

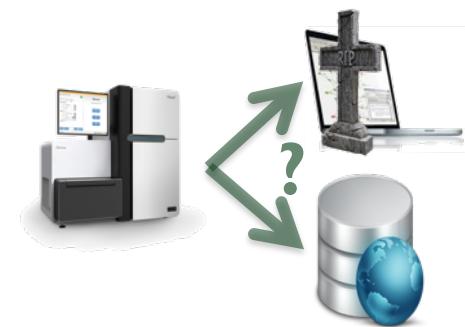
---

# Research Data Management

Niclas Jareborg, NBIS  
[niclas.jareborg@nbis.se](mailto:niclas.jareborg@nbis.se)

*Introduction to NGS course, 2017-10-27*

- 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.

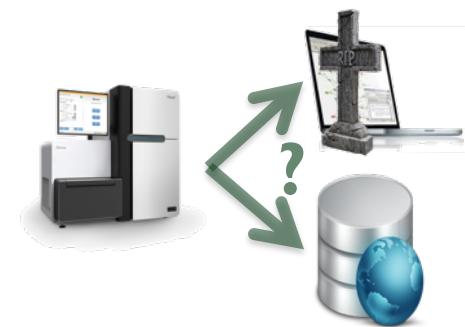
### REFERENCES

1. O. M. Lönnstedt, P. Eklöv, *Science* **352**, 1213 (2016).
2. J. Berg, *Science* **354**, 1242 (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

- *The practice of providing on-line access to scientific information that is free of charge to the end-user and that is re-usable.*
  - Does not necessarily mean unrestricted access, e.g. for sensitive personal data
- 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]



Propositionens huvudsakliga innehåll

I propositionen presenteras regeringens syn på forskningspolitiken inriktning i ett långtgående perspektiv, med särskilt fokus på seningen 2017–2020. Syftet är att ge en nationell handlingsplan för forskningspolitiken och dess tillämpning i praktiken.

En utgångspunkt är att varan den fria forskningens betydelse för landet och dess medborgare. Detta innebär att forskning och utveckling ska bidra till att skapa värde för samhället och till att öka konkurrenskraften. För att detta ska ske krävs att forskningen är tillgänglig för alla och att den används till förmån för samhället och dess medborgare. En annan utgångspunkt är att forskningen och utvecklingen avser att samordna resurserna.

Forskningspolitiken ska i framtiden utgöra en integrerad del i nationell strategisk planering och styrkas genom ökad anslag för forskning och utveckling. Detta ska göras genom att förstärka nationell forskning, att utveckla samhälls- och handindustriell forskning och att utveckla forskningsmiljöer. Detta ska göras genom att utveckla nationell forskning och att utveckla forskningsmiljöer. Detta ska göras genom att utveckla nationell forskning och att utveckla forskningsmiljöer.

Regeringen har i budgetbeslutet för 2017 lämnat förlag och överlämnat till riksdagen för att bli lag. Den 24 oktober 2016 beslutade riksdagen om att godkänna propositionen. Den 24 oktober 2016 beslutade riksdagen om att godkänna propositionen. Den 24 oktober 2016 beslutade riksdagen om att godkänna propositionen.

Sammanfattningsvis är det viktigt att se att propositionen är en del av en strukturering av nationell forskningspolitiken och att den är en del av en strukturering av nationell forskningspolitiken. Sammanfattningsvis är det viktigt att se att propositionen är en del av en strukturering av nationell forskningspolitiken och att den är en del av en strukturering av nationell forskningspolitiken.

# Why Open Access?

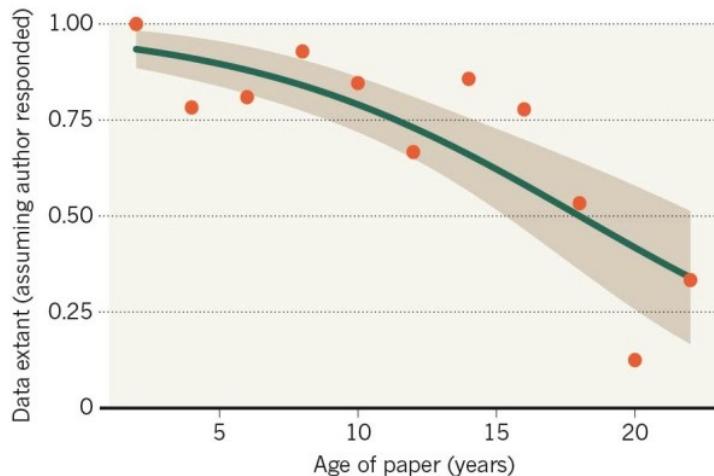
- Democracy and transparency
  - Publicly funded research data should be accessible to all
  - Published results and conclusions should be possible to check by others
- Research
  - Enables others to combine data, address new questions, and develop new analytical methods
  - Reduce duplication and waste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

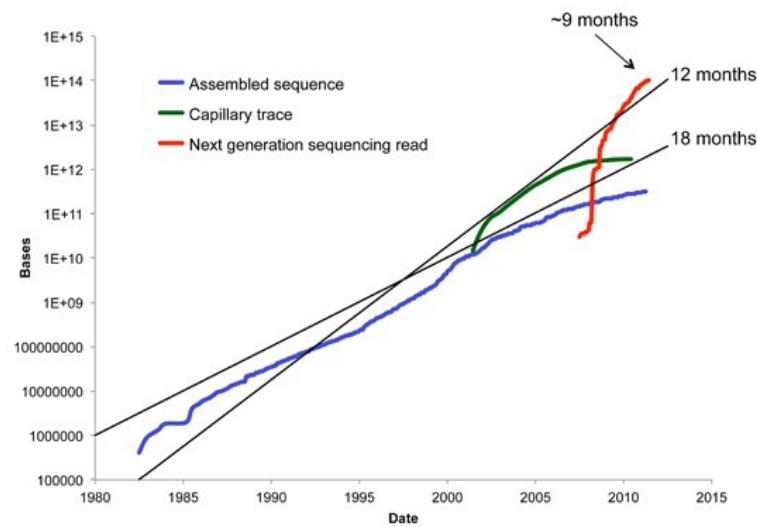
As research articles age, the odds of their raw data being extant drop dramatically.



Nature news, 19 December 2013



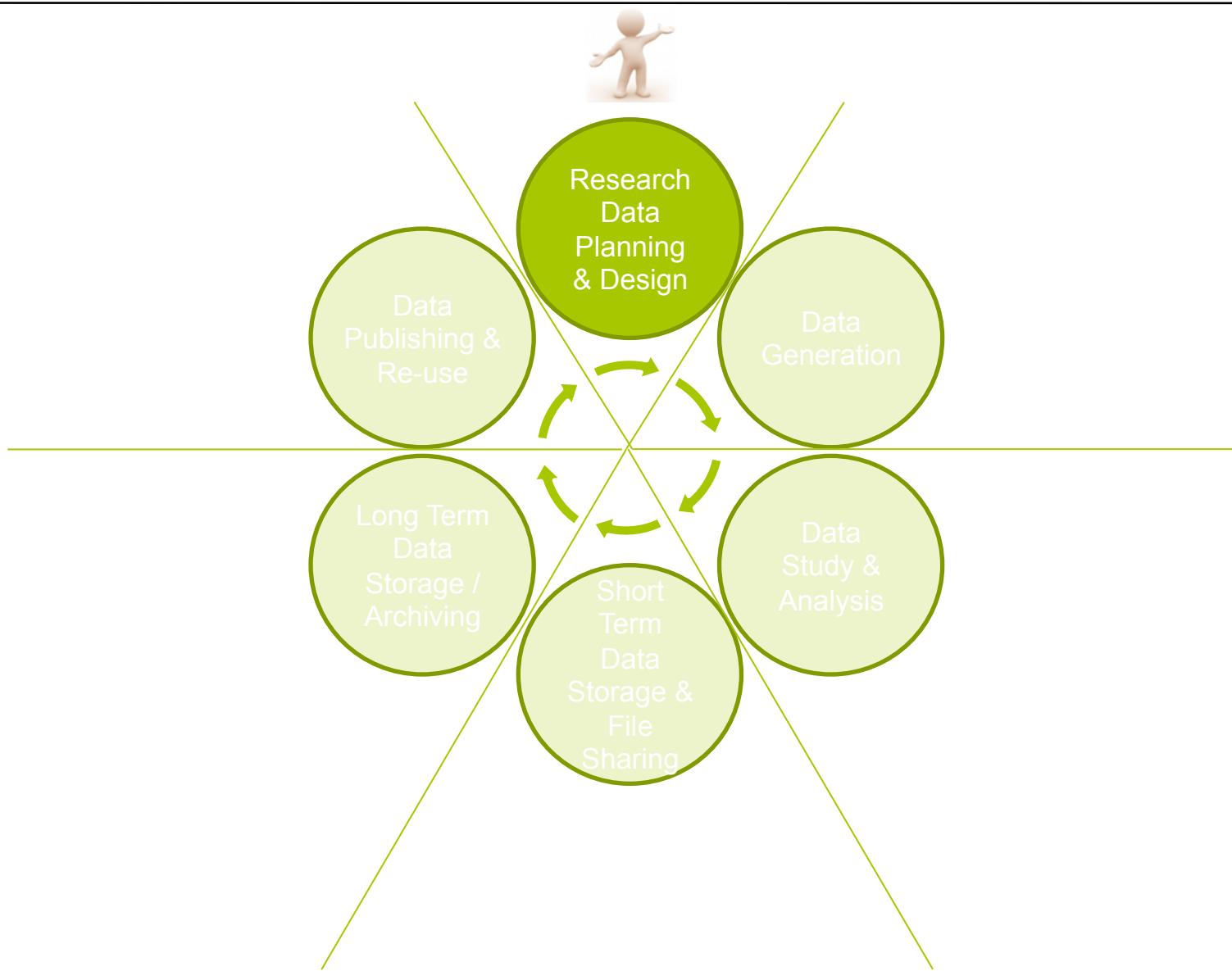
*'Oops, that link was the laptop of my PhD student'*



- 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

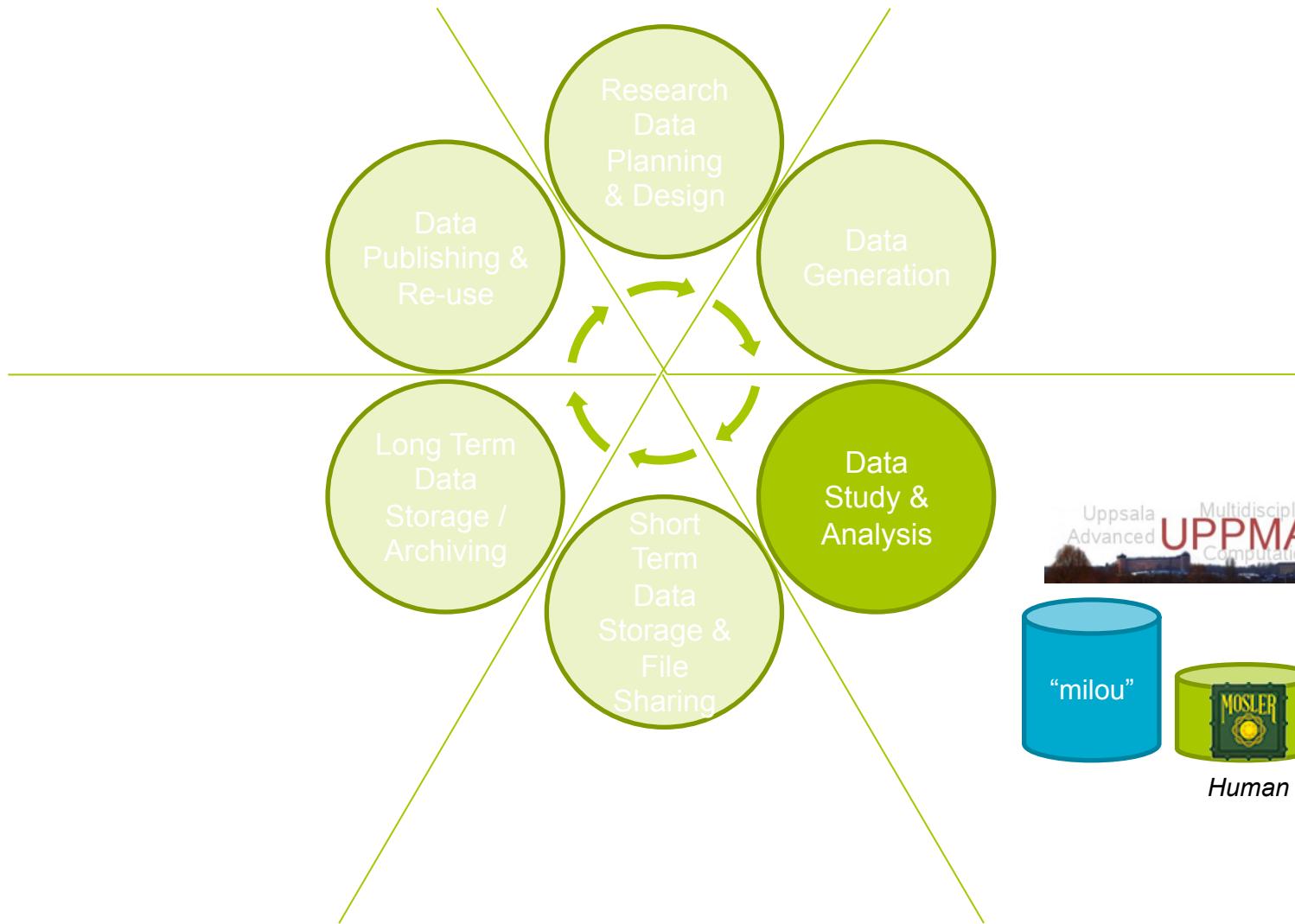
Slide stolen from Barend Mons





- Data Management planning
  - Data types
    - Sizes, were to store, etc
  - **Metadata**
    - Study, Samples, Experiments, etc
    - Use standards!
- *Data Management Plans*
  - Will become a standard part of the research funding application process
  - What will be collected?, Size?, Organized?, Documented?, Stored and preserved?, Disseminated?, Policies?, Budget?





Human derived data

- 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



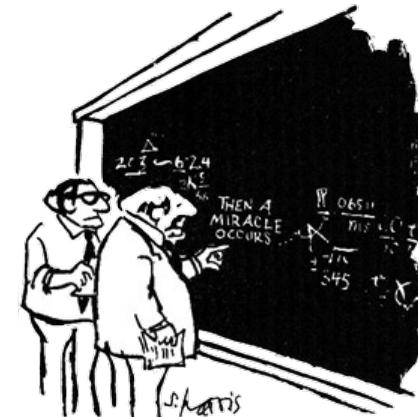
Trevor A. Branch  
@TrevorABranch

 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



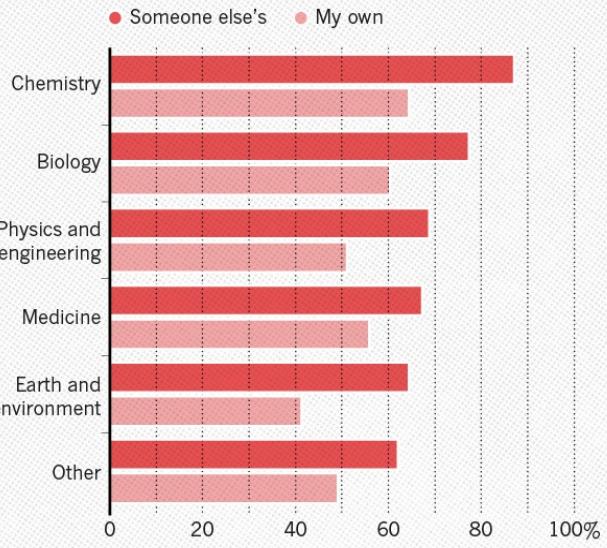
- Structuring data for analysis
  - Poor organizational choices lead to significantly slower research progress.
  - It is critical to make results reproducible.



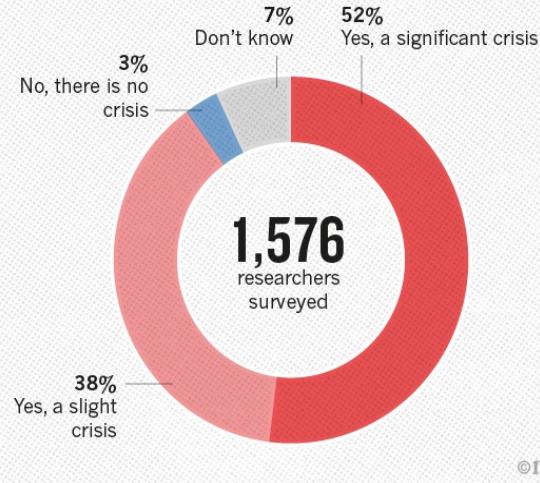
# A reproducibility crisis

## 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 science<sup>1</sup>.

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 irreproducible<sup>2</sup>.

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

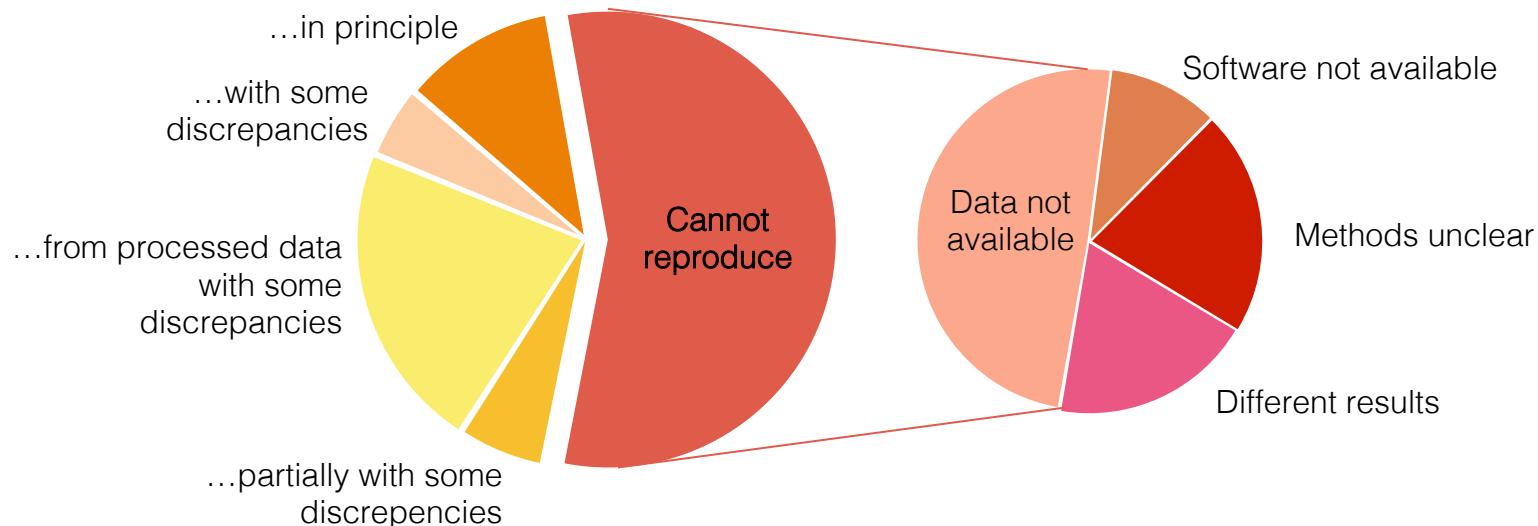
[1] "1,500 scientists lift the lid on reproducibility". Nature. 533: 452–454

[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?

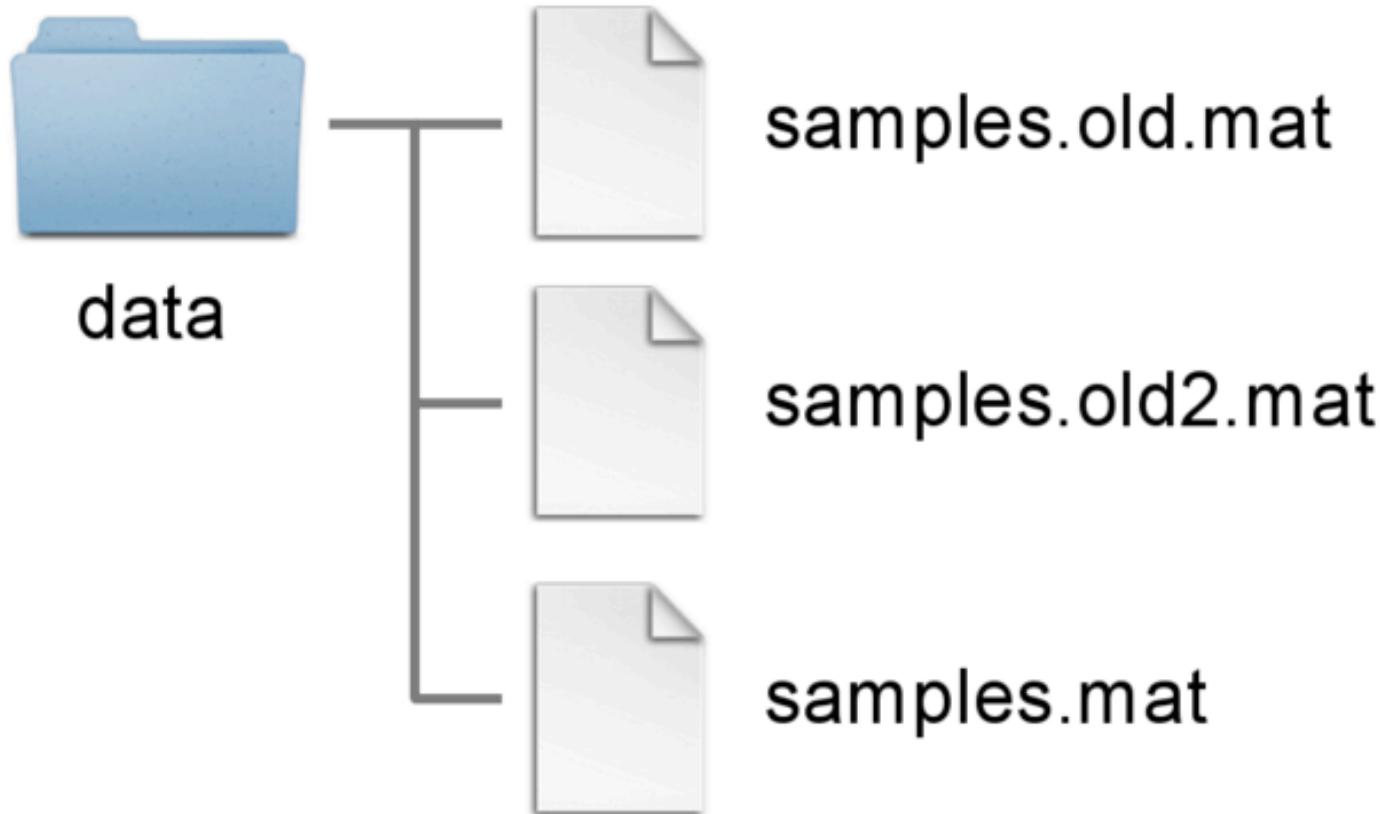
- Primarily for ones own benefit!  
Organized, efficient, in control.  
Dynamic team members.
- Transparent what has been done
- Some will be interested in parts of the analysis. Make it easy to redo, then adapt to own data.

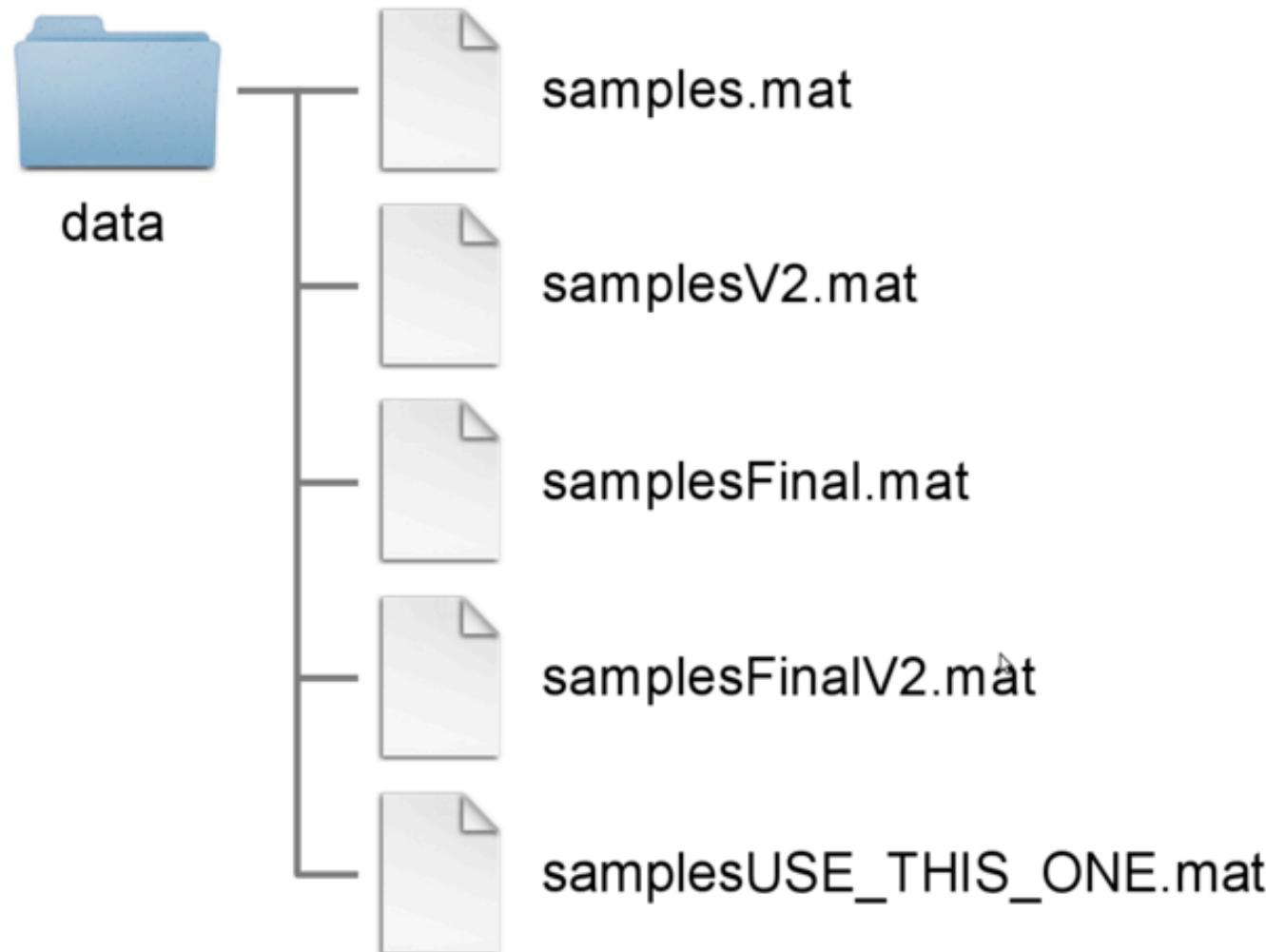


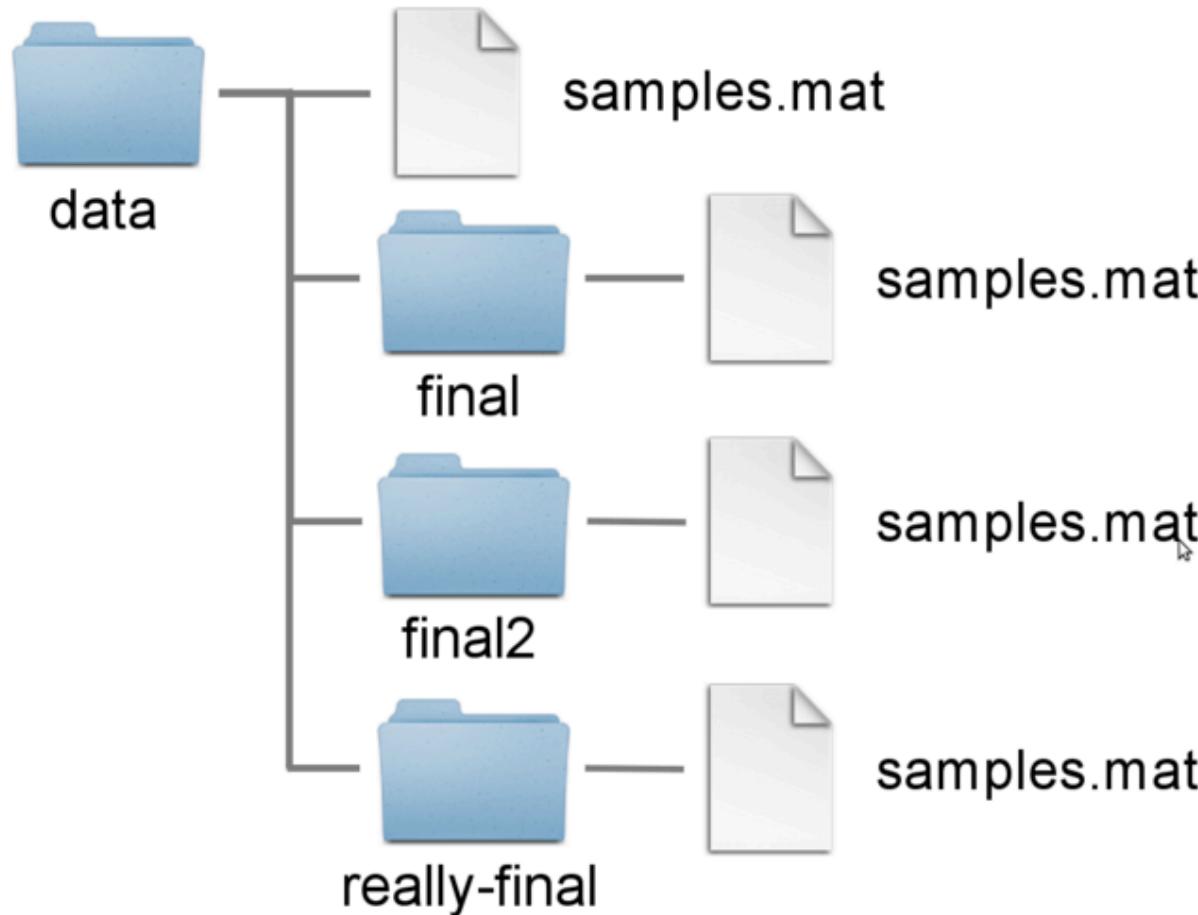
data

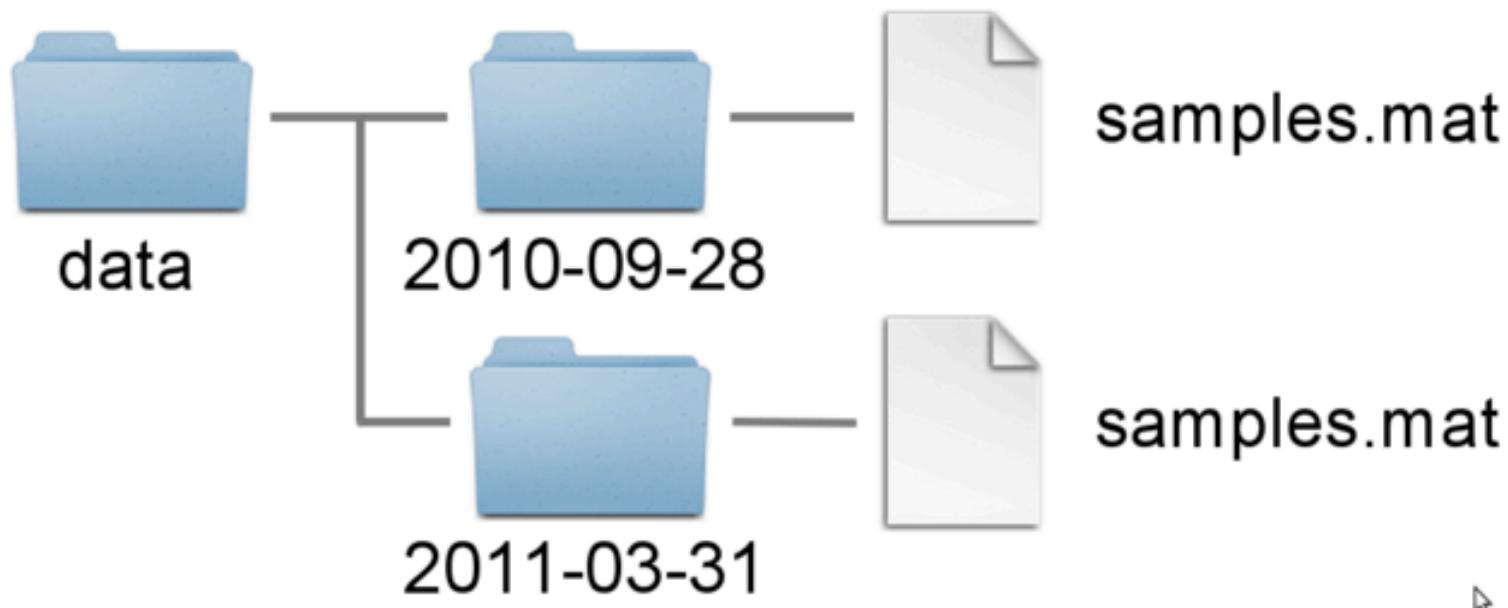
samples.mat





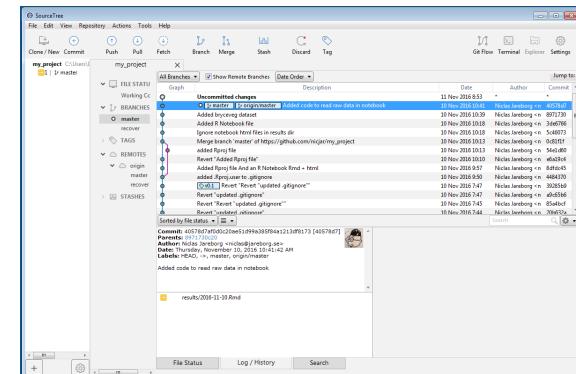






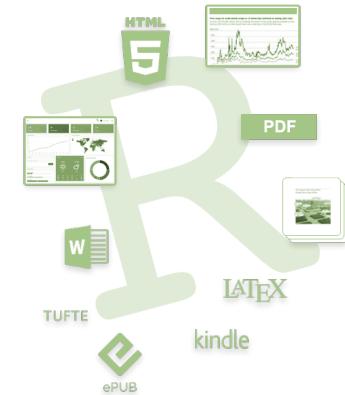
- There is a **folder for the raw data**, which do not get altered, or intermixed with data that is the result of manual or programmatic manipulation. I.e., derived data is kept separate from raw data, and **raw data are not duplicated**.
- **Code is kept separate from data.**
- Use a **version control 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 .x/sx
- 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 .x/sx
- 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, ...



- Need context → document **metadata**
  - How was the data generated?
  - From what was the data generated?
  - What where the experimental conditions?
  - Etc
- Use standards
  - Controlled vocabularies / Ontologies
  - *Not straight-forward...*

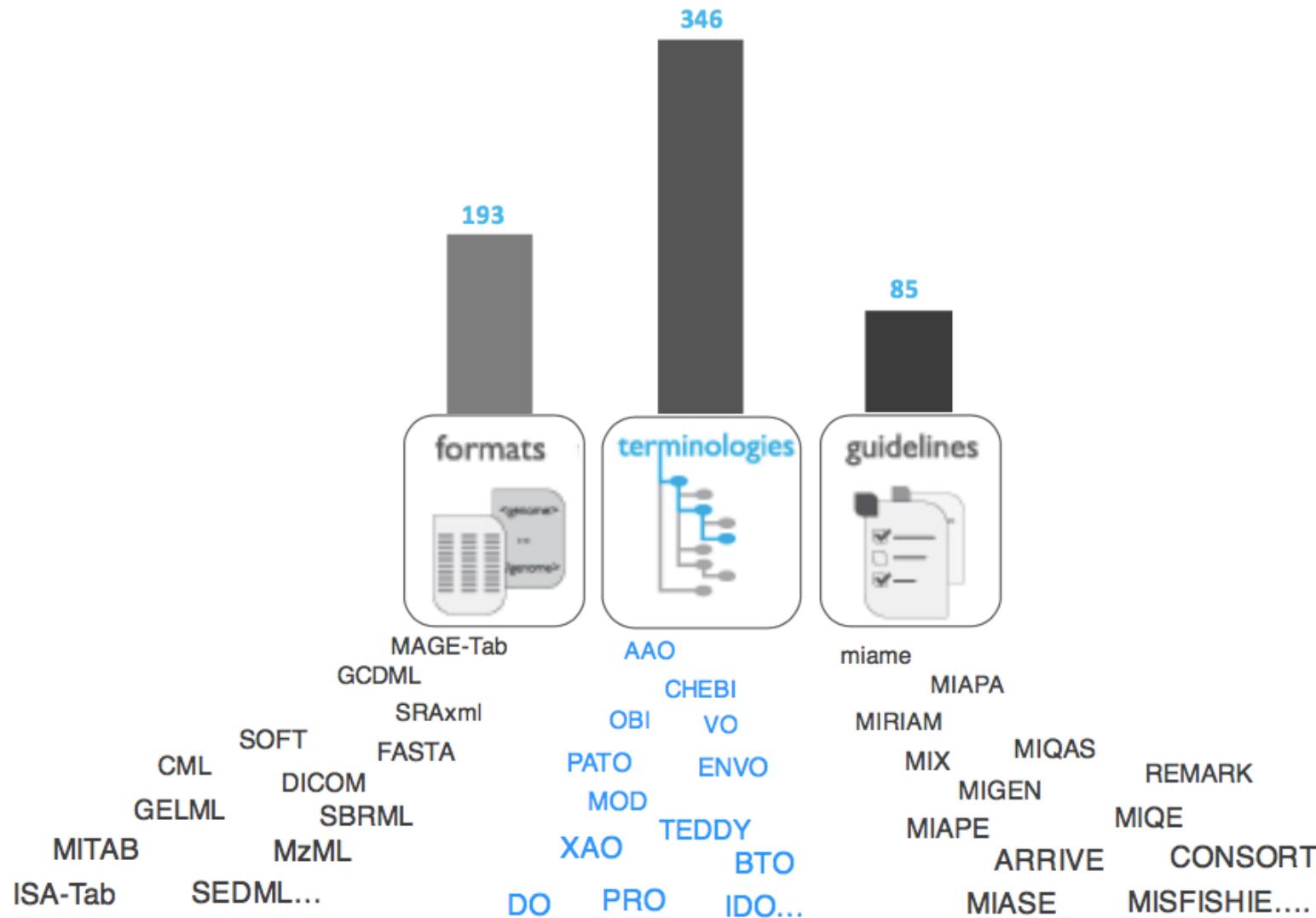
**Human Phenotype Ontology**

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	<a href="http://purl.obolibrary.org/obo/HP_0004808">http://purl.obolibrary.org/obo/HP_0004808</a>		
database_cross_reference	MeSH:D015470 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		

Jump To:

- All
  - Clinical modifier
  - Mode of inheritance
  - Mortality/Aging
  - Phenotypic abnormality
    - Abnormality of blood and blood-forming tissues
      - Abnormal bleeding
      - Abnormal thrombosis
      - Abnormality of bone marrow cell morphology
      - Abnormality of coagulation
      - Abnormality of leukocytes
      - Abnormality of thrombocytes
      - Extramedullary hematopoiesis
      - Hematological neoplasm
    - Leukemia
      - Acute leukemia
        - Acute lymphoblastic leukemia
        - Acute megakaryocytic leukemia
        - Acute monocytic leukemia
        - Acute myeloid leukemia
        - Acute myelomonocytic leukemia
        - Acute promyelocytic leukemia
        - Biphenotypic acute leukaemia
      - Chronic leukemia
      - Lymphoid leukemia
      - Myeloid leukemia
      - Myeloproliferative disorder
    - Lymphoma
      - Lymphoma
      - Lymphoproliferative disorder
      - Malignant eosinophil proliferation
      - Multiple myeloma
      - Myelodysplasia
      - Plasmacytoma
    - Abnormality of connective tissue
    - Abnormality of head or neck
    - Abnormality of limbs
    - Abnormality of metabolism/homeostasis

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

Standards  Databases  Policies  Collections/Recommendations

**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

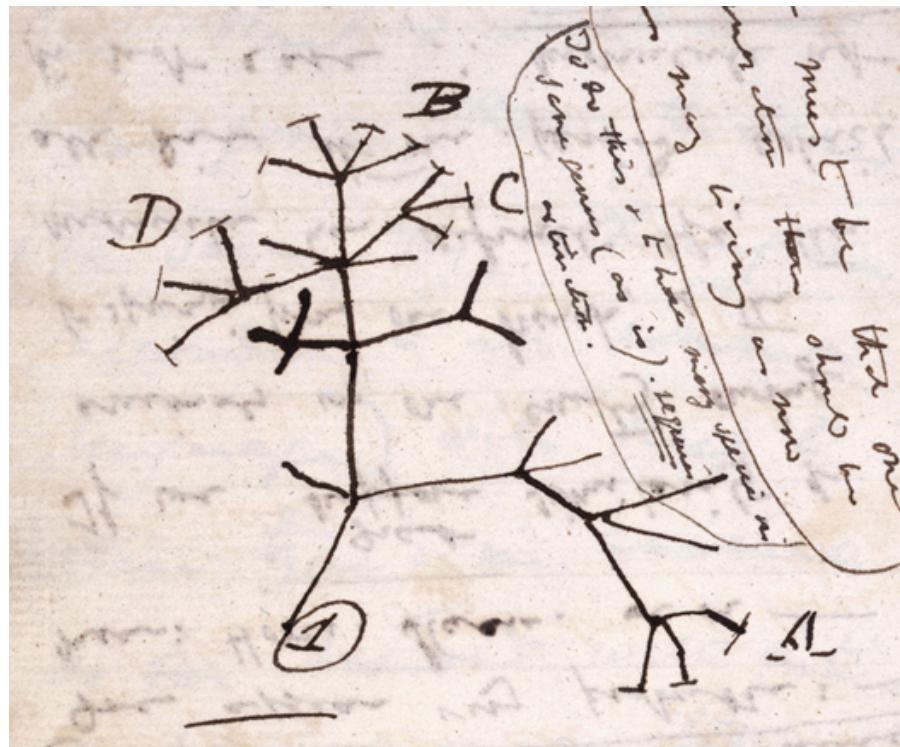
[View all](#)

### 97 Policies

Funder	22
Journal	68
Society	3

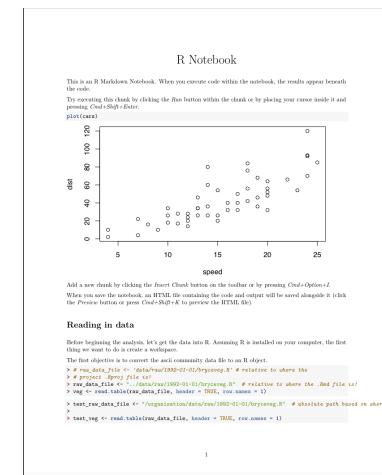
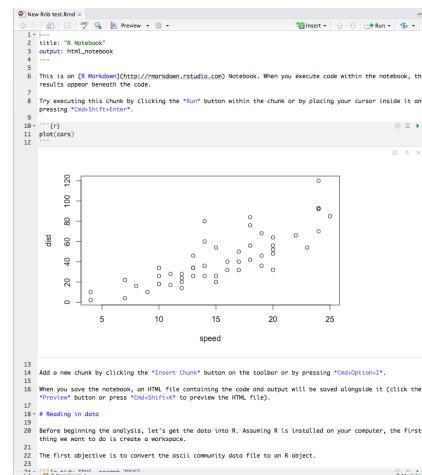
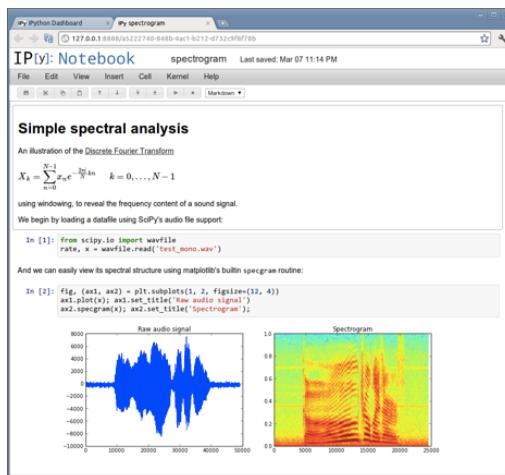
[View all](#)

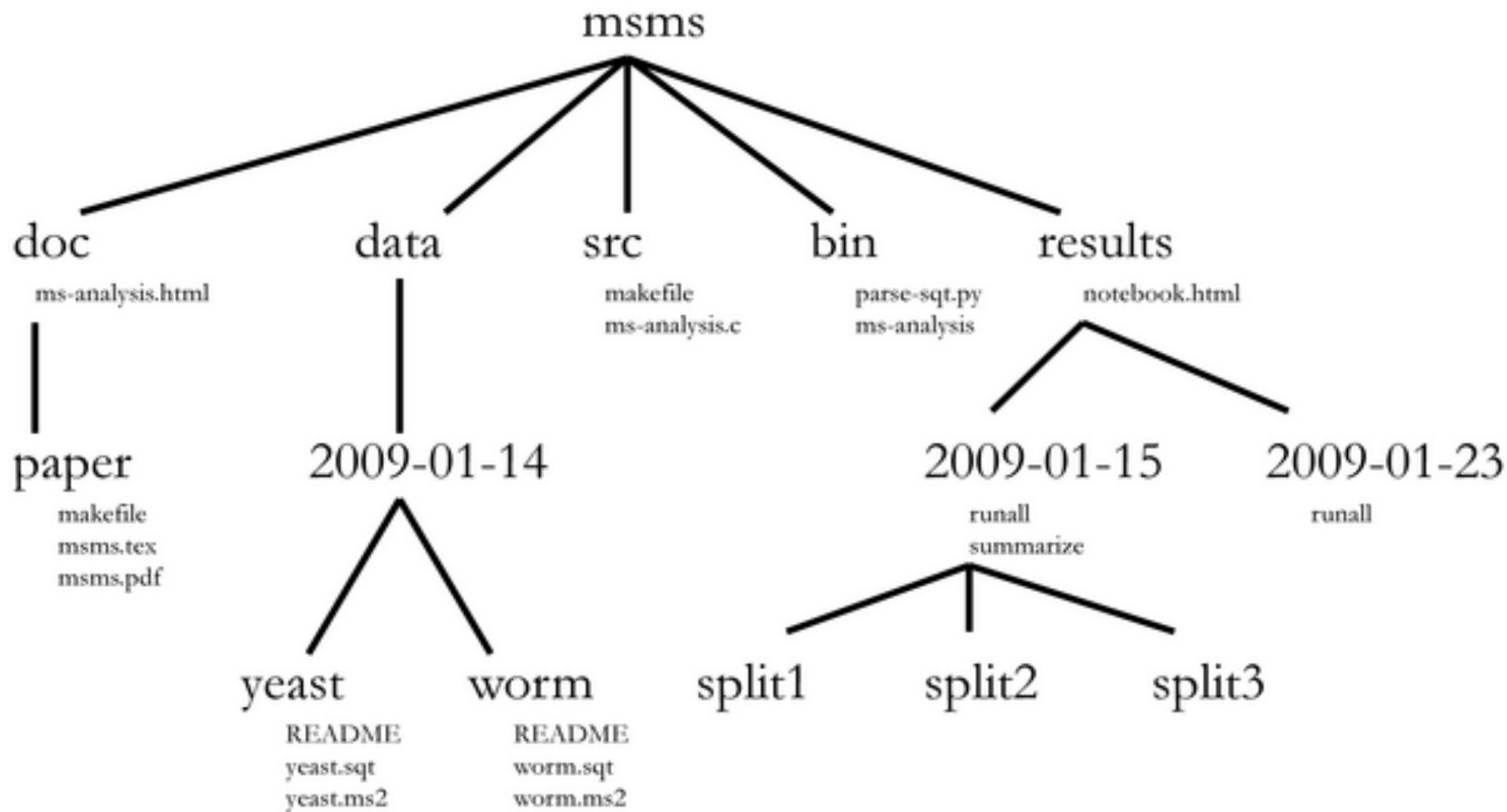
- Why?
  - You have to understand what you have done
  - **Others should be able to reproduce what you have done**



- 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

- 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





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>

```
bin <-----# Binary files and executables (jar files & proj-wide scripts etc)
conf <-----# Project-wide configuraiton
doc <-----# Any documents, such as manuscripts being written
experiments <----# The main experiments folder
    2000-01-01-exa <-# An example Experiment
        audit <----# Audit logs from workflow runs (higher level than normal logs)
        bin <----# Experiment-specific executables and scripts
        conf <----# Experiment-specific config
        data <----# Any data generated by workflows
        doc <----# Experiment-specific documents
        log <----# Log files from workflow runs (lower level than audit logs)
        raw <----# Raw-data to be used in the experiment (not to be changed)
        results <---# Results from workflow runs
        run <----# All files rel. to running experiment: Workflows, run confs/scripts...
        tmp <----# Any temporary files not supposed to be saved
    raw <-----# Project-wide raw data
    results <-----# Project-wide results
    src <-----# Project-wide source code (that needs to be compiled)
```

From Samuel Lampa's blog: <http://bionics.it/posts/organizing-compbio-projects>

- There's no perfect set-up
  - Pick one! e.g.
    - <https://github.com/chendaniely/computational-project-cookie-cutter>
    - <https://github.com/Reproducible-Science-Curriculum/rr-init>
    - <https://github.com/nylander/pTemplate>
    - ...
- Communicate structure to collaborators
- Document as you go
- Done well it might reduce post-project explaining



# Reproducible research for bioinformatics projects

Leif Väremo ([leif.varemo@scilifelab.se](mailto:leif.varemo@scilifelab.se))  
Rasmus Ågren ([rasmus.agren@scilifelab.se](mailto:rasmus.agren@scilifelab.se))  
Bioinformatics long-term support (WABI)

## Everything can be a project

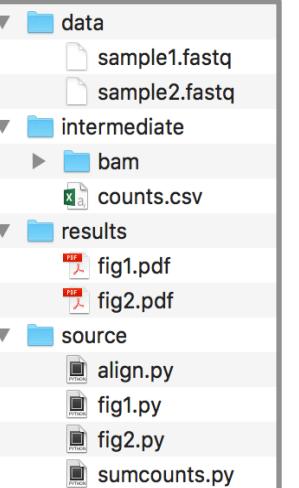
Divide your work into distinct projects and keep all files needed to go from raw data to final results in a dedicated directory with relevant subdirectories (see example).

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.



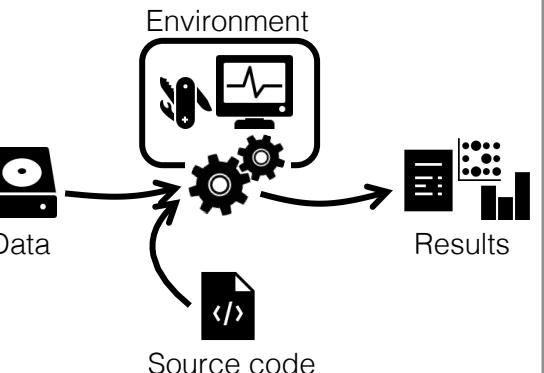
<https://git-scm.com/>



## 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.

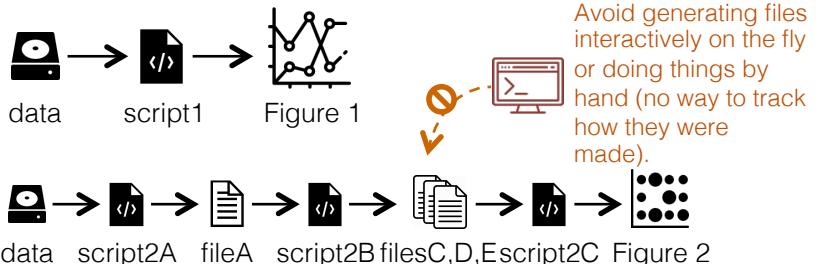


## Treasure your data

- 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)



## 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 environment 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 publically available Docker container!
- Singularity is an alternative to Docker which runs better on HPC clusters.



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



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

- 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 OSF dashboard for a project titled "My fabulous project". The top navigation bar includes links for My Dashboard, Browse, Help, and Settings. The main content area displays basic project metadata: Contributors (Niclas Jareborg), Date created (2016-03-16 03:04 PM), Last Updated (2016-03-16 03:08 PM), Category (Project), Description (No description), and License (No license). Below this, there are four main sections: "Wiki" containing a "Welcome" page with the text "This is a test project to check out functionality"; "Citation" with a link to osf.io/85f7h; "Components" showing contributions for "Data files" (1 contribution by Jareborg) and "Code" (5 contributions by Jareborg); and "Tags" with categories like Data management and Testing.

# Personal data



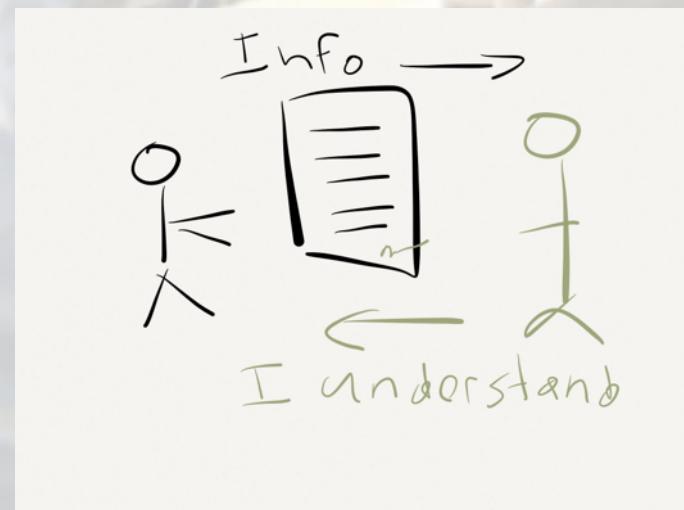
- Personal Data Act (*Personuppgiftslagen (PUL)*)
- Act concerning the Ethical Review of Research Involving Humans (*Lag om etikprövning av forskning som avser människor*)



- All kinds of information that is directly or indirectly referable to a natural person who is alive constitute personal data
- Sensitive data
  - It is **prohibited** to process personal data that discloses *ethnic origin, political opinions, religious or philosophical convictions, membership of trade unions*, as well as personal data relating to **health** or **sexual life**.
  - Sensitive personal data can be handled for **research purposes** if person has given **explicit consent**
- The Data Inspection Board (*Datainspektionen*) is the supervisory authority under the Personal Data Act
- May 2018: **General Data Protection Regulation (GDPR)**
  - New Swedish
    - **Data Protection Act (Dataskyddslag)**
    - **Research Data Act (Forskningsdatalag)**

- The (legal) person that decides why and how personal data should be processed is called the **controller of personal data** (*personuppgiftsansvarig*)
  - e.g. the employing university
- The controller of personal data can delegate processing of personal data to a **personal data assistant** (*personuppgiftsbiträde*)
  - e.g. UPPMAX/Uppsala university
- A **personal data representative** (*personuppgiftsombud*) is a natural person who, on the assignment of the controller, shall ensure that personal data is processed in a lawful and proper manner
- Obligation to report handling of personal data to the Data Inspection Board
  - Or, notify the Board of the named representative

- 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
- *“The controller is liable to implement technical and organizational measures to protect the personal data. The measures shall attain an appropriate level of security.”*

- **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**
  - e-Infrastructure for working with sensitive data for academic research
    - Developed & operated by NBIS
  - Inspired by Norwegian solution (TSD)
  - Designed to look like UPPMAX clusters
    - UPPMAX modules
    - UPPMAX can assist with installing custom tools
  - 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*



- High-performance computing in a virtualized environment (OpenStack)
  - Each project environment is isolated from all other projects
    - Separated private networks and file systems
    - No internet access
    - No root access
- Only accessible over remote Linux desktop (ThinLinc) via a web dashboard
- **2-factor authentication for login**
- **Restricted data transfer in/out**
  - Via a file gateway
  - Project members can transfer IN / only PI allowed to transfer out
  - Not possible to copy/paste out



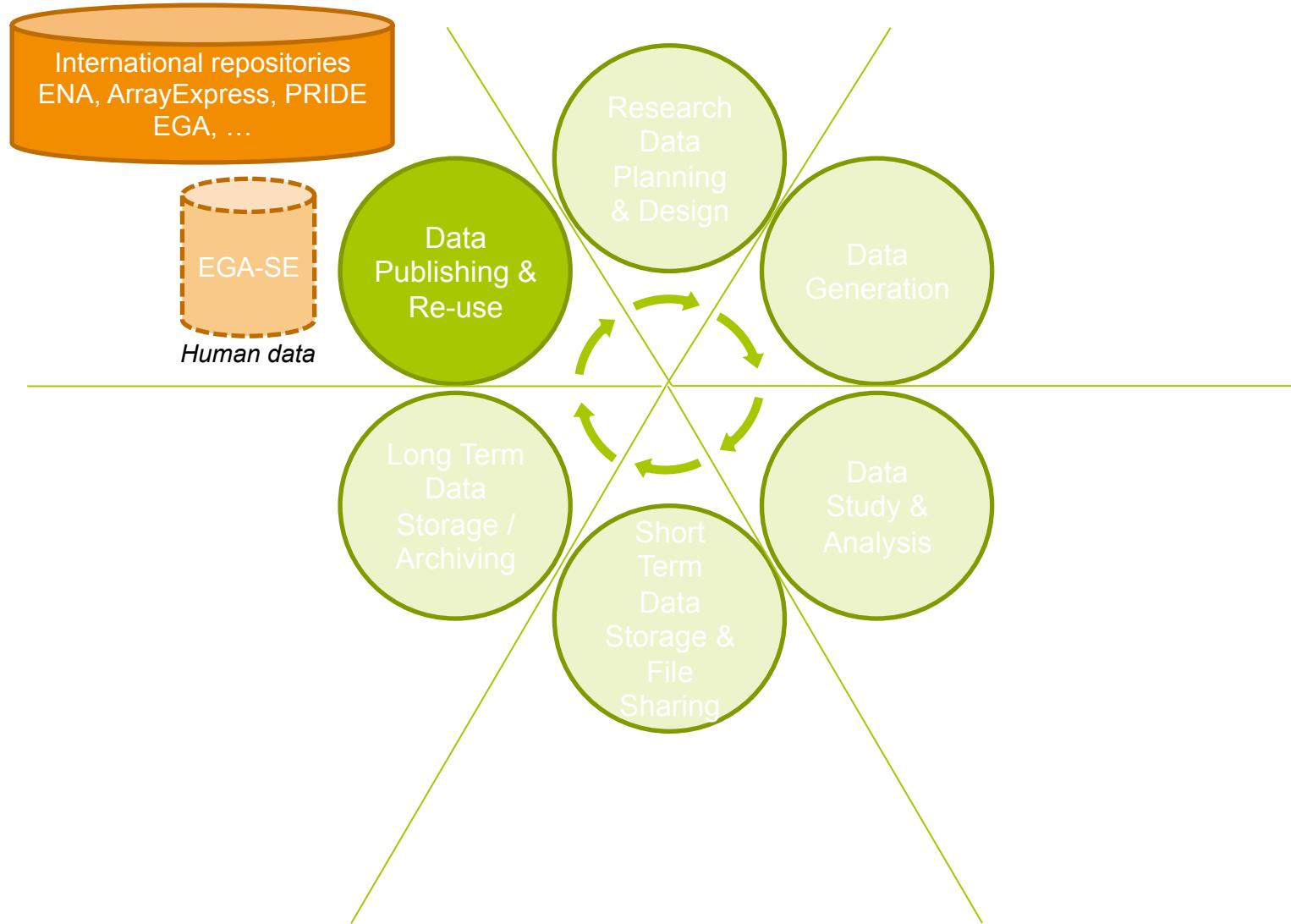
# Tryggve – collaboration for sensitive biomedical data

- Project aims to strengthen Nordic biomedical research by facilitating use of **sensitive data in cross-border projects**
- Collaborators and funders are NeIC and ELIXIR Nodes in Denmark, Finland, Norway and Sweden
- Project will build on strong existing capacities and resources in Nordic countries



1. Technical development
  - Building blocks: Secure systems in Den, Fin, Nor & Swe
2. Interoperability of systems
  - Data transfer service – *sFTP beamer*
  - Portable software installations – *docker containers*
  - Shared computing resources – *Mosler-ePouta*
  - Investigate common authentication and authorization mechanisms
3. Process development
  - Knowledge-sharing (e.g. IT security, administrative processes, harmonizing user agreements)
  - Code of Conduct
4. Legal framework
  - Assessing relevant legislation
  - Analyzing legal requirements in use cases
5. **Use cases**
  - **Implement and support concrete use cases to facilitate cross-border research, and to connect project to actual user demands.**
6. Communication and outreach

[https://wiki.neic.no/wiki/Tryggve\\_Getting\\_Started](https://wiki.neic.no/wiki/Tryggve_Getting_Started)



- *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



- 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 electronic us for help. *Floppy disks.* We can read M or 5½-in diskettes written on MS-DOS so that the submitted data be written as flat t in a format specific to a given word |

Growth of the GenBank Database  
 October 1982 to August 1987

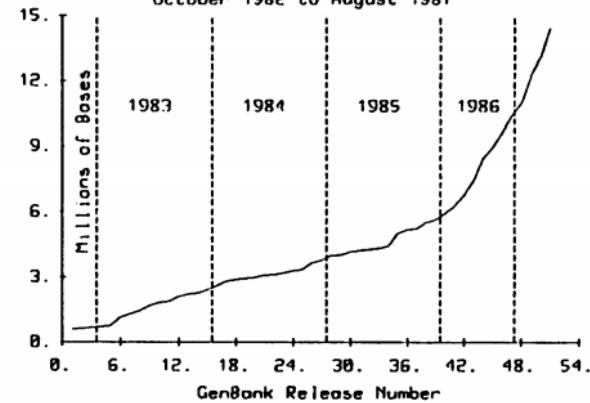


Figure 1.

"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

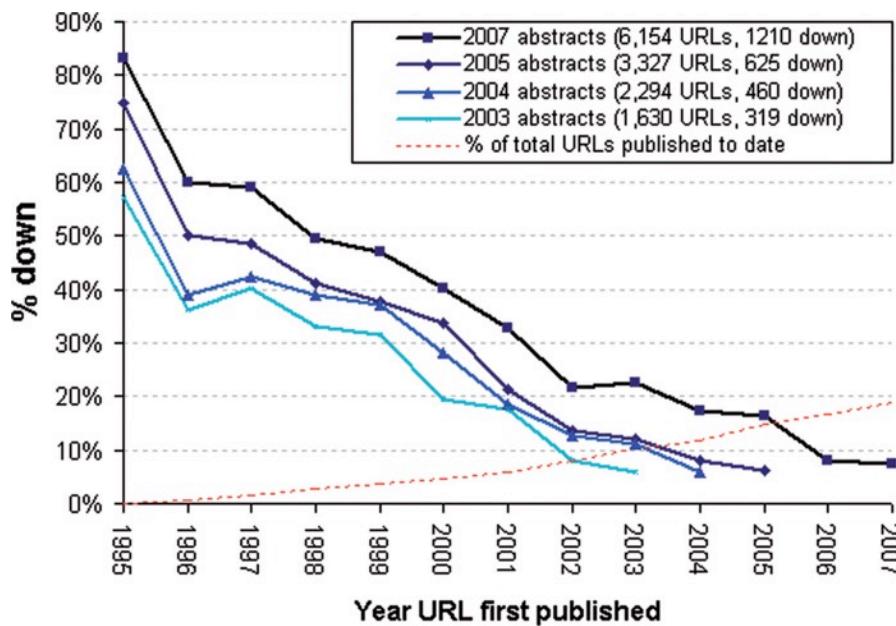


## URL decay in MEDLINE—a 4-year follow-up study

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>

[www.nature.com/scientificdata/](http://www.nature.com/scientificdata/)

**SCIENTIFIC DATA**

**OPEN**

SUBJECT CATEGORIES  
 » Research data  
 » Publication characteristics

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

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 SEPT



- 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

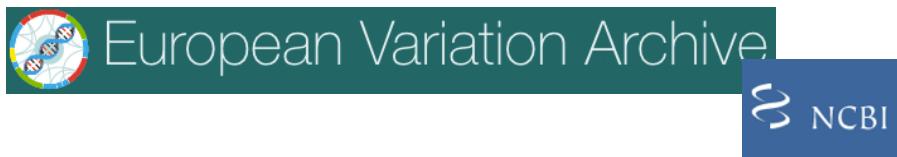
### EOSC Declaration

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 commons, 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 will actively support their implementation in the respective capacities:

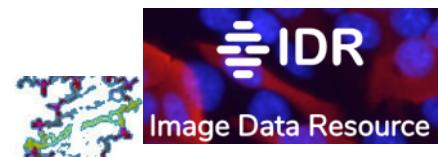
#### Data culture and FAIR data

- [Data culture] European science must be grounded in a common culture of data stewardship, so 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 will enable long-term reuse for science and for innovation of data created by research activities: no 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 they need substantial support from the European Commission and the Member States.
- [Data stewardship] Researchers need the support of adequately trained data stewards. The European Commission and Member States should invest in the education of data stewards via 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, both





dbSNP  
Short Genetic Variations



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

Deposition Database	Data type	International collaboration framework <sup>1</sup>	Deposition Database	Data type	International collaboration framework <sup>1</sup>
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 <a href="#">How countries join</a> ):	
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

---

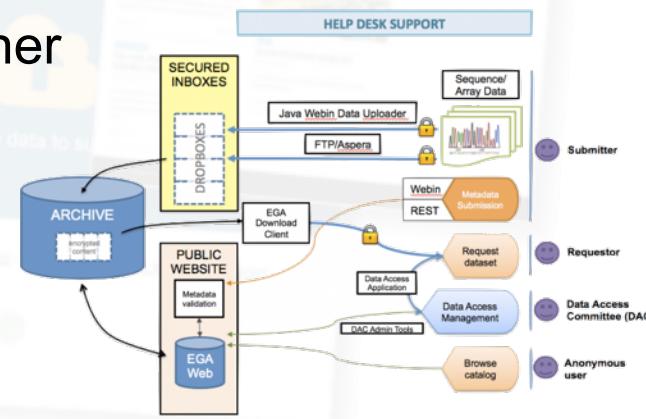
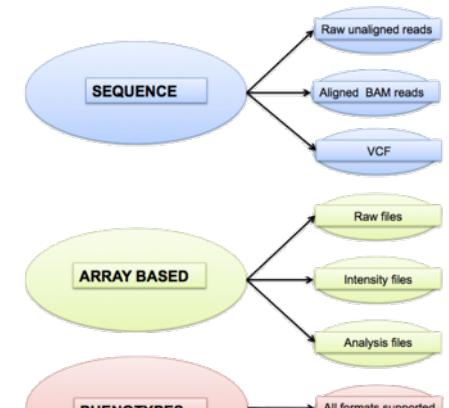
- NIH funded research
  - Only 12% of articles from NIH-funded research mention data deposited in international repositories
  - Estimated 200000+ “invisible” data sets / year

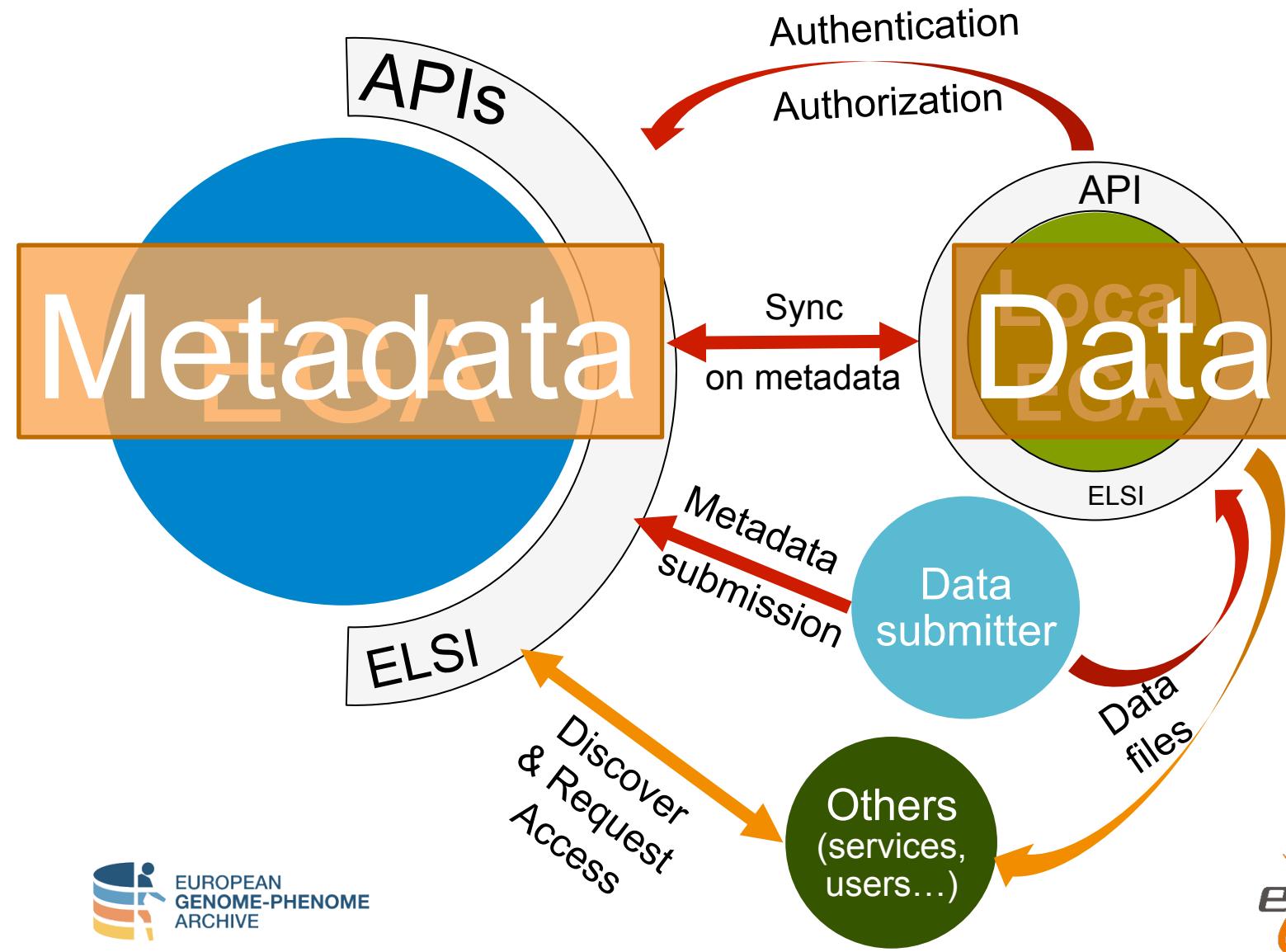
*Read et al. “Sizing the Problem of Improving Discovery and Access to NIH-Funded Data: A Preliminary Study” (2015)*

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

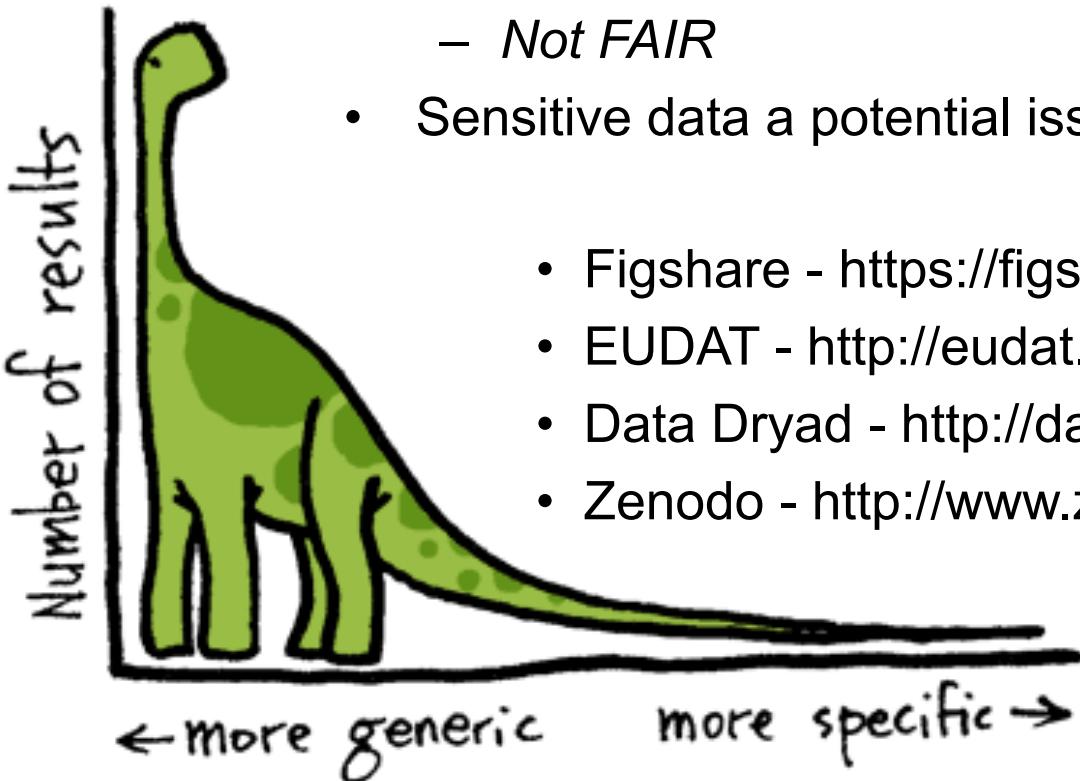
- **EGA – European Genome-phenome Archive**
  - Repository that promotes the distribution and sharing of **genetic and phenotypic data** consented for specific approved uses but **not fully open, public distribution.**
  - All types of sequence and genotype experiments, including case-control, population,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

European  
genome-phenome  
archive



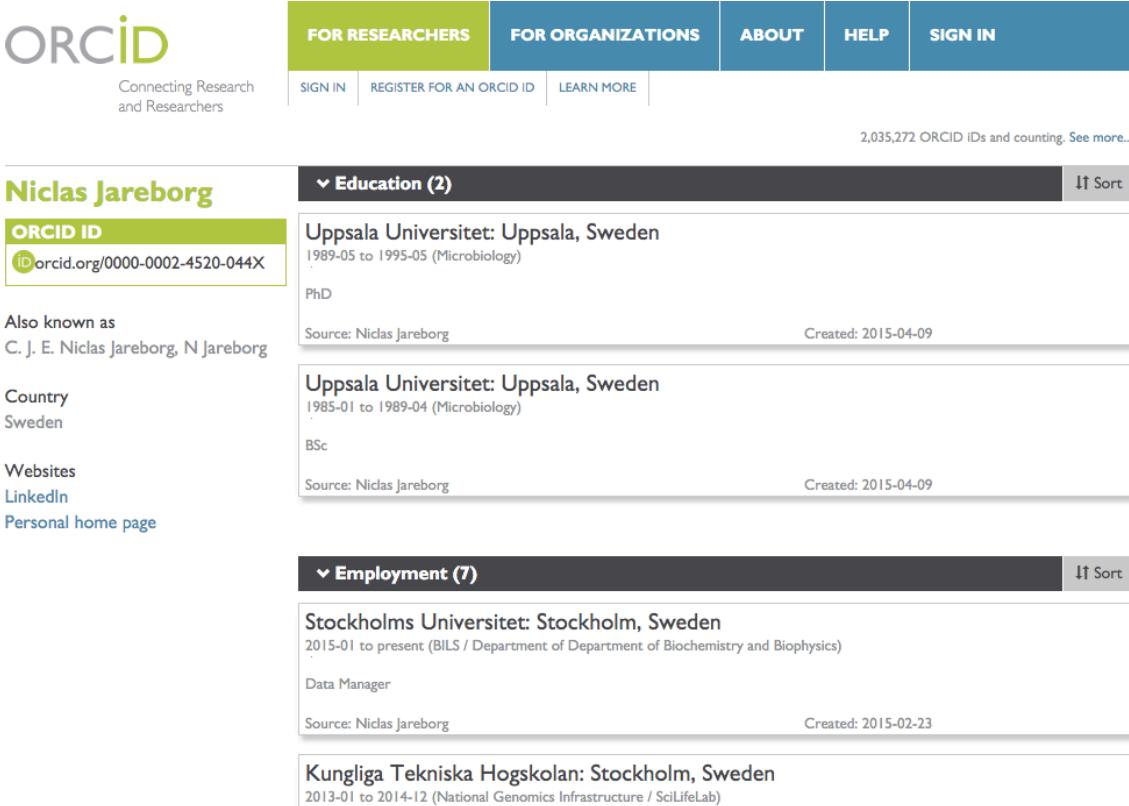


- Research data that doesn't fit in structured data repositories
- Data publication – persistent identifiers
- Metadata submission – not tailored to Life Science
  - *Affects discoverability*
  - *Not FAIR*
- Sensitive data a potential issue



- Figshare - <https://figshare.com/>
- EUDAT - <http://eudat.eu/>
- Data Dryad - <http://datadryad.org/>
- Zenodo - <http://www.zenodo.org/>

- 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 tabs: FOR RESEARCHERS (highlighted in green), FOR ORGANIZATIONS, ABOUT, HELP, and SIGN IN. Below the navigation bar, there are links for SIGN IN, REGISTER FOR AN ORCID ID, and LEARN MORE. A statistic at the top right says "2,035,272 ORCID IDs and counting. See more..."

The main content area displays Niclas Jareborg's profile information. It includes his ORCID ID (ID.orcid.org/0000-0002-4520-044X) and sections for Education and Employment.

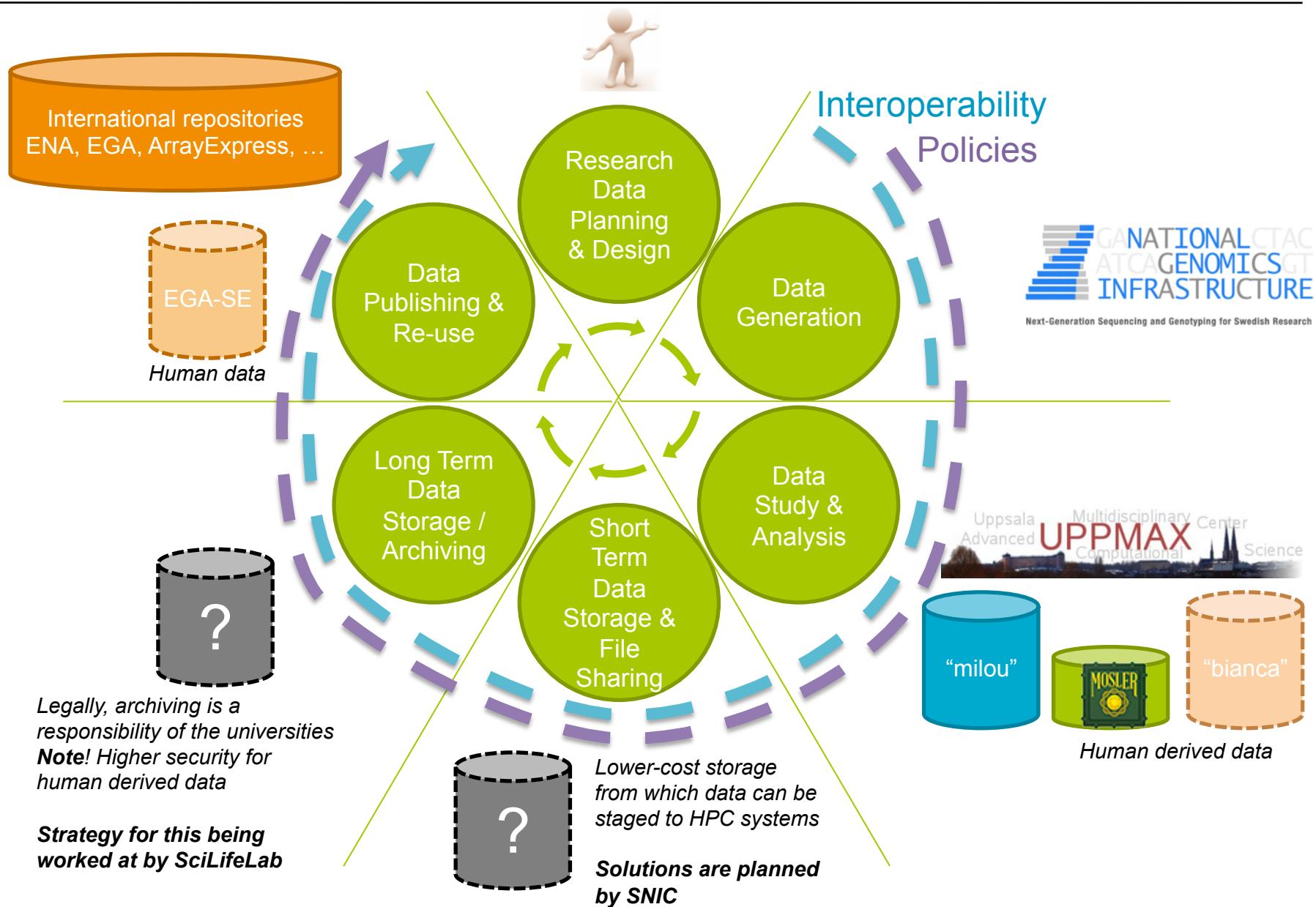
**Education:**

- 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.

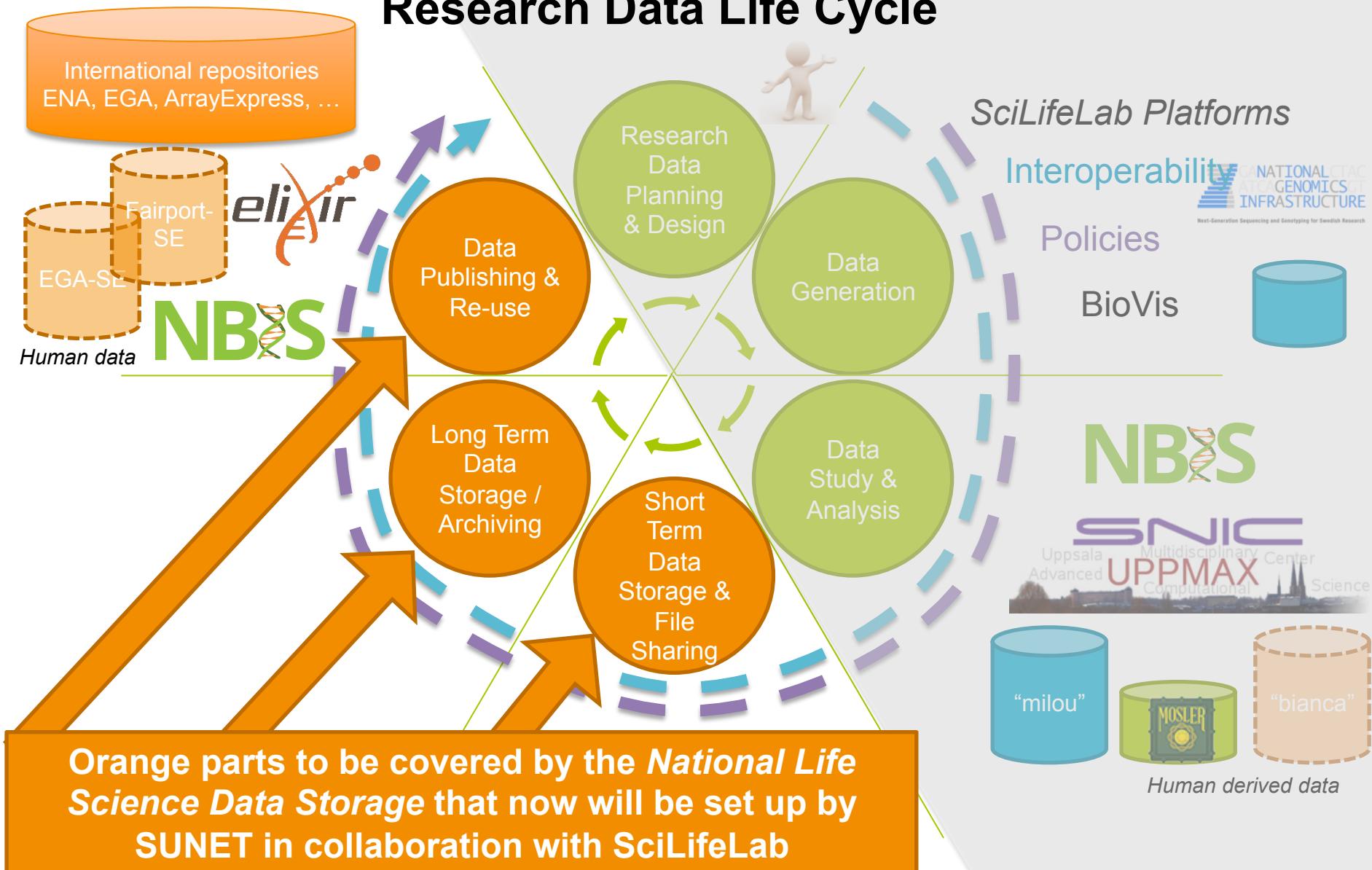
**Employment:**

- Stockholms Universitet: Stockholm, Sweden (2015-01 to present, BILS / Department of Biochemistry and Biophysics) - Data Manager. Source: Niclas Jareborg. Created: 2015-02-23.
- Kungliga Tekniska Högskolan: Stockholm, Sweden (2013-01 to 2014-12, National Genomics Infrastructure / SciLifeLab)

- Project planning
  - Metadata
  - File formats
  - Licensing
  - *Data Management Plans*
- Data analysis
- Data publication and submission
  - Automate submissions to public repositories
  - Metadata
  - Licensing

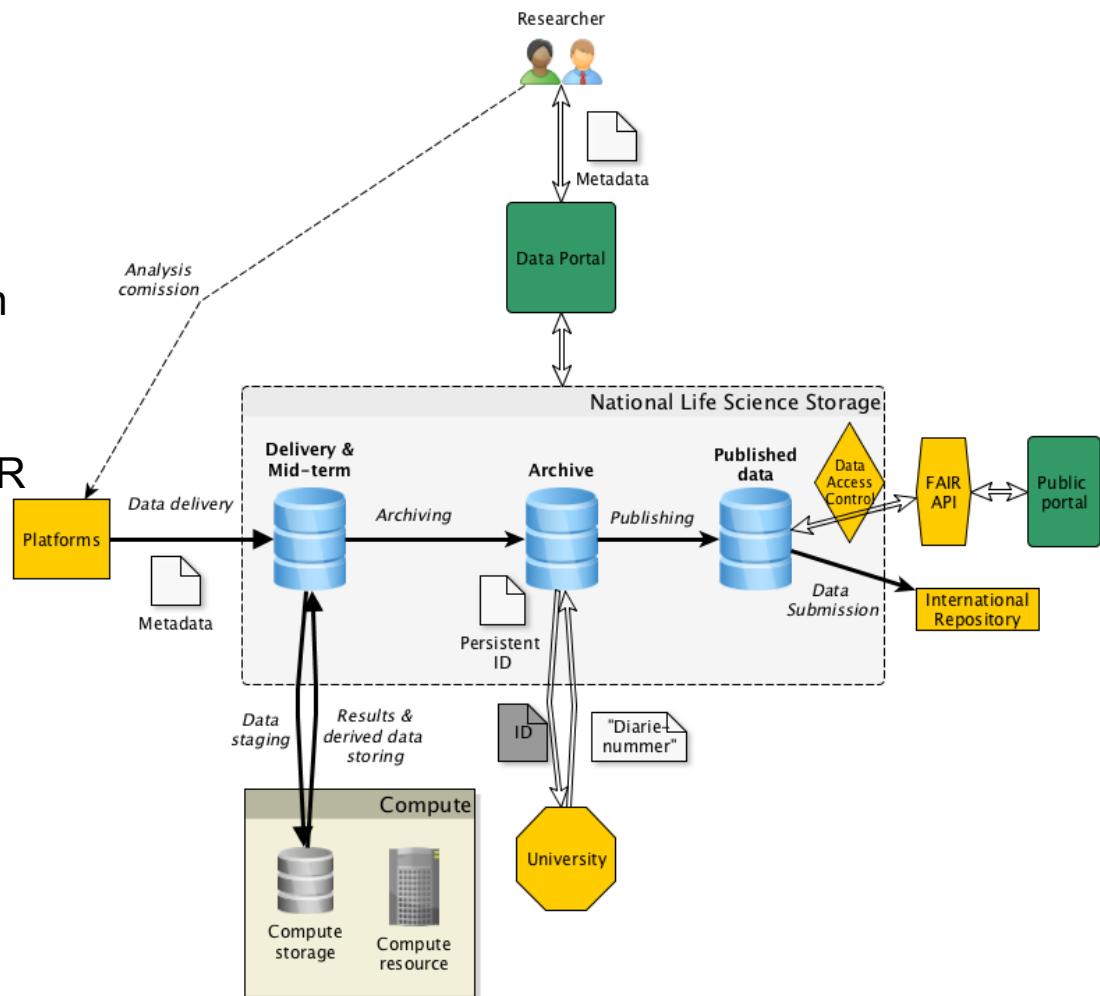


# Research Data Life Cycle



## Components

- “Active Data” storage
- Data staging to HPC resources
- Archiving
  - Offer a solution to the universities’ legal obligation (*possible funding stream*)
- Data publication
  - Making SciLifeLab data FAIR
- User-friendly interface to manage the data life cycle process
- Support the SciLifeLab Data Office way of working



- Research Data Management, EUDAT -  
<http://hdl.handle.net/11304/79db27e2-c12a-11e5-9bb4-2b0aad496318>
- Barend Mons – FAIR Data
- Antti Pursula – Tryggve <https://wiki.neic.no/wiki/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)
- Samuel Lampa - <http://bionics.it/posts/organizing-compbio-projects>
- Reproducible Science Curriculum –  
<https://github.com/Reproducible-Science-Curriculum/rr-init>
- Leif Väremo -  
[https://bitbucket.org/scilifelab-lts/reproducible\\_research\\_example/src](https://bitbucket.org/scilifelab-lts/reproducible_research_example/src)