# *From data management to data science: current trends & future challenges*
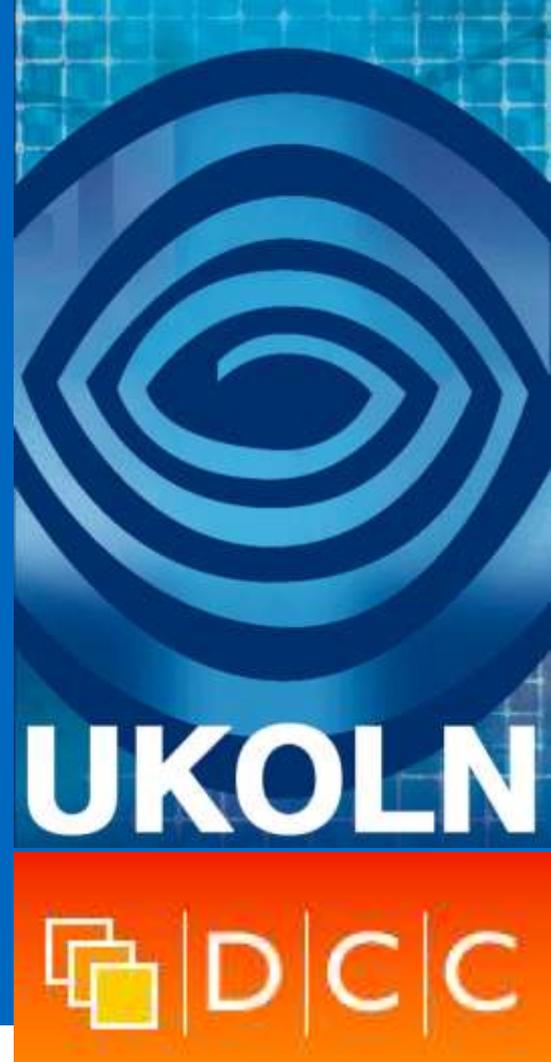
**Dr Liz Lyon,**

Associate Director, UK Digital Curation Centre,
Director, UKOLN, University of Bath, UK

DST Seminar, Pretoria South Africa, November 2012

**UKOLN**

**D|C|C**

**UKOLN is supported by:**

**JISC** Microsoft Research

**www.ukoln.ac.uk**

UNIVERSITY OF **BATH**

THE SUNDAY TIMES
**UNIVERSITY OF THE YEAR** 2011-12

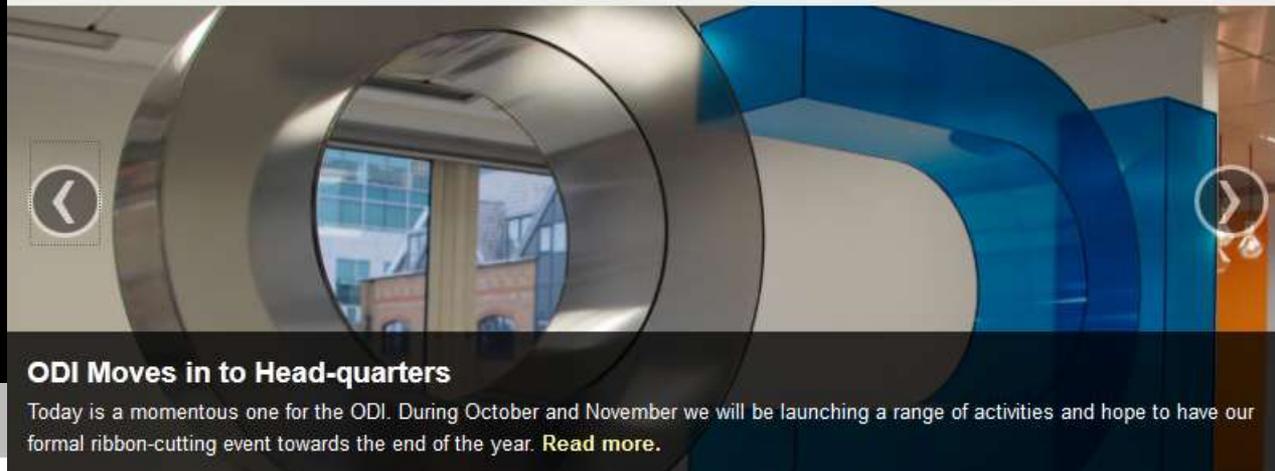A centre of expertise in digital information management

UKOLN

# Overview

- Trends, Launches, Readiness
- UK Social science data
- Institutions, data and the DCC University of Bath experience
- Data scientists: a new breed

UK government supports open data

Science as an
open enterprise

June 2012

THE
ROYAL
SOCIETY

Royal Society Report

Science as an Open
Enterprise

June 2012

10 Recommendations

Cites the DCC

*http://royalsociety.org/uploadedFiles/Royal_Society_Content/policy/projects/sape/2012-06-20-SAOE.pdf*

# Report sound-bytes

- "intelligently open data"
- "Scientists ….. are increasingly turning to their university libraries and institutional repositories for support for their data….."
- "familiarity with … tools and principles of data management should be an integral part of the training of scientists in the future…. "
- "The skills of data scientists are crucial in supporting the data management needs of researchers and of institutions."

McKinsey&Company

McKinsey Global Institute

May 2011

Big data: The next frontier
for innovation, competition,
and productivity

Implications of "Big Data" and data science for organisations in all sectors

Predicts a shortage of 190,000 data scientists by 2019

*http://www.mckinsey.com/Insights/MGI/Research/Technology_and_Innovation/Big_data_The_next_frontier_for_innovation*

# Data : from Big to Broad (Jim Hendler)

**BROAD data**

**Tetherless World Constellation**

- 4th context: Broad Data
  - The huge amount of freely available, but widely varied, Open Data on the World Wide Web (Structured and Semi-structured)
    - Example: The extended Facebook OGP graph (the part outside Facebook's datasets)
    - Example: The growing linked open data cloud of freely available RDF linked data
    - Example: More than 710,000 datasets that are available on the Web free from governments around the world

# Research Data Alliance

**Vision**

Researchers around the world sharing and using research data without barriers

*http://rd-alliance.org//*

## Currently involved

The individuals currently involved in working to bring the Research Data Alliance into being are:

- Fran Berman, Professor of Computer Science, Rensselaer Polytechnic Institute
- Juan Bicarregui, Acting Director e-Science, STFC Rutherford Appleton Laboratory
- Leif Laaksonen, Collaboration Director, CSC Finland
- Beth Plale, Director Data to Insight Center, Professor of Computer Science, Indiana University Bloomington
- Andrew Treloar, Director of Technology, Australian National Data Service
- Ross Wilkinson, Executive Director, Australian National Data Service
- Peter Wittenburg, The Language Archive, Max Planck Institute for Psycholinguistics
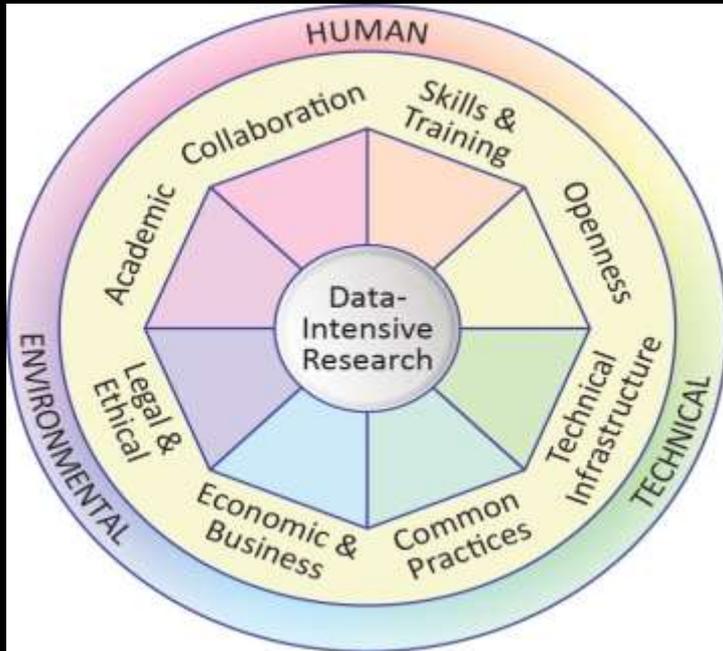- John Wood, Secretary General of the Association of Commonwealth Universities

**RD-Alliance**
forum

RESEARCH DATA ALLIANCE
**CANDIDATE WORKING GROUP**
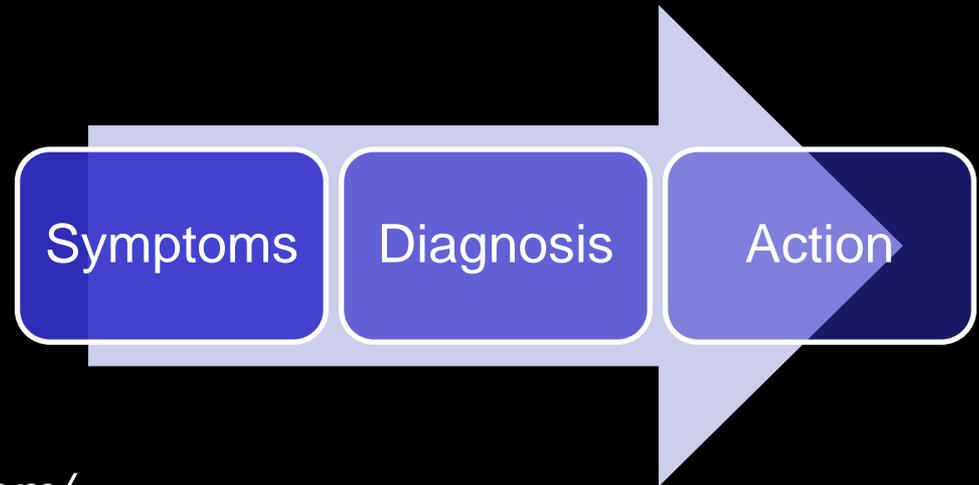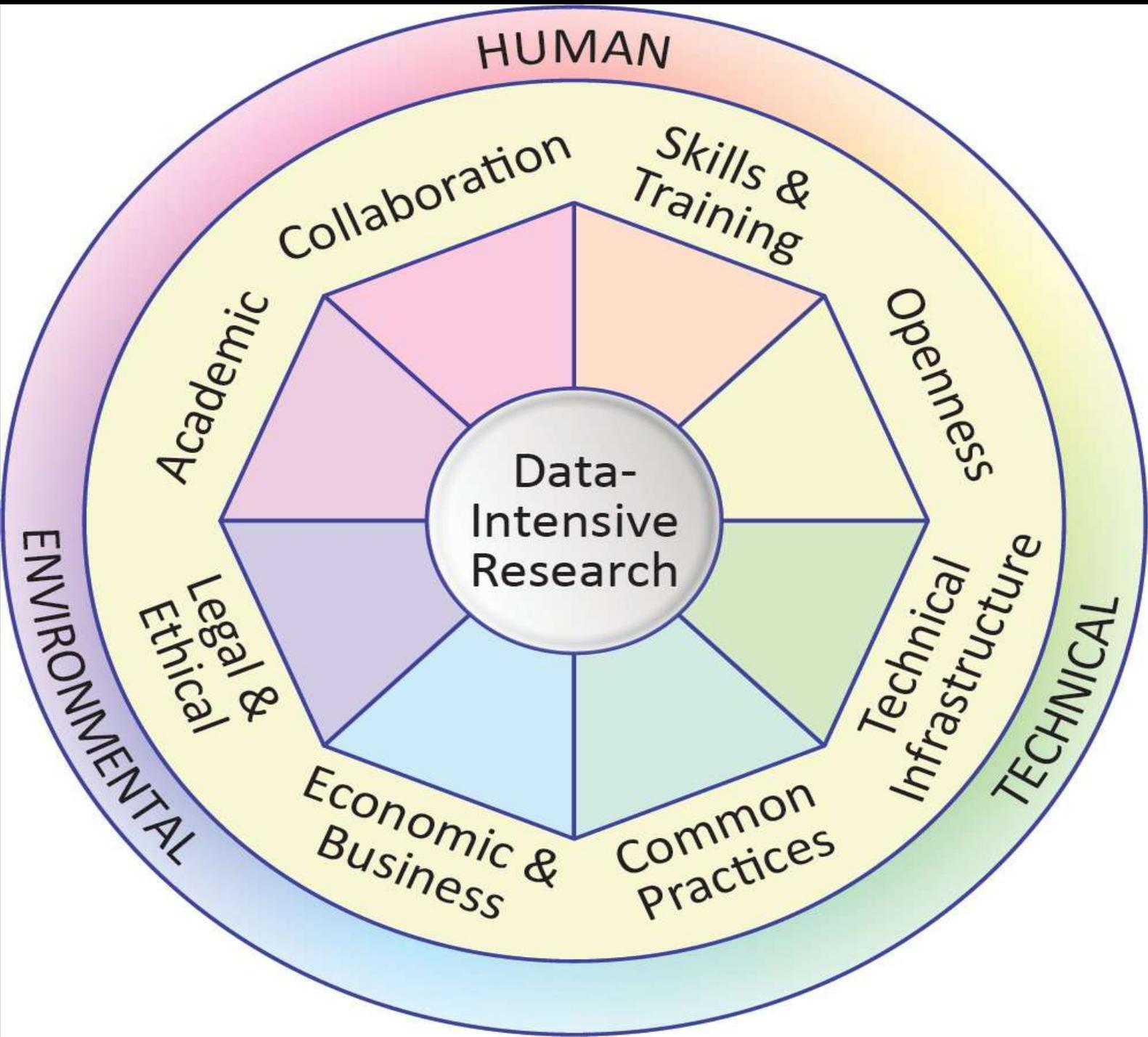**CASE STATEMENT GUIDELINES V.1**

*OCTOBER, 2012*

*RD-Alliance Launch March 2013, Gothenburg, Sweden*
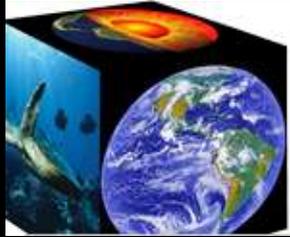
# Community Capability Model Framework CCMF



- Understanding the data habitat :
- Funder, institution, PI

| Symptoms | Diagnosis | Action |
|----------|-----------|--------|

*http://communitymodel.sharepoint.com/*

CCMF 8 Capability Factors

HUMAN

Collaboration

Skills & Training

Academic

Openness

Data-Intensive Research

ENVIRONMENTAL

Legal & Ethical

Technical Infrastructure

TECHNICAL

Economic & Business

Common Practices

# Applying CCMF: geosciences

EarthCube

- Disciplinary readiness for Cyber-Infrastructure (CI)
- NSF geoscience initiative
- Cross-Domain Interoperability Roadmap
- Applied the CCMF model

**EarthCube ROADMAP**

PREPARED BY
**CROSS-DOMAIN INTEROPERABILITY TEST BED GROUP**

Version 1.1    August 2012

*https://www.dropbox.com/s/0oqk5ostahfokbg/interop_roadmap_master8_Aug16.pdf*

# UK social science:  data ready



COMING SOON

**UK Data Service**

A new national data service for social and economic data

- UKDS Launch
- >40 yrs data archive
- Policy mandate to deposit datasets
- Admin data linkage



SECURE DATA SERVICE
enabling the research community



E·S·R·C
ECONOMIC & SOCIAL RESEARCH COUNCIL

Economic and Social Research Council
Shaping Society

# April 2011 - EPSRC Letter to VCs

**EPSRC**
Pioneering research
and skills

Engineering and Physical Sciences Research Council

> EPSRC expects all those institutions it funds
> - to develop **a roadmap** that aligns their policies and processes with EPSRC's **expectations** by **1st May 2012**;
> - to be fully compliant with these **expectations** by **1st May 2015**.

- Awareness of regulatory environment
- Data access statement
- Data policies and processes
- Data storage
- Structured metadata descriptions
- DOIs for data
- Data securely preserved for a minimum of 10 years

# How ready are institutions?

# Research360 Project



The Research360 Institutional Research Lifecycle Concept

JISC Managing Research Data Programme

Focus on academia-industry data

Focus on meeting EPSRC requirements

# 10 RDM challenges for institutions

1. *Priorities, risks, benefits?*
2. What research data do we have?
3. Where is it stored?
4. What state is it in?
5. How long should we keep it?
6. Do we have a data strategy?
7. How much will it cost ?
8. Do we have a data policy?
9. Who has responsibility for RDM?
10. Are researchers and services staff "data-aware" and "data-savvy"?

# 1. Risk: where is your data?

2. Reputation
3. Quality
4. Scale
5. Partnerships
6. Funding

Getting attention….

# Data Requirements at your Univ?


Data Asset Framework


BOS Bristol Online Surveys


CARDIO
collaborative assessment of research data infrastructure and objectives


D|C|C

http://www.dcc.ac.uk/

If research data lies at the heart of your organisation, you need to know that you have adequate infrastructure, staff skills and resources, and senior management support in place to ensure that your data is effectively managed for validation, reuse and evidential purposes.

**CARDIO enables you to:**

✓ collaboratively assess data management requirements, activity, and capacity at your institution

✓ build consensus between data creators, information managers and service providers

✓ identify practical goals for improvement in data management provision and support;

✓ identify operational inefficiencies and opportunities for cost saving;

✓ make a compelling case to senior managers for investment in data management support

University of Bath Research Data Survey results



- 210 respondents (3.5% response rate) : PIs, ROs, postgrads
- ***Some preliminary findings:***
- Most have not had to produce a data management plan (81%)
- Much data is confidential, anonymised, under non-disclosure agreements, commercially sensitive, DPA, encrypted
- Data is also in non-digital form: lab notebooks, interviews
- Researchers store data on Univ Bath shared filestore ☺
- They also use Dropbox, USB sticks, home computers ☹
- Data loss: accidental deletion, hardware failure, obsolescence
- Open data is not the norm – often shared informally
- Lack of recognition for data sharing and reuse is an issue

# NERC Data Value Checklist

**http://www.dcc.ac.uk/**

## Checklist

**Mandatory criteria:** These are mandatory criteria and answering 'Yes' to one or more of the questions below will automatically result in selection for retention.

| Legal/statutory considerations | Yes | No |
|---|---|---|
| Is there a legal or legislative reason for NERC to retain the data? | | |
| Is there any obvious reason why the data may be used in litigation, public enquiries, police investigations or any report or paper that could be legally challenged? | | |
| Are there any financial or contractual obligations that require us to retain the data? | | |

**Important criteria:** These are primary criteria and answering 'Yes' to at least one of the questions from each section below should probably result in selection for retention.

| Policy | | |
|---|---|---|
| Are the data a result of full or partial NERC funded activities? | | |
| Do the data fall within the selected Data Centre's Collection Policy? If no – refer to NERC Data Coordinator or pass to the correct data centre. | | |
| **Scientific or historic value** | | |
| Are the data a unique unrepeatable measurement of the environment? | | |
| Do the data have a broad geographical or temporal extent that makes them useful to others? | | |
| Do the data have historic value i.e. do they represent a landmark in scientific discovery? | | |
| Do the data include changes in processing methods, new standards or set any precedents? | | |
| Do the data support current projects or trends in science? | | |
| Are the data likely to meet the future needs/direction of the scientific community? | | |
| Do the data contribute to a pre-existing collection? | | |
| Is there potential for re-use of the data? | | |
| Are the data likely to be cited or referenced in a publication? | | |

**Supporting criteria:** These are important criteria and answering 'Yes' to the majority of the questions below should result in selection for retention.

| Origin | | |
|---|---|---|
| Do the data have their original integrity? | | |
| Would the data be costly to reproduce? | | |
| Will this become the reference copy of the data? | | |
| **Condition** | | |
| Do the data have relevant metadata available? | | |
| Are there proportionally more valuable data than non-valuable data within the collection? | | |
| Can the data be ingested into the Data Centre without significant additional processing? (reboxing, sifting, conversion etc) | | |
| Are the data in a suitable condition for addition to the collections? i.e. Readable, Undamaged, | | |

- Research Council data requirements
- Institutional Roadmaps for EPSRC
- *http://www.bath.ac.uk/rdso/University-of-Bath-Roadmap-for-EPSRC.pdf*

# Alignment with EPSRC Expectations



**Where we are now**

**Roles and Responsibilities:**
Who's responsible

**Objectives**:
Where we need to be

**How we're going to get there**

**Milestones:**
When it will be done by

- Finding the balance between 'gold standard' and 'good enough'
- Ensuring that all stakeholders were satisfied

# Roadmap → Strategic planning ↓

## Operational Plan

- Detailed activities
- Defined roles and responsibilities
- Timescale for implementation
- Ensures infrastructure development
- Ensures compliance with funder requirements

## Strategic Plan

- Aligns data management with existing institutional strategic aims
- Identifies the benefits of research data management
- Recommendations for high level next steps

## Business Case

- Benefits of investing in data management
- Risks associated with not investing in data management
- Options for different levels of investment
- Rational for recommendations

Transition from project-based to integrated infrastructure

| Institution | Policy name | Date released |
|---|---|---|
| University of Oxford | Statement of commitment to Research Data Management (formal policy forthcoming through the DaMaRO project) | 2010 |
| University of Edinburgh | Research Data Management Policy | 16 May 2011 |
| University of Northampton | Research Data Policy | June 2011 |
| University of Hertfordshire | Data Management Policy (see s.7 on research data and the appendix 'Guide to RDM') | 1 Sept 2011 |
| University of Warwick | Research Data Management Policy | 7 Nov 2011 |
| Glyndwr University | Policy on the Management of and Access to Research Data | 20 December 2011 |
| University of Southampton | Research Data Management Policy | February 2012 |
| University of East London | Research Data Management Policy for UEL | 15 March 2012 |
| Brunel University | Research Data Management Vision | 20 March 2012 |
| Queen Mary, University of London | Research Data Management Policy | 7 June 2012 |
| University of Sheffield | Research Data Management Policy | July 2012 |

**Draft policies**

University of Leeds - via the RoaDMaP project
Timeline of developments including draft policy text

University of Lincoln - via the Orbital project
Blog post with link to the draft policy text

University of Manchester - via the MiSS project
'Towards a Research Data Management Policy' document outlining progress

University of Exeter - Open Access and Research Data Management Policy

University of Exeter - Open Access and Research Data Management Policy for PGR Students

# Institutional data policy development

- Aspirational?
- Pragmatic?
- Emergent?
- High-level?
- With teeth?

# Developing a research data policy: reflections from Bath

- Keep it as succinct as possible
- It does need to have teeth
- Consult widely – many stakeholders
- Essential to discuss with the legal office
- Definitions of terms are helpful
- Provide detail in supporting procedures
- Research Data Steering Group approval
- PVC Research lead
- Roles and Responsibilities (again)

| Role | Responsibilities | Requirements | Relationships |
|---|---|---|---|
| Director Information Services / CIO University Librarian | To lead and co-ordinate data informatics support | Appropriate LIS structure in place | PVC Research, Deans, Associate Deans, Faculty/School Directors of Research, IT Director, Director Research Support |
| | | Library staff with data informatics & research data management skills | Other key institutional stakeholders |
| | | Institutional repository with content links to underlying research data | Open Access Publishers |
| Data librarian / Data scientist / Liaison /Subject / Faculty Librarian | To deliver expert data informatics advice and guidance to research staff | Knowledge of data management planning and data audit and assessment tools | DTCs, post-grads, PIs |
| | To facilitate access to datasets for PIs, research staff, postgraduate and undergraduate students | Knowledge of selection and appraisal, metadata standards and schema, data formats, domain ontologies, identifiers, data citation, data licensing | DCC |
| | | | DataCite |
| | | Knowledge of appropriate disciplinary data centres, | Data centre staff |

*Advocacy*

*Full mapping : Informatics Transform, IJDC Current issue, 2012*

# RDM Challenges for institutions

10. Are researchers and professional services staff "data-aware" & "data-savvy"?

# Training cascade

Research360
Managing data across the institutional research lifecycle

# Supporting resources

Research Data Management Website

Researcher Development Framework
1-2 hour training sessions

Virtual Training module eg MANTRA @ Edinburgh

DCC Briefing Paper on managing academic-industry research data

DCC How To Guide on managing academic-industry research data

Data Storage Guidelines

# Data Management Planning





Institutional DMP template

Research360 is developing:

- Institutional template for **DMP***online*

- Guidelines for postgraduate students

- Support for researchers during grant application process:

  - Internal peer review college

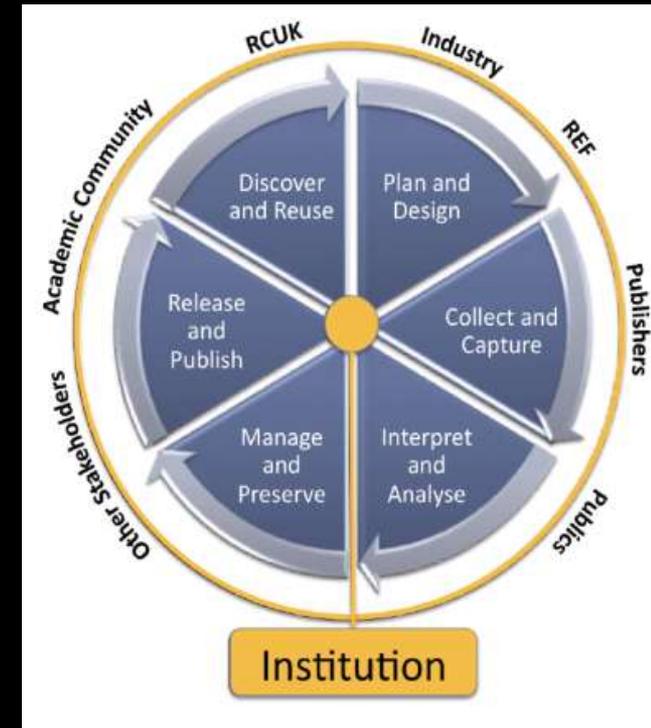# Institutional data scientist role

- **Co-ordination and Collaboration**
  - Liaison / subject librarians
  - Repository manager
  - IT/Computing Services
  - Research Support & Development Office
  - Doctoral Training Centres
  - Researchers
- **Advocacy**
- **Training**



*Liz Lyon, Informatics Transform, IJDC Current Issue, 2012*

Research360
Managing data across the institutional research lifecycle

# Family of data scientist roles

- ***data engineer -*** focus on software development, coding, programming, tools
- ***data analyst –*** focus on business/scientific analytics and statistics e.g. R, SAS, Excel to support researchers and modellers, business
- ***data librarian –*** focus on advocacy, research data management / informatics in a university / institute
- ***data steward –*** focus on long term digital preservation, repositories, archives, data centres
- ***data journalist –*** focus on telling stories and news

Jer Thorp: Hope / Crisis, NYT Word Frequency

New York Times Data Artist in Residence, Jer Thorp Joins Stellar Cast of Speakers at TEDxVancouver 2011
Posted by TEDxVancouver Team on October 17th, 2011 · No Comments

The New York Times

*Infrastructure, Intelligence, Innovation: driving the Data Science agenda*
8th International Digital Curation Conference, Amsterdam, 14-16 January 2013

# Thank you.

*DCC Resources can be downloaded from*

*http://www.dcc.ac.uk*

*Slides (including What is a data scientist?) at*

*http://www.ukoln.ac.uk/ukoln/staff/e.j.lyon/presentations.html*

*Informatics Transform paper at*

*http://www.ijdc.net/index.php/ijdc/article/view/210/279*