# UK Physical Sciences Data Infrastructure (PSDI) initiative

*24th November 2023 - NFDI4Chem Stammtisch*

Dr Nicola Knight
&
Dr Samantha Pearman-Kanza
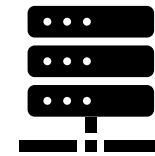
https://www.psdi.ac.uk/

# Aim(s) of PSDI

Support Data as a major driver of research in Physical Sciences

**PSDI** will provide
A data infrastructure that
**connects existing**
experimental and computational facilities
within Physical Sciences and beyond

# *Building Bridges*

▶ Sustaining data resources beyond lifespan of individual research projects

# PSDI: filling a Gap in Provision

- **Other countries** have initiatives underway **in this domain**, e.g.
  - USA: Materials Genome Initiative
  - Japan: NIMS
  - European data infrastructures, such as E-CAM, MaX and NOMAD
  - German National Research Data Infrastructure (NFDI)

UK catch up
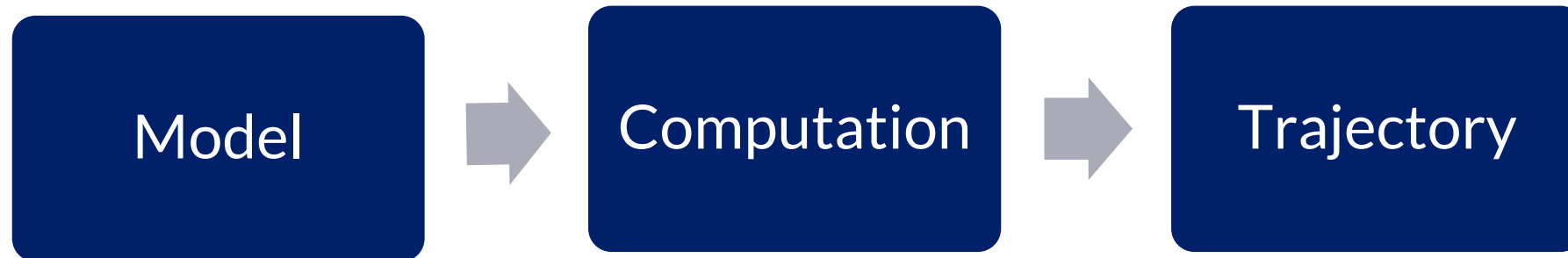
- **Other domains** have initiatives underway **in the UK**, e.g.
  - EBI in Life Sciences
  - NERC Data centres in Environmental Science
  - UK Data Archive in Social Science

Physical Sciences catch up

**We are building a UK, Physical Science, Data Infrastructure**

- Supporting Chemistry, Materials and related disciplines
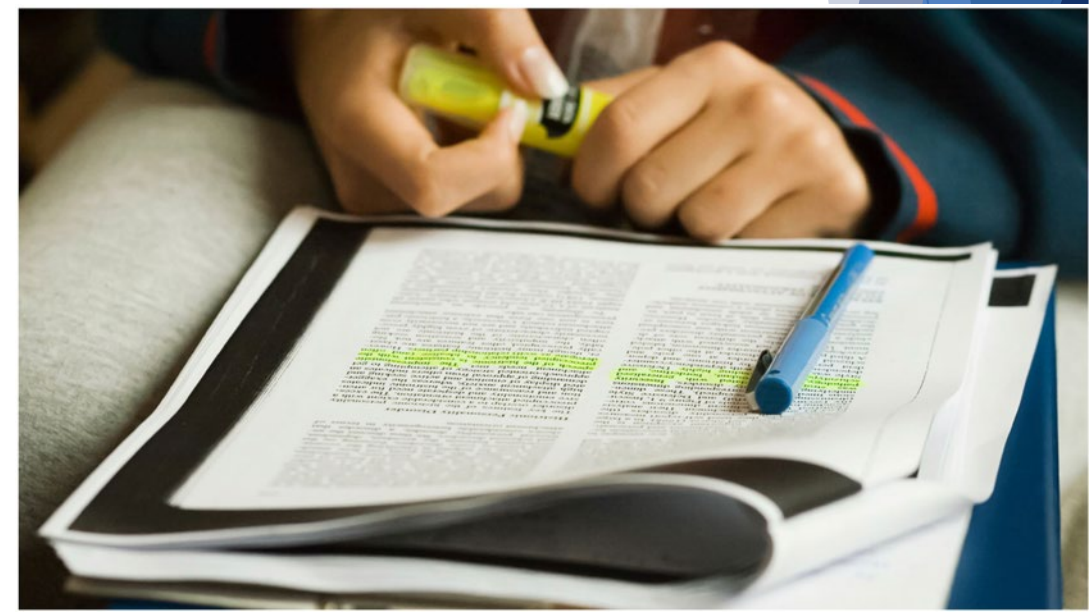- Traversing to and interfacing with Life, Medical, Engineering and Environmental Sciences through federated systems

# An Example:
# Biomolecular Simulations

| Model | → | Computation | → | Trajectory |
|:-----:|:-:|:-----------:|:-:|:----------:|

- Run 10s of simulations to generate data
- Apply know-how to extract science from data
- Publish paper

But ….

- Paper does not include all details needed to **repeat** simulation
- Citations do not give **credit** for *all* resources used

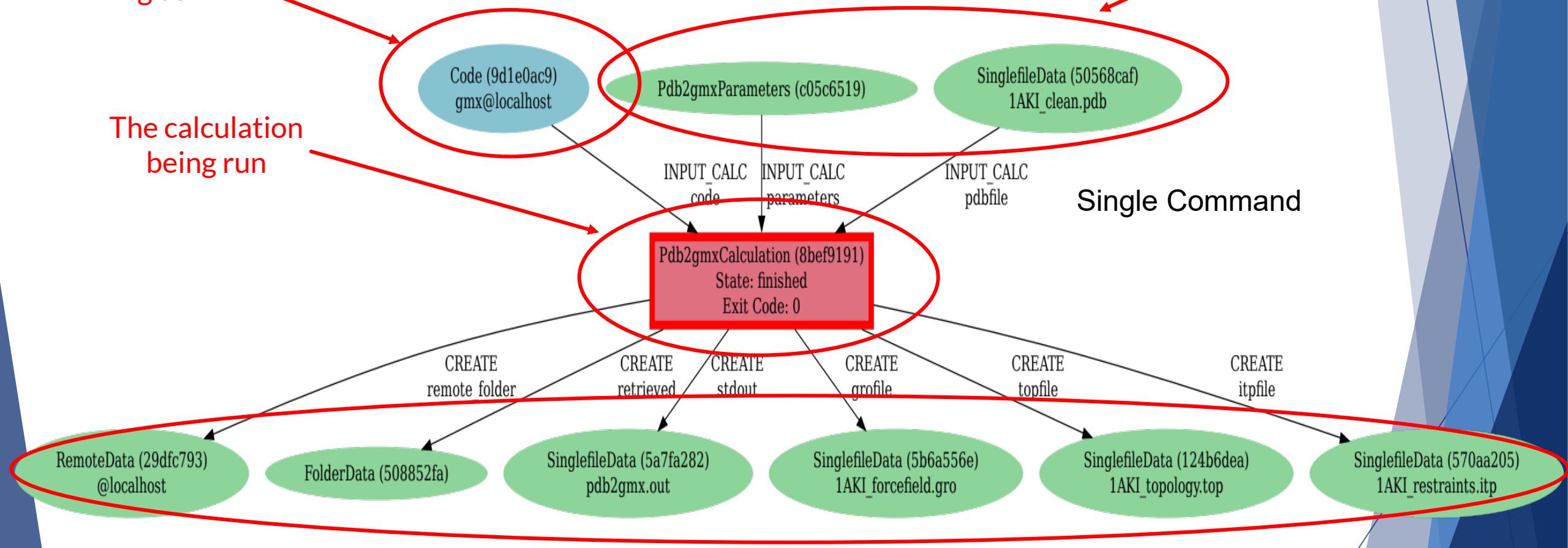# Provenance map of a Single Command in a Simulation

# PSDI PathFinder on
# Research Process Orchestration

Main aim is to improve data practices in domain – align with FAIR principles

▶ Prototype tools to **capture full data provenance** for model creation, simulation and analytics (FAI**R**)

▶ Prototype infrastructure tools to **store, access, find and share** data (**F**AIR)

▶ **Collect** and Integrate existing small scale, disparate data sources

▶ Maintain **compatibility** with other data initiatives (EBI, EU and US)

▶ Link **computational and experimental** data sources

▶ "**I**" (FA**I**R) **Integrations** *not yet in scope* of this pathfinder (excellent projects in CCPBioSim)

James Gebbie & Jas Kalayan

# Process Orchestration PathFinder:
# User Environment Prototype

- Building on GROMACS software (70% of users in UK HPC Biosim Consortium

- Designed to mimick working with native package (command line driven)

- Simple to install and setup our plugin "`pip install aiida-gromacs`" – available through AiiDA

**Normal command:**

gmx pdb2gmx -f prot.pdb -ff oplsaa -water spce -o prot.gro -p prot.top -i prot.itp

**Capture provenance with AiiDA:**

gmx_pdb2gmx -f prot.pdb -ff oplsaa -water spce -o prot.gro -p prot.top –i prot.itp

# Pilot Phase at a Glance



**8 Case Studies**

**> 400 People**

**> 30 Engagement Activities**
- Workshops
- Interviews
- Survey
- Focus Groups
- Group Discussions
- Presentations

**> 50 Organisations**
- Franklin
- Catalysis Hub
- CCPs
- Turing
- National Research Facilities
- HPC
- Central Facilities
- Royce

CS1: Data and simulation driven understanding of **catalytic** activity

CS2: Simulations driven **materials discovery**

CS3: Combining data sources in **Materials Physics**

CS4: **Spectroscopy** data infrastructure

CS5: **Data curation** and availability at instrument-based facilities

CS6: Process Recording and **Electronic Laboratory Notebooks**

CS7: Data **trust**, sharing & **preservation**

CS8: The **role of structure** in Physical Sciences data management

# Pilot Recommendations

13 recommendations in 4 areas:

## Connecting existing infrastructures

3 Recommendations: connecting existing research data services, beyond the lifespan of individual projects, co-operation and co-creation between all stakeholder organisations

## Best Use of Data

4 Recommendations: developing a toolkit for publishing, access to provenanced data, tools for reproduceable data processing, support for transforming data to knowledge

## Best Use of People

4 Recommendations: co-ordination for community activities and input, community training and support, professionalisation for data roles, governance structure for PSDI

## Best Use of Technology

2 Recommendations: services to connect existing provision (data and services), adopt existing technologies

Full recommendations at: https://www.psdi.ac.uk/the-pilot/recommendations
Outputs available via www.psdi.ac.uk and PSDI zenodo community

# Current Work –
# Platform, Pathfinders and Hub

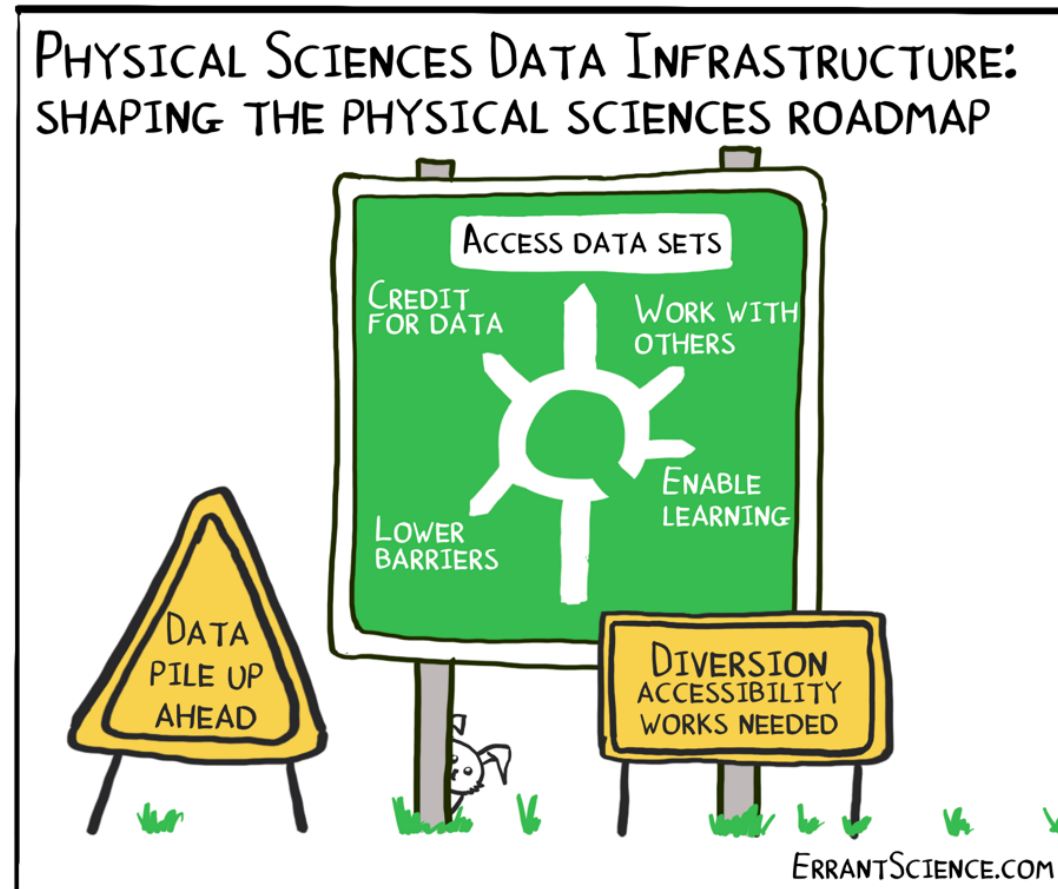- **Platform**
  - Requirements Analysis
  - Capacity Planning
  - System Architecture design
  - Component testing
  - Beginning Build

- **"Pathfinders"**
  - PF1: Experimental data capture
  - PF2: Process Recording
  - PF3: Building Data Collections
  - PF4: Process Orchestration
  - PF5: Data to Knowledge
  - PF6: CCP-NC Database
  - PF7: Reproducible Computational Workflows

- **Hub:** Communications, Governance, Planning,…

# PSDI Hub
# Core Activities & Services

Management, Governance &coordination

Core data infrastructure components

Communications and Engagement

Training

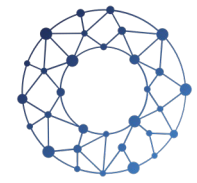# International Collaboration



Research and data is not bounded by international borders!

Alignment with other ongoing and developing international projects
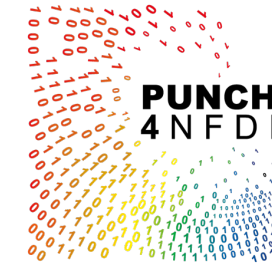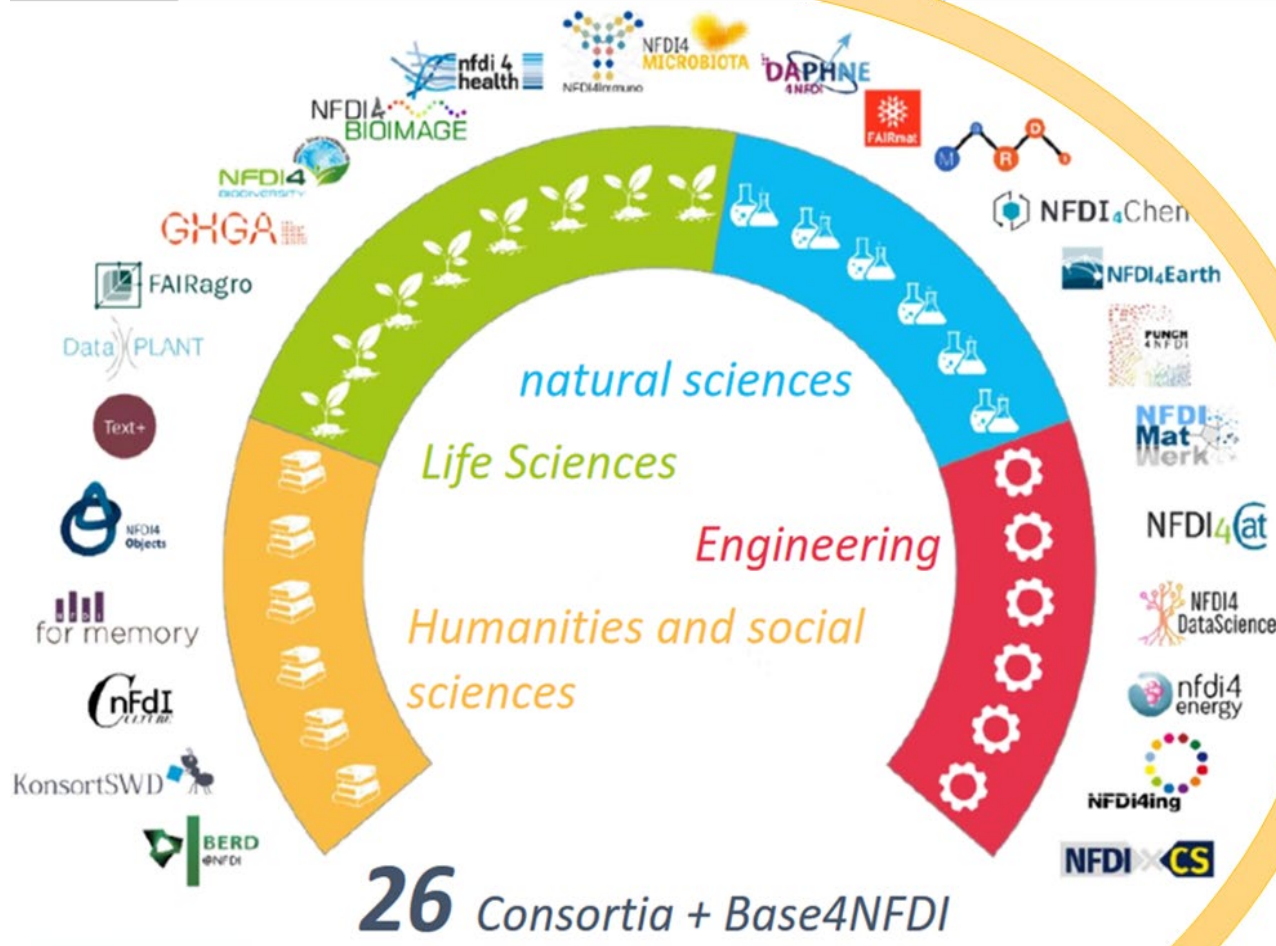
CODATA, RDA, WorldFAIR engagement (among others)

How might we align / collaborate with NFDI?

# PF1: Experimental data capture

Lead: Abraham Nieva de la Hidalga

## Goal

Improve data publishing practices to promote better use of this valuable resource

## Requirements

- track data provenance
- reference all data used
- link data to results
- generate publish ready data

## Proposed approaches

### Create FAIR data objects

- Enable reproducibility and replicability
- Promote data reuse



FAIR Data Object

Publication(s) — SCHOLIX — Link Pubs-Data
Data — PROV-O — Provenance
Organizations — EXPO — Experiment
Researchers — ChEBI — Chemical entities

### Scientific workflows

- Create custom processing/analysis tools
- Combine tools into workflows
- Share and publish workflows
- Generate FAIR digital objects

Galaxy PROJECT

# PF1 - Opportunities for collaboration

**NOMAD** — Data management platform

**voc4cat** **reac4cat** — Vocabulary and ontology tools tailored to the needs of a research community

**Lab Motion** — ELN for documenting custom experimental workflows

**ROCK-IT** — Automate beamline experiments and accelerate operation

# PF2: Process Recording

Lead: Dr Samantha Pearman-Kanza

Investigating routes for recording research process, as well as developing metadata and ontology layers to enable processing and analysis tools

Research focus/service areas:

▶ Process recording tools

  ▶ ELNS & generic notebooks

  ▶ Investigating the data trail from Lab Notebook to Thesis/Paper to Supplementary Information

▶ Exemplars for FAIR data/software/research

▶ Data format conversion service

▶ Converting paper lab notebooks into machine-readable data using Data Revival

▶ Metadata & semantics research

## ELN Finder

The ELN Finder helps you to search and select a suitable Electronic Lab Notebook (ELN) for your purposes.

- More than 40 filter criteria available.
- Filter criteria clearly divided into categories.
- Result list of the identified ELN tools displayed in an overview.
- Brief descriptions of the individual tools included.

🔍 Find ELNs

### Chemistry File Format Conversion Database

Convert from:
cml: Chemical Markup Language

Convert to:
pdb: Protein Data Bank

Conversion success:
Open Babel: complete

Converter details:

Open Babel
Comprehensive converter
https://openbabel.org/docs/dev/Command-line_tools/babel.html Visit website

**Webinar recording available on @PSDI_UK YouTube focusing on this PF**

# PF2: Process Recording – Collaborations & NFDI Alignment

Current Collaborations:

▶ ELNFinder – enabling scientists to choose between ELNs

NFDI Alignment:

▶ NFDI4Chem

  ▶ FAIR data publishing – this aligns with our research for exemplars on FAIR data/software/research

  ▶ Chemotion – opportunities to use Chemotion for case studies – aligns with ELN research

  ▶ Terminology Service – aligns with semantics research

▶ FAIRMAT – FAIR data
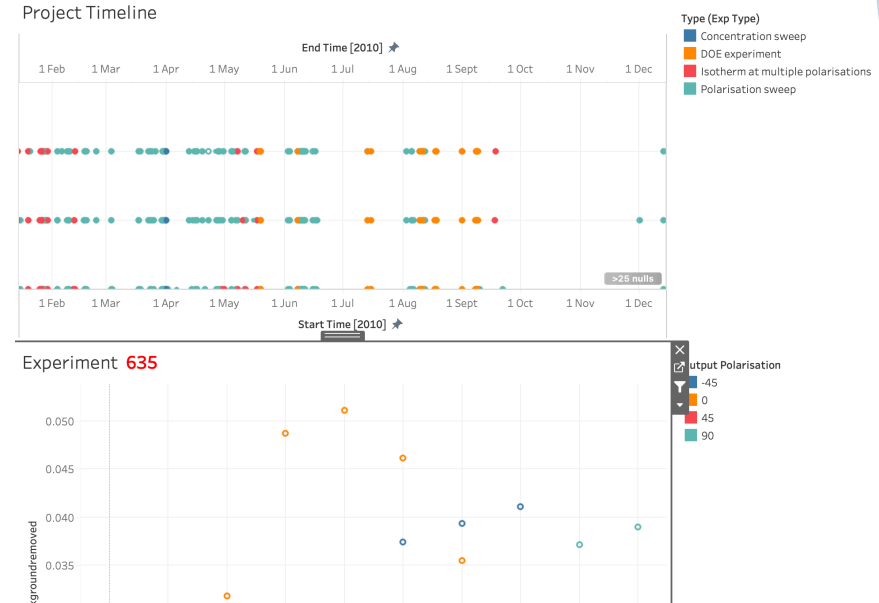
▶ DAPHNE - metadata

# PF3: Building Data Collections

Lead: Professor Jeremy Frey

Explore and develop methods to **build**, **store**, **manage** and access collections for types of data, such as:

- ▶ institutional data
- ▶ facilities data
- ▶ legacy data
- ▶ orphaned data

Several use cases are being worked on:

- ▶ Multiple data types (legacy/active/paper/pdf/electronic/structured)

- ▶ How best to manage, curate and store data

- ▶ Working with a range of tools and technologies

- ▶ These use cases will enable us to provide guidelines on data practices e.g.
  - ▶ Database creation
  - ▶ Chemical identifiers
  - ▶ Data publishing

- ▶ Production of high quality curated datasets

- ▶ Investigating the available repositories (institutional and domain based)

Legacy Second Harmonic Generation data now curated and on Tableau

AMR data curated from spreadsheets, stored in FAIR database & Sharepoint Site

## NFDI Alignment:

- **NFDI4Chem**
  - ▶ Chemistry repositories – this aligns with our research on available repositories, and guidelines towards data publishing

- **DAPHNE**
  - ▶ Community repositories

- **NFDI4Cat**
  - ▶ Data Management in Catalysis – links to guidelines on data

# PF4: Biomolecular Simulations

Lead: Dr James Gebbie-Rayet

▶ Exemplar shown earlier in presentation

▶ This pathfinder aims to establish tools and an infrastructure prototype for capturing the full data provenance for biomolecular simulations, starting from experimental input data through to eventual publications.

▶ Most closely aligned with **NFDI4Chem** and **FAIRMAT**

▶ Will be working with MDDB (Molecular Dynamics Data Bank) so there will be collaboration through EBI with European partners

Webinar recording available on @PSDI_UK YouTube focusing on this PF (publishing soon)

# PF5: Data to Knowledge
## Lead: Dr Alin Marin Elena

- **Design and deploy hardware infrastructure** to host both training data for machine learnt interatomic potentials and the potentials

- Proof of concept database & app:
  - installable centrally within PSDI
  - blueprints for local installation
  - users can interrogate, download and deposit data
  - API and web interface

- Advanced search features: elastic search techniques

- Integrate into ML workflows



Alin Elena, Elliott Kasoar, Federica Zanca, collaboration with Gábor Csányi

Upcoming Webinar focusing on this PF
Dec 14th 1400 UTC
www.psdi.ac.uk/event/webinar-psdi-pf5/

# PF6: CCP-NC Database

Lead: Dr Sathya Sai Seetharaman

## Collaborative Computational Project for NMR Crystallography (https://www.ccpnc.ac.uk/)

- Supports the multidisciplinary experimental NMR community with computational tools

- Recognised by International Union of Crystallography

- Created .magres file format for unified representation of crystal NMR parameters.
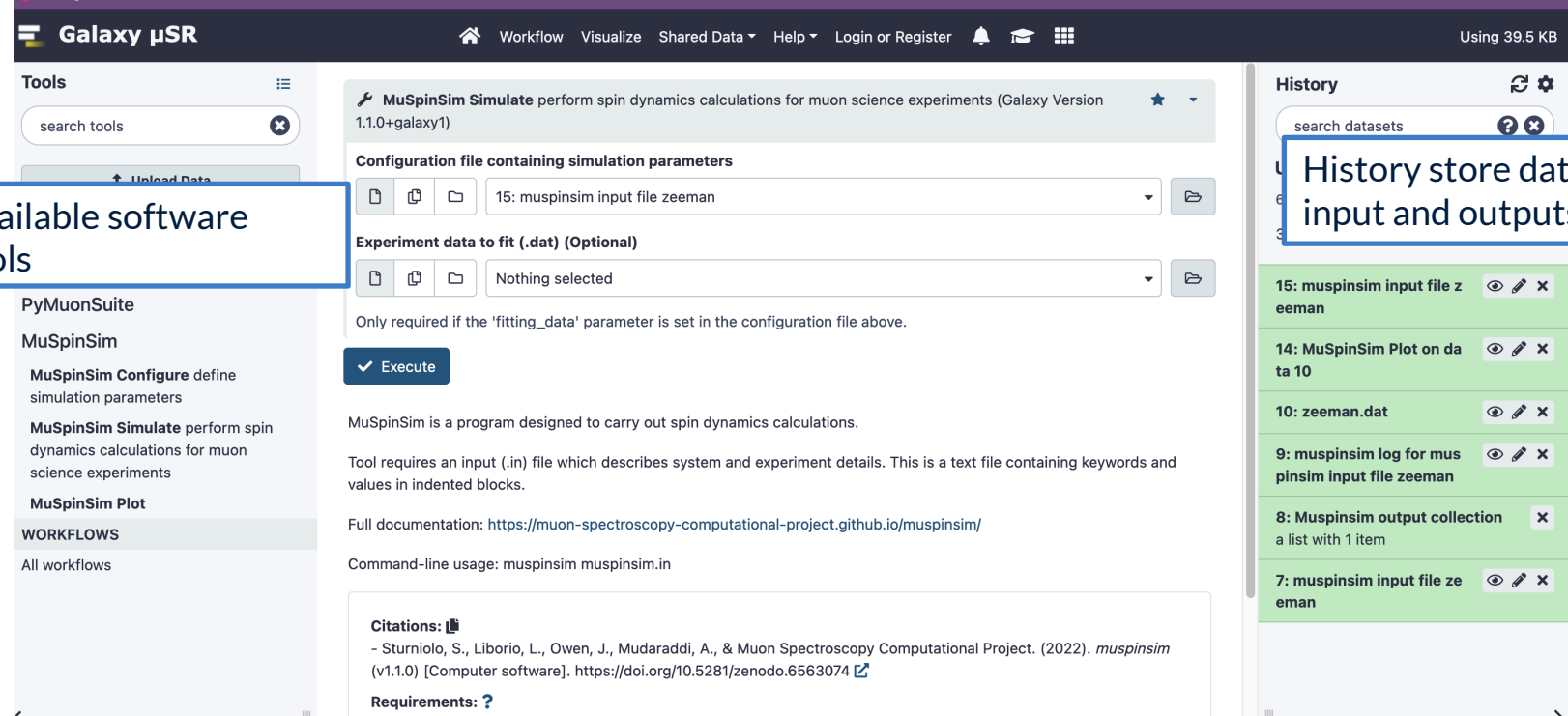
## Pathfinder Objectives

- Improving the current CCP-NC magres database (https://www.ccpnc.ac.uk/database/)

- Development of a state-of-the art database (version 2)

# CCP-NC – FAIRmat

**Productive discussion with NOMAD about potential collaboration:**

- Magres database support

    - Initial discussions about magres file parser support

    - Add functionality for QE/CASTEP -> magres workflows

    - CCP-NC magres database – support to legacy data

- Adding NMR specific searchability within NOMAD

- Potential extension of services to NMR experimental community

- Data visualisation support in-line with CCP-NC tool standards

- Potential CCP-NC – NOMAD collaborative development

# PF7: Reproducible Computational Workflows

**PSDI** PHYSICAL SCIENCES DATA INFRASTRUCTURE

Lead: Dr Leandro Liborio

**Galaxy** PROJECT

https://galaxyproject.org/



Available software tools

History store data files: input and outputs

## Muon Experiments

- Muon experiments are performed at the Rutherford Appleton Lab, STFC, UK.
- Develop software tools for interpretation of those muon experiments.
- Tools based on computing simulations, i.e.: DFT.
- Created associated Galaxy tools and Galaxy instance.
- Use the Galaxy platform to manage the workflows resulting from the tools.

## X-ray Absorption Spectroscopy (XAS) Experiments

- XAS experiments are performed at the Rutherford Appleton Lab, STFC, UK. Catalysis-related experiments.
- Software tools for processing experimental data already available.
- Created associated Galaxy tools and Galaxy instance.
- Use the Galaxy platform to manage the workflows resulting from the tools

# Potential Collaborations/Alignment Between PF7 and NFDI

## NFDI4Cat

▶ Present galaxy tools as a complementary method for processing workflows. Currently working with A. Nieva de la Hidalga (Pathfinder 1) on galaxy for catalysis experiments.
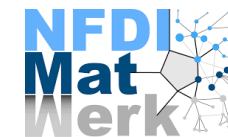
## DAPHNE4NFDI

▶ Task area 3 from DAPHNE4NFDI refers to "Infrastructure for Data and Software Reuse". Galaxy is an open, web-based platform for accessible, reproducible, and transparent computational research.

▶ We are collaborating with colleagues from Oak Ridge National Lab on Galaxy tools for neutron science.

▶ We are working with the Diamond Light Source on Galaxy tools for x-ray experiments.
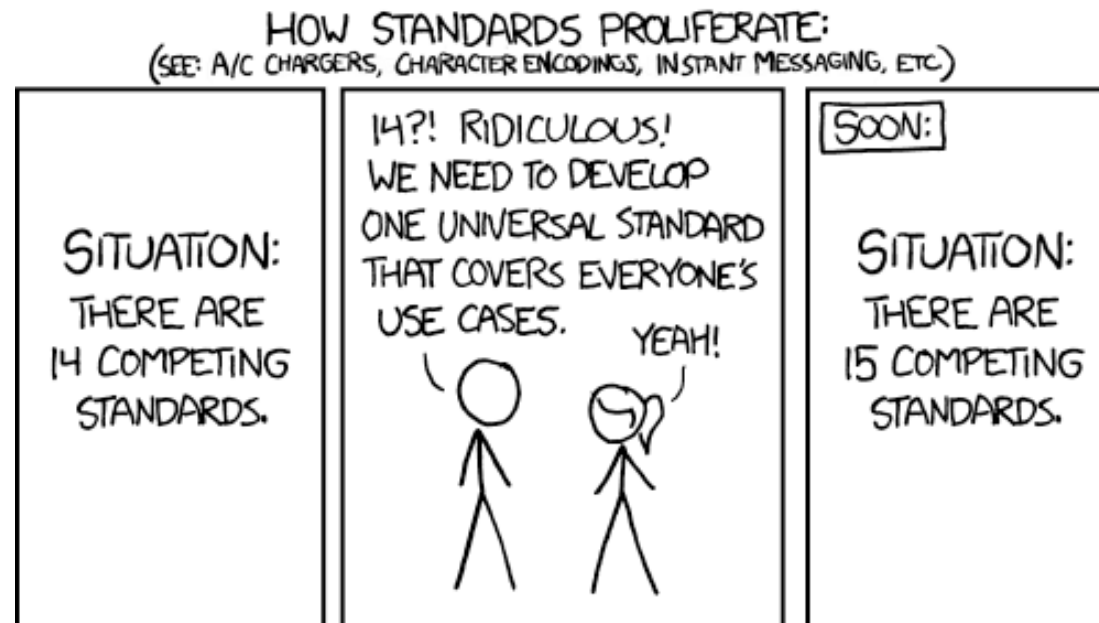
## NFDI MatWerk

▶ IUC05 Digital infrastructure and workflows for labs.

▶ PP13 Tomography and Microstructure-based Modelling.

*"a workflow for the transfer of tomography data to related multi-scale simulations will be established".*

▶ IUC15 Method- and scale-bridging workflows and data structures for tomography.

*"Materials tomography methods and resulting data vary strongly depending on the method used, the experimental approach and the workflow for post-processing. Currently, there is no established protocol which would allow to conduct all necessary steps in a well-defined manner. The resulting data from different methods are therefore not interconnected and workflows are intransparent."*

# Cross Project Topics

▶ Best practices

▶ Skills / Training

▶ Standards

   ▶ Files formats

   ▶ Metadata

▶ Publishing / Sharing

▶ Semantics / Ontologies

XKCD – Standards, https://xkcd.com/927/ -Creative Commons Attribution-NonCommercial 2.5 License

# Enabling communication

www.psdi.ac.uk    @PSDI_UK    @PSDI_UK    PSDIUK

# Any Questions?

Please do contact our researchers directly

They just love to talk about our work!