

# Linking raw data with scientific workflow and software repository: some early experience in PanData-ODI

Erica Yang, Brian Matthews

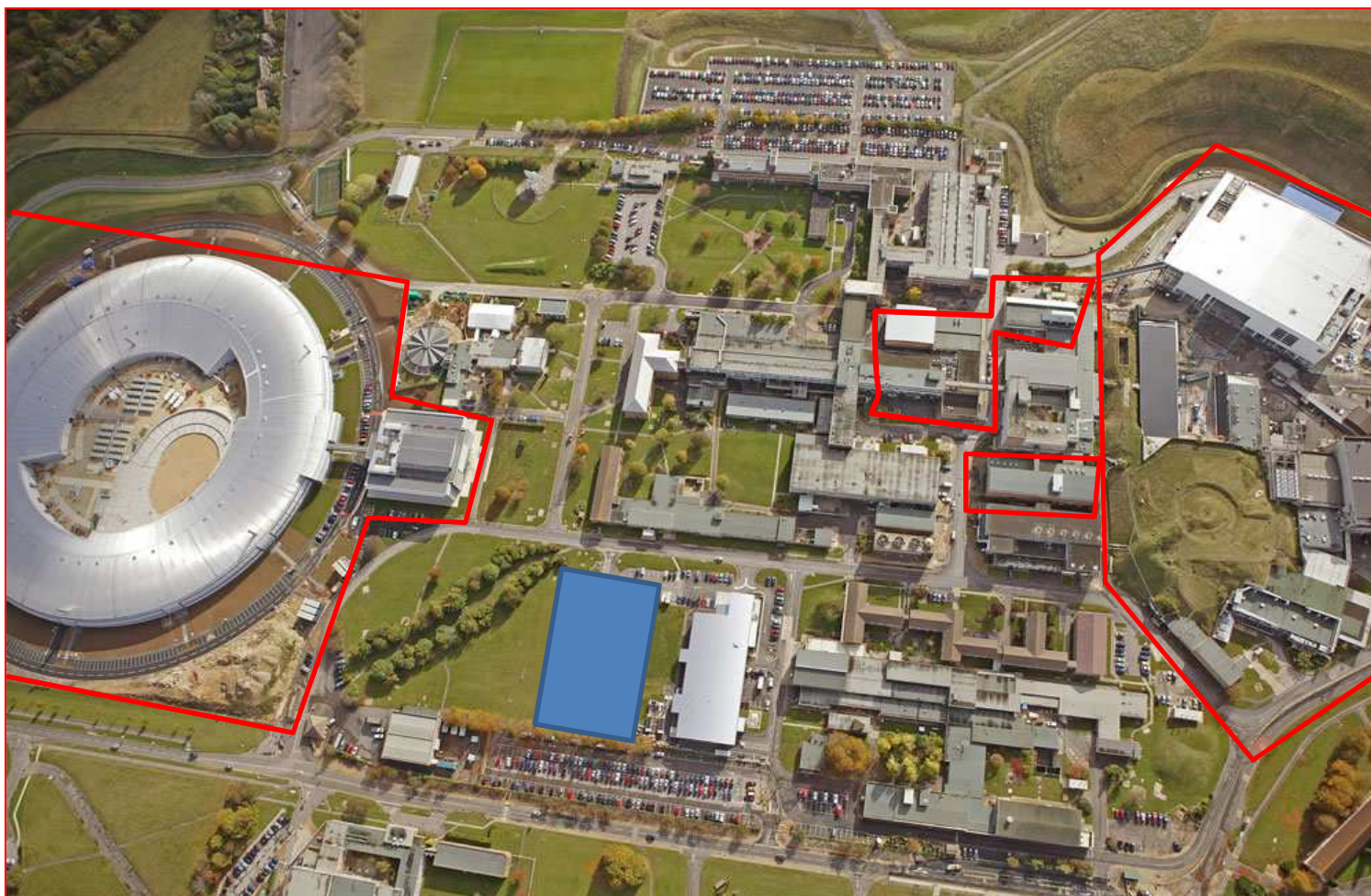
Scientific Computing Department (SCD)

Rutherford Appleton Laboratory (RAL)  
Science and Technology Facilities Council (STFC), U. K.

[erica.yang@stfc.ac.uk](mailto:erica.yang@stfc.ac.uk)

# BACKGROUND

# STFC Rutherford Appleton Laboratory



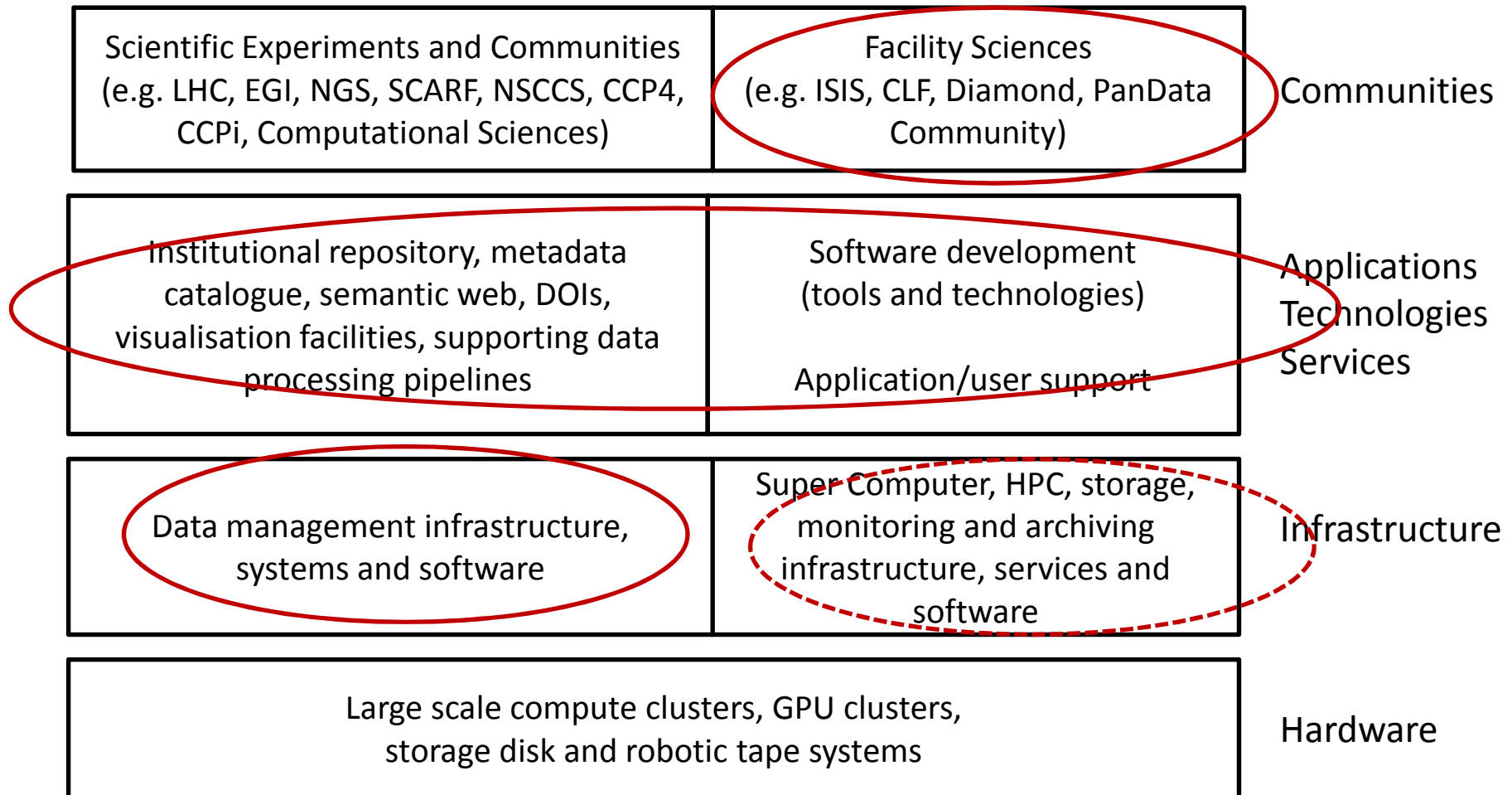
# What we do

## Scientific Computing Department ...

Scientific Experiments and Communities (e.g. LHC, EGI, NGS, SCARF, NSCCS, CCP4, CCPi, Computational Sciences)	Facility Sciences (e.g. ISIS, CLF, Diamond, PanData Community)	Communities
Institutional repository, metadata catalogue, semantic web, DOIs, visualisation facilities, supporting data processing pipelines	Software development (tools and technologies)  Application/user support	Applications Technologies Services
Data management infrastructure, systems and software	Super Computer, HPC, storage, monitoring and archiving infrastructure, services and software	Infrastructure
Large scale compute clusters, GPU clusters, storage disk and robotic tape systems		Hardware

# What we do

Scientific Computing Department ...



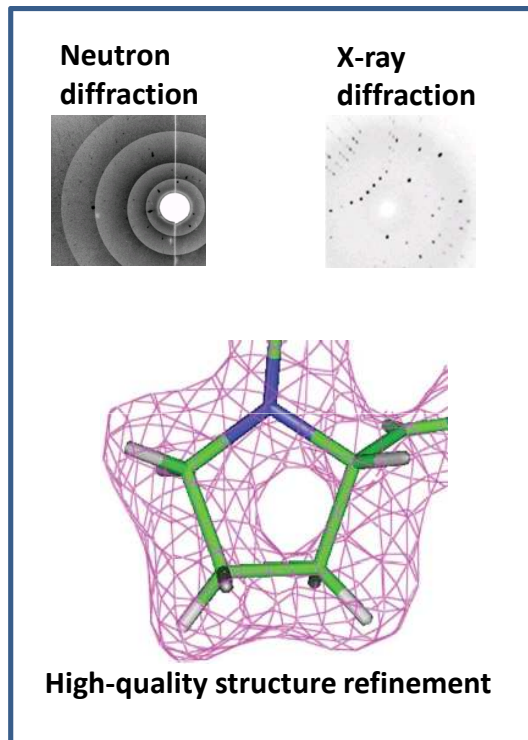
**I am a computer scientist ...**

# PANDATA-ODI

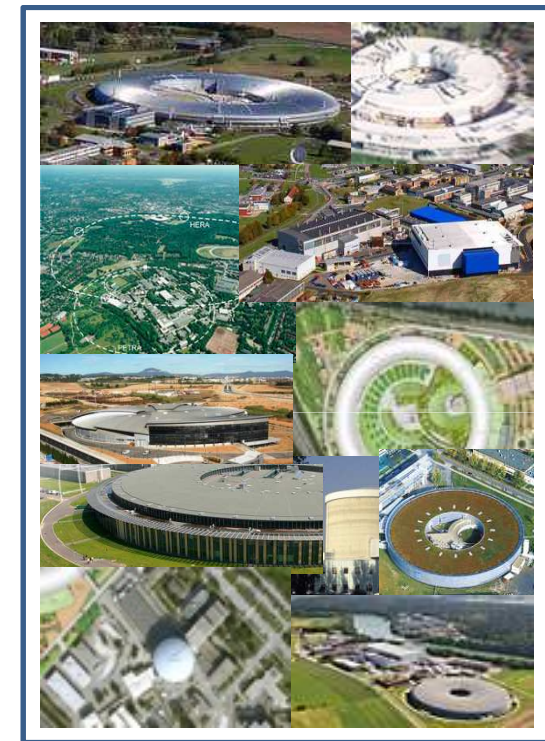


# PaN-data ODI– an Open Data Infrastructure for European Photon and Neutron laboratories

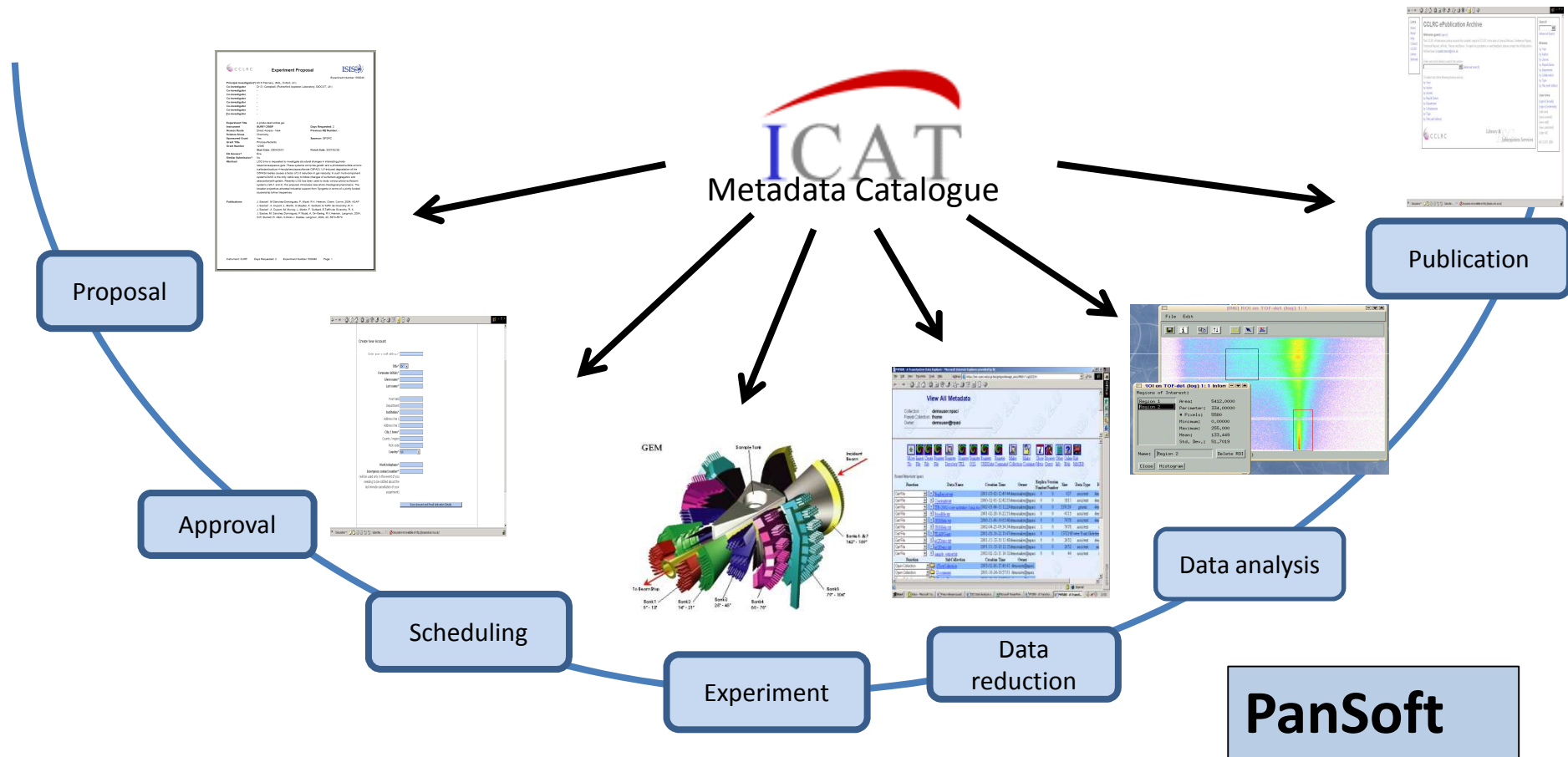
Federated data catalogues supporting cross-facility, cross-discipline interaction at the scale of atoms and molecules



- Unification of data management policies
- Shared protocols for exchange of user information
- Common scientific data formats
- Interoperation of data analysis software
- **Data Provenance WP: Linking Data and Publications**
- Digital Preservation: supporting the long-term preservation of the research outputs



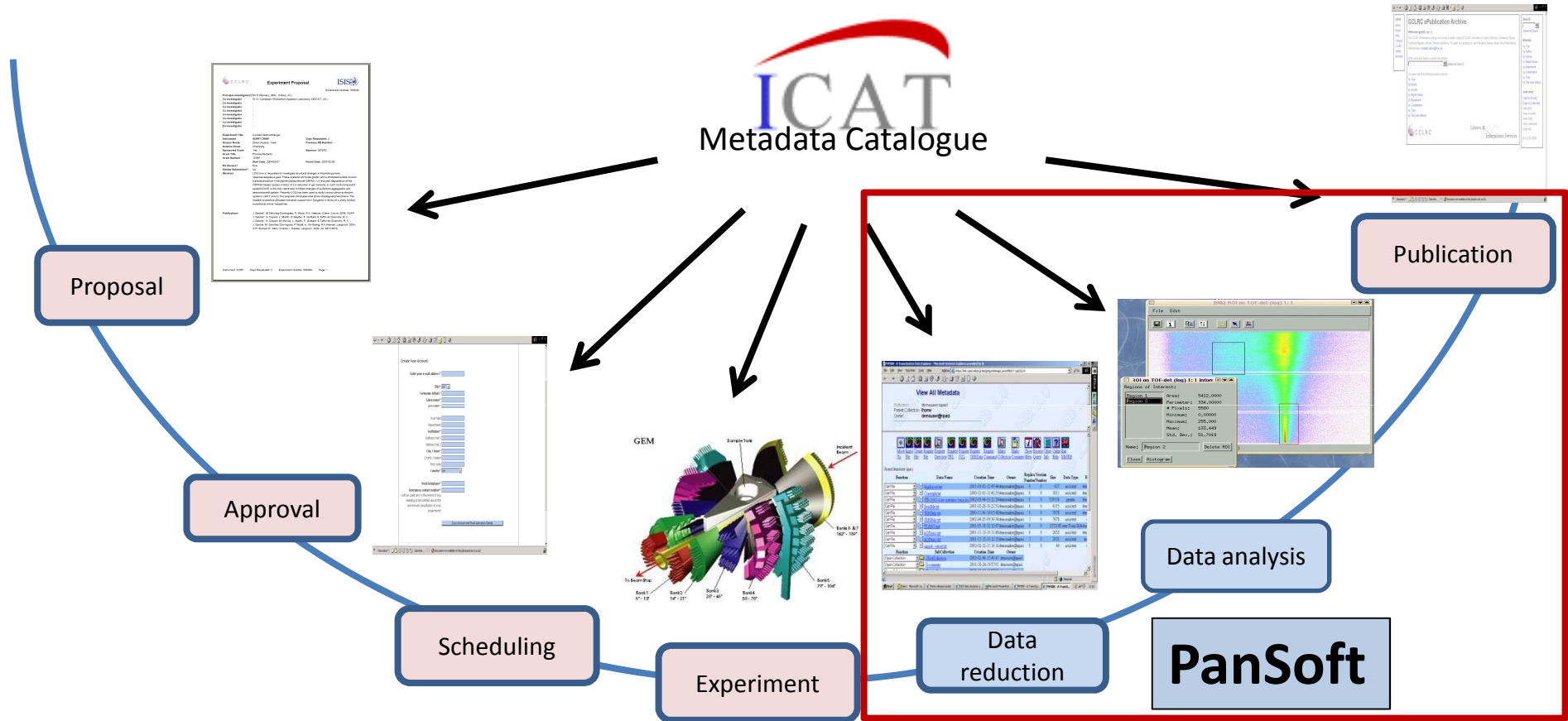
# (Facility) Data Continuum



Services on top of ICAT:  
DOI, TopCAT, Eclipse ICAT Explorer ...



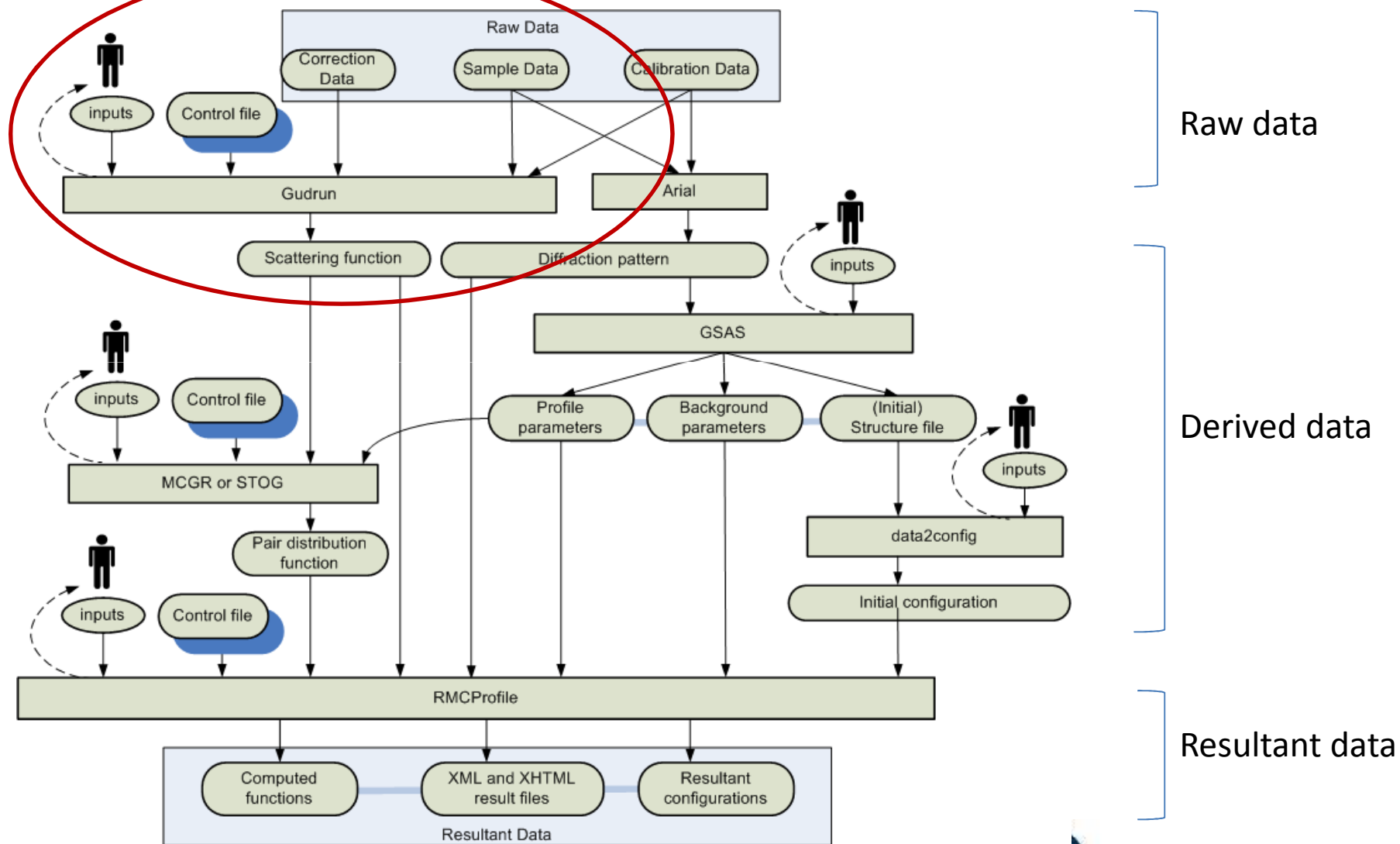
# Data Continuum



 Well developed and supported

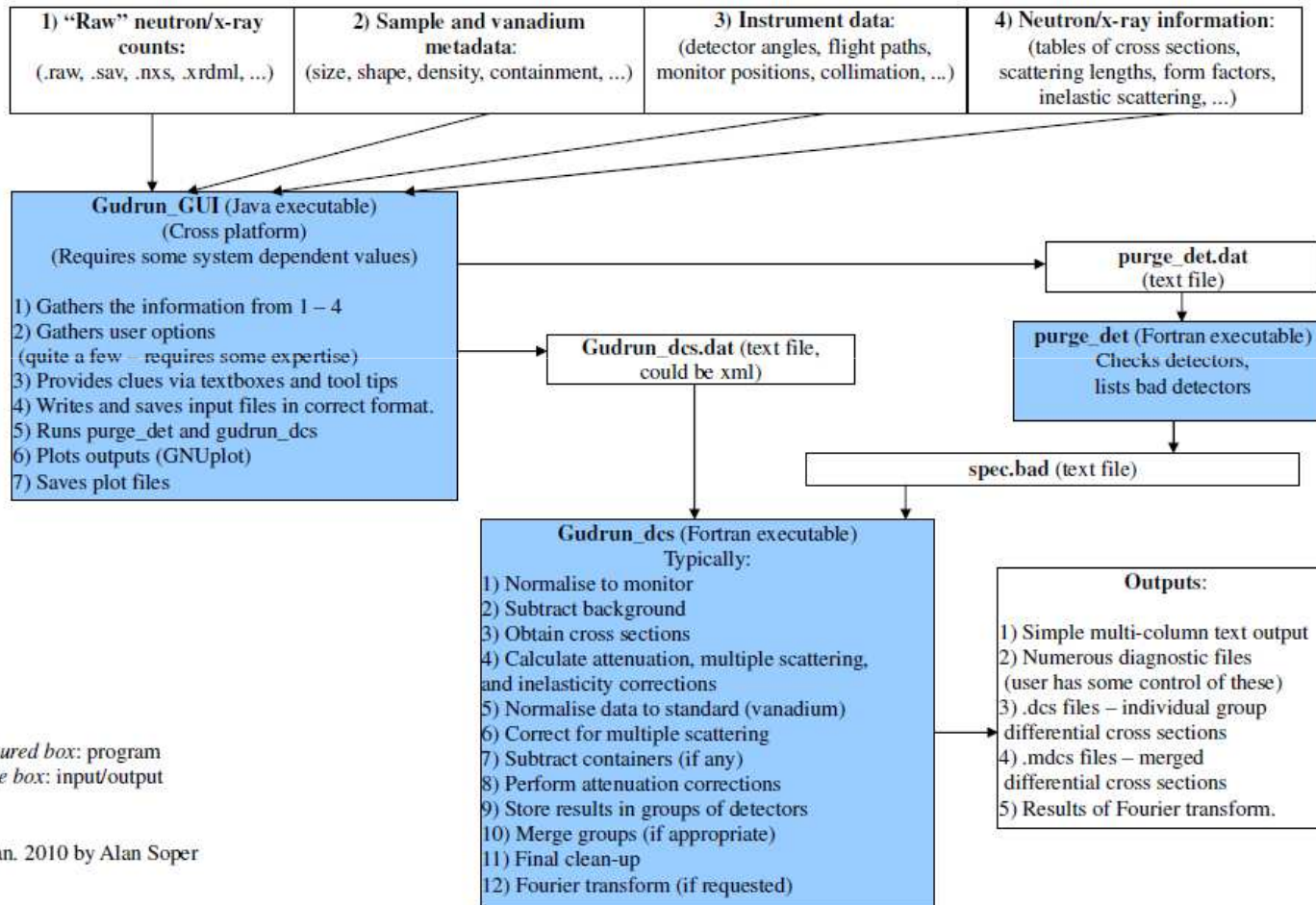
- These are with users.
- Traditionally, these, although very useful for data *citation, reuse and sharing*, are very difficult to capture!
- Practices vary from individuals to individuals, and from institutions to institutions

# Prior Experience



# Can be difficult to capture ...

Gudrun flow diagram



# Data Provenance: Case Studies (so far)

	Description	Facility	Likely stages of continuum
1.	<b>Automated SANS2D reduction and analysis</b>	ISIS	<ul style="list-style-type: none"> <li>- Experiment/Sample preparation</li> <li>- Raw data collection</li> <li>- Data reduction</li> <li>- Data Analysis</li> </ul>
2.	<b>Tomography reconstruction and analysis</b>	Manchester Uni./DLS	<ul style="list-style-type: none"> <li>- Raw data</li> <li>- Reconstructed data</li> <li>- Analysed data</li> </ul>
3.	EXPRESS Services	ISIS	<ul style="list-style-type: none"> <li>- Experiment preparation</li> <li>- Raw data collection</li> <li>- Data Analysis to standard final data product</li> </ul>
4.	Recording publications arising from proposal	ISIS	<ul style="list-style-type: none"> <li>- Proposal system</li> <li>- Raw Data collection</li> <li>- Publication recording</li> </ul>
5.	DAWN + iSpyB	DLS, ESRF	<ul style="list-style-type: none"> <li>- Experiment preparation</li> <li>- Data analysis steps</li> </ul>

**These case studies have given us unique insights into today's facilities ...**

**Looking for more case studies ...**

# TODAY'S FACILITIES – FACTS



# Diamond and ISIS (the story so far ...)

- Diamond:
  - ~ 290TB and 120 million files [1]
  - In SCD data archive
  - Largest file: ~120GB – Tomography beamlines [3]
- ISIS [2]:
  - ~16TB and 11 million files
  - In ISIS data archive & SCD data archive (on going)
  - Largest file: ~16GB – the WISH instrument

# Diamond Tomography

- Up to 120 GBs/file every 30 minutes
- 6,000 TIFF images/file
- Up to 200 GBs/hr
- ~5 TBs/day
- 1-3 days/experiment

How many images are there for each experiment?

# ISIS

- Peak data copy rate (off instrument) 100MB/S
  - Expecting to reach 500MB-1GB/S (but not soon).
- Expect to grow to 10TB/cycle in 3-5 years
- Become interested in centrally hosted services (WISH)

# WHAT DOES IT MEAN?

# It means ...

- Due to the volume, it is not cost effective to transfer the (raw + reconstructed) data back to the home institutions **or elsewhere to process**
  - The network bandwidth to universities mean that it will take a long time to transfer ...
  - So, users have to physically take data back home on storage drive ...
- It is impossible for users to do the reconstruction or analysis on their own computer/laptop.
  - How much RAM do you have on your laptop? And how big is the file from WISH? The Mantid story ...
- It is expensive to re-do the analysis back to the home institutions because of
  - Lack of hardware resources
  - Lack of metadata (a large number of files often mean that there is not much useful metadata)
  - Lack of expertise (e.g. parallel processing, GPU programming)
  - (Assuming software is open source ...)
- Facilities become interested, again, in centralised computing services, **right next to the data**
  - The ISIS WISH story ...
  - Diamond GPU cluster vs. SCD GPU cluster (directly linked to the data archive)

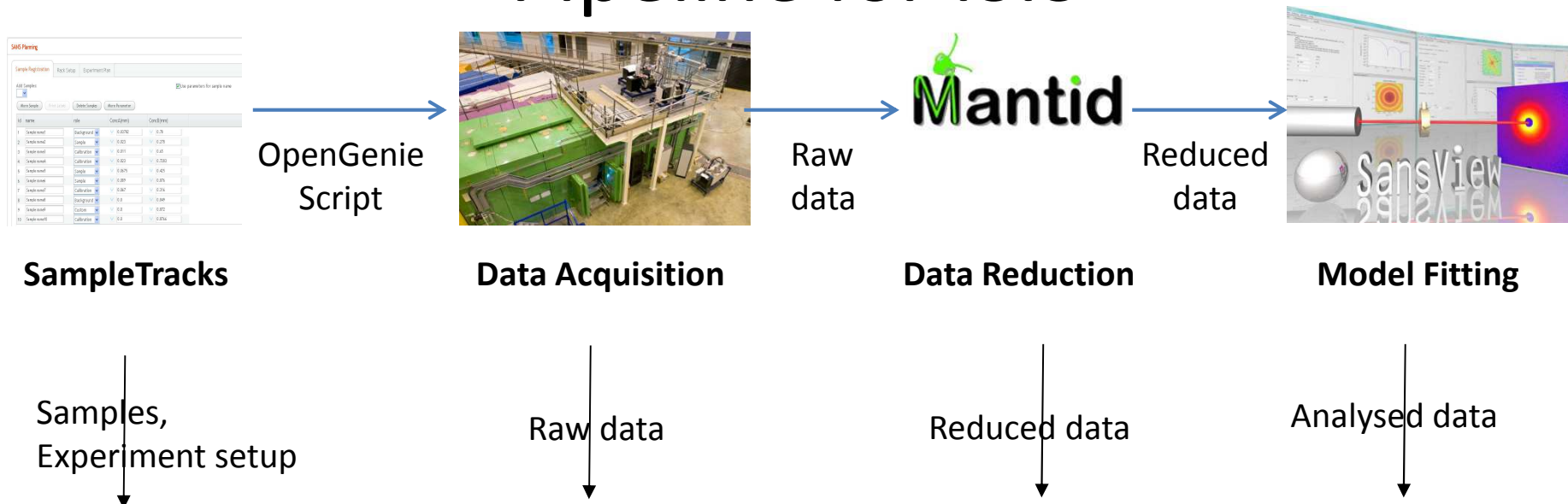


# Users now ask for remote services

- Users are interested in remote data analysis services
  - “... Of course this would mean a *step change in the facilities provided and the time users spend at the facility*. ...”
  - In response, we are developing: “TomoDM” ...
- Then, what benefits can remote services bring?
  - Systematic recording of data continuum, thus allowing the recording of scientific workflows, software, and data provenance (the first four categories data as defined in the “Living Publication” [4])
  - Drive data processing (reduction & analysis) with data provenance
  - It is not only possible to create bi-directional links between raw data and publications, it is also possible to systematically create pair-wise bi-directional links between raw, derived, resultant data and publications.

# TWO EXAMPLES

# SRF – Automated Data Processing Pipeline for ISIS



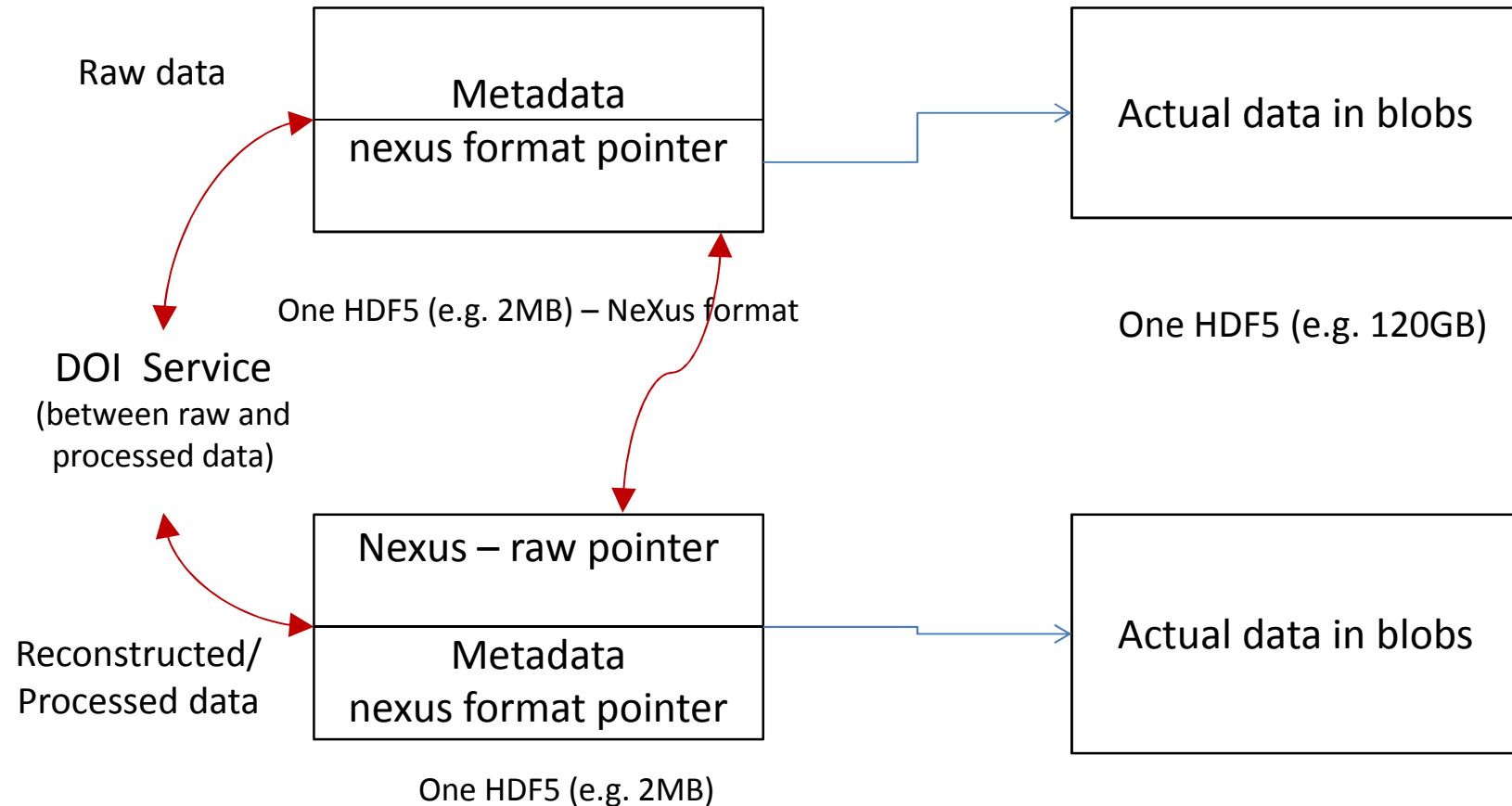


**ICAT Data Catalogue**



**Blog Posts in LabTrove**

# Links between metadata and files



# But ...

- If we can capture all types of data
  - Samples
  - Raw data
  - Reduced data
  - Analysed data
  - Software
  - Data provenance (i.e. Relationships between datasets)
- Facilities operators do not normally
  - Perform validation and fixity checking of data files (volume vs. cost)
  - *Actually MUCH cheaper to simply back up all images without regard for quality, relevance or even duplication, than it is to build an infrastructure for automatically analyzing millions of image files to determine which are “worth keeping”. [1]*
  - But lifetime of data on tape vs. Lifetime of tape
- Will DOIs good enough for data citation?



# STFC EXPERIENCE OF USING DATACITE DOIS

# STFC DOI landing page



Log in | Sign up

Search

About STFC

Business & Innovation

Funding & Grants

Sites & Facilities

Our Research

Public & Schools

How we operate

Collaborate with STFC

Calls, rules & statistics

A guide to STFC

Overview of programmes

Engaging the public



Data collected on the CRISP instrument at the ISIS facility

## ISIS Data

RB920486.

Investigation title: Electric field effect on the interfacial uncompensated spins in the Co/BiFeO<sub>3</sub>/STO exchange bias system.

Creator: *Steinke, N J*

DOI: 10.5286/ISIS.E.24079627

Date of Experiment: Fri Jul 23 08:52:43 BST 2010

Publisher: STFC ISIS Facility

Data format: RAW/Nexus

Select the data format above to find out more about it.



DOWNLOAD

download the dataset

## Data Citation

The recommended format for citing this dataset in a research publication is as:  
[author], [date], [title], [publisher], [doi]

For Example:

Steinke, N J. et al; (2010): RB920486, STFC ISIS Facility, doi:10.5286/ISIS.E.24079627

Abstract

# Behind the landing page

## Datafile Window

Download		
File Name	File Location	File Size
<b>Dataset Name: Default (20 Items)</b>		
CSP80476_He_Level.bt	Wisinst\$Instruments\$IN...	0.001 MB
CSP80476_height.bt	Wisinst\$Instruments\$IN...	0.003 MB
CSP80476_ICPdebug.bt	Wisinst\$Instruments\$IN...	0.091 MB
CSP80476_ICPevent.bt	Wisinst\$Instruments\$IN...	0.012 MB
CSP80476_ITC_Set_Point.bt	Wisinst\$Instruments\$IN...	0 MB
CSP80476_ITC_Temperat...	Wisinst\$Instruments\$IN...	0.006 MB
CSP80476_Linear_Det_He...	Wisinst\$Instruments\$IN...	0 MB
CSP80476_Moderator_Te...	Wisinst\$Instruments\$IN...	0.001 MB
CSP80476_phi.bt	Wisinst\$Instruments\$IN...	0.001 MB
CSP80476_psi.bt	Wisinst\$Instruments\$IN...	0.001 MB
CSP80476_res.bt	Wisinst\$Instruments\$IN...	0.002 MB
CSP80476_s1.bt	Wisinst\$Instruments\$IN...	0 MB
CSP80476_s2.bt	Wisinst\$Instruments\$IN...	0 MB
CSP80476_s3.bt	Wisinst\$Instruments\$IN...	0.005 MB
CSP80476_s4.bt	Wisinst\$Instruments\$IN...	0 MB
CSP80476_sample.bt	Wisinst\$Instruments\$IN...	0 MB

Page 1 of 100

20/02/2009 16:03
20/02/2009 16:03
20/02/2009 16:03
20/02/2009 16:03
20/02/2009 16:03
20/02/2009 16:34
20/02/2009 16:05
20/02/2009 16:14
20/02/2009 16:43
20/02/2009 16:03
20/02/2009 16:14
20/02/2009 16:14
20/02/2009 16:14
20/02/2009 16:03
20/02/2009 18:12
20/02/2009 18:12

Displaying 1 - 20 of 1995

# Citation Issues

- At what granularity should data be made citable?
  - If single datasets are given identifiers, what about collections of datasets, or subsets of data?
  - What are we citing? Datasets or an aggregation of datasets
- STFC link DOI to experiments on large facilities which contain many data files
- Other organisations DOI usage policies use different granularities.
- Can there be a common, or discipline common DOI usage policy for data?
- Citation of different types of object: software, processes, workflows ...
- Private vs. Public datasets/DOIs

# Summary

- Infrastructure services
  - Compute
  - Storage
  - Archiving + preservation (validation, integrity)
  - Remote data processing
  - DOIs
- It is already happening
  - DLS is already collecting processed data (e.g. iSpyB)
  - But, interesting issues remind to be resolved when lots of data become available ...



# Acknowledgements

Project (s)	Institution(s)	People
PanData-ODI (case study)	Diamond, U.K.	Mark Basham, Kaz Wanelik
PanData-ODI: (case study)	Manchester University, U.K. Harwell Research Complex, RAL Diamond, U.K.	Prof. Philip Withers, Prof. Peter Lee
SCAPE, SRF, PanData-ODI	Scientific Computing Department, STFC	David Corney, Shaun De Witt, Chris Kruk, Michael, Wilson
I2S2, SRF	Chemistry, Southampton University, UK	Prof. Jeremy Frey, Dr. Simon Coles
I2S2	Cambridge University, UK	Prof. Martin Dove
I2S2, PanData-ODI, SRF, SCAPE	ISIS, STFC	Tom Griffin, Chris Morton-Smith, Alan Soper, Silvia Imberti, Cameron Neylon, Ann Terry, Matt Tucker

I2S2, SRF are funded by JISC, UK  
PanData-ODI, SCAPE are funded by EC

# References

1. Colin Nave (Diamond), “Survey of Future Data Archiving Policies for Macromolecular Crystallography at Synchrotrons”, distributed via “dddwg-bounces@iucr.org”, July 2012.
2. Chris Morton-Smith (ISIS), “ISIS Data rates and sizes – up to March 2012”, May 2012.
3. Mark Basham (Diamond), “HDF5 Parallel Reading and Tomography Processing at DLS”, PanData-ODI DESY meeting, Hamburg, Germany, Feb. 2012.
4. John R Helliwell, Thomas C. Terwilliger, Brian McMahon, “The Living Publication”, April 2012.

