# PaNdata ODI 1st Open Workshop

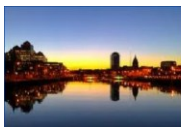Dublin 24-25/28th of March 2014

Co-located with the RDA 3rd plenary at Croke Park

https://indico.desy.de/event/1stow

## Booklet of collected presentations

PART 3: PaNSIG

pandata ODI

# Photon and Neutron Science Interest Group meeting

# 26th and 28th of March at Croke Park / RDA

| Thomas Proffen/ Amber Boehnlein | Photon and Neutron Science and facilities - the computing challenges |
|---|---|
| Tom Griffin: | Data Management solutions – ICAT |
| Brian Matthews: | Data Management solutions – Tardis |
| Brian Matthews/ Steve Androulakis | Potential areas for future collaboration and development of common standards and services within RDA |
| Adam Farquhar: | DataCite |
| Alun Ashton: | Data analysis issues and frameworks |
| Eugen Wintersberger: | HDF5 |
| Brian Matthews: | PaNdata |
| PaNSIG chairs: | Developing a plan of activities for PaNSig |

pandata ODI

# Computing Challenges for Photon and Neutron Facilities

# Driving Factors: Computing

- Meeting the current science goals for major initiatives such as energy research and materials discovery require improvements in computational tools and techniques.

- Science is driving source upgrades
  - Brighter and more precise sources drive detector development
  - Detector development drives computing needs—volume, and complexity

- Science also directly drives the computing needs
  - Simulations
  - Data Analysis and analytics

# Facility Challenges

- O(10000) users
  - Diversity of science, needs, skills, longevity

- O(100) beam lines

- O(10) different imaging techniques

- Operational constraints

- Can Do/Make Do Culture

- Source and detector upgrades
  - Challenging in themselves;
  - Environment is changing

# Context

It is relatively easy to make a list of common areas of interest for data and computing

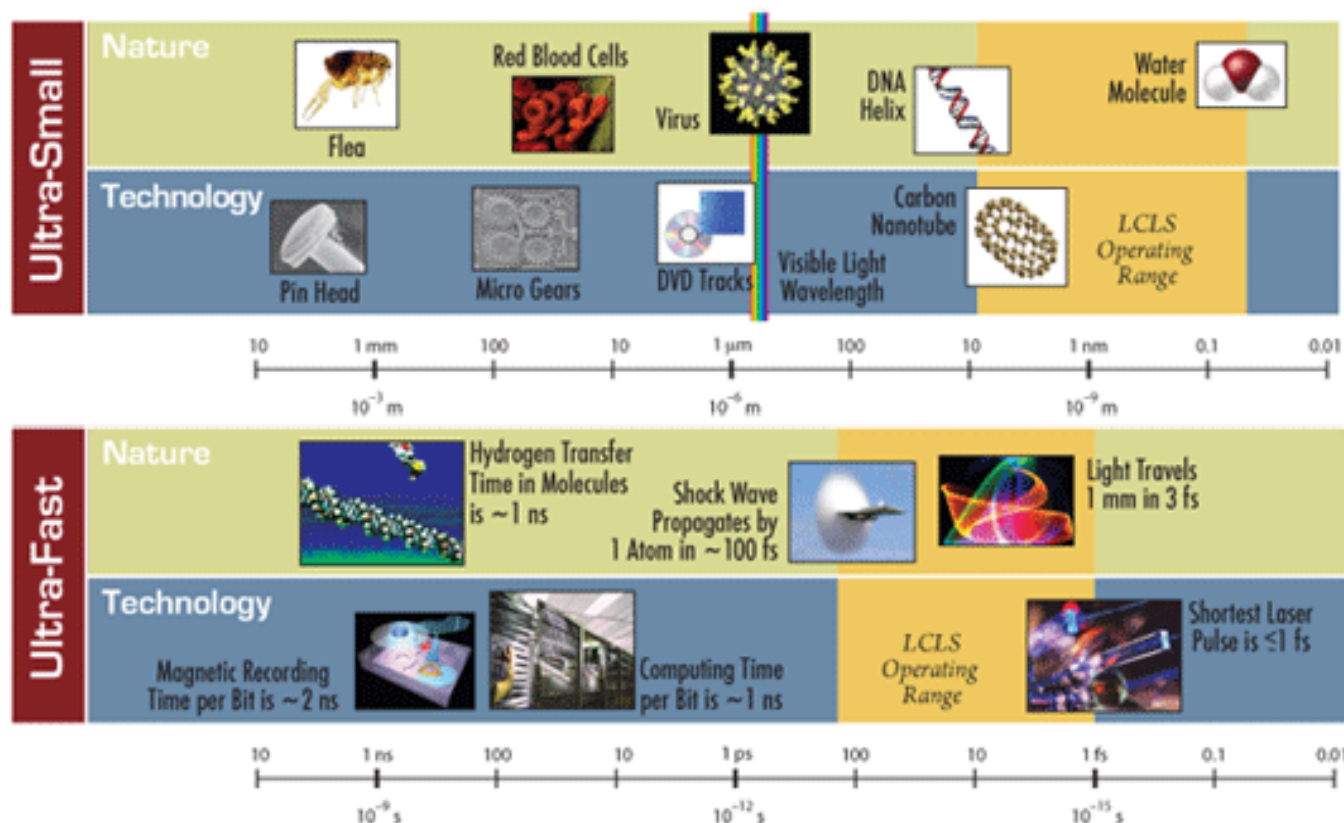Can we collaborate across facilities on those areas productively?

With limited resources, what are good investments?

How do we plan and make the case?

Who, what, where, when and how?

# LCLS Data Rates & Volumes

Vary for different instruments/experiments

Three operation modes (X-Ray pulse rate):

**30 Hz** (first runs on LCLS, AMO/CAMP, Oct-Dec 2009)
**60 HZ** (second series of runs, May 2010 -> on)
**120 Hz** (Fall 2010)

Instruments:

**Fall 2009**: AMO/CAMP
**Summer 2010**: AMO/CAMP, XPP, SXR
**Fall 2012:** all 6+1 instruments (station #2 for CXI)

First runs at AMO/CAMP (30 Hz):

up to **180 MB/s**, **3 TB/day**, **~100 TB** of raw data recorded

At full capacity (120 Hz):

**Up to 1.5 GB/s** for CXI (event size exceeds 10 MB)
**30 TB/shift**
**Up to 3 instruments can take data in parallel (**approaching this condition**)**
**1-2 PB** of raw data per year

# The Data Challenge view from 2012

**CURRENTLY STORED**   **APPROX. PER ANNUM DATA COLLECTION RATE**   **PROJECTED DATA COLLECTION RATE IN 2017**

+15/yr   LHC

8.2   DØ EXPERIMENT (FERMILAB)

+2.5/yr   LCLS

+2/yr   APS

+1/yr   ALS

.7   BaBar EXPERIMENT (SLAC)

.2   LIBRARY OF CONGRESS

.05   HUBBLE SPACE TELESCOPE

+20/yr   EURO XFEL

+15/yr   LCLS

+12/yr   APS

+11/yr   NSLS-II

+8/yr   ALS

+8/yr   ESS (2019)

+2/yr   SNS

7

# Data Volumes Growing everywhere

## Data Rates and Data Volumes

| Application | Current data rates [GB/h] | | Future data rates [GB/h] | |
|---|---|---|---|---|
| | Peak | Average | Peak | Average |
| Protein Crystallography | 500 | 50 | 500 | 200 |
| Coherent Diffraction Imaging | 500 | 50 | 4000 | 400 |
| Tomography | 700 | 50 | 800 | 200 |
| Spectroscopy | 450 | 45 | 18000 | 1800 |
| Small Angle Scattering | 1400 | 140 | 14400 | 4200 |
| Grain Mapping | 140 | 80 | 800 | 300 |

DESY: CFEL/PETRA III+/FLASH → 1.6 PB/year
Cern: ATLAS 100 MB/s → 3 PB/year

- Increasing source brilliance
- Faster detectors with more pixels (e.g. Pilatus)
- New experiment scenarios: Fast time resolved, scanning
- Use of simultaneous Detectors

DESY

HELMHOLTZ | GEMEINSCHAFT

*Photo-system II, Nick Sauter (LBL):* A single experiment in 2013 generated **120TB** over four days. Data was relayed to NERSC via Esnet, analysis required **135K CPU hours** of computing.

LCLS-II will require > **100M CPU hours** per experiment. Higher resolution and advanced image analysis could grow computational complexity. Some algorithms scale an M*NlogN for M images of N pixels.



From : Thu Feb 28 13:25:12 2013    To : Fri Mar 1 13:25:12 2013

To site    From site

**Total traffic**   *Tip: Double Click to Zoom-In and [SHIFT] Double click to Zoom-Out*

All NERSC Traffic

Traffic split by : **'Autonomous System (origin)'**

nersc-SLAC:3671

Photosystem II X-Ray Study

9

# A starter wish list…

Easy to use data management and processing frameworks  that scale with data rate

Conditions data storage and management; meta data

Data validation and fast feedback at collection

Data volume reduction/compression techniques

Development and availability of algorithm and algorithm tools

Community developed simulations and simulation tool kits for beamlines and detectors

Computational science

Compute and storage hardware platforms appropriate for the task**S**

# Data management

- Life is easier if the data is managed from point of origin

- Detector readouts can have proprietary readouts

Is standardization possible?
Data containers/formats
Metadata
Workflows/tool kits
Visualization tools
Curation/cataloging



Naively, some standardization
 would make good use of resources and expertise,
simplify life, open up many possibilities

# Rapid feedback

Time is money
    Beam time, experimenter time; instrument scientist time
Experiment simulations
Advance preparation of
analysis tools

Could provide better use of
beam time, shorten time to
publication and refine where to spend resources on
computing improvements

A reconstructed image of the Photosystem I complex. (Image courtesy Raimund Fromme, Arizona State University.)

Three-dimensional rendering of the X-ray diffraction pattern for the Photosystem I protein, reconstructed from more than 15,000 single nanocrystal snapshots taken at the LCLS. (Image courtesy Thomas White, DESY.)

# High Throughput XRD for Material Discovery

**Strategy:**
- Ternary and Quaternary Combinatorial Libraries
- Screen for Properties
  - HT Spectroscopy
  - New Materials
- Determine Structural Phase Diagram
  - Sam Webb
  - HT XRD | HT EXAFS
  - New Phase Diagrams
- Composition-Structure-Property Relations
- Phase Transitions
  - Metastable & Near Equilibrium
- Input for Atomistic Material Modeling
  - New Theory for (Metastable) Materials

Resistivity map

Ni
Fe
Co
FIG. 5-78. Electrical resistivity of annealed Fe-Co-Ni alloys.

Phase map

Ni
α
Fe
Co

Ni
Fe
Co

properties — structures
New composition
properties — structures

Prototype setup @ SSRL BL 1-5

Laser Ht2D XRD
Fluorescence Detector
Gauge
Detector
X-rays

**Currently ~ 2000xrd+mca/day**
exp~30sec + readout ~10sec

**Moving Forward with HT-XRD :**

Enhance Capabilities:
- Auto Alignment and Calibrations
- Low Z Detection Sensitivity
- Rapid Annealing

IncreaseThroughput:
- New X-ray Focusing Optics: exp ~1-5sec
- New 2D detector: readout <1 sec
  - → 20-50,000 xrd+mca/day
  - Need Automation and Robotics

Data Management:

JCAP
JOINT CENTER FOR
ARTIFICIAL PHOTOSYNTHESIS

John Gregoire
Ichiro Takeuchi
Matt Kramer
Apurva Mehta

New Composition-
Structure-Property
Relations

$_{40}Ga_{60}$

New Materials

Scientific
Communities

New Theory of
(metastable) Materials

Hubs

Hubs

SSRL

...ction

tabase

CD

# COMPUTING CHALLENGES FOR PHOTON AND NEUTRON FACILITIES

Thomas Proffen, ORNL

Amber Boehnlein, SLAC

# Creating a Pathway for Scientific Discovery

- Accelerating discovery in materials science
- Enhancing predictive capabilities



Data and Communications in Basic Energy Sciences:
Creating a Pathway for Scientific Discovery
Report of a Workshop Linking Experimental User Facility Needs
with Advances in Data Analysis and Communications
June 2012

- Theory and analysis components should be integrated seamlessly within experimental workflow.

- Move analysis closer to experiment – future possibility of experiment steering.

- Match data management access and capabilities with advancements in detectors and sources.

# Driving Factors: Computing

- Meeting the current science goals for major initiatives such as energy research and materials discovery require improvements in computational tools and techniques.

- Science is driving source upgrades
    - Brighter and more precise sources drive detector development
    - Detector development drives computing needs—volume, and complexity

- Science also directly drives the computing needs
    - Simulations
    - Data Analysis and analytics

# Facility Challenges

- O(10000) users
  - Diversity of science, needs, skills, longevity

- O(100) beam lines

- O(10) different imaging techniques

- Operational constraints

- Can Do/Make Do Culture

- Source and detector upgrades
  - Challenging in themselves;
  - Environment is changing

# Context

- It is relatively easy to make a list of common areas of interest for data and computing

- Can we collaborate across facilities on those areas productively?

- With limited resources, what are good investments?

- How do we plan and make the case?

- Who, what, where, when and how?

-

# Neutron Data Life Cycle

**Feedback**

### DAS
- Neutron events
- Events from sample environment
- Other triggers

### Reduction
- Corrected reduced data (histograms, S(Q,E), ..)
- Merging, reconstruction of data
- Instrument/technique dependent
- Need for 'real' time reduction

### Analysis
- Multi dimensional fitting
- Advanced visualization
- Comparison to simulation / feedback
- Field dependent, large variety of approaches

### Simulation Modeling
- Multitude of techniques (DFT, MD, ..)
- Advanced simulation of experiments
- 'Refinement' using experimental data
- Multiple experiments / probes

# User Facility
- Variety of experiments, topics, methods and 'computer literacy' of users are significant challenge.

# Example: NOMAD Diffractometer



| Raw Data: up to 10^12 events per second | **Acquisition** | Translated Data: Gigabytes to Terabytes | **Reduction** | Reduced Data: e.g. Powder Diffraction Pattern | **Analysis & Simulation** | Analysis: PDF, MD simulation, etc. |

Feedback guiding changes to the experiment setup

Data captured and stored on multiple systems at the beamline

After completion of a "run" data is aggregated on a single system, translation begins

Once data is aggregated reduction begins using a workstation

Analysis and Simulation using mid-scale compute

# Improving Productivity = Changing the Workflow

# ADARA is enabling real-time feedback from experiment, analysis and computational steering

- Leverages our multi-disciplinary capabilities at ORNL coupling Neutron Sciences Directorate with Computing and Computational Sciences Directorate.

- The ADARA Project lets us stream data to computational resources and provide live feedback from experiment in real-time $S(Q,E)$.

- Provides a high performance data backplane for reduction, analysis, and coupling with simulation forming the basis for future work to integrate experiment and simulation.

- Prototype running on HYSPEC instrument. Deployment to other beamlines in 2013/2014.

# ORNL has launched the Center for Accelerating Materials Modeling (CAMM)

- The CAMM will integrate materials modeling/simulation (MD/DFT) directly into the chain for neutron scattering data analysis, **offline** and **online** (in near real time)

- Developing workflows for refinement, integration of MD codes, **neutron scattering corrections** ..

- The CAMM is working with ORNL's Materials Science and Technology Division to study coarse grained MD simulations of polymers PEO-AA (CNMS), *ab-initio* MD simulations for ferroelectrics/thermoelectrics



Example: *ab-initio* MD simulations for ferroelectrics/thermoelectrics. Focus on *width* of dispersions

---

**The Center for Accelerating Materials Modeling (CAMM)**

- *Partnership between ORNL's Neutron Sciences, Physical Sciences and Computing and Computational Sciences Directorates*
- *ORNL SEED money and DOE funds provided to study force field refinement from quasi-elastic and inelastic neutron scattering data*
- *CAMM formed in response to BES proposal call for* Predictive Theory and Modeling

# A New Kind of Laser - The LCLS creates X-ray pulses that can capture images of atoms and molecules in motion

# LCLS Data Rates & Volumes

- Three operation modes (X-Ray pulse rate):
  - **30 Hz** (first runs on LCLS, AMO/CAMP, Oct-Dec 2009)
  - **60 HZ** (second series of runs, May 2010 -> on)
  - **120 Hz** (Fall 2010)
- Instruments:
  - **Fall 2009**:       AMO/CAMP
  - **Summer 2010**:   AMO/CAMP, XPP, SXR
  - **Fall 2012:** all 6+1 instruments (station #2 for CXI)
- First runs at AMO/CAMP (30 Hz):
  - up to **180 MB/s**, **3 TB/day**, **~100 TB** of raw data recorded
- At full capacity (120 Hz):
  - **Up to 1.5 GB/s** for CXI (event size exceeds 10 MB)
  - **30 TB/shift**
  - **Up to 3 instruments can take data in parallel (**approaching this condition**)**
  - **1-2 PB** of raw data per year

Vary for different instruments/experiments

# The Data Challenge view from 2012



Legend:
- CURRENTLY STORED
- APPROX. PER ANNUM DATA COLLECTION RATE
- PROJECTED DATA COLLECTION RATE IN 2017

Left panel:
- +15/yr — LHC
- 8.2 — DØ EXPERIMENT (FERMILAB)
- +2.5/yr — LCLS
- +2/yr — APS
- +1/yr — ALS
- .7 — BaBar EXPERIMENT (SLAC)
- .2 — LIBRARY OF CONGRESS
- .05 — HUBBLE SPACE TELESCOPE

Right panel:
- +20/yr — EURO XFEL
- +15/yr — LCLS
- +12/yr — APS
- +11/yr — NSLS-II
- +8/yr — ALS
- +8/yr — ESS (2019)
- +2/yr — SNS

# Data Volumes Growing everywhere

## Data Rates and Data Volumes

| Application | Current data rates [GB/h] | | Future data rates [GB/h] | |
|---|---|---|---|---|
| | Peak | Average | Peak | Average |
| Protein Crystallography | 500 | 50 | 500 | 200 |
| Coherent Diffraction Imaging | 500 | 50 | 4000 | 400 |
| Tomography | 700 | 50 | 800 | 200 |
| Spectroscopy | 450 | 45 | 18000 | 1800 |
| Small Angle Scattering | 1400 | 140 | 14400 | 4200 |
| Grain Mapping | 140 | 80 | 800 | 300 |

DESY: CFEL/PETRA III+/FLASH → 1.6 PB/year
Cern: ATLAS 100 MB/s → 3 PB/year

- Increasing source brilliance
- Faster detectors with more pixels (e.g. Pilatus)
- New experiment scenarios: Fast time resolved, scanning
- Use of simultaneous Detectors

# LCLS Diffract & Destroy

*Photo-system II, Nick Sauter (LBL):* A single experiment in 2013 generated **120TB** over four days. Data was relayed to NERSC via Esnet, analysis required **135K CPU hours** of computing.

LCLS-II will require > **100M CPU hours** per experiment. Higher resolution and advanced image analysis could grow computational complexity. Some algorithms scale an M*NlogN for M images of N pixels.

From : Thu Feb 28 13:25:12 2013    To : Fri Mar 1 13:25:12 2013

To site    From site

**Total traffic**    *Tip: Double Click to Zoom-In and [SHIFT] Double click to Zoom-Out*

All NERSC Traffic

13:25 Feb 28    17:36 Feb 28    21:46 Feb 28    01:57 Mar 01    06:08 Mar 01    10:19 Mar 01

Traffic split by : **'Autonomous System (origin)'**

**nersc-SLAC:3671**

Photosystem II X-Ray Study

13:25    17:36    21:46    01:57    06:08    10:19

# A starter wish list…

- Easy to use data management and processing frameworks that scale with data rate
- Conditions data storage and management; meta data
- Data validation and fast feedback at collection
- Data volume reduction/compression techniques
- Development and availability of algorithm and algorithm tools
- Community developed simulations and simulation tool kits for beamlines and detectors
- Computational science

- Compute and storage hardware platforms appropriate for the taskS

# Data management

- Life is easier if the data is managed from point of origin
- Detector readouts can have proprietary readouts

- Is standardization possible?
  - Data containers/formats
  - Metadata
  - Workflows/tool kits
  - Visualization tools
  - Curation/cataloging



- Naively, some standardization
 would make good use of resources and expertise,
simplify life, open up many possibilities

# Rapid feedback

- Time is money
  - Beam time, experimenter time; instrument scientist time
- Experiment simulations
- Advance preparation of analysis tools



- Could provide better use of beam time, shorten time to publication and refine where to spend resources on computing improvements

A reconstructed image of the Photosystem I complex. (Image courtesy Raimund Fromme, Arizona State University.)

Three-dimensional rendering of the X-ray diffraction pattern for the Photosystem I protein, reconstructed from more than 15,000 single nanocrystal snapshots taken at the LCLS. (Image courtesy Thomas White, DESY.)

# High-Throughput Pipeline:
## Robotics, Automation & Machine Learning

John Gregoire
Ichiro Takeuchi
Matt Kramer
Apurva Mehta

PCA, Density based clustering, Genetic optimization

XRD dendrograms

Fluorescence

Pd

Fe  Ga

Feature clustering

New Composition-Structure-Property Relations

**Hubs**

**Hubs**

Library Production

Data Collection    Data Archiving

New Materials

Property Screening

**SSRL**

Scientific Communities

New Theory of (metastable) Materials

- Auto-alignment and Calibration

- Robotic Library Changer

- Library Tracking

- Multiple Data Stream Archiving
  - Hdf5?

New structure, not in Database
Ortho Fe/Pd silicide
BCC Fe

New structure, not in Database

FCC FePd

Hex FeGa

FCC

Theory

Fe  Machine Learning

Database

Resistivity mapping (from Bozarth)

# ICAT

Facility Metadata Catalogue

Tom Griffin, STFC ISIS Facility

# Overview

- What is ICAT?
- Uses and Benefits
- Installation
- Data Integration
- Interfaces
- About the ICAT project
- Future developments

# What is ICAT?

ICAT is a database, with a well defined API that provides a uniform interface to experimental data and a mechanism to link all aspects of research from proposal through to publication.

# What is ICAT?



*Example ISIS Proposal*

*GEM – High intensity, high resolution neutron diffractometer*

*H2-(zeolite) vibrational frequencies vs polarising potential of cations*

*B-lactoglobulin protein interfacial structure*

ICAT

### Proposals

Once awarded beamtime at ISIS, an entry will be created in ICAT that describes your proposed experiment.

### Experiment

Data collected from your experiment will be indexed by ICAT (with additional experimental conditions) and made available to your experimental team

### Analysed Data

You will have the capability to upload any desired analysed data and associate it with your experiments.

### Publication

Using ICAT you will also be able to associate publications to your experiment and even reference data from your publications.

Science & Technology
Facilities Council

# Benefits

- Links proposal to data to analysis to publication
- Makes data searchable
- Remote data access
- Assign DOIs to data
- Increased opportunities for data sharing and re-use
- Implements a data policy
- Provenance (Creation, ownership, history)
- Integration with applications

# Where is it used?

- STFC (main developer)
  - ISIS pulsed neutron and muon facility
  - CLF (UK national laser facility)
- Diamond
  - UK national synchrotron
- ILL
  - European reactor neutron source
- Oak Ridge National Laboratory (USA)
  - Spallation Neutron Source
- Currently being rolled out at: ESRF, ALBA, SOLEIL, other PaNdata partners

# Technical

# Installation

- Relational database - Oracle, MySQL
- Java application server - Glassfish
- Support available

# Components

- ICAT is a modular system
  - Authentication
    - Plain text simple database – testing and development
    - LDAP/Active Directory
    - Custom user office connections
  - Data Server
    - Defines an interface – can be implemented as appropriate
    - 'Disk only' implementation (ISIS)
    - Tape/disk cache (Diamond)

**Science & Technology**
Facilities Council

# Flexibility

- Access to data is defined through a 'rules' system
    - E.g. Grant read access to data when a user is a Co-Investigator
    - Grant full control to data when a user is the Principal Investigator
    - Grant full control to all data on an instrument when the user is the instrument scientist for that instrument.
- Facility specific: Instruments (beamlines) , Instrument Scientists, Datafile types, parameters
    - E.g. wavelength, total_proton_count, etc

**Science & Technology**
Facilities Council

# Interfaces

- ICAT exposes a web service (SOAP) API
    - Available for client applications
    - Enables integration with data analysis applications such as DAWN and Mantid.
    - RESTful interface planned
- Default web interface – 'TopCAT'
    - Allows basic data browsing
    - Able to search many ICAT's in parallel

**Science & Technology**
Facilities Council

# Interfaces - TopCAT

# Interfaces - TopCAT

# Interfaces - Mantid

# Interfaces - Mantid

# Interfaces - DAWN

# Interfaces - DAWN

# **Populating ICAT with Data**

- ICAT has a SOAP API so data can be pushed in from most languages

- Typically metadata is imported from two sources
  - User office
  - Experiment data files

- User office link code tends to be bespoke

- Can be simple (experiment title, people) or more complex (abstracts, samples, links between related work etc)

**Science & Technology**
Facilities Council

# Populating ICAT with Data

- Data files
  - Tools exist to assist ingesting nexus files
  - Ingest custom/specific data formats will require code to extract the metadata from them
- Can be done using the API or 'XML Ingest' layer
- XML Ingest defines a simple schema that describes a set of datafiles and inserts them into an ICAT
- Prototype available, but still work-in-progress
- Will simplify data ingestion.

**Science & Technology**
Facilities Council

# The ICAT Project

# The ICAT Project

- ICAT is an open source project
- Currently managed by STFC
- Released under BSD and Apache licenses (permissive)
- [www.icatproject.org](www.icatproject.org)
- [http://code.google.com/p/icatproject](http://code.google.com/p/icatproject)
- Monthly collaboration meetings (telephone)
- Annual developer meetings (face-to-face)
- Steering committee

Science & Technology
Facilities Council

# Future Developments

- Under active development
- Releases every 6 months
- Current roadmap
  - Clustered deployment
  - Improved documentation
  - Data import/export/migration
  - Non-relational databases (hybrid)
  - Interface improvements to TopCAT
  http://icatproject.org/releases/road-map/
- Tell us what you want to see…..

**Science & Technology**
Facilities Council

# Summary

- ICAT is a mature solution for facility metadata management

- Enables remote data access  and linking between proposals, data, analysis and publications

- Flexible architecture enables  integration with differing facility systems and requirements

- API enables integration with other software (analysis)

- Open source project with an active collaboration

# Questions?

Thanks to: Steve Fisher, Brian Matthews, Alistair Mills, Alun Ashton, Antony Wilson

# Store.Synchrotron.org.au

**Steve Androulakis**

# Store.Synchrotron.org.au

• Store.Synchrotron is a system that captures all macromolecular beamline data, available online to all non-commercial Australian Synchrotron users. It was developed by Monash University in a strategic, ongoing partnership.

• Data is **immediately shareable** by the researcher on the web and able to be **published**.

• The service operates on the Australian NeCTAR Research compute cloud in a scalable setup able to withstand load. (http://nectar.org.au/)

• We're actively opening access to raw data behind high-impact research publications under CC BY licenses. Six institutions have opened data so far.

• Built on MyTardis – an open source, Australian made data management platform used all over Australia in proteomics, genomics, electron microscopy, medical imaging, astrophysics, quantum physics and more.

• Visit store.synchrotron.org.au for access / contact steve.Androulakis@monash.edu
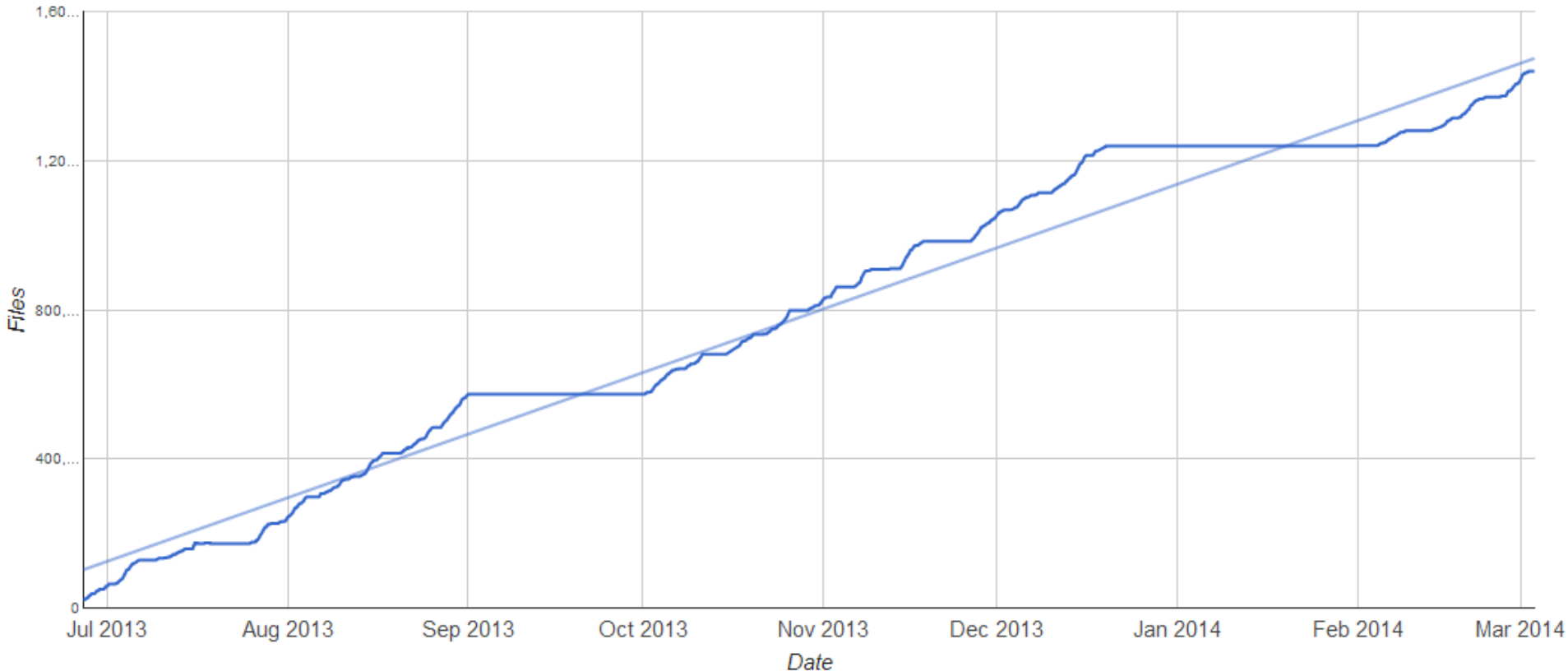
MONASH University     Australian Synchrotron     nectar

# Real-time instrument data capture



Capture began June 2013. As of March 2013, it has captured over 18 terabytes of data in over 1.5 million raw diffraction images.

Presented on the Synchrotron's Control PCs in a slideshow.

# How can we contribute to PaNSig?

• Store.Synchrotron.org.au is a scalable cloud service that was created to withstand network bandwidth, CPU, disk IO. We have considerable engineering experience in this area and can provide advice and guidance.

• We have worked with Synchrotron users in order to have their raw data released publicly (http://tardis.edu.au/syncpublish/) and have a basic model for this process.

• The Australian National Data Service (ANDS.org.au) is funding an 18 month project to refine the Store.Synchrotron.org.au process of collecting quality metadata at time of collection, and time of publication. Many of its outputs would be relevant to PaNSig.

• Monash University are in the process of increasing its 'instrument integration' of this kind using the same underlying system as Store.Synchrotron. This includes expansion to more beamlines at the Australian Synchrotron and likely the Australian neutron source ANSTO. Experiences here with different research communities, policies, practices, technology and culture should be broadly useful to this group.

# Developing a plan of activities for RDA PaNSig

RDA Plenary, Dublin, March 2014

# Co-chairs

Currently:

- Amber Boehnlein , SLAC, US
- Frank Schluenzen, DESY, DE
- Brian Matthews,  STFC, UK

# Topics for Discussion

- Purpose
- Technical activities: RDA WGs etc
- Interaction within RDA
- Future Meetings
- Interactions between meetings
- Members

# Purpose

- Data related issues of science applications associated with large scale source facilities
- Source facilities share a number of issues in their data handling.
- These could include:
  - Scalability of data volumes and data access rates
  - Standardization of (meta-)data and vocabularies
  - Data, cataloguing, publishing, discovery, sharing, transfer and access, policies
  - Data analysis tools and frameworks supporting workflows and provenance
  - Interaction with the data handling practices and standards within different communities.

# Technical Activities

- Spinning out RDA technical working groups
  - Specific
  - Time-bound (c. 18 months)
  - Clear outcome (document, format, code etc)
  - Need commitment
  - **Topics** ?


- Any documents as a group ?
  - capturing community views.
- Future projects?

# Interaction with other RDA groups

- We will be requested to comment on draft recommendations from other RDA working groups
  - Will notify group and request if people want to review documents
- Are there specific groups we want to interact with?
  - Domain specific groups: Structural Biology,  Materials Data
  - Technical infrastructure:  Persistent IDs, Big data Analytics, cloud computing, metadata, data publication, …..

# Future meetings

- RDA Plenary
  - 4th RDA Plenary: Amsterdam, 22-24 September, 2014
  - 5th RDA Plenary: USA, March 2015
- NoBugs 2014 ?
- Other events?
- Opportunities to establish a PaN Computing event similar to CHEP ??

# Interactions in between f2f

- Keeping up momentum is important

  - Website/Wiki (RDA Mediated)
  - mailing list (RDA Mediated)
  - phone/video meetings?

- Subgroups on particular topics?

# Other members

- Who should we be aiming at ?
- PaNS outside US/EU/AU ?

# Possible immediate activities?

- Online technical forum
- Counting users
- Collecting views data usage
- Project/software "market"
- Common terms
- Other RDA Groups: Structural Biology, Materials, Persistent ID …
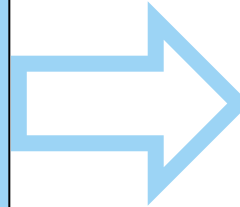- Meetings : RDA the Hague

# Other discussion?

# DataCite –
# Persistent links to scientific data

# DOI names for citations

**URLs are not persistent**

- (e.g. Wren JD: **URL decay in MEDLINE- a 4-year follow-up study**. Bioinformatics. 2008, Jun 1;24(11):1381-5).

**Digital Object Identifiers (DOI names) offer a solution**

- Mostly widely used identifier for scientific articles
- Researchers, authors, publishers know how to use them
- Put datasets on the same playing field as articles



The page cannot be found

The page you are looking for might have been removed, had its name changed, or is temporarily unavailable.

Please try the following:

- If you typed the page address in the Address bar, make sure that it is spelled correctly.
- Open the httpd.apache.org home page, and then look for links to the information you want.
- Click the ⇦ Back button to try another link.
- Click 🔍 Search to look for information on the Internet.

HTTP 404 - File not found
Internet Explorer

**Dataset**
Yancheva et al (2007). Analyses on sediment of Lake Maar. PANGAEA.
doi:10.1594/PANGAEA.587840

DataCite

# DataCite members

1. Technische Informationsbibliothek (TIB)
2. Canada Institute for Scientific and Technical Information (CISTI),
3. California Digital Library, USA
4. Purdue University, USA
5. Office of Scientific and Technical Information (OSTI), USA
6. Library of TU Delft, The Netherlands
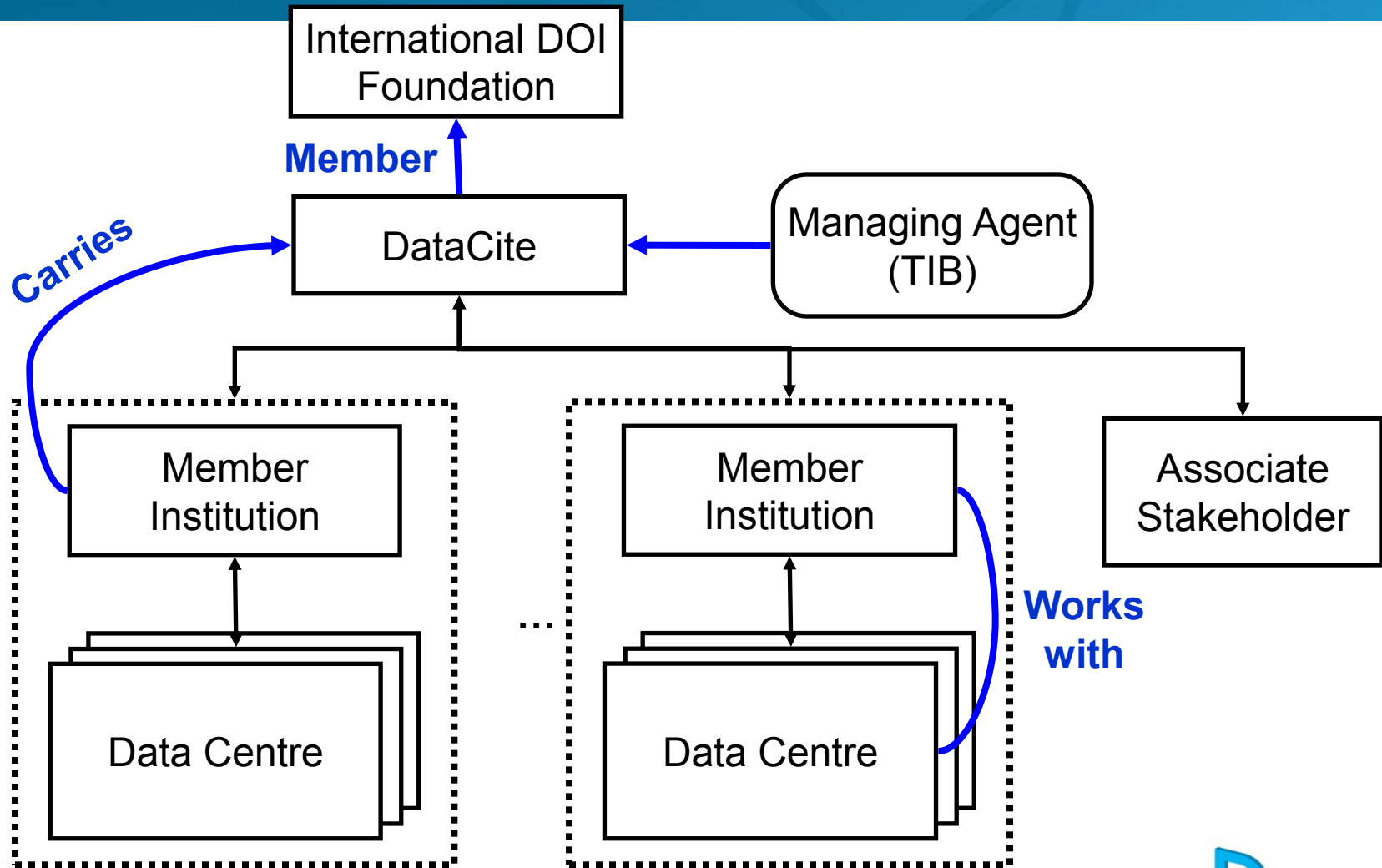7. Technical Information Center of Denmark
8. The British Library
9. ZB Med, Germany
10. ZBW, Germany
11. Gesis, Germany
12. Library of ETH Zürich
13. L'Institut de l'Information Scientifique et Technique (INIST), France
14. Swedish National Data Service (SND)
15. Australian National Data Service (ANDS)
16. Conferenza dei Rettori delle Università Italiane (CRUI)
17. National Research Council of Thailand (NRCT)
18. The Hungarian Academy of Sciences
19. University of Tartu, Estonia
20. Japan Link Center (JaLC)
21. South African Environmental Observation Network (SAEON)
22. European Organisation for Nuclear Research (CERN)

**Affiliated members:**

1. Digital Curation Center (UK)
2. Microsoft Research
3. Interuniversity Consortium for Political and Social Research (ICPS
1. Korea Institute of Science and Technology Information (KISTI)
5. Bejiing Genomic Institute (BGI)
6. IEEE
7. Harvard University Library
8. World Data System (WDS)
9. GWDG

# DataCite structure

# DataCite in 2014

Over 3,200,000 DOI names registered so far.

290 data centers.

8,000,000 resolutions in 2013.

DataCite Metadata schema published (in cooperation with all members) http://schema.datacite.org

DataCite MetadataStore
http://search.datacite.org

# OAI and Statistics

OAI Harvester

http://oai.datacite.org

DataCite statistics (resolution and registration)

http://stats.datacite.org

# ODIN project

ORCID and DataCite interoperability network. Funded under FP7

[http://www.odin-project.eu](http://www.odin-project.eu)

Claim your DataCite DOI with your ORCID profile:

[http://datacite.labs.orcid-eu.org/](http://datacite.labs.orcid-eu.org/)

# 2012: STM, CrossRef and DataCite Joint Statement

1. To improve the availability and findability of research data, the signers encourage authors of research papers to **deposit researcher validated data in <u>trustworthy and reliable Data Archives</u>**.

2. The Signers encourage Data Archives to **enable bi-directional <u>linking</u> between datasets and publications** by using established and community endorsed unique persistent identifiers such as database <u>accession codes</u> and <u>DOI's</u>.

3. The Signers encourage publishers and data archives to make visible or increase **<u>visibility of these links</u>** from publications to datasets and vice versa

DataCite

# Example

**The dataset:**

Storz, D et al. (2009):

*Planktic foraminiferal flux and faunal composition of sediment trap L1_K276 in the northeastern Atlantic*.

http://dx.doi.org/10.1594/PANGAEA.724325

**Is supplement to the article:**

Storz, David; Schulz, Hartmut; Waniek, Joanna J; Schulz-Bull, Detlef; Kucera, Michal (2009): *Seasonal and interannual variability of the planktic foraminiferal flux in the vicinity of the Azores Current.*

Deep-Sea Research Part I-Oceanographic Research Papers, **56(1),** 107-124,

http://dx.doi.org/10.1016/j.dsr.2008.08.009

DataCite

# Cooperation

MoU with ORCID

Agreement with Re3Data and DataBib to include their service in 2016

MoU with RDA to become organisational affiliate

Joint Declaration of Data Citation Principles

https://www.force11.org/datacitation

# 2014 Annual conference

## DataCite Annual Conference 2014
### 25-26 August 2014, Inist-CNRS, Nancy, France
(Just after the IFLA World Library and Information Congress in Lyon)

## Giving value to data:
## advocacy, guidance, services

**Coming to Nancy:**

- *by train:*
1h30 from Paris
4 hours from Lyon

- *by plane:*
direct link from Paris Airport

Paris    Nancy

Lyon

DataCite

# Thank you!

# Data analysis issues and frameworks

Alun Ashton

Group Leader, Data Analysis Software

# Diamond Light Source



- EM Facility
- SFX XFEL UK Hub

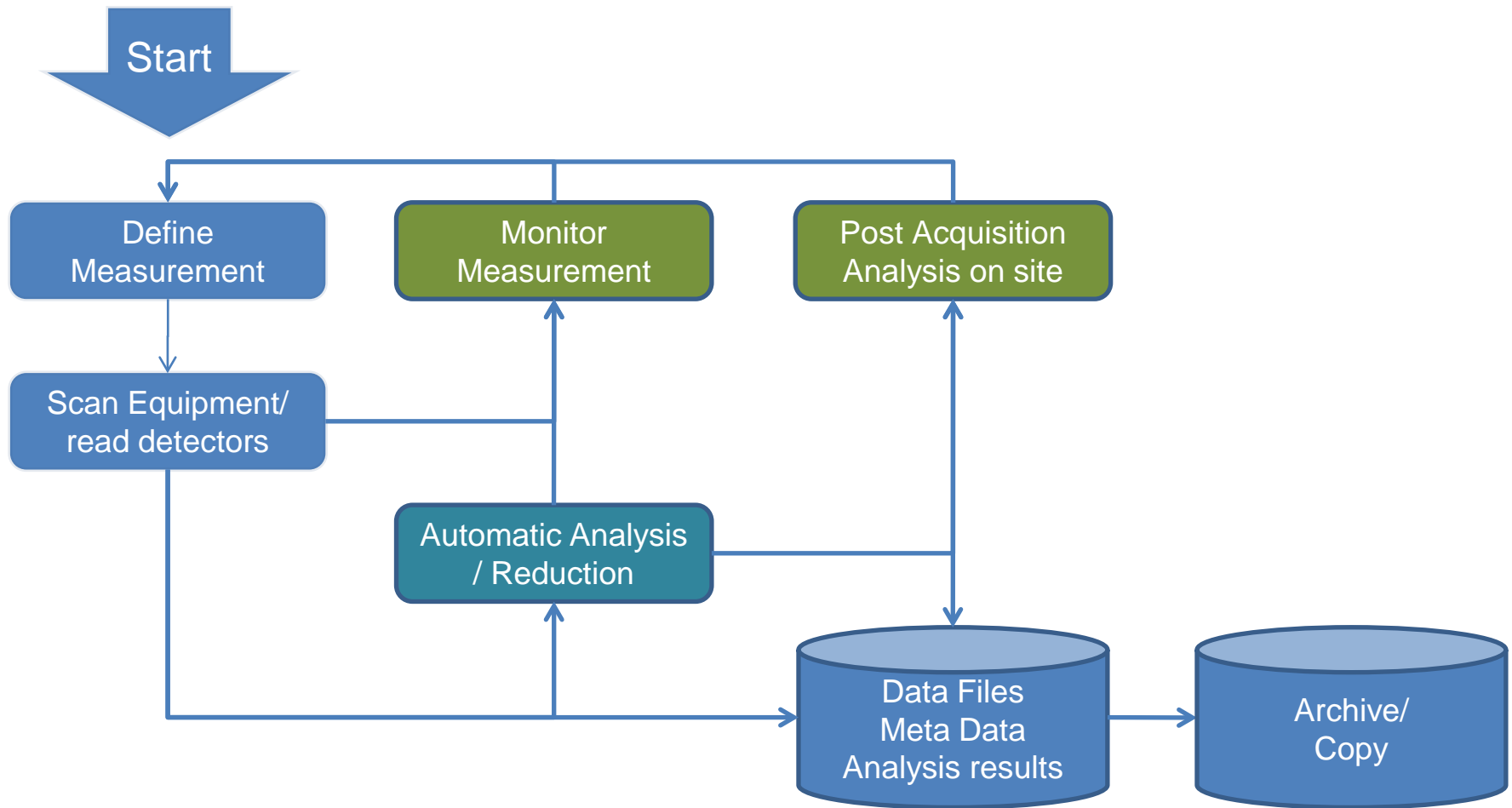# Data Analysis – Experimental Challenges Diversity: 33 beamlines > 120 experiment techniques

# Data Analysis – Experimental Challenges

| BEFORE | IMMEDIATE | SHORT TERM | LONG TERM |
|---|---|---|---|
| (FROM DLS/USER'S INSTITUTION) | (DURING EXPERIMENTS) | (BEFORE THE USER GOES HOME) | (FROM DLS/USER'S INSTITUTION) |
| Simulations<br><br>Processing of older datasets | "Real time" data processing, analysis and visualisation – to make experimental decisions | Data reduction and processing – Users go home with clean data free of instrument artefacts.<br><br>Preliminary data analysis – helpful, but may require significant processing power and know-how | Detailed analysis – from data to information.<br><br>Incorporating results from other techniques.<br><br>Experiments:<br><br>➢ Provide parameters for a model.<br><br>➢ Test/verify a model or theory.<br><br>➢ Show where a new theory or model is required. |

# Data Analysis – Experimental Challenges

# The 'Customer'

# Data Analysis – Experimental Challenges

How facility users want to analyse their data (sample):

- Command line, interactive analysis
- 'scripts'
- Black box
- Wizards
- GUI
- Automatically
- Someone else do it
- Don't you look at it

- Quickly
- Estimates
- Properly
- Publication quality
- The way I used to do it
- Use the new stuff
- On all the data
- On some of the data

diamond

# Data Analysis – Experimental Challenges

How facility users want to analyse their data (sample):

- Only when I am at the facility
- When I am at home
- On a web page
- On my laptop
- On your computers

- Yesterday
- Next year
- By the student
- By the expert
- Securely
- Shared data
- ASCII not binary...

diamond

# Other Factors

- Existing codes and expertise
- Plethora of data structures/file formats/data rates
- Available computing resources
  - Enough computing
  - Quick access to Data
- Information management
- Collaborations

diamond

# How we are dealing with the problem in Diamond

examples

diamond

# Beamline Computing/Software Support

- Data Analysis Group: – Alun Ashton
- Data Acquisition Group: – Paul Gibbons
- Scientific Computing Team: - Greg Matthews
- Beamline Controls Group: - Nick Rees
- User Office development team: – Bill Pulford, Ben Peacock

# Acquisition and Analysis in the Eclipse Framework
## (www.eclipse.org)

**GDA** software for science

**DAWN** SCIENCE

Client server technology

Communication with EPICS and hardware

Scan mechanism

www.opengda.org

## Acquisition

Jython and Python

Visualisation

Communication with external analysis

Analysis tools

Data read, write, convert

Metadata structure

Workflows

www.dawnsci.org

## Analysis

DAWN is a collection of generic and bespoke 'views' collated into 'perspectives'.

The perspectives and views can be used in part or whole in either the GDA or DAWN.

All core technologies open source

**diamond**

An open source collaboration with a core of generic tools with local extensions and implementations.

- Visualisation of scientific data
- Workflow tool for the treatment/manipulation of scientific data
- Integrated python environment

# Image stack analysis



Thanks to: Gareth Nisbet, I16

# 2D Diffraction Processing

# BioSAXS experiment

# SAXS Data Reduction

# Automated Data Reduction

# Tomography Reconstruction GUI

# ARPES Angle-Resolved Photo-Emission Spectroscopy

# ARPES

# ARPES

# File Formats x n

Diamond has a policy of, where feasible, to standardise on file formats, the choice being NeXus/HDF5

**Green:** predominantly using NeXus.

**Orange:** Mixed NeXus and other formats or considering NeXus in the next 12 months.

Beamline labels around the ring: I02, I03, I04-1, I04, I05, I06, I07, I09, I10, I11, I12, I13, I15, I16, B16, I18, B18, I19, I20, B21, I22, B22, B23, I24

Files can be generated by Detector, EPICS or Data Acquisition

pandata  NeXus  IUCr

# ISPyB information Flow: Display of Results

Data tracking and status checking of data acquired:

- Automatic plot generation of Scattering curves, Kratky, Guinier and P(r) plots

| Macromolecule | Concentration | Scattering | Kratky | Guinier | P(r) | Frames (Averaged/Total) | | Guinier | | | | Gnom | | Porod | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Rg (nm) | Points | Quality (%) | I(0) | Rg (nm) | Dmax (nm) | Volume (nm3) | MM (kD) Vol. est. |
| bsa | 5.80 mg/ml | | | | | bbsa (10 of 10) bsa (10 of 10) bbsa (10 of 10) | | 3.13 nm | 44 - 87 (43) | 78.37 | 68.25 ± 4.964017e-2 | 3.20 nm | 10.94 | 118.64 | 59.3 - 79.1 |

- Display of 1D curves, averaged curves and subtracted curves

- 3D Models displayed using WebGL



- Concentration Effects

A prototype version of ISPyB for BioSAXS (ISPyBB) database is now available at BM29@ESRF and is being ported to P12@EMBL

**Visits** | Experiments (i02) | Cancel     **Visits** | Experiment details (nt5814-2, i02)

🔍 lysozyme ⊗

| | | |
|---|---|---|
| 🟡 | **LysoCT200_99bcryo** | 1.66 Å |
| | 120511_hc1_work/lysozymeCT200/ | |
| AM | P1 | 40 x 77 x 78 Å |

| | | |
|---|---|---|
| 🟢 | **LysoCT200_99bcryo** | 1.66 Å |
| | 120511_hc1_work/lysozymeCT200/ | |
| AM | RSymm = 0.115 | I/Sigma = 20.2 |

| | | |
|---|---|---|
| 🟢 | **LysoCT200_80cryo** | 1.66 Å |
| | 120511_hc1_work/lysozymeCT200/ | |
| AM | RSymm = 0.070 | I/Sigma = 24.3 |

| | | |
|---|---|---|
| 🟢 | **LysoCT200_80cryo** | 2.20 Å |
| | 120511_hc1_work/lysozymeCT200/ | |
| AM | P4 | 76 x 76 x 37 Å |

| | | |
|---|---|---|
| 🟢 | **LysoCT200_80cryo** | 2.20 Å |
| | 120511_hc1_work/lysozymeCT200/ | |
| AM | P4 | 77 x 77 x 37 Å |

| | | |
|---|---|---|
| 🟢 | **LysoCT200_80cryo** | 2.20 Å |
| | 120511_hc1_work/lysozymeCT200/ | |
| AM | P4 | 78 x 78 x 38 Å |

| | |
|---|---|
| Omega end | 300.0° |
| Rotation per image | 0.30° |
| **Exposure time** | 1.000 s |
| **Beamsize X** | 50.0 μm |
| **Beamsize Y** | 25.0 μm |
| **Transmission** | 20.0% |
| **Resolution** | 1.7 Å |

**Autoprocessing results**

🟢    **Autoprocessing successful**

**1 - fast_dp results** ❯
/dls_sw/apps/fast_dp_refactor/svn/src/fast_dp.py -j 0 -J 36 /dls/i02/data/2012/nt5814-2/120...

**2 - xia2 results** ❯
xia2 -min_images 3 -3dii -atom s -blend -project nt5814v2 -crystal LysoCT20099bcryo1 -is...

**Comments**

(2247,9,164) EDNAStrategy1: subWedge:1Aperture: Large

**Associated images**

Visits      Status      Settings

# Archive and Reprocessing Service

# Collaboration experience

- Collaborating Inside the facility is often as challenging as with external groups.

- 'customer' involvement essential

- Management enthusiasm and support essential

- The scope of a collaborative project must fit in with the ambitions of your facility.

- Collaborating is expensive and ideally all parties must have adequate and commensurate resources though a centre of mass can produce results quicker, but at the cost of less well resourced partners.

- Although service providers are often happy to collaborate, our customers might not....

diamond

# Current status/volumes in Diamond ICAT = 285,198,074 files



**Diamond Total Data**

This image shows the diffraction pattern of an RNase 4 crystal at 1.6A, from Prof K. Ravi Acharya's group (at University of Bath) - the image that took DLS total data catalogued and archived > 1 Petabyte.

10/03/14
1,020,040

# Novel developments in HDF5

**An overview on new developments in the HDF5 library**

Eugen Wintersberger
Dublin, 28.03.2014

# HDF5 key features

→ Transparent compression of individual datasets

→ platform independent (Windows, Linux, AIX, Solaris, OSX, OpenVMS)

→ library is implemented in C (easy to interface)

→ bindings to many languages: C++, Fortran, Python, Java, Perl, R, Go

→ Supported by commercial products like: Matlab, Mathematica, IDL

**Easy parsing:** data access via library → it is not about how data is stored but rather where (path).

# Recently developed features already available in 1.8.12

# How compression works in HDF5

## Writing compressed data



Compression algorithms
implemented as shared libraries.

## Reading compressed data

# Writing precompressed data ...

Sometimes compression can be done better in hardware (FPGAs)
→ reduced bandwidth requirements for network and disk I/O!



Feature funded
by DECTRIS

This features allows bypassing the filter chain and write compressed data directly.

**Reading the data is the same as for conventionally written data!**

# External filter API

Detector vendors or facilities may want to use custom compression algorithms.

Possible with all HDF5 1.8.X versions but: reading software needs to be recompiled => commercial applications do not have access to the data.



Feature funded by DESY

**External filters can be loaded at runtime on demand.**

**Requirement:**
Detector vendors and facilities provide libraries for all target platforms.

# New features addressing concurrent data access

## Single Writer – Multiple Reader = SWMR

**Current situation**: cannot read data from a file while it is written by two different processes!

### SWMR would make this possible!



**Status**: A prototype is available – missing funding for library integration.

# Multithreading support

## Option 1: Make HDF5 fully threadsafe!

➔ would give best performance

➔ major rewrite of 300K lines of C-Code

➔ costs of 4-6 FTEs

➔ significant future maintenance efforts

**Very unlikely to happen without extensive funding!**

## Option 2: Use threads within the library!

➔ use multiple threads for compression

➔ asynchronous (non-blocking) I/O

➔ costs ~1.5 FTEs

➔ lower maintenance effort

➔ could improve continuously

DESY

# Remote access to data stored in an HDF5 file

## Open Source Project for a Network Data Access Protocol

**Advantages**

➜ Already existing and actively developed protocol

➜ Data access via HTTP requests where URLs encode the information of what shall be retrieved from a data source

➜ Independent of the dataformat (other formats can easily be supported)

➜ Data can be accessed with any software that can dereference URLS (webbrowser, Excel, Libreoffice,...)

➜ Matlab 2012a provides native support for OPeNDAP

➜ C++, Java, and Python libraries available

www.opendap.org

# HDF5 data server

Serialization of HDF5 library calls rather than using HTTP and URLs



**Status:** only as a design draft – no implementation will be done without funding.

# Licenses and funding issues

**Currently the HDF5 group is entirely project funded!**

… they are looking for new funding (licensing) models

➔ the core library will always remain open source

➔ dual licensing would be an alternative

➔ some features may become commercial products in future

**As a user community we should take funding issues of the HDF5 group serious and think about how we can improve their financial backup!**

# Conclusions: what the RDA could do

➜ Establish an HDF5 special interest group as a hub for HDF5 related work

➜ collect feature requests from various user communities

➜ organize funding activities if a new feature should be implemented

➜ manage filters (compression algorithms) which should be available for HDF5 and provide hosting resources

➜ Provided solutions based on HDF5 for certain selected use-cases

➜ Organize workshops around HDF5 (maybe as satellite events around meetings)

Brian Matthews

STFC

# PaNdata

- Photon and Neutron Data Infrastructure
- Established in 2007 with 4 facilities
  - now standing at 13
  - With "friends" around the world
- Combined Number of Unique Users more than 35000 in 2011
- Combines Scientific and IT staff from the collaborating facilities
- European Framework 7 Projects
  - PaNdata-Europe: SA, 2009-11
  - PaNdata-Open Data Infrastructure, IP, 2011-14

  – Guestimates
  - Investment > €4.000.000.000*
  - Running costs > €500.000.000/yr*
  - Publications > 10.000/yr*
  - RCosts/Publication ~ €50.000*%
  - Data volume >> 10PB/yr*

# Counting Users

| | ALBA | BER II | DESY | DLS | ELETTRA | ESRF | FRM-II | ILL | ISIS | LLB | SINQ | SLS | SOLEIL | neutron | photon | all |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Number of Users shared between facilities | | | | | | | | | | | | | | | | |
| ALBA | 773 | 7 | 61 | 58 | 51 | 281 | 2 | 51 | 13 | 5 | 10 | 77 | 105 | 69 | 400 | 773 |
| BER II | 7 | 1563 | 115 | 46 | 27 | 179 | 157 | 383 | 198 | 98 | 191 | 62 | 36 | 580 | 329 | 1563 |
| DESY | 61 | 115 | 4197 | 137 | 222 | 851 | 116 | 255 | 113 | 62 | 95 | 315 | 188 | 469 | 1294 | 4197 |
| DLS | 58 | 46 | 137 | 4407 | 102 | 810 | 30 | 267 | 399 | 33 | 52 | 229 | 192 | 546 | 1130 | 4407 |
| ELETTRA | 51 | 27 | 222 | 102 | 3167 | 433 | 11 | 77 | 35 | 20 | 18 | 179 | 367 | 141 | 900 | 3167 |
| ESRF | 281 | 179 | 851 | 810 | 433 | 10287 | 139 | 900 | 369 | 190 | 174 | 963 | 1286 | 1313 | 3586 | 10287 |
| FRM-II | 2 | 157 | 116 | 30 | 11 | 139 | 1095 | 347 | 137 | 89 | 161 | 33 | 29 | 509 | 259 | 1095 |
| ILL | 51 | 383 | 255 | 267 | 77 | 900 | 347 | 4649 | 731 | 301 | 395 | 156 | 222 | 1518 | 1347 | 4649 |
| ISIS | 13 | 198 | 113 | 399 | 35 | 369 | 137 | 731 | 2880 | 89 | 233 | 94 | 56 | 936 | 745 | 2880 |
| LLB | 5 | 98 | 62 | 33 | 20 | 190 | 89 | 301 | 89 | 1235 | 74 | 39 | 151 | 391 | 323 | 1235 |
| SINQ | 10 | 191 | 95 | 52 | 18 | 174 | 161 | 395 | 233 | 74 | 1219 | 224 | 31 | 590 | 415 | 1219 |
| SLS | 77 | 62 | 315 | 229 | 179 | 963 | 33 | 156 | 94 | 39 | 224 | 3827 | 399 | 371 | 1470 | 3827 |
| SOLEIL | 105 | 36 | 188 | 192 | 367 | 1286 | 29 | 222 | 56 | 151 | 31 | 399 | 4568 | 394 | 1817 | 4568 |
| neutron | 69 | 1563 | 469 | 546 | 141 | 1313 | 1095 | 4649 | 2880 | 1235 | 1219 | 371 | 394 | 10023 | 2334 | 10023 |
| photon | 773 | 329 | 4197 | 4407 | 3167 | 10287 | 259 | 1347 | 745 | 323 | 415 | 3827 | 4568 | 2334 | 25336 | 25336 |
| all | 773 | 1563 | 4197 | 4407 | 3167 | 10287 | 1095 | 4649 | 2880 | 1235 | 1219 | 3827 | 4568 | 10023 | 25336 | 33025 |

http://pan-data.eu/Users2012-Results

# PaN-Data Integration

*Common data environment, common user experience*

## Shared Data Policy Framework

Federated User Authentication 

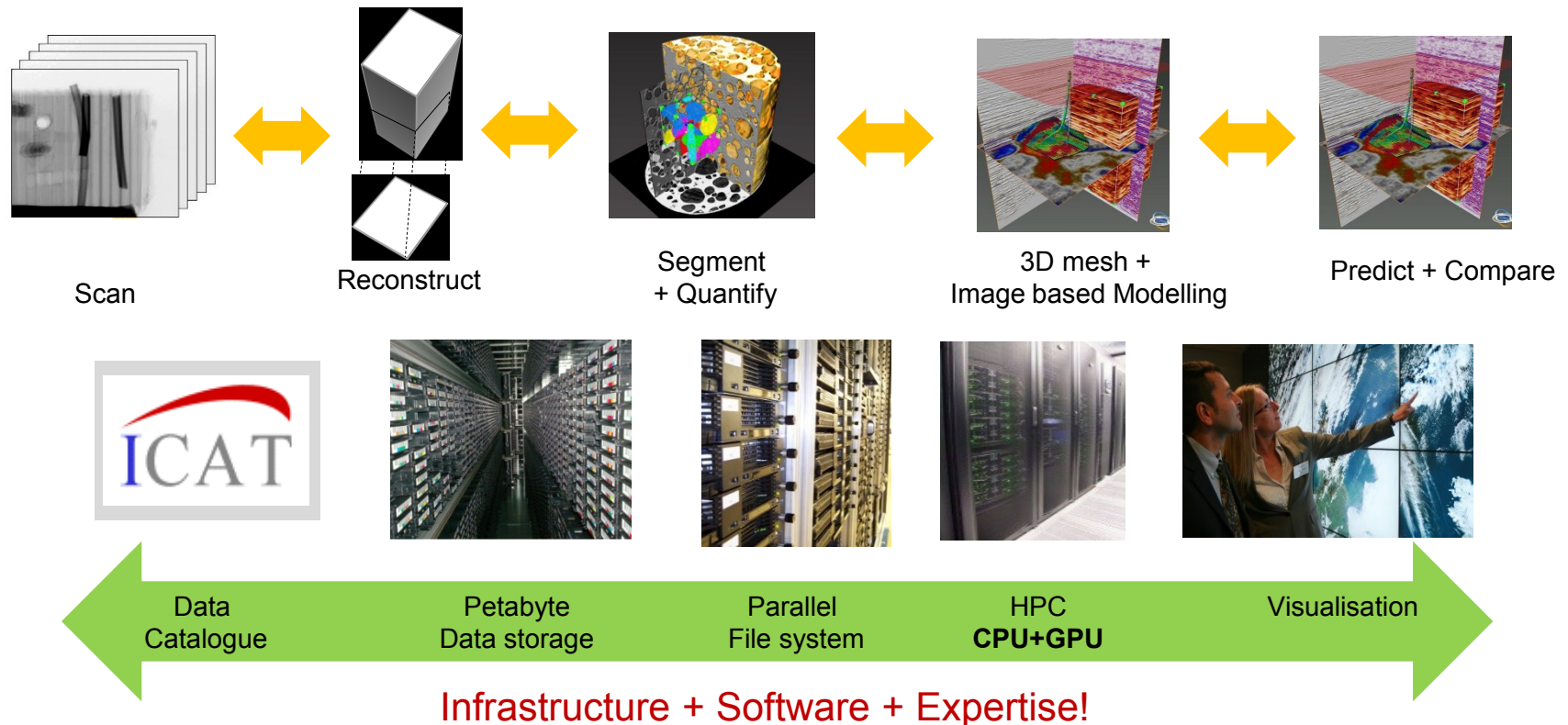Federated Data Catalogue 

Common Data Format 

# Towards the Future

- ## Provenance

  - Integrating context, analysis and publication into the record

- ## Preservatiom

  - Long-term need for archiving and curating data
  - Persistence Identifiers, itegtrity, context,
  - Costs and Benefits of data preservation

- ## Scalability

  - Managing high data rates and volumes
  - Parallel file stores

# Infrastructure for managing data flows



Scan

Reconstruct

Segment + Quantify

3D mesh + Image based Modelling

Predict + Compare

Data Catalogue

Petabyte Data storage

Parallel File system

HPC **CPU+GPU**

Visualisation

Infrastructure + Software + Expertise!

- **Tomography**: Dealing with high data volumes – 200Gb/scan, ~5 TB/day (one experiment)
- **MX**: high data volumes, smaller files, but a lot more experiments
- Hard to move the data – needs to be handled at the facility?