



**Building an
Open Data Infrastructure for Research:
*Turning Policy into Practice***

Juan Bicarregui

Head of Data Services Division
STFC Department of Scientific Computing

Overview

1. The Policy Context

- OECD
- EC/NSF/...
- G8+5
- RCUK
- Royal Society
- G8

2. PaNdata

- Photon and Neutron Open Data Infrastructure

3. The Research Data Alliance

- Fostering Collaboration on a global scale

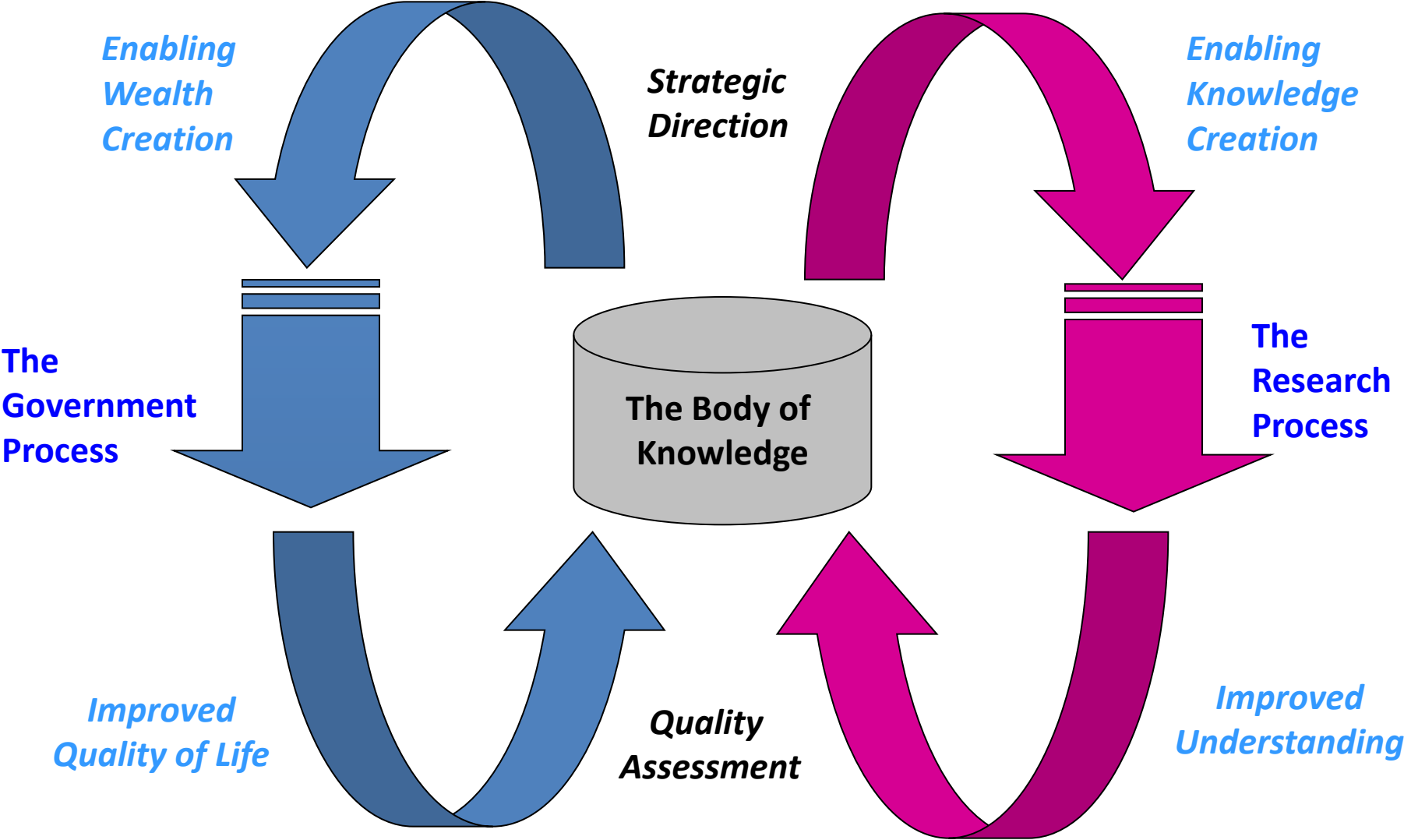
1. The Policy Context

- OECD, 2004-2006
 - Principles and *Guidelines for Access to Research Data from Public Funding*
- EC, 2007-2012
 - Recommendation on access to and preservation of scientific information
- G8+5, 2011-2012
 - Global Research Infrastructure Sub Group on Data
- Research Councils UK, 2011
 - Joint Principles on Data
- Royal Society, 2011-2012
 - Science as an Open Exercise
- G8 Ministerial Statement, 2013
 - Grand Challenges, Global Research Infrastructures,
 - Open Scientific Research Data, Open Access

The views expressed herein are the personal views of the author and do not necessarily reflect the views of the policy makers

Economic Impact

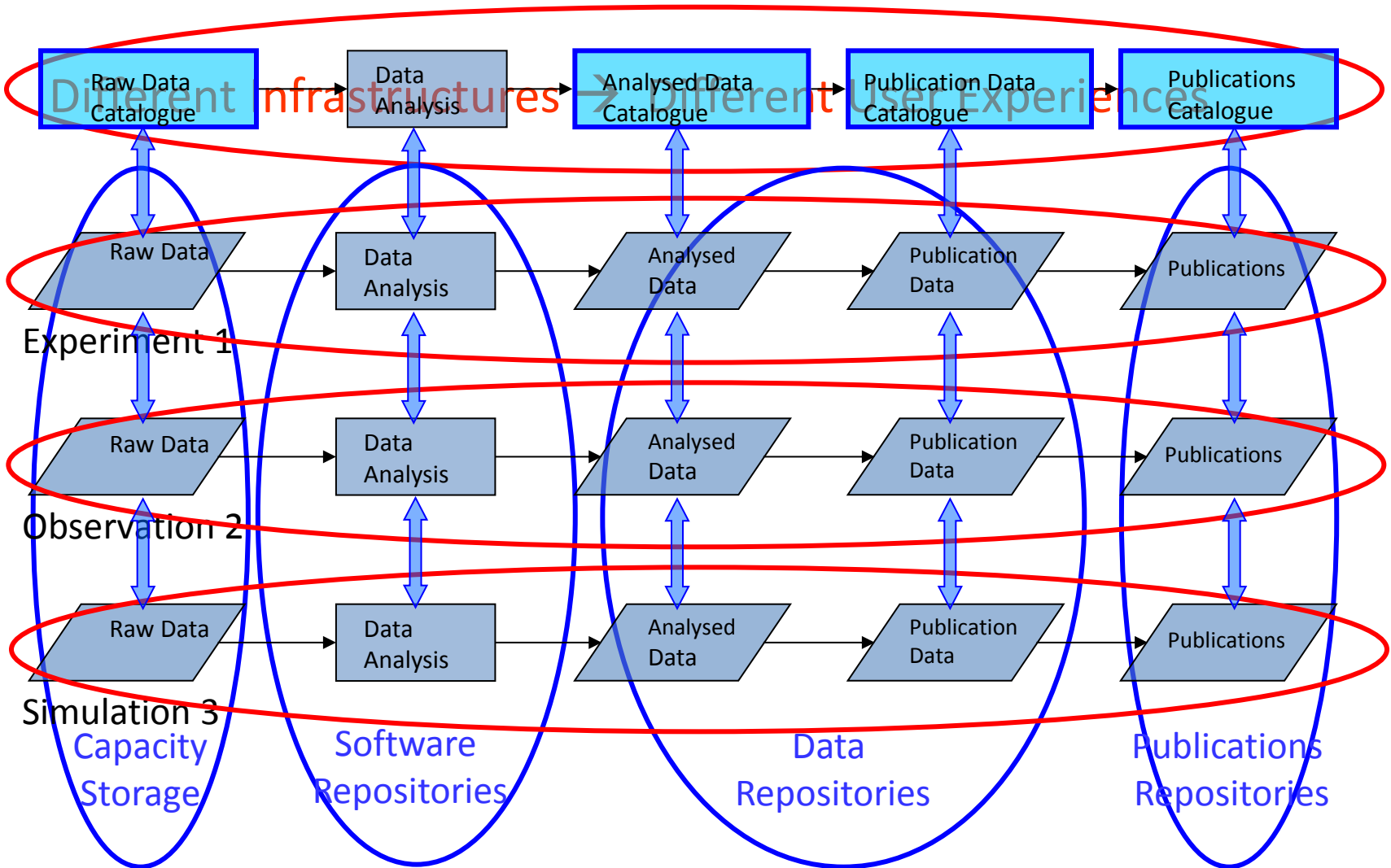
The Innovation Lifecycle



Aggregation of Knowledge lies at the heart of the innovation lifecycle

Technology Sharing

Single Infrastructure → Single User Experience



Open Science

the researcher acts through ingest and access

Research Environment

Archival

Provenanced Research

the researcher shouldn't have to worry about the information infrastructure

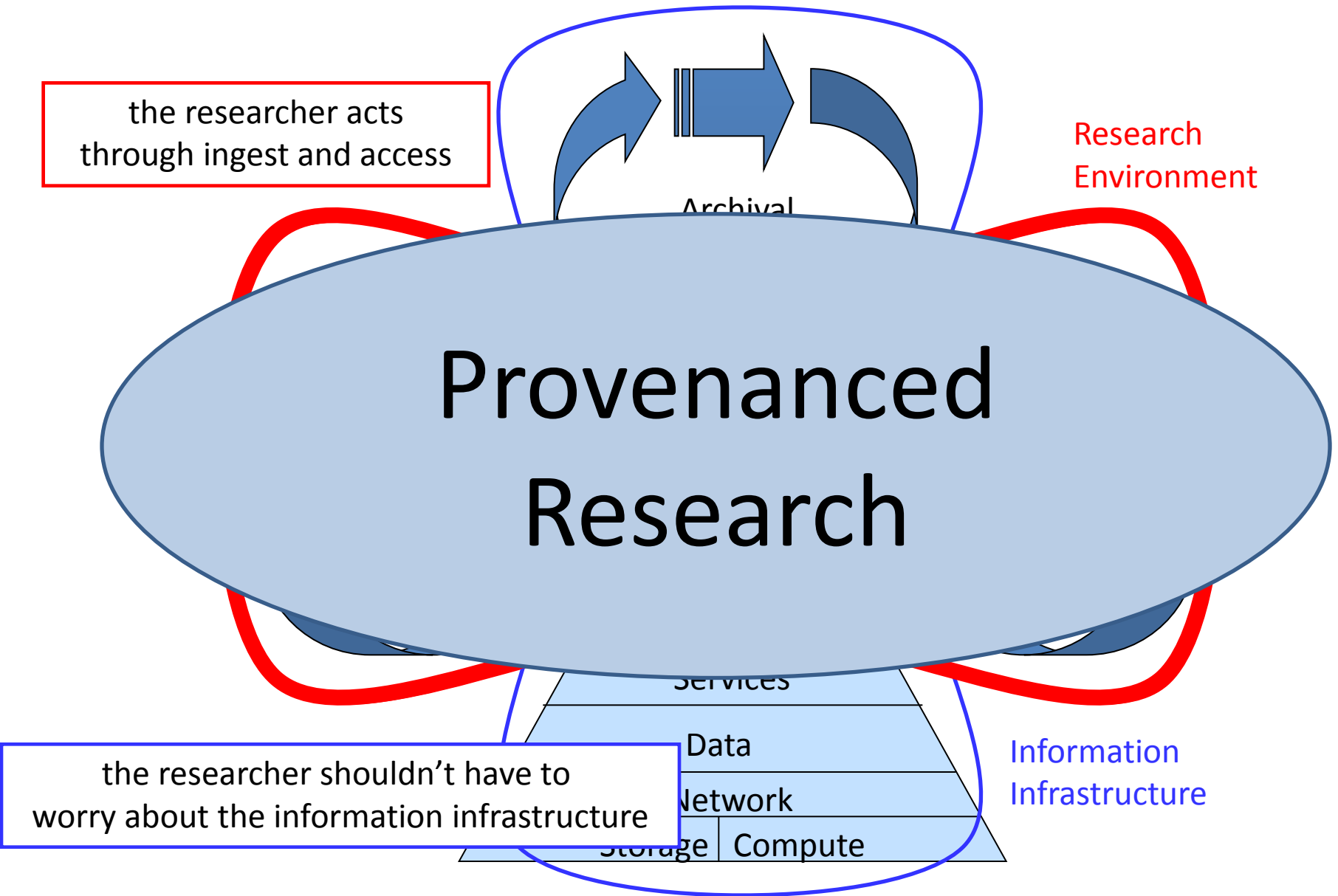
Information Infrastructure

Services

Data

Network

Storage Compute



RCUK principles: Data are a Public Good

Publicly funded research data are a public good, produced in the public interest, which should be made openly available with as few restrictions as possible in a timely and responsible manner that does not harm intellectual property.

Public good – is [nonrival](#) and [non-excludable](#) [wikipedia]

consumption by one does not reduce availability for others

no one can be effectively excluded from using

Research Data

recorded factual material commonly retained by and accepted in the scientific community as necessary to validate research findings

As few restrictions as possible

Later (distinguish registration from restriction)

Timely

Later (discipline specific)

Responsible

Later (maximising access does not necessarily maximising research benefit)

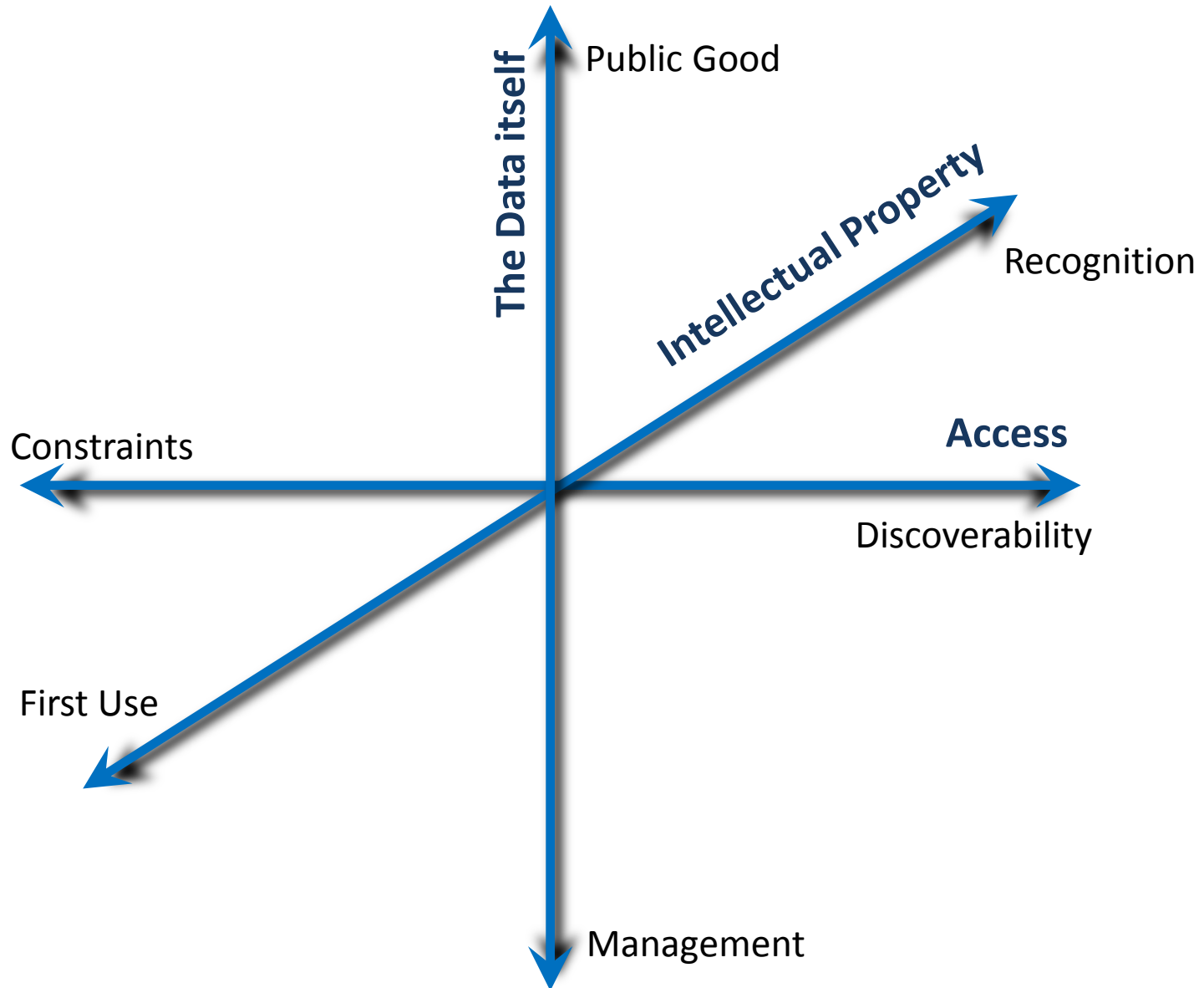
Intellectual Property

Later (balance contribution from sharing and from primary research)

RCUK Principles on Data Policy

- 2) Data should be managed
- 3) Data should be discoverable
- 4) There may be constraints
- 5) Originators may have first use
- 6) Reusers have responsibilities
- 7) Data sharing is not free

3 Dimensions of policy





Overview

1. The Policy Context

- OECD
- EC/NSF/...
- G8+5
- RCUK
- Royal Society
- G8

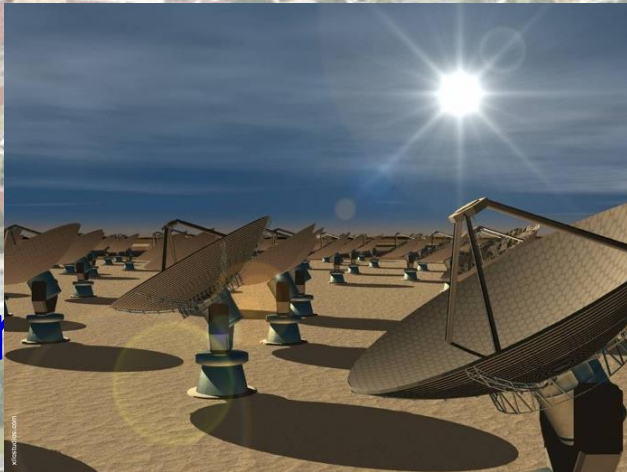
2. PaNdata

- Photon and Neutron Open Data Infrastructure

3. The Research Data Alliance

- Fostering Collaboration on a global scale

What is STFC?



Square Kilometre Array

- Synchrotron Radiation Source
- Lasers
- Space Science
- Particle Physics



Large Hadron Collider



Daresbury Laboratory



Data Management

Communications



ESRF & ILL, Grenoble

The PaNdata Collaboration

- Established 2007 with 4 partners
- Expanded since to 13 organisations
(see next slide)
- Aims:
 - *“...to construct and operate a shared data infrastructure for Neutron and Photon laboratories...”*

2007	2008	2009	2010	2011	2012	2013	2014
	EDNS (4)						
		EDNP (10)					
				PaNdataEurope(11)			
						Pandata ODI(11)	

PaN-data Partners

PaN-data bring together 13 major European Research Infrastructures

ISIS is the world's leading pulsed spallation neutron source

ILL operates the most intense slow neutron source in the world

PSI operates the Swiss Light Source, SLS, and Neutron Spallation Source, SINQ, and is developing the SwissFEL Free Electron Laser

HZB operates the BER II research reactor the BESSY II synchrotron

CEA/LLB operates neutron scattering spectrometers from the Orphée fission reactor

JCNS Juelich Centre for Neutron Science

ESRF is a third generation synchrotron light source jointly funded by 19 European countries

Diamond is new 3rd generation synchrotron funded by the UK and the Wellcome Trust

DESY operates two synchrotrons, Doris III and Petra III, and the FLASH free electron laser

Soleil is a 2.75 GeV synchrotron radiation facility in operation since 2007

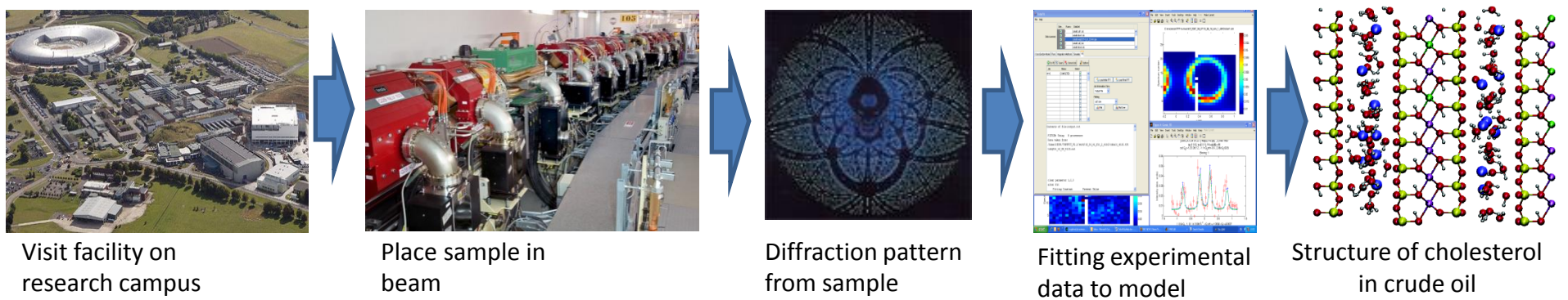
ELETTRA operates a 2-2.4 GeV synchrotron and is building the FERMI Free Electron Laser

ALBA is a new 3 GeV synchrotron facility due to become operational in 2010

MaxLab, Max IV Synchrotron

PaN-data is coordinated by the STFC Department of Scientific Computing

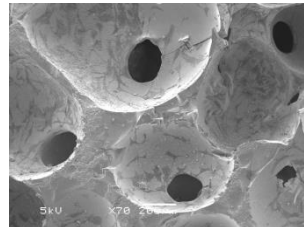
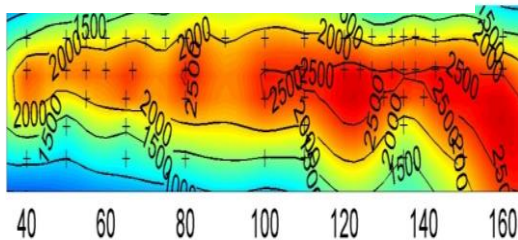
The Science we do - Structure of materials



- Over 30,000 user visitors each year:
 - physics, chemistry, biology, medicine,
 - energy, environmental, materials, culture
 - pharmaceuticals, petrochemicals, microelectronics

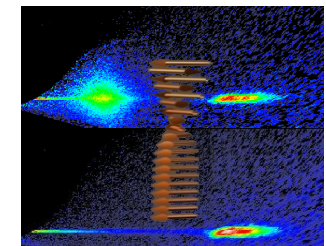
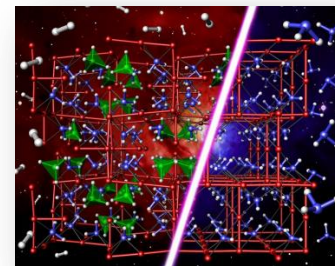
- Over 5,000 high impact publications per year
 - But so far no integrated data repositories
 - Lacking sustainability & traceability

Longitudinal strain in aircraft wing



Bioactive glass for bone growth

Hydrogen storage for zero emission vehicles



Magnetic moments in electronic storage

PaN-data Standardisation

PaN-data Europe is undertaking 5 standardisation activities:

1. Development of a **common data policy** framework
2. Agreement on protocols for shared **user information exchange**
3. Definition of standards for common **scientific data formats**
4. Strategy for the interoperation of **data analysis software** enabling the most appropriate software to be used independently of where the data is collected
5. **Integration and cross-linking** of research outputs completing the lifecycle of research, linking all information underpinning publications, and supporting the long-term preservation of the research outputs



Standards from PaNdata Support Action

users

data

s/w

Integ

uCat

dCat

PaNdata ODI Service Activities

vLabs

PaNdata ODI Service Releases

Rel 1

Jun 2013

Rel 2

Sep 2013

Rel 3

Dec 2013

Rel 4

Mar 2014

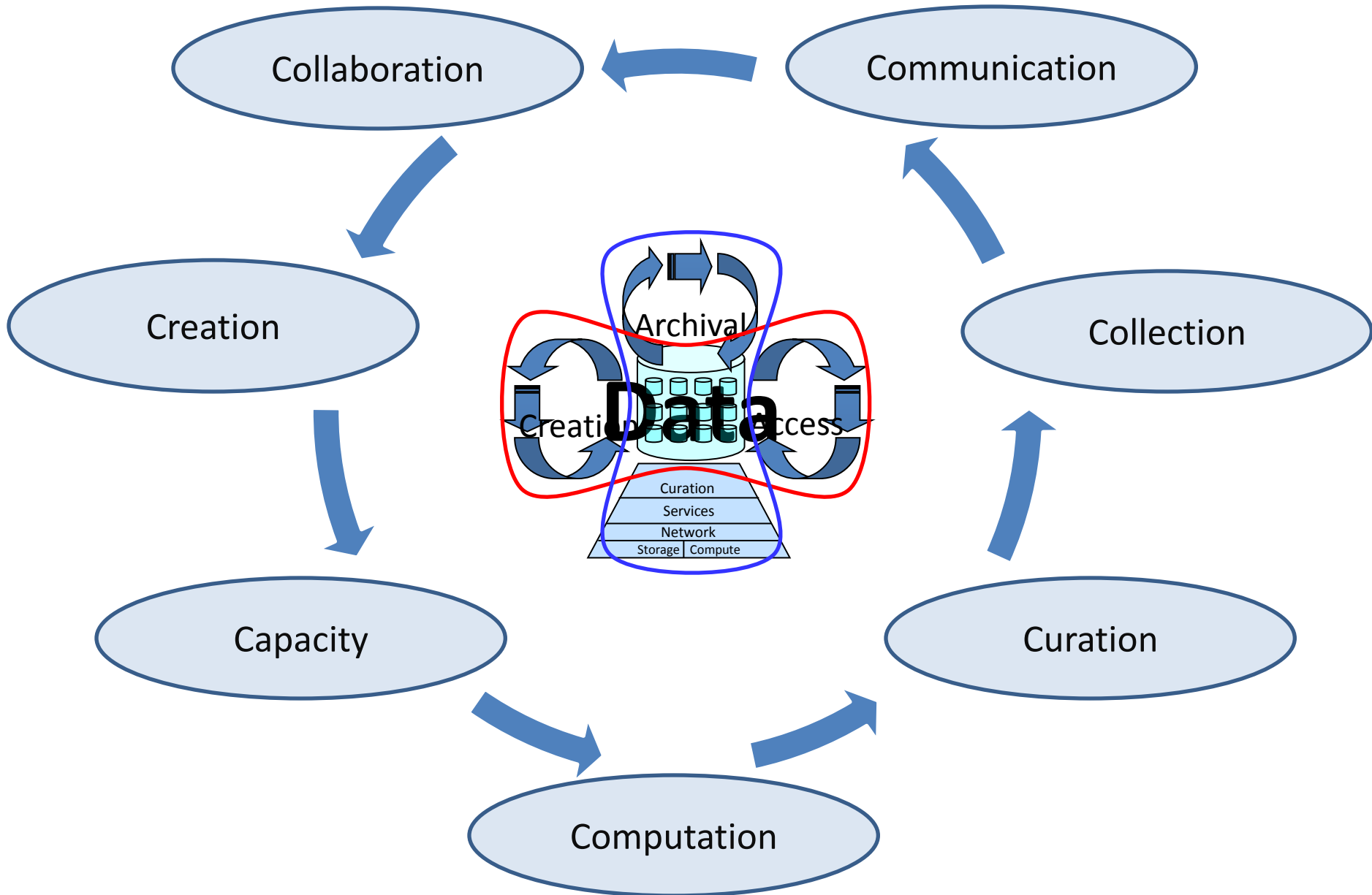
Prov

Pres

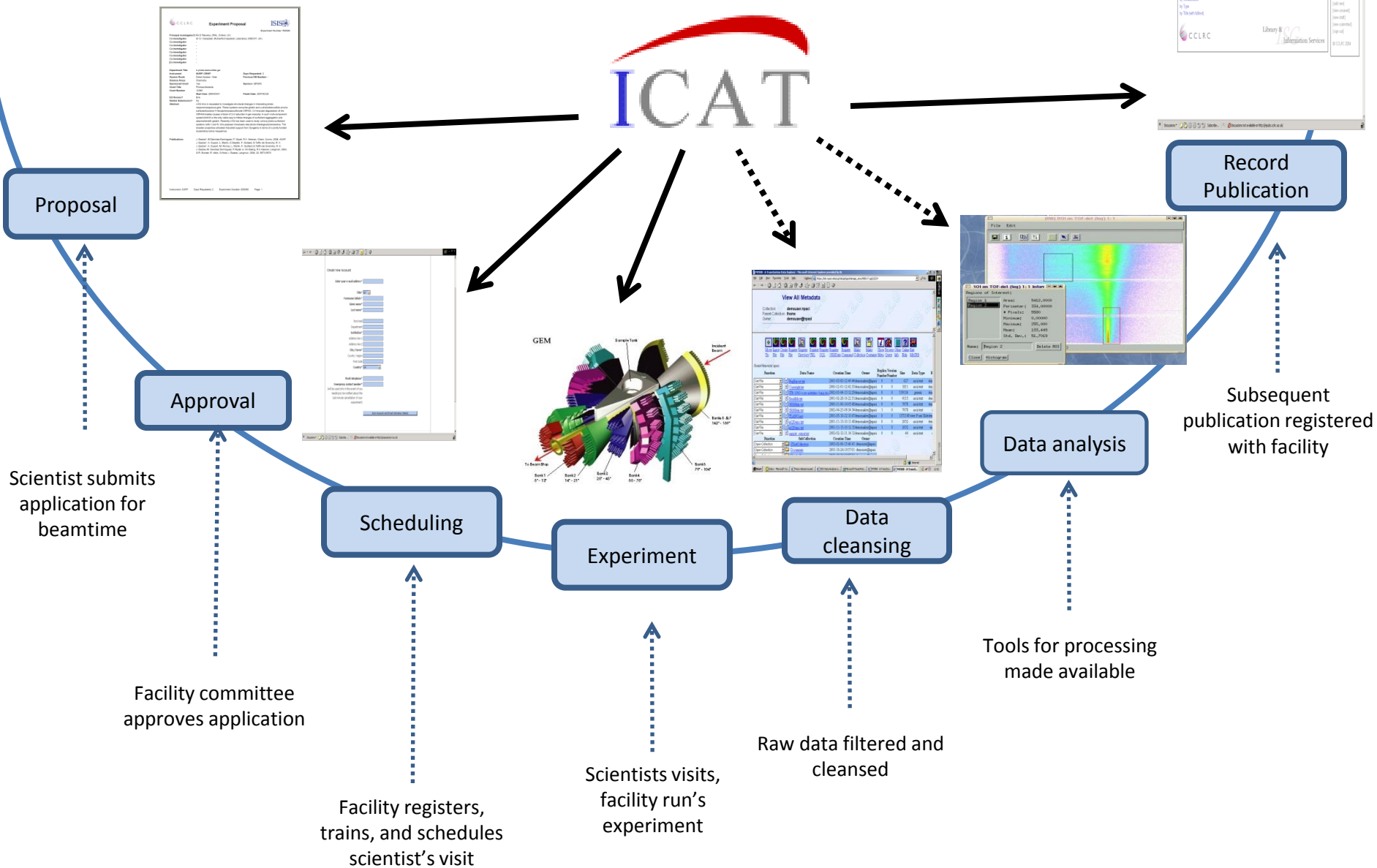
Scale

PaNdata ODI Joint Research Activities

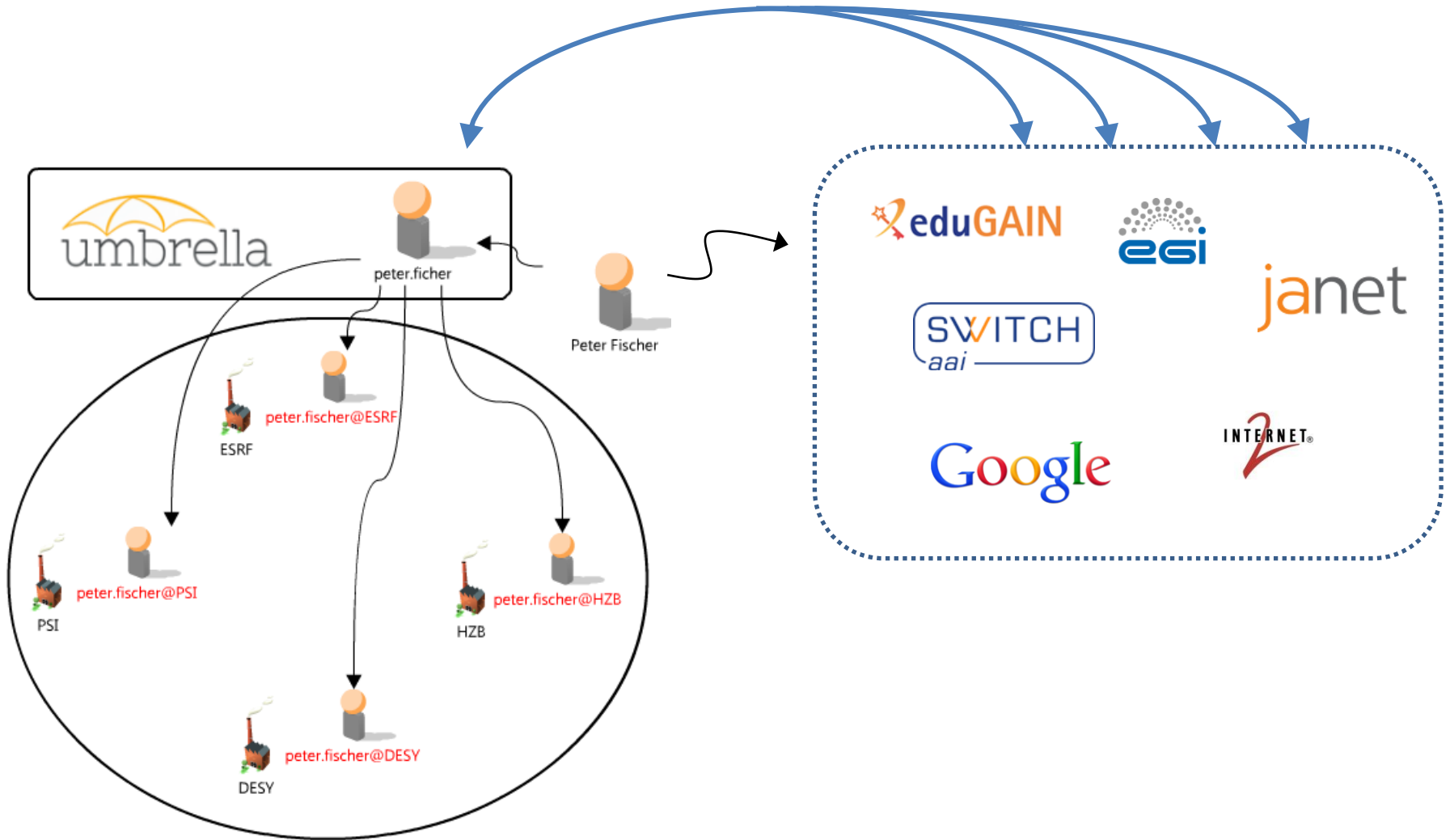
The 7 C's



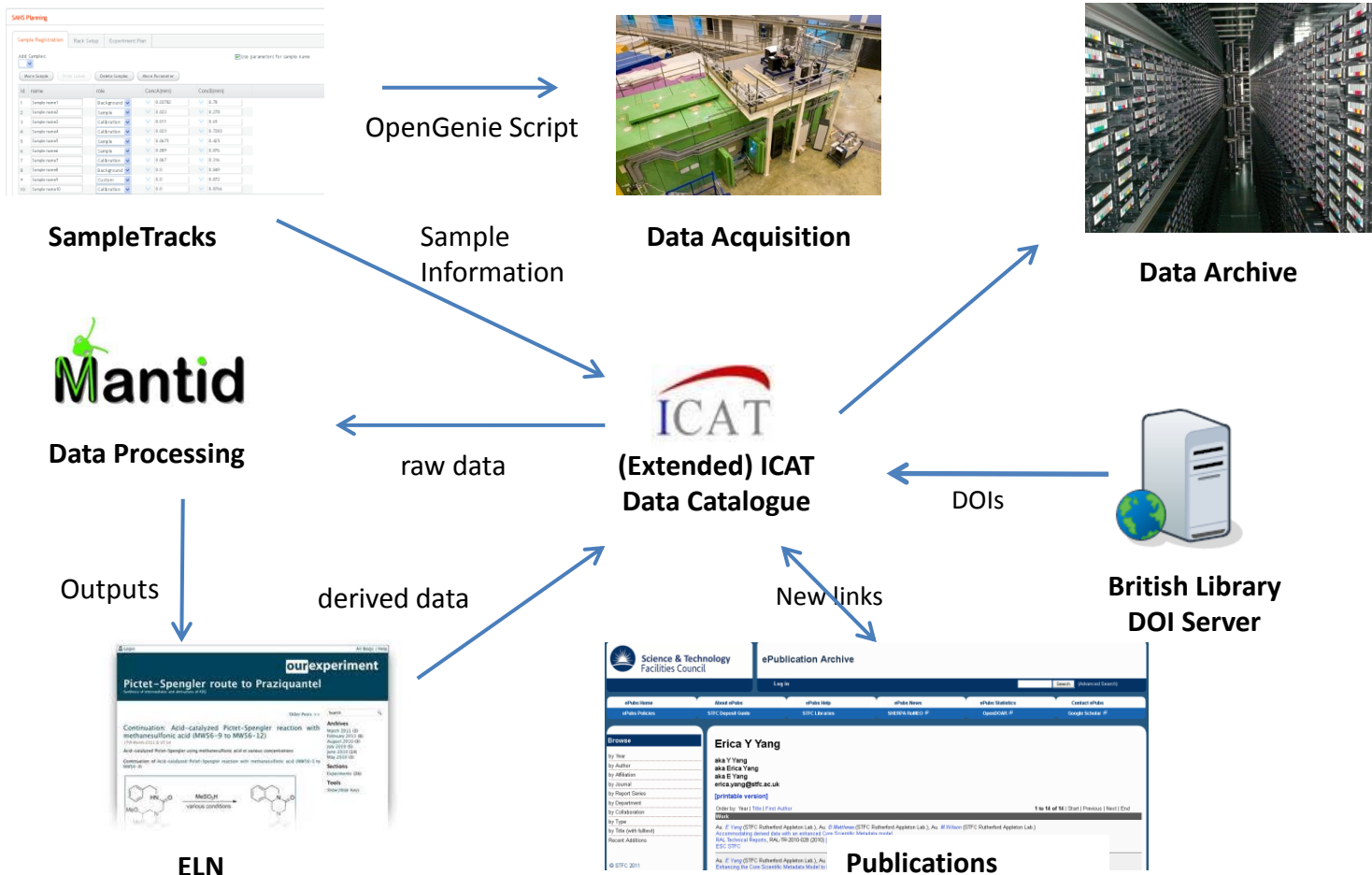
Metadata Collection



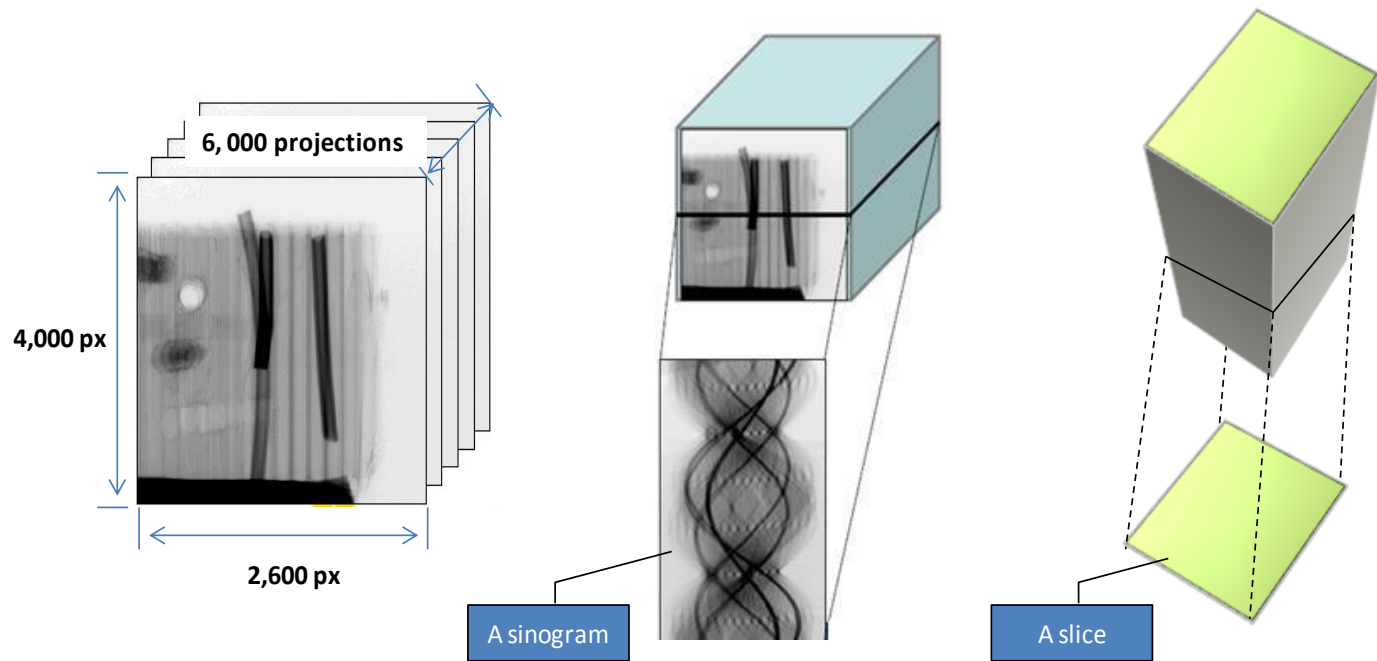
Authentication



Provenance: SANS2d: Experiment coordination



Tomographic Reconstruction



Raw images

(2,600 X 4,000 px each image/projection)
(Total: 6,000 images/projections)

Sinograms
from the projections

(Total: 4,000 sinograms
Each sinogram: 2,600 x 4,000 x 1 px)

Reconstructed slices
in a 3d volume

(Total: 4,000 slices)



~100Gb per 3D image - ~40 mins on 16 GPU cluster
~10 TB per experiment" - ~3 days on site
~ 1PB per year (per beamline)

Working on using the Emerald (376 GPUs)

ESRF example: Amber inclusion

[Prioriphora schroederhohenwarthi](#)

- Xray imaging of 1mm *Prioriphora* (scuttle fly) from Cretaceous period
- found at [Archingeay-Les Nouillers](#) in opaque amber



Overview

1. The Policy Context

- OECD
- EC/NSF/...
- G8+5
- RCUK
- Royal Society
- G8

2. PaNdata

- Photon and Neutron Open Data Infrastructure

3. The Research Data Alliance

- Fostering Collaboration on a global scale

3. The Research Data Alliance

New international organization

Currently supported by:

EU

NSF

Australian National Data Service

To accelerate data-driven innovation
through research data sharing and exchange.

Infrastructure, Policy, Practice and Standards

Research Data Alliance

Vision and Purpose

Vision

*Researchers around the world
sharing and using research data without barriers.*

Purpose

*... to accelerate international
data-driven innovation and discovery
by facilitating research data
sharing and exchange,
use and re-use,
standards harmonization, and discoverability.
...through the development and adoption of
infrastructure, policy, practice, standards, and other
deliverables.*

RDA Principles

Openness

- Membership is open to all interested organizations,
- all meetings are public,
- RDA processes are transparent, and
- all RDA products are freely available to the public;

Consensus

- The RDA moves forward by achieving consensus and
- resolves disagreements through appropriate voting mechanisms;

Balance

- The RDA is organized on the principle of balanced representation for individual organizations and stakeholder communities;

Harmonization

- The RDA works to achieve harmonization across standards, policies, technologies, tools, and other data infrastructure elements;

Voluntary

- The RDA is not a government organization or regulatory body and, instead, is a public body responsive to its members; and

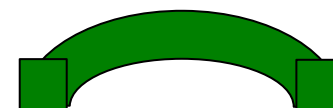
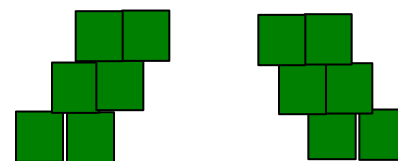
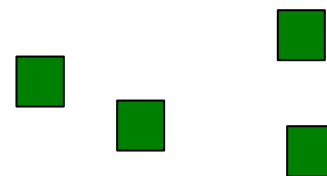
Non-profit

- RDA is not a commercial organization and will not design, promote, endorse, or sell commercial products, technologies, or services.

“Building Bridges”



- Bridges to the future
 - data preservation
- Bridges to research partners
- Bridges across disciplines
- Bridges across regions
- Bridges to integration
 - to solve new problems
- Bridges across communities



RDA role

Two bridges we can build:

- Connecting Data
- Connecting People

What kind of organisation do we need to do this?

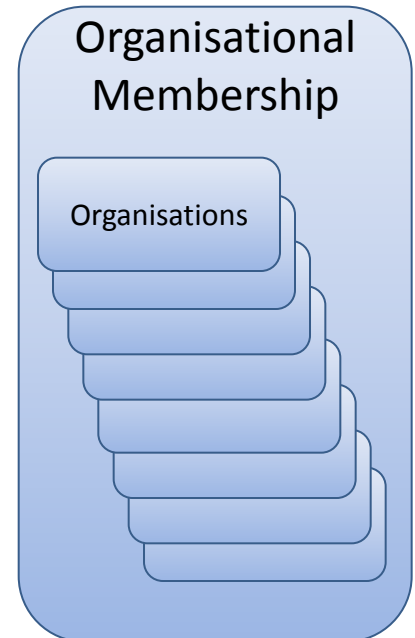
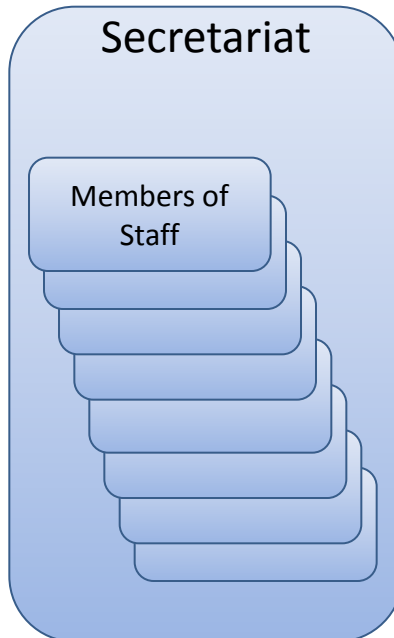
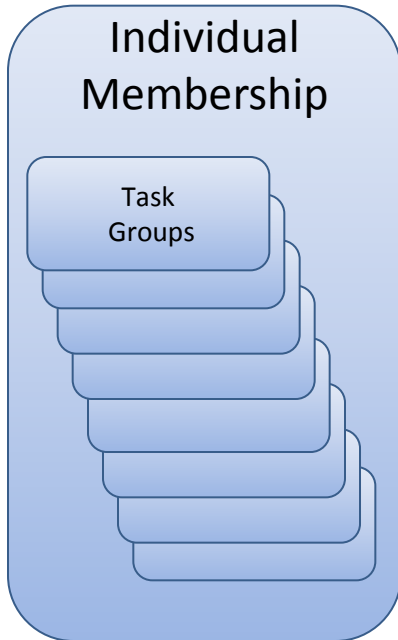
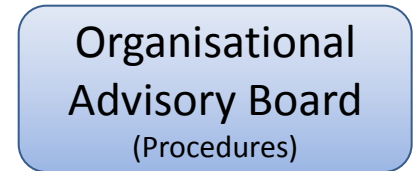
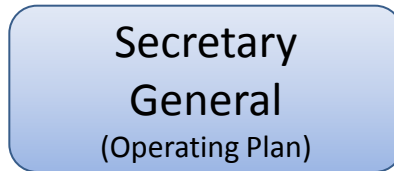
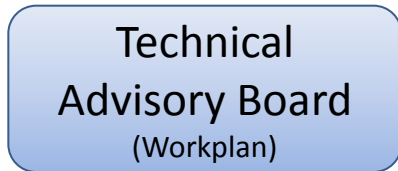








RDA Bodies



Technical Domain

Administrative Domain

Procedural Domain

Data Practitioners Domain

Working Groups and Interest Groups

- Carry out work of RDA
- Reach consensus on
- May suggest BoFs about new topics
- Open to all but...
 - some commitment expected

Plenary

- Open to all persons involved in RDA
- Hears and comments on reports from WGs
- Suggests new IGs and WGs
- Hears candidates for TAC

Technical Advice Committee

- advise on WG work activities
- Interacting directly with working groups
- advise on new WGs and new BoFs
- Give implementation suggestions to strategic direction from council

Administrative Domain

Administration and Management Team

- Implement strategic direction set by council
- Supports the activities of the RDA
 - Arrange plenary meetings
 - Run the on-line for a
 - Manage documents
- Convene nominating committees for
 - Council and TAC
- Monitor and controls finances
- Prepare reports for
 - Council, funders,....

Council

- Set
- Fin
- Approve new WGs (TAC advised)
- control balanced WG approach

- use for all kinds of activities, open to all RDA members

Online Open Interaction Fora

Example RDA Working Groups

- Data Citation
- Data Foundation and Terminology
- Data Type Registries
- Metadata Standards
- PID Information Types
- Practical Policy
- Standardisation of Data
- ...

The screenshot displays the RDA (Research Data Alliance) website. At the top, the RDA logo and tagline 'Research Data Sharing without barriers' are visible, along with social media icons for Facebook, Twitter, Pinterest, and Google+. A navigation menu includes 'Home', 'About', 'Working and Interest Groups', 'News & Events', 'Plenary Meetings', and 'RDA in the Press'. The main content area is titled 'All Working Groups' and lists several groups with brief descriptions and 'READ MORE' links:

- Community Capability Model WG:** The RDA Community Capability Model (CCM) Working Group (WG) collects, validates and publishes a range of data-centric "capability profiles" to enhance inter- and intra-domain interoperability and catalyse RDA data-sharing goals.
- Data Citation WG:** The RDA Working Group on Data Citation (WG-DC) aims to bring together a group of experts to discuss the issues, requirements, advantages and shortcomings of existing approaches for efficiently citing subsets of data. The WG-DC focuses on a narrow field where we can contribute significantly and provide prototypes and reference implementations.
- Data Foundation and Terminology WG:** The Data Foundation and Terminology Working Group describes a basic abstract data organization model which can be used to derive a reference data terminology that can be used across communities and stakeholders to better synchronize conceptualization, enable better understanding within and between communities and stimulate tool building.
- Data Type Registries WG:** The Data Type Registries Working Group will compile a set of use cases for data type use and management. It will identify and distinguish among existing type registry efforts and their potential interaction with this group and formulate a data model and expression for types. It will also design a functional specification for type registries and propose a federation strategy among multiple type registries at both technical and organizational levels.
- Metadata Standards WG:** The RDA Metadata Working Group identifies and endorses metadata solutions for addressing data management challenges. Metadata is crucial for the discovery, access, preservation, exchange, manipulation, and use/re-use of scientific data.
- PID Information Types WG:** (Partially visible at the bottom)

On the right side, there is a vertical navigation menu with a green background, listing: Home, About, Working and Interest Groups, Working Groups Process, Goals and Outcomes, Case Statements, Working Groups (highlighted), Interest Groups, RDA Groups, News & Events, Plenary Meetings, and RDA in the Press. Below this menu, there is a 'News & Events' section featuring a photo of science ministers and the headline 'G8 science ministers meet in London'. A 'Who is involved?' section features a photo of Professor John Wood CBE, FREng, Secretary-General of the Association of Commonwealth Universities and High Level Expert Group on Scientific Data Information Chair & European Research Area Board Chair. At the bottom right, there is an 'RDA Newsletter' section and a 'Tag Cloud'.

Some Risks

- Standardisation is easy, I've done it a hundred times
(apologies to Mark Twain)
- Two easy ways to standardise:
 - The Imperial model
 - The Esperanto model
- Justify need, define benefit, involve stakeholders
- Make a small steps and reassess
- “Never generalise from one example”

Supporting Projects

Three projects supporting RDA through its first phase:

- RDA/Europe (previously iCordi) EC Project
- RDA/US NSF Project
- Support in Australia through ANDS

Steering Group setting it up:

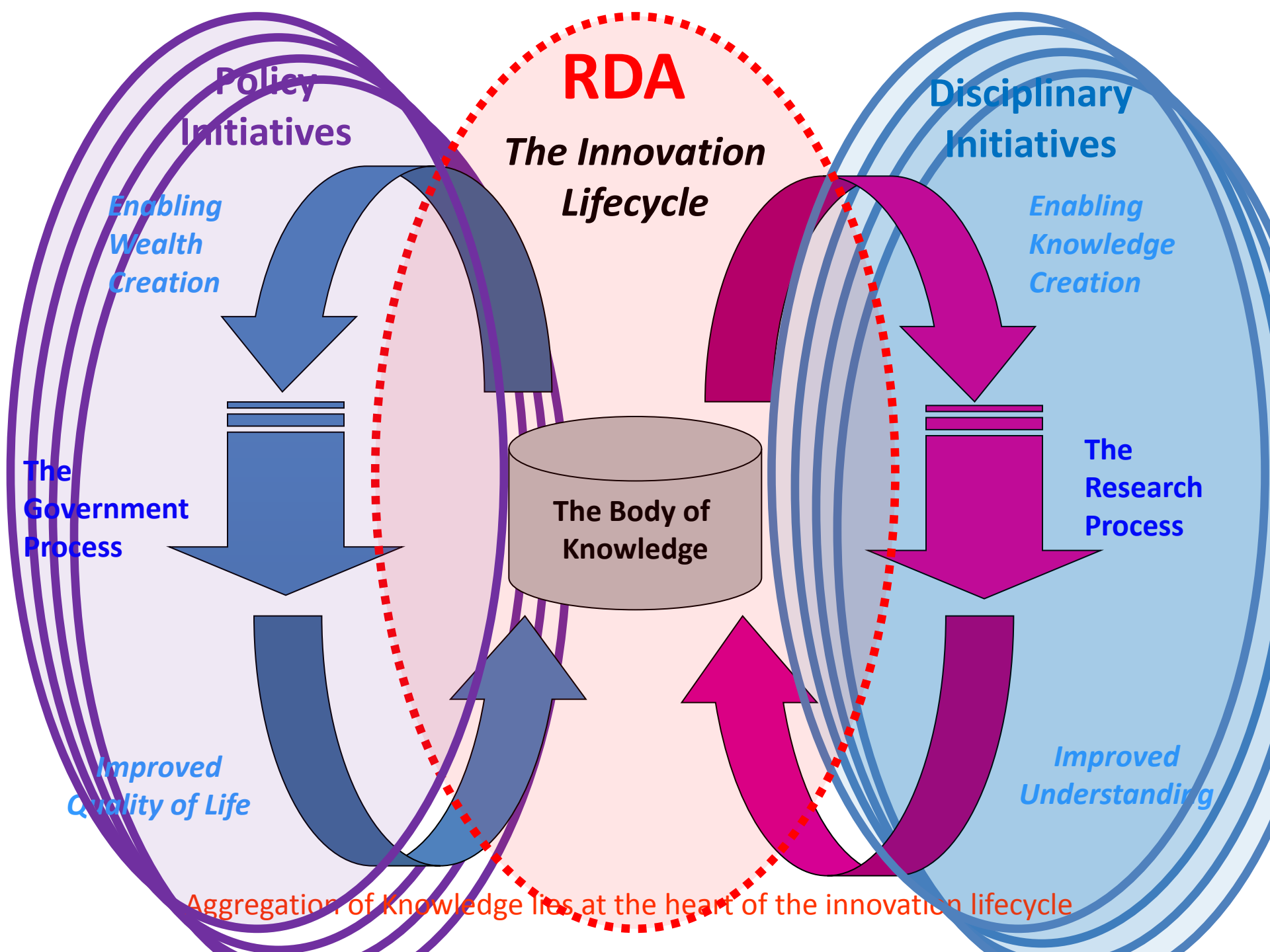
- US – Fran Berman, Beth Plale
- EU – Leif Laaksonen, Peter Wittenburg, Juan Bicarregui
- Australia – Ross Wilkinson, Andrew Treloar
- TAB to be elected at 2nd Plenary
- First Organisational Assembly at 2nd Plenary

RDA Status in June 2013

- Pre-launch meetings in Munich and Washington September 2012,
 - ~200 Delegates
- Various Workshops eg through eIRG, IDCC,
- Launch and First Plenary, March 2013, Guttenberg,
 - ~250 participants
- Currently, 8 Working Groups and 14 Interest Groups
- **Second Plenary, September 16-18 2013, Washington**
- **Third Plenary, March 26-28, 2014, Dublin**
- **Fourth Plenary, TBD**

Please get involved by registering and participating in the discussions:

Website: rd-alliance.org/



Overview

1. The Policy Context

- OECD
- EC/NSF/...
- G8+5
- RCUK
- Royal Society
- G8

2. PaNdata

- Photon and Neutron Open Data Infrastructure

3. The Research Data Alliance

- Fostering Collaboration on a global scale

Thank You

A silhouette of a satellite dish on a tower against a sunset sky with clouds. The dish is mounted on a lattice tower and is pointed towards the upper right. The sky is filled with soft, golden clouds, and the sun is low on the horizon, creating a warm, orange glow. The overall scene is peaceful and evocative of communication and technology.

www.rcuk.ac.uk/research/Pages/DataPolicy.aspx

www.pan-data.eu

www.rd-alliance.org

The End