# WP6: JRA on Provenance
# Overview

**Brian Matthews, Erica Yang**

Scientific Applications Group
E-Science Centre
STFC Rutherford Appleton Laboratory

brian.matthews@stfc.ac.uk

# WP 6 : JRA on provenance

- Start M7 (April 2012), Finish M30
  - STFC (Lead)    18 SM
  - ILL                           6 SM
  - ELETTRA                  12 SM

# *WP6  : managing the data continuum*

- The Provenance JRA
  - Extends the repository of information about an experiment
  - Tracking and logging the data analysis steps it links all the data artefacts
  - Records the  data continuum
  - tracking of provenance of data
  - from proposal to publication.

- In general
  - A large and complex task
  - Establishing Science benefit

# Objective

- To develop a conceptual framework, which can record and recall the data continuum, and especially the analysis process.

- To provide a software infrastructure which implements that model to record analysis steps hence enabling the tracing of the derivation of analysed data outputs.

# Tasks

- **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**
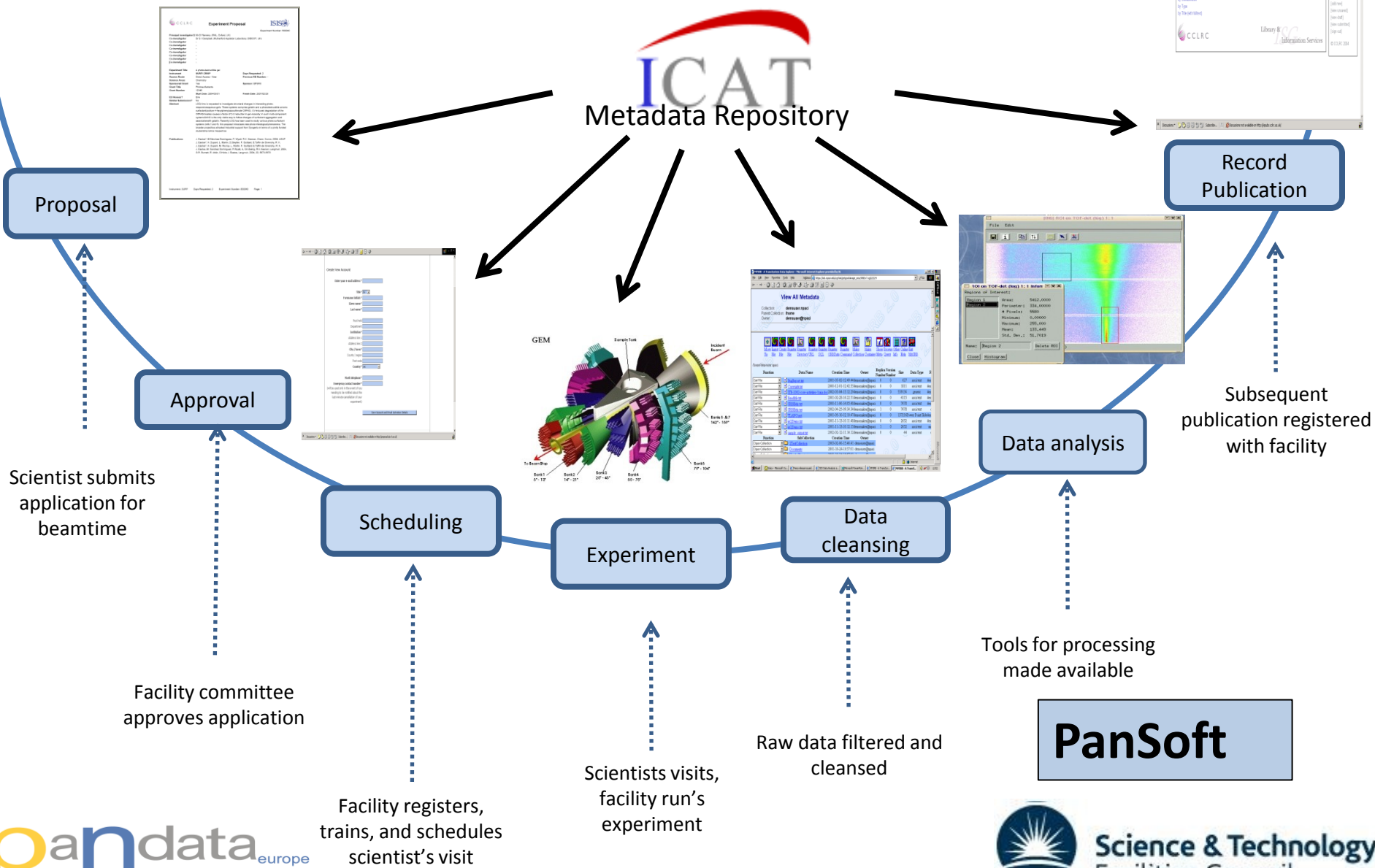- **Task 6: Evaluation**

# Deliverables

- **Deliverables and month of delivery**

- **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)** .

# Requirements

- Explore some case studies e.g.
  - ISIS – SNS (see later) (ISIS)
  - Express services (ISIS, DLS?)
  - DAWN (ESRF/DLS)
  - Directly Programming Data Analysis Kit (DPDAK) (DESY)
  - ISPyB (DLS/ESRF)
  - Publication linking (e.g. DLS+IUCr, ISIS)
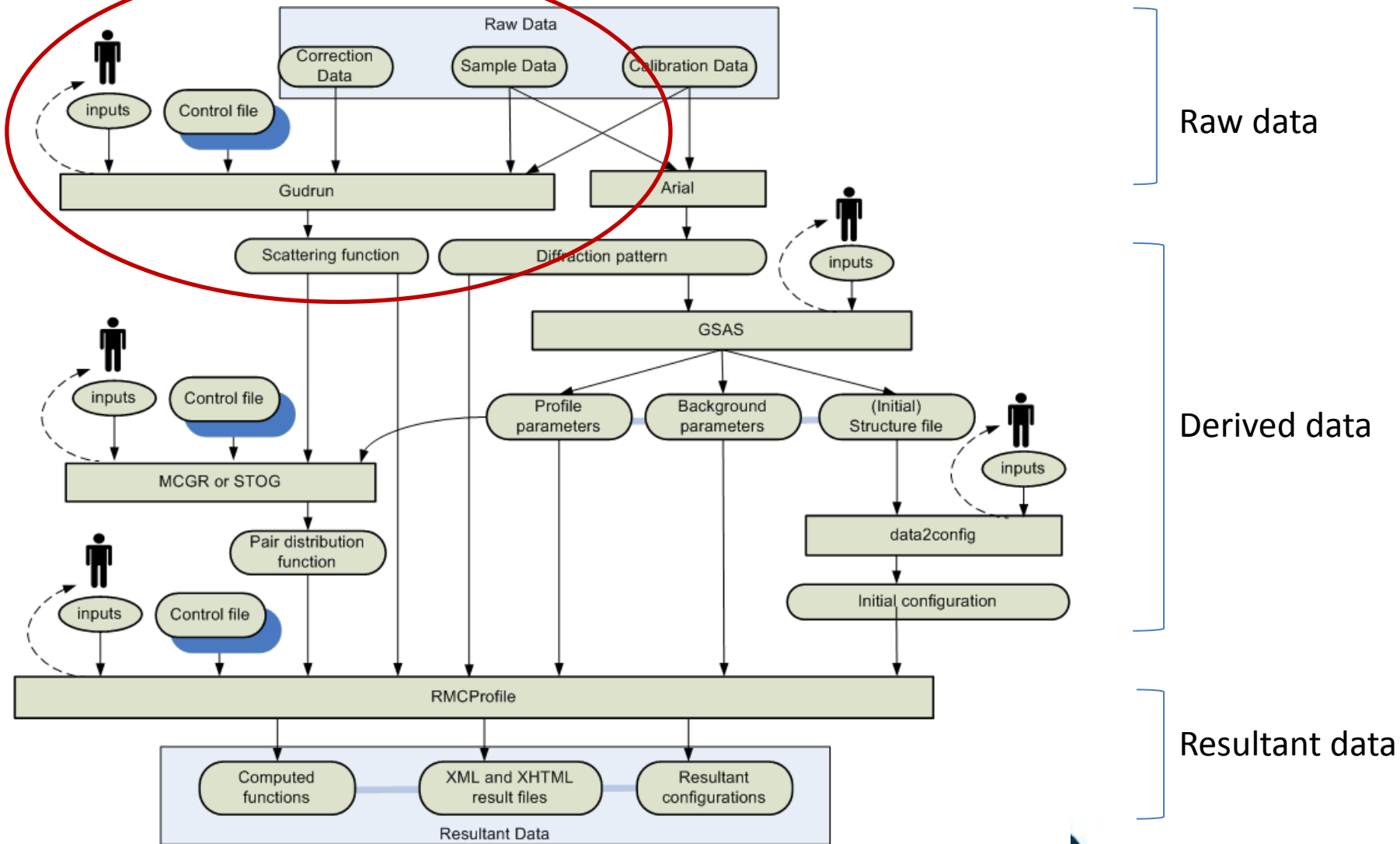- Work with Virtual Laboratories

# Data Continuum



Metadata Repository

Proposal

Record Publication

Approval

Scheduling

Experiment

Data cleansing

Data analysis

Scientist submits application for beamtime

Subsequent publication registered with facility

Facility committee approves application

Tools for processing made available

Raw data filtered and cleansed

**PanSoft**

Facility registers, trains, and schedules scientist's visit

Scientists visits, facility run's experiment

pandata europe

Science & Technology Facilities Council

# Capturing data Provenance for Science: A Use Case and Next Steps

Erica Yang, Brian Matthews

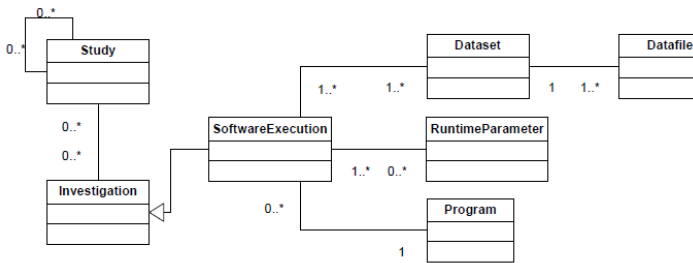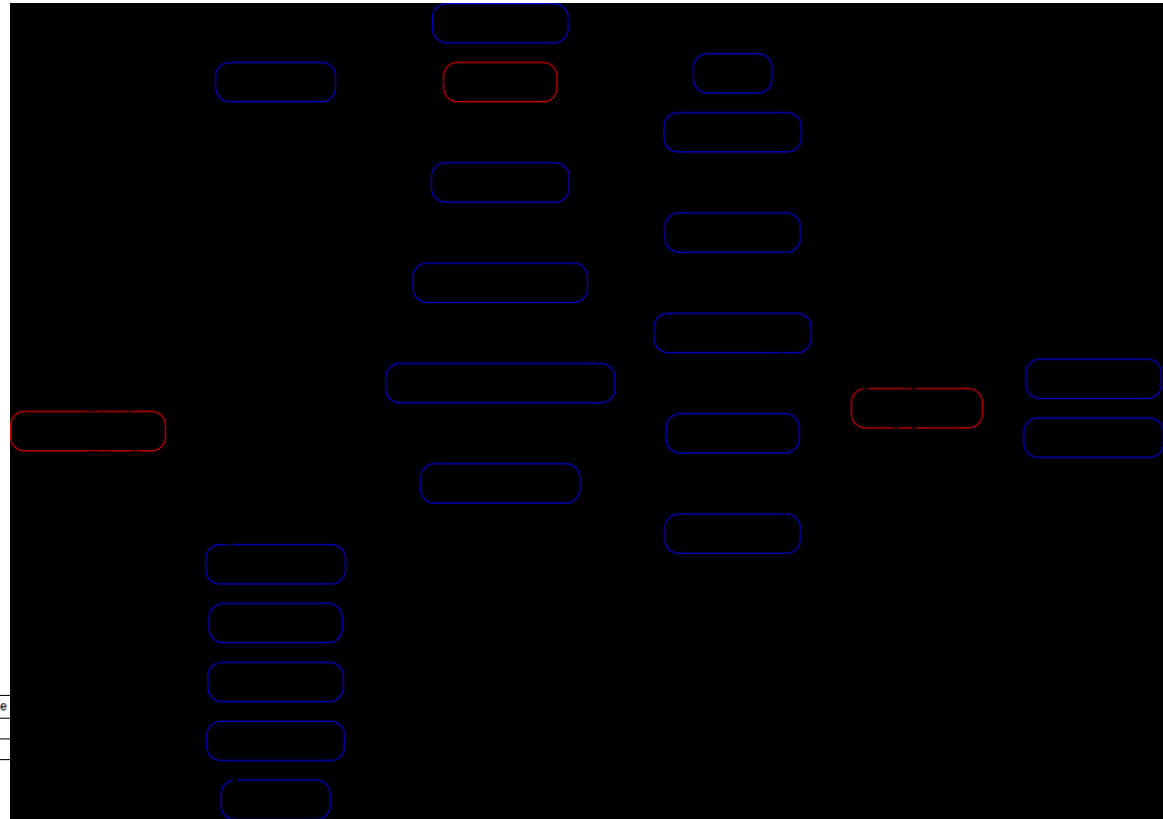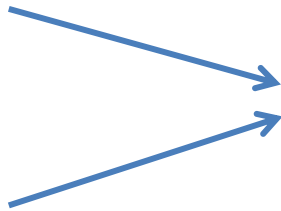Scientific Information Group

STFC e-Science

# Prior Experience



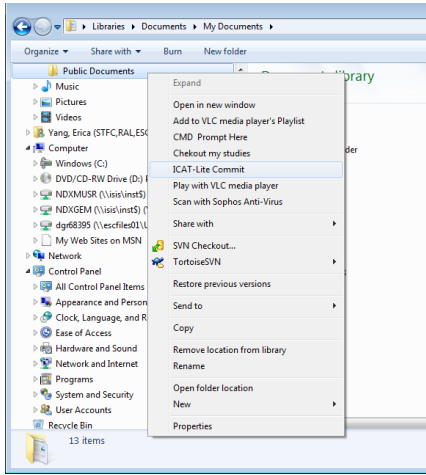Credits: Martin Dove, Erica Yang (Nov. 2009)

# Prior Experience

**Gudrun flow diagram**



25 Jan. 2010 by Alan Soper

Credits: Alan Soper (Jan. 2010)

# Prior Experience

# What we have learned

- Flexibility is the key to manage data provenance
  - Allow "Mix and match" of data processing trials
  - Support forward and backwards tracing of data provenance (e.g. raw<->derived<->paper)
- Researchers are hesitated to change their well established software/practice. "Why would I change?"
  - Need to demonstrate the benefits!

# Use Case: Leveraging Data Provenance for Data Reduction

# An Overview



**SampleTracks**

OpenGenie Script

**Data Acquisition**

**Data Archive**

Sample Information

**Mantid**

**Data Processing**

raw data

**(Extended) ICAT Data Catalogue**

DOIs

**British Library DOI Server**

Outputs

derived data

New links

**ELN**

**Publications**

**Science & Technology Facilities Council**

# Sample Registration

# Experiment Planning – Runs: Detailed View

## Autofill Runs - Configurations

| Run | | Sample | Configuration | Length | Thick | Background |
|---|---|---|---|---|---|---|
| 1 | ❌ | Sample name1: ConcB/0.78(mm)-ConcA/0.03782(mm) | SANS 6m | 40 | 1.0 | |
| 2 | ❌ | Sample name2: ConcA/0.023(mm)-ConcB/0.278(mm) | SANS 6m | 40 | 1.0 | Sample name1: ConcB/0.78(mm)-ConcA/0.03782(mm) |
| 3 | ❌ | Sample name3: ConcB/0.65(mm)-ConcA/0.011(mm) | SANS 6m | 40 | 1.0 | Sample name1: ConcB/0.78(mm)-ConcA/0.03782(mm) |
| 4 | ❌ | Sample name4: ConcB/0.7283(mm)-ConcA/0.023(mm) | SANS 6m | 40 | 1.0 | Sample name1: ConcB/0.78(mm)-ConcA/0.03782(mm) |
| 5 | ❌ | Sample name5: ConcB/0.425(mm)-ConcA/0.0675(mm) | SANS 6m | 40 | 1.0 | Sample name1: ConcB/0.78(mm)-ConcA/0.03782(mm) |
| 6 | ❌ | Sample name6: ConcA/0.089(mm)-ConcB/0.876(mm) | SANS 6m | 40 | 1.0 | Sample name8: ConcA/0.0(mm)-ConcB/0.849(mm) |
| 7 | ❌ | Sample name7: ConcA/0.067(mm)-ConcB/0.316(mm) | SANS 6m | 40 | 1.0 | Sample name8: ConcA/0.0(mm)-ConcB/0.849(mm) |
| 8 | ❌ | Sample name8: ConcB/0.849(mm)-ConcA/0.0(mm) | SANS 6m | 40 | 1.0 | |
| 9 | ❌ | Sample name9: ConcB/0.872(mm)-ConcA/0.0(mm) | SANS 6m | 40 | 1.0 | Sample name8: ConcA/0.0(mm)-ConcB/0.849(mm) |
| 10 | ❌ | Sample name10: ConcA/0.0(mm)-ConcB/0.8766(mm) | SANS 6m | 40 | 1.0 | Sample name1: ConcB/0.78(mm)-ConcA/0.03782(mm) |
| 11 | ❌ | Sample name1: ConcB/0.78(mm)-ConcA/0.03782(mm) | TRANS | 10 | 1.0 | |
| 12 | ❌ | Sample name2: ConcA/0.023(mm)-ConcB/0.278(mm) | TRANS | 10 | 1.0 | |
| 13 | ❌ | Sample name3: ConcB/0.65(mm)-ConcA/0.011(mm) | TRANS | 10 | 1.0 | |
| 14 | ❌ | Sample name4: ConcB/0.7283(mm)-ConcA/0.023(mm) | TRANS | 10 | 1.0 | |
| 15 | ❌ | Sample name5: ConcB/0.425(mm)-ConcA/0.0675(mm) | TRANS | 10 | 1.0 | |
| 16 | ❌ | Sample name6: ConcA/0.089(mm)-ConcB/0.876(mm) | TRANS | 10 | 1.0 | |
| 17 | ❌ | Sample name7: ConcA/0.067(mm)-ConcB/0.316(mm) | TRANS | 10 | 1.0 | |
| 18 | ❌ | Sample name8: ConcB/0.849(mm)-ConcA/0.0(mm) | TRANS | 10 | 1.0 | |
| 19 | ❌ | Sample name9: ConcB/0.872(mm)-ConcA/0.0(mm) | TRANS | 10 | 1.0 | |
| 20 | ❌ | Sample name10: ConcA/0.0(mm)-ConcB/0.8766(mm) | TRANS | 10 | 1.0 | |

WriteScript

# Instrument Control Script

```
# A faked up script example for the SRF use case sandbox

SETSCRIPTNAME THIS_PROCEDURE()

Sample_par width=8 height=8 geometry="Disc"

DO_TRANS

MOVE pos="T9" thick=1 uAhr=10 title="Sample_Example_TRANS_A1-01"
MOVE pos="T10" thick=1 uAhr=10 title="Background_Sample_Example_TRANS_A2-01"
MOVE pos="T11" thick=1 uAhr=10 title="Sample2_Example_TRANS_A3-01"
MOVE pos="T12" thick=1 uAhr=10 title="Directbeam_A0-00"

DO_SANS

MOVE pos="T9" thick=1 uAhr=40 title="Sample_Example_SANS_A1-02"
MOVE pos="T10" thick=1 uAhr=40 title="Background_Sample_Example_SANS_A2-02"
MOVE pos="T11" thick=1 uAhr=40 title="Sample2_Example_TRANS_A3-02"
```

# Flexible Data Model

- Arbitrary number and types of sample parameters, for different types of <u>experiments</u>

- Arbitrary number and types of sample configurations, for different types of <u>instruments</u>

- Designed to work with existing instrument control system, data acquisition system, data catalogue (ICAT), and data processing software

- Designed to allow capturing and tracking of data provenance from samples, to raw and reduced data, and to publications, and backwards

# Next Steps

- Firming up the use cases
  - Revisiting the requirements and existing approaches for data provenance within the PanData facilities
  - Modelling the data continuum
    - Commonalities across existing approaches
    - Common processes across facilities
  - Work with other WPs/partners, e.g.
    - WP4: Data Catalogue Service
    - WP5: Use cases in Virtual Laboratories
    - Directly Programming Data Analysis Kit (DPDAK)
    - Nexus Application Definitions
  - Flagship demonstrations to demonstrate the benefits
    - Express services

# Benefits

- Showcase that data provenance can directly improve research productivity

- Improve facility operational efficiency

- Follow the "non-intrusive" principle to capture and catalogue metadata, designing as part of researchers' existing workflow