# *Building Capacity and Capability for Data : Requirements, Challenges, Opportunities*

**Dr Liz Lyon,** Associate Director, UK Digital Curation Centre
Director, UKOLN, University of Bath, UK

Horizon2020 Workshop Brussels, May 2012

UKOLN

DCC

**UKOLN is supported by:**

JISC

**www.ukoln.ac.uk**

A centre of expertise in digital information management

UNIVERSITY OF BATH

THE SUNDAY TIMES
UNIVERSITY OF THE YEAR 2011-12

# Running order…..

- Data landscape snapshots
- Roles and responsibilities
- Skills and competencies
- Gaps and opportunities

"The ability to take data - to be able to understand it, to process it, to extract value from it, to visualise it, to communicate it - that's going to be a hugely important skill in the next decades."

*Hal Varian, Chief Economist, Google*

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

"Big Data"
Data scientist

Data Science Revealed community survey

About how much time do you spend on the following activities (% A lot)

| Activity | Big Data | Normal Data |
|---|---|---|
| Acquiring new data sets | 48% | 27% |
| Parsing data sets | 50% | 21% |
| Filtering and organizing data | 58% | 34% |
| Mining data for patterns | 52% | 23% |
| Applying advanced algorithms to solve analytical problems | 48% | 22% |
| Representing data visually | 54% | 31% |
| Telling a story with data | 50% | 27% |
| Interacting with data dynamically | 58% | 30% |
| Making business decisions based on data | 60% | 35% |

■ Big Data  ■ Normal Data

theguardian

News | Sport | Comment | Culture | Busine

News > Datablog

DATABLOG
Facts are sacred

# What is a data scientist?

It's the job of the moment. But what exactly is a data scientist?

# Other data-related roles?

# Data journalist
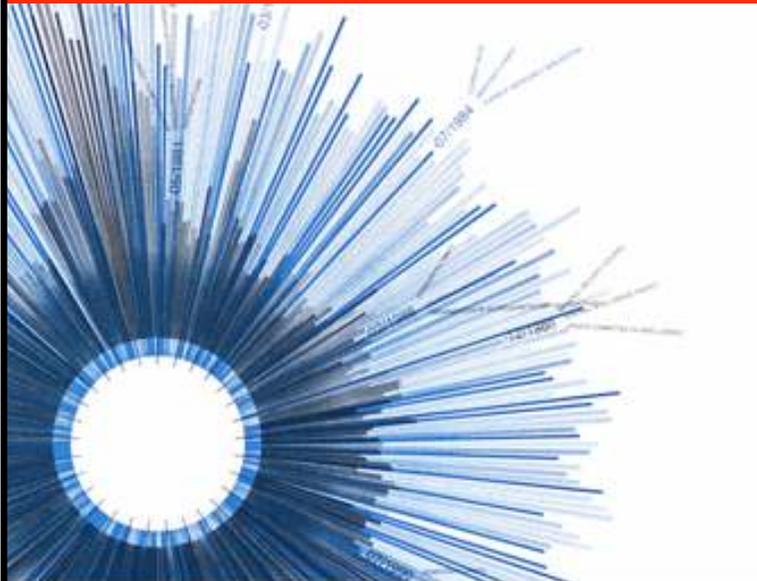
New York Times Data Artist in Residence, Jer Thorp Joins Stellar Cast of Speakers at TEDxVancouver 2011

1 · No Comments »

The New York Times

Jer Thorp: Hope / Crisis, NYT Word Frequency
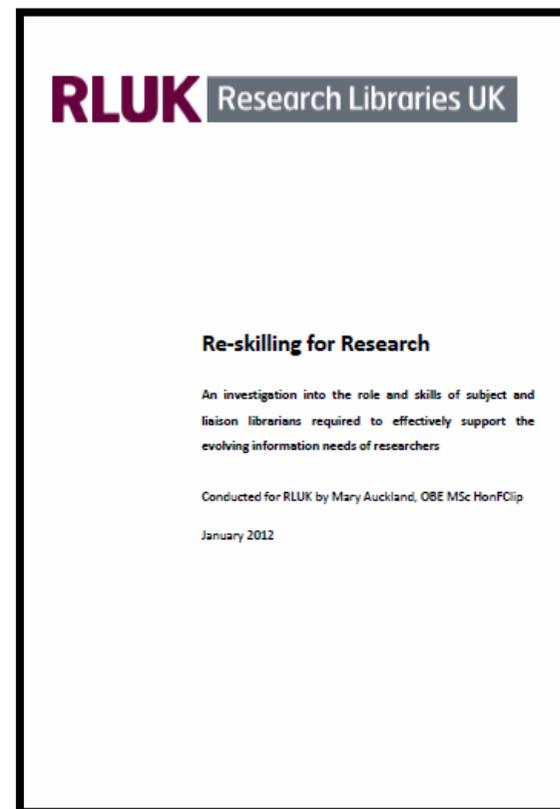
| Position | Location |
|---|---|
| Science Data Librarian | Stanford |
| Data Management Librarian | Oregon State |
| Social Sciences Data Librarian | Brown |
| Data Curation Librarian | Northeastern |
| Data Librarian | New South Wales |
| Research Data Management Co-ordinator | Sydney |
| Research Data & Digital Curation Officer | Cambridge |
| Data Services Librarian | Iowa |
| Data Analyst | ANDS |
| Institutional Data Scientist | Bath |

DCC

RLUK/Mary Auckland:
Reskilling for Research
9 areas are skill gaps
for subject librarians

Sheila Corrall: Libraries,
Librarians and Data
Many action exemplars

**RLUK** Research Libraries UK

**Re-skilling for Research**

An investigation into the role and skills of subject and
liaison librarians required to effectively support the
evolving information needs of researchers

Conducted for RLUK by Mary Auckland, OBE MSc HonFClip

January 2012

2012: Libraries in review

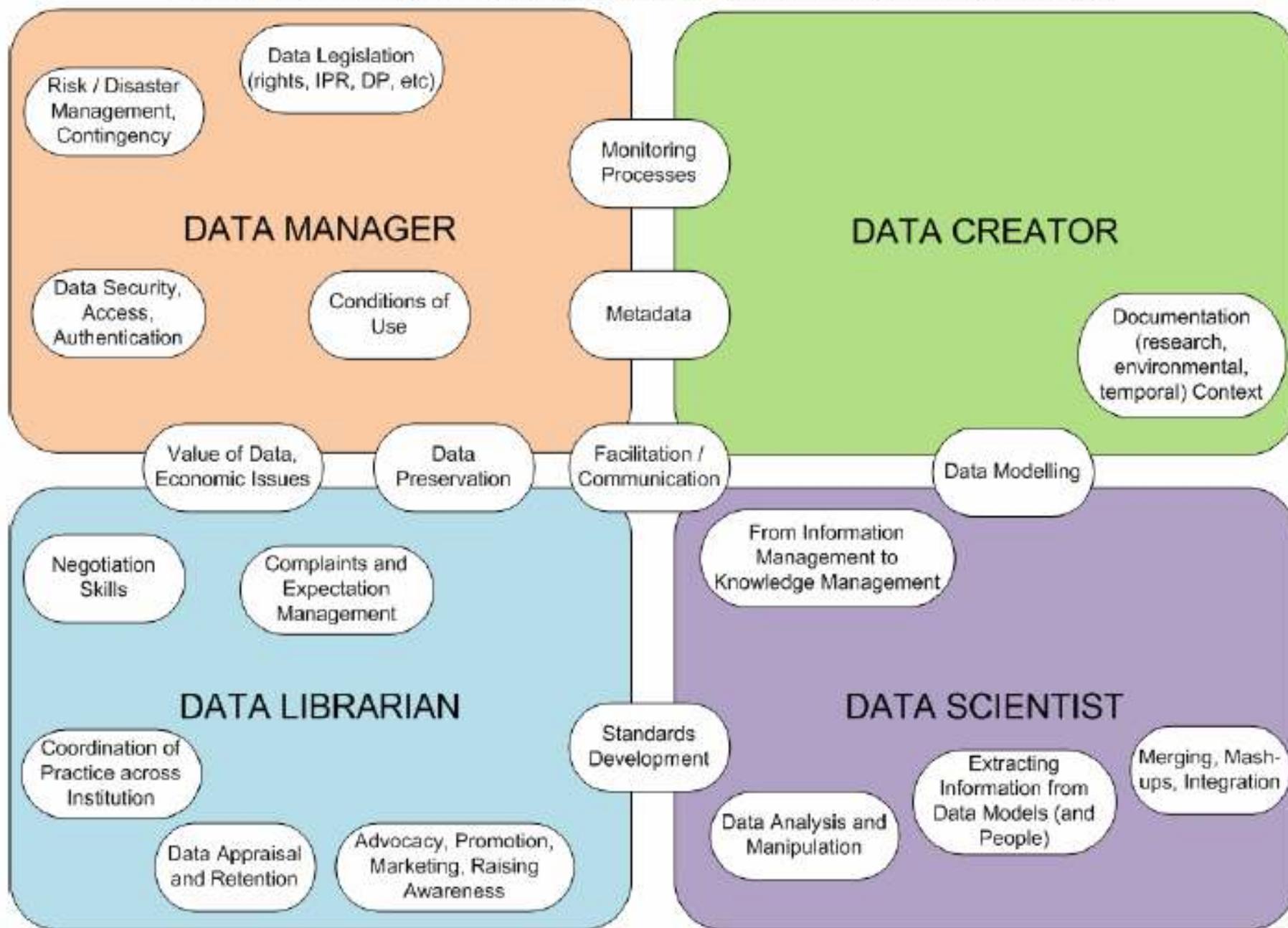| Skill gap | 2-5 years | Now |
|---|---|---|
| Preserving research outputs | 49% | 10% |
| Data management & curation | 48% | 16% |
| Comply with funder mandates | 40% | 16% |
| Data manipulation tools | 34% | 7% |
| Data mining | 33% | 3% |
| Metadata | 29% | 10% |
| Preservation of project records | 24% | 3% |
| Sources of research funding | 21% | 8% |
| Metadata schema, discipline standards, practices | 16% | 2% |

*Data from RLUK/Mary Auckland: Reskilling for Research 2012*

"Very few librarians are likely to have specialist scientific or medical knowledge - if you train as a research scientist or a medic, you probably won't become a librarian."

# CORE SKILLS FOR DATA MANAGEMENT

A follow-up from the second DCC Research Data Management Forum (November 2008)

## DATA MANAGER

- Data Legislation (rights, IPR, DP, etc)
- Risk / Disaster Management, Contingency
- Monitoring Processes
- Data Security, Access, Authentication
- Conditions of Use
- Metadata

## DATA CREATOR

- Documentation (research, environmental, temporal) Context

- Value of Data, Economic Issues
- Data Preservation
- Facilitation / Communication
- Data Modelling

## DATA LIBRARIAN

- From Information Management to Knowledge Management
- Negotiation Skills
- Complaints and Expectation Management
- Coordination of Practice across Institution
- Standards Development
- Data Appraisal and Retention
- Advocacy, Promotion, Marketing, Raising Awareness

## DATA SCIENTIST

- Data Analysis and Manipulation
- Extracting Information from Data Models (and People)
- Merging, Mash-ups, Integration

- Leadership & co-ordination
- Strategy and planning
- Policy
- Legal and ethical (FoI, Data Protection)
- Advocacy (data informatics)
- Data repositories
- Data storage
- Data analysis
- Data visualisation
- Data mining
- Data modelling
- Data licensing
- Training….

WANTED
GlassGiant.com

# University data roles?

- Roles (7 listed)
- Responsibilities
- Requirements
- Relationships

| 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 |

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

1. Director IS/CIO/University Librarian
2. Data librarians /data scientist /liaison/subject/faculty librarians
3. Repository managers
4. IT/Computing Services
5. Research Support/Innovation Office
6. Doctoral Training Centres
7. PVC Research
8. *+ Public Engagement Office*

**Data roles**

| 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 |

*Leadership*

*Advocacy*

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

# April 2011 - EPSRC Letter to VCs



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

- Leadership
- Co-ordination
- Pan-institutional perspective
- Operational plan
- Wider strategic alignment



UNIVERSITY OF BATH

# University of Bath Roadmap for EPSRC

Compliance with Research Data Management Expectations

28th April 2012, Version 1.1
Authors: Dr Liz Lyon, UKOLN, & Dr Catherine Pink, UKOLN
Status: Submitted to Research Data Steering Group     5th April 2012
     Approved, with amendments, Research Data Steering Group    17th April 2012
     Submitted to Vice-Chancellor's Group (VCG)     23rd April 2012
     Submitted to VCG with revisions     30th April 2012
     Approved, with amendments, by VCG     30th April 2012

Acknowledgement

We would like to acknowledge the leadership of Monash University in the area of research data management. The Monash University Research Data Management Strategy and Strategic Plan 2012-2015, released under a CC-BY licence, was highly influential in the development of this document.

1

# Advocacy and support

- *Data requirements:* legacy data
- *Data management plans*: tools
- *Informatics*: disciplinary metadata schema, standards, formats, identifiers, ontologies
- *Citation:* links to publications
- *Reuse*: tracking your data

JISC

DCC because good research needs good data

# Understanding Data Requirements



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

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

# Data management plans

| Role | Responsibilities | Requirements | Relationships |
|---|---|---|---|
| Repository managers | To ensure research papers have persistent links to underlying research data | Knowledge of persistent identification mechanisms and publisher requirements | Data librarians / Data scientists / Liaison /Subject / Faculty Librarians |
| IT / Computing Services | To provide data storage infrastructure and guidance | Knowledge of data storage options including cloud-based services | EduServ data centre. Cloud service providers<br><br>National data centres |
| Research & Development Support Office / Research & Innovation Services | To provide RIM/CRIS capability for research outputs | Provision for non-textual outputs such as datasets, software and program code, gene sequences, models | Research funding bodies<br><br>Data scientists / Liaison /Subject / Faculty Librarians |

*Discovery*

*Storage*

*CRIS*

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

A Digital Curation Centre 'working level' guide

## How to Cite Datasets and Link to Publications

Alex Ball (DCC) and Monica Duke (DCC)

# How to cite data

Helping you to find, access, and reuse data

DataCite

# Using DOIs

# How to track impact

## total·Impact

### Uncover the invisible impact of research.

Create a collection of research objects you want to track. We'll provide you a report of the total impact of this collection. You can peruse a sample report or check out the most recently shared reports.

| **Collect research objects** | | **Create report** |
|---|---|---|

### Paste object IDs,
Add one DOI, PubMed ID, URL, or other supported identifier per line:

```
10.1371/journal.pcbi.1000361
20334632
2BAK
GSE2109
10.5061/dryad.1295
http://www.carlboettiger.info/research
/lab-notebook
http://www.slideshare.net/phylogenomics
/eisenall-hands
```

Add to collection

### ...or pull object IDs from existing collections.
▸ Mendeley profiles
▸ Mendeley groups
▸ Slideshare accounts
▸ Dryad dataset authors
▸ PubMed grants
▸ GitHub users
▸ GitHub organizations

Something missing on import?
See a list of current limitations.

### Name your collection:
my collection

get my metrics!

… or fetch a quick collection based on your Mendeley contacts and public groups »

*http://total-impact.org/*

- **Storage:** file-store, cloud, data centres, funder policy
- **Access:** embargoes, FoI

CRIS integration, CERIF and data

| Role | Responsibilities | Requirements | Relationships |
|------|------------------|--------------|---------------|
| Faculty Doctoral Training Centres | To supply training to new-entrant researchers and PIs | Knowledge of data management planning and data audit and assessment tools  Training programmes and modules | Deans & Associate Deans, PIs  Data librarian / Data scientist / Liaison / Subject / Faculty Librarians |
| PVC Research | To develop institutional research policy and code of practice | Understanding of data management compliance implications, risks including legal and ethical issues, and sustainability challenges | Deans & Associate Deans  Key service directors  Research & Development Support Office / Research & Innovation Services |
| Public Engagement Unit | To facilitate citizen participation in the research process | Understanding of open science methodologies and infrastructure | PVC Research Director, Communications Deans & Associate Deans, PIs  The Media |

*Training*

*Policy*

*Participation*

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

Institutional data policy development

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

# Doctoral Training Centres: Research360 Project @ Bath

JISC projects

DCC resources

- Leadership & co-ordination
- Strategy and planning
- Policy
- Legal and ethical (FoI, Data Protection)
- Advocacy (data informatics)
- Data repositories
- Data storage
- Data analysis
- Data visualisation
- Data mining
- Data modelling
- Data licensing
- Training….

WANTED
GlassGiant.com

# Gaps?  Opportunities??

Analyse LIS entry qualifications & increase STEM entrants

Target
- Biologists
- Chemists
- Mathematicians

*Lyon, Informatics Transform, IJDC 2012*

# Gaps?  Opportunities??

Define core components of data informatics and data science

- Metadata (discovery, preservation)
- Domain ontologies
- Visualisation e.g. VisTrails
- Workflow e.g. Taverna
- Analysis e.g. R

*Lyon, Informatics Transform, IJDC 2012*

http://www.flickr.com/photos/50542505@N08/5723947474/

# Data scientist flavours?

- Analysis, mining, modelling
- Informatics, advocacy, training
- Repositories, preservation
- Visualisation, simulations

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

# Thank you!

Informatics Transform  article

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

Slides

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

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