

# *The Informatics Transform: Re-engineering Libraries for the Data Decade*

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

VALA2012, Melbourne, Australia



This work is licensed under a Creative Commons Licence  
Attribution-ShareAlike 2.0

UKOLN is supported by:

**JISC**

[www.ukoln.ac.uk](http://www.ukoln.ac.uk)

A centre of expertise in digital information management



***“Data is the new oil.”***

Andreas Weigend, Stanford (ex Amazon)

***“The future belongs to companies and people that turn data into products”***

Mike Loukides, O'Reilly Media



Data...

<http://www.google.co.uk/imgres?q=illumina+bgi&hl=en&client=firefox-a&hs=JI2&rls=org.mozilla:en-GB:official&biw=1366&bih>



# Oceans: last unmapped frontier?



<http://www.wired.com/wiredscience/2011/09/ocean-sensor-network/>



<http://bohemiaadventures.blogspot.com.au/2010/06/bering-sea-day-dutch-harbor.html>

## NEWS POLITICS

[Home](#) [World](#) [UK](#) [England](#) [N. Ireland](#) [Scotland](#) [Wales](#) [Business](#) [Politics](#) [Health](#)  
[Entertainment & Arts](#)5 December 2011 Last updated at  
21:221.3K [Share](#) [f](#) [t](#) [e](#)

## Everyone 'to be research patient', says David Cameron



..using personal  
data for research

"Let me be clear, this does not threaten privacy, it doesn't mean anyone can look at your health records, but it does mean using anonymous data to make new medical breakthroughs.

# Share your genome data?

- Buy a DTC kit
- Join a project

Personal Genome Project



## Personal Kit

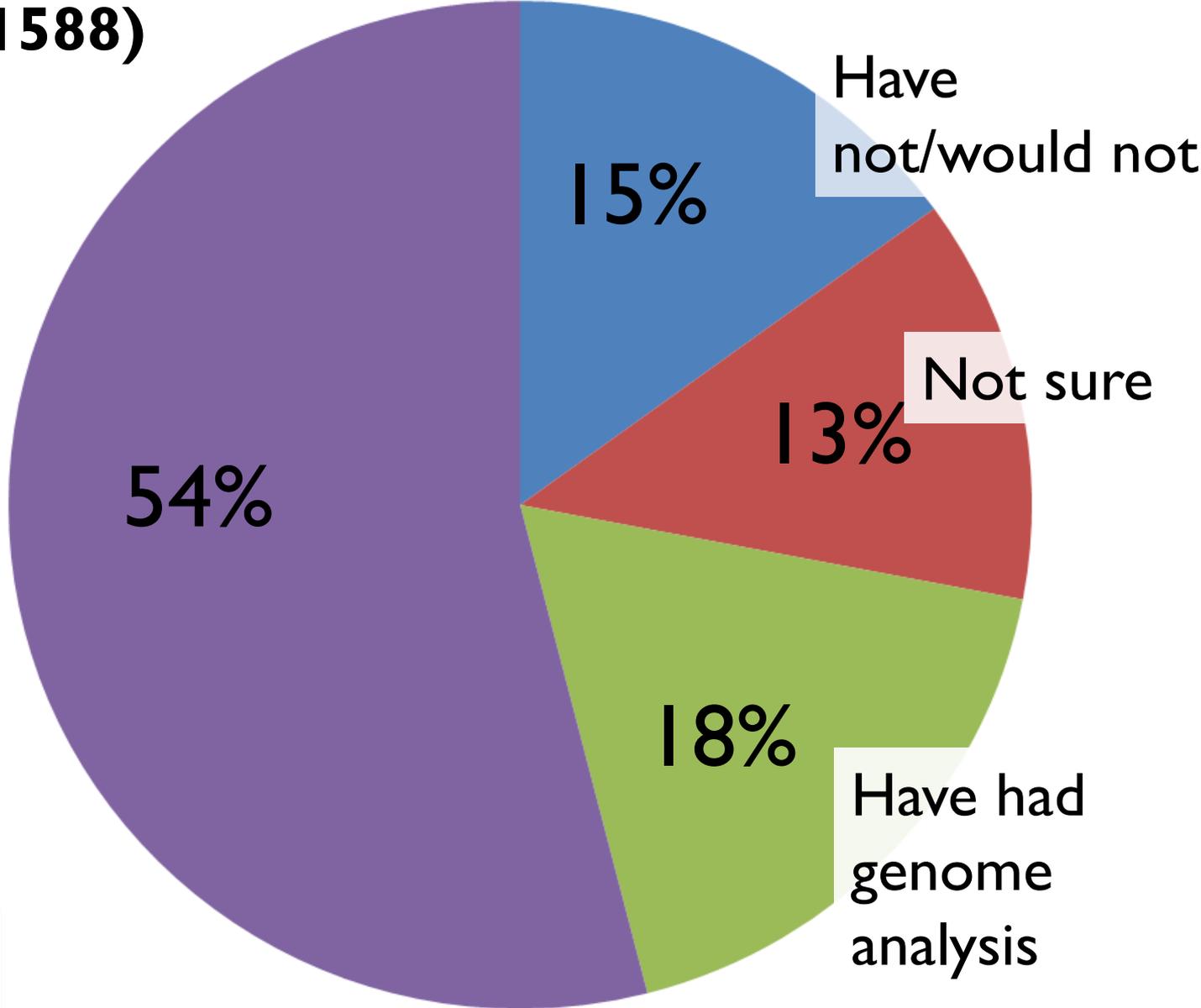
Order a kit for yourself today. Requires 1-year subscription to 23andMe's [Personal Genome Service®](#) billed at \$9 per month.

\$99 (1-year subscription required)

add to cart



**In a recent 2011 survey, *Nature* asked its readers whether they had, or would consider, a genome analysis (n=1588)**



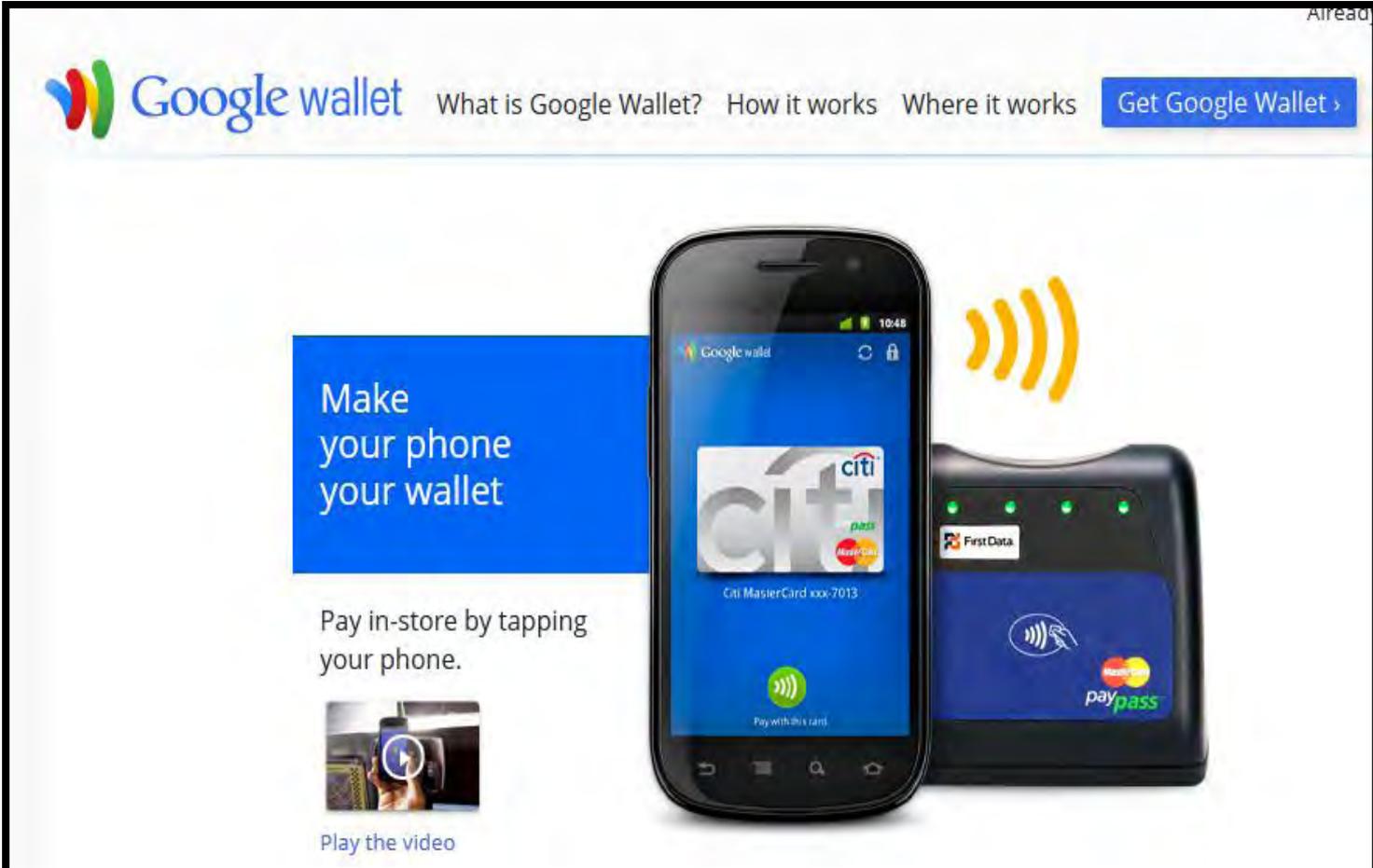
Would if given the opportunity

# Consumer data...

## Companies Are Gearing Up To Track Every Retail Transaction Through Your Smartphone

Michael Griffin

Already



The screenshot shows the Google Wallet website interface. At the top left is the Google Wallet logo. To its right are navigation links: "What is Google Wallet?", "How it works", and "Where it works". A blue button labeled "Get Google Wallet" is on the far right. The main content area features a blue box on the left with the text "Make your phone your wallet". Below this is a video player with a play button and the caption "Play the video". To the right of the video is a smartphone displaying the Google Wallet app with a Citi MasterCard card. Further right is a payment terminal with a "First Data" logo and a "paypass" logo. A yellow NFC symbol is positioned between the phone and the terminal.

**Google wallet** What is Google Wallet? How it works Where it works [Get Google Wallet](#)

Make your phone your wallet

Pay in-store by tapping your phone.

Play the video

*One in every nine people on Earth is  
on Facebook*

30billion pieces of content are shared on  
Facebook each month

The Facebook logo, consisting of the word "facebook" in white lowercase letters on a blue rectangular background.

*People upload 3000images to Flickr every  
minute*

The Flickr logo, featuring the word "flickr" in blue lowercase letters with a pink "r", followed by "from YAHOO!" in purple text.

Google+ has > 25million users

The Google+ logo, featuring the word "Google" in its multi-colored font followed by a plus sign.

*From 20 Social Media Statistics (Jeffbullas)*

# ...and conversations

## #Twitter and Tweets

there are  
**200,000,000**  
registered Twitter users

that's over 1,650 Tweets every second  
there are  
**one BILLION**  
new Tweets posted every week

every minute there are  
**138,888**  
new Tweets

just  
**5%**  
of Twitter users create 75% of the content  
that's a few very busy Tweeters...

not bad in just 5 years...  
almost  
**88%**  
of people have awareness of Twitter

there are up to  
**180,000,000**  
new tweets posted every day

as many as  
**52%**  
of users update their status every day  
daily updates make the difference...

want a retweet?  
**5PM** is the best time to get retweeted

there are  
**450,000**  
new Twitter accounts created every day  
that's 5.2 every second, every day...

addictive personalities...  
**24%**  
of all users check Twitter several times a day

search and ye shall find...

***“Data is the new oil.”***

Andreas Weigend, Stanford (ex Amazon)

***Data is more like soup –  
its messy and you don't  
know what's in it....***



Kyle Machulis

**“DIY”**

**Human  
physiology  
data**

# “Herculean” and “Heroic”

13 December 2011 Last updated at 17:20

15K

## LHC: Higgs boson 'may have been glimpsed'

By Paul Rincon

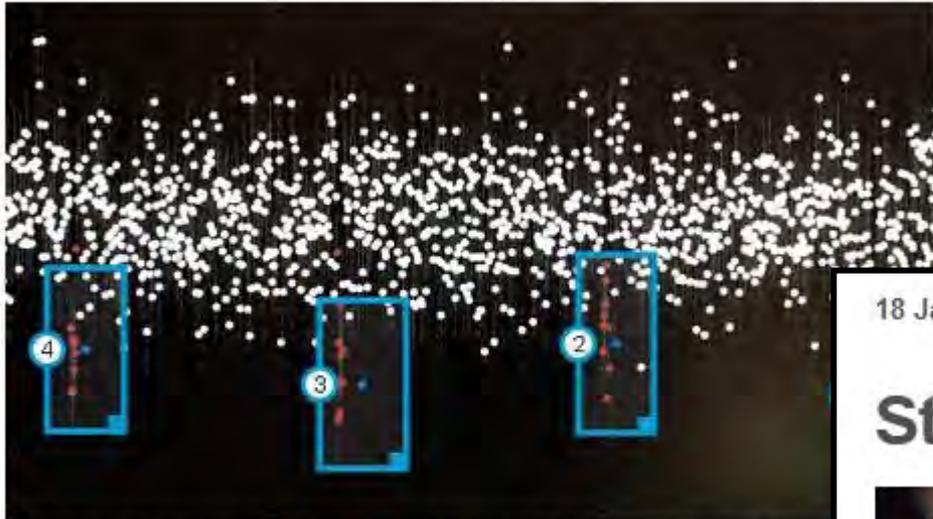
Science editor, BBC News website, Geneva



Two teams at the LHC have seen hints of what may well prove to be the Higgs.

Particle  
physics  
data

## Volunteers wanted for planet hunt



Time-lapsed images of a single star show dips in brightness as a plane

Members of the public are being asked to join the hunt for planets that could support life.

# “Crowd-sourced” astronomy

BBC Mobile

# NEWS

18 January 2012 Last updated at 22:16

251

## Stargazing viewer in planet coup



By Jonathan Amos

Science correspondent, BBC News



Amateur astronomer Chris Holmes from Peterborough stumbled upon SPH10066540

Researchers need help to  
manage their data.

This is a really exciting  
opportunity for libraries.....

**With a bit of re-engineering**



# 1. Leadership

(Getting attention...)

Six reasons why you should care  
about managing your research data

# 1. Risk: where is your data?



# 2. Reputation : data access, FOI

University told to hand over tree ring data - April 15, 2010



theguardian

[News](#) | [Sport](#) | [Comment](#) | [Culture](#) | [Business](#) | [Money](#) | [Life & style](#)

[News](#) > [Society](#) > [Smoking](#)

## Tobacco firm demands university's research on children and smoking

Stirling University fighting attempt by Philip Morris to gain access to research under freedom of information laws

Severin Carrell, Scotland correspondent  
guardian.co.uk, Thursday 1 September 2011 15.02 BST  
[Article history](#)



Philip Morris International, which makes Marlboro cigarettes, has asked for Stirling University's research on teenagers and smoking. Photograph: Paul Sakuma/AP

# 3. Quality: data gold standard



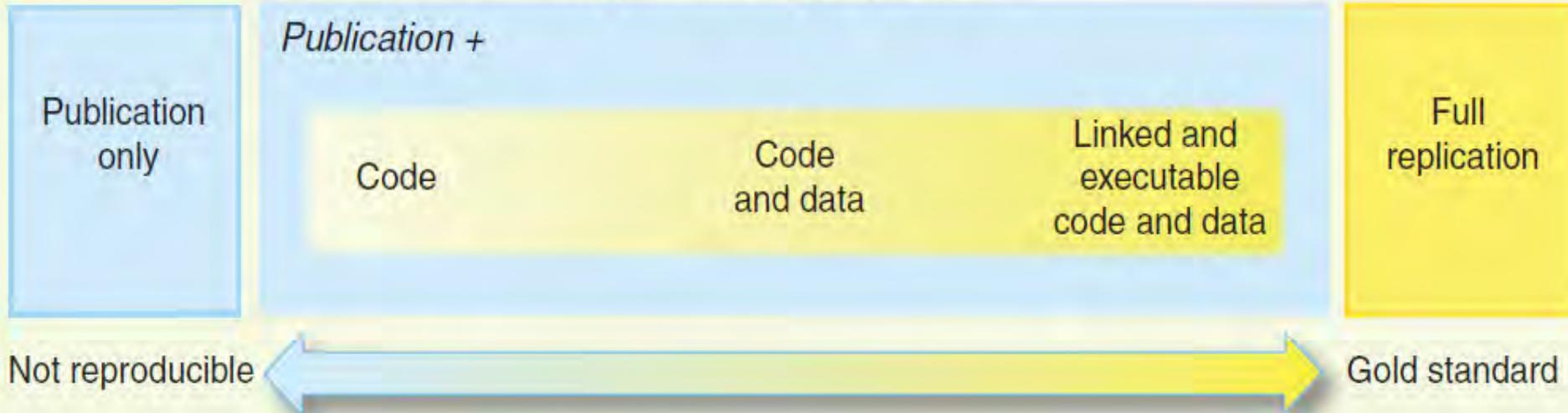
## Data Replication & Reproducibility

PERSPECTIVE

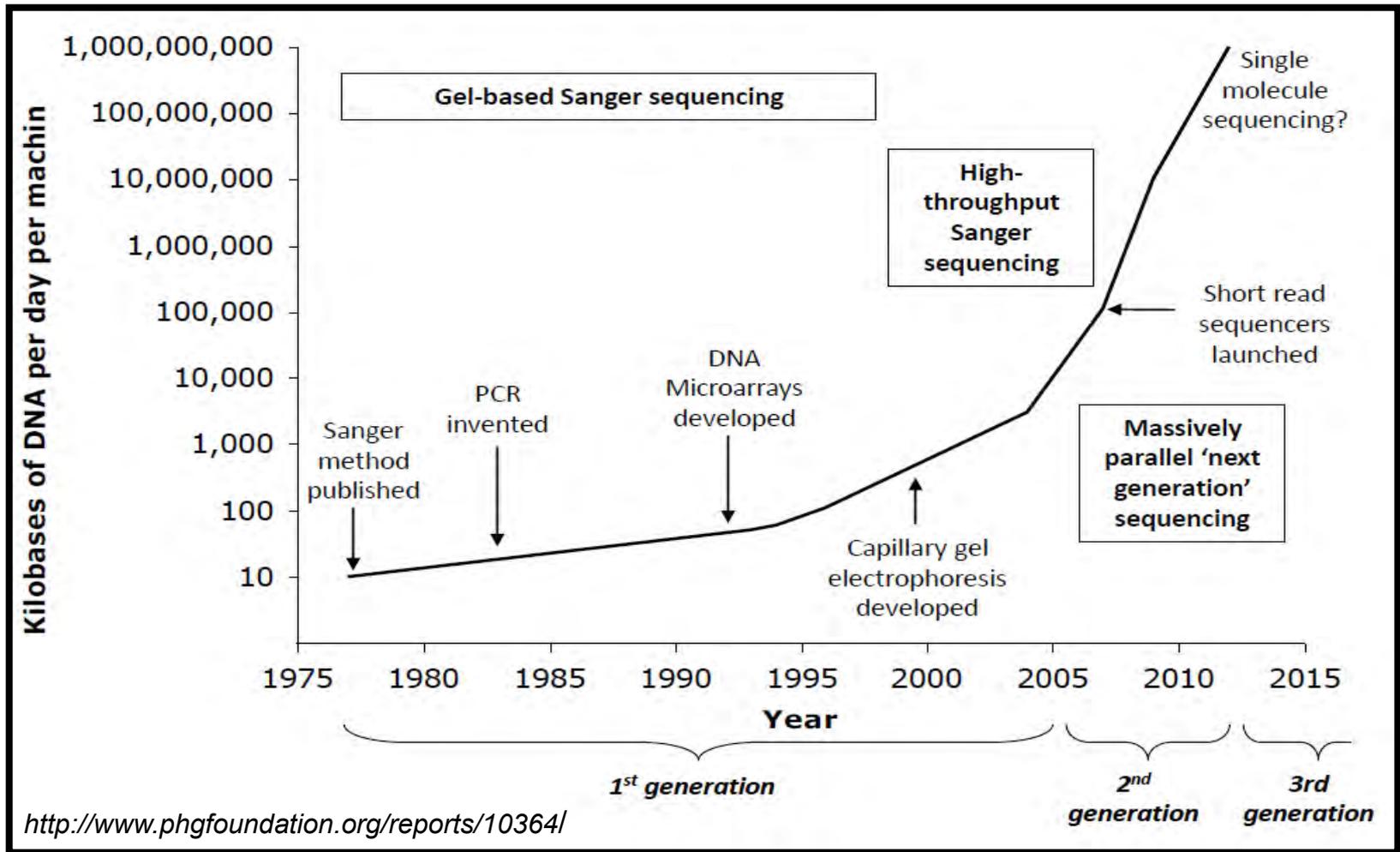
# Reproducible Research in Computational Science

Roger D. Peng

## Reproducibility Spectrum



# 4. Scale: an explosion of data



"A single sequencer can now generate in a day what it took 10 years to collect for the Human Genome Project"

# The New York Times 5.Partnerships

## Sharing of Data Leads to Progress on Alzheimer's

By GINA KOLATA

Published: August 12, 2010

Alzheimer's Disease Neuroimaging Initiative: a unique (open) \$60M partnership between NIH, FDA, universities and drug companies.

*“It was unbelievable. Its not science the way most of us have practiced in our careers. But we all realised that we would never get biomarkers unless all of us parked our egos and intellectual property noses outside the door and agreed that all of our data would be public immediately.”*

*Dr John Trojanowski, University of Pennsylvania*

# 6. Funding

**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 **1<sup>st</sup> May 2012**;
  - to be fully compliant with these **expectations** by **1<sup>st</sup> May 2015**.

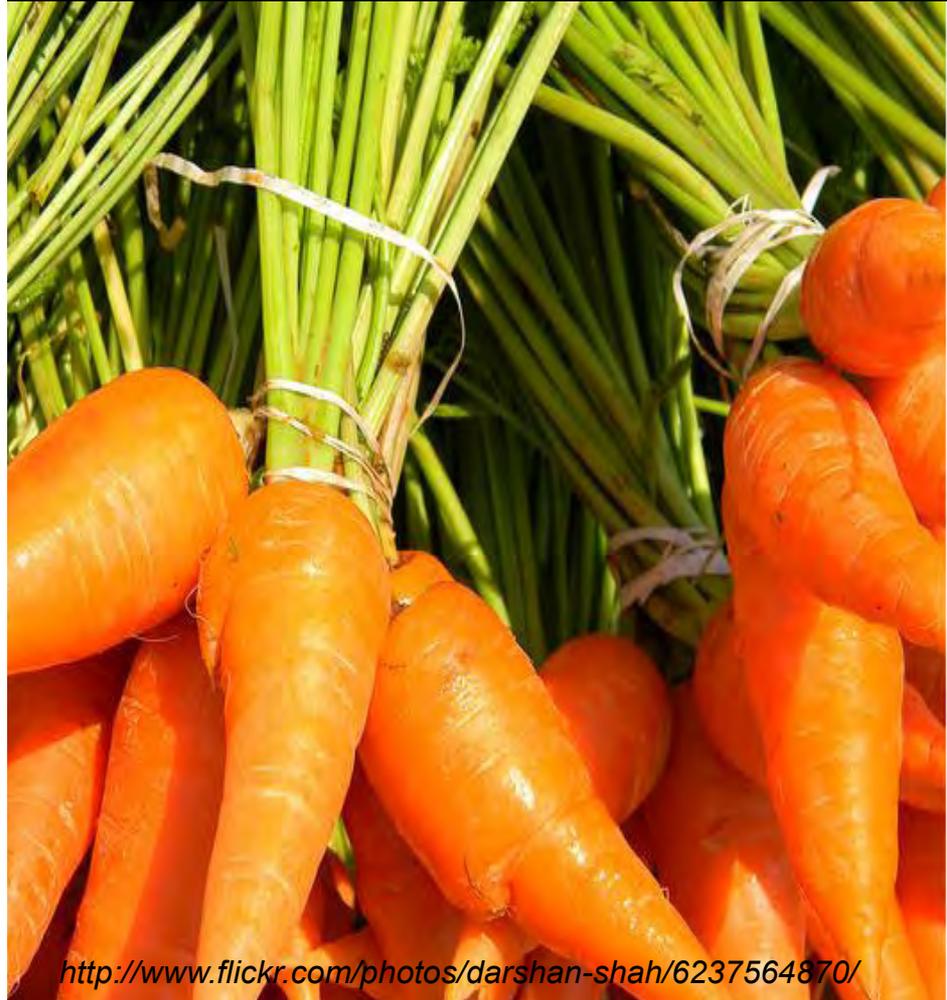
- Awareness of regulatory environment
- Data access statement
- Policies and processes
- Data storage
- Structured metadata descriptions
- DOIs for data
- Securely preserved for a minimum of 10 years



search\_ID:cs14846

WARM RECEPTION OF MR. PECKSNIFF BY HIS VENERABLE FRIEND  
Old Martin, with his burning indignation crowded into one vehement  
burst, struck him down upon the ground with a well-directed blow.  
*Chuzzlewit, chap. lii*

# Sticks ...and Carrots



# 2. Research Data Management services

(Providing tools & support)

# Understanding Data Requirements



```
001 0001 10
1001 0001 10
1 1000 0001
010 101 00
1001 100 00
1001 1000 00
0 1001 0001
1001 1000 00
1001 0001
```

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.

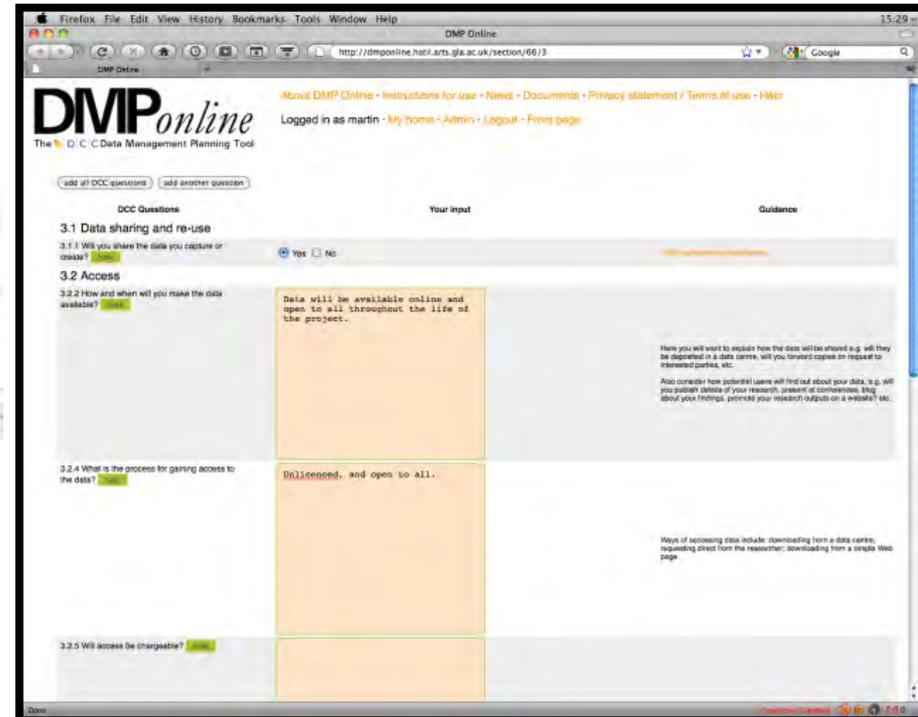
## CARDIO enables you to:

- ✓ 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



# Data management plans

**DMP***online*  
The  D C C Data Management Planning Tool



**DMP**Tool

Guidance and Resources for your Data Management Plan

- **Advocacy & Training**
  - **Informatics:** disciplinary metadata schema, standards, formats, identifiers, ontologies
  - **Storage:** file-store, cloud, data centres, funder policy
  - **Access:** embargoes, FOI

# What data to keep

A Digital Curation Centre 'working level' guide



## How to Cite Datasets and Link to Publications

Alex Ball (DCC) and Monica Duke (DCC)



Digital Curation Centre, 2011.  
Licensed under Creative Commons Attribution 2.5 Scotland:  
<http://creativecommons.org/licenses/by/2.5/scotland/>



A Digital Curation Centre and Australian  
National Data Service 'working level' guide

## How to Appraise & Select Research Data for Curation

Angus Whyte (DCC) and Andrew Wilson (ANDS)



Digital Curation Centre, Australian National Data Service 2010.  
Licensed under Creative Commons BY-NC-SA 2.5 Scotland:  
<http://creativecommons.org/licenses/by-nc-sa/2.5/scotland/>

# How to cite data

## How to License Research Data

Alex Ball (DCC)

DRAFT: 29 OCTOBER 2010



Digital Curation Centre, 2010.

Licensed under Creative Commons BY-NC-SA 2.5 Scotland:

<http://creativecommons.org/licenses/by-nc-sa/2.5/scotland/>

# Data Licensing

Bespoke licences  
Standard licences  
Multiple licensing  
Licence mechanisms



# Tools 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

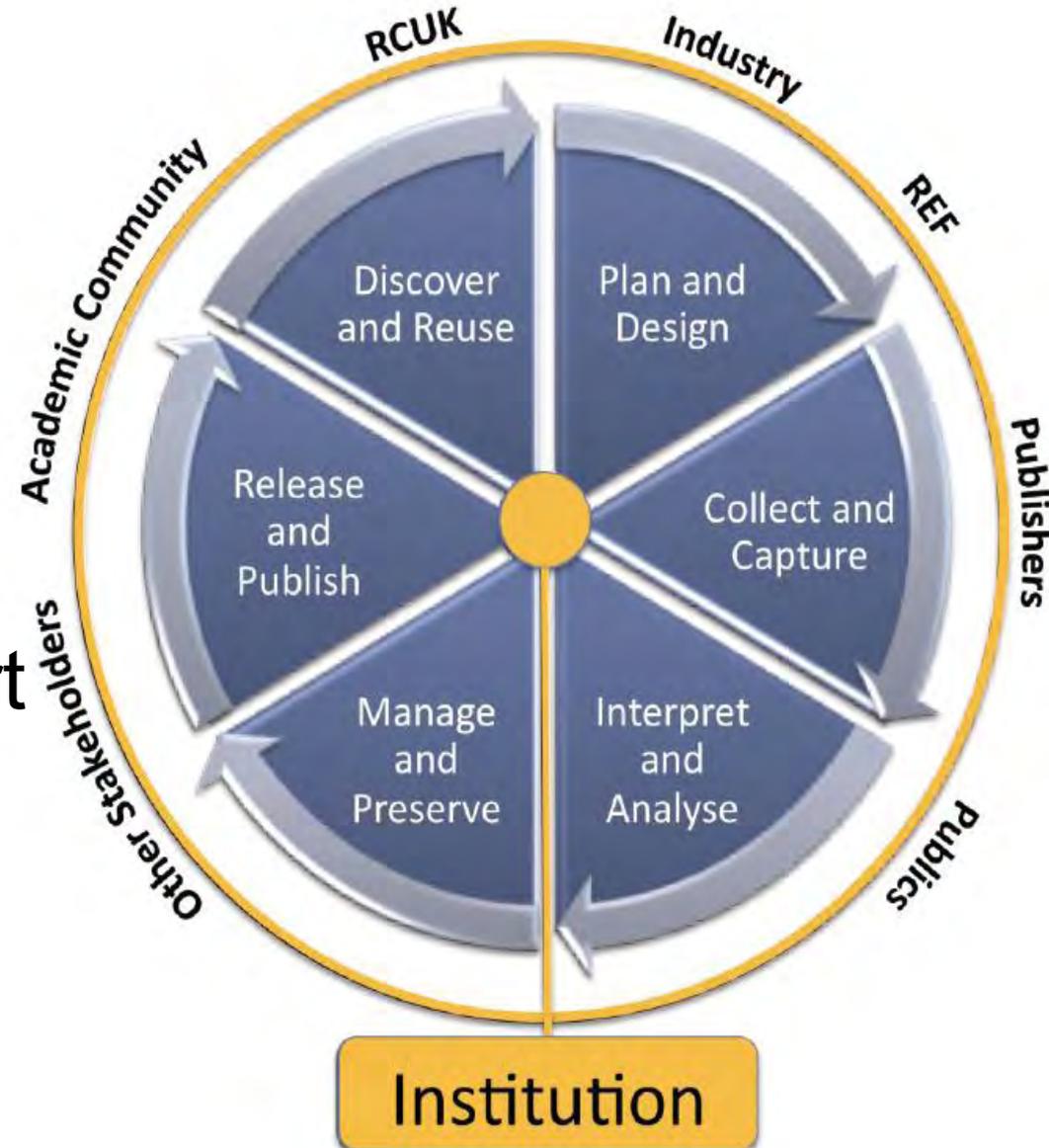
get my metrics!

... or fetch a quick collection based on your [Mendeley contacts](#) and [public groups](#) »

<http://total-impact.org/>

- **Partnership approach**

- UKOLN-DCC
- Library
- IT services
- Research Support Office
- Doctoral Training Centres



# Partnership approach

## Library & institutional stakeholders

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

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

# Data roles

*Liz Lyon, Informatics Transform,  
Ariadne Issue 68, 2012*

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

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

# 3. Developing data informatics capacity & capability

(Acquiring the skills....)



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



Sheila Corral: Libraries,  
Librarians and Data  
Many action exemplars

2012: Libraries in review

<b>Skill gap</b>	<b>2-5 years</b>	<b>Now</b>
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*

# Pause for reflection....



- Skills shortage for data informatics?
- Reposition LIS curriculum?
- LIS entry requirements?
- Get credit for informatics work?

# Play for action....



## 1. Define core components of data informatics

- Visualisation e.g. VisTrails
- Workflow e.g. Taverna
- Analysis e.g. R

**“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.”**

# Play for action....



## 2. Analyse LIS entry qualifications & increase STEM entrants

### Target

- Biologists
- Chemists
- Mathematicians



**ISB** International Society  
for Biocuration



**SCONUL**

Society of College, National and University Libraries



COUNCIL OF AUSTRALIAN UNIVERSITY LIBRARIANS



Chartered Institute of  
Library and Information  
Professionals

*Research Information Network*



# Let's get together

**vala**



Australian  
Library and  
Information  
Association

**RLUK** Research Libraries UK



**ARL**



**ASSOCIATION OF RESEARCH LIBRARIES**

[www.arl.org](http://www.arl.org)

Search ARL



# Play for action....



## 3. International Data Informatics Working Group to explore promotion, recognition & reward

- Global awareness campaign
- Career incentives
- Benchmark good practice

<b>Position</b>	<b>Location</b>
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

# 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*

# Data artist?



Jer Thorp: Hope / Crisis, NYT Word Frequency

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

EMC<sup>2</sup>

## Data Science Revealed community survey

<http://www.emc.com/collateral/about/news/emc-data-science-study-wp.pdf>

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

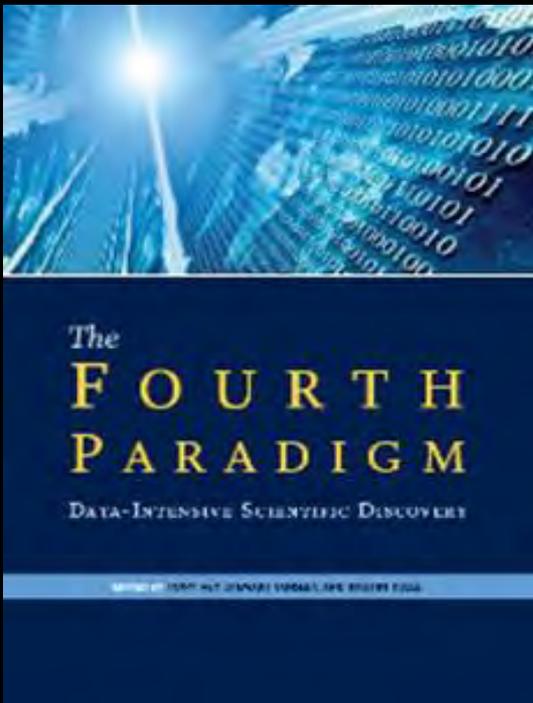


For a University, research data is a key element of “Big Data”.

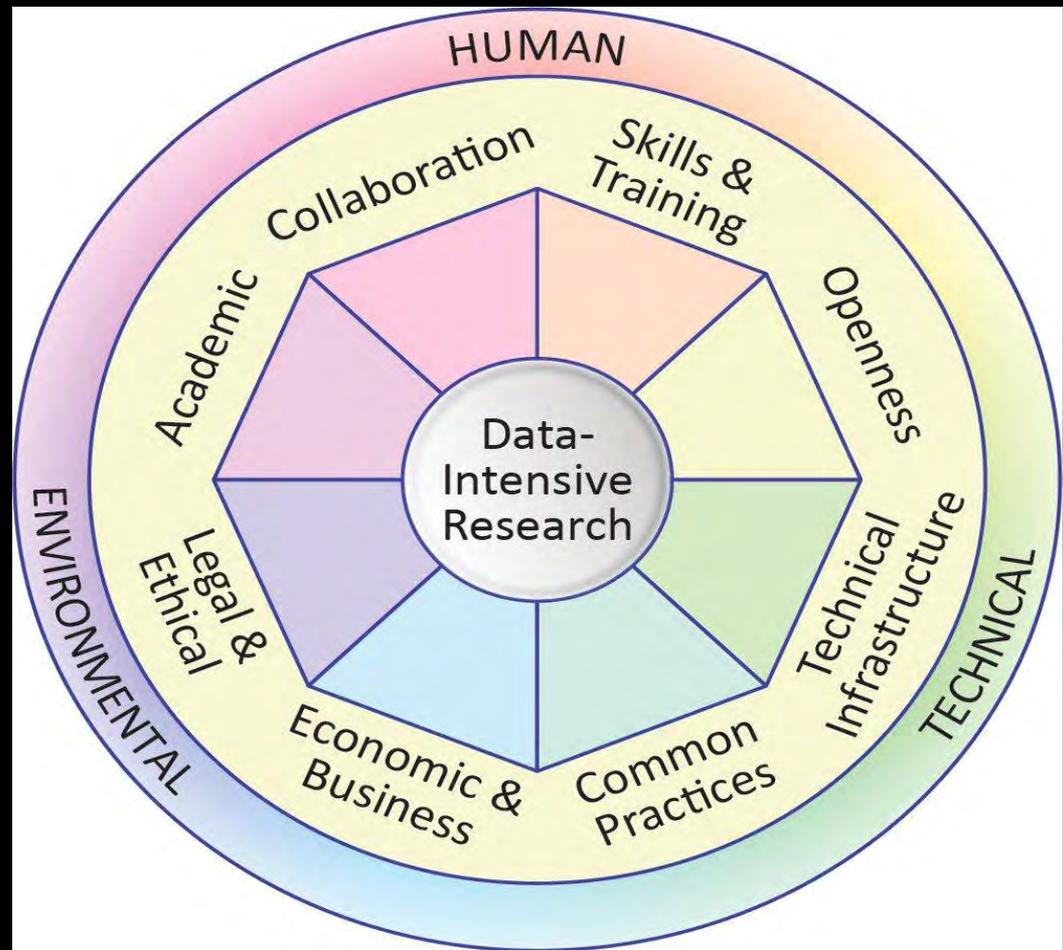
Managing research data effectively will give business advantage.

# Data-intensive research

- Intelligence
- Decision-making
- Planning
- Investment
- Capacity
- Capability



# Community Capability Model Framework CCMF



- Research Funders
- Institutions
- Research leaders/PIs

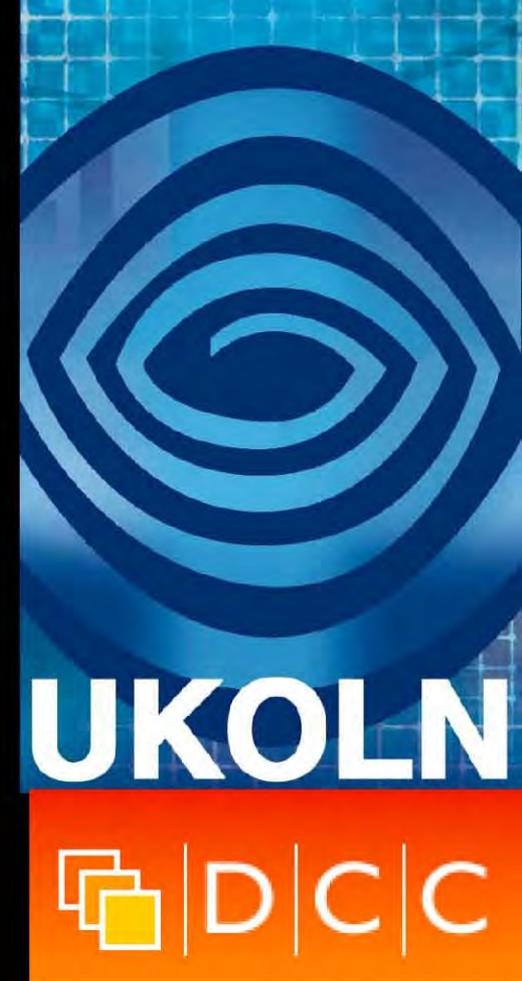
**“The ability to take data -  
to be able to understand it,  
to process it, to extract  
value from it, to visualise  
it, to communicate it -  
that’s going to be a hugely  
important skill in the next  
decades.”**

*Hal Varian, Chief Economist, Google*

**Libraries are on a data journey -  
the Informatics Transform is the  
first step in a new direction...**



# Thank you!



Slides

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

DCC <http://www.dcc.ac.uk>