# DATA PAPERS AND THEIR APPLICATIONS: DATA DESCRIPTORS IN *SCIENTIFIC DATA*

LCPD Workshop 12th September 2014

Iain Hrynaszkiewicz
Head of Data and HSS Publishing, Open Research
Nature Publishing Group & Palgrave Macmillan

macmillan
Science and Education

*Transforming Learning and Discovery*

http://www.nature.com/sdata/

# *Scientific Data*

## Scope

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.

## Editorial team

Managing Editor (Andrew Hufton)
Editorial Curator (Victoria Newman)
Honorary Academic Editor (Susanna Sansone, Oxford)
Advisory Panel and Editorial Board

## Open access article processing charge

$1,000 USD / £650 GBP / €750 for each accepted article.

# The 'Data Descriptor' article

*Detailed descriptions of the methods and technical analyses supporting the quality of the measurements. Does not contain tests of new scientific hypotheses*

**Sections:**
- Title
- Abstract
- Background & Summary
- Methods
- Technical Validation
- Data Records
- Usage Notes
- Figures & Tables
- References
- Data Citations

**Data Records**

All the samples used in this study are summarized in Table 1. Consistent identifiers are used in Tables 2 and 3 to allow mapping between the proteomic and transcriptomic data outputs.

**Data Record 1**
The raw data, peaklists (.mgf), ProteomeDiscoverer result files (.msf) and ProteomeDiscoverer workflow files (.xml) have been uploaded to ProteomeXchange (http://www.proteomexchange.org/) with the following accession number PXD000134 (ref. 67; Table 2).

**Data Record 2**
Microarray data are available at the NCBI Gene Expression Omnibus (GEO) database under the accession numbers GSE26451 (ref. 68) and GSE26453 (ref. 69; Table 3).

**Data Record 3**
The peptide and protein identification data sets have been annotated by The Global Proteome Machine at http://gpmdb.thegpm.org/

**Data Record 4**
The peptide and protein identification data sets have been annotated by the StemCellOmicsRepository (SCOR) at http://scor.chem.wisc.edu/

**Data Citations**
67. Low, T. Y. *et al.* ProteomeXchange: PXD000134 (2013).
68. Chin, A. *et al.* Gene Expression Omnibus: GSE26451 (2011).
69. Chin, A. *et al.* Gene Expression Omnibus: GSE26453 (2011).
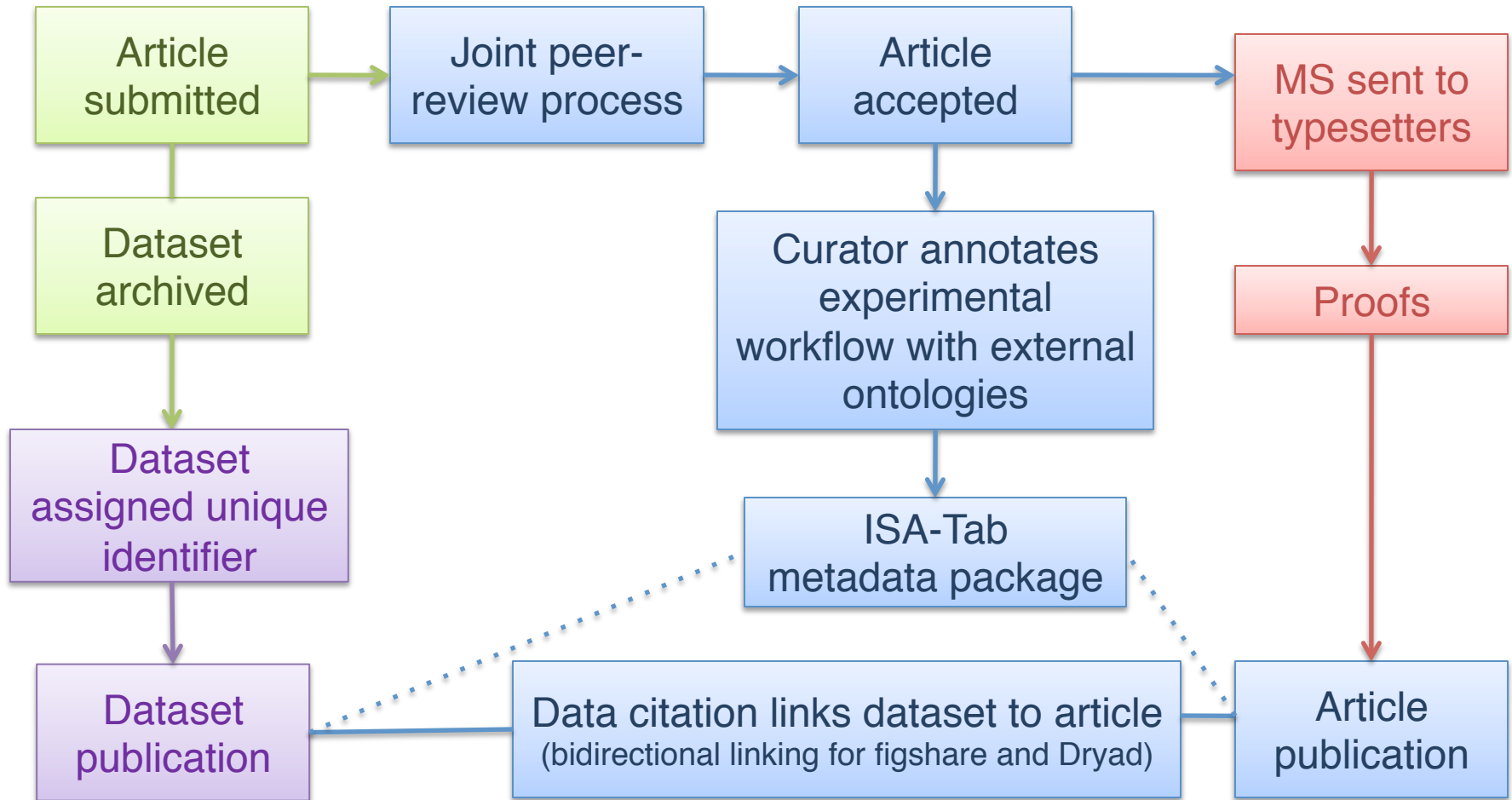
# Peer review at *Scientific Data*

**Focuses on:**

- Completeness (can others reproduce?)

- Consistency (were community standards followed?)

- Integrity (are data in the best repository?)

- Experimental rigour technical quality
  (were the methods sound?)

**Does not focus on:**

- Perceived impact/importance

- Size/complexity of data

# *Scientific Data* workflow overview



Green: author; Purple: repository; Blue: SciData; Red: production

# SCIENTIFIC DATA

## SCIENTIFIC DATA
# Investigation file creation tool *beta*

## Welcome to the Investigation File Creator

This tool will help you to create an Investigation File, a component of the ISA-Tab–based structured metadata included with all manuscripts published in Scientific Data.

### Why am I doing this?

Annotating your data with detailed metadata enables their discovery and reuse, increasing the likelihood that others will build on your research. Creating metadata files to submit with your Data Descriptor manuscript may substantially speed the time to publication should your work be accepted.

### Complete Your Investigation File Below…

**Don't have an ISA-Tab File?**

You can still create an Investigation File from scratch.

**Start**

**Have an ISA-Tab File?**

Import Investigation File (.txt format)

Choose file    No file chosen

**Continue**

### Before You Start

- Completing the narrative portions of your draft Data Descriptor

- Creating tables to describe the samples and assays in your study

- Depositing your data in a public repository

- Collecting details (e.g. Title, DOI/PubMed ID, authors) on any publications related to your data

will all help you create your Investigation File. Please note that template files to describe your samples (Study Files) and assays (Assay Files) can be downloaded from within the Investigation File Creation Tool.

It is not mandatory to complete any sections within the Investigation File or columns within the Study and Assay Files. Partially completed files can be saved and downloaded from the "Thank You" tab in the tool. Downloaded files can also be reopened and edited in the tool. Please direct your questions to scientificdata@nature.com, citing your manuscript tracking number where possible.

Support email address:
scientificdata@nature.com

nature publishing group npg

# The 'Data Descriptor' article



Article or *narrative* component
(PDF and HTML)

Experimental metadata or *structured* component
(in-house curated, machine-readable formats)

Subject terms: Reproductive biology · Data publication and archiving · Biological anthropology

| | |
|---|---|
| **Design Type(s)** | observation design · Demographics · data integration · longitudinal animal study |
| **Measurement Type(s)** | phenotypic profiling |
| **Technology Type(s)** | phenotype characterization |
| **Factor Type(s)** | |
| **Sample Characteristic(s)** | Otolemur garnettii garnettii · Galago moholi · Cheirogaleus medius · Eulemur rubriventer · Eulemur rufus · Eulemur sanfordi · Eulemur · Hapalemur griseus griseus · Lemur catta · Microcebus murinus · Mirza coquereli · Propithecus coquereli · Daubentonia madagascariensis · Varecia · Varecia rubra · Varecia variegata variegata · Eulemur albifrons · Eulemur collaris · Eulemur coronatus · Eulemur fulvus · Eulemur flavifrons · Eulemur macaco · Eulemur mongoz · Loris tardigradus · Nycticebus coucang · Nycticebus pygmaeus · Perodicticus potto · multi-cellular organism |

**Figures at a glance**

Displaying 3 of 4 figures | Figures index



Figure 1   Figure 2   Figure 3

# Stem Cells



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

# Neuroscience



- **New Dataset**
- Data in OpenfMRI
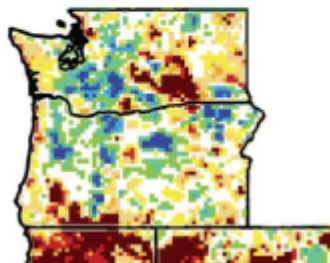- Source code in GitHub
- *Big Data*

**Code in GitHub**

# Environmental



- **New Dataset**
- Data in **fig**share
- Code in **fig**share
- Integrated **fig**share data viewer
- Cited in *Science*

# The right licence

**Data Descriptor article:** Licensed under one of two Creative Commons licenses, by author choice:



**Metadata:** released under the **CC0 waiver** to maximize reuse and aid data miners

**Data:** depends on public repositories. Partner repositories figshare and Dryad both use the **CC0 waiver**.

# Repository criteria

http://www.nature.com/sdata/data-policies/repositories

1. Broad support and recognition within their scientific community

2. Ensure long-term persistence and preservation of datasets

3. Provide expert curation

4. Implement relevant, community-endorsed reporting requirements

5. Provide for confidential review of submitted datasets

6. Provide stable identifiers for submitted datasets

7. Allow public access to data without unnecessary restrictions
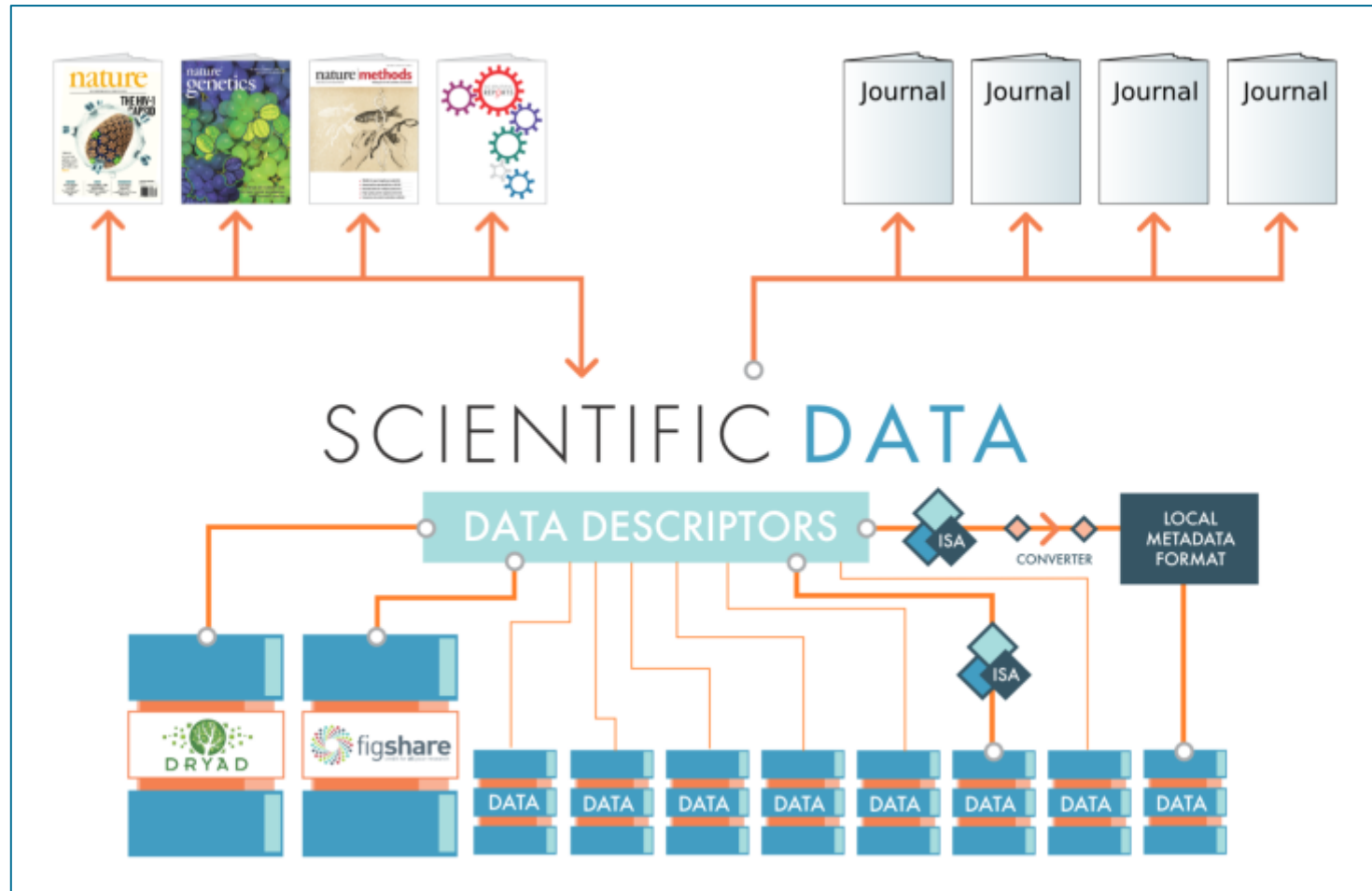
# Thank you

For more information
please contact

**IAIN HRYNASZKIEWICZ**
**Head of Data and HSS Publishing,**
**Open Research**

M: +44 (0)7814 290576
T: +44 (0)207 0146753
E: iain.hrynaszkiewicz@nature.com

macmillan
Science and Education

*Transforming Learning and Discovery*

SCIENTIFIC DATA

## Get Credit for Sharing Your Data

Publications will be indexed and citeable.

## Open-access

Creative Commons licenses (CC-BY/CC-BY-NC) for the main Data Descriptor. Each publication supported by CCO metadata.

## Focused on Data Reuse

All the information others need to reuse the data; no interpretative analysis, or hypothesis testing

## Peer-reviewed

Rigorous peer-review focused on technical data quality and reuse value

## Promoting Community Data Repositories

Not a new data repository; data stored in community data repositories