# From My Data to Our Data
## Creating a culture of data reuse

**Kevin Ashley**
**Digital Curation Centre**
**www.dcc.ac.uk**
**@kevingashley**
**Kevin.ashley@ed.ac.uk**

# My home – the DCC

because good research needs good data

Home | Digital curation | About us | News | Events | Resources | Training | Projects | Community | Tailored support

▷ Mission – to increase capability and capacity for research data services in UK institutions

▷ Not just a UK problem – an international one

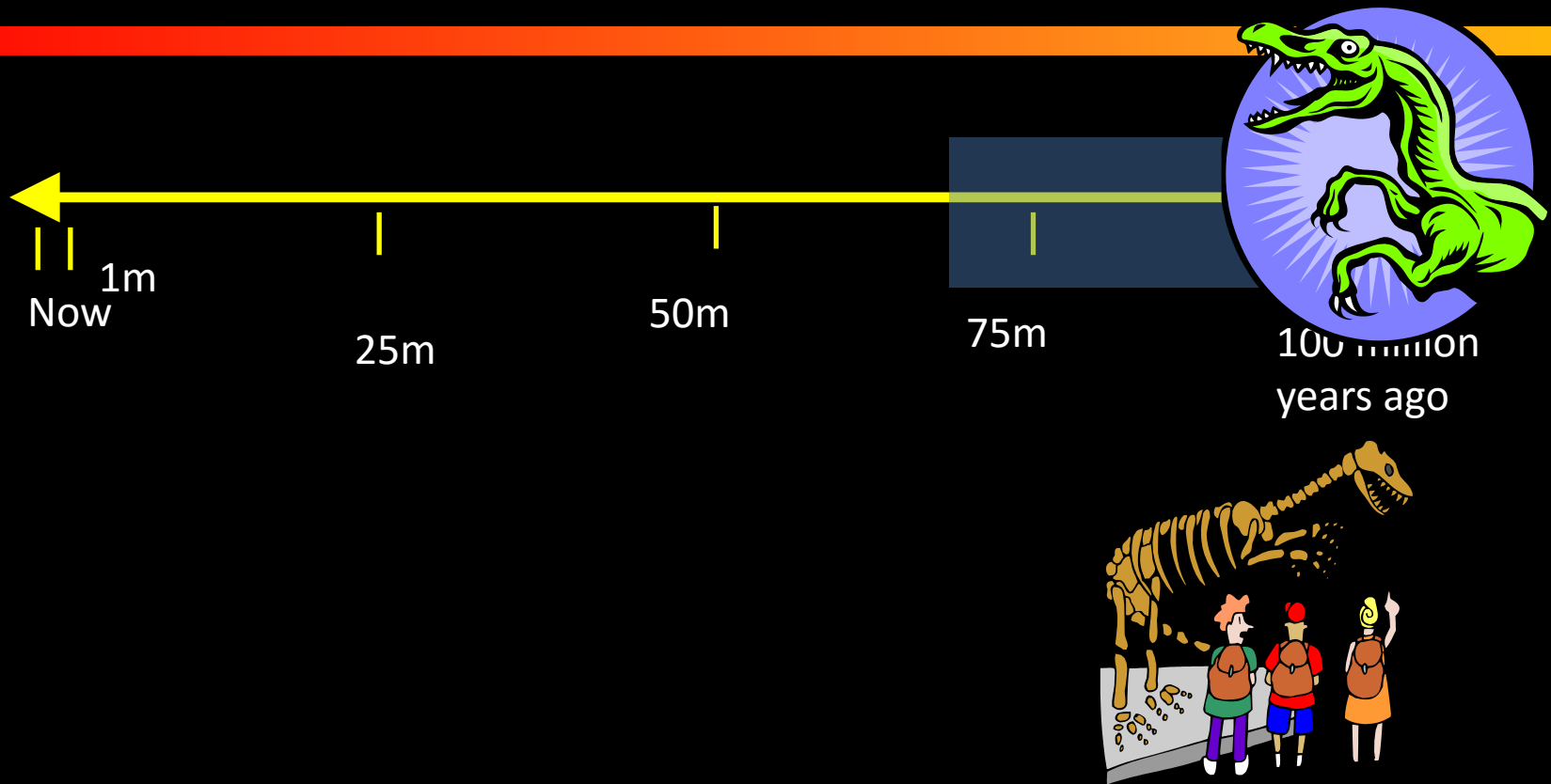▷ Training, shared services, guidance, policy, standards, futures

Latest news | Next events

**IDCC17 - Call for Papers**
22 August, 2016 | in DCC News

**New H2020 DMP guidelines**
12 August, 2016 | in DCC News

**Research data policy briefing welcomes UK Concordat**
5 August, 2016 | in DCC News

12th International Digital Curation Conference
Edinburgh, 20 - 23 February 2017

Call for Papers | Dates | Submissions

How can the DCC help you?

**About us**
We are a world-leading centre of expertise
...p-guidelines information curation...

Editor's choice

**DMPonline & DMPTool roadmap**
Sarah Jones on the recent reciprocal visits between our teams...

Recent blog posts

**Getting our ducks in a row**
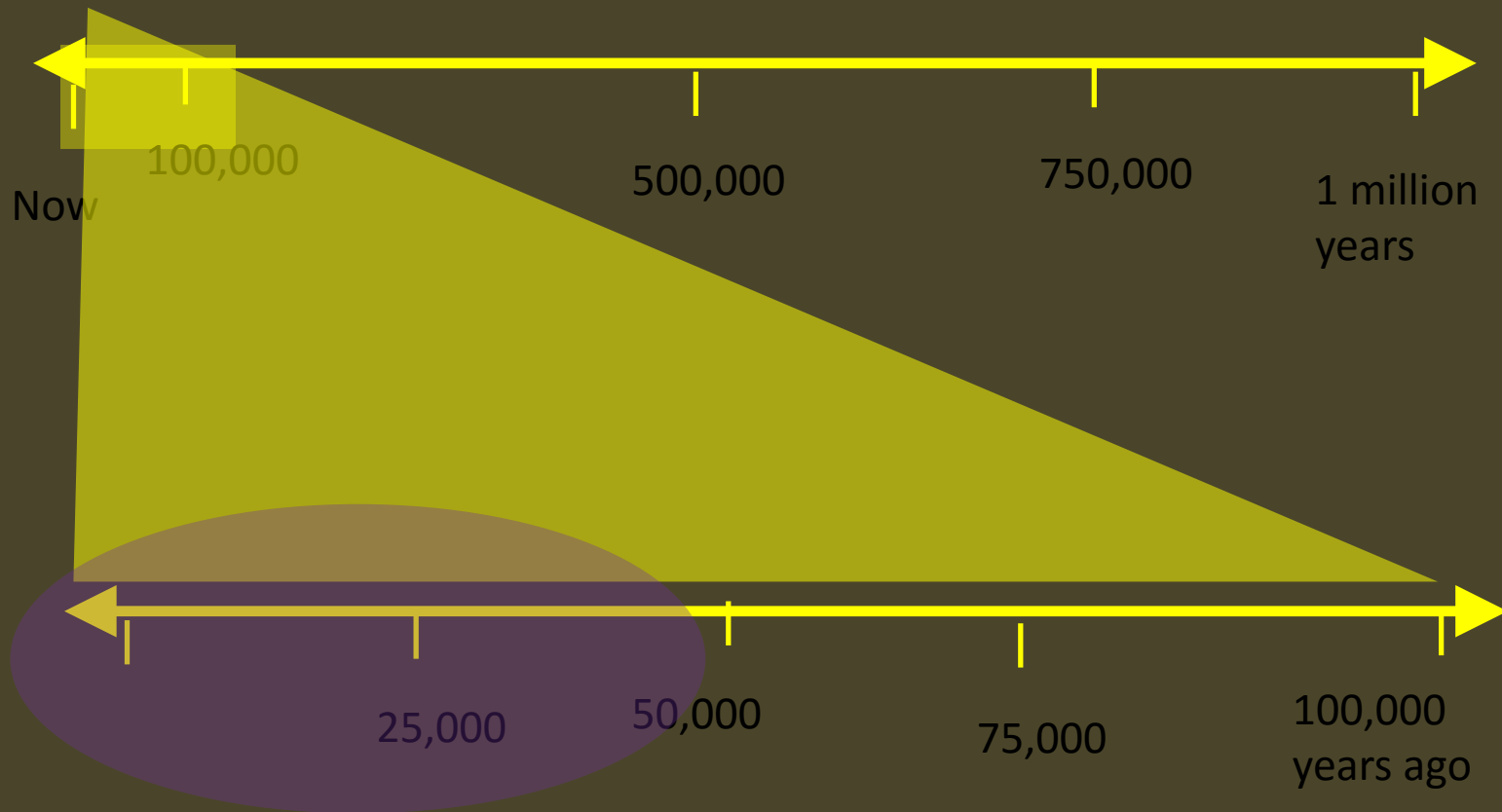
**FOSTER in Scandinavia**

# **Data reuse stories**

▷ The palaeontologist who saved years of work with archaeological data

# What a paleontologist looks at

1m

Now

25m

50m

75m

100 million years ago

Now

100,000

1m

Now

25m

500,000

50m

750,000

75m

1 million years

100 million years ago

# What an archaeologist looks at

100,000

Now

500,000

750,000

1 million
years

25,000

50,000

75,000

100,000
years ago

# Data reuse stories

▷ The palaeontologist who saved years of work with archaeological data

▷ The 19th-century ships logs that help us model climate change

# The Old weather project

Data for research, not from research

# Data reuse stories

▷ The palaeontologist who saved years of work with archaeological data

▷ The 19th-century ships logs that help us model climate change

▷ The 'noise' from research radar that mapped dust from Eyjafjallajökull

# Data reuse - messages

Often your data tells stories that your publications do not

Discipline-bounded data discovery doesn't give us all we need or want

Not all data comes from other researchers

One person's noise is another person's signal

# What is data curation ?

▷ "Maintaining, preserving and adding value to research data throughout its lifecycle"

▷ More than preservation:

» Active management – dealing with change

▷ Less than preservation:

» Lifecycle sometimes involves destruction

▷ Sometimes, not always, about publication or citation

▷ Always about sharing in some way

# What is research data management?

Plan

Create

Discover and Reuse

Use

Deposit and Publish

Appraise

"the active management and appraisal of data over the lifecycle of scholarly and scientific interest"

"an explicit process covering the creation and stewardship of research materials to enable their use for as long as they retain value."

**Data management is part of good research practice**

# Why care?

▷ Data is expensive – an investment

▷ Reuse:

» More research

» Teaching & Learning

» Planning

▷ Impact – with or without publication

▷ Accountability

▷ Legal & regulatory requirements

# Why does this matter?

▷ Research quality
  » How close can we get to the truth?

▷ Research speed
  » How quickly can we get to the truth?

▷ Research finance
  » How much does the truth cost?

▷ Improving one or more of these is of interest to all actors:

▷ Researchers as data creators

▷ Researchers as data reusers

▷ Research institutions

▷ Funders – hence government and society

**Centres like these provide a return on investment of between 400% and 1200%**

http://www.jisc.ac.uk/whatwedo/programmes/di_directions/strategic directions/badc.aspx

# Integrity – not without data



...dies on inte...
...ed 1976; n...

...ng le...
...obab...
...es
...55 p...
» Poldermans -



The case for open data: the Duke clinical trials

9 May, 2011 | in Blogs
By: Kevin Ashley

A recent story in the Times Higher Educational Supplement, backed up by leader comment, provides a highly readable summary of a long and complex case of flawed clinical research and the difficulties encountered by those trying to expose the flaws. It also provides a strong argument for being open with data and code at an early stage, even where sensitive data is involved.

Since this research involved cancer chemotherapy, the lives of people and their quality of life whilst undergoing treatment potentially depended on the truth of the research findings. As the article shows, falsifying the findings would have been far easier and quicker had the original data, and the methods used to analyse it, been made available from the outset. Expensive clinical trials could have been avoided. Potentially, better treatments could have been brought to trial more quickly once the false promise of this particular intervention was clear.

It's often felt that whilst some subjects may be prime candidates for openness with data, those involving human subjects, and in particular clinical medicine, present too many ethical and regulatory challenges. Examples such as this show that such a position is wrong. Even if ethical and regulatory barriers exist, wider ethical issues - the avoidance of unnecessary human suffering being one - demand that we be as open as possible with clinical data. In this case, no identifying information needed to be released to allow others to validate or invalidate this work. Even when the inclusion of identifying information is inescapable, data can still be open in the sense that its existence is public and it is made available to anyone who can satisfy the

Most Read    Site Comments

▶ Re-skilling for Research - observations on an RLUK report
▶ Re-engineering Libraries for the Data Decade
▶ New book: Managing Research Data
▶ 'What's New' Issue 42: February 2012
▶ How can we evaluate data repositories? Pointers from DryadUK

Incremental project

home    why this matters    news    get involved    comments    organisations    about    contact

**It's time all clinical trial results are reported.**
Patients, researchers, pharmacists, doctors and regulators everywhere will benefit from publication of clinical trial results. Wherever you are in the world please sign the petition:
**Thousands of clinical trials have not reported their results; some have not even been registered.**
Information on what was done and what was found in these trials could be lost forever to doctors and researchers, leading to bad treatment decisions, missed opportunities for good medicine, and trials being repeated.
**All trials past and present should be registered, and the full methods and the results reported.**
We call on governments, regulators and research bodies to implement measures to achieve this.

**Sign the petition**

First Name **
Last Name **
Email **
Country
Occupation
I signed this because... (add your comment for the wall here)

"The case for open data: the Duke Clinical Trials "– blog post, Kevin Ashley, http://www.dcc.ac.uk/news/case-open-data-duke-clinical-trials
"Lies, Damned Lies and Research Data: Can Data Sharing Prevent Data Fraud?" – Doorn, Dillo, van Horik,  IJDC 8(1); doi:10.2218/ijdc.v8i1.256

# Why manage research data –The selfish view

- To make research easier!

- To stop yourself drowning in irrelevant stuff

- In case you need the data later

- To avoid accusations of fraud or bad science

- To comply with the law or regulations

- To share data so others can use and learn from it

- To get credit for producing the data

- Because it's a condition of research funding

# Data loss

Digital data are fragile and susceptible to loss for a wide variety of reasons

▷ Natural disaster

▷ Facilities infrastructure failure

▷ Storage failure

▷ Server hardware/software failure

▷ Application software failure

▷ Format obsolescence

▷ Legal encumbrance

▷ Human error

▷ Malicious attack

▷ Loss of staffing competencies

▷ Loss of institutional commitment

▷ Loss of financial stability

▷ Changes in user expectations



Image CC BY-NC-SA 2.0 by Dave Hill
https://www.flickr.com/photos/dmh65
0/4031607067

# Data is variable

▷ Not always textual

▷ Not always tabular

▷ Not always fixed – continual change

▷ Not always clearly authored – think of archival provenance

▷ Not always associated with publication

▷ Often with indistinct boundaries

▷ Multi-dimensional and non-linear

# Data reuse from Hubble



HST Publication Statistics

# New research with old data



A network meta-analysis offers a wider picture than a single traditional meta-analysis

**700 trials of advanced breast cancer treatment**

Quantitative synthesis allowing to combine direct and indirect information and allowing to estimate all possible pair-wise comparisons between treatments

Size of each node proportional to the

- ▷ Synthesis allows new analyses
- ▷ Research that cannot be done with any one of these datasets

# Make data citable

- ▷ Making data available increases citations
- ▷ Everyone – academic, funder, institution – loves citations
- ▷ Want evidence?
  - » Alter, Pienta, Lyle – 240%, social sciences *
  - » Piwowar, Vision – 9% (microarray data)†
  - » Henneken, Accomazzi – 20% (astronomy) #

# Edwin Henneken, Alberto Accomazzi, (2011) Linking to Data - Effect on Citation Rates in Astronomy. http://arxiv.org/abs/1111.3618

* Amy Pienta, George Alter, Jared Lyle, (2010) The Enduring Value of Social Science Research: The Use and Reuse of Primary Research Data. http://hdl.handle.net/2027.42/78307

† Piwowar H, Vision TJ. (2013) Data reuse & the open data citation advantage. PeerJ PrePrints 1:e1v1 http://dx.doi.org/10.7287/peerj.preprints.1v1

# Improve your research impact

**REPRODUCIBLE RESEARCH FOR SCIENTIFIC COMPUTING**

## Code Sharing Is Associated with Research Impact in Image Processing

*In computational sciences such as image processing, publishing usually isn't enough to allow other researchers to verify results. Often, supplementary materials such as source code and measurement data are required. Yet most researchers choose not to make their code available because of the extra time required to prepare it. Are such efforts actually worthwhile, though?*

Vandewalle (2012) DOI: 10.1109/MCSE.2012.63

# Barriers to Data and Code Sharing in Computational Science

Survey of Machine Learning Community, NIPS (Stodden, 2010):

| Code | | Data |
|---|---|---|
| 77% | Time to document and clean up | 54% |
| 52% | Dealing with questions from users | 34% |
| 44% | Not receiving attribution | 42% |
| 40% | Possibility of patents | - |
| 34% | Legal Barriers (ie. copyright) | 41% |
| - | Time to verify release with admin | 38% |
| 30% | Potential loss of future publications | 35% |
| 30% | Competitors may get an advantage | 33% |
| 20% | Web/disk space limitations | 29% |

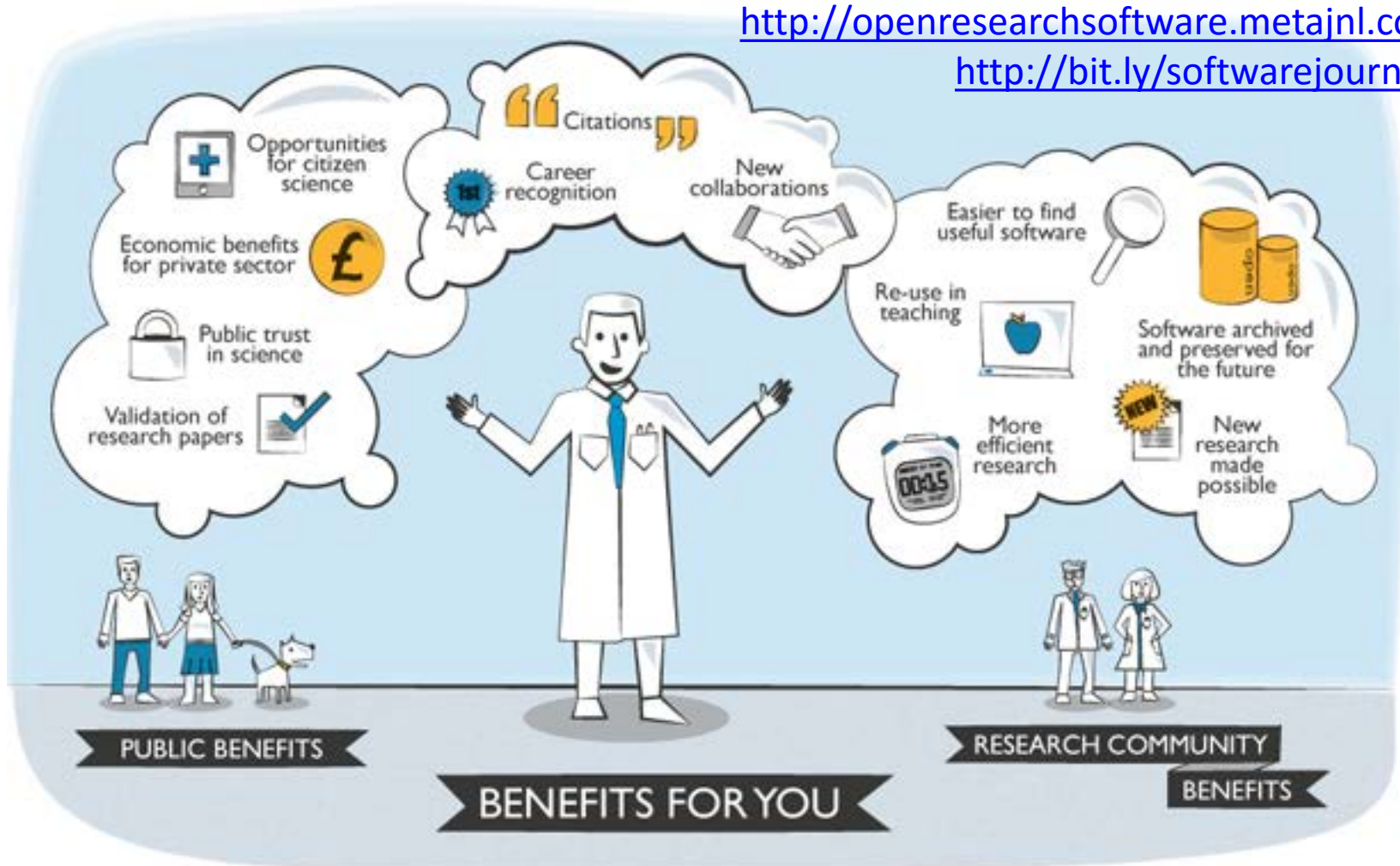Victoria Stodden, AMP 2011  http://www.stodden.net/AMP2011/,

Special Issue Reproducible Research Computing in Science and Engineering July/August 2012, 14(4)

Howison and Herbsleb (2013) "Incentives and Integration In Scientific Software Production" CSCW 2013.

http://openresearchsoftware.metajnl.com
http://bit.ly/softwarejournals

Slide: Neil Chue Hong

# In which j

By Neil Chue Hong.

Until there is a radi
way that academic
principal record of
is still the peer-rev
Given that softwar
part of doing scien
the question we ar
*where can I publis*
*primarily focused*
*software?*

**Engineer**

- Adv
- Coa
- Ren

**Humaniti**

- Dig
- Jou
- Jou

**Image pr**

- Ima
- Insi

## General Journals

- Con
- Con
- Jou
- Jou
- Nat
- Scie
- SIA
- Sof

**Life Scienc**

- Ameri
- Artifici
- Artifici
- Bioinf
- Bioinf
- Bioph
- BMC
- BMC
- BMC
- Bone
- Curre
- Datab
- eLife (
- Epide
- Evolut
- F1000
- Fronti

### Informatics, Mathematics and Statistics

- ACM Transactions on Mathematical Software
- The Archive of Numerical Software
- Future
- Journa
- Journa
- Journa
- Journa
- Journa
- Know
- LMS J
- The M
- Mathe
- Nume
- The R

- BMC Bioinformatics
- BMC Systems Biology
- BMC Source Code for Biology and Medicine
- Bone
- Current Protocols in Bioinformatics
- Database: The Journal of Biological Databases and Curation
- eLife (Tools and Resources) [*example*]
- Epidemiology
- Evolutionary Bioinformatics
- F1000 Research
- Frontiers in Neuroinformatics
- Gigascience
- Methods in Ecology and Evolution
- Nature Methods [*example*]
- Neurocomputing
- Neuroinformatics
- Nucleic Acids Research (special issues)
- PeerJ [*example*]
- PLoS Computational Biology: Software collection
- PLoS ONE
- Trends in Parasitology

### Physical Sciences and Geosciences

- Communications in Computational Physics
- Computer Physics Communications
- Computers and Geosciences
- Geoscientific Model Development
- International Journal of Quantum Chemistry
- Journal of Chemical Theory and Computation
- Journal of Computational Chemistry (special articles - software news and updates)
- Molecular Simulation
- Wiley Interdisciplinary Reviews: Computational Molecular Science (Software Focus) [*example*]

### Acknowledgements

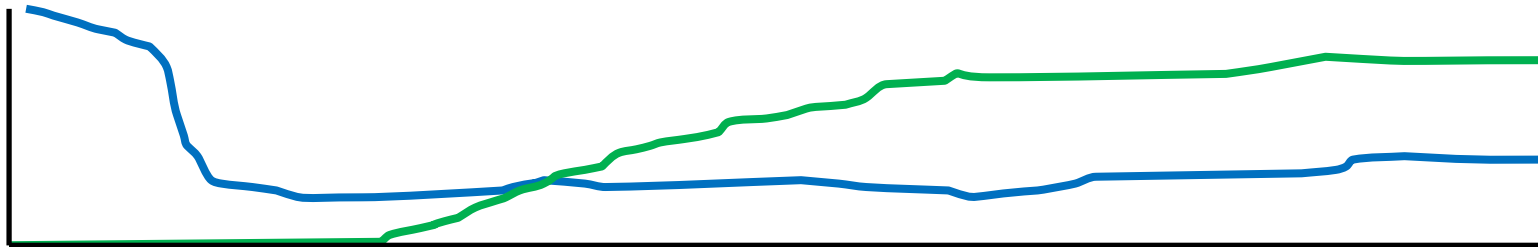# PRACTICE AROUND THE WORLD – AND COSTS

# UKRDS – vision and business case

▷ Assume need for data preservation within each university

▷ Linked by common national services:

» Discovery

» Data Management Planning

» Permanent Identifiers

▷ Working with international data infrastructure

▷ £5m (5.57m Euro) over 5 years – investment then repaid by increased efficiency

# A simplified business case

# Where should data go and when?

▷ Where a national or international subject repository exists – use it

▷ Where there is no repository – the university is responsible

▷ It takes responsibility for data once active use is finished
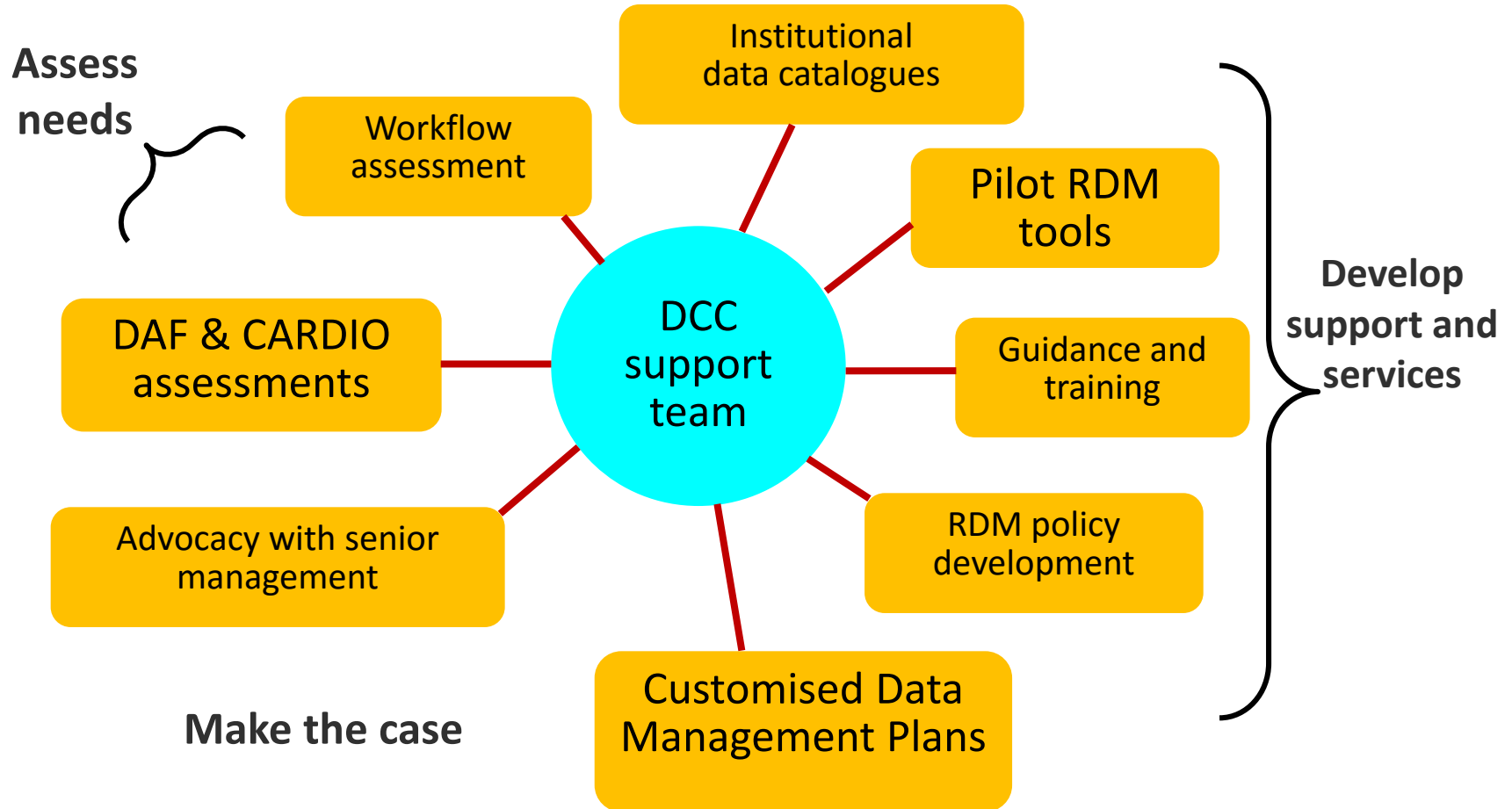
▷ Not all data is kept for ever

# Funder policy helps the change

▷ NERC (environment) – 2010

▷ ESRC (social science) – 1990s

▷ EPSRC (engineering, physical sciences) – 2013

▷ European Commission – 2012

▷ Most require data management plans

▷ Policies influence researchers and their universities

# The source of my lessons

▷ In 2011, DCC began working closely with 20 UK universities to develop research data management (RDM) services

▷ Putting guidance, learning into practice

▷ Since expanded to > 60 universities and other organisations around the world

# DCC 'institutional engagement'



**Assess needs**

Workflow assessment

Institutional data catalogues

Pilot RDM tools

DCC support team

DAF & CARDIO assessments

Guidance and training

Advocacy with senior management

RDM policy development

**Develop support and services**

**Make the case**

Customised Data Management Plans

**...and support policy implementation**

# Some institutional roles

▷ Leadership – coordinate action

▷ Audit – who has what, where does it go?

▷ Advice on access – data, wherever it is

▷ Preservation – permanence

▷ Citability

▷ Data/publication linking

▷ Promoting data in teaching

▷ Selection

▷ Education – early career researchers

# Acquire research data skills

DCC - NTU Seminar - CC-BY

Without senior management attention and researcher involvement, your initiative will fail

# Who (in the UK) is leading RDM work?



**RESEARCHERS**

**Library**

**Research Office**

**IT**

# Research data management services cannot involve the library alone

*"I just back everything up onto data sticks. I didn't even know you could back-up to servers".*

*"Departments don't have guidelines or norms for personal back-up and researcher procedure, knowledge and diligence varies tremendously. Many have experienced moderate to catastrophic data loss"*

Incremental Project Report, June 2010

# Researchers need to know your services exist

# Institutional support



http://www.dcc.ac.uk/resources/how-guides/how-develop-rdm-services

# Goals for the university

▷ Each university can provide a safe home for data that is discoverable

▷ Each university has RDM skills in library, IT, research support

▷ Each university is training new researchers in RDM skills

▷ Where sensible, universities work together to provide services

▷ Each university uses national, international services where appropriate

# Is there a better home for my data?

An international
service –
building on work
by Purdue, DCC,
Biomedcentral
and others

# Data discovery around the world

▷ Research Data Australia

▷ UK data registry pilot & Gateway2Research

▷ Research Data Netherlands

▷ World Data System



Still focussed on discovery of whole datasets – we need to move to discovery of what's inside them

# Australia

▷ Significant long-term funding for national services & support

▷ ANDS – data discovery, software, skills, coordination, advice

▷ National storage & HPC infrastructure

▷ Financial incentives for universities to use common services in standard ways

# Canada

▷ Like Australia, national research funding but province-level university funding

▷ Two initiatives – one from universities (PROTAGE), one from federal level (RDC)

▷ Data management planning, discovery, skills, repositories

▷ Common tools for local deployment

# Netherlands

- ▷ Strong national data repository – DANS

- ▷ One cooperative service – 3TU (now 4TU)

- ▷ Combined to produce RDNL – back office tech services, front office liaison

- ▷ National Dataverse instance

- ▷ Some shared services now being proposed by SURF-SARA

# Portugal

- ▷ Existing national publication repository
- ▷ Extending to cover research data
- ▷ One provider, university-branded front-ends
- ▷ Copying other aspects of UK model – e.g. regular meetings for professional staff, funders, other stakeholders

▷ Scale means many initiatives

▷ Much is NSF project-funded

▷ Some existing university collaborations

▷ Similar spread to UK – 3 or 4 tiers from Ivy League to small institutions

▷ Very complex funding model

▷ Technology is useful – models less so

# What about the cost?

▷ A rough guide – 5% of total project cost on data curation

▷ May not all fall to original research group

▷ How to pay depends on funder and university costing models

▷ Benefits to society & industry are proven

▷ Automation and simplification of many processes is helping

# SOME FINAL MESSAGES

# Should all data be open?

▷ NO

▷ Many reasons – most to do with human subjects

▷ But data existence should always be open

▷ Allows discovery & negotiation on use

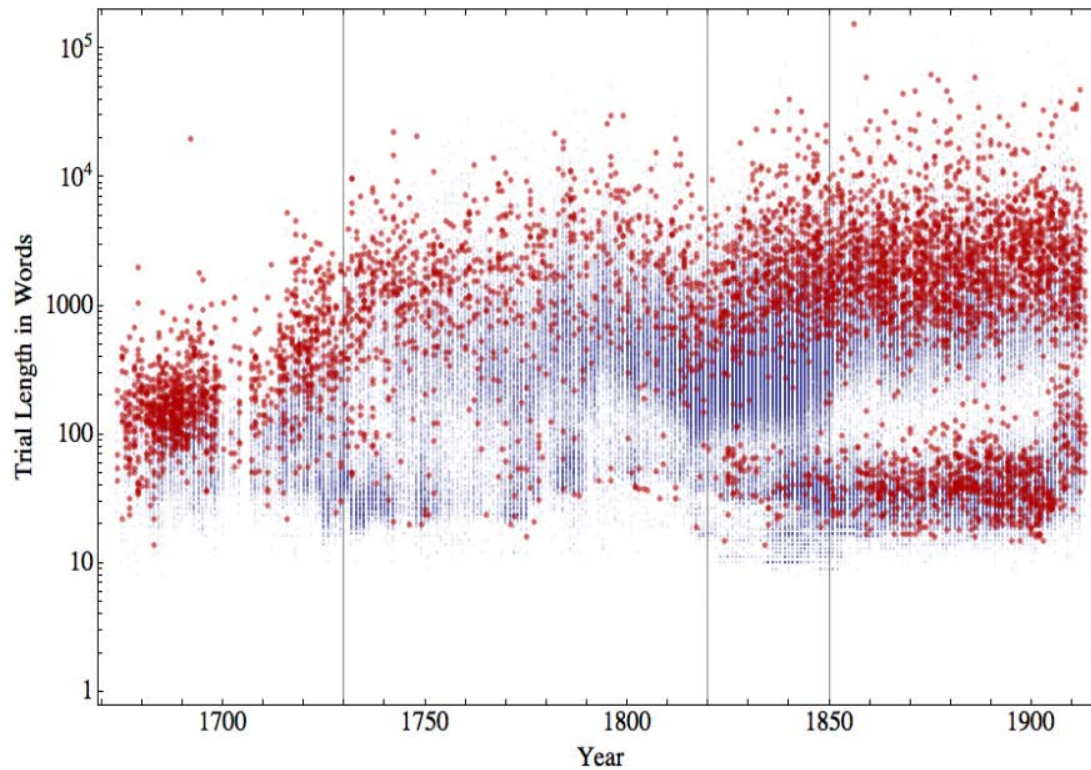▷ Avoids pointless replication

# Data isn't all about numbers

▷ Data can be words, images, sound, video…

▷ Anything which can be analysed to provide insight

▷ Some examples from Old Bailey online – 300 years of English court records

# "loveless marriage"



Credit: "Criminal Intent" – Cyril Briquet, Dan Cohen, Frederick Gibbs, Tim Hitchcock, Jamie McLaughlin, Geoffrey Rockwell, Joerg Sander, Robert Shoemaker, John Simpson, Stefan Sinclair, Sean Takats, William J. Turkel
http://criminalintent.org/

Distribution of trial lengths in words for 'killing' displayed in red; all other trials in grey. 'Killing' includes all trials tagged as including the offences of, 'Infanticide', 'murder', 'petty treason', 'manslaughter', and 'killing: other', by the Old Bailey online.



## Length of trial – killing – from Old Bailey Online

Credit: "Criminal Intent" – Cyril Briquet, Dan Cohen, Frederick Gibbs, Tim Hitchcock, Jamie McLaughlin, Geoffrey Rockwell, Joerg Sander, Robert Shoemaker, John Simpson, Stefan Sinclair, Sean Takats, William J. Turkel

http://criminalintent.org/

# Traditional skills can win



**Data made available before paper was published – result was immediate impact**

...ets it wr...

...s us why

...y, Ryan; King...

...: The Parabl...

Analysis", http://dx.doi.org/10.7...

UNF:5:BJh9WzZQNEeSEpV3EWs...

[Distributor] V1 [Version]

Gking.harvard.edu/data

Tools to make data findable & reusable

Researcher-friendly: incremental approach to metadata

# Some messages for you

▷ Some things we need to know about data:
- » When/where/what is it about?
- » Who owns it
- » What rights apply
- » What it is derived from & how
- » What software may be associated
- » What data management plan applies
- » How do I gain access ?
- » Where is it ?
- » When was/will it be destroyed?

# My messages to researchers

▷ Sharing is difficult

▷ Reusing is difficult

▷ Both are key to advancing science, and advancing your own career

▷ Your data can live longer than your findings

▷ All this can be easier than you think

# The value of data in astronomy

▷ Zij star catalogues – 8$^{th}$ century onwards
▷ Abd Al-Rahman Al-Sufi – book of fixed stars
▷ Abu-Mahmud al-Khujandi
  » meridian transits of sun
  » calculate earth's angle of tilt
  » Different to earlier Indian/Greek calculations (but data lost)
▷ Lots of work refinin                motion
▷ Kepler –  data from
▷ Modern day – chine         us measure changes

**The old theories are discredited**

**The old data has value**

**Darwin had something to say about this**

False facts are highly injurious to the progress of science, for they often endure long; but false views, if supported by some evidence, do little harm, for every one takes a salutary pleasure in proving their falseness

# Our message to researchers

▷ The credit belongs to you

▷ The data belongs to all of us

▷ Share, and we all reap the benefits



Store your valuable data.
Show it to the world.
Share it with others.

http://datacentrum.3tu.nl

# Excuses – and responses

▷ "People will ask questions"
  » So use a data centre or repository
▷ "It will be misinterpreted"
  » Stuff happens. Also, openness encourages correction
▷ "It's not interesting"
  » Let others be the judge – your noise is my signal
▷ "I might get another paper out of it"
  » Up to a point. We might get more research out of it
▷ "I don't have permission"
  » A real problem. But solvable at senior level
▷ "It's too bad/complicated" –see above
▷ "It's not a priority"
  » Unfortunately, funders are making it so. But if you looked at the evidence, it would be your priority as well

See e.g. Carly Strasser's blog:
http://datapub.cdlib.org/2013/04/24/closed-data-excuses-excuses/