

# The role of journals in (and publishers) in data sharing

Open in practice, Reading  
30<sup>th</sup> March 2017

Iain Hrynaszkiewicz

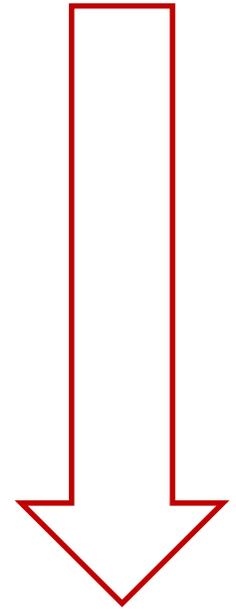
**SPRINGER NATURE**



## Not all research data are open data

Different levels of openness in research data publishing:

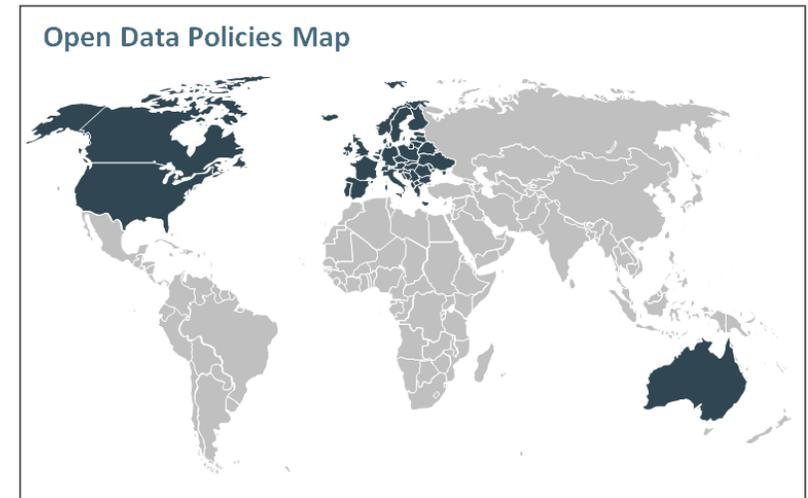
1. Accessible only to an individual researcher/group
2. Accessible to others on (reasonable) request
3. Published as electronic supplementary material
4. Deposited in a general or institutional data repository (e.g. figshare)
5. Deposited in a subject/community specific data repository



**More open**

# Research funders' policies

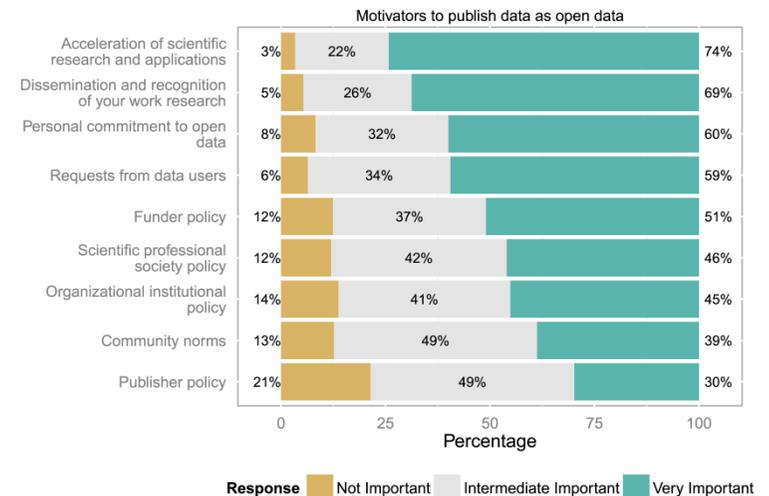
- More than 40 research funders globally have policies or mandates that require data archiving as a condition of grants (1), including:
  - National Science Foundation (NSF)
  - National Institutes of Health (NIH)
  - Wellcome Trust
  - Bill and Melinda Gates Foundation
  - Research Councils UK
    - Medical Research Council
    - BBSRC
    - ESRC
- **Some of these require data to be linked to publications including:**
  - Engineering and Physical Sciences Research Council (EPSRC)



(1) Hahnel, M: Global funders who require data archiving as a condition of grants. *figshare*.  
<https://dx.doi.org/10.6084/m9.figshare.1281141.v1> (2015)

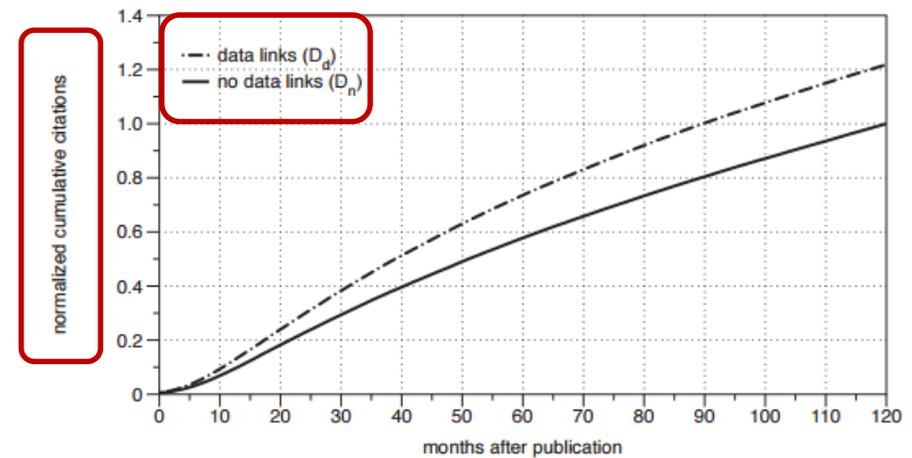
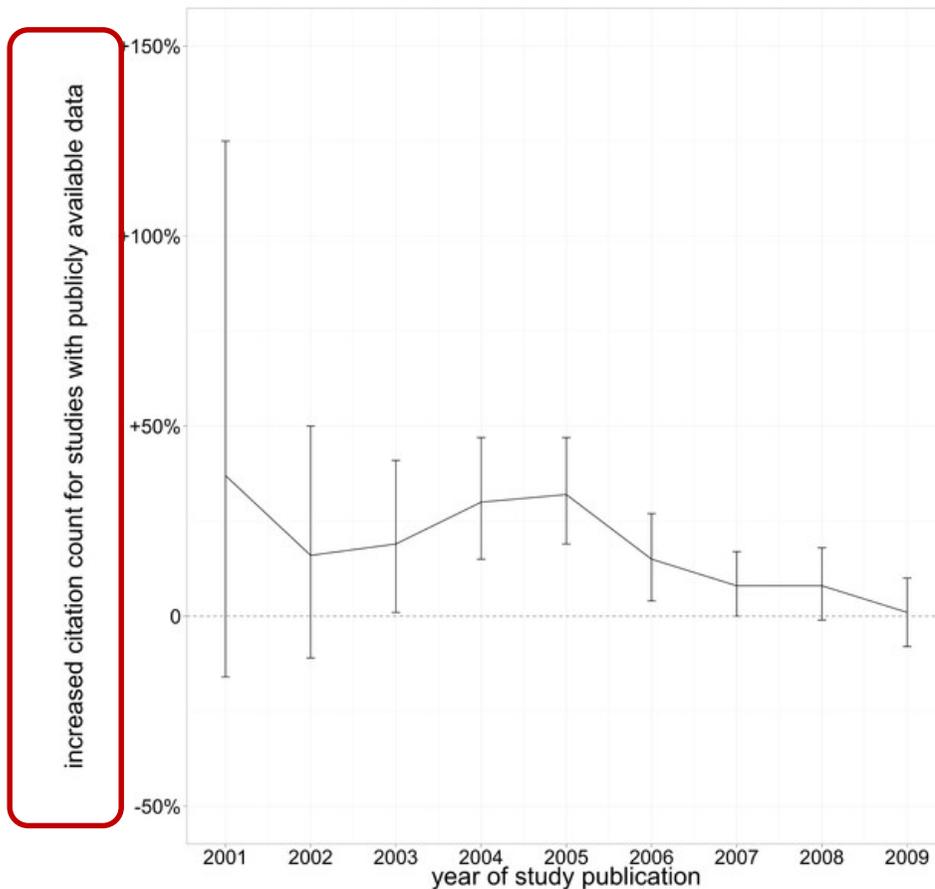
# What motivates researchers to share data?

- **97%** - to accelerate research and its applications<sup>1</sup>
- **96%** - increased visibility and discovery of their research data<sup>1,2</sup>
- **95%** - increased usability of their research data<sup>2</sup>
- **>90%** - credit mechanism for deposit of data<sup>1,2</sup>
- **88%** - to comply with funder policy<sup>1</sup>



1. Schmidt et al. (2016). PLoS ONE 11(1): e0146695. doi:10.1371/journal.pone.0146695 (n=1248) (& image credit, CC BY)
2. Nature Publishing Group (2014): Data publication survey - raw data. figshare. <http://dx.doi.org/10.6084/m9.figshare.1234052> (n=387)

# Publicly available data/links to data may = more citations

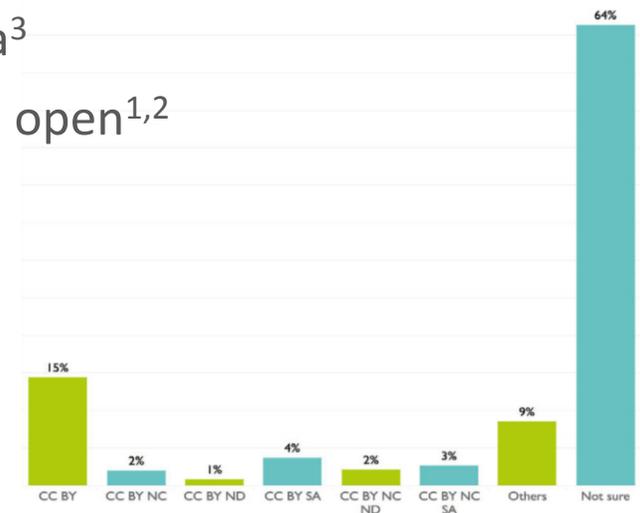


Source: <http://arxiv.org/pdf/1111.3618v1.pdf>

Source: <https://doi.org/10.7717/peerj.175/fig-2>

# What challenges do researchers face?

- **64%** unsure about open licensing of research data<sup>1</sup>
- **56%** do not use a metadata standard<sup>2</sup>
- **54%** would like more guidance complying with funder policies<sup>1</sup>
- **54%** do not have enough time to make data available<sup>1</sup>
- **45%** unaware of a repository for some of their data<sup>3</sup>
- **39%** uncertain about meeting costs of making data open<sup>1,2</sup>



1. Treadway et al. (2016). figshare. <https://dx.doi.org/10.6084/m9.figshare.4036398.v1> (n= 2061) (& image credit CC BY)
2. Tenopir et al. (2011). PLoS ONE 6(6): e21101. doi:10.1371/journal.pone.0021101 (n=1315)
3. Nature Publishing Group (2014). figshare. <http://dx.doi.org/10.6084/m9.figshare.1234052> (n=387)

# What are journals & publishers doing about it?

- **Content types** e.g. data articles and journals
- **Credit and incentives** e.g. data citation and data articles
- **Encouraging reuse** e.g. open licenses
- **Data quality** e.g. data peer review
- **Data discoverability** e.g. linking data to publications; supporting repositories
- **Raising awareness** e.g. editorials, outreach
- **Guidance and policy** e.g. information for authors, policy harmonization
- **Technology** e.g. platform developments, repository integration and other partnerships

# What are journals & publishers doing about it?

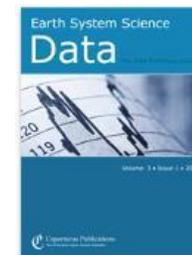
- **Content types** e.g. data articles and journals
- **Credit and incentives** e.g. data citation and data articles
- **Encouraging reuse** e.g. open licenses
- **Data quality** e.g. data peer review
- **Data discoverability** e.g. linking data to publications; supporting repositories
- **Raising awareness** e.g. editorials, outreach
- **Guidance and policy** e.g. information for authors, policy harmonization
- **Technology** e.g. platform developments, repository integration and other partnerships

SCIENTIFIC DATA 



(GIGA)<sup>n</sup>  
SCIENCE 

  
BMC  
Research Notes



SPRINGER NATURE

# Scientific Data

nature.com > scientific data

a natureresearch journal

MENU

SCIENTIFIC DATA

Search E-alert Submit My account

Data Descriptor | 07 June 2016 | OPEN

## Evaluating vertebrate fossils with X-ray computed tomography

Data Descriptor | 07 June 2016 | OPEN

### Spatializing 6,000 years of global urbanization from 3700 BC to AD 2000

Meredith Reba, Femke Reitsma & Karen C. Seto

Data Descriptor | 07 June 2016 | OPEN

### Extensive sequencing of seven human genomes to characterize benchmark reference materials

Justin M. Zook, David Catoe [...] Marc Salit

Announcement

### The #scidata16 draft programme and call for lightning talks

June 7 | by Mat Astell

scientificdata updates

# The data descriptor article

**SCIENTIFIC DATA**

Home | Archive | About | For Authors | For Referees | Advisory & Editorial Board | Data Policies

Home > Data Descriptors > Data Descriptor

SCIENTIFIC DATA | DATA DESCRIPTOR **OPEN**

## Time-resolved gene expression profiling during reprogramming of C/EBP $\alpha$ -pulsed B cells into iPS cells

Bruno Di Stefano, Samuel Collombet & Thomas Graf

Affiliations | Contributions | Corresponding authors

Scientific Data 1, Article number: 140008 | doi:10.1038/sdata.2014.8  
Received 21 February 2014 | Accepted 22 April 2014 | Published online 27 May 2014

PDF | ISA tab | Citation | Reprints | Rights & permissions | Article metrics

**About Scientific Data**  
Scientific Data is an open-access, peer-reviewed publication for descriptions of scientifically valuable datasets. Our primary article-type, the Data Descriptor, is designed to make your data more discoverable, interpretable and reusable.

E-alert | RSS | Facebook | Twitter

**Associated Links**  
Nature | Article  
C/EBP $\alpha$  paises B cells for into induced pluripotent stem cells by Bruno Di Stefano et al

- **Associated Nature Article**
- Data at **figshare** & NCBI GEO
- Integrated **figshare** data viewer

## Data Citations

[Abstract](#) • [Background and Summary](#) • [Methods](#) • [Data Records](#) • [Technical Validation](#) • [Usage Notes](#) • [Additional information](#) • [References](#) • [Data Citations](#) • [Acknowledgements](#) • [Author information](#)

1. Di Stefano, B., Collombet, S., & Graf, T. *Gene Expression Omnibus GSE46321* (2014).
2. Di Stefano, B., Collombet, S., & Graf, T. *Gene Expression Omnibus GSE52396* (2014).
3. Di Stefano, B., Collombet, S., & Graf, T. *Figshare* <http://dx.doi.org/10.6084/m9.figshare.939408> (2014).

**NCBI GEO - Accession Display**

Series: GSE52396

Series	Status	Publ
GSE52396	Public	C/EBP $\alpha$
		Misc
		Expn

**figshare**

AIRProbes\_AIRReplicates.xls

	A	B	C	D
1	GeneName	AP1_1	AP1_2	AP1_1
2	ProbeID	6.66662102	6.15802008	6.58027107
3	Accession	6.17811101	6.17120214	6.27110804
4	Accession	12.18803304	12.18803304	12.27111101
5	Accession	12.28236607	12.28236607	12.27111101
6	Accession	12.10233304	12.20237108	12.28236607
7	Accession	6.48027107	6.48027107	6.17120214
8	Accession	6.25188033	6.15802008	6.27110804

# What are journals & publishers doing about it?

- **Content types** e.g. data articles and journals
- **Credit and incentives** e.g. data citation and data articles
- **Encouraging reuse** e.g. open licenses
- **Data quality** e.g. data peer review
- **Data discoverability** e.g. linking data to publications; supporting repositories
- **Raising awareness** e.g. editorials, outreach
- **Guidance and policy** e.g. information for authors, policy harmonization
- **Technology** e.g. platform developments, repository integration and other partnerships



**DC<sup>1</sup>**

*Data Citation Principles*



# What is a data citation?

- *A reference made to data in the same way as researchers routinely provide a bibliographic reference to journal articles and books*
- When public datasets have Digital Object Identifiers (DOIs), **very similar to citing a journal article**
- Include the minimum information recommended by DataCite and follow Nature style i.e.
  - authors, title, publisher (repository name), identifier, year

Creator (author)

Title

77. Di Stefano, B., Collombet, S. & Graf, T. Time-resolved gene expression profiling during reprogramming of C/EBP $\alpha$ -pulsed B cells into iPS cells. *figshare* [https://dx.doi.org/10.6084/m9.figshare.939408\\_D1](https://dx.doi.org/10.6084/m9.figshare.939408_D1) (2014).

[+ Show context](#)

Publisher (repository name)

Identifier (DOI)

Publication year

SPRINGER NATURE

## Data are increasingly recognised as a 1st class object

HEFCE and the REF: <http://blog.hefce.ac.uk/2015/09/01/opening-up-research-data/> *“Via the REF, we fully recognise data as an equally valid form of research output, and, through our open access policy for the next REF, we plan to reward research environments that deliver open access to a wider set of outputs than just journal articles and conference papers.”*

National Science Foundation (US) *“For all new grant applications from 14 January, the US National Science Foundation (NSF) asks a principal investigator to list his or her research “products” rather than “publications” in the biographical sketch section. This means that, according to the NSF, a scientist’s worth is not dependent solely on publications. Data sets, software and other non-traditional research products will count too.”*

<http://www.nature.com/nature/journal/v493/n7431/full/493159a.html>

# What are journals & publishers doing about it?

- **Content types** e.g. data articles and journals
- **Credit and incentives** e.g. data citation and data articles
- **Encouraging reuse** e.g. open licenses
- **Data quality** e.g. data peer review
- **Data discoverability** e.g. linking data to publications; supporting repositories
- **Raising awareness** e.g. editorials, outreach
- **Guidance and policy** e.g. information for authors, policy harmonization
- **Technology** e.g. platform developments, repository integration and other partnerships



# What are journals & publishers doing about it?

- **Content types** e.g. data articles and journals
- **Credit and incentives** e.g. data citation and data articles
- **Encouraging reuse** e.g. open licenses
- **Data quality** e.g. data peer review
- **Data discoverability** e.g. linking data to publications; supporting repositories
- **Raising awareness** e.g. editorials, outreach
- **Guidance and policy** e.g. information for authors, policy harmonization
- **Technology** e.g. platform developments, repository integration and other partnerships

SCIENTIFIC DATA 

# Making data peer review a reality

## Peer review at *Scientific Data* focuses on:

- Completeness (can others reproduce?)
- Consistency (were community standards followed?)
- Integrity (are data in the best repository?)
- Experimental rigour and technical quality (were the methods sound?)

## Does not focus on:

- Perceived impact/importance
- Size/complexity of data

# What are journals & publishers doing about it?

- **Content types** e.g. data articles and journals
- **Credit and incentives** e.g. data citation and data articles
- **Encouraging reuse** e.g. open licenses
- **Data quality** e.g. data peer review
- **Data discoverability** e.g. linking data to publications; supporting repositories
- **Raising awareness** e.g. editorials, outreach
- **Guidance and policy** e.g. information for authors, policy harmonization
- **Technology** e.g. platform developments, repository integration and other partnerships



# Helping researchers find the right repository

- Recommended repositories list (>80 repositories)
- <http://www.springernature.com/gp/group/data-policy/repositories>

## *What makes a good data repository?*

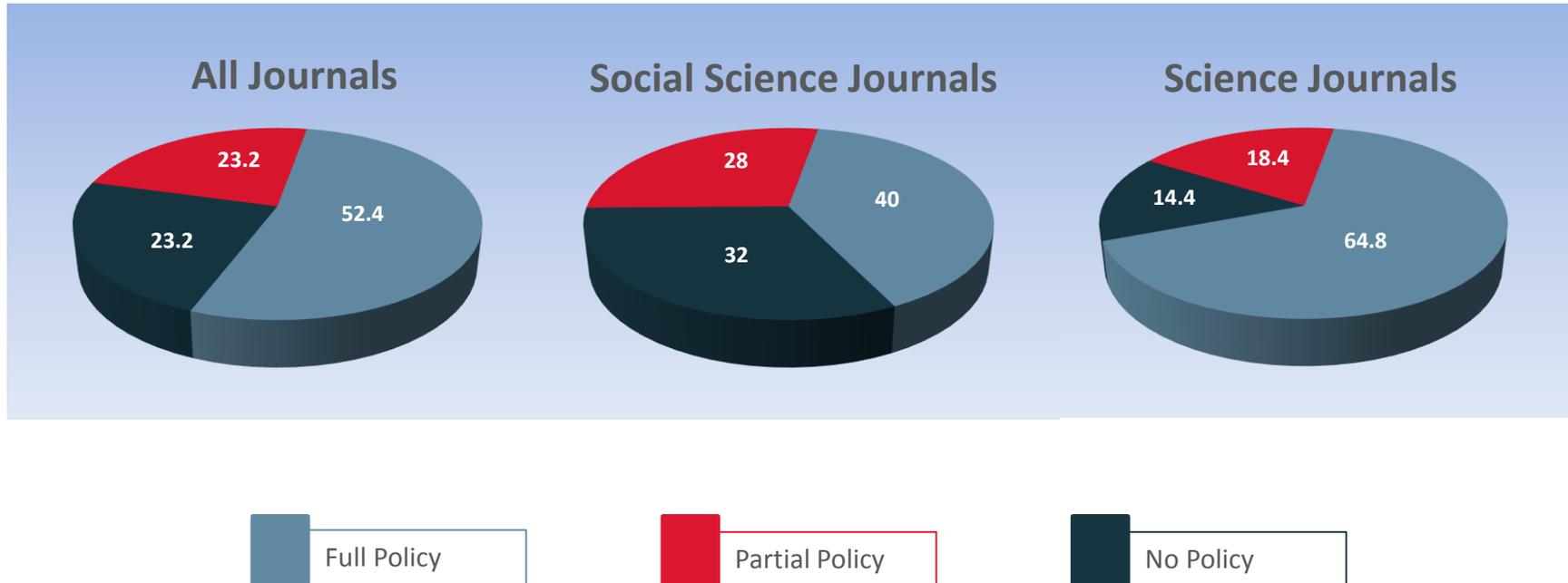
Data repositories for data supporting peer-reviewed publications generally should:

- I. Ensure long-term persistence and preservation of datasets
- II. Be recognized by a research community or research institution
- III. Provide deposited datasets with stable and persistent identifiers, such as Digital Object Identifiers (DOIs)
- IV. Allow access to data without unnecessary restrictions
- V. Provide a clear license or terms of use for deposited datasets

# What are journals & publishers doing about it?

- **Content types** e.g. data articles and journals
- **Credit and incentives** e.g. data citation and data articles
- **Encouraging reuse** e.g. open licenses
- **Data quality** e.g. data peer review
- **Data discoverability** e.g. linking data to publications; supporting repositories
- **Raising awareness** e.g. editorials, outreach
- **Guidance and policy** e.g. information for authors, policy harmonization
- **Technology** e.g. platform developments, repository integration and other partnerships

# How many journals have a research data policy?

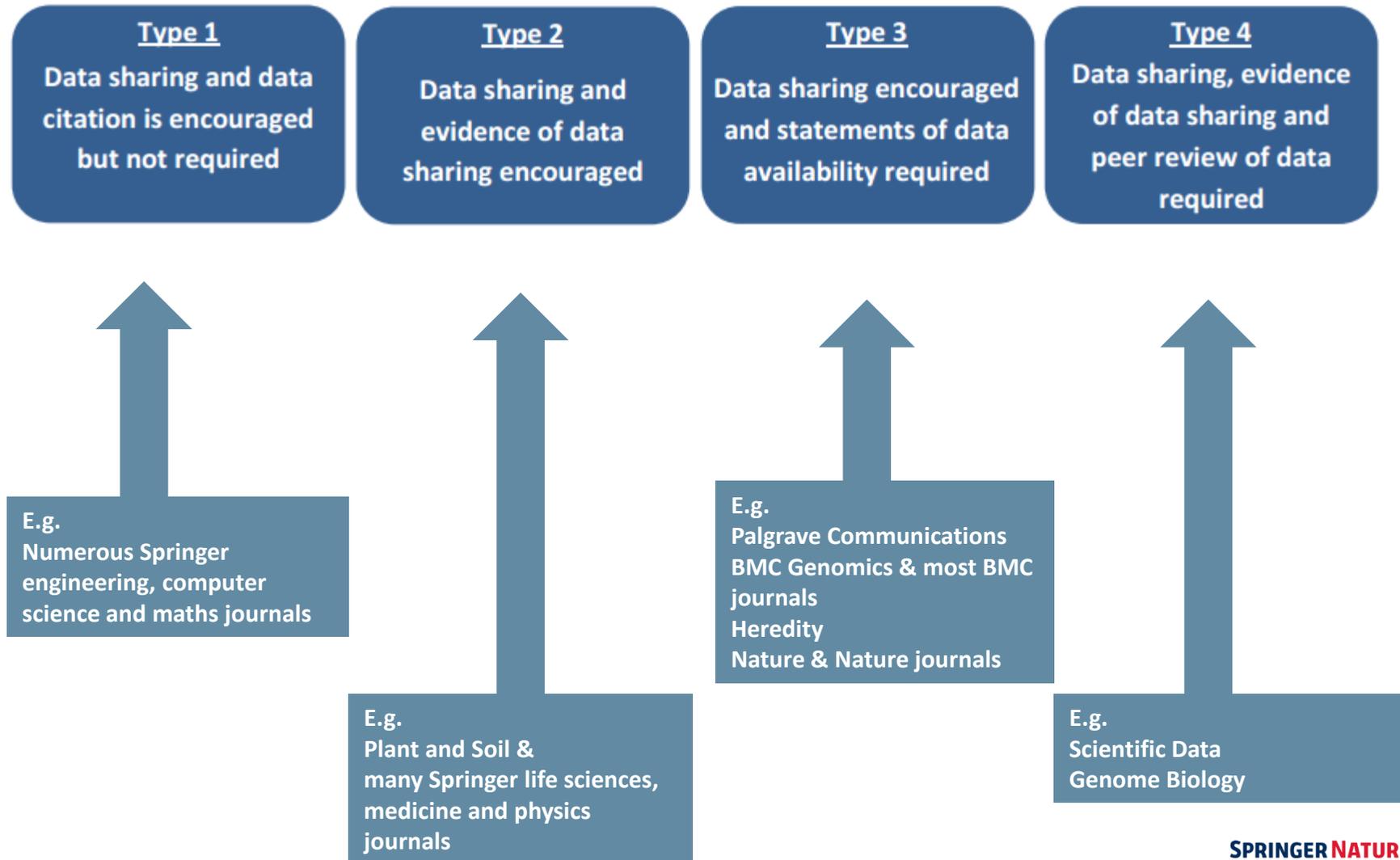


Data source: Linda Naughton, JISC Journal Research Data Policy Bank project presentation (n = 250)

***“The evidence shows that the current research data policy ecosystem is in critical need of standardization and harmonization”***

-- Naughton, L. & Kernohan, D., (2016). Making sense of journal research data policies. Insights. 29(1), pp.84–89. DOI: <http://doi.org/10.1629/uksg.284>

## Policy Types



# Research Data Support helpdesk @Springer Nature

## Support for editors:

- Identifying and implementing a data policy
- Identifying data repositories for their audience(s)
- Dealing with peer review of sensitive/human data
- Good practice for data-literature integration

## Support for authors:

- Information on the data policy of their target journal(s)
- Identifying and using data repositories
- Compliance with funders' and institutions' data sharing policies
- Data reporting standards

**SPRINGER NATURE** About Us Responsible Business Careers Media Contact

---

### Research Data Support Helpdesk

Springer Nature provides a research data policy support service for authors and editors, which can be contacted at [researchdata@springernature.com](mailto:researchdata@springernature.com)

This service provides advice on research data policies of funders, institutions and journals and on finding research data repositories. It is independent of journal, book and conference proceedings editorial offices and does not advise on specific manuscripts.

**Scope of the service**

**For Editors:**

- Help identify a research data policy type, from Springer Nature's standard data policy framework, for their journal
- Identifying relevant data repositories for their audience, that can be recommended to authors
- Advice on implementation and management of research data policies in journals
- Educational materials and support for better integration of data and literature

**For authors:**

- Information on the research data policy of their target journal(s)
- Help identifying suitable data repositories for their research data
- Information on data reporting standards for different research communities
- Advice on funders' and institutions' data sharing and open data policies\*

## What is a data availability statement (DAS)?

A statement about where data supporting the results reported in the article can be found

- The datasets generated during and/or analysed during the current study are available in the [NAME] repository, [PERSISTENT WEB LINK TO DATASETS].
- The datasets generated during and/or analysed during the current study are available from the corresponding author on reasonable request.
- All data generated or analysed during this study are included in this published article (and its supplementary information files).

Required by many journals/publishers e.g. PLOS, Royal Society, Nature, BioMed Central, BMJ, Hindawi

<http://www.reading.ac.uk/reas-DataAccessStatements.aspx>

# Data availability statements

- Guidance on and published examples of data availability statements

<http://www.springernature.com/gp/group/data-policy/data-availability-statements>

Statement type/description	Template/example text	Published example
Data generated during the study are subject to a data sharing mandate and available in a public repository that does not issue datasets with DOIs	[Data type e.g. "Sequence"] data that support the findings of this study have been deposited in [repository name e.g. "GenBank"] with the [primary] accession codes [list accession codes with links e.g. "KP253039 <a href="#">↗</a> "]	<p><a href="#">BMC Biology <a href="#">↗</a></a></p> <p><a href="#">Nature Communications <a href="#">↗</a></a></p>
Data available in a public (institutional,	The [data type] data that support the findings of this study are available in	<a href="#">Nature Communications <a href="#">↗</a></a>

## Future directions

- A research data policy for every journal
  - Expect to provide a data availability statement to publish in a reputable journal
- More and better links between datasets/repositories and articles
  - Impact/reuse information for research datasets
- More services from publishers to support data sharing
- Integration of larger datasets into online article pages
- Information about data sharing required during submission
  - More peer review visibility for data
- More evidence on costs and benefits of sharing data supporting publications

# Making supplementary data more discoverable

- Discoverability/visibility of articles
- Improve author service
  - Storage of larger datasets and files
  - Better metrics for data files
- Improve reader service
  - Better display/integration of a variety of file types e.g. video, large images
- Improve editor and peer reviewer service
  - Easier access to supporting data files

<https://blogs.biomedcentral.com/bmcblog/2016/12/15/making-research-more-accessible-with-figshare/>



# Figshare & Springer Nature

The screenshot shows the Springer Nature Figshare website. At the top left is the Figshare logo. A search bar is located in the top center with the text "search on figshare". To the right of the search bar are links for "Browse", "Upload", "Sign up", and "Log in". The main banner features a colorful abstract background with a pushpin on the right and the "SPRINGER NATURE" logo in the center. Below the banner is a navigation bar with "NEW", "POPULAR", "CATEGORIES", and "SEARCH" options. The main content area displays a grid of research items, each with a thumbnail, title, author, and "today" status.

Discover research from **Springer Nature** ▾

NEW POPULAR CATEGORIES SEARCH 🔍

**COLLECTION**

Collection: Nanomaterial datasets to advance tomography in scanning...  
Héctor Abruña ▾ today

**DATA SET**

3D reconstruction of hyperbranched Co<sub>2</sub>P nanoparticle u...  
Richard Robinson ▾ today

3D reconstruction of platinum nanoparticles on carbon nanofibre ...  
RobertHovden ▾ today

3D reconstruction of platinum nanoparticles on carbon nanofibre ...  
Chien-Chun Chen ▾ today

**DATA SET**

3D reconstruction of tungsten at atomic resolution using electron to...  
Rui Xu ▾ today

3D reconstruction of Pt-Cu nanocatalysts using electron tomog...  
Deli Wang ▾ today

Electron tomography of hyperbranched Co<sub>2</sub>P nanoparticle  
Richard Robinson ▾ today

Variable dose electron tomography series of platinum nanoparticles on...  
RobertHovden ▾ today

<https://springernature.figshare.com/>

# Figshare “widget” on all BMC & SpringerOpen journals

Globally distributed root endophyte *Phialocephala subalpina* links pathogenic and saprophytic lifestyles  
Showing 1/9: 12864\_2016\_3369\_MOESM8\_ESM.pdf

36 views 0 shares 0 downloads

**ab initio repeat library**

- RepeatScout (l-mer-table, repeatscout)
- Filter step 1 (MS, low complexity)
- blastelust (-S90 -L 0.9 -b F -p F)

**other supporting evidence**

- 454 ESTs *P. subalpina* (TE and non-TE with >4 matches on genome)
- Sanger-sequenced ESTs (TE and non-TE with >4 matches on genome)
- manually annotated TEs

Filter out non-TE sequences (e.g. HET domain proteins)

Use raw library and map against genome

Get consensus for complete TEs and classify them (especially highly abundant TEs)

manually-curated TE library for repeat masking prior to annotation

figshare 1 / 9 < > ≡ 🔍 Share Download

Additional file 1: Mapping statistics for the assembled 454 ESTs. (PDF 197 kb)

Additional file 2: Validation of low identity gene models. (PDF 190 kb)

Additional file 3: Putative secondary metabolite clusters identified in the *P. subalpina* genome (XLSX 12 kb)

Additional file 4: General statistics on the presence of InterPro accessions in the analyzed genomes. (XLSX 10 kb)

Additional file 5: CAZyme modules related to PCWD used for principle component analysis. (XLSX 12 kb)

Additional file 6: General genome statistics for the species included in the present study. (XLSX 15 kb)

Additional file 7: PCA analysis based on InterPro accessions. Placement of the 13 ascomycete species and *P. subalpina* in PCA based on InterPro accessions described in Soanes et al. 2008 to be enriched in pathogens. (PDF 359 kb)

Additional file 8: Strategy to construct the manually-curated repeat library. (PDF 114 kb)

Additional file 9: Overview of the annotation strategy applied. (PDF 120 kb)

<http://bmcbgenomics.biomedcentral.com/articles/10.1186/s12864-016-3369-8>

## References

1. Sieber TN. Fungal root endophytes. In: Waisel Y, Eshel A, Kafkafi U, editors. Plant Roots: The hidden half. 3rd ed.

# The role of journals & publishers: in summary

- Create incentives and increase motivations
- Remove barriers to data sharing
- Support community specific solutions
- Acknowledge community differences
- Implement/support funder and community policy requirements
- Help improve quality and value
- Increase transparency (to increase reproducibility)
- Provide clear, consistent guidance and support
- Make it easier to do the right thing
- Supporting transition to open access



# Thank you

[iain.hrynaszkiewicz@nature.com](mailto:iain.hrynaszkiewicz@nature.com)

[researchdata@springernature.com](mailto:researchdata@springernature.com)

@iainh\_z