

WP6: JRA on Provenance

Brian Matthews

Scientific Computing Department, STFC, Didcot, UK

brian.matthews@stfc.ac.uk

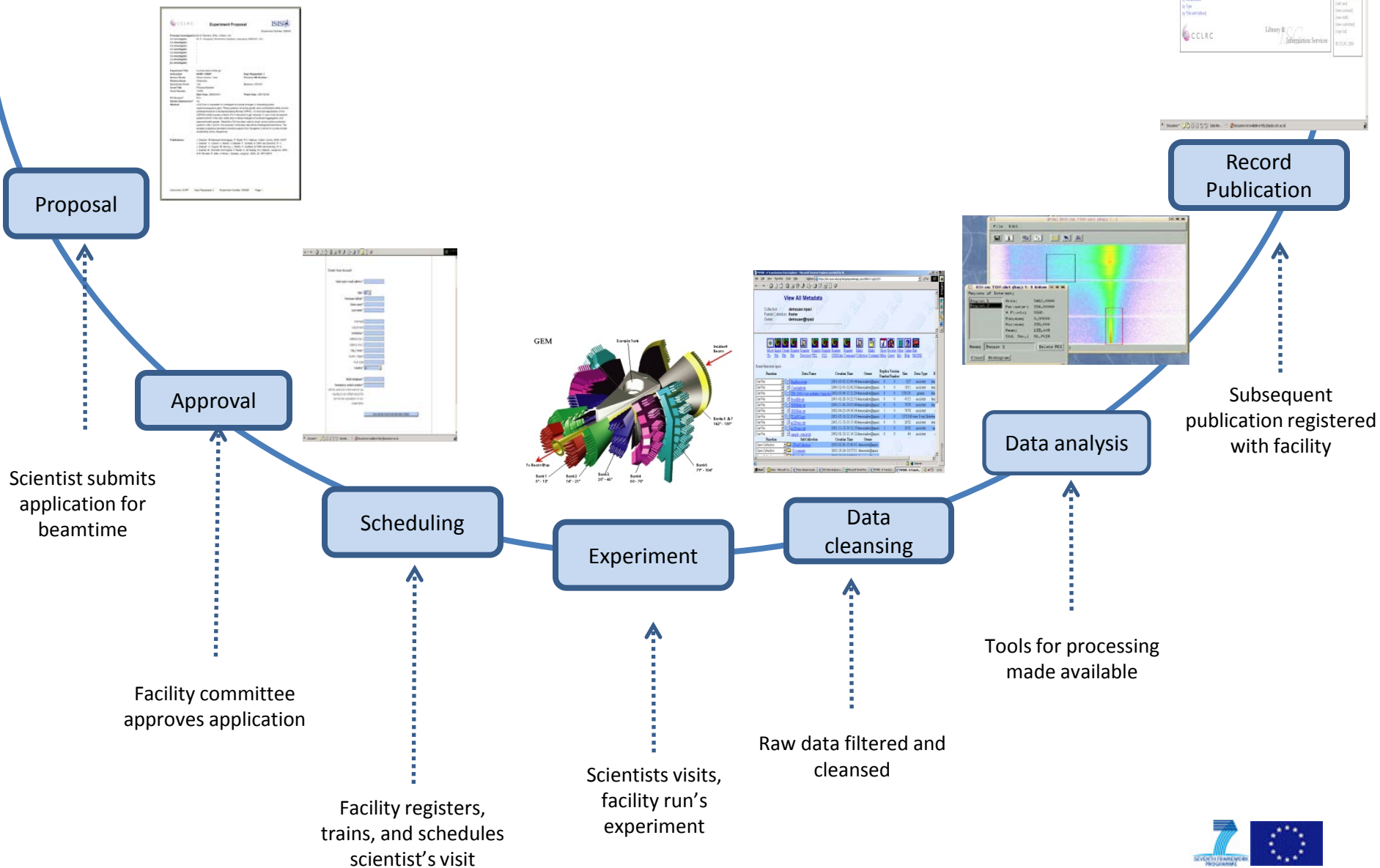
To research and develop a conceptual framework, defined as a metadata model, which can record the analysis process, and to provide a software infrastructure which implements that model to record analysis steps hence enabling the tracing of the derivation of analysed data outputs.

- To develop a conceptual framework, which can record and recall the data continuum.
- To provide a software infrastructure which implements that model to record analysis steps
- In general
 - A large and complex task
 - Establishing Science benefit
- Start M7 (April 2012), Finish M30

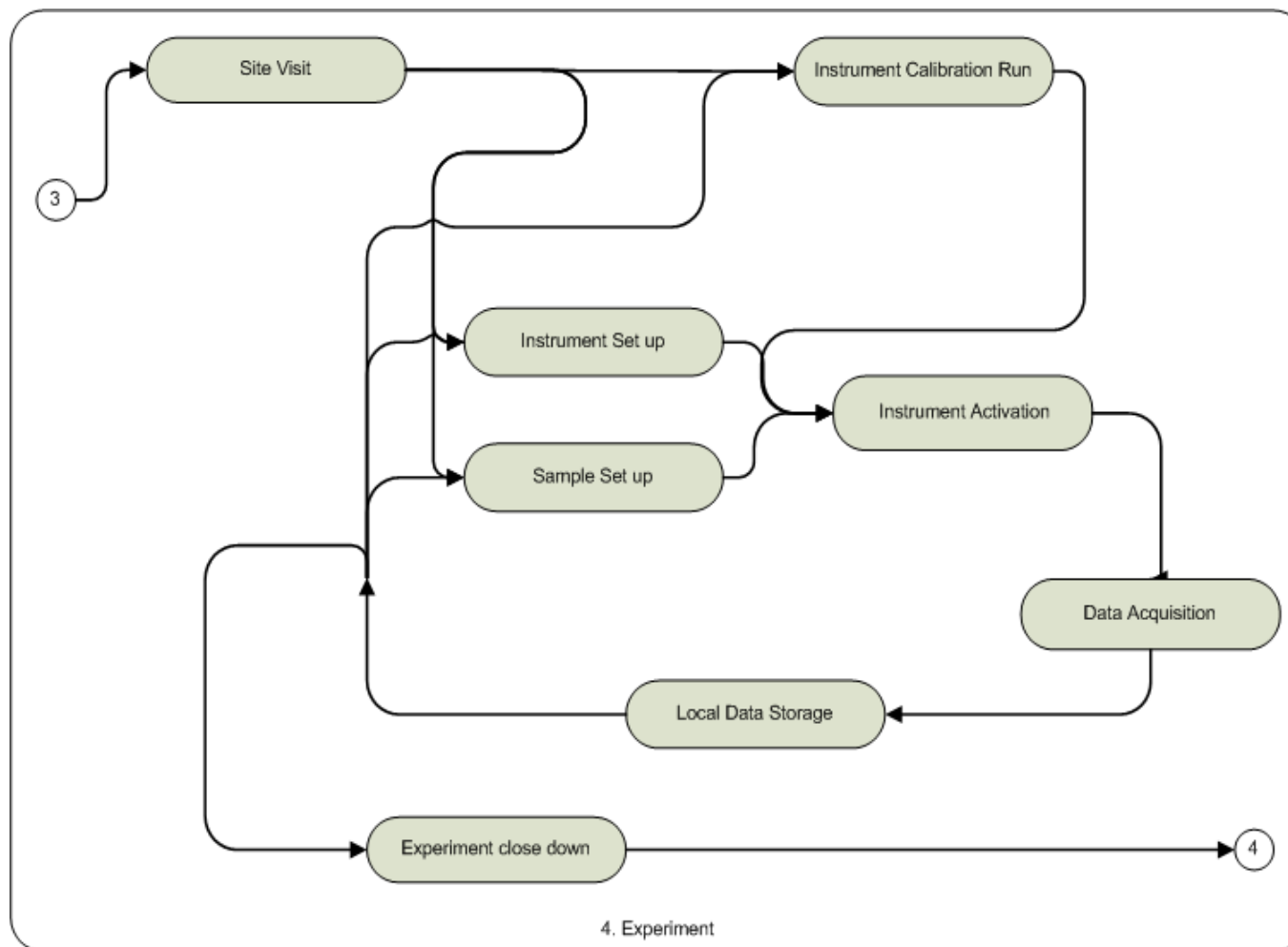
— STFC (Lead)	18 SM
— ILL	6 SM
— ELETTRA	12 SM

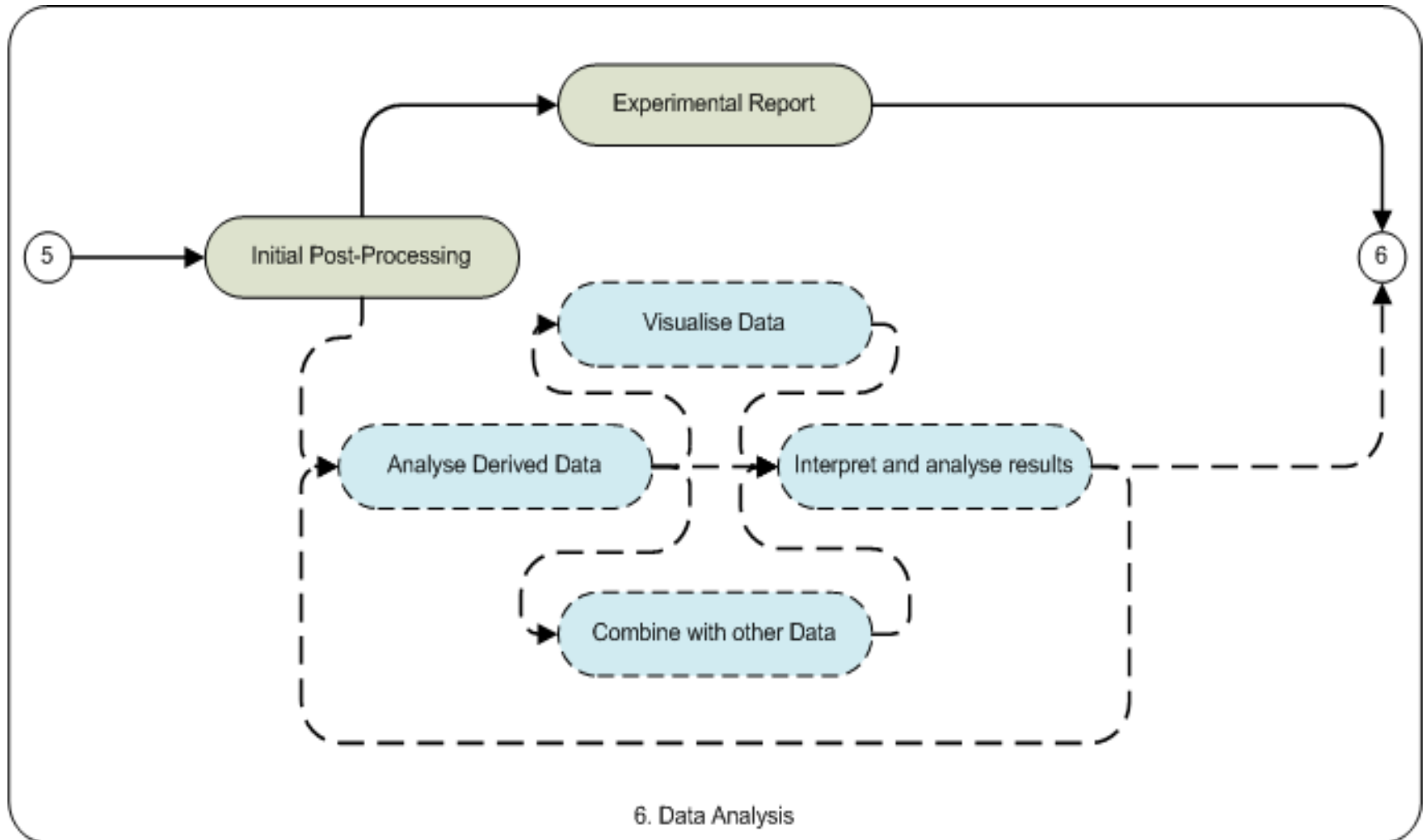
- Task 1: Requirements for Provenance
 - Task 2: Modelling the data continuum
 - Task 3: Ontologies for specific instruments/techniques
 - Task 4: Tool Support for the Data Continuum
 - Task 5: Tracing the Data Continuum
 - Task 6: Evaluation
- D6.1: Model of the data continuum in Photon and Neutron Facilities (M12)
 - D6.2: Common ontology definition and definition of tools to support the use of provenance for Photon and Neutron Facilities (M18)
 - D6.3: Tools for building research objects in Photon and Neutron Facilities (M24)
 - D6.5: Evaluation report on provenance management in Photon and Neutron Facilities (M30)

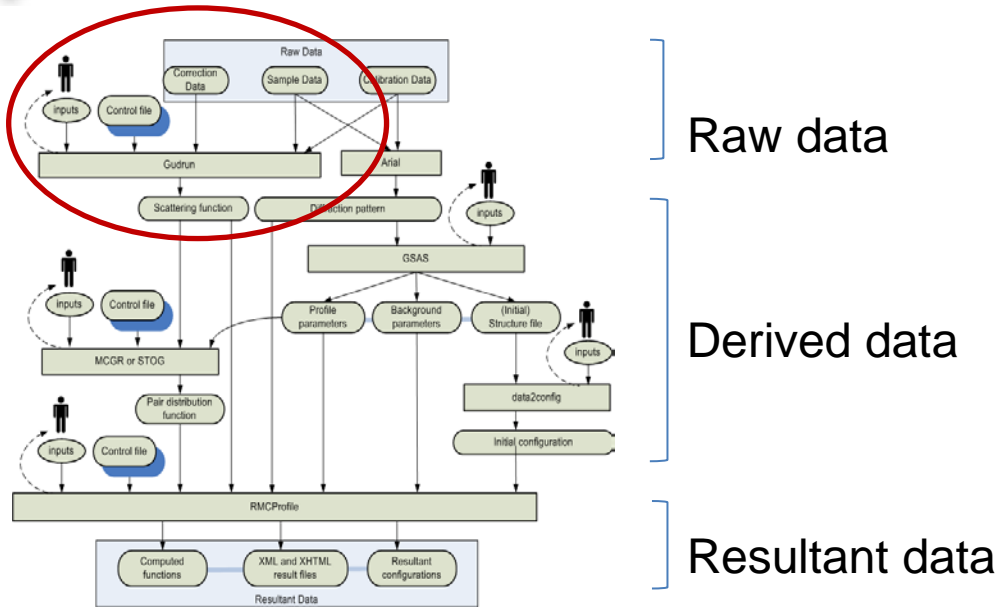
- Define Data Continuum
 - Define the stages of the facilities process we would need to support for end-to-end coverage
 - Consider who is involved
 - Consider the data and metadata involved
- Costs and Benefits
- Use cases
 - first sketch
 - Derive requirements



Stage	Actors	Information Systems	Metadata Types
1. Proposal	User Office, Principal Investigator	User office systems, User registration and management, User identity, Proposal systems	User identity, funding sources, institutional information, project description, experiment description, prior art
2. Approval	User Office, Approval Panel	User Office Systems, Proposal Systems	User identity, funding sources, experiment description prior art
3. Scheduling	User Office, Experimental Team, Instrument scientist	User Office Systems, Proposal Systems, H&S systems, Scheduling, Sample tracking	User identity, Sample information, Instrument information, Experiment planning
4. Experiment	Experimental Team, Instrument scientist	Sample tracking, Instrument control, Environmental monitoring, Data Acquisition systems, Data Management systems	User identity, Sample information, Instrument information, Experiment planning, Environmental parameters Calibration information
5. Data Storage	Scientific Team, Instrument scientist, Data infrastructure team	Data Acquisition systems, Data Management systems, Data storage systems, Archival Systems	User identity, Data formats, Data set information, File identifiers
6. Data Analysis	Scientific Team, Instrument scientist,	Data Storage systems, analysis software packages,	User identity, Data formats, Data set information, File identifiers Instrument parameters ...
7. Publication	Scientific Team, Instrument scientist, User Office, Library	User office systems Research Output tracking systems Library systems	User Identity Proposal information Publication information







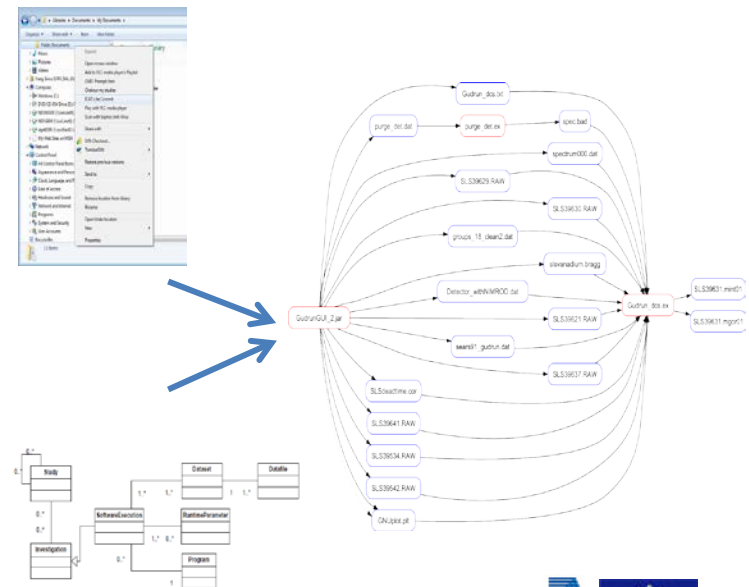
Issues:

1. Complexity of dependencies
2. Data volumes
3. Valuable data amongst noise
4. Software tracking
5. Distributed analysis
6. Workflow
7. Integration of tools

Credits: Martin Dove, Erica Yang (Nov. 2009)

Researchers are hesitated to change their well established software/practice.
“Why would I change?”

Need to demonstrate the benefits!



- Analyse practical situations where we can gain benefit from recording provenance
- A number of use cases
 - SANS2d (ISIS)
 - Co-ordinating and automating “near to experiment” processes
 - TwinMic x-ray spectromicroscopy beam line (ELETTRA)
 - Co-ordinating multiple experimental runs
 - Tomography beam lines (DLS)
 - Managing high volumes of data
 - Express services (ISIS)
 - Automating “standardised” experiments.
 - Publication linking (ISIS)
 - Tracing research outputs for impact analysis
- Would explore some of these in more depth in the project – related to V Labs

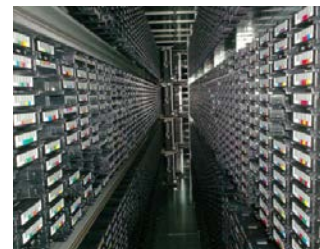
OpenGenie Script

ID	Name	URI	Concentration	Condition
1	Sample 1	uri1	0.1 M	25 °C
2	Sample 2	uri2	0.2 M	25 °C
3	Sample 3	uri3	0.3 M	25 °C
4	Sample 4	uri4	0.4 M	25 °C
5	Sample 5	uri5	0.5 M	25 °C
6	Sample 6	uri6	0.6 M	25 °C
7	Sample 7	uri7	0.7 M	25 °C
8	Sample 8	uri8	0.8 M	25 °C
9	Sample 9	uri9	0.9 M	25 °C
10	Sample 10	uri10	1.0 M	25 °C

OpenGenie Script



Data Acquisition



Data Archive

SampleTracks

Sample Information

Mantid
Data Processing

raw data

ICAT
(Extended) ICAT
Data Catalogue

DOIs



British Library
DOI Server

- Small Angle
Scattering VLab

- WP4, WP6, WP7

Outputs

derived data

New links



Science & Technology
Facilities Council

ePublication Archive

Log In

Erica Y Yang

aka Y Yang
aka Erica Yang
aka E Yang
erica.yang@stfc.ac.uk

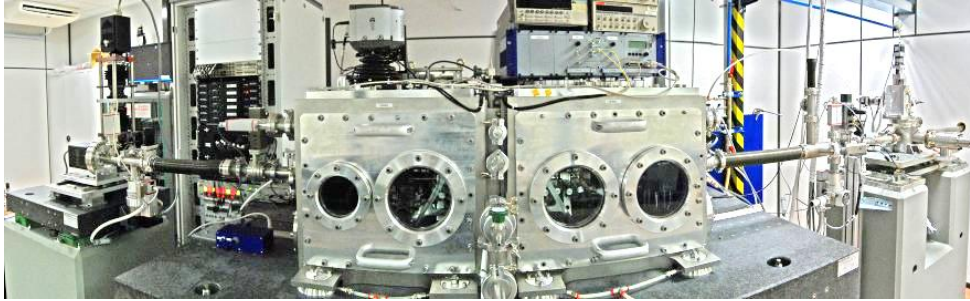
[printable version]

Order by: Year | Title | Author

1 to 14 of 14 | Data Processing | Sort | Edit

1. Y. Yang (STFC) Subscribed Application Log...
2. Y. Yang (STFC) Subscribed Application Log...
3. Y. Yang (STFC) Subscribed Application Log...
4. Y. Yang (STFC) Subscribed Application Log...
5. Y. Yang (STFC) Subscribed Application Log...
6. Y. Yang (STFC) Subscribed Application Log...
7. Y. Yang (STFC) Subscribed Application Log...
8. Y. Yang (STFC) Subscribed Application Log...
9. Y. Yang (STFC) Subscribed Application Log...
10. Y. Yang (STFC) Subscribed Application Log...

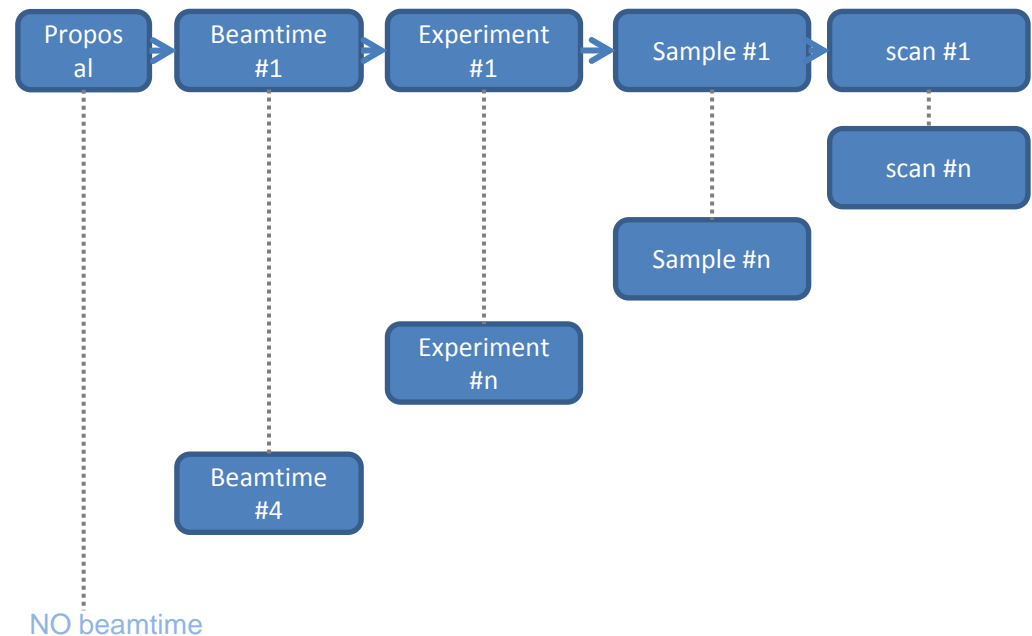
Publications

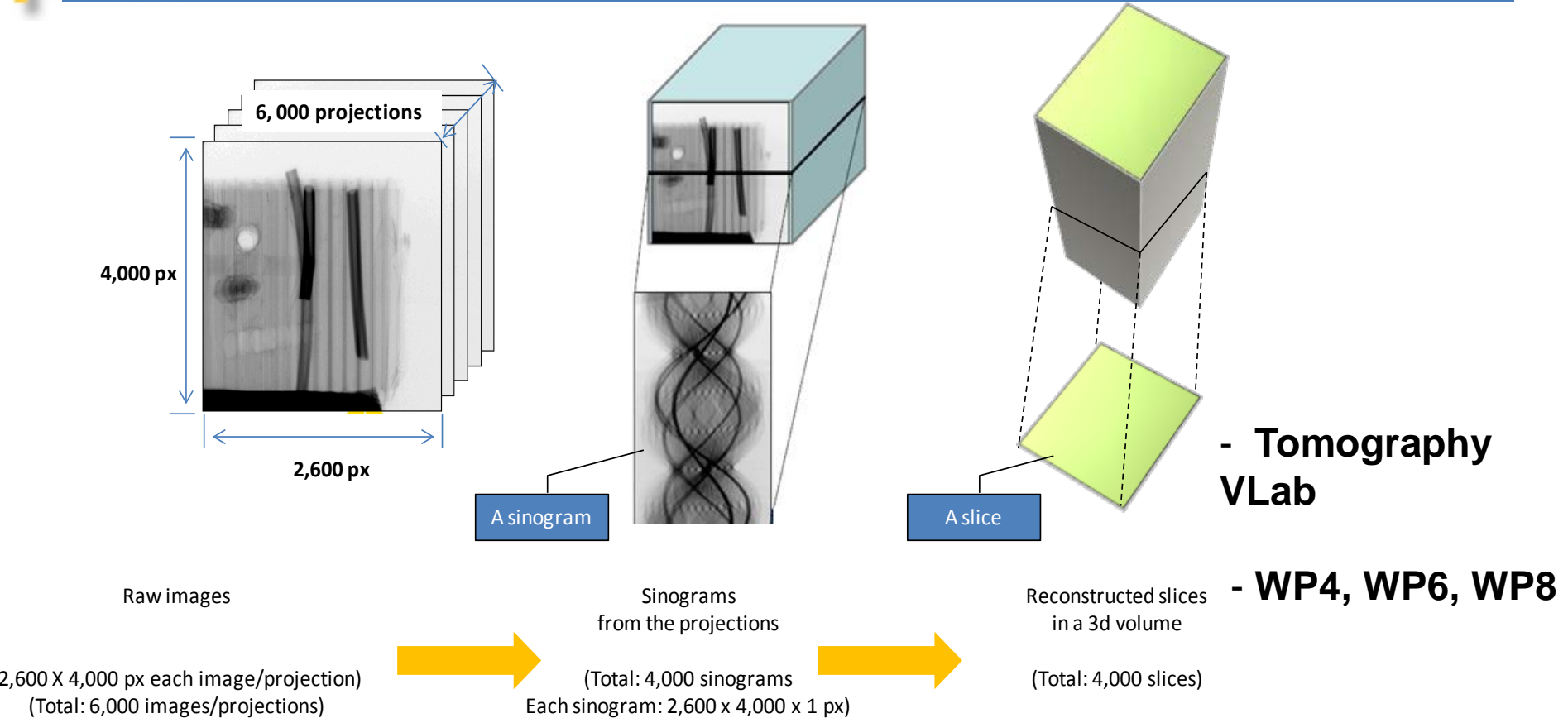


Beam-time Proposal has lots of meta-information

- 1 LT proposal may grant up to 4 beam-times

- x-ray spectromicroscopy
- 1 instrument → 2 modes: Scanning & Full field
- **Modified ICAT to support multiple beam-times**
 - WP4, WP6





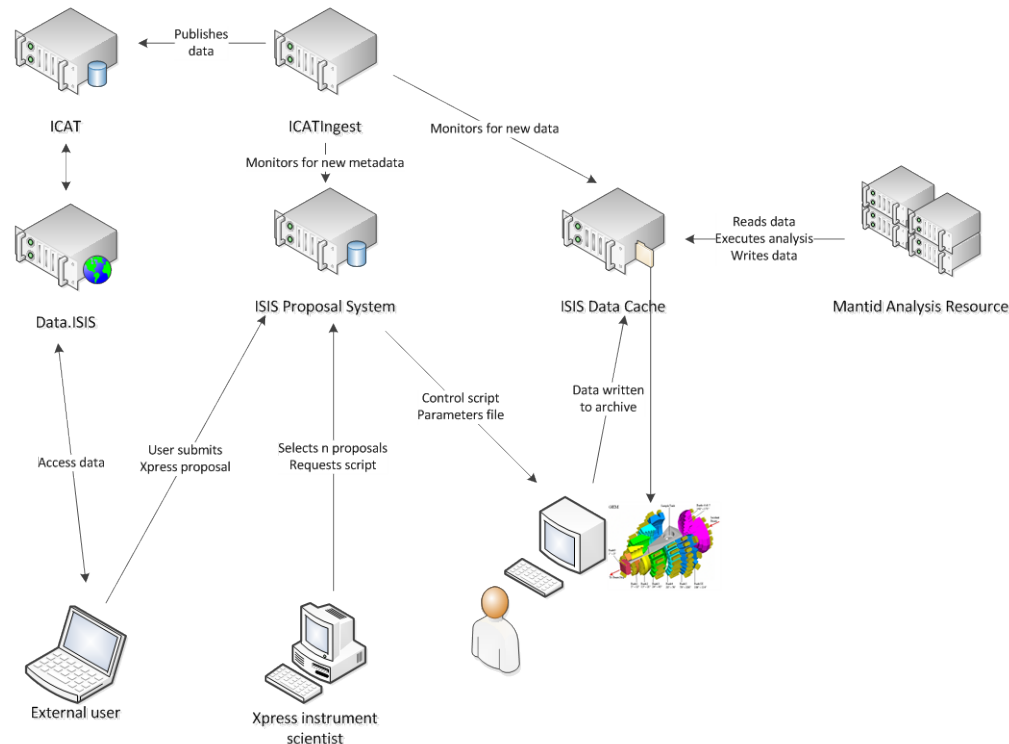
Dealing with high data volumes – 200Gb per experiment in a reconstruction

- hard to move the data – needs to be handled at the facility



Credit: Dr. Mark Basham, Diamond,

- Powder Diffraction
- Experiment by courier
 - Facility staff carry out complete experiment
 - Return Fully reduced and corrected high-quality data to user
- Suitable for automation
- An example of a service which is common in facilities

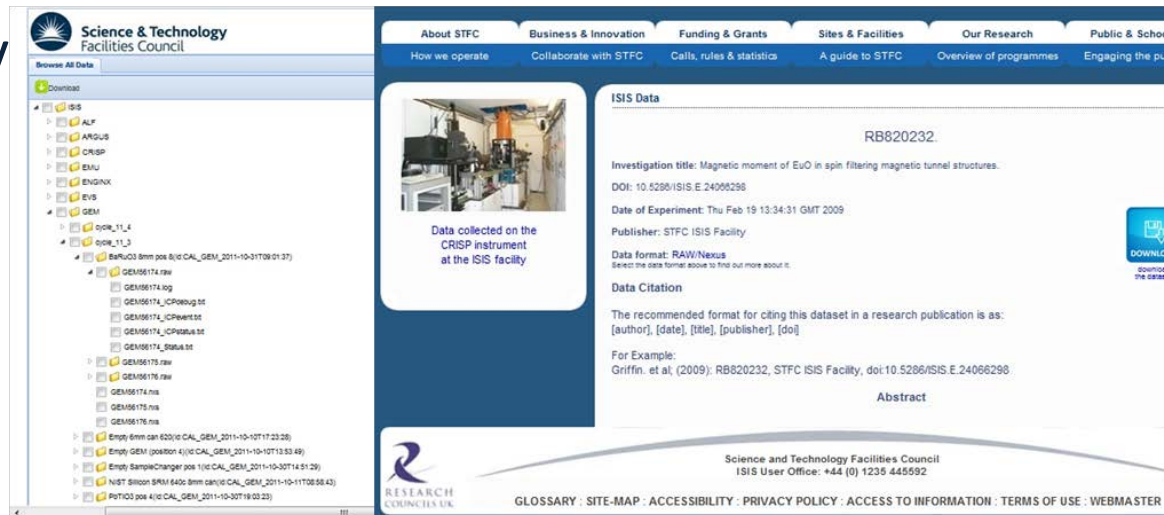


- Powder Diffraction VLab

- WP4, WP6, WP7



- **WP4, WP6, WP7**



- Controlled Vocabulary/Ontology
- Definition of tools which support Provenance

Agreed terms to be used to describe a number of different aspects of Facilities science

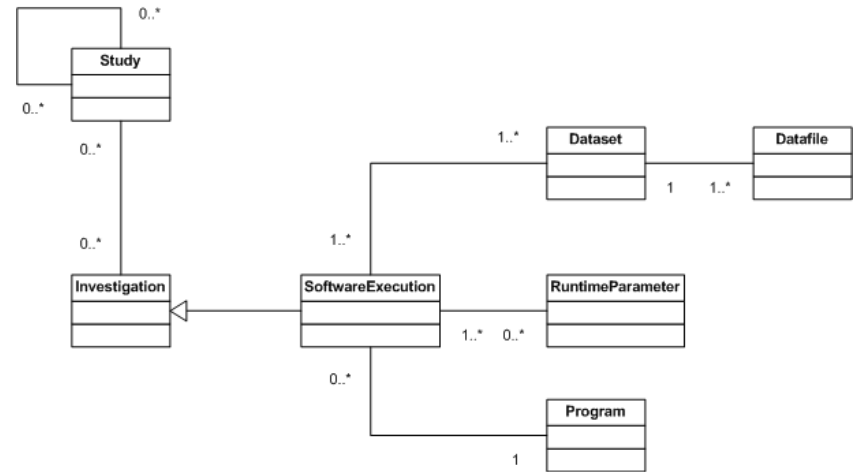
- Facility and Facility Type
- Analysis technique
- Instrument and Instrument Type
- Environmental parameters
- Sample information
- Measured parameters

Can use in tools:

- Catalogue and Search tools
- Linking between tools
- Increase precision and reduce ambiguity
- Hard work to get agreed
 - needs community buy-in

Instrument	----- CLF Laser
--- Neutron Instrument	----- Astra-Gemini
----- ISIS Neutron Instrument	----- Artemis
----- Alf	Scientific Technique
----- Crisp	--- Neutron Technique
----- Engin-X	----- Neutron Diffraction
----- GEM	----- Neutron Spectroscopy
----- LOQ	----- Neutron Reflectometry
----- Merlin	----- Small Angle Neutron Scattering
----- OSIRIS	--- Muon Technique
----- SANS	----- Muon Spectroscopy
--- Muon Instrument	
----- ISIS Muon Instrument	
----- Argus	
----- Deva	
----- Emu	
----- Hifi	
----- MuSR	
--- Pulsed Laser Instrument	

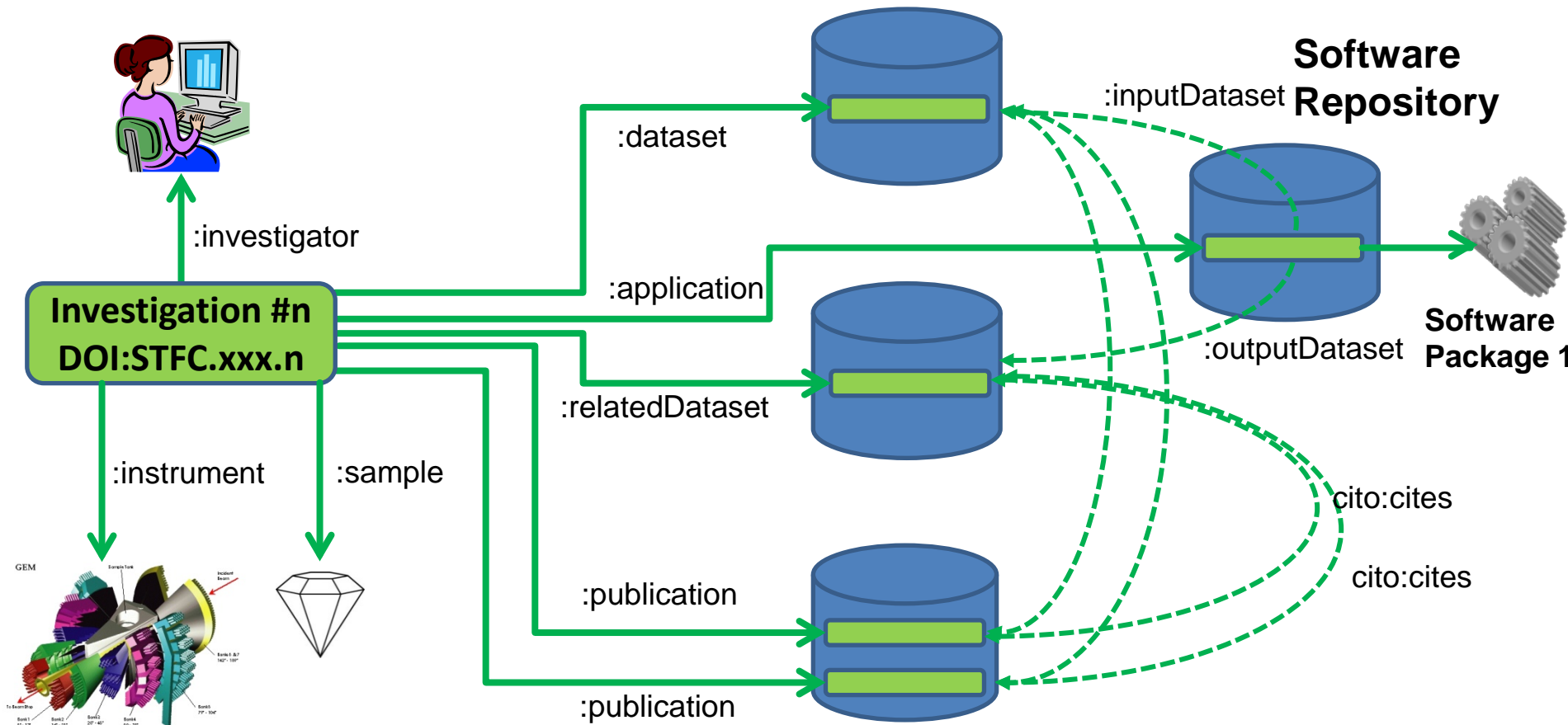
- Modified ICAT to support:
 - Derived data
 - Software, jobs
 - Linking between these
- Modification to the metadata model
- Some early prototypes
- Further develop and use PanSoft for reference software catalogue
 - Would need an API onto PaNSoft
 - Use common information model and vocabularies



Would like to see:

- Modified ICAT
 - Manage Applications and jobs
 - Manage dependencies between (raw and analysed) datasets
 - Manage link to Publications
 - Using controlled terms for Keywords
- ICAT Job Portal
- Modified TopCat for front end
 - Can display “research objects”
 - investigations with all linked components
- PanSoft
- Common Metadata model (extended CSMD)
- Controlled vocabulary
 - Instruments, techniques, parameters ...

Linking the software application into the research object



- Own metadata format (CSMD)
- OAI-ORE
- W3C Prov ontology
- Assume that the software is in a repository



WP6: Provenance

Brian Matthews

Scientific Computing Department, STFC, Didcot, UK