

Data FAIRification using RStudio workflows

Brendan Palmer & Darren Dahly,
Clinical Research Facility - Cork & School of Public Health,
University College Cork
 [@B_A_Palmer](https://twitter.com/B_A_Palmer) [@statsepi](https://twitter.com/statsepi)

Where to begin...



Don't do what Donny Dont does!



"In short, peer review misses all the hard stuff, and a worrying amount of the easy stuff"

James Heathers,
Northwestern University

#datathugs



Brian Wansink: The grad student who never said no

"Every day we would scratch our heads, ask "Why," and come up with another way to reanalyze the data with yet another set of plausible hypotheses. Eventually we started discovering solutions"

Credibility crisis

2005

PLOS MEDICINE

BROWSE PUBLISH ABOUT SEARCH advanced search

OPEN ACCESS ESSAY

Why Most Published Research Findings Are False

John P. A. Ioannidis
Published: August 30, 2005 • <https://doi.org/10.1371/journal.pmed.0020124>

68,436 Save	3,184 Citation
2,813,238 View	10,483 Share

2016

nature International weekly journal of science

Search Go Advanced search

Home News & Comment Research Careers & Jobs Current Issue Archive Audio & Video For Authors

Archive Volume 533 Issue 7604 News Feature Article

NATURE | NEWS FEATURE

1,500 scientists lift the lid on reproducibility

Survey sheds light on the 'crisis' rocking research.

Monya Baker

2018

THE IRREPRODUCIBILITY CRISIS OF MODERN SCIENCE

Causes, Consequences, and the Road to Reform



DAVID RANDALL AND CHRISTOPHER WELSER
NATIONAL ASSOCIATION OF SCHOLARS
APRIL 2018
ISBN: 978-0-9986635-5-5



REFLECTIONS

ON THE

DECLINE OF SCIENCE IN ENGLAND,

AND ON

SOME OF ITS CAUSES.

BY

CHARLES BABBAGE, ESQ.

LUCASIAN PROFESSOR OF MATHEMATICS IN THE UNIVERSITY OF CAMBRIDGE,
AND MEMBER OF SEVERAL ACADEMIES.

LONDON:

PRINTED FOR B. FELLOWES, LUDGATE STREET;
AND J. BOOTH, DUKE STREET, PORTLAND PLACE.

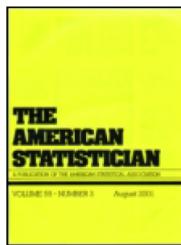
1830

The winds of change

CONSORT 2010

The CONSORT (CONsolidated Standards of Reporting Trials) 2010 guideline is intended to improve the reporting of parallel-group randomized controlled trial (RCT), enabling readers to understand a trial's design, conduct, analysis and interpretation, and to assess the validity of its results. This can only be achieved through complete adherence and transparency by authors.

CONSORT 2010 was developed through collaboration and consensus between clinical trial methodologists, guideline developers, knowledge translation specialists, and journal editors (see [CONSORT group](#)). CONSORT 2010 is the current version of the guideline and supersedes the 2001 and 1996 versions. It contains a 25-item [checklist](#) and [flow diagram](#), freely available for viewing and [downloading](#) through this website.



The American Statistician



ISSN: 0003-1305 (Print) 1537-2731 (Online) Journal homepage: <http://amstat.tandfonline.com/loi/utas20>

The ASA's Statement on *p*-Values: Context, Process, and Purpose

Ronald L. Wasserstein & Nicole A. Lazar



The American Statistician

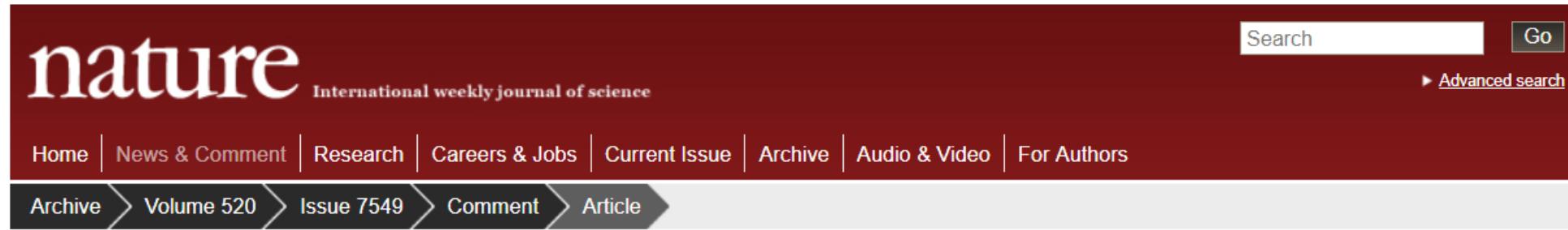


ISSN: 0003-1305 (Print) 1537-2731 (Online) Journal homepage: <https://www.tandfonline.com/loi/utas20>

Moving to a World Beyond "*p* < 0.05"

Ronald L. Wasserstein, Allen L. Schirm & Nicole A. Lazar

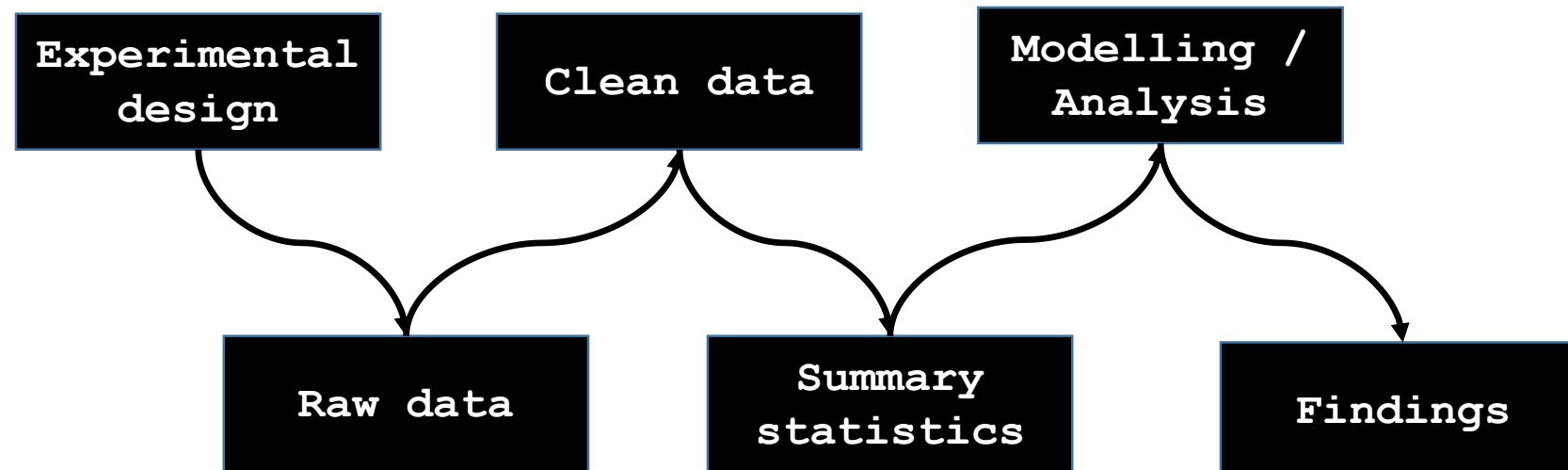
p-values should not define a study



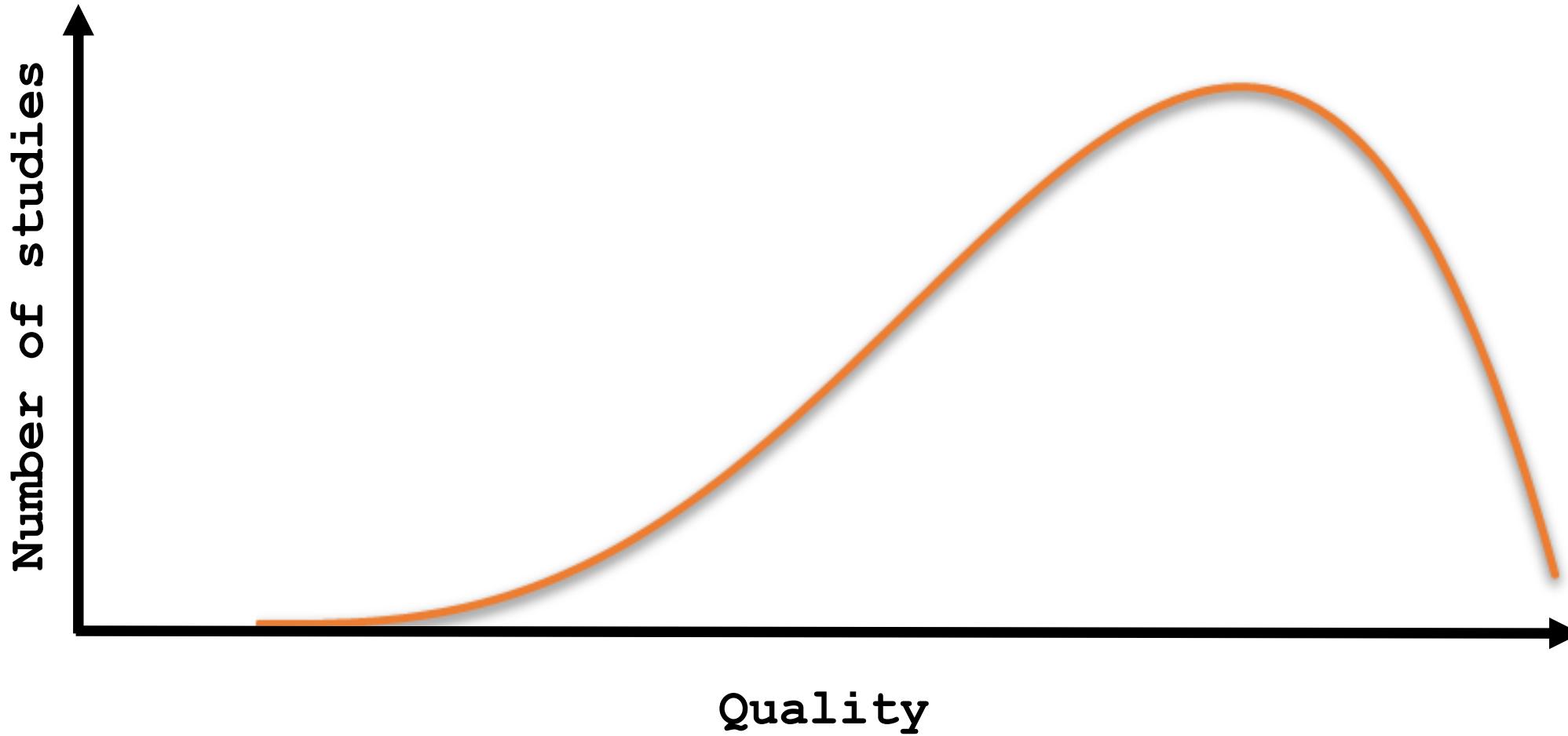
Statistics: *P* values are just the tip of the iceberg

Jeffrey T. Leek & Roger D. Peng

28 April 2015



Today



The butterfly has started flapping its wings



Why Plan S [10 Principles](#) Funders & support Implementation About Contact

"After 1 January 2020 scientific publications on the results from research funded by public grants provided by national and European research councils and funding bodies, must be published in compliant Open Access Journals or on compliant Open Access Platforms."



EUROPEAN COMMISSION
Directorate-General for Research & Innovation

H2020 Programme

Guidelines on
FAIR Data Management in Horizon 2020

Home > Funding > Policies and principles > **Open Research**

Open Research

The HRB is committed to ensuring that its funded research is open, accessible and usable, so it can have the greatest possible impact.

There is a fundamental shift across Europe towards making research more transparent, collaborative, accessible and efficient. The Open Science movement is a strategic priority for the European Commission in research and innovation policy and an EU high-level Expert Group, the [Open Science Policy Platform](#) (OSPP 2016–2018) has been established to consider key implementation areas.

Funding schemes
EU funding support
Manage a grant
Funding awarded
Evaluation
GDPR guidance for researchers
Policies and principles
EU legislation
Gender
Good research practice
Open Research

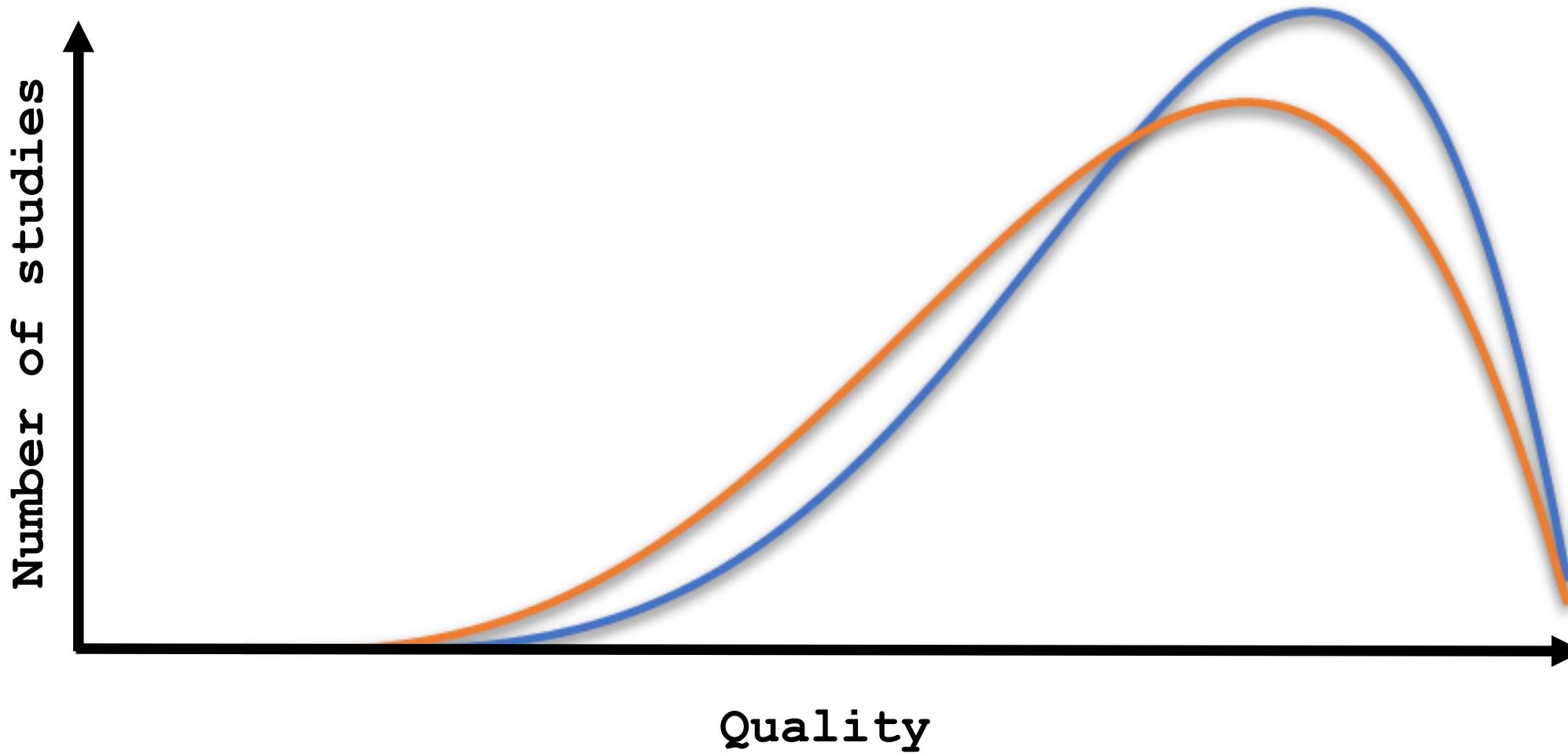


Funding Engagement Events Research News SFI Research Centres

→ Science Foundation Ireland joins DORA

14th February 2019, Dublin – Science Foundation Ireland has become a signatory to the San Francisco Declaration of Research Assessment (DORA), making a formal commitment to assessing the quality and impact of research through means other than journal impact factors.

Tomorrow



FAIR is a part of your life now!

Data Guidelines

1. Background
 - 1.1 Open Data Policy
 - 1.2 Fair Data Principles

2. Share Your Data in 3 Steps
 - 2.1 Prepare Your Data for Sharing
 - 2.2 Select a Repository
 - 2.3 Add a Data Availability Statement to Your Manuscript
 - 2.4 Linking your datasets to your article

Some types of data benefit from visualization within the article. Wellcome Open Research welcomes the submission of manuscripts featuring [Plot.ly interactive figures](#) and [Code Ocean compute capsules](#). For further detail, please [contact us](#).



Research Data Management

Good data governance and stewardship are key components of good research practice. In this regard, Science Foundation Ireland supports that research data should be Findable, Accessible, Interoperable and Reusable (FAIR)*. Appropriate data management and data sharing are fundamental to all stages of the research process and support high quality, reproducible research. As such, access to research data arising in whole or in part from SFI funding should be as open as possible.



FAIR Data Management

Describe the approach to data management that will be taken during and after the project, including who will be responsible for data management and data stewardship. The word limit is 500 words.



Social Research Ethics Committee (SREC) ETHICS APPROVAL FORM

✉ srec@ucc.ie

<https://www.ucc.ie/en/research/about/ethics/>

⁴ Data management should follow the FAIR guiding principles (Findability, Accessibility, Interoperability & Reusability). See, for example, Wilkinson, M. D. et al. (2016) The FAIR Guiding Principles for Scientific Data Management and Stewardship. Full text: <http://www.nature.com/articles/sdata201618>. It is required that all staff and student researchers store those data which are required to replicate research findings, and the information required to enable re-use of data. Details of the UCC policy on research data storage can be found in section 8 of the Code of Research Conduct (2016): <https://www.ucc.ie/en/media/research/researchatucc/documents/UCCCodeofResearchConduct.pdf>. SREC advises against storing research data on non UCC approved cloud-based storage services. Physical data must be stored in a locked cabinet and you must specify who has permission to access this data.



A set of Digital Object Compliance principles that describes the properties of digital objects that enables them to be findable, accessible, interoperable and reproducible (FAIR).

What are the FAIR data principles

www.nature.com/scientificdata

SCIENTIFIC DATA

Amended: Addendum

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

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 exemplar 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 governmental agencies are beginning to require data management and stewardship plans for data generated in publicly funded experiments. Beyond proper collection, annotation, and archival, data stewardship includes the notion of ‘long-term care’ of valuable digital assets, with the goal that they should be discovered and re-used for downstream investigations, either alone, or in combination with newly generated data. The outcomes from good data management and stewardship, therefore, are high quality digital publications that facilitate and simplify this ongoing process of discovery, evaluation, and reuse in downstream studies. What constitutes ‘good data management’ is, however, largely undefined, and is generally left as a decision for the data or repository owner. Therefore, bringing some clarity around the goals and desiderata of good data management and stewardship, and defining simple guideposts to inform those who publish and/or preserve scholarly data, would be of great utility.

This article describes four foundational principles—Findability, Accessibility, Interoperability, and Reusability—that serve to guide data producers and publishers as they navigate around these obstacles, thereby helping to maximize the added-value gained by contemporary, formal scholarly digital publishing. Importantly, it is our intent that the principles apply not only to ‘data’ in the conventional sense, but also to the algorithms, tools, and workflows that led to that data. All scholarly digital research objects—from data to analytical pipelines—benefit from application of these principles, since all components of the research process must be available to ensure transparency, reproducibility, and reusability.

There are numerous and diverse stakeholders who stand to benefit from overcoming these obstacles: researchers wanting to share, get credit, and reuse each other’s data and interpretations; professional data publishers offering their services; software and tool-builders providing data analysis and processing services such as reusable workflows; funding agencies (private and public) increasingly

Correspondence and requests for materials should be addressed to B.M. (email: barend.mons@dtls.nl).
*A full list of authors and their affiliations appears at the end of the paper.

- A minimal set of community agreed guiding principles and practices to ensure that research data is:
 - **F**indable
 - **A**ccessible
 - **I**nteroperable
 - **R**eusable
- Initially developed by Dutch Tech Centre for the Life Sciences
- Reviewed and refined through multi-stakeholder practitioner groups, including Force11 and the Research Data Alliance

What are the FAIR data principles



- F**indable - Assign persistent IDs
- Machine readable descriptions to support structured searches



- A**ccessible - Retrievable using a standard protocol
- Metadata available, even if data aren't
- Authentication and authorization procedure

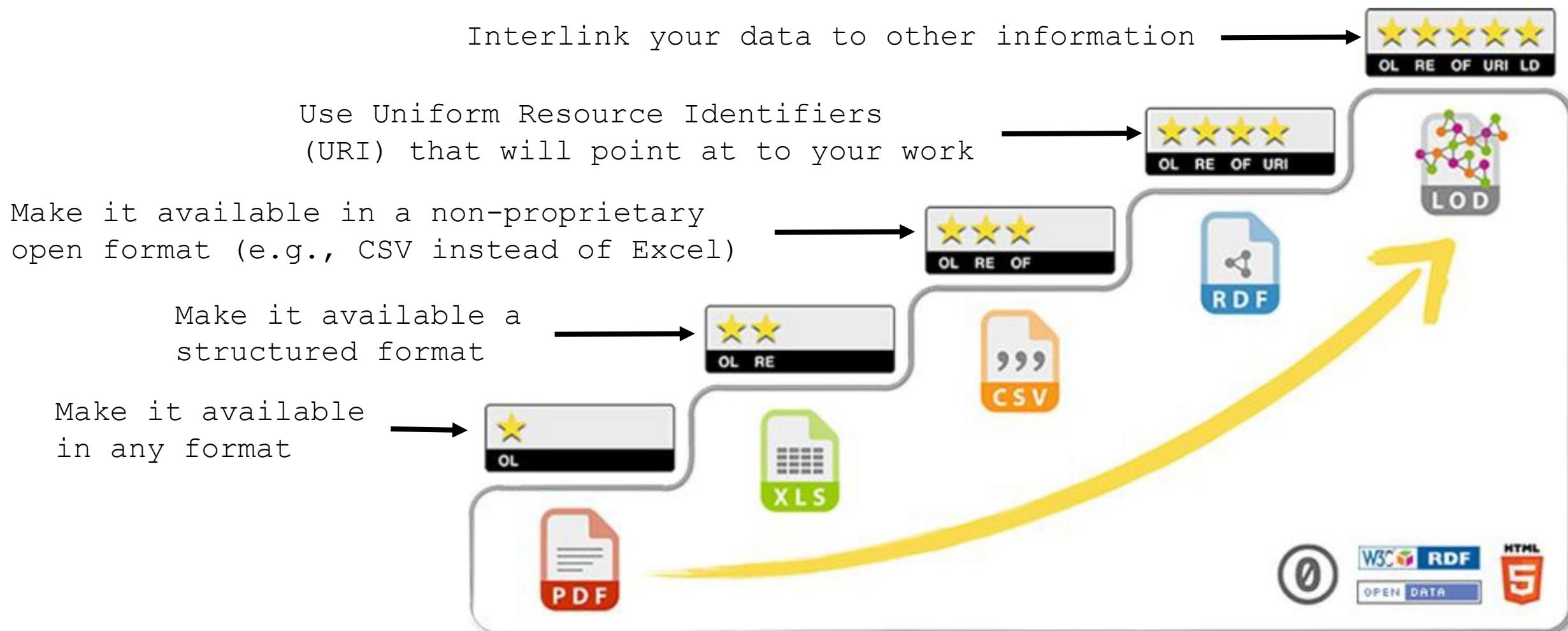


- I**nteroperable - Uses standard (FAIR) vocabularies
- Linked to other resources



- R**eusable - Clear licenses
- Provenance
- Meets domain-relevant community standards

A path towards FAIR



Incorporating FAIR into your routine workflow

F1000



Your go-to guide to making your data Findable, Accessible, Interoperable, and Reusable (FAIR)

So that you and others can get the most out of your data, it is important that you adhere to the [FAIR principles](#) to ensure your data are **Findable, Accessible, Interoperable, and Reusable** – whilst making your data openly available where it is safe to do so. This is no small task, so here are some ideas to help you get started:

1

Start with a management plan

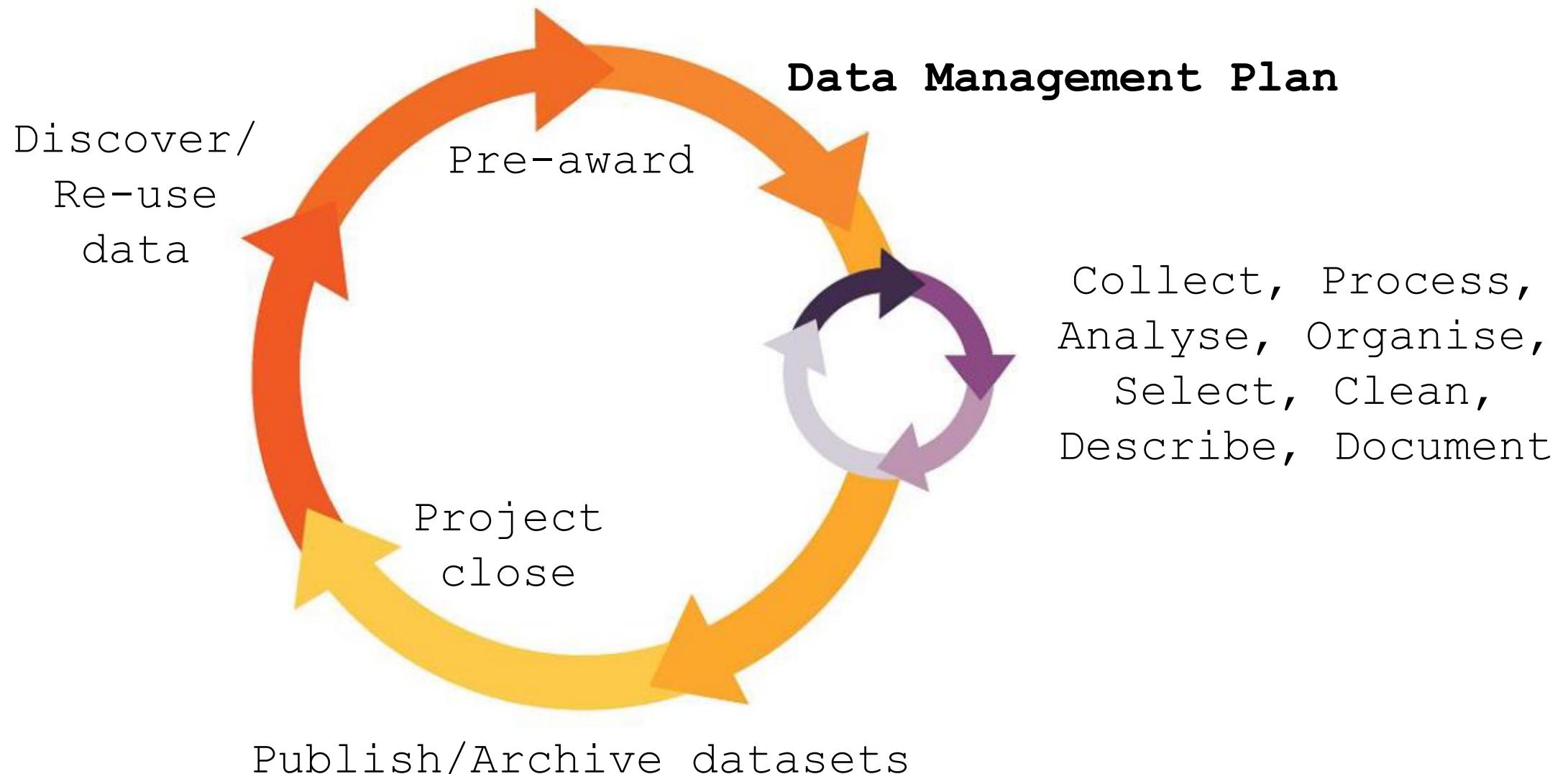
An output management plan (OMP) is a useful starting point for collecting or creating data, software, research materials, and intellectual property. Creating an OMP before you begin your research, and updating it throughout the research cycle, will help ensure that your outputs are as open and **FAIR** as possible when your project is complete.

Some funders require grant-holders to produce a plan as part of their application for funding, and/or after funding has been secured.

You should consider:

- What outputs you will be creating or collecting, and how these will be documented
- What ethical or legal requirements, if any, apply to the outputs
- How you will organise, store, secure, and share the outputs
- What resources are required and who is responsible

Incorporating FAIR into your routine workflow



Taking those first steps



Smart Data Management Plans for FAIR Open Science
For Serious Researchers and Data Stewards

How to use the Data Stewardship Wizard

We offer several options how to use the Data Stewardship Wizard, each suited for different use case.

Demo Instance	Researchers Instance	Self-hosted instance	Instance hosted by us
<p><i>For exploring the DSW features</i></p> <ul style="list-style-type: none">• Easy to sign up and use• A shared instance with other users• Not for serious usage	<p><i>For individual researchers</i></p> <ul style="list-style-type: none">• Easy to sign up and use• Ready to use Knowledge Models• Privacy and stability	<p><i>For organizations</i></p> <ul style="list-style-type: none">• All the DSW features available• Your own instance• You need to host and run the instance by yourself	<p><i>For organizations</i></p> <p>We offer managing the DS Wizard instance for interesting projects that want to use it seriously but don't want to run it by themselves.</p>

Current Phase

Before Submitting the DMP ▾

Chapters

I. Design of experiment 6

II. Data design and planning 7

III. Data Capture/Measurement 3

IV. Data processing and curation 11

V. Data integration 7

VI. Data interpretation 3

VII. Information and insight 8

More

Summary Report

IV. Data processing and curation

1 Workflow development

It is likely that you will be developing or modifying the workflow for data processing. There are a lot of aspects of this workflow that can play a role in your data management, such as the use of an existing work flow engine, the use of existing software vs development of new components, and whether every run needs human intervention or whether all data processing can be run in bulk once the work flow has been defined.

 Desirable: *Before Submitting the Proposal* a. This has been arranged b. More guidance is desired ⚙ Clear answer

2 How will you make sure to know what exactly has been run?

 Desirable: *Before Submitting the DMP* a. Explore ⚙ Clear answer

2.a.1 Will you keep results together with all processing scripts or workflows including documentation of the versions of the tools that have been run?

 Desirable: *Before Submitting the DMP* a. NoReusability b. YesReusability

2.a.2 Will you make use of the metadata fields in your output data files to register how the data was obtained?

File formats like VCF (for genetics) and TIFF (for images) have possibilities to document metadata in the file header. It is a good idea to use work flow tools that use these fields to document what was done to obtain the data.



Exploring standards in your field



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

HOW CAN WE HELP?

We guide consumers to discover, select and use these resources with confidence, and producers to make their resource more discoverable, more widely adopted and cited.



Research data facilitators, librarians, trainers

Use FAIRsharing to provide a foundation on which to create or enrich educational lectures, training and teaching material, and to plug into data management planning tools...
[\[read more\]](#)

Help building metadata

[PURPOSE](#)[RESEARCH](#)[TOOLS | TRAINING](#)[COMMUNITY](#)[ABOUT US](#)[Try CEDAR Now!](#)[YouTube](#)[SHARE](#)

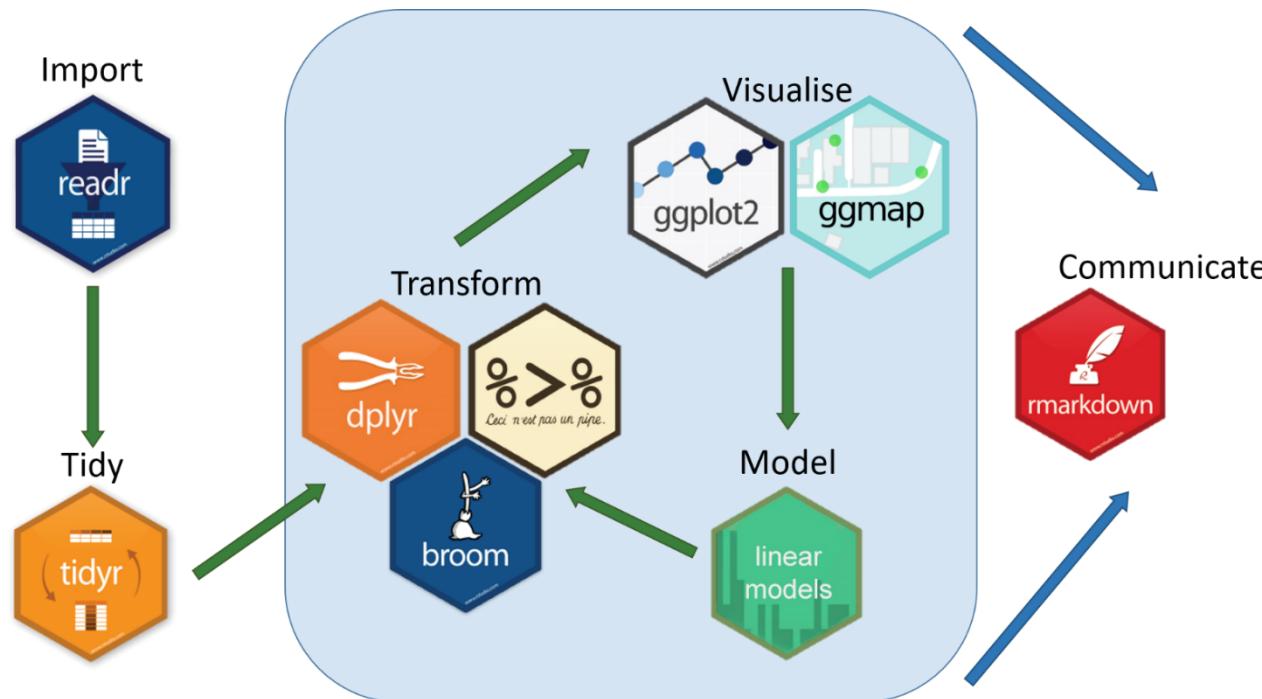
Better data for better science

[Home](#) › [Tools | Training](#)

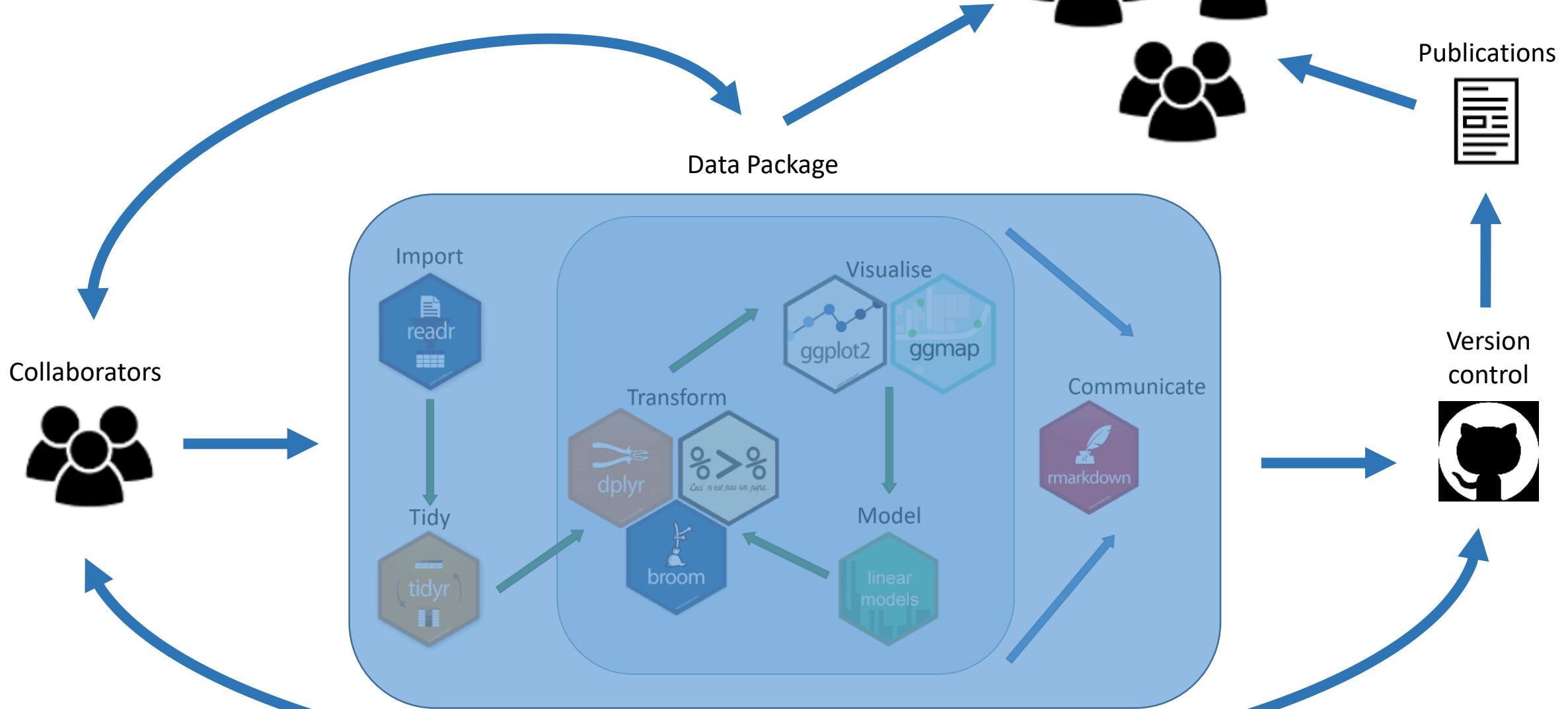
CEDAR Metadata Tools

When a biomedical scientist needs to upload her data and enter corresponding metadata into a repository, she is faced with a formidable task. Not only does she need to navigate and fill out many forms to enter (and re-enter!) information, and make sure everything is cross-referenced correctly, but the metadata frequently end up stored in an ad hoc manner, in a non-standard format, and using non-standard terminology. As a result, finding or reusing the metadata, or understanding the underlying experiments, becomes extremely hard, if not impossible. But when the scientist uses CEDAR, our tools can make describing laboratory studies—or metadata for any other biomedical content—much easier.

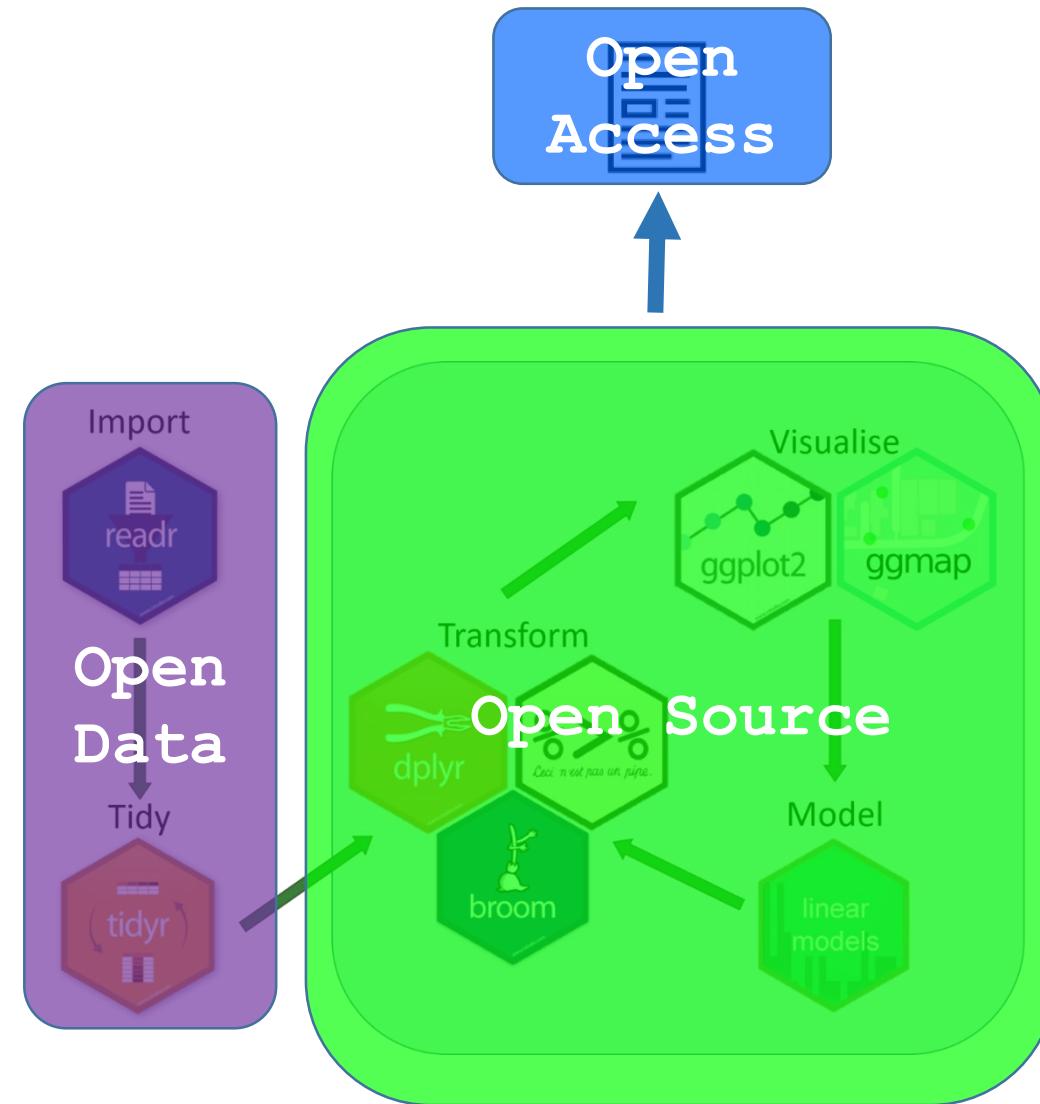
Putting the pieces together using R



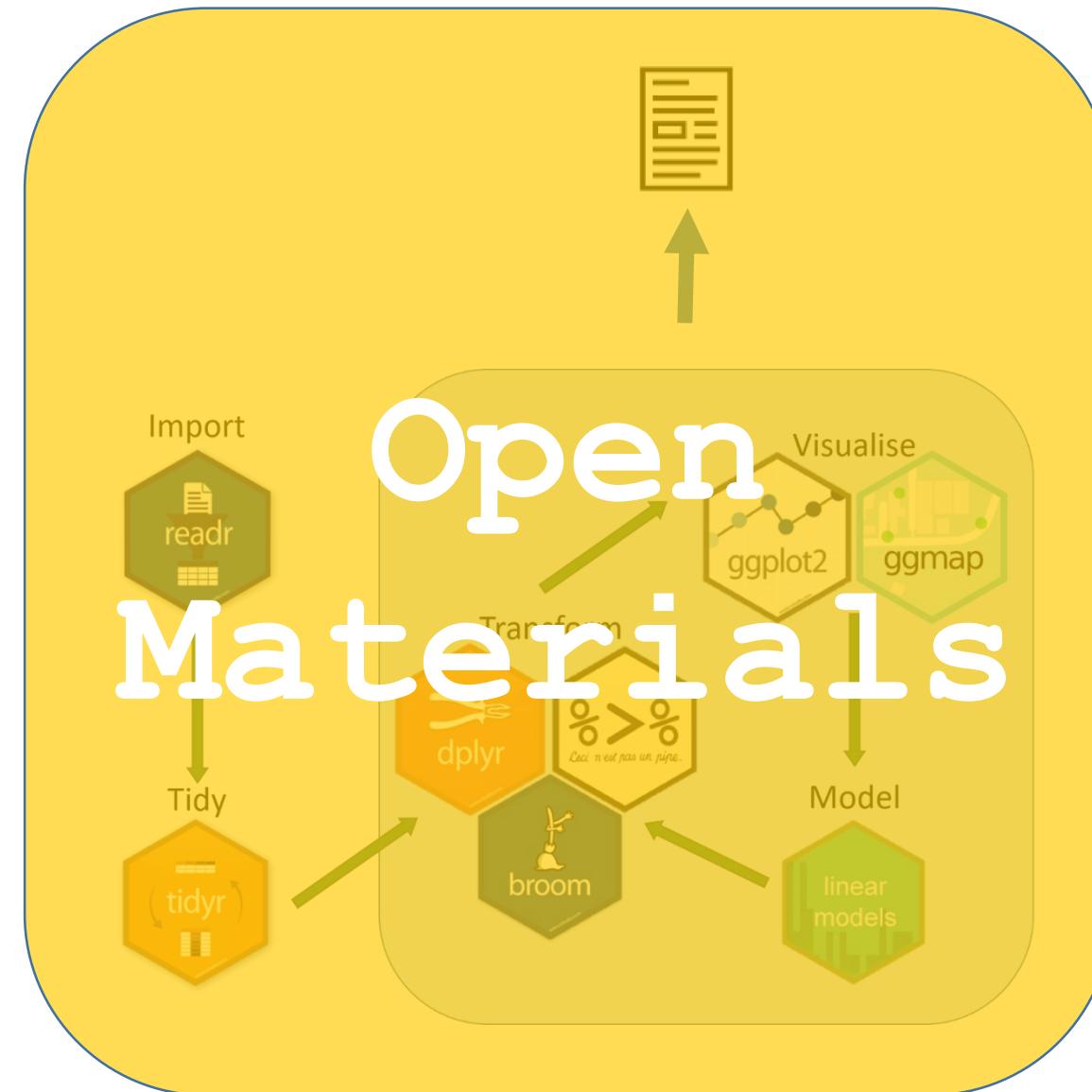
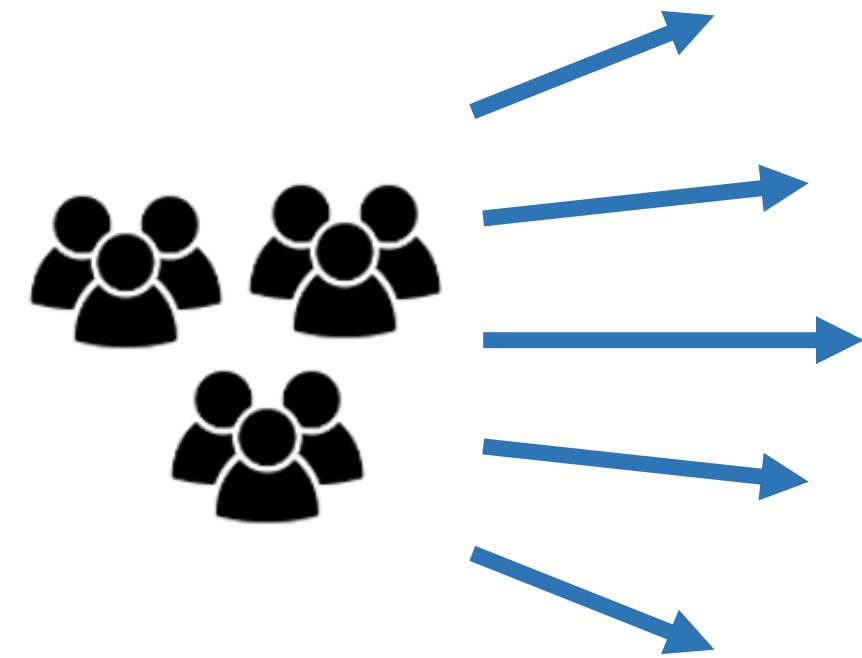
The bigger picture



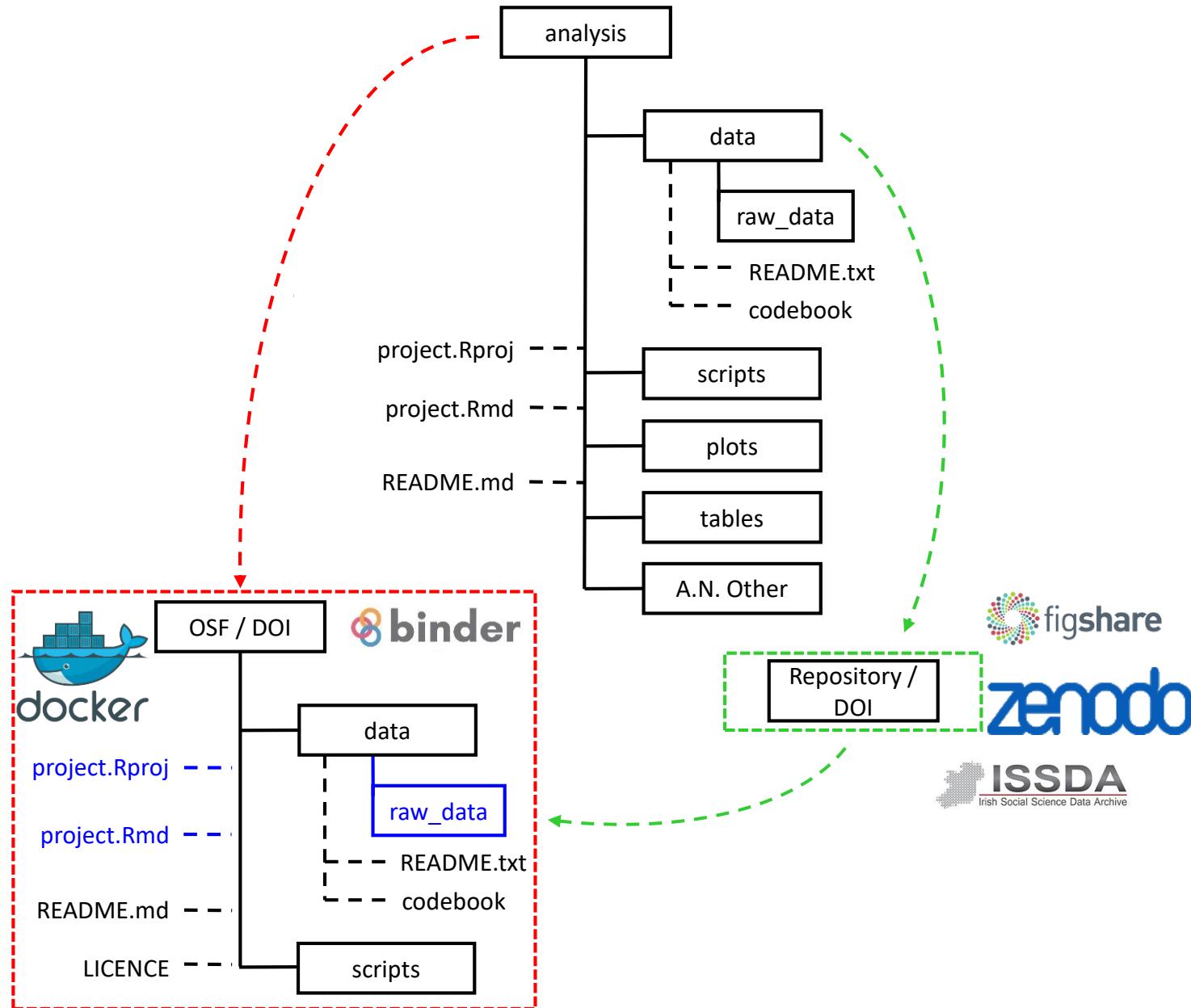
The ‘Open Science’ picture



The ‘Open Science’ picture



What does this allow us to do?



Our real life experiment



- UV light has potential to change the secondary metabolite composition (colour) of bronze/red lettuce
- Experimental setup:
 - 3 lettuce varieties
 - 3 UV filter conditions
 - 3 week duration

Real data comes with real problems

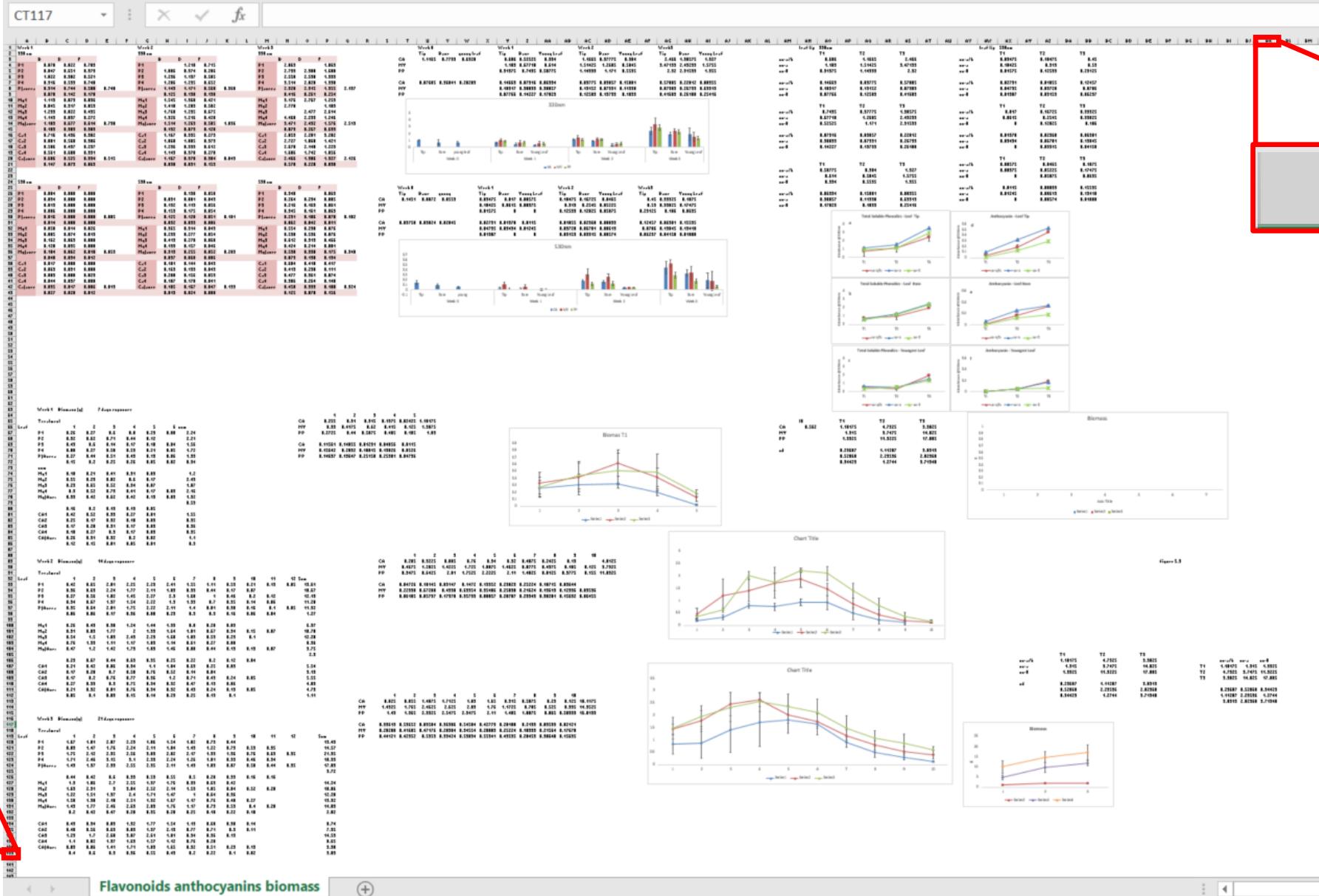
Raw Data wk 1-3 Lettuce Exp 1 - Excel

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
1	Week 1						Week 2						Week 3					
2	330 nm						330 nm						330 nm					
3		B	D	F				B	D	F				B	D	F		
4	P1	0.870	0.822	0.703			P1						1	2.869		1.069		
5	P2	0.847	0.651	0.379			P2						2	2.739	2.380	1.688		
6	P3	1.022	0.902	0.521			P3	1.236	1.197	0.585			P3	2.558	2.538	1.333		
7	P4	0.916	0.599	0.748			P4	1.206	1.295	0.652			P4	3.514	2.028	1.330		
8	P(average)	0.914	0.744	0.588	0.748		P(average)	1.149	1.171	0.560	0.960		P(average)	2.920	2.315	1.355	2.197	
9		0.078	0.142	0.170				0.125	0.138	0.190				0.416	0.261	0.254		
10	My1	1.119	0.873	0.896			My1	1.545	1.360	0.421			My1	3.176	2.767	1.259		
11	My2	0.845	0.917	0.853			My2	1.418	1.203	0.502			My2	2.778		1.183		
12	My3	1.299	0.822	0.435			My3	1.768	1.295	0.675			My3		2.477	2.614		
13	My4	1.149	0.097	0.272			My4	1.326	1.216	0.420			My4	4.460	2.233	1.246		
14	My(average)	1.103	0.677	0.614	0.798		My(average)	1.514	1.269	0.505	1.096		My(average)	3.471	2.492	1.576	2.513	
15		0.189	0.389	0.309				0.192	0.073	0.120				0.879	0.267	0.693		
16	Ca1	0.716	0.496	0.382			Ca1	1.167	0.935	0.273			Ca1	2.853	2.201	3.202		
17	Ca2	0.881	0.568	0.386			Ca2	1.060	1.005	0.373			Ca2	2.727	1.860	1.421		
18	Ca3	0.586	0.437	0.237			Ca3	1.296	0.993	0.612			Ca3	2.678	2.140	1.229		
19	Ca4	0.561	0.600	0.331			Ca4	1.143	0.978	0.278			Ca4	1.606	1.742	1.856		
20	Ca(average)	0.686	0.525	0.334	0.515		Ca(average)	1.167	0.978	0.384	0.843		Ca(average)	2.466	1.986	1.927	2.126	
21		0.147	0.073	0.069				0.098	0.031	0.159				0.578	0.220	0.890		
22																		
23																		
24	530 nm						530 nm						530 nm					
25		B	D	F				B	D	F				B	D	F		
26	P1	0.004	0.000	0.000			P1		0.138	0.050				P1	0.340		0.069	
27	P2	0.034	0.000	0.000			P2		0.091	0.081	0.043			P2	0.264	0.234	0.085	CA
28	P3	0.019	0.000	0.000			P3		0.132	0.119	0.056			P3	0.216	0.163	0.061	MY

File Home Insert Page Layout Formulas Data Review View Tell me what you want to do...

Normal Page Break Preview Custom Layout Views Workbook Views

Ruler Formula Bar Gridlines Headings Zoom 100% Zoom to Selection Window New Arrange Freeze All Panes Hide Synchronous Scrolling Reset Window Position Window Switch Windows Macros Macros



Take small steps to enact big changes

THE AMERICAN STATISTICIAN
2018, VOL. 72, NO. 1, 2–10
<https://doi.org/10.1080/00031305.2017.1375989>



OPEN ACCESS



Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

ABSTRACT

Spreadsheets are widely used software tools for data entry, storage, analysis, and visualization. Focusing on the data entry and storage aspects, this article offers practical recommendations for organizing spreadsheet data to reduce errors and ease later analyses. The basic principles are: be consistent, write dates like YYYY-MM-DD, do not leave any cells empty, put just one thing in a cell, organize the data as a single rectangle (with subjects as rows and variables as columns, and with a single header row), create a data dictionary, do not include calculations in the raw data files, do not use font color or highlighting as data, choose good names for things, make backups, use data validation to avoid data entry errors, and save the data in plain text files.

ARTICLE HISTORY

Received June 2017
Revised August 2017

KEYWORDS

Data management; Data organization; Microsoft Excel; Spreadsheets

Less stress, more success

	A	B	C	D	E	F	G	H	I	J	K	L
1	id	week_no	filter_nam	treatment	replicate_no	flavonoids	biomass	variety	date	investigator		
2	1	0	ptp	nofilter	1	1.061	0.39	cos	2019/04/01	Darren Dahly		
3	2	0	ptp	nofilter	2	1.1805	0.42	cos	2019/04/01	Darren Dahly		
4	3	0	ptp	nofilter	3	1.0345	0.62	cos	2019/04/01	Darren Dahly		
5	4	0	ptp	nofilter	4	1.094	0.63	cos	2019/04/01	Brendan Palmer		
6	1	0	my	nofilter	1	1.061	0.39	cos	2019/04/01	Brendan Palmer		
7	2	0	my	nofilter	2	1.1805	0.42	cos	2019/04/01	Brendan Palmer		
8	3	0	my	nofilter	3	1.0345	0.62	cos	2019/04/01	Brendan Palmer		
9	4	0	my	nofilter	4	1.094	0.63	cos	2019/04/01	Brendan Palmer		
10	1	0	ca	nofilter	1	1.061	0.39	cos	2019/04/01	Brendan Palmer		
11	2	0	ca	nofilter	2	1.1805	0.42	cos	2019/04/01	Brendan Palmer		
12	3	0	ca	nofilter	3	1.0345	0.62	cos	2019/04/01	Brendan Palmer		
13	4	0	ca	nofilter	4	1.094	0.63	cos	2019/04/01	Darren Dahly		
14	5	1	ptp	filter	1	0.87	0.76	cos	2019/04/08	Darren Dahly		
15	6	1	ptp	filter	2	0.847	0.95	cos	2019/04/08	Darren Dahly		
16	7	1	ptp	filter	3	1.022	0.95	cos	2019/04/08	Darren Dahly		
17	8	1	ptp	filter	4	0.916	0.95	cos	2019/04/08	Darren Dahly		
18	9	1	my	filter	1	1.119	1.55	cos	2019/04/08	Darren Dahly		
19	10	1	my	filter	2	0.845	3.16	cos	2019/04/08	Darren Dahly		
20	11	1	my	filter	3	1.299	4.9	cos	2019/04/08	Brendan Palmer		
21	12	1	my	filter	4	1.149	5.5	cos	2019/04/08	Brendan Palmer		
22	13	1	ca	filter	1	0.716	5.5	cos	2019/04/08	Brendan Palmer		
23	14	1	ca	filter	2	0.881	7.94	cos	2019/04/08	Brendan Palmer		
24	15	1	ca	filter	3	0.586	8.71	cos	2019/04/08	Brendan Palmer		
25	16	1	ca	filter	4	0.561	8.71	cos	2019/04/08	Brendan Palmer		
26	17	2	ptp	filter	1	0	14.45	cos	2019/04/15	Brendan Palmer		
27	18	2	ptp	filter	2	1.006	2.14	cos	2019/04/15	Brendan Palmer		
28	19	2	ptp	filter	3	1.236	1.86	cos	2019/04/15	Brendan Palmer		
29	20	2	ptp	filter	4	1.206	1.2	cos	2019/04/15	Brendan Palmer		
30	21	2	mv	filter	1	1.545	2.45	cos	2019/04/15	Brendan Palmer		

data

dictionary

values



Less stress, more success

Less stress, more success

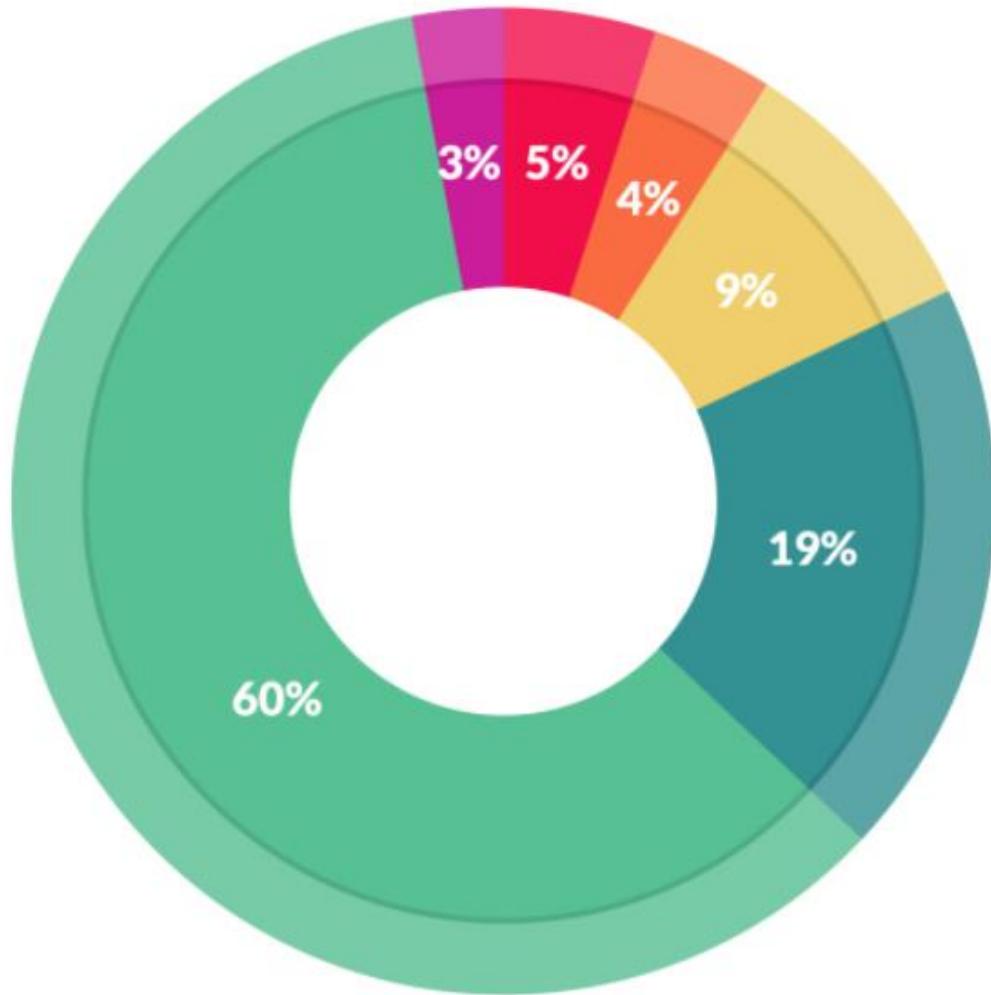
1	A	B	C	D	E	F	G	H	I	J	K	L
2	1	0	ptp	nofilter	1	1.061	0.39	cos	2019/04/01	Darren Dahly		
3	2	0	ptp	A	B	C	D	E				
4	3	0	ptp	1	field_name	data_type	data_format	example	standard_units	description		
5	4	0	ptp	2	id	numeric	integer	23	NA	Unique identifier applied to each observation		
6	1	0	my	3	week_no	numeric	integer					
7	2	0	my	4	filter_name	character	NA					
8	3	0	my	5	treatment	character	NA					
9	4	0	my	6	replicate_no	numeric	integer					
10	1	0	ca	7	flavonoids	numeric	double					
11	2	0	ca	8	biomass	numeric	double					
12	3	0	ca	9	variety	character	NA					
13	4	0	ca	10	date	date	YYYY/MM/DD					
14	5	1	ptp	11	investigator	character	Firstname Lastname					
15	6	1	ptp	12								
16	7	1	ptp	13								
17	8	1	ptp	14								
18	9	1	my	15								
19	10	1	my	16								
20	11	1	my	17								
21	12	1	my	18								
22	13	1	ca	19								
23	14	1	ca	20								
24	15	1	ca	21								
25	16	1	ca	22								
26	17	2	ptp	23								
27	18	2	ptp	24								
28	19	2	ptp	25								
29	20	2	ptp	26								
30	21	2	mv	27								
		data	dictionary	28								
				29								
				30								

The screenshot shows a data entry interface with two tabs: 'data' and 'dictionary'. The 'data' tab displays a grid of experimental data. The 'dictionary' tab provides a detailed schema for each column, including field_name, data_type, data_format, example, standard_units, and description. A tooltip for 'id' indicates it is a unique identifier applied to each observation.

Below the tabs, there are navigation buttons for the data grid: back, forward, data, dictionary, values, and a plus sign.

Less stress, more success

Resources are being wasted by not doing this



What data scientists spend the most time doing

- *Building training sets: 3%*
- *Cleaning and organizing data: 60%*
- *Collecting data sets; 19%*
- *Mining data for patterns: 9%*
- *Refining algorithms: 4%*
- *Other: 5%*

Try it out yourself

Screenshot of a GitHub profile page for Brendan Palmer (@bapalmer). The profile features a photo of a puppet with orange hair and a green jacket, and includes pinned projects related to R and reproducible research.

Brendan Palmer
bapalmer

[Edit profile](#)

Overview Repositories 14 Projects 0 Stars 1 Followers 12 Following 10

Pinned Order updated. Customize your pins

- RSS_Belfast_2019**
Data FAIRification using R/RStudio workflows
R
- R-A_Hitchhikers_Guide_to_Reproducible_Research**
A 3-day R course given in University College Cork that encompasses various elements off reproducible research facilitated through RStudio projects, the R tidyverse language and reporting using R Ma...
HTML ★ 2
- RCR**
Section of the UCC Reproducible Conduct of Research digital badge dedicated to exposing researchers to reproducible research practices.
HTML
- lunchtime_sessions**
Short 1 hour introductions to R-related topics such as creating R projects, using GitHub through RStudio and more
HTML ★ 1

Putting the final pieces into place

Make Your Code Citable Using GitHub and Zenodo: A How- to Guide

By [Open Science MOOC](#) on July 24, 2018

