes increasingly difficult to keep track of all the parts and how they fit together. Snakemake is a workflow management system that keeps track of how your files tie together, from raw data and scripts to final figures. If anything changes (script code, parameters, software version, etc) it will know what parts to rerun in order to have up to date and reproducible results.



Snakemake

<https://snakemake.readthedocs.io/>

Connect your results with the code

Rmarkdown and Jupyter notebooks blur the boundaries between code and its output. They allow you to add non-code text (markdown) to your code. This generates a report containing custom formatted text, as well as figures and tables together with the code that generated them.

R Markdown

<http://rmarkdown.rstudio.com/> <http://jupyter.org/>



Master your dependencies

- Full reproducibility requires the possibility to recreate the system that was originally used to generate the results.
- Conda is package, dependency, and env-ironment manager that makes it easy to install (most) software that you need for your project.
- Your environment can be exported in a simple text format and reinstalled by Conda on another system.

CONDA

<https://conda.io>

For the advanced

- Conda cannot always completely recreate the system, which is required for proper reproducibility.
- A solution is to package your project in an isolated Docker container, together with all its dependencies and libraries.
- A vision is that every new bioinformatics publication is accompanied by a publicly available Docker container!
- Singularity is an alternative to Docker which runs better on HPC clusters.



<https://www.docker.com/>



<http://singularity.lbl.gov/>

NBIS Reproducible research course

Search docs

- Welcome
- About
- The course
- Schedule
- Travel info
- Feedback

Tutorials

Introduction to the tutorials:

- Introduction to the tutorials
- The case study
- Setup
- For Mac / Linux users
- For Windows users
- The tutorials

Conda

Snakemake

Git

Jupyter

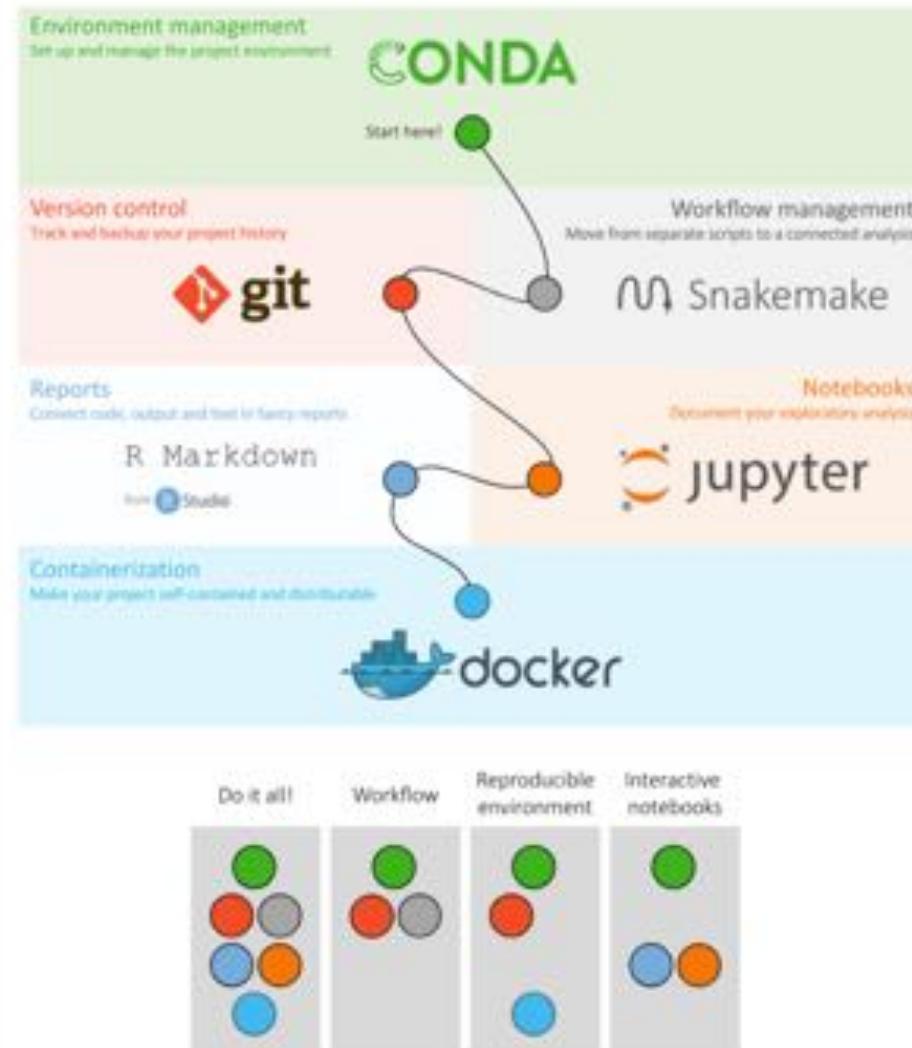
R Markdown

Docker

Take down

Documentation

Read the Docs



- Open Science Framework – <http://osf.io>
 - Organize research project documentation and outputs
 - Control access for collaboration
 - 3rd party integrations
 - Google Drive
 - Dropbox
 - GitHub
 - External links
 - Etc
 - Persistent identifiers
 - Publish article preprints

The screenshot shows the Open Science Framework (OSF) web interface. At the top, there's a navigation bar with links for 'My Dashboard', 'Browse', 'Help', and user profile information. Below the header, the title 'My fabulous project' is displayed. To the right of the title are buttons for 'Private', 'Make Public', and other sharing options. The main content area is divided into several sections: 'Wiki' (containing a 'Welcome' page with a note about being a test project), 'Citation' (with a DOI link), 'Components' (listing 'Data files' and 'Code' components, each with a progress bar indicating contributions from 'Janeborg'), and 'Tags' (with a list including 'Data management', 'Testing', and a placeholder 'add a tag'). On the left side, there's a sidebar with a tree view of the project structure, showing nodes like 'Project: My fabulous project', 'OSF Storage', 'Component: Data files', 'Component: Code', and 'GitHub: nicasj/fresco (master)'.

Personal data



- GDPR – General Data Protection Regulation (*Dataskyddsförordningen*) + others
- Act concerning the Ethical Review of Research Involving Humans (*Lag om etikprövning av forskning som avser människor*)



- All kinds of information that is directly or indirectly referable to a natural person who is alive constitute personal data
- To process personal data:
 - *All processing of personal data must fulfil the **fundamental principles** defined in the Regulation.*
 - Decide a **purpose** and stick to it
 - Only collect data that is needed
 - Don't collect more data than necessary
 - Don't use data for another incompatible purpose
 - Erase data when no longer needed
 - Ensure that data is correct and updated
 - Protect collected data – confidential and intact
 - Identify the **legal basis** for data processing before it starts
 - Inform in a transparent and honest way
- The Data Inspection Board (*Datainspektionen*) is the Swedish Data Protection Agency
 - Changing name - *Integritetsskyddsmyndigheten*

- **Consent**
- To be able to fulfil contract with data subject
- Legal obligation
- Necessary in order to protect the vital interests of the data subject
- **Public interest**
- Necessary for the purposes of the legitimate interests pursued by the controller

- Special categories (*Sensitive data*)
 - ... **racial or ethnic origin**, [...] **genetic data**, [...], data concerning **health** ... Art. 9 (1)
 - Processing is **prohibited** unless...
 - **explicit consent** is given Art. 9 (2)a
 - processing is necessary for **scientific research** in accordance with Article 89(1) based on Union or *Member State law* which shall be proportionate to the aim pursued, respect the essence of the right to data protection and provide for suitable and specific measures to safeguard the fundamental rights and the interests of the data subject. Art. 9 (2)j
 - Member State specific conditions and *limitations possible* for processing of health & genetic data Art. 9 (4)
 - **Sweden**
 - Consent?
 - Public interest → Ethical review necessary (often includes consent)

- The (legal) person that decides why and how personal data should be processed is called the **Controller** (*personuppgiftsansvarig*)
 - e.g. the employing university
 - Controller responsible for
 - Has to ensure the **rights of the individuals**
 - **Take measures** to ensure that the Regulation is followed, and be able to **show** that it is
 - **Privacy by Design** as standard
 - Keep a **register of processing**
 - Apply **security measures** when processing data
 - **Report** personal data **breaches** to the Data Protection Authority
 - Perform *Impact Assessments* and consult Data Protection Authority (when necessary)
 - Appoint **Data Protection Officer**
- The controller of personal data can delegate processing of personal data to a **Processor** (*personuppgiftsbiträde*)
 - e.g. UPPMAX/Uppsala university
 - Joint responsibility with Controller

- **A Data Protection Officer (*dataskyddssombud*)**
 - The natural person that is responsible for ensuring that the organization/company adheres to the GDPR
 - Educate
 - Audit
 - Contact point between organization and Data Protection Agency

GU

<https://medarbetarportalen.gu.se/projekt-process/aktuella-projekt/dataskyddsforordning>

KI

<https://ki.se/medarbetare/gdpr-pa-karolinska-institutet>

KTH

<https://intra.kth.se/anstallning/anstallningsvillkor/att-vara-statligt-an/behandling-av-person/dataskyddsforordningen-gdpr-1.800623>

LiU

<https://insidan.liu.se/dataskyddsforordningen/anmalan-av-personuppgiftsbehandling?l=sv>

LU

<https://personuppgifter.blogg.lu.se>

SU

<https://www.su.se/medarbetare/organisation-styrning/juridik/personuppgifter/dataskydds%C3%B6rningen>

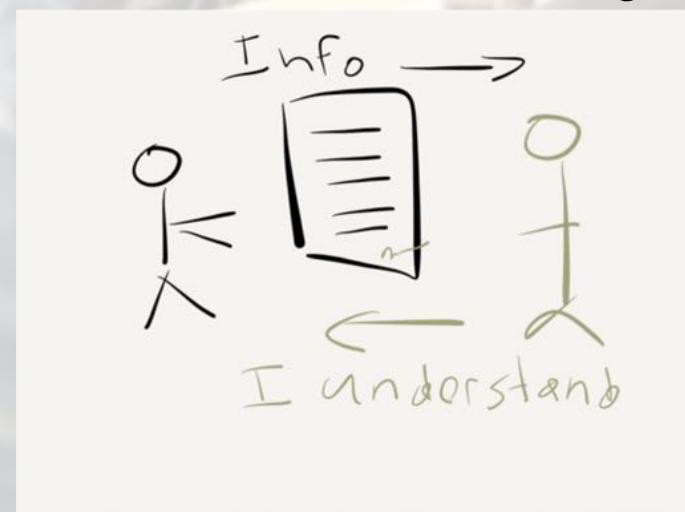
UmU

<https://www.aurora.umu.se/regler-och-riktlinjer/juridik/personuppgifter/>

UU

<https://mp.uu.se/web/info/stod/dataskyddsforordningen>

- Research that concerns studies of biological material that has been taken from a living person and that can be traced back to that person may only be conducted if it has been approved subsequent to an ethical vetting
- Informed consent
 - The subject must be informed about the purpose or the research and the consequences and risks that the research might entail
 - The subject must consent



- The genetic information of an individual is personal data
 - **Sensitive** personal data (as it relates to health)
 - Explicitly defining in GDPR
 - Even if *anonymized / pseudonymized*
 - In principle, **no** difference between WGS, Exome, Transcriptome or GWAS data
- Theoretically possible to identify the individual person from which the sequence was derived from the sequence itself
 - The more associated metadata there is, the easier this gets
 - Gymrek et al. “Identifying Personal Genomes by Surname Inference”. Science 339, 321 (2013); DOI:10.1126/science.1229566

- **Bianca**
 - Swedish Research Council funded - SNIC Sens project
 - Implemented by SNIC/UPPMAX
 - 3200 cores / 1 PB
 - Opened april 2017 <https://uppmax.uu.se/resources/systems/the-bianca-cluster/>
- **Mosler (nearing end of life)**
 - e-Infrastructure for working with sensitive data for academic research
 - Developed & operated by NBIS
 - Inspired by Norwegian solution (TSD) <https://nbis.se/infrastructure/mosler.html>
 - Designed to look like UPPMAX clusters
 - Implementation project completed Nov 2015
 - “Pilot-size system”
 - 24 nodes, 270 TB
- Provide users with a compute environment for sensitive data, with an *appropriate level of security*





Nordic Collaboration for Sensitive Data



NordForsk



<https://wiki.neic.no/tryggve>

Tryggve vision

Tryggve2 develops and facilitates access to secure e-infrastructure for sensitive data, suitable for hosting large-scale cross-border biomedical research studies

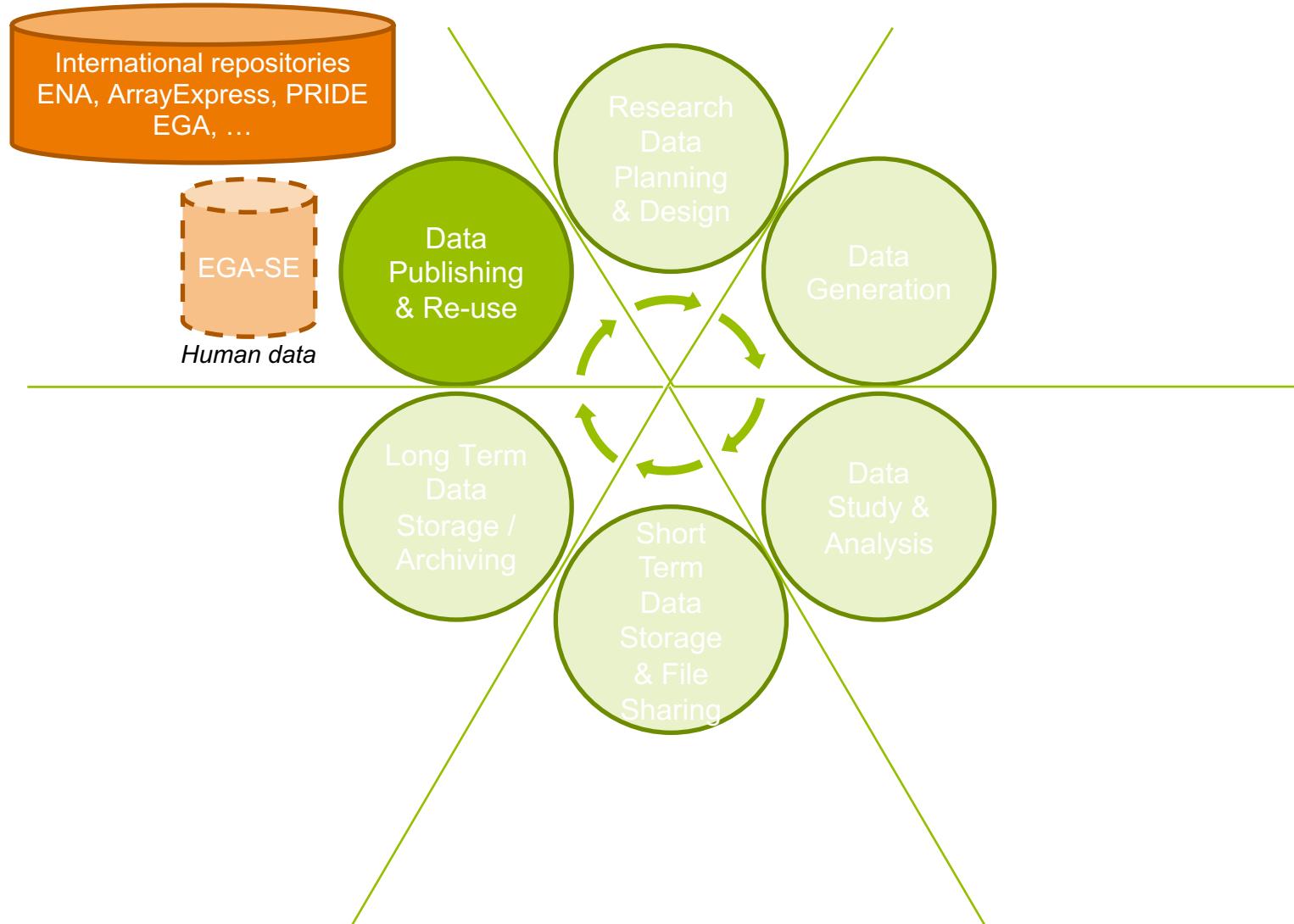


Tryggve major deliverables

1. Sensitive data archiving
2. Production quality processing services
3. Homogenized user experience
 - User mobility
 - Workflow mobility
 - Data mobility
4. **Nordic use cases**
 - **Research**
 - **Infrastructure development**
5. ELIXIR AAI
6. IT Security
7. ELSI Topics

<https://neic.no/tryggve/usecase/>





- *Research Data Publishing is a cornerstone of Open Access*
- Long-term storage
 - Data should not disappear
- Persistent identifiers
 - Possibility to refer to a dataset over long periods of time
 - Unique
 - e.g. DOIs (Digital Object Identifiers)
- Discoverability
 - Expose dataset metadata through search functionalities
- *Strive towards uploading data to its final destination already at the beginning of a project*



- DNA sequence databases: Genbank and *EMBL db* 1982
- Protein structures: *PDB* 1969

Proc. Natl. Acad. Sci. USA
 Vol. 86, p. 408, January 1989
 Data Submission

1989

Submission of data to GenBank

CHRISTIAN BURKS AND LAURIE J. TOMLINSON

Theoretical Biology and Biophysics Group T-10, MS K710, Los Alamos National Laboratory, Los Alamos, NM 87545

In response to both the ever-increasing rate of determining nucleotide sequences (1) and the growing trend among journals to allow articles to appear that describe the results of determining a sequence without explicitly presenting the sequence (1), GenBank* (2-5) and a number of the journals that publish nucleotide sequence data are working together to promote the direct, timely submission of nucleotide sequence data to GenBank. The policy being established by the PROCEEDINGS is described in the editorial on p. 407; here, we will provide a brief summary, in the context of this policy, of

Electronic file transfer. Files can be network to the network GenBank submit above. This address—in most cases with can be reached from various networks, ARPANET, USENET, JANET, JUNET, etc. / work or system expert how to send electro us for help. *Floppy disks.* We can read 8 or 5½-in diskettes written on MS-DOS s that the submitted data be written as flat t in a format specific to a given word |

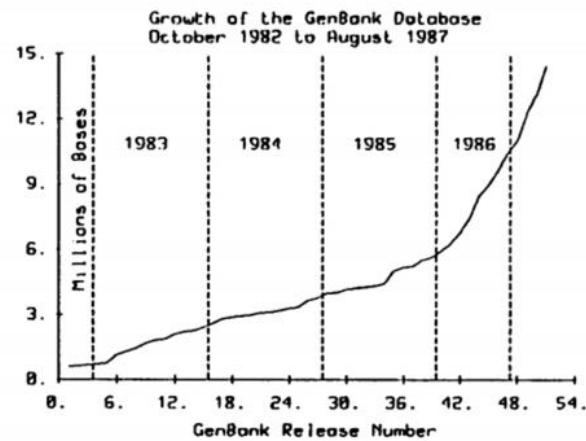


Figure 1.

"The author will provide the accession number to the PROCEEDINGS [PNAS] office to be included in a footnote to the published paper."

Bilofsky & Burks (1988)
Nucleic Acids Research v16 n5

Bermuda Principles for sharing DNA sequence data

- Automatic release of sequence assemblies larger than 1 kb (preferably within 24 hours).
- Immediate publication of finished annotated sequences.
- Aim to make the entire sequence freely available in the public domain

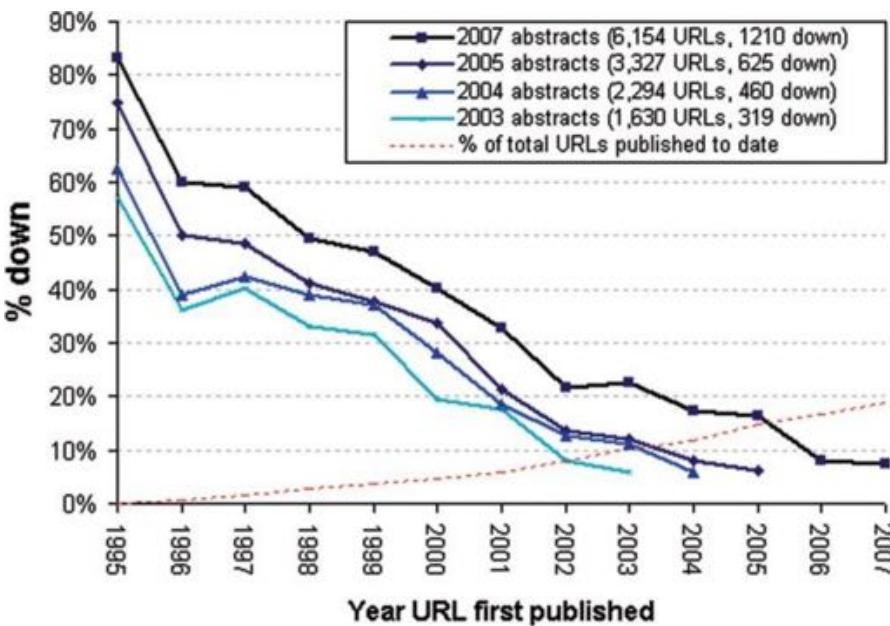


URL decay in MEDLINE—a 4-year follow-up study

Jonathan D. Wren^{*}
 Author Affiliations

^{*}To whom correspondence should be addressed.

Received January 22, 2008.
 Revision received March 11, 2008.
 Accepted April 6, 2008.



- Link rot – more 404 errors generated over time
- Reference rot* – link rot plus content drift i.e. webpages evolving and no longer reflecting original content cited

* Term coined by Hiberlink <http://hiberlink.org>

- 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>



SCIENTIFIC DATA 

OPEN **Comment: The FAIR Guiding Principles for scientific data management and stewardship**

Mark D. Wilkinson et al.*

Received: 10 December 2015
Accepted: 12 February 2016
Published: 15 March 2016

There is an urgent need to improve the infrastructure supporting the reuse of scholarly data. A diverse set of stakeholders—representing academia, industry, funding agencies, and scholarly publishers—have come together to design and jointly endorse a concise and measurable set of principles that we refer to as the FAIR Data Principles. The intent is that these may act as a guideline for those wishing to enhance the reusability of their data holdings. Distinct from peer initiatives that focus on the human scholar, the FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals. This Comment is the first formal publication of the FAIR Principles, and includes the rationale behind them, and some exemplary implementations in the community.

Supporting discovery through good data management
Good data management is not a goal in itself, but rather is the key conduit leading to knowledge discovery and innovation, and to subsequent data and knowledge integration and reuse by the community after the data publication process. Unfortunately, the existing digital ecosystem surrounding scholarly data publication prevents us from extracting maximum benefit from our research investments (e.g., ref. 1). Partially in response to this science funders, publishers and

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

G20 HANGZHOU SUMMIT

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

HANGZHOU, CHINA 4-5 SEPTEMBER 2016



- European Open Science Cloud – EOSC
 - *Enable trusted access to services, systems and the re-use of shared scientific data across disciplinary, social and geographical borders.*
 - FAIR principles are a cornerstone of EOSC



EUROPEAN COMMISSION

DIRECTORATE-GENERAL FOR RESEARCH & INNOVATION

The Director-General

Brussels, 10 July 2017

EOSC Declaration

RECOGNISING the challenges of data driven research in pursuing excellent science;

GRANTING that the vision of European Open Science is that of a research data community widely inclusive of all disciplines and Member States, sustainable in the long-term,

CONFIRMING that the implementation of the EOSC is a process, not a project, by its nature iterative and based on constant learning and mutual alignment;

UPHOLDING that the EOSC Summit marked the beginning and not the end of this process, one based on continuous engagement with scientific stakeholders, the European Commission

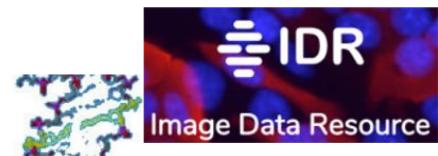
PROPOSES that all EOSC stakeholders consider sharing the following intents and actively support their implementation in the respective capacities:

Data culture and FAIR data

- [Data culture] European science must be grounded in a common culture of data stewardship that research data is recognised as a significant output of research and is appropriately curated throughout and after the period conducting the research. Only a considerable cultural change can enable long-term reuse for science and for innovation of data created by research activities in disciplines, institutions or countries must be left behind.
- [Open access by-default] All researchers in Europe must enjoy access to an open-by-default, efficient and cross-disciplinary research data environment supported by FAIR data principles. Open access must be the default setting for all results of publicly funded research in Europe, allowing for proportionate limitations only in duly justified cases of personal data protection, confidentiality, IPR concerns, national security or similar (e.g. 'as open as possible and as closed as necessary').
- [Skills] The necessary skills and education in research data management, data stewardship and data science should be provided throughout the EU as part of higher education, the training system and on-the-job best practice in the industry. University associations, research organisations, research libraries and other educational brokers play an important role but need substantial support from the European Commission and the Member States.
- [Data stewardship] Researchers need the support of adequately trained data stewards. European Commission and Member States should invest in the education of data stewardship career programmes delivered by universities, research institutions and other trans-European agents.
- [Rewards and incentives] Rewarding research data sharing is essential. Researchers who make research data open and FAIR for reuse and/or reuse and reproduce data should be rewarded.



dbSNP
Short Genetic Variations



- Best way to make data FAIR
- Domain-specific metadata standards

Deposition Database	Data type	International collaboration framework ¹	Deposition Database	Data type	International collaboration framework ¹
ArrayExpress	Functional genomics data. Stores data from high-throughput functional genomics experiments.		PDBe	Biological macromolecular structures.	wwPDB
BioModels	Computational models of biological processes.		PRIDE	Mass spectrometry-based proteomics data, including peptide and protein expression information (identifications and quantification values) and the supporting mass spectra evidence.	The ProteomeXchange Consortium
EGA	Personally identifiable genetic and phenotypic data resulting from biomedical research projects.	European Bioinformatics Institute and the Centre for Genomic Regulation		Pending incorporation into a Node Service Delivery Plan (see How countries join):	
ENA	Nucleotide sequence information, covering raw sequencing data, contextual data, sequence assembly information and functional and taxonomic annotation.	International Nucleotide Sequence Database Collaboration	BioSamples	BioSamples stores and supplies descriptions and metadata about biological samples used in research and development by academia and industry.	NCBI BioSamples database
IntAct	IntAct provides a freely available, open source database system and analysis tools for molecular interaction data.	The International Molecular Exchange Consortium	BioStudies	Descriptions of biological studies, links to data from these studies in other databases, as well as data that do not fit in the structured archives.	
MetaboLights	Metabolite structures and their reference spectra as well as their biological roles, locations and concentrations, and experimental data from metabolic experiments.		EVA	The European Variation Archive covers genetic variation data from all species.	dbSNP and dbVAR
			EMDB	The Electron Microscopy Data Bank is a public repository for electron microscopy density maps of macromolecular complexes and subcellular structures.	

<https://www.elixir-europe.org/platforms/data/elixir-deposition-databases>

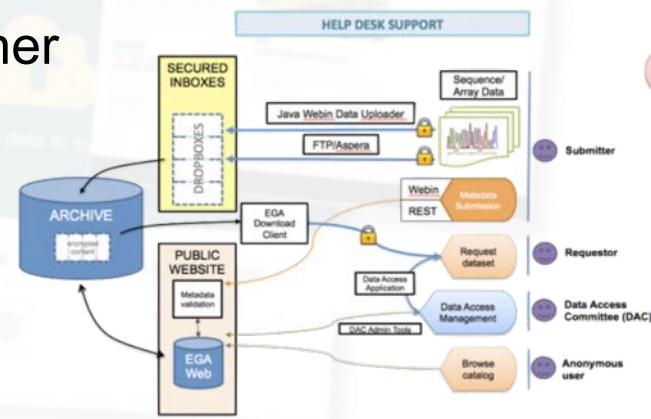
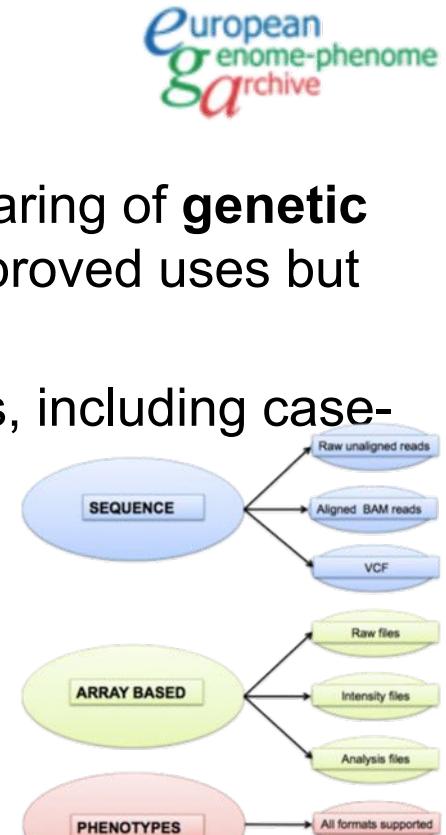
Surprisingly few submit to international repositories

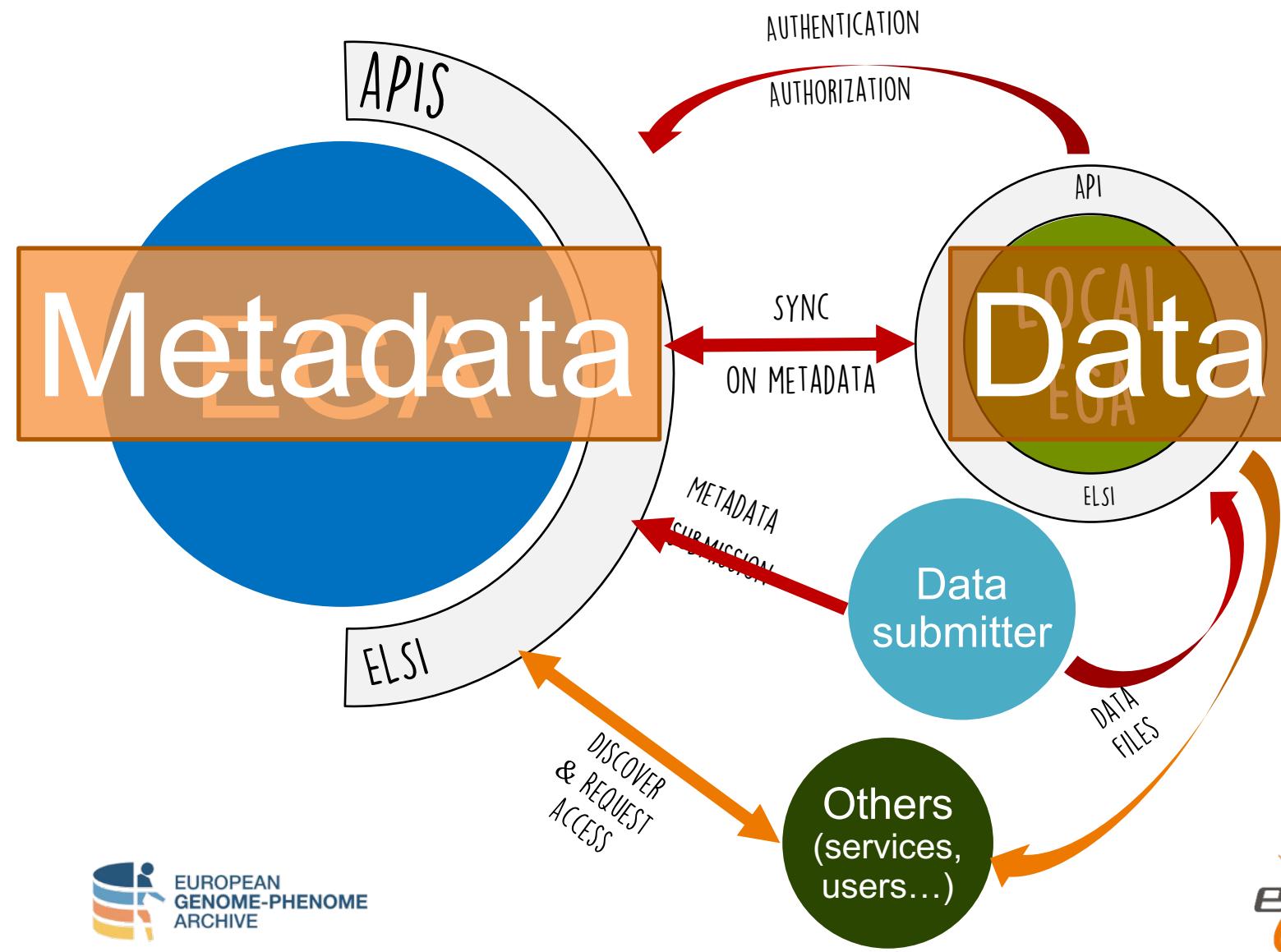
- NIH funded research
 - Only 12% of articles from NIH-funded research mention data deposited in international repositories
 - Estimated 200000+ “invisible” data sets / year

Read et al. “Sizing the Problem of Improving Discovery and Access to NIH-Funded Data: A Preliminary Study” (2015)

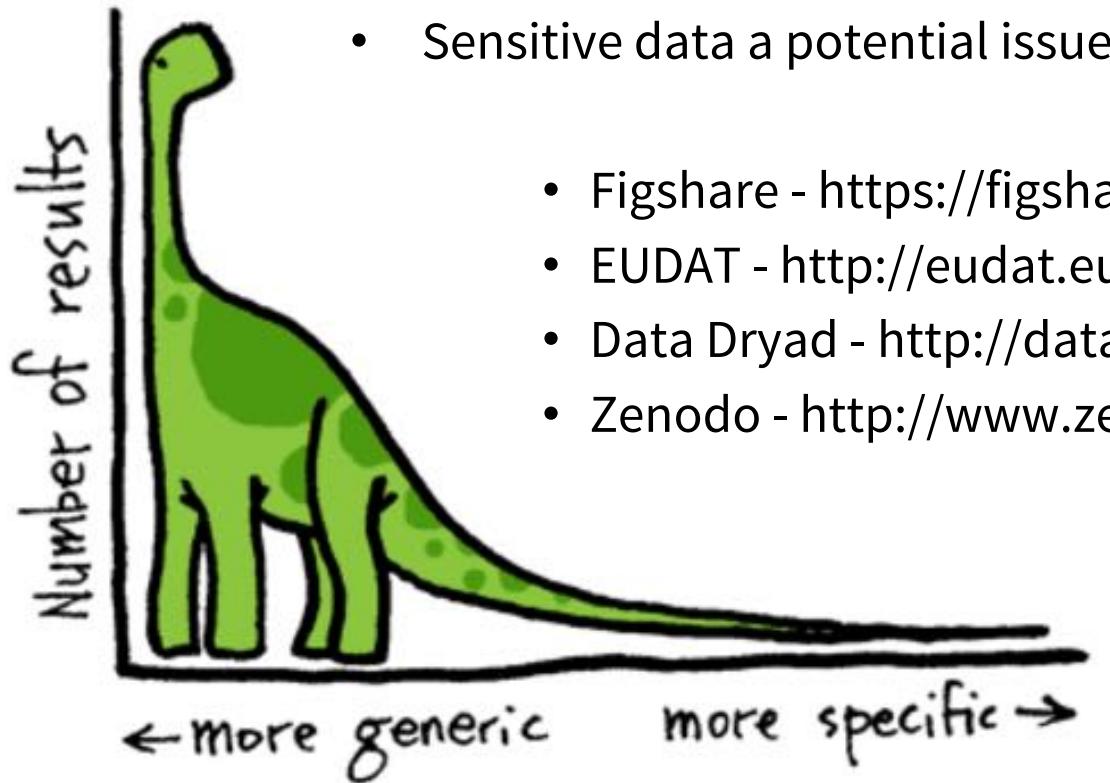
PLoS ONE 10(7): e0132735. doi: 10.1371/journal.pone.0132735

- **EGA – European Genome-phenome Archive**
 - Repository that promotes the distribution and sharing of **genetic and phenotypic data** consented for specific approved uses but **not fully open, public distribution**.
 - All types of sequence and genotype experiments, including case-control, population, and family studies.
- Data Access Agreement
 - Defined by the data owner
- Data Access Committee – DAC
 - Decided by the data owner





- Research data that doesn't fit in structured data repositories
- Data publication – persistent identifiers
- Metadata submission – not tailored to Life Science
 - *Affects discoverability*
 - *(Less) FAIR*
- Sensitive data a potential issue



- ORCID is an open, non-profit, community-driven effort to create and maintain a registry of unique researcher identifiers and a transparent method of linking research activities and outputs to these identifiers.
- <http://orcid.org>
- Persistent identifier for you as a researcher

The screenshot shows the ORCID profile page for Niclas Jareborg. At the top, there's a navigation bar with links for 'FOR RESEARCHERS', 'FOR ORGANIZATIONS', 'ABOUT', 'HELP', and 'SIGN IN'. Below the navigation bar, it says '3,035,272 ORCID IDs and counting' with a 'See more...' link. The main content area displays two education entries under the heading 'Education (2)'. Each entry includes the institution ('Uppsala Universitet: Uppsala, Sweden'), the period ('1989-05 to 1995-05 (Microbiology)' or '1985-01 to 1989-04 (Microbiology)'), the degree ('PhD' or 'BSc'), the source ('Niclas Jareborg'), and the creation date ('Created: 2015-04-09'). Below these, there are two employment entries under the heading 'Employment (7)'. Each entry includes the institution ('Stockholms Universitet: Stockholm, Sweden' or 'Kungliga Tekniska Hogskolan: Stockholm, Sweden'), the period ('2015-01 to present (BALS / Department of Biochemistry and Biophysics)' or '2013-01 to 2014-12 (National Genomics Infrastructure / SciLifeLab)'), the role ('Data Manager'), the source ('Niclas Jareborg'), and the creation date ('Created: 2015-03-23'). On the left side of the profile, there are sections for 'Also known as' (listing 'C. J. E. Niclas Jareborg, N Jareborg'), 'Country' (listing 'Sweden'), and 'Websites' (listing 'LinkedIn' and 'Personal home page').

Niclas Jareborg

ORCID ID
orcid.org/0000-0002-4520-044X

Also known as
C. J. E. Niclas Jareborg, N Jareborg

Country
Sweden

Websites
LinkedIn
Personal home page

Education (2)

Uppsala Universitet: Uppsala, Sweden
1989-05 to 1995-05 (Microbiology)
PhD
Source: Niclas Jareborg
Created: 2015-04-09

Uppsala Universitet: Uppsala, Sweden
1985-01 to 1989-04 (Microbiology)
BSc
Source: Niclas Jareborg
Created: 2015-04-09

Employment (7)

Stockholms Universitet: Stockholm, Sweden
2015-01 to present (BALS / Department of Biochemistry and Biophysics)
Data Manager
Source: Niclas Jareborg
Created: 2015-03-23

Kungliga Tekniska Hogskolan: Stockholm, Sweden
2013-01 to 2014-12 (National Genomics Infrastructure / SciLifeLab)

- Consider doing a Data Management Plan for your project
 - How do you ensure that your research output is FAIR?
- Consider submitting "raw data" to public repositories as early as possible
- Organize project metadata from the start
 - In ways that makes it easy to submit to public repositories
 - Use available standards
- Pick a thought-through file and folder structure organization for your computational analyses
- Strive for reproducibility
 - Data & Code
- Be aware that there are legal aspects to processing human data
- *Ask for help if you need it!*

- Research Data Management, EUDAT - <http://hdl.handle.net/11304/79db27e2-c12a-11e5-9bb4-2b0aad496318>
- Barend Mons – FAIR Data
- Antti Pursula – Tryggve <https://neic.no/tryggve/>
- Noble WS (2009) [A Quick Guide to Organizing Computational Biology Projects. PLoS Comput Biol 5\(7\): e1000424. doi:10.1371/journal.pcbi.1000424](https://doi.org/10.1371/journal.pcbi.1000424)
- Reproducible research
 - Reproducible Science Curriculum – <https://github.com/Reproducible-Science-Curriculum/rr-init>
 - Leif Väremo & Rasmus Ågren
 - https://bitbucket.org/scilifelab-lts/reproducible_research_example/src
 - https://nbis-reproducible-research.readthedocs.io/en/course_1803
- GDPR
 - Datainspektionen – <https://www.datainspektionen.se/lagar-regler/dataskyddsforordningen/>
 - Regina Becker, ELIXIR Luxemburg
- ... and probably others I have forgotten