


# ***Scientific Data***

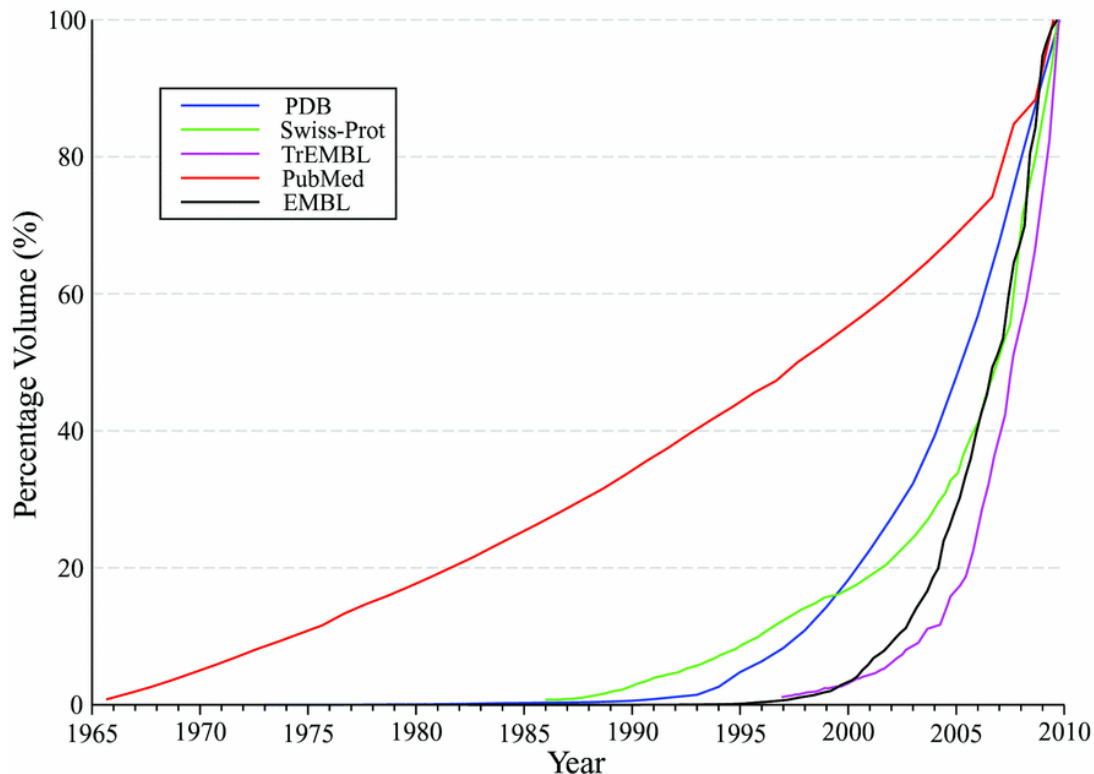
**An open-access, peer-reviewed platform for data-focused publications**



Andrew Hufton  
Managing Editor, Data  
Nature Publishing Group

# Data, data, data

Depositions of datasets in archives continue to grow, surpassing journal articles in biomedical research



Growth of biomedical research publications (red; current total >19 million), alongside the accumulation of research data, including nucleic acid sequences (black; current total ~163 million), computer-annotated protein sequences (magenta; current total 9 million), manually annotated protein sequences (green; current total 500,000) and protein structures (blue; current total 60,000)

Source: Biochemical Journal 2009 424, 317-333 - Teresa K. Attwood, Douglas B. Kell and others.

# An illustrative example....



## DNA Structure 1953

1 Page  
2 Authors  
1 Figure  
no data

Year	Number of people
2000	100
2001	102
2002	104
2003	106
2004	108
2005	110
2006	112
2007	114
2008	116
2009	118
2010	120

**62 Pages, 150 Authors,  
49 Figure, 27 tables**



# Existing challenges

- Data producers do not necessarily get appropriate credit for their work
- The peer review process at research journals is not ideal for ensuring data release and data standards
- Data and info about datasets often ends in supp. material
- Valuable datasets are not released



# A Role for Data Publications

*GigaScience*

*F1000R*

*ZooKeys*

*Ecological Archives*

**and soon *Scientific Data*...**

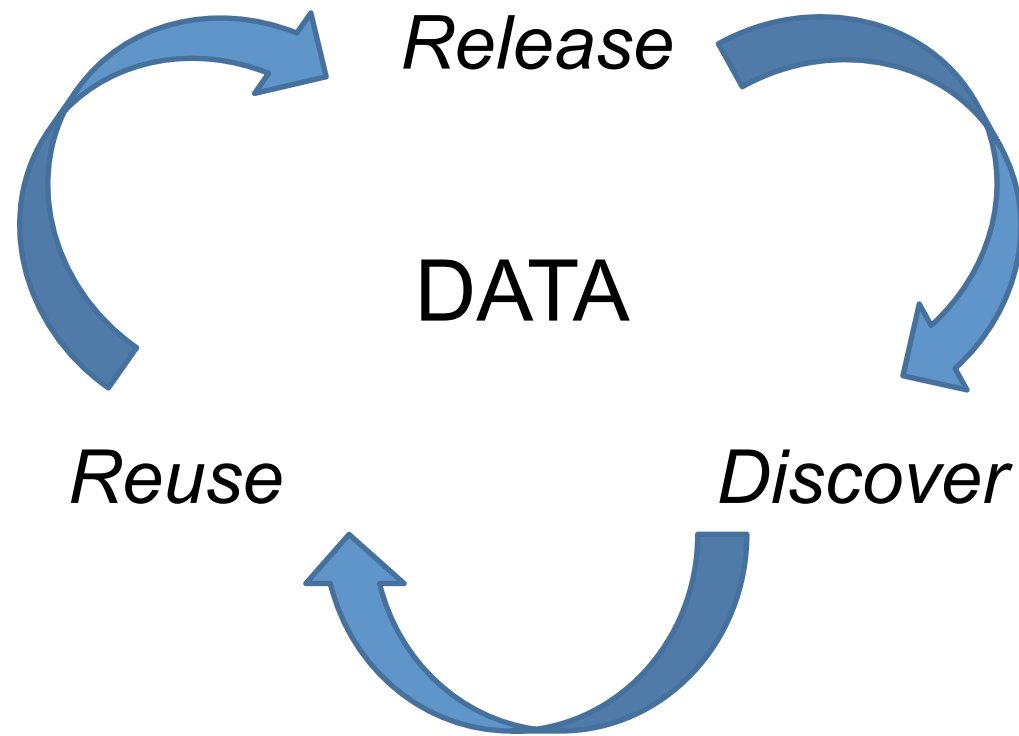




Year	Number of people in the workforce
2000	100
2001	98
2002	100
2003	98
2004	100
2005	102
2006	104
2007	106
2008	107
2009	108
2010	107

- 
- | Year | Number of people in the workforce |
|------|-----------------------------------|
| 2000 | 100                               |
| 2001 | 98                                |
| 2002 | 100                               |
| 2003 | 98                                |
| 2004 | 100                               |
| 2005 | 102                               |
| 2006 | 104                               |
| 2007 | 106                               |
| 2008 | 107                               |
| 2009 | 108                               |
| 2010 | 107                               |







**If we want scientists to **release** their data, we need to provide a credit mechanism**



If we want scientists to **release** their data, we need to provide a credit mechanism

## PUBLICATION



**If want released data to be reusable,  
we need critical evaluation to verify  
experimental rigor and the  
completeness of the description**



**If want released data to be reusable,  
we need critical evaluation to verify  
experimental rigor and the  
completeness of the description**

## PEER REVIEW



If we want scientists to be able to find datasets that will accelerate their research, then datasets need to be searchable & discoverable



If we want scientists to be able to find datasets that will accelerate their research, then datasets need to be searchable & discoverable

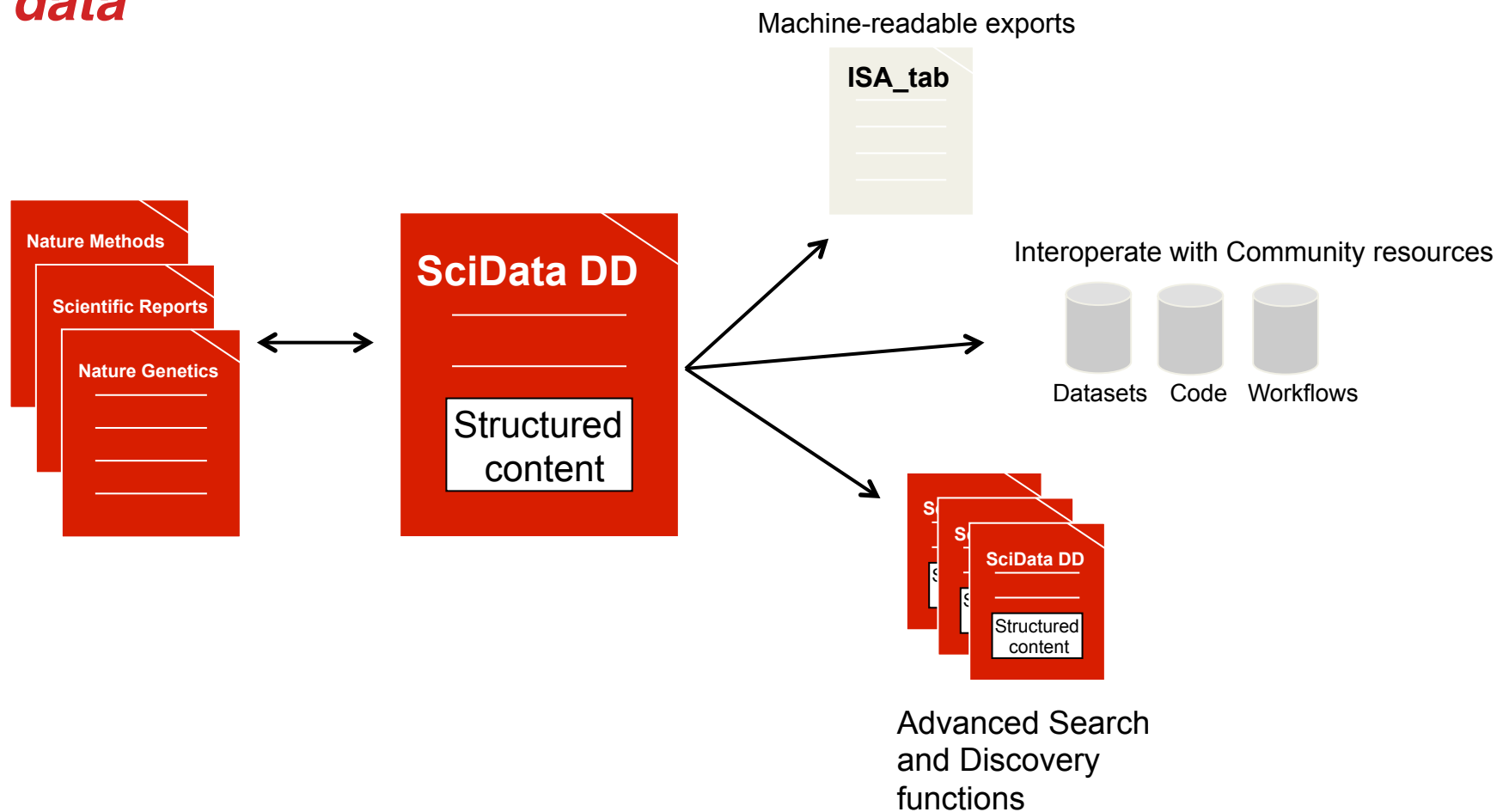
STANDARDIZED,  
CURATED DESCRIPTIONS



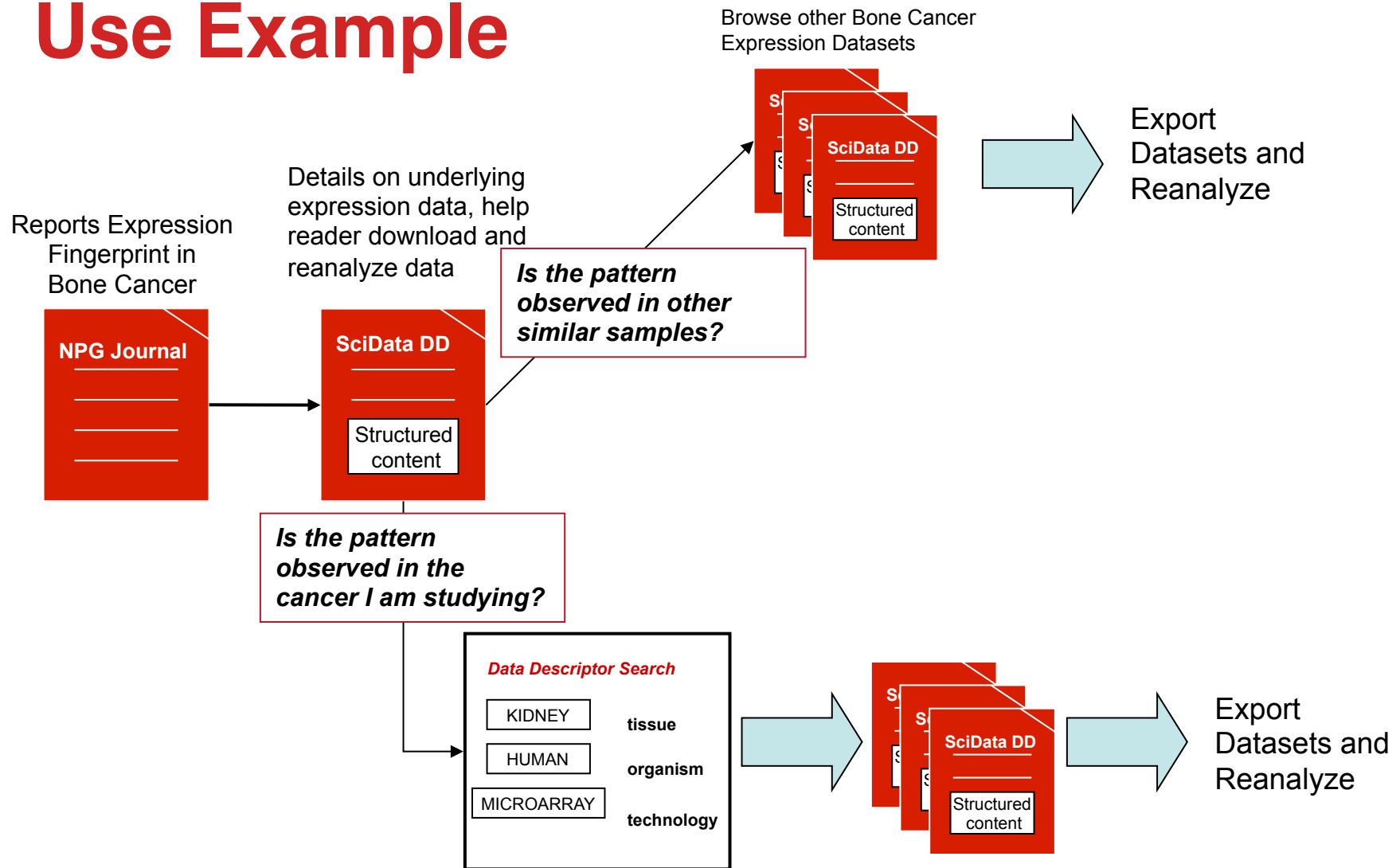


# The Data Publication

*seamless integration between research findings and data*



# Use Example



# Scientific Data Concept Overview

**Credit:** citeable, peer-reviewed mechanism for recording credit for dataset creation.

**Rapid peer-review and publication** evaluate the technical quality of the procedures used to generate the data, the reuse value of the resulting datasets, the completeness of the data description, and alignment of with existing community standards.

**Reuse:** complete and standardized descriptions help others discover and reuse your data, and discover your related research articles.

**Diverse data-types:** *Scientific Data* will initially focus on experimental and observational datasets from the life, biomedical and environmental science communities, but will be open to content from a wide range of experimental disciplines.

**Cooperation** with a broad set of public, community-recognized repositories will promote existing standards, and **integration** with generalist scientific repositories like Dryad and FigShare will make it easy to deposit diverse data-types.

**Open access:** Data Descriptors will be released freely to the public under an open access license.

**Complementary:** Data Descriptors can be used to describe both datasets that are associated with a traditional research article, as well as standalone datasets.

# Focused Scope

complements both journal articles and repository records

## Includes

- Highly detailed, reproducible methods descriptions
- Quality control & technical validation experiments
- Searchable, machine-readable meta-data

## Does Not Include

- In depth analysis or tests of hypotheses
- New scientific conclusions
- Exploratory analysis (e.g. clustering)

# Thanks!

Andrew Hufton [andrew.hufton@nature.com](mailto:andrew.hufton@nature.com)  
Managing Editor, Data

