# Connecting data repositories and publishers for data publication

Sarah Callaghan, Fiona Murphy, Jonathan Tedds, Rob Allan, John Kunze, Rebecca Lawrence, Matthew S. Mayernik, Angus Whyte, and the PREPARDE project team

#preparde

sarah.callaghan@stfc.ac.uk @sorcha_ni

JISC · University of Leicester · British Atmospheric Data Centre — NATIONAL CENTRE FOR ATMOSPHERIC SCIENCE NATURAL ENVIRONMENT RESEARCH COUNCIL · WILEY-BLACKWELL · University of Reading · D|C|C · University of California CDL California Digital Library · F1000 FACULTY of 1000 POST-PUBLICATION PEER REVIEW · NCAR NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

# Why link data and publications?

- Data is the foundation of science – without it we can't test our assertions or reproduce our results

- The Internet allows us to link things to other things quickly and easily

- But there are still serious problems to address when it comes to linking data to the scientific record:
  - Data persistence
  - Data and metadata quality
  - Attribution and credit for data producers
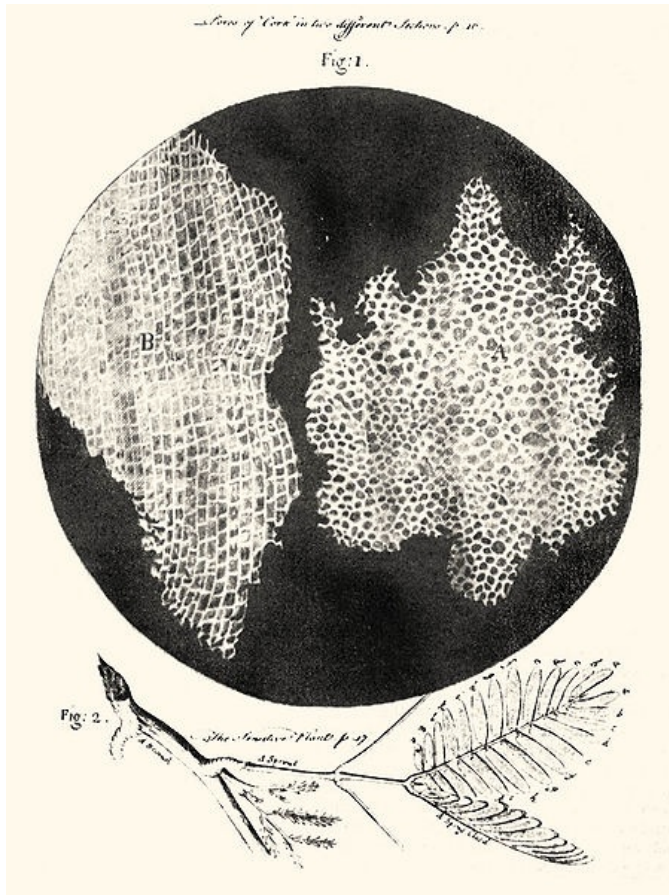  - … and many more



Engraving of printer using the early Gutenberg letter press during the 15th century.
Date            unknown - estimate 16th - 19th century
http://commons.wikimedia.org/wiki/File:Gutenberg_press.jpg

# Historically, journals have always published data



Suber cells and mimosa leaves. Robert Hooke, Micrographia, 1665



The Scientific Papers of William Parsons, Third Earl of Rosse 1800-1867

# But now… the Data Deluge

*"the amount of data generated worldwide…is growing by 58% per year; in 2010 the world generated 1250 billion gigabytes of data"*

The Digital Universe Decade – Are You Ready?
IDCC White Paper, May 2010

A lot of people are creating a lot of data, and we're only going to get more of it.

If this is a data deluge – time to start building boats!

Figure 1: The Digital Universe 2009 – 2020
*Growing by a Factor of 44*

2009
0.8 ZB*

2020
35 ZB*

*Zettabyte = 1 trillion gigabytes

Source: IDC Digital Universe Study, sponsored by EMC, May 2010

# Using citations to link research outputs

- We already have a working method for linking between publications which is
  - commonly used
  - understood by the research community
  - used to create metrics to show how much of an impact something has (citation counts)
  - applied to digital objects (digital versions of journal articles)

- We can extend citation to other things like
  - data
  - code
  - multimedia

 And the best bit is, we don't need to teach researchers a new method of linking – they cite like they normally would!

http://www.flickr.com/photos/anton41/6588935181/

# Reasons for citing and publishing data

- Pressure from (UK) government to make data from publicly funded research available for free.
  - Scientists want attribution and credit for their work
  - Public want to know what the scientists are doing

- Research funders want reassurance that they're getting value for money
  - Relies on peer-review of science publications (well established) and data (not done yet!)

- Allows the wider research community to find and use datasets, and understand the quality of the data

- Extra incentive for scientists to submit their data to data centres in appropriate formats and with full metadata

http://www.evidencebased-management.com/blog/2011/11/04/new-evidence-on-big-bonuses/

# PREPARDE: Peer REview for Publication & Accreditation of Research Data in the Earth sciences

- Lead Institution: University of Leicester
- Partners
  - British Atmospheric Data Centre (BADC)
  - US National Centre for Atmospheric Research (NCAR)
  - California Digital Library (CDL)
  - Digital Curation Centre (DCC)
  - University of Reading
  - Wiley-Blackwell
  - Faculty of 1000 Ltd
- Project Lead:          Dr Jonathan Tedds  (University of Leicester, jat26@le.ac.uk)
- Project Manager:    Dr Sarah Callaghan  (BADC, sarah.callaghan@stfc.ac.uk )
- Length of Project:   12 months
- Project Start Date:  1st July 2012
- Project End Date:    31st June 2013

# *Geoscience Data Journal*, Wiley-Blackwell and the Royal Meteorological Society

- Partnership formed between Royal Meteorological Society and academic publishers Wiley Blackwell to develop a mechanism for the formal publication of data in the Open Access *Geoscience Data Journal*

- GDJ publishes short data articles cross-linked to, and citing, datasets that have been deposited in approved data centres and awarded DOIs (or other permanent identifier).

- A data article describes a dataset, giving details of its collection, processing, software, file formats, etc., without the requirement of novel analyses or ground breaking conclusions.
  - the when, how and why data was collected and what the data-product is.

Volume 1, Number 1 — May 2012
RMetS — Geoscience Data Journal — Open Access
Editor-in-Chief Dr. Rob Allan
WILEY Open Access

RMetS
Royal Meteorological Society

WILEY-BLACKWELL

**The traditional online journal model**

1) Author prepares the paper using word processing software.

Word processing software with journal template

2) Author submits the paper as a PDF/Word file.

A Journal (Any online journal system)

PDF | PDF | PDF | PDF | PDF

Data ?

3) Reviewer reviews the PDF file against the journal's acceptance criteria.

**Overlay journal model for publishing data**

1) Author prepares the data paper using word processing software and the dataset using appropriate tools.

Word processing software with journal template

2a) Author submits the data paper to the journal.

2b) Author submits the dataset to a repository.

Data Journal (Geoscience Data Journal)

html | html | html | html

Data | Data
BADC

Data | Data
BODC

3) Reviewer reviews the data paper and the dataset it points to against the journals acceptance criteria.

RMetS — Royal Meteorological Society

## Geoscience Data Journal

Open Access

Data Paper

### On the South Atlantic Convergence Zone affecting southern Amazonia in austral summer

Fabien C. Lamaze[1,*], Dany Garant[2], Louis Bernatchez[1]

Article first published online: 23 OCT 2012

DOI: 10.1002/asl.401

Total views since publication: 25   view chart

Additional Information (Show All)

How to Cite | Author Information | Publication History | Funding Information

Issue

Evolutionary Applications

Early View (Online Version of Record published before inclusion in an issue)

Abstract | Article | References | Supporting Information | Cited By

Get PDF (504K)

**Keywords:**
brook charr; gene expression; hybridization; introgression; quantitative PCR; stocking

### Dataset                                           Jump to…

GBS (Global Broadcast Service), doi:10.1029/2007RS003793

### Abstract                                          Jump to…

Time series analysis of the average rainfall over a target area in southern Amazon Basin showed a spectral peak at 11 day period. An objective method for defining the South Atlantic Convergence Zone (SACZ) is used to identify 28 episodes affecting southern Amazon Basin during the 10 summers in the period 1999–2010. The 28-episode composite precipitation anomalies show significant positive values over the target area. The convergence of moisture over the target area in the SACZ composites is about 35% stronger than the climatological value. Copyright © 2012 Royal Meteorological Society

### 1. Introduction                                   Jump to…

One of the regional scale meteorological systems that affect the weather over a major part of the South American tropics is the South Atlantic Convergence Zone (SACZ). This system somewhat plays the same role for the South American monsoon (Vera et al., 2006; Carvalho et al., 2010) as does the monsoon trough for the South Asian summer monsoon (Keshavamurty and Awade, 1970). The cloud band associated with SACZ extends from the Amazon Basin (Amazonia) to the South Atlantic subtropics, over a stretch of 4000 km or more (Kodama, 1992; Satyamurty et al., 1998) and affects many regions of Brazil with intense rainfall. Some SACZ events are especially intense over interior South American continent (Carvalho et al., 2002, 2004; Muza and Carvalho, 2006), strongly affecting Amazonia.

# Data paper mock-up

Dataset citation is first thing in the paper and is also included in reference list (to take advantage of citation count systems)

JISC · University of Leicester · British Atmospheric Data Centre — NATIONAL CENTRE FOR ATMOSPHERIC SCIENCE NATURAL ENVIRONMENT RESEARCH COUNCIL · WILEY-BLACKWELL · University of Reading · DCC · University of California CDL California Digital Library · F1000 FACULTY of 1000 POST-PUBLICATION PEER REVIEW · NCAR NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

# PREPARDE topics

Example steps/workflow required for a researcher to publish a data paper

3 main areas of interest (in orange)
1. Workflows and cross-linking between journal and repository
2. Repository accreditation
3. Scientific peer-review of data

- Division of area of responsibilities between
  - *repository controlled* processes
  - *journal controlled* processes

# Data repository workflows

- Data centre and journal workflows captured
  - Workflows are very varied! No one-size fits all method
  - Can have multiple workflows in the same data centre, depending on interactions with external sources ("Engaged submitter"/ "Data dumper" / "Third party requester")

# Repository Workflow – NCAR Comp. & Info. Systems Lab Research Data Archive (RDA)

**Data Ingest**

**Check with data provider for changes to files**

**Notification to provider/user community**

**Data Preparation:**
- Automated file collection.
- Check integrity of file receipts.
- Compare bytes and checksums (if available) with original data providers.

**Processing:**
- Validate files – using software, read the full content of every file.
- Pull out metadata.
- Identify errors and metadata holes.
- Do time-series checks.
- Check metadata against internal standard/expectation.
- If necessary, filter data or fix metadata.

**Embargo**

**Archive (Tape-based)**

**Access Development Phase**

**Online Data (Most Demanded)**

**Publish Metadata – User GUIs**

**Distribute metadata**

Not ok

Ok

Contact data provider

Errors found

**Metadata Database**
- Spatial info
- Temporal info
- Global Change Master Directory (GCMD) keywords
- Parameters
- Format table relationships

**Remote backup**

GCMD

NCAR CDP

BADC

... OAI-PMH

# Journal workflow

**Geoscience Data Journal**
**Data Paper workflow**

Data set in repository digested and suitable for DOI (passed technical review at data centre)

Confirmation and DOI sent to author[1]

Author writes Data Paper about the data set, including DOI and submits to GDJ editorial office[2]

Data Paper is reviewed, assessing the following criteria:

1) does it meet the journal's editorial guidelines? Eg data set has a DOI, paper is in scope for the journal

If 'no' then Data Paper is rejected[3]

2) scientific review of data set. Eg is it accurate in its methods of data acquisition, statistical info and error calculations etc? Is the data scientifically useful?

If data set does not pass scientific review then Data Paper is rejected. Author should correct/add to the data set and resubmit it to the data centre as a new version. Once the new version of the data set has been ingested the author could submit a new Data Paper (indicating the new data set version number in bibliographic details)[3]

3) review of Data Paper. Eg does the paper adequately describe the data set? Does it explain the data acquisition methodology etc?

Revision required to Data Paper (ie data set ok but not sufficiently described in the paper) → author revises Data Paper and then sends revised version for further review

Data Paper passes review and is accepted for publication in GDJ

Data Paper goes through production at publisher (undergoes copyediting and typesetting, checking of reference details, formatting to journal style, necessary coding and tagging added to allow cross linking, citation and discoverability)

XML info (provided by data centre, via author) is used to populate the data set tagging within the article → this appears both fully tagged in the 'data set' section at the start of the article and as a normal reference within the reference list

Data Paper published online in Wiley Online Library (first in Early View, then later within an issue). The Data Paper is assigned its own DOI

Citation of the data set by the Data Paper is registered in ISI and other indexing services

Data Paper details are sent to the relevant data centre for them to add a cross link from the data set to the Data Paper[4]

Data Paper (and its authors) can accrue citations by other articles citing the data set[5]
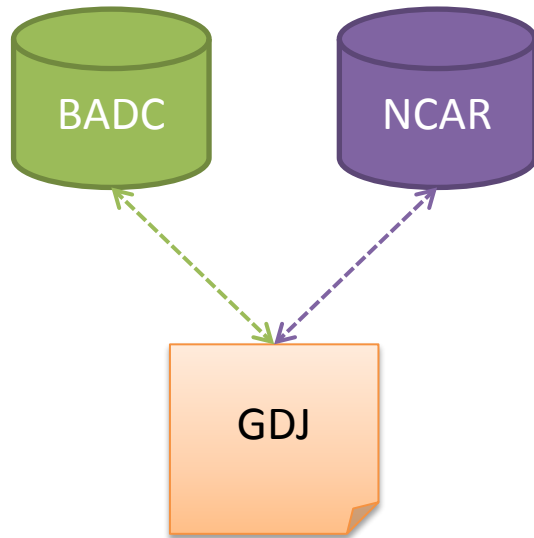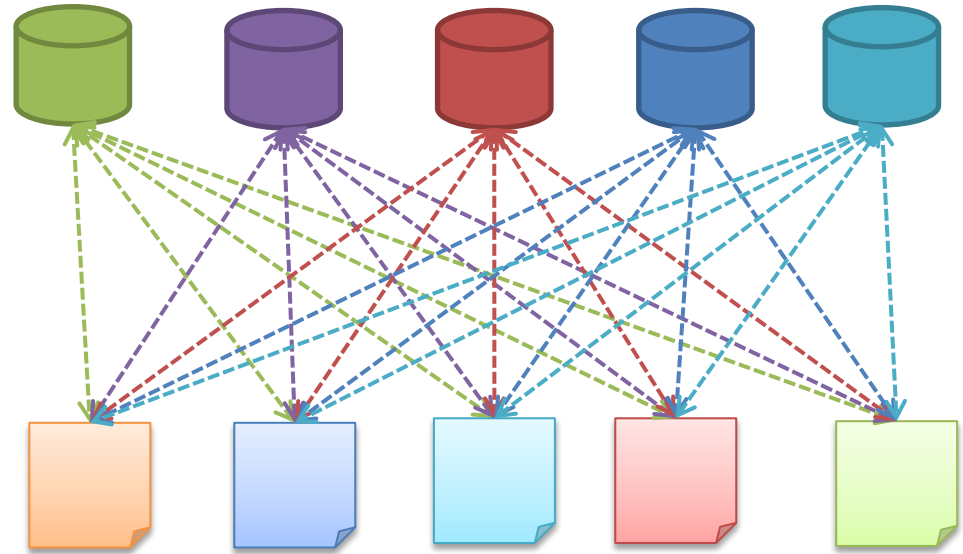
- Work on comparisons and identification of cross-linking points is continuing.
  - Aim is to minimise effort needed to submit data paper by taking advantage of already submitted metadata.
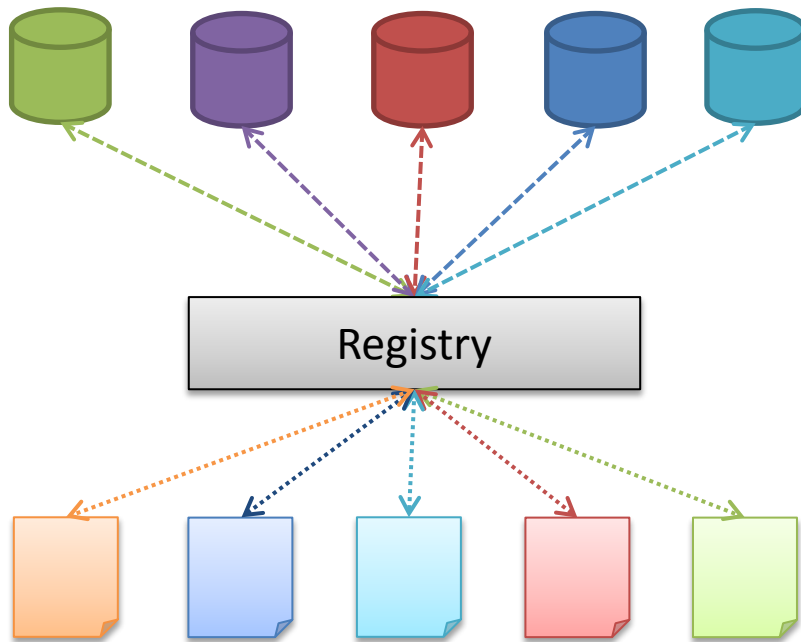
# Cross-linking



This is what we have to focus on for PREPARDE – demonstrate cross linking between GDJ and BADC (and maybe NCAR)

Unfortunately this direct cross-linking isn't scaleable!

Need for off-the shelf solutions that can work across multiple research domains

# Cross-linking – the ideal situation

Registry could provide other functions as well as being an intermediary between journals and data repositories like:

- Certify data centres are "trustworthy"
- Administer linking mechanism
- Provide search and metrics functions

Disadvantages:

- Single point of failure
- Difficulty of standardisation across different research domains

Could OpenAIRE be this registry?
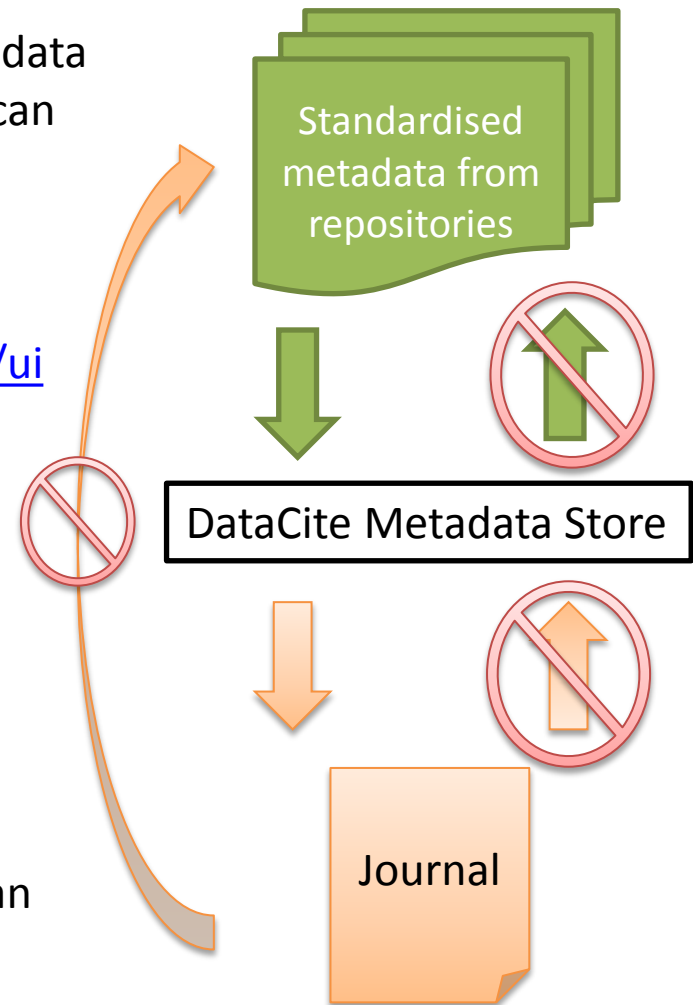
# Do we have a start?

DataCite have standardised a set of bibliometric metadata that have to be submitted before a DOI for a dataset can be minted by a repository.

This metadata is then made openly available via the DataCite metadata search: http://search.datacite.org/ui

Given a DOI, a journal can then easily find the DOI standard metadata.

DataCite also have a content resolver http://data.datacite.org/static/index.html

What's missing is the return link, where the journal can let the repository know that a dataset has been cited (directly or via DataCite)
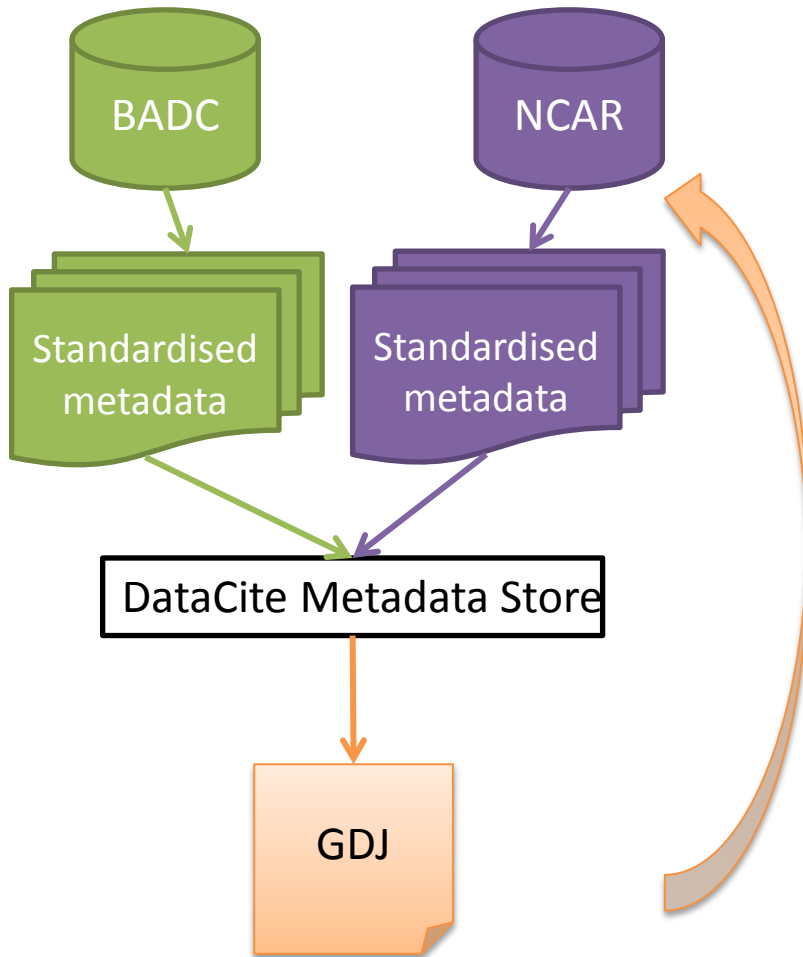
Standardised metadata from repositories

DataCite Metadata Store

Journal

# DataCite Metadata Schema

http://schema.datacite.org/

| ID | Property |
|----|----------|
| DataCite Mandatory Properties | |
| 1 | Identifier (with type attribute) |
| 2 | Creator (with name identifier attributes) |
| 3 | Title (with optional type attribute) |
| 4 | Publisher |
| 5 | PublicationYear |

| ID | Property |
|----|----------|
| DataCite Optional Properties | |
| 6 | Subject (with schema attribute) |
| 7 | Contributor (with type and name identifier attributes) |
| 8 | Date (with type attribute) |
| 9 | Language |
| 10 | ResourceType (with description attribute) |
| 11 | AlternateIdentifier (with type attribute) |
| 12 | RelatedIdentifier (with type and relation type attributes) |
| 13 | Size |
| 14 | Format |
| 15 | Version |
| 16 | Rights |
| 17 | Description (with type attribute) |

JISC · University of Leicester · British Atmospheric Data Centre · NATIONAL CENTRE FOR ATMOSPHERIC / NATURAL ENVIRONMENT RESEARCH · F1000 POST-PUBLICATION PEER REVIEW · NCAR NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

# MOLES: Metadata Objects for Linking Environmental Sciences v3.4



http://proj.badc.rl.ac.uk/moles/browser/branches/V3.4/MODEL/Diagrams/MOLES3.4Summary.png

# What PREPARDE is going to do



We already have a link from the GDJ data article to the data repository – thanks to the DOI.

GDJ can also pull the standard DOI metadata attached to that DOI from the DataCite metadata store

We need to figure out a way so GDJ can inform the repository that their dataset has been cited/published – bearing in mind scaling issues!

Might have to start with a manual work-around.

# Tell us what you think

Workshop on cross-linking between data centres and publishers planned for May 2013 at Rutherford Appleton Laboratory, UK

Workshop on peer-review of data planned for March 2013 at the British Library

Always happy to get input from others!



Image Credit: http://bit.ly/9H4qBX

Project website: http://proj.badc.rl.ac.uk/preparde/wiki
Project blog: http://proj.badc.rl.ac.uk/preparde/blog

# Thanks! Any questions?

Yor data organishun is harrible



Pweez to stand by...we seemz to be ekzpeerinsing teknikul difikulteez.