# Digital Preservation

# Data Citation

Stefan Pröll

# Outline

- Introduction and Motivation
  - Why should we reference?
- Persistent Identifiers
  - Isn't a URL enough?
- Citing Datasets
  - Best Practices
- Future Research
  - Dynamic Datasets

# Why Should We Cite?

- Science is a collaborative approach

"If I have seen further, it has been by standing on the shoulders of giants."

??

- Giving credit to peers and acknowledge their work

# Benefits of Citation

- Enables reproducibility

- Enhances transparency

- Serves as documentation

- Settles the context

- Allows exact identification of results

- Track the impact of research

# Citations Help Detecting Scientific Misconduct

RESEARCH ARTICLE

## How Many Scientists Fabricate and Falsify Research? A Systematic Review and Meta-Analysis of Survey Data
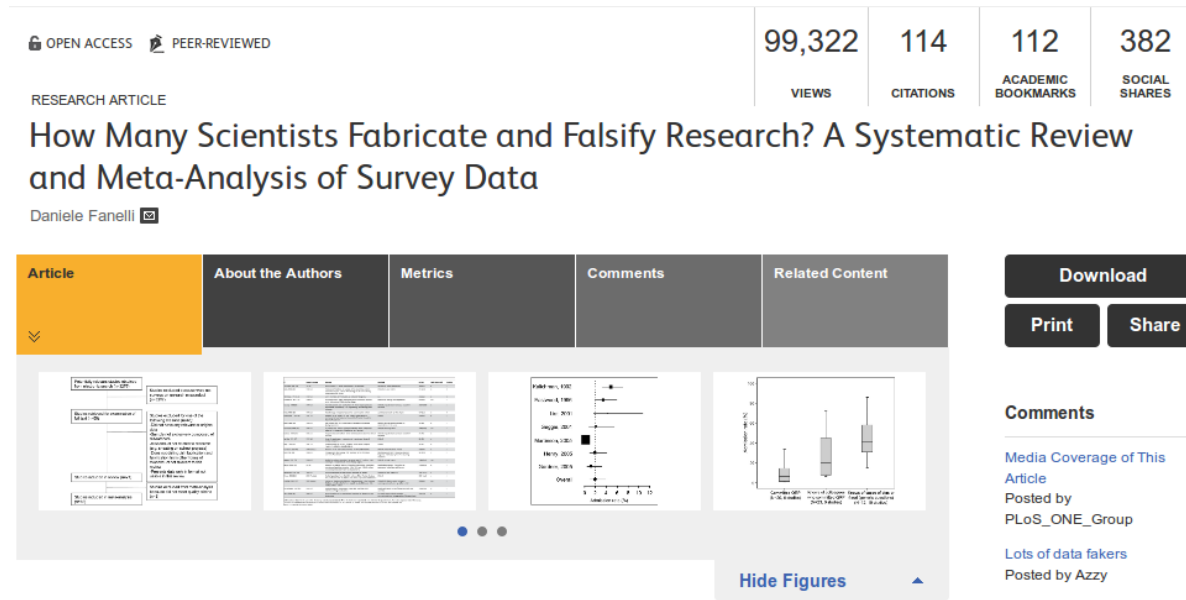
Daniele Fanelli

| Article | About the Authors | Metrics | Comments | Related Content |

Download    Print    Share

Comments

Media Coverage of This Article
Posted by PLoS_ONE_Group

Lots of data fakers
Posted by Azzy

Interestingly enough...
Posted by SQLserver

ADVERTISEMENT

- Abstract
- Introduction
- Methods
- Results
- Discussion
- Supporting Information
- Acknowledgments
- Author Contributions
- References

Reader Comments (4)

Hide Figures

## Abstract

The frequency with which scientists fabricate and falsify data, or commit other forms of scientific misconduct is a matter of controversy. Many surveys have asked scientists directly whether they have committed or know of a colleague who committed research misconduct, but their results appeared difficult to compare and synthesize. This is the first meta-analysis of these surveys.

To standardize outcomes, the number of respondents who recalled at least one incident of misconduct was calculated for each question, and the analysis was limited to behaviours that distort scientific knowledge: fabrication, falsification, "cooking" of data, etc... Survey questions on plagiarism and other forms of professional misconduct were excluded. The final sample consisted of 21 surveys that were included in the systematic review, and 18 in the meta-analysis.

Source: http://www.plosone.org

# Advance in Research by Sharing

# Paper Publications

- Referencing research papers is well established

# Data Publications

- Data is an essential part of research
  - Majority of papers is based upon research data
  - Needed for validation and reproduction of experiments
- Challenges
  - Encourage researchers to share
  - Different data formats
  - Potentially large storage size
  - Who maintains it?

# Stakeholders

# Data Center Policies
## Criteria for Assessing Value of Data

- Relevance to mission

- Scientific value

- Uniqueness

- Potential for redistribution

- Non-Replicability

- Costs

- Documentation

- ….

http://www.dcc.ac.uk/resources/how-guides/appraise-select-data

# Data Citation Requirements

- Unique identification
- Identify subsets and complete dataset
- Machine readable metadata
- Human readable metadata
- Citation metrics

# Elements of Data Citation

- Classical bibliographic details:
  - Author, date, edition
  - Publisher, version

- Specific details:
  - Feature name, resource type
  - Unique numeric fingerprint (hash)
  - <u>Persistent identifier</u>
  - <u>Location</u>

# Digital Object Life Cycle



Archive

Create

Revised

Edit

Destroy

Cited

Published

# Identifiers

- Identifier is a <u>symbol</u> that uniquely <u>identifies</u> a <u>digital object</u>. Can be dependent on the context.

- Are URLs identifiers?

- Are URLs persistent?
    - Can URLs be mapped to the digital object life cycle?

Sources: http://ands.org.au/guides/persistent-identifiers-working.html

# Existing Unique Identifier Models

- **Traditional Mechanisms**
  - International Standard Serial Number (ISSN)
    - Unique eight-digit number
    - Identifiers periodical publications
    - Can be encoded as URN
  - International Standard Book Number (ISBN)
    - Unique commercial book identifier barcode
    - 13 (since 2007) or 10 digits with checksum
    - ISBN-10: 3836217155
    - ISBN-13: 978-3836217156

secure
sba-research.org

# Unique Identifiers for Digital Objects

- Locator Based Mechanisms
  - Uniform Resource Identifier (URI)
  - Uniform Resource Locator (URL)
  - Uniform Resouce Name (URN)
  - National Bibliographic Numbers (NBNs)
- Delegating Methods
  - DOI
  - The Handle System
  - Digital Object Identifier
  - Persistent URL (PURL)
  - Archival Resource Key (ARK)

# URL

- Uniform Resource Locator
- Particular form of URIs
- Addressing documents
- Can only be used to locate resources
- Depend on DNS information
- URLs are not persistent
- <u>Problems</u>
  - Links may break
  - No mechanism to handle broken links

# URN

- Uniform Resource Name
- Combination of namespace identifier (NID) and a namespace specific string (NSS)
- Naming scheme for URNs:
- urn: <NID> :<NSS>
- Example: urn:isbn:0451450523

# URN (2)

- Main functions of a URN
  - Global scope of names
  - Global uniqueness
  - Persistence
  - Scalability
  - Legacy support
  - Extensibility
  - Independence
  - Resolution
- Persistence of identifiers
- Possibility to resolve them

# Persistency?

- Standard URLs are not forever
  - Network locations
  - Not suitable for the long term
  - Link rot: half of the links in publications are not available after 5 to 7 years
- Persistent identifiers need to be maintained
  - Proactively
  - Throughout the whole lifecycle
- Can a URL be persistent?

# Digital Object Identifier (DOI)

- Identifier scheme administered by the International DOI Foundation

- Consists of three parts:

http://dx.doi.org/10.1016/S0169-7552(98)00110-X

| Resolver Service | Prefix (Assigning Body) | Suffix (Resource) |

# DOI

- Publisher (organizations) register and get a unique ID (Prefix)

- Resource gets an ID (Suffix) which is unique within the prefix

- Resolver services maintain the link between the endpoint (e.g. URL) and the resource
  - http://dx.doi.org/

# Handle

- Distributed persistent naming system
- Conforms to URN framework
- Used by <u>DOI</u> (Digital Object Identifier) system
- Persistent identifier consists of two parts:
  - Naming authority
  - Name (must be unique string to the authority)
- Digital objects on the Internet can be assigned, managed and resolved by handles.
- Resolved by global handle service

# Handle (2)

- Main points
  - Handles are unique and persistent
  - Handle system supports internationalization
  - Operations on handle system have to be authorized
- Syntax:
  <Handle Naming Authority> ,/' <Handle Local Name>
- Example:
  - 10.1045/january2013-burns
- Available Services:
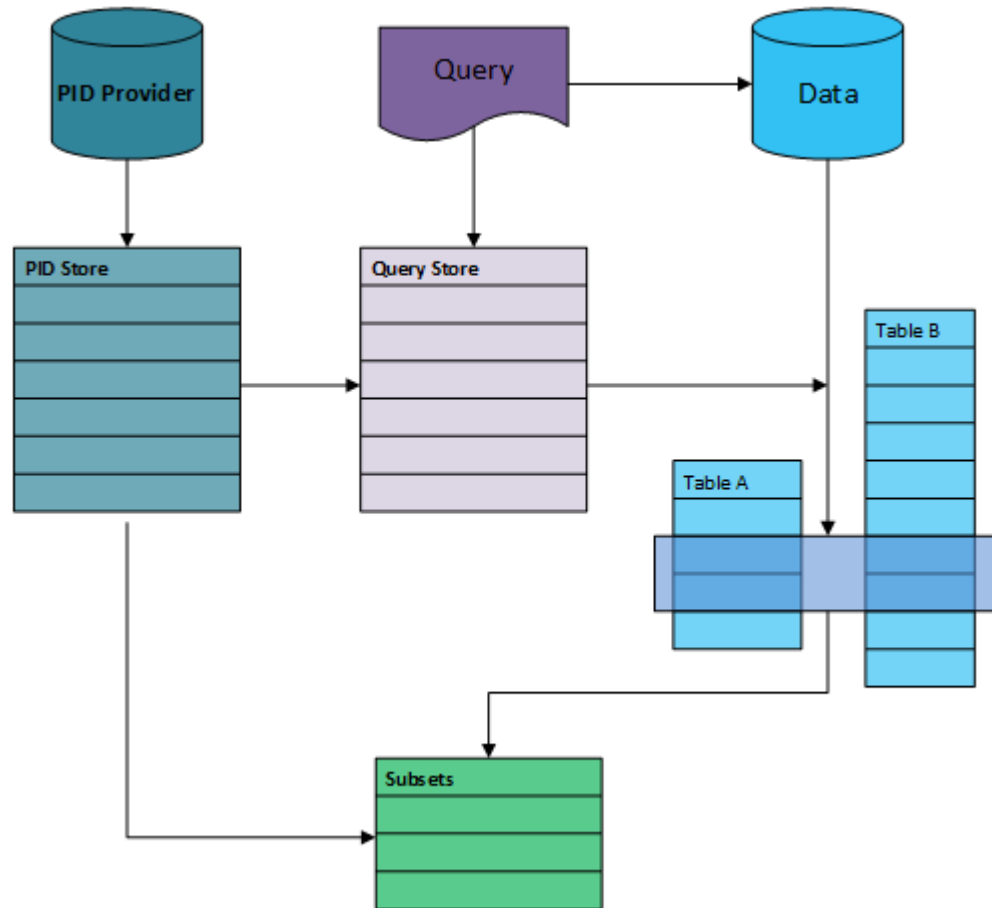  - http://hdl.handle.net

# Current Challenges

- Granularity
  - How to define subsets?
  - Should individual data records be cited?
  - Assign each database row a DOI?
- Identify people (authors)?
- Dynamic data
  - How to treat evolving data?

# Referencing Dynamic Subsets

- What is needed:
  - Uniquely identifiable data records
  - Time stamps of data
  - Versioned data, considering markings of deleted, altered or inserted data records
  - Precise query language for constructing subsets
  - Persistent query store that keeps queries and the timestamp of their issuing
  - An identification mechanism for queries, that enables access

secure
sba-research.org

# Citing Dynamic Data in Databases

# Referencing Dynamic Subsets

- High Level Requirements
  - Dynamic data
    - Queries need to be stored
    - Temporal data and queries
  - Assemble subsets
  - Scalability is enabled
  - Implementation is transparent
  - Machine actionable

# Literature and Links

- http://www.dcc.ac.uk/resources/how-guides/appraise-select-data
- http://www.dcc.ac.uk/resources/how-guides/cite-datasets
- http://dl.acm.org/citation.cfm?doid=602421.602422
- http://ands.org.au/guides/persistent-identifiers-working.html
- http://hdl.handle.net/
- http://dx.doi.org/

# Thank you for your attention.

sproell@sba-research.org