



Science and  
Technology  
Facilities Council

Scientific Computing



# Computational Biology

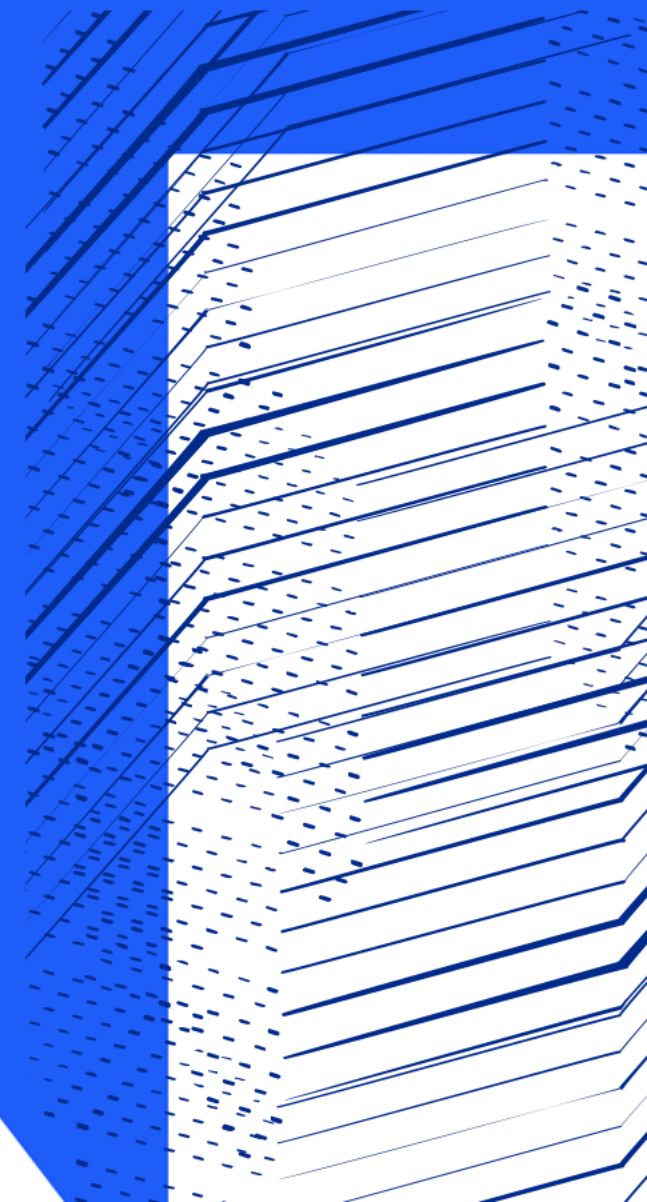
Martyn Winn 20th May 2024

SCD Computational Biology theme

Current facility activities

- CCP4 & CCP-EM
- Molecular simulation and modelling
- Data & machine learning
- Computed tomography

Future opportunities



# Our history ....

- 1990s - CCP4 team at Daresbury in Synchrotron Radiation Department
- Ca. 2005 Molecular Dynamics simulation activity begun
- CCP4 move to Diamond House (2009) then Research Complex (2010)
- CCP-EM team created (2012) – help curate eBIC software
- 2020 Computed Tomography team added
- Jan 2024 Computational Biology Theme created



# Computational Biology Theme

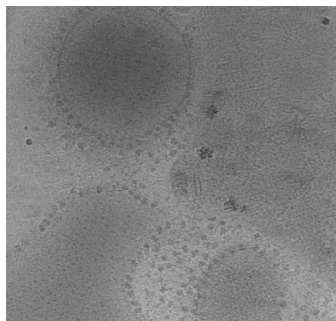
## Macromolecular Crystallography

Methods development  
Software suite



## Molecular and Cellular Electron Microscopy

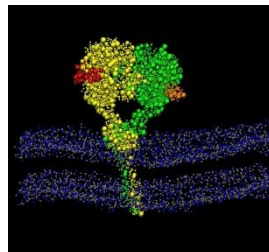
Methods development  
Software suite



## Biomolecular simulation

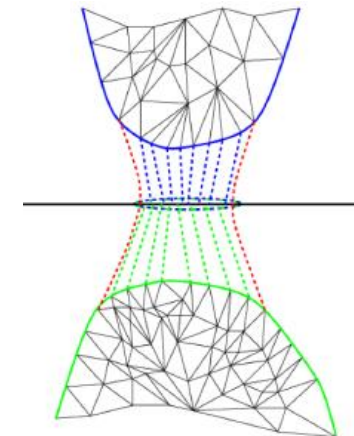
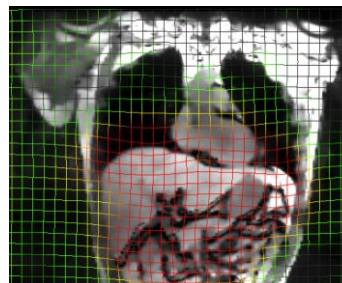
Dynamics / *in silico* expts  
Tools and HPC access

CCPBioSim



## Computed Tomography and Imaging

Software libraries  
CCPi – materials  
CCPSyneRBI – medical imaging



## Elsewhere in SCD ...

**Materials** – MD, mesoscale,  
QM/MM, electronic structure

**Engineering** – meshes, finite  
element

**Maths** - algorithms

**SciML** – machine learning

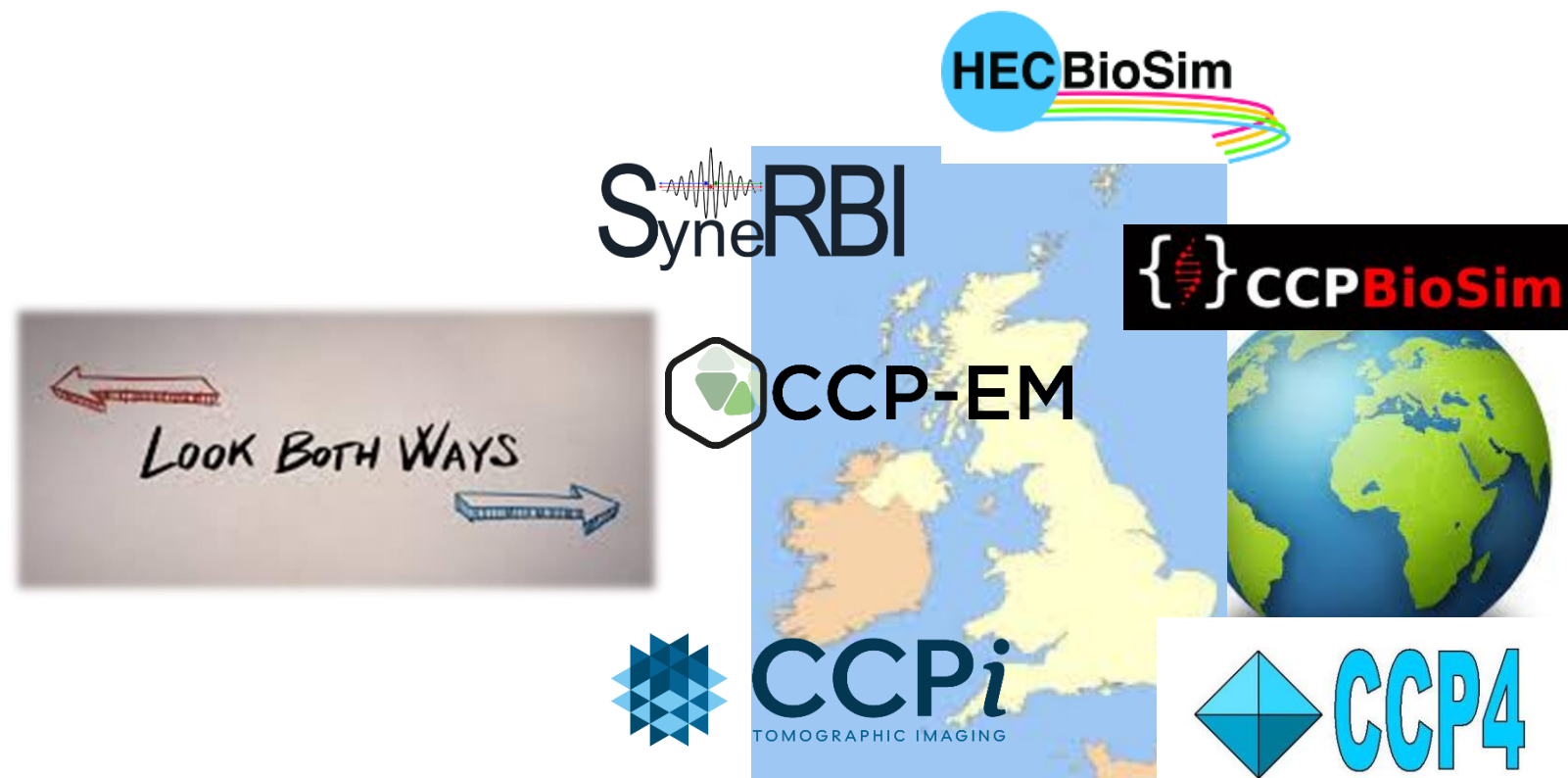
**DAaaS** - services

**Cloud** – hardware backend

# Harwell facilities



# UK and global communities



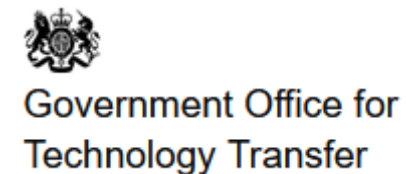
# Our interests ....

- Largely molecular / structural biology
  - Also wider imaging and informatics
- Closely tied to experimental data
  - Modelling for interpretation and in silico experiments
- Grown organically over 30 years
- Variety of funding, mostly external

- Development of **software**, long term and stable
- Development of methods, though not deep theory
- **Data / metadata** handling important

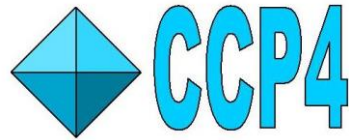


In practice, software libraries and frameworks to support multiple workflows



	Diffraction	Microscopy	Tomography	Scattering
1. Expt.	CCP4, DIALS	CCP-EM	CCPi	MD for neutron reflectometry
2. Sim.	Electron diffraction	AlphaFold, Parakeet	Digital twins	Molecular Dynamics
3. Data	CCP4Cloud	CCP-EM Pipeliner	CIL	Provenance for MD
4. AI/ML	Solvent content prediction	2D/3D classification	De-noising	Learned potentials

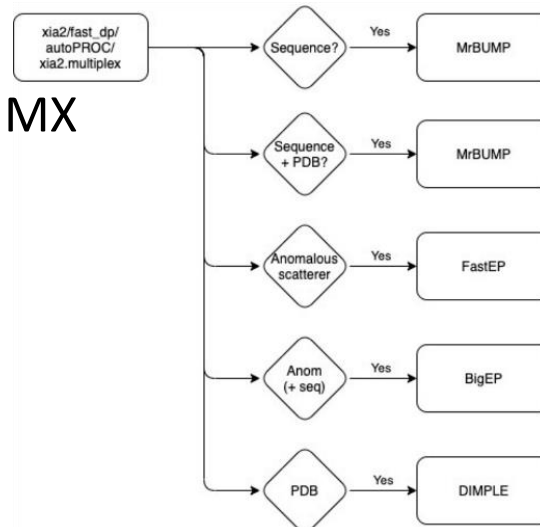
# CCP4 and CCP-EM



## CCP4

- Macromolecular crystallography
- Originally X-ray, now neutron and electron
- Diffraction pattern to structure
- Large downloadable software suite
- Joint workshop with Diamond every Nov/Dec; Study Weekend
- Software used by Diamond MX
  - DIALS
  - MrBUMP
  - Crank2
  - etc

[www.ccp4.ac.uk](http://www.ccp4.ac.uk)



## CCP-EM

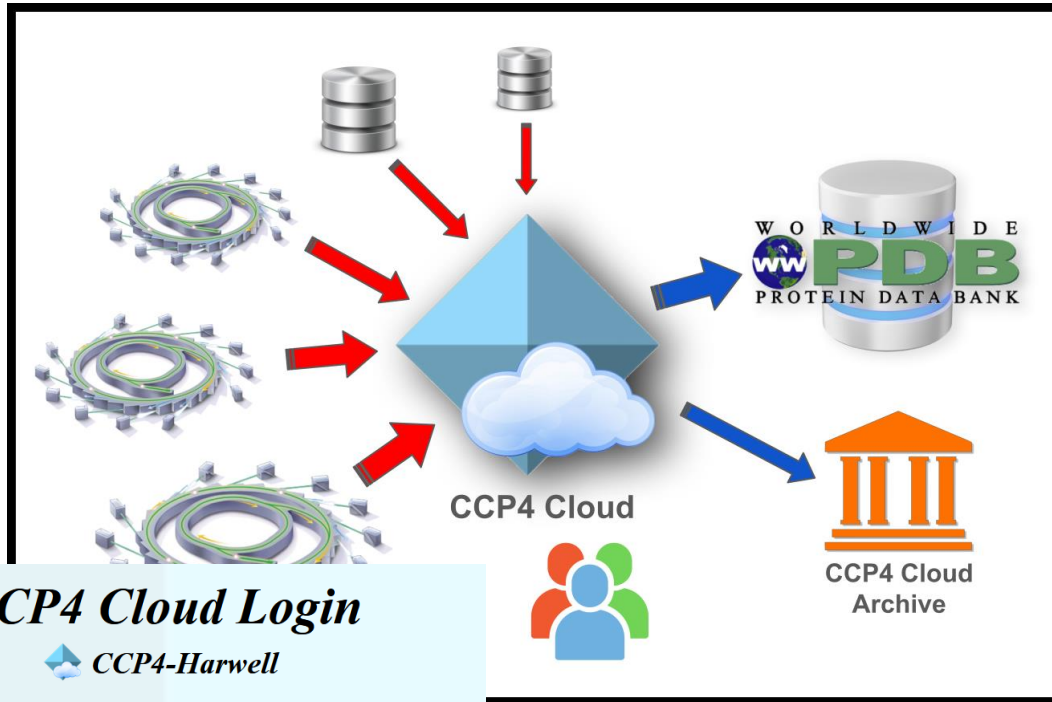
- Cryogenic electron microscopy and tomography
- From frames/micrographs to map/structures
- Downloadable software suite
- Software used by eBIC
  - CCP-EM suite
  - Relion
  - Pipeliner based automated workflows

[www.ccpem.ac.uk](http://www.ccpem.ac.uk)

# CCP4Cloud



Jools Wills



## CCP4 Cloud Login

CCP4-Harwell

Login name:   
Password:

- 
- 
- 

**NOTE:** For using **Coot**, **DUI**, **iMosflm** and similar tasks, install [the CCP4 Software Suite](#) and start CCP4 Cloud with this icon:



- Online solution for crystallographic computation  
diffraction images → structure refinement →  
deposition in the PDB
- Main service at RAL, backed by SCD Cloud, but  
other instances have been created.
- **Issue:** upload of diffraction images difficult  
⇒ usage typically starts with reduced data.

- With Diamond, working on DataLink between data collected at MX  
beamlines and CCP4 Cloud.

- Pull-type link from CCP4Cloud.
- Push-type link (using Node JS API).
- Globus Endpoint.



# Diffraction across the facilities

- Diffraction experiments on MX beamlines (X-ray), eBIC (electrons), ISIS (neutrons).
- Workflow similar in each case. Differences due to particle properties (e.g. flat Ewald sphere for electrons) or experimental setup (time-of-flight for neutrons).
- Processing of diffraction images to spot intensities carried out by DIALS software.
  - Originally for X-rays, replacing old CCP4 software
  - Adapted to neutron TOF diffraction experiments, starting Oct 2020
  - Adapted to ED from 2018 onwards
- DIALS core well-maintained.
- Collaboration Diamond, Lawrence Berkeley, CCP4.
- Customisations for particular experiment types.



Gwyndaf  
Evans



David  
Waterman



<https://dials.diamond.ac.uk/>

# Extending DIALS to Polychromatic Experiments



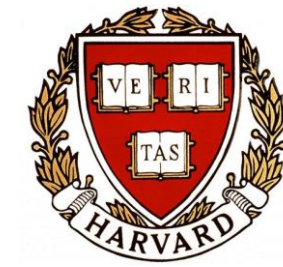
SXD  
LMX



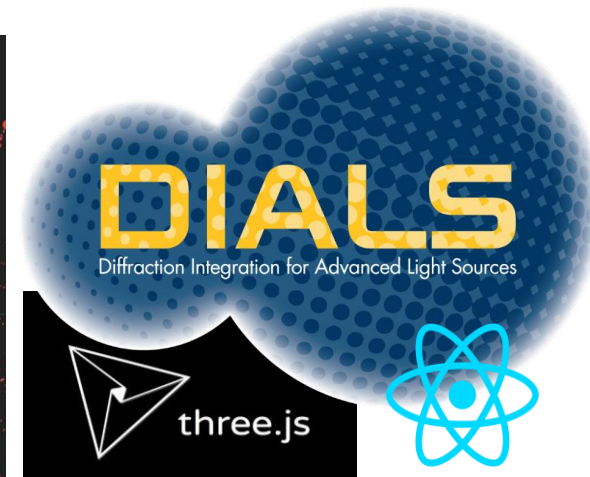
MANDI  
TOPAZ



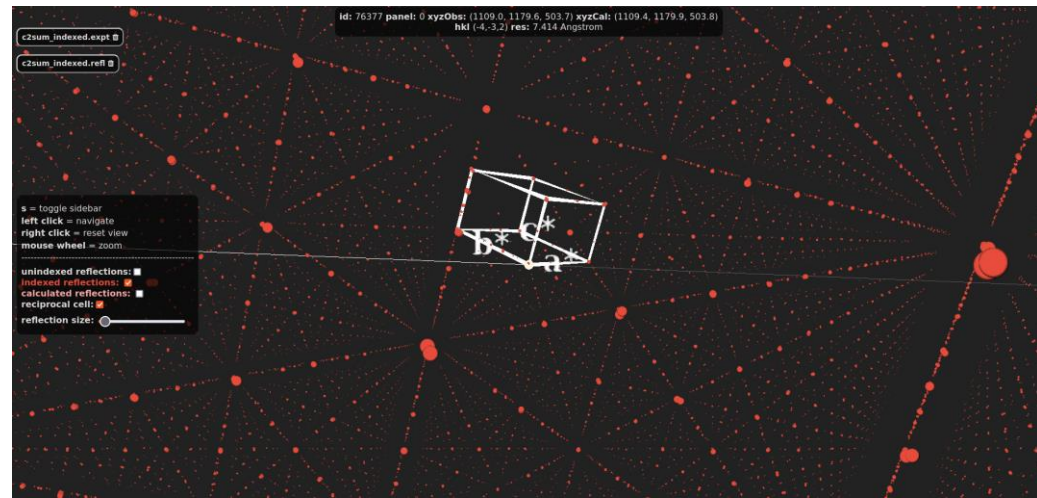
NMX



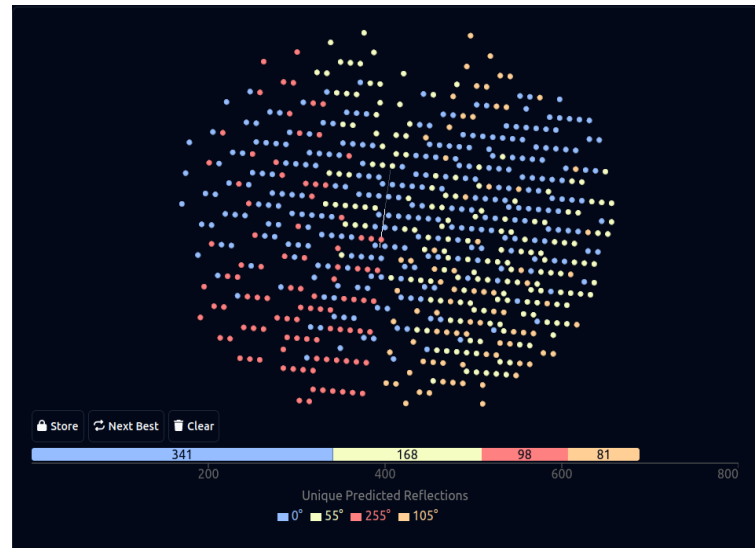
X-ray Laue



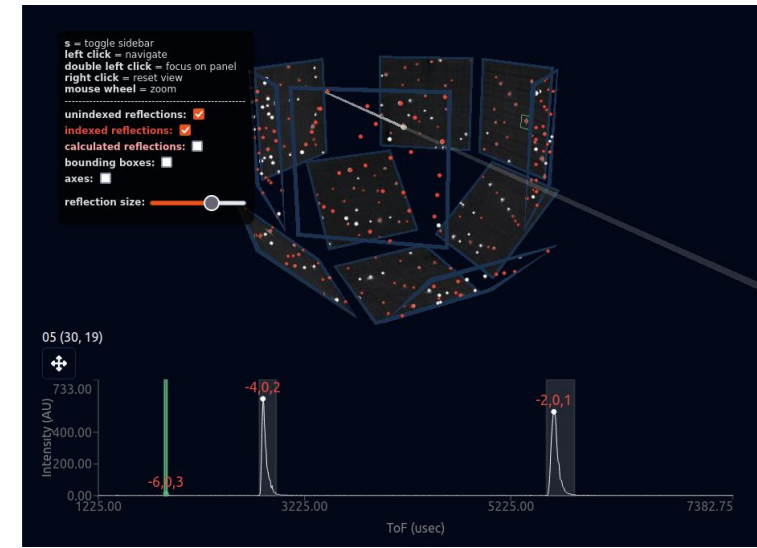
New browser-based GUI



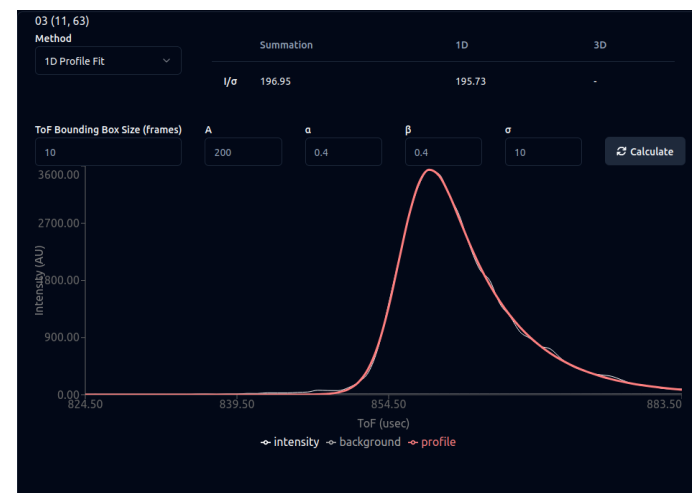
Reciprocal Lattice Viewer



Experiment Planner



Experiment Viewer



Integration Profiler

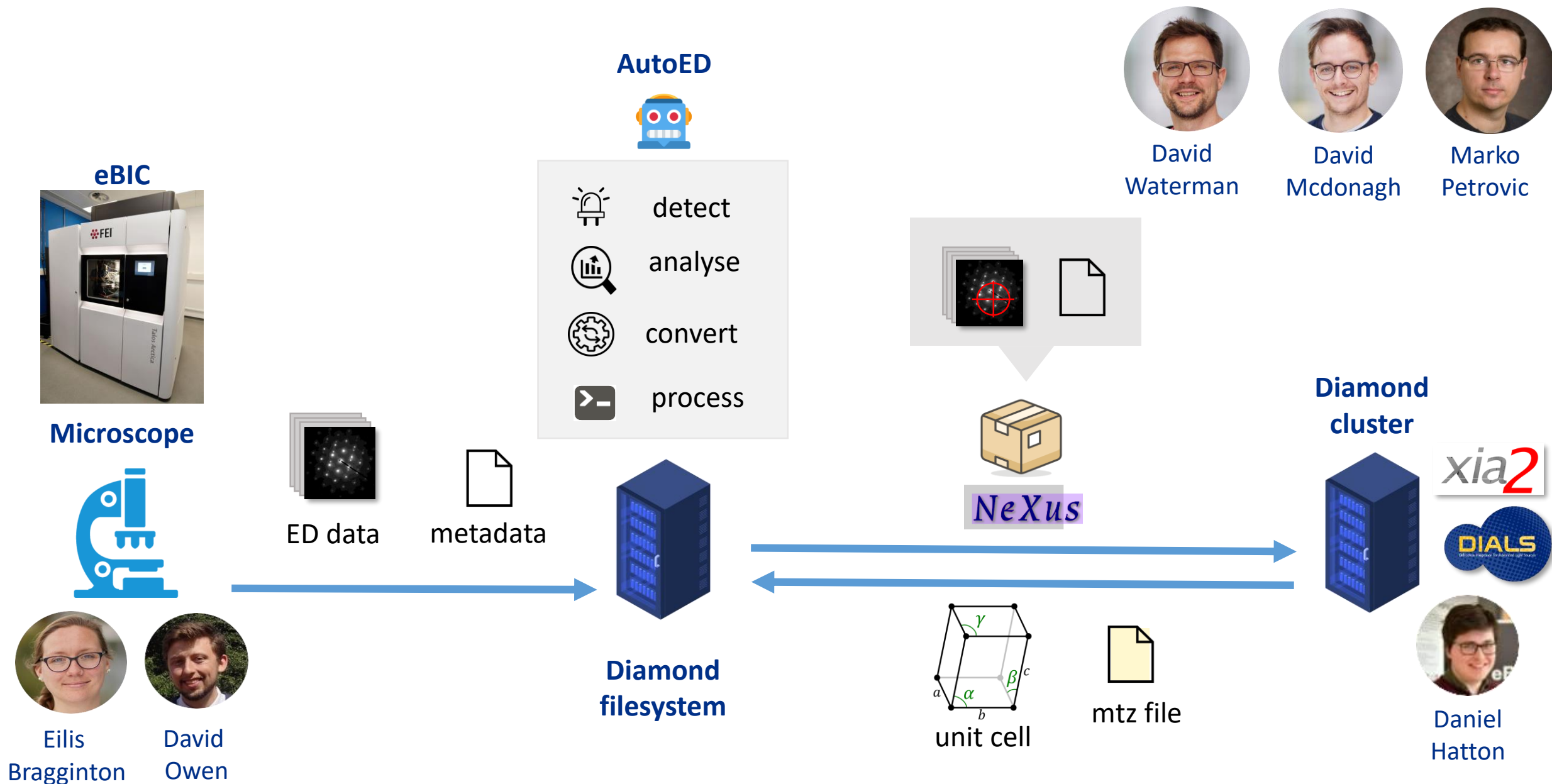


David Mcdonagh



Silvia Capelli

# Auto-processing of electron diffraction data



# CCP-EM Pipeliner



- **Business logic layer decoupled from tasks and UI**
- Python 3
- MPL 2.0 licence
- <https://gitlab.com/ccpem/ccpem-pipeliner>
- Directed Acyclic Graph data flow
  - Metadata tracking
  - Jobs have input/output nodes
- Plugin architecture
  - 100-500 lines per app
  - Open - supports external apps
- CLI / UI / scripting APIs



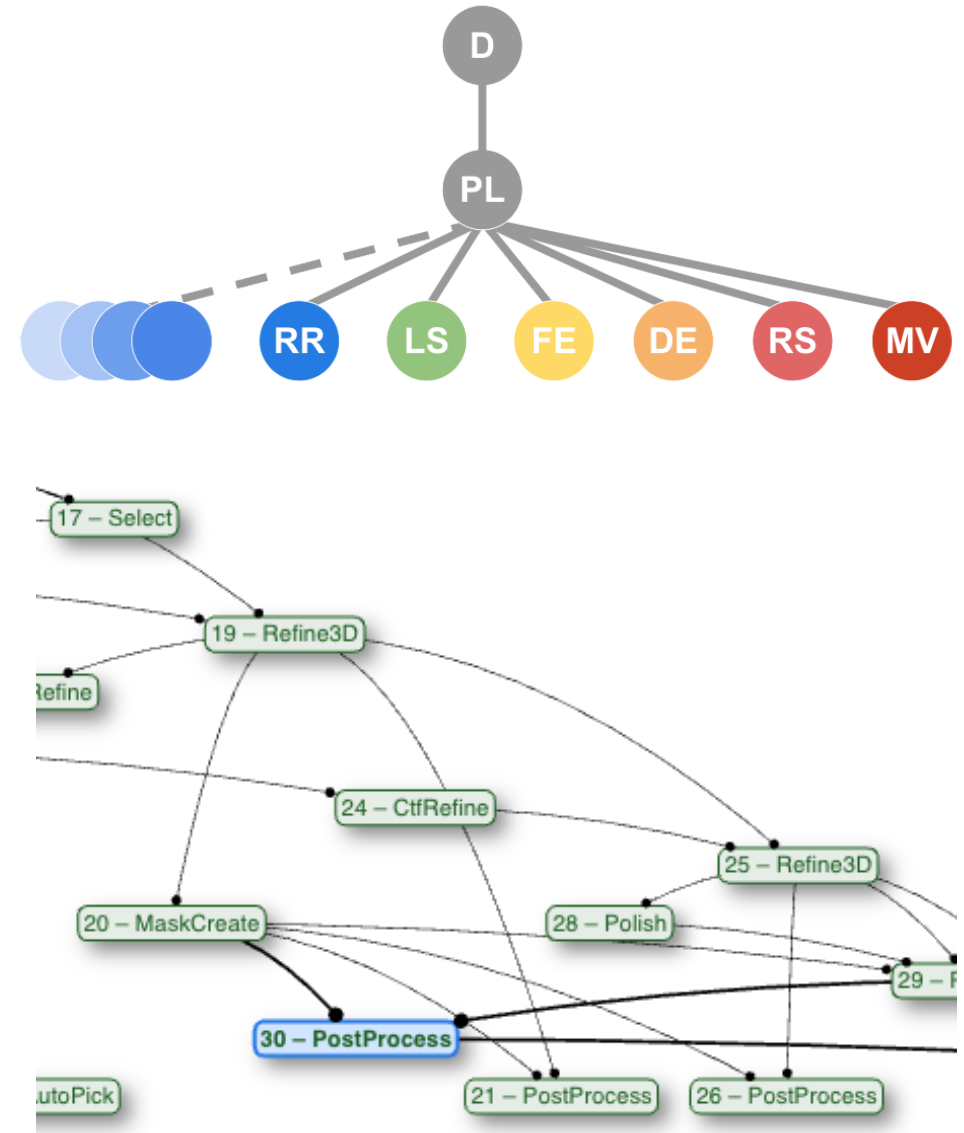
Matt  
Iadanza



Sjors  
Scheres



Gerard  
Bricogne

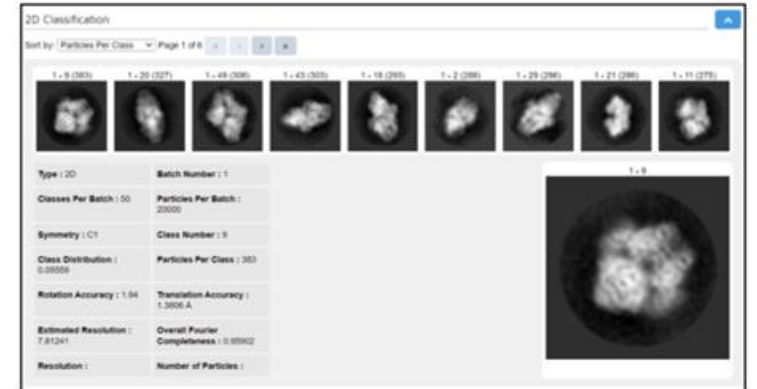
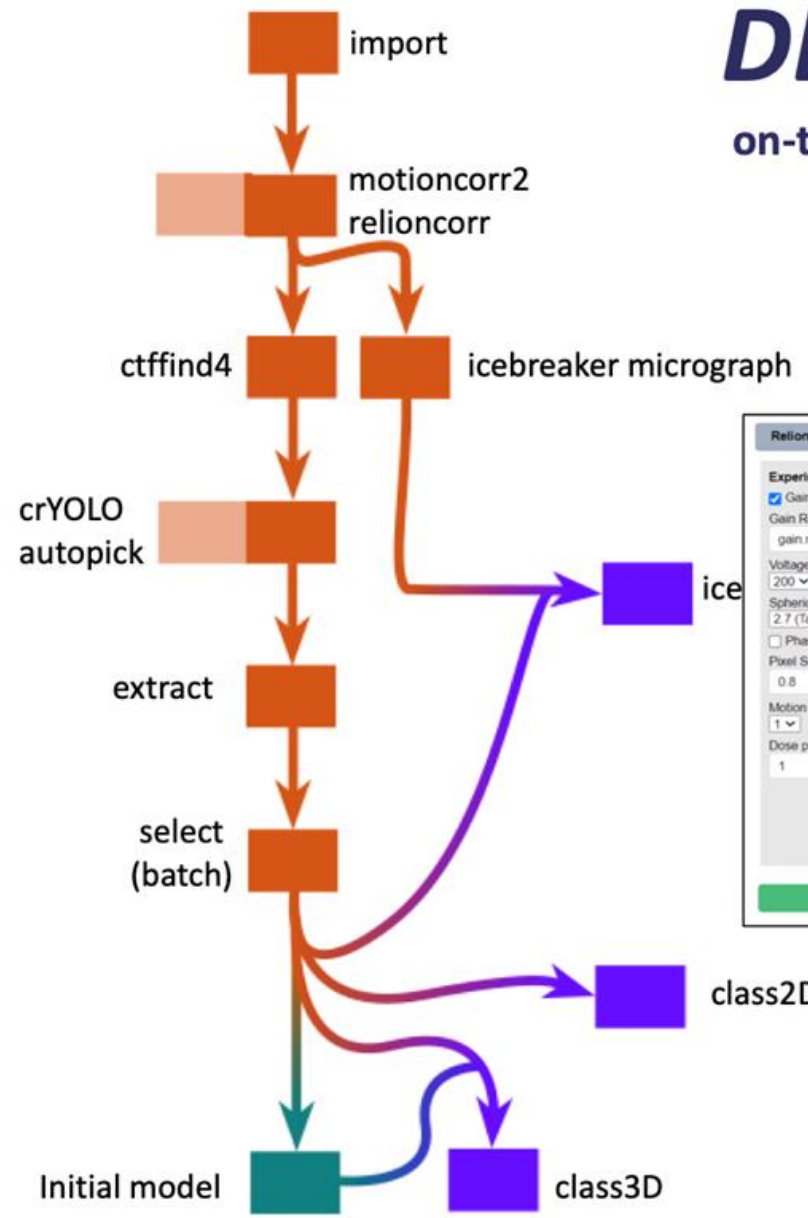


# DLS/eBIC Processing

on-the-fly batched preprocessing and evaluation



ISPyB



Dan Hatton  
Diamond Light Source



# Deposition and archiving

## Project archiving

### Full Archive

Preserve the project and all associated files



Project directory structure



All input files, except raw data



All results files

### Script archive

Preserve and repeat the workflow



Project directory structure



All input files, except raw data



Script to re-run the project



Matt  
Iadanza



Kyle  
Morris

Prepare EMPIAR deposition

**RUN** **JOB INFO** **RESET PARAMETERS**

EMPIAR

Job alias: deposit micrographs

Main

Job to create deposition from: \* PostProcess/job030

Deposit raw micrographs (if available)?  Yes  No (i)

Deposit corrected micrographs (if available)?  Yes  No (i)

Deposit particles (if available)?  Yes  No (i)

Deposit polished particles (if available)?  Yes  No (i)

Entry title \* required\* complexX

Use tracking to generate EMPIAR deposition files for raw data leading to the selected job.

Automates deposition via empiar-depositor at EBI.

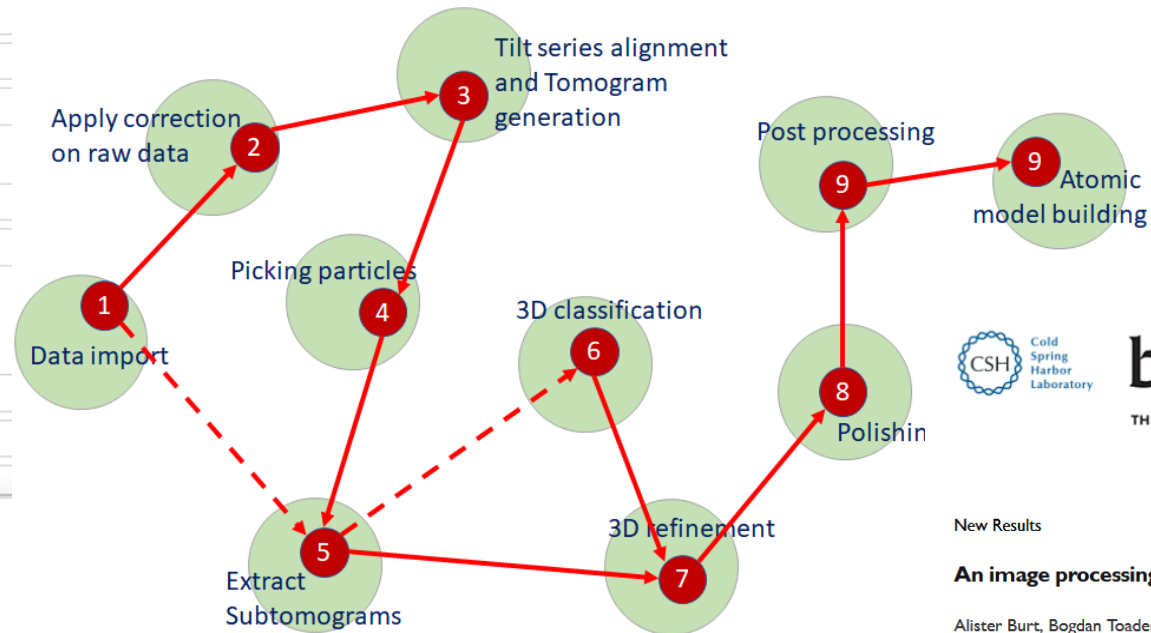


Under development: generating EMDB & PDB deposition files for final map/model submission to OneDep

# cryoET pipelines

- First pipeline based on Relion5.
- Extend for other software and options.
- Collaboration with developer community, including eBIC and RFI.

The screenshot shows the CCP-EM Doppio web interface for a project titled 'Import raw tilt images'. The interface includes a navigation menu with 'PROJECT', 'JOBS', 'NODES', and 'NEW JOB'. A search bar is present for filtering jobs. The main panel displays various input fields and options for importing raw tilt images, such as 'Tilt image files', 'mdoc files', 'Dose rate per tilt-image', 'Is dose rate per movie frame?', 'Tilt axis angle (deg)', 'MTF file', 'Invert defocus handedness?', 'Movies already motion corrected?', 'Pixel size (Angstrom)', 'Voltage (kV)', and 'Spherical aberration (mm)'. There are buttons for 'RUN', 'JOB INFO', and 'RESET OPTIONS'.



Rangana Warshamanage



bioRxiv  
THE PREPRINT SERVER FOR BIOLOGY

New Results

Follow this pr

**An image processing pipeline for electron cryo-tomography in RELION-5**

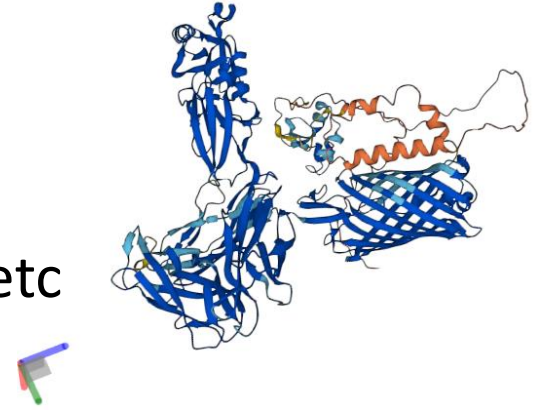
Alister Burt, Bogdan Toader, Rangana Warshamanage, Andriko von Kugelgen, Euan Pyle, Jasenko Zivanov, Dari Kimanius, Tanmay A.M. Bharat, Sjors Scheres

doi: <https://doi.org/10.1101/2024.04.26.591129>

# Molecular simulation and modelling

## Structural models:

- Phasing in crystallography via Molecular Replacement
- Interpretation of cryoEM maps, SAS, etc
- Often trimming to high confidence regions
- From experimental structures, homology modelling, AlphaFold, etc
- Genome-wide modelling



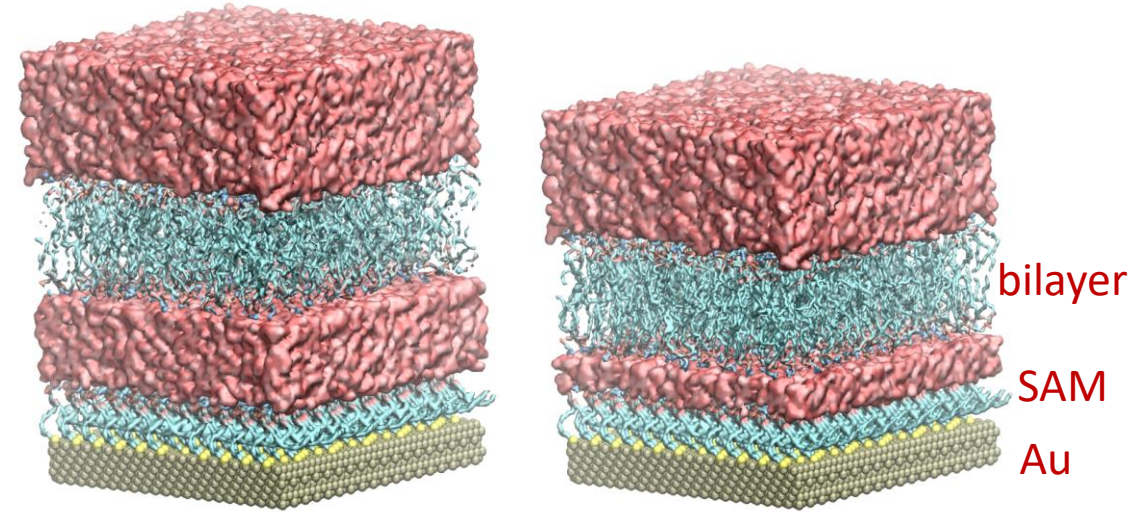
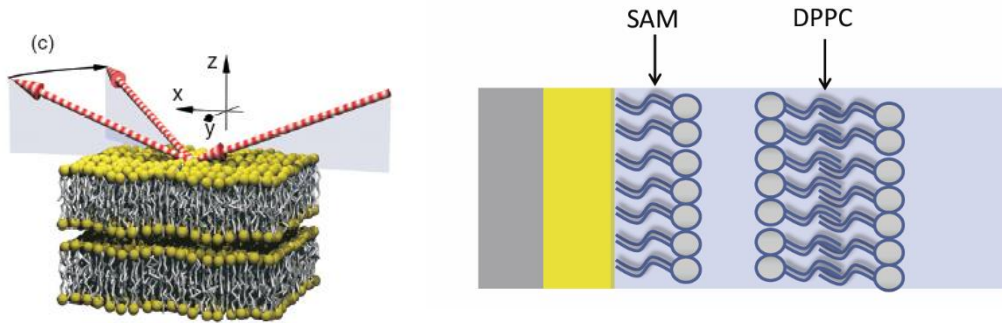
## Molecular simulation:

- Previously a largely separate discipline
- Now used for generating conformations, understanding dynamics
- *In silico* experiments





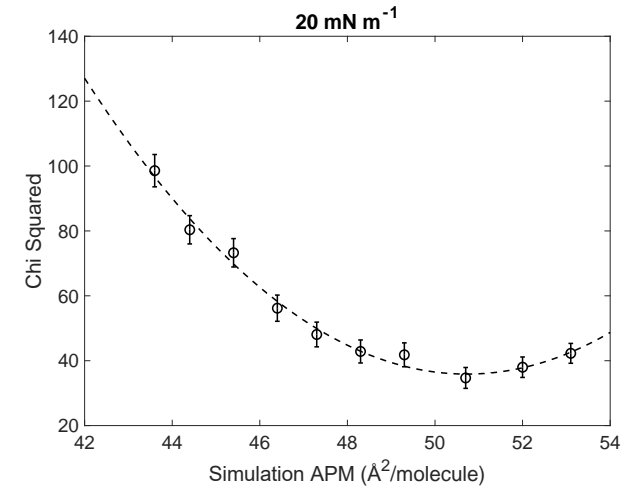
# Neutron reflectometry



- NR considers z-profile of layered structures.
- Traditionally modelled with set of continuum layers, but more accurate models from MD (prior knowledge of chemistry).

(a) lipid bilayers, (b) monolayers at air-water interface, (c) self-assembled monolayers.

- Fit to experimental data (**Rascal/RAT**) very sensitive to choice of simulation cell.



*The Analysis of Neutron Reflectivity from Langmuir Monolayers of Lipids Using Molecular Dynamics Simulations: The Role of Lipid Area.*

Arwel V. Hughes<sup>a</sup>, Valeria Losasso<sup>b</sup> and Martyn Winn<sup>c</sup>.  
(submitted)



Valeria Losasso

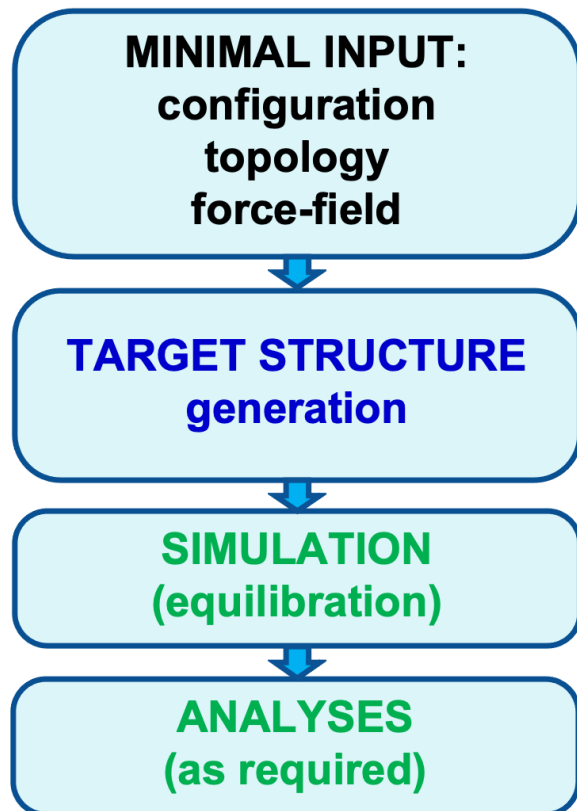


Arwel Hughes

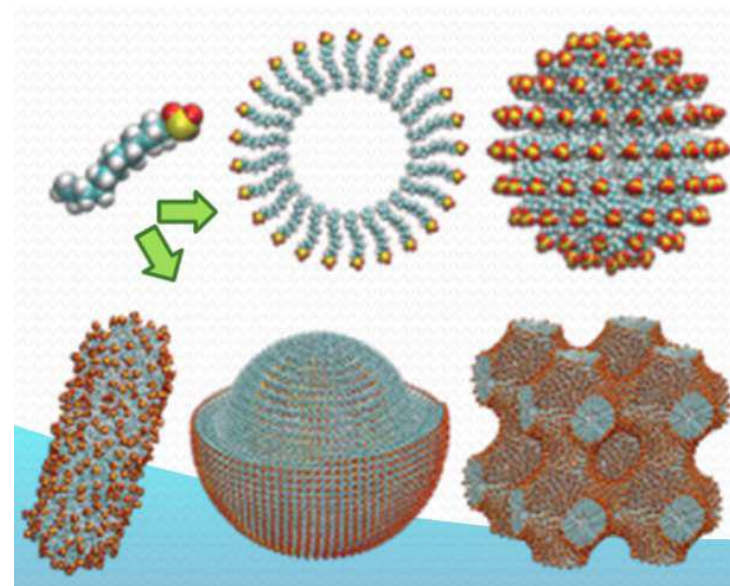
# Shapespyer

Toolkit and workflows for molecular simulations of nanostructures.

Developed for soft matter SANS.



1. Generate molecular aggregates in various geometries (layered structures being added)
2. Setup and run MD in Gromacs (NAMD being added)
3. Cluster analyses, radii of gyration, hydration layer, cavity occupation



Andrey  
Brukhno



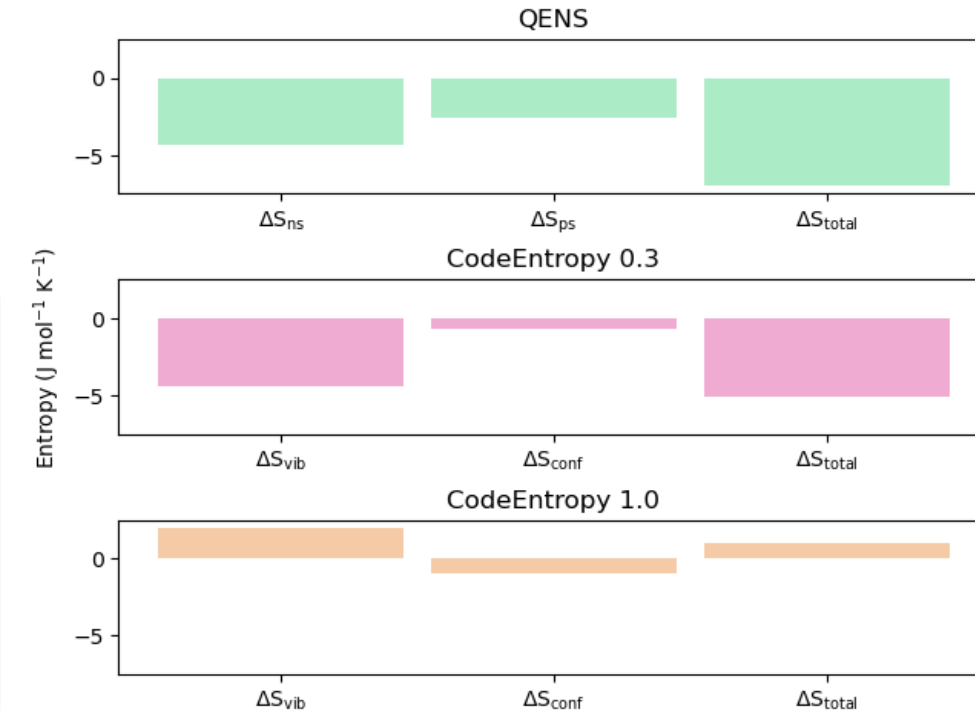
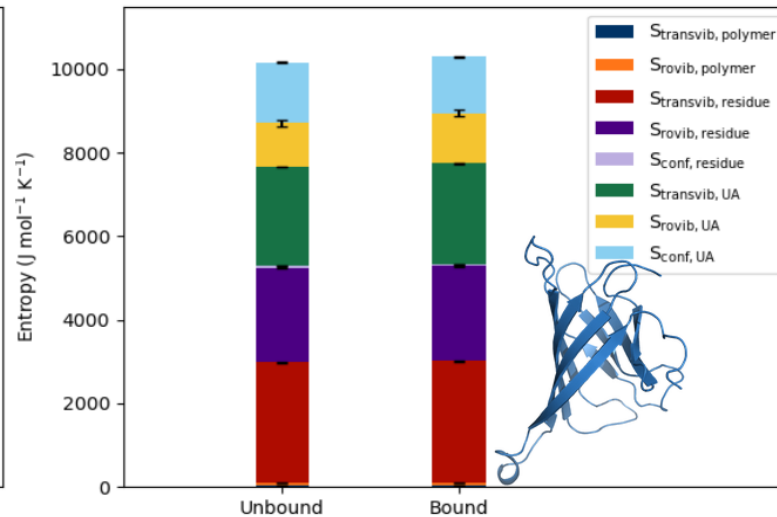
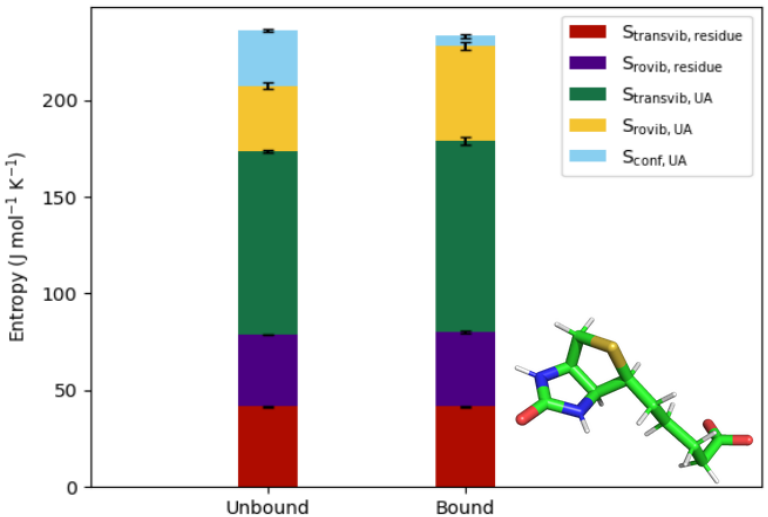
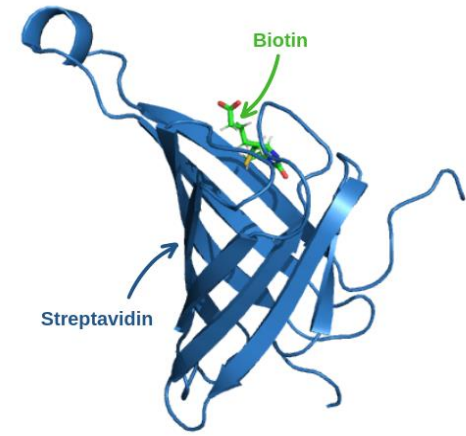
James  
Douch



Tim  
Snow

# Binding Entropy of Protein-Ligand Systems

- Molecular dynamics (MD) simulations of the streptavidin-biotin system were computed on the JADE2 HPC using GROMACS
- **CodeEntropy** - a python package for computing the entropy of macromolecular systems from forces sampled from MD simulation trajectories using multiscale cell correlation (MCC) [1]
- **Quasielastic neutron scattering (QENS)** data obtained by Sarter et al. [2] was compared with the calculated binding entropies

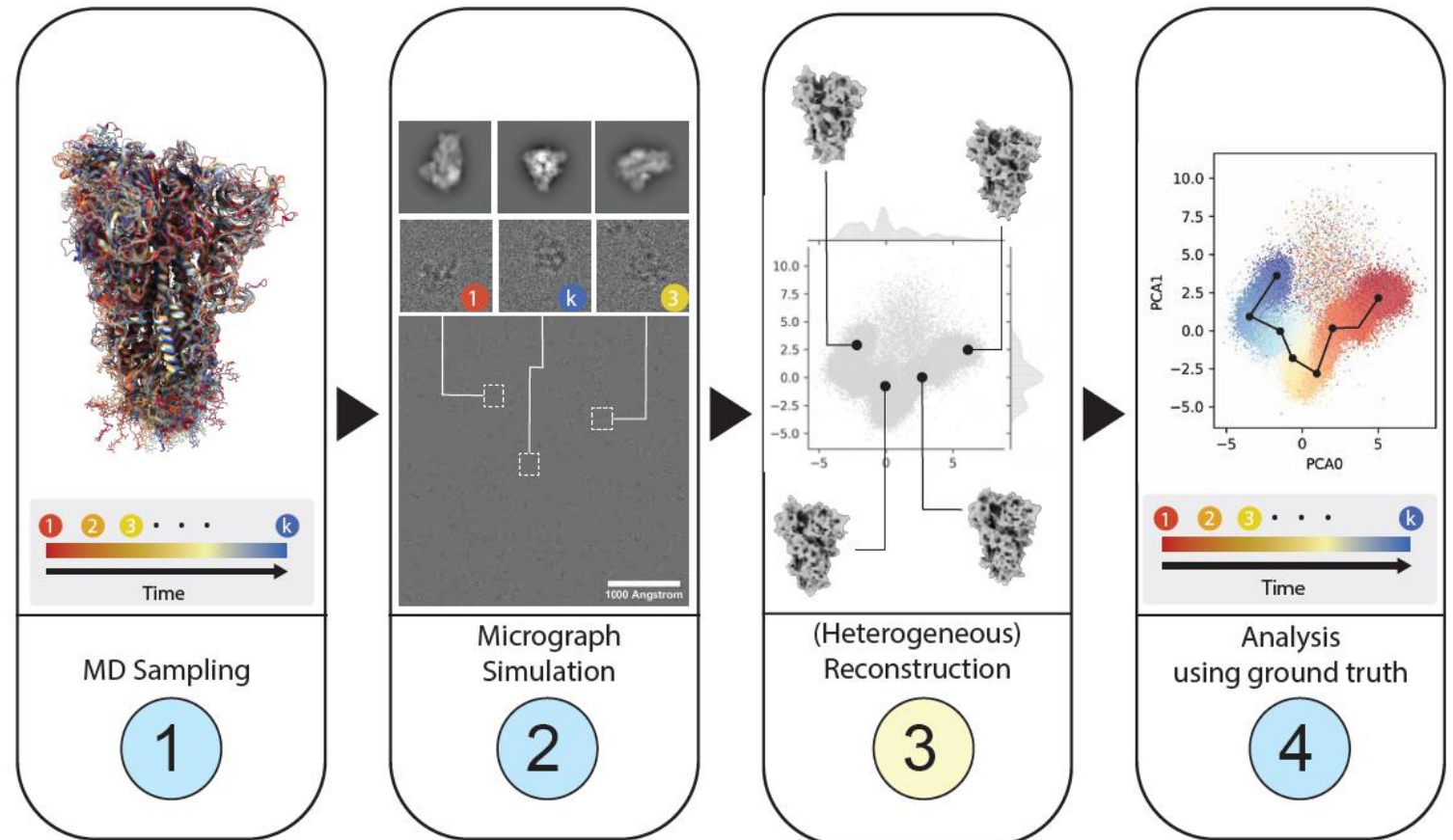


[1] Chakravorty, Arghya; Higham, Jonathan and Henchman, Richard H. (2020). *J. Chem. Inf. Model.*, **60**, 5540–5551.

[2] Sarter, Mona; Niether, Doreen; Koenig, Bernd. W; Lohstroh, Wiebke; Zamponi, Michaela; Jalarvo, Niina H.; Wiegand, Simone; Fitter, Jörg and Stadler, Andreas M. (2020) *J. Phys. Chem. B*, 2020, **124**, 324–335.

# MD and cryoEM

- Molecular dynamics explores conformational states of a protein / complex.  
*Ergodicity, simplified sample*
- CryoEM can separate conformational states into discrete or continuous classes  
*Signal to noise limit*



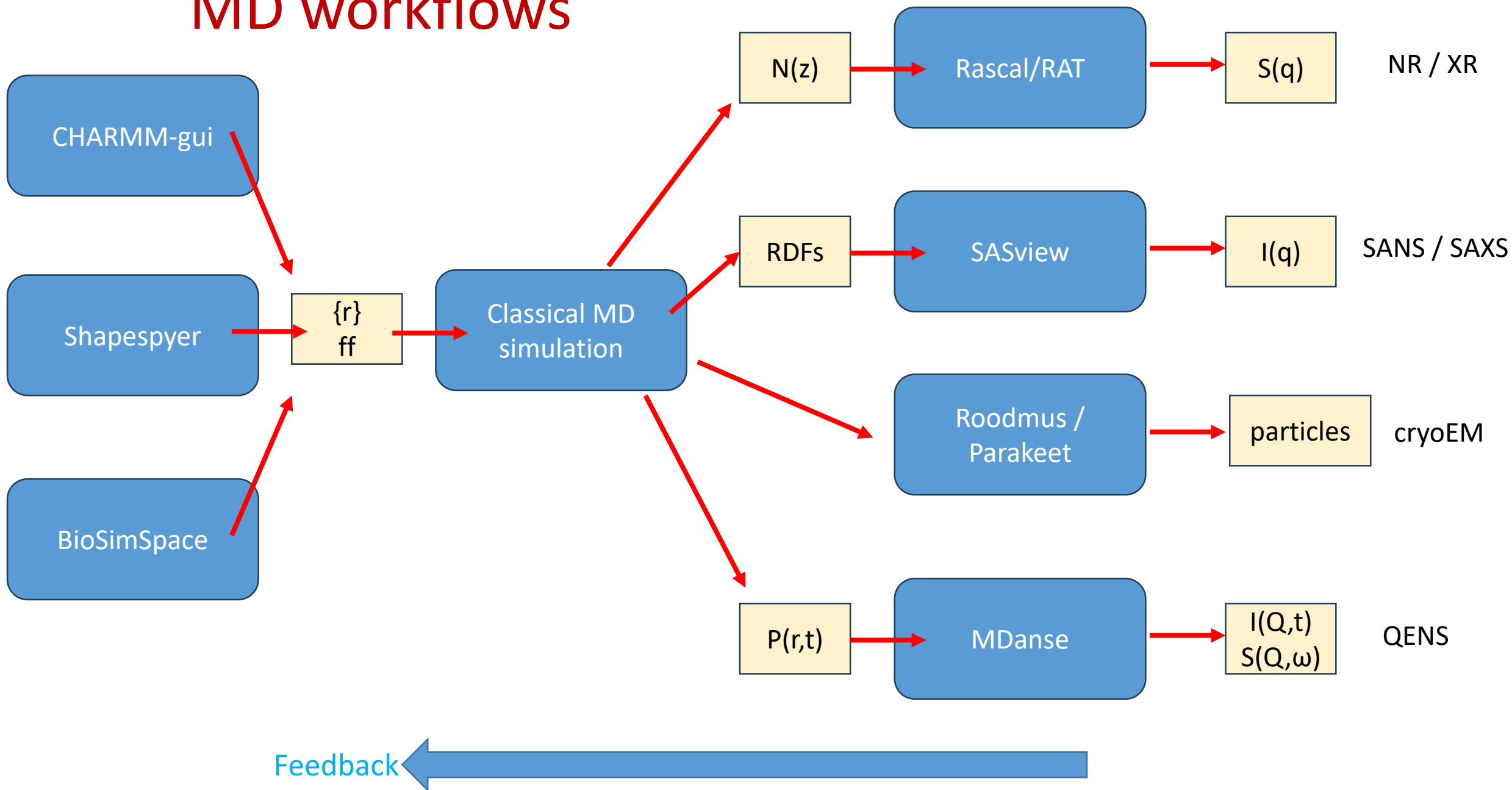
Aim to benchmark heterogeneous reconstruction algorithms (e.g. CryoDRGN and 3DFlex) using synthetic particle sets from MD where ground truth known.

**Roodmus: A toolkit for benchmarking heterogeneous electron cryo-microscopy reconstructions**

Maarten Joosten,<sup>a†</sup> Joel Greer,<sup>b†</sup> James Parkhurst,<sup>c,d</sup> Tom Burnley<sup>b\*</sup> and Arjen J. Jakobi<sup>a\*</sup>

Delft, SCD, RFI

# MD workflows



# Data / metadata management

*We do not handle data storage or data transfer.  
Better experts in SCD!*

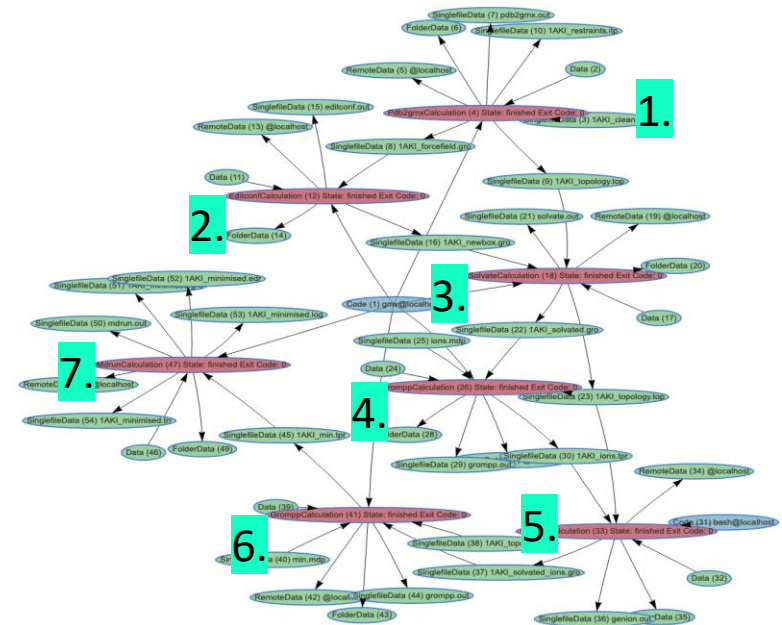
## Our workflows:

- Handle file format conversions
- Store and track metadata
- Connect to external data stores
- Deposit in international repositories

Good data management a pre-requisite for ....



**PSDI**  
PHYSICAL SCIENCES  
DATA INFRASTRUCTURE



<https://aiida-gromacs.readthedocs.io>

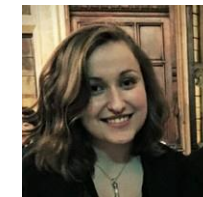
# Machine learning

- Early adopters of ML for interpretation of scientific data
- Looking to adapt established models to our datasets / processing
- Have worked closely with SciML group and Alan Lowe's group at ATU

*> 50% of any project is preparing / organising the training data ...*

- Many methods coded in the Macromolecular Machine Learning Toolbox  
<https://gitlab.com/ccpem/ml-protein-toolbox>
- Besides our own tools, incorporate 3<sup>rd</sup> party tools e.g. AlphaFold, Blush (Relion 5)
- Working on integrating CryoDRGN

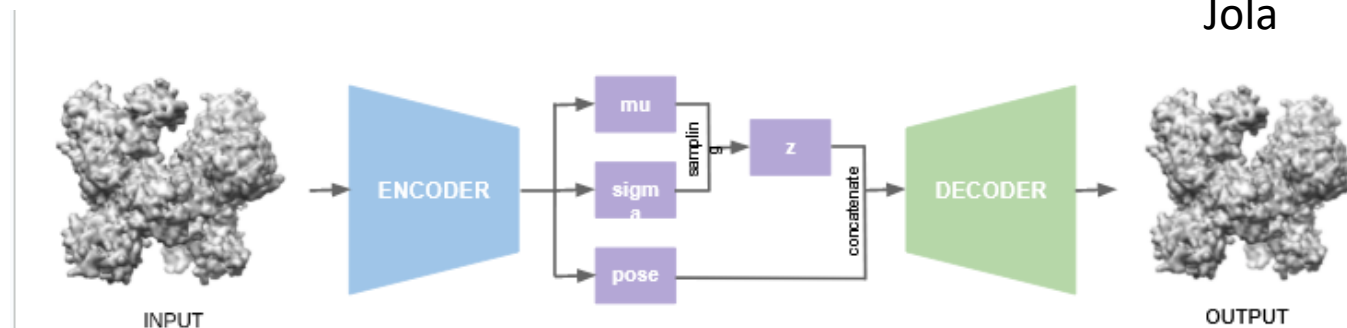
# Affinity-VAE



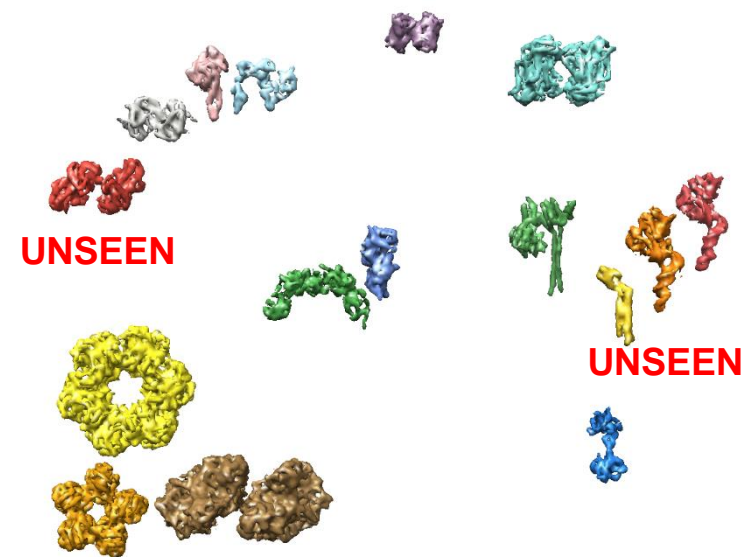
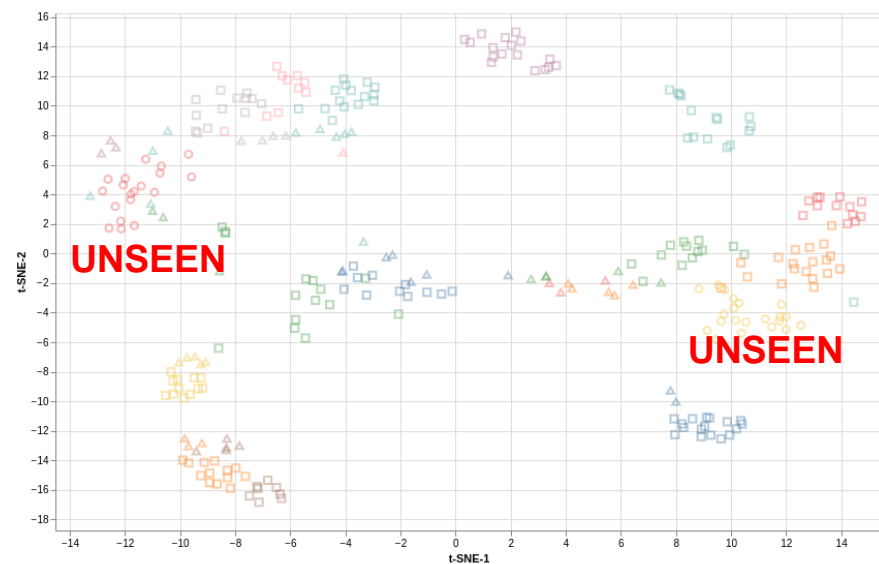
Jola

## Identification of molecules in cryoET

- beta-VAE optimisation function plus affinity regulariser
- Based on similarity of two 3D objects
- Helps to disentangle of inter-class variation from intra-class (rot, trans)

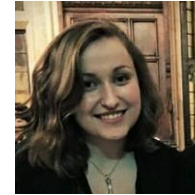


CCP-EM, ATI, RFI, SciML





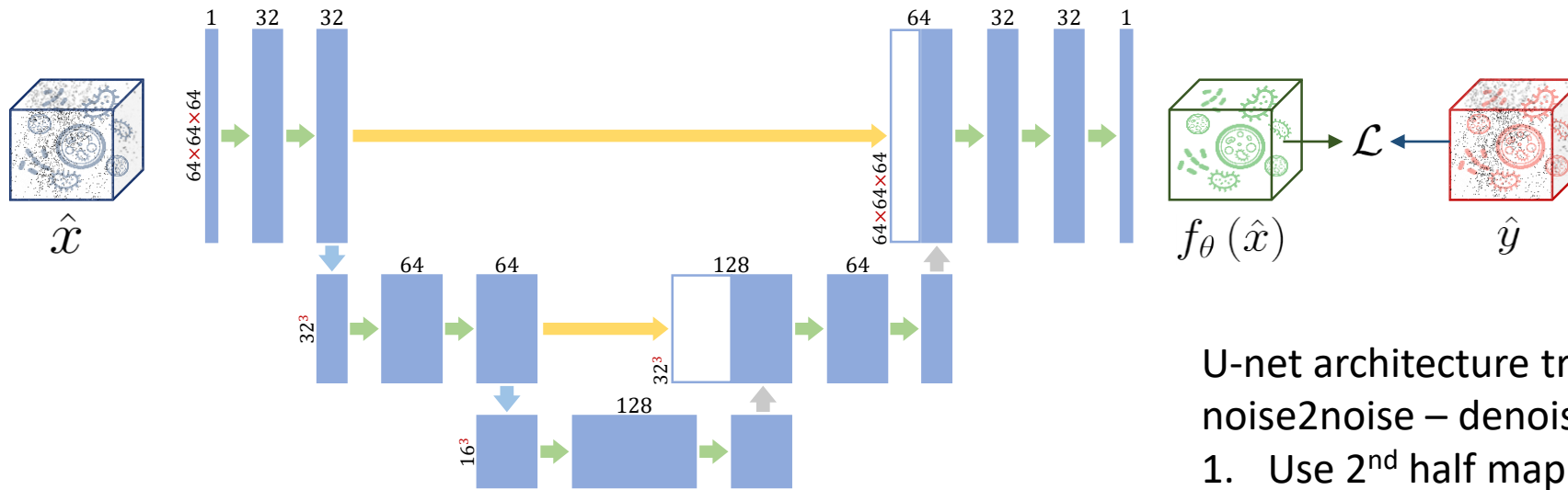
# Denoising of electron tomograms



Jola



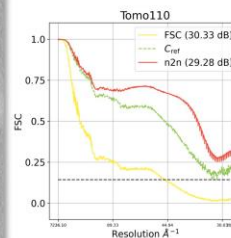
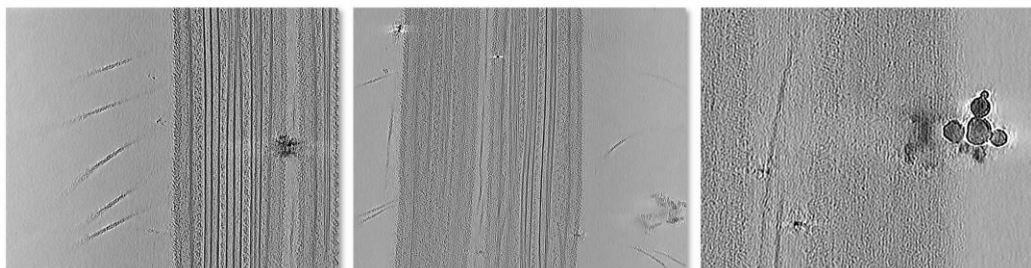
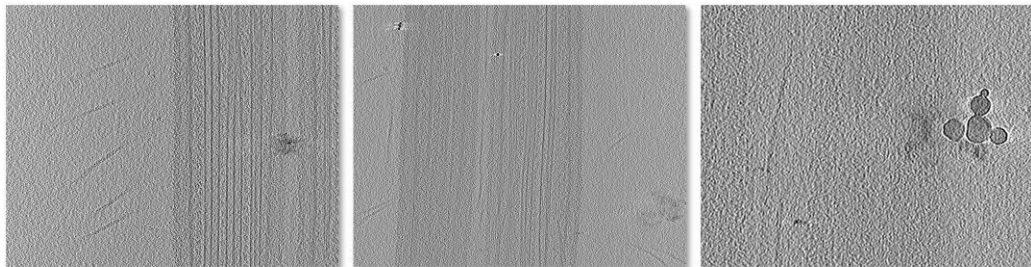
Ding



U-net architecture trained with 64 x 64 x 64 patches.  
noise2noise – denoising without clean data:

1. Use 2<sup>nd</sup> half map as noisy target
2. Randomly remove high frequency components to give corrupted target

Pros and cons in efficacy, speed and data requirements.



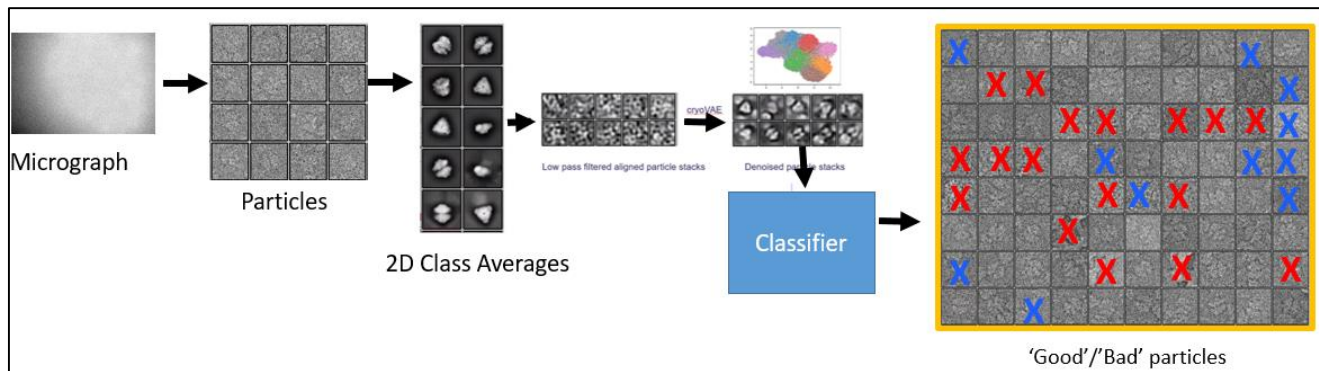
Feeds into tomography projects with eBIC and RFI.

# cryoDANN: automated evaluation of particles



Sony

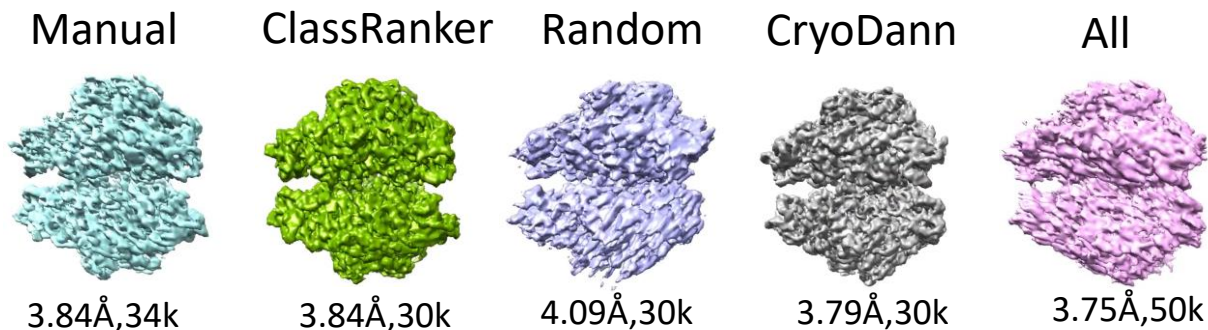
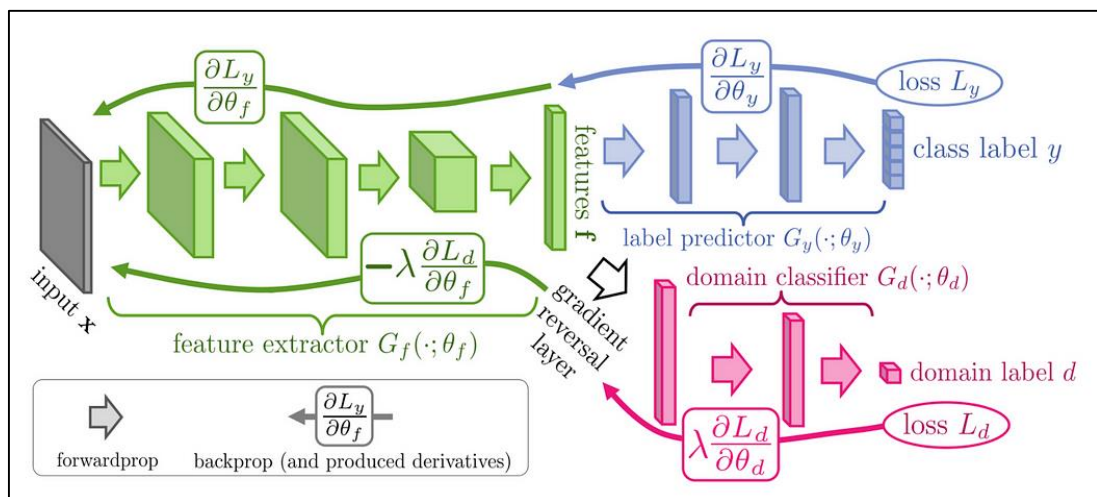
Objective: Annotating particle images



- Joint project with Diamond/eBIC, SciML
- ML based method for an automated evaluation of good/bad particle selection
- Tested the classifier on EMPIAR datasets
- Installed and being tested at eBIC (Yuriy Chaban, Dan Hatten)

**Aim: Automate the particle selection step in data processing step**

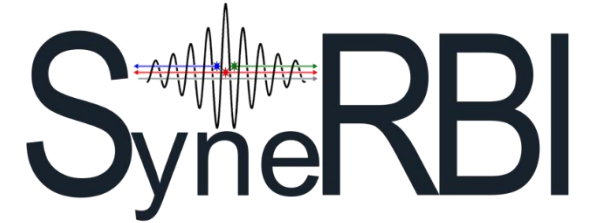
CryoDANN:



Results: EMPIAR-10547

# Computed tomography

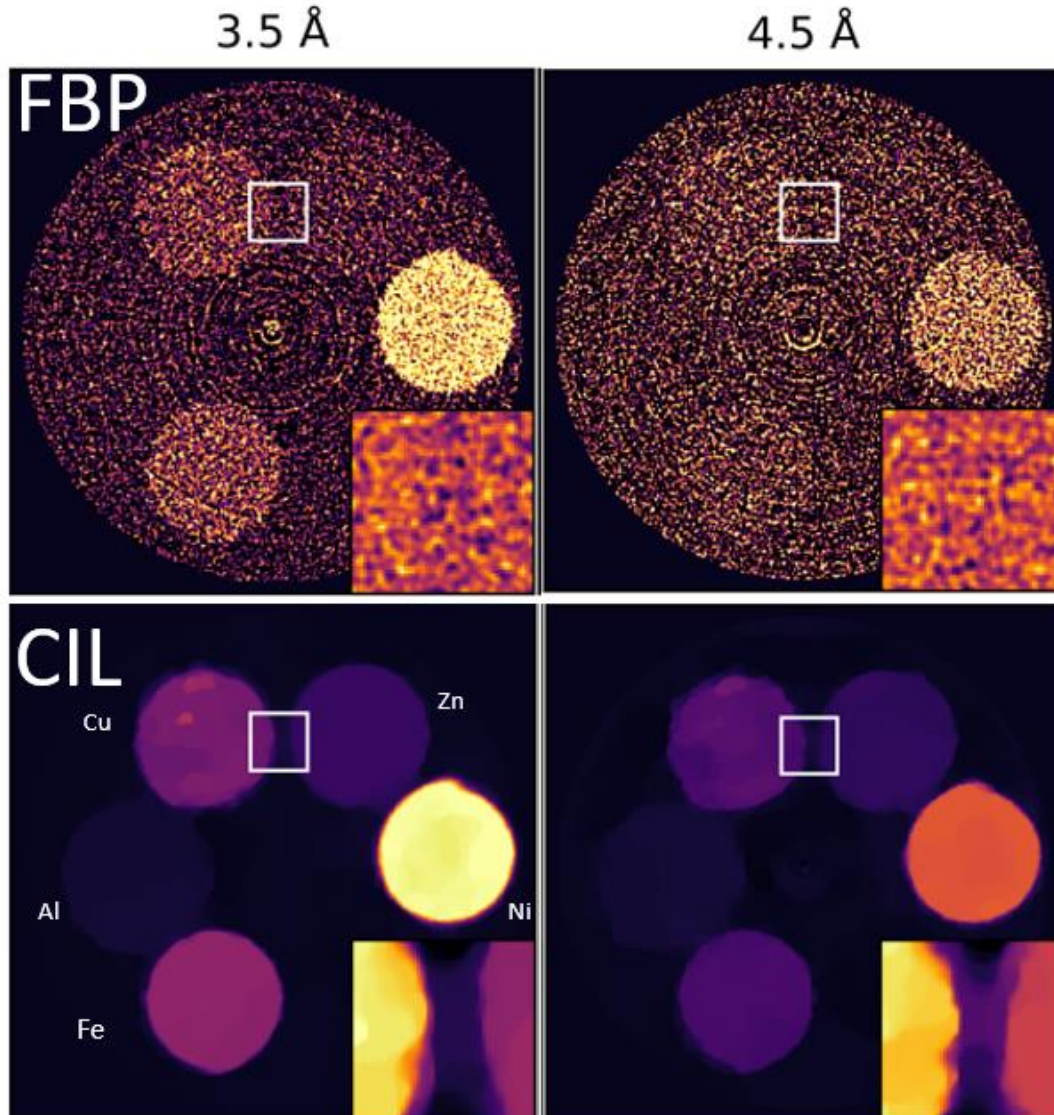
- CCPi mainly materials, but formally in the Comp Bio theme and not covered elsewhere
  - CCPSyneRBI is medical imaging
  - On-going collaborations with ISIS, Diamond and EPAC
  
  - Most software developments centred on the Core Imaging Library (CIL )
- <https://ccpi.ac.uk/cil/>
- Underpinning for inverse problems



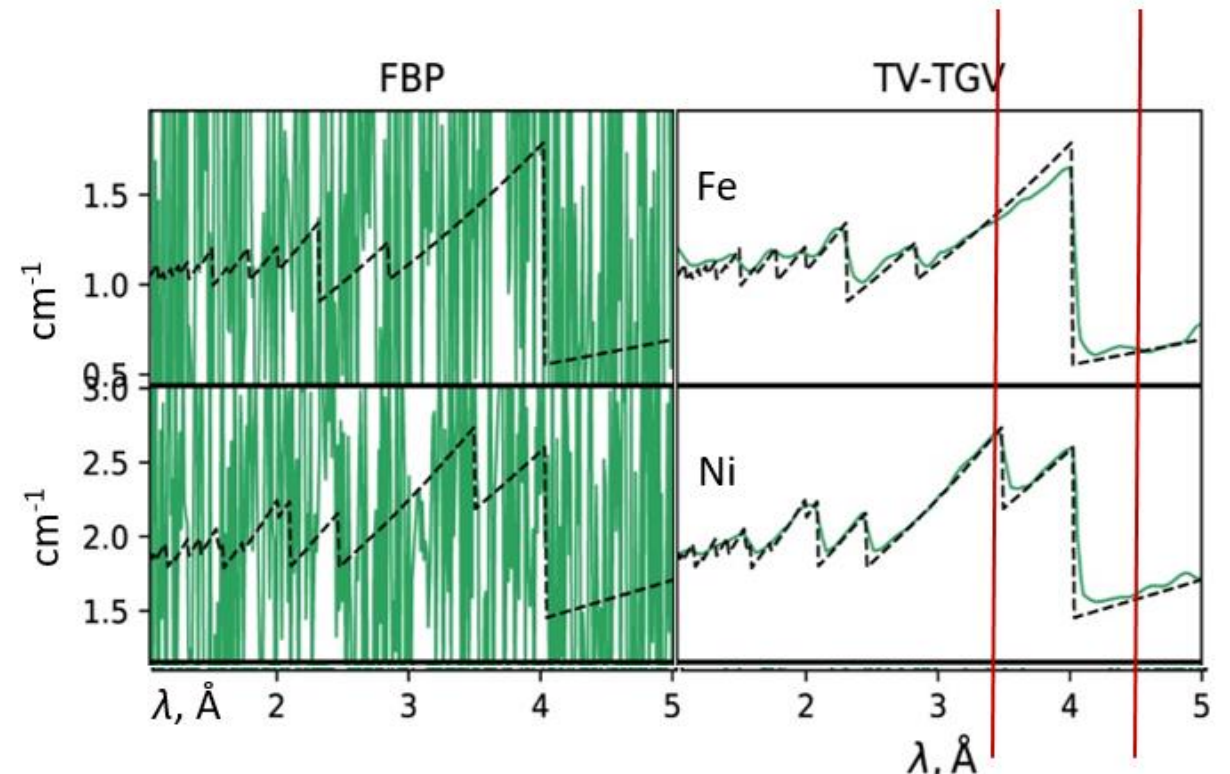
<https://ccpi.ac.uk/>

<https://www.ccpsynerbi.ac.uk/>

# Energy-resolved neutron CT



- Data from ISIS/IMAT
- Proposed spatio-spectral TV-TGV regularization
- Enables clear identification of Bragg edges in 3D



Ametova et al. 2021: *Crystalline phase discriminating neutron tomography using advanced reconstruction methods*, J. Physics D, <https://doi.org/10.1088/1361-6463/ac02f9>

# Digital Volume Correlation

Alignment / interpretation of 4D data.

Collaboration with **Brian Bay** (Oregon).

Porting, parallelisation, GUI.

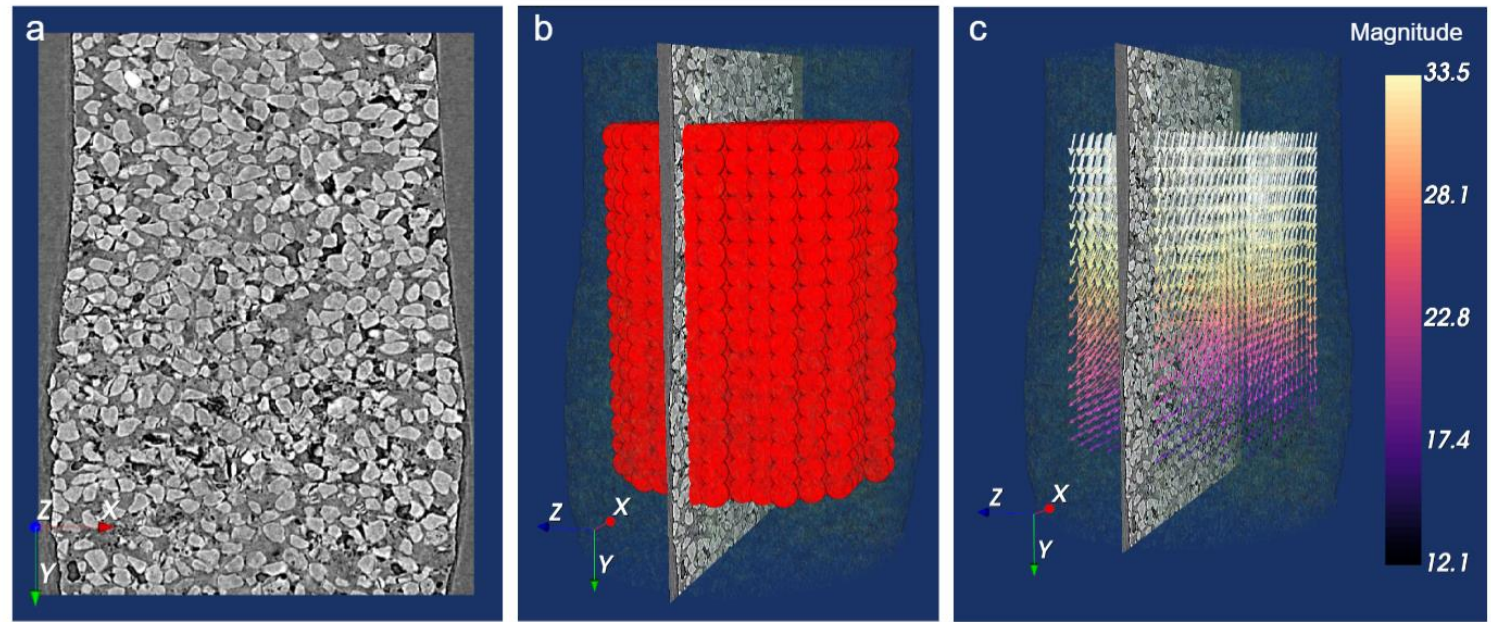


Figure 1: CT reconstruction of magma data. a) slice view, b) region of interest with point cloud and spherical subregions, c) displacement vector display in unit of pixels.

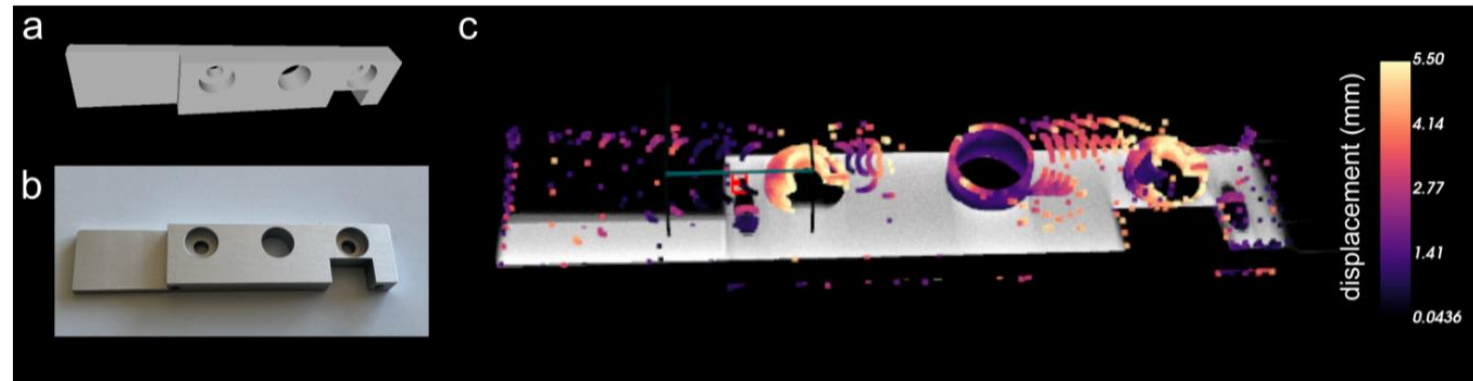
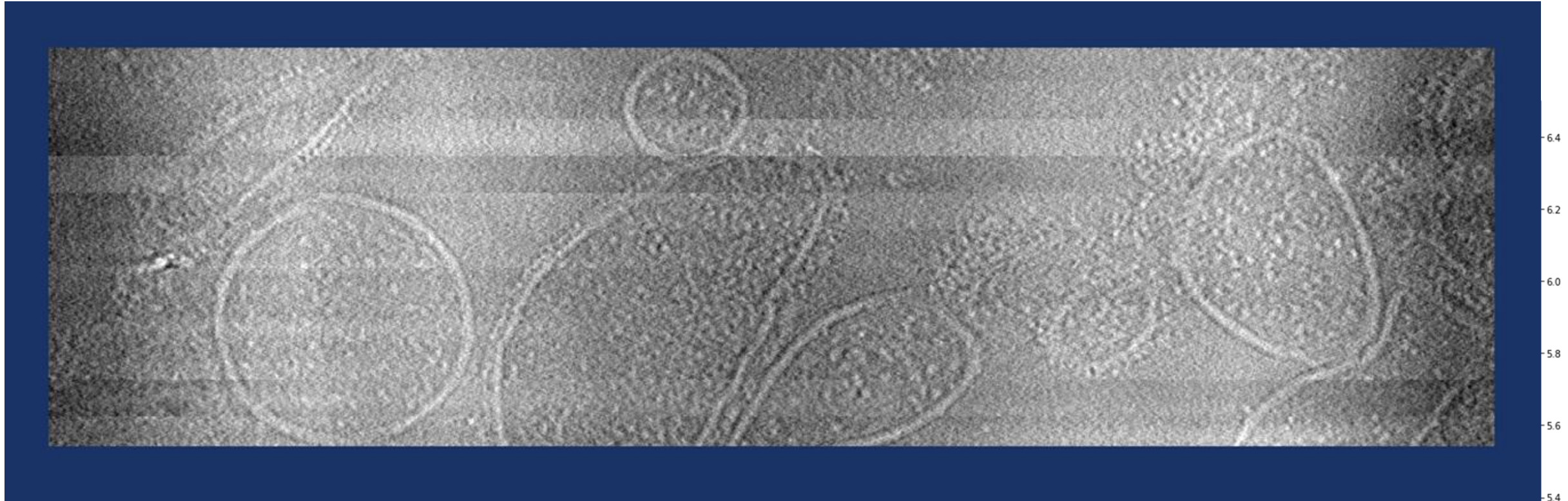


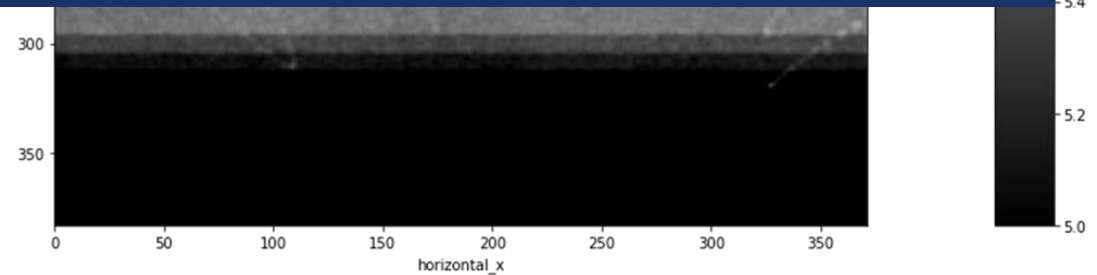
Figure 2: Mechanical component of an in-situ tensile testing machine. a) Original CAD model, b) component manufactured from the CAD model, c) displacement field between the 3D volumes of the experimental CT scan and of the virtual CT scan of the CAD model.

# Reconstructions in CIL for cryoET

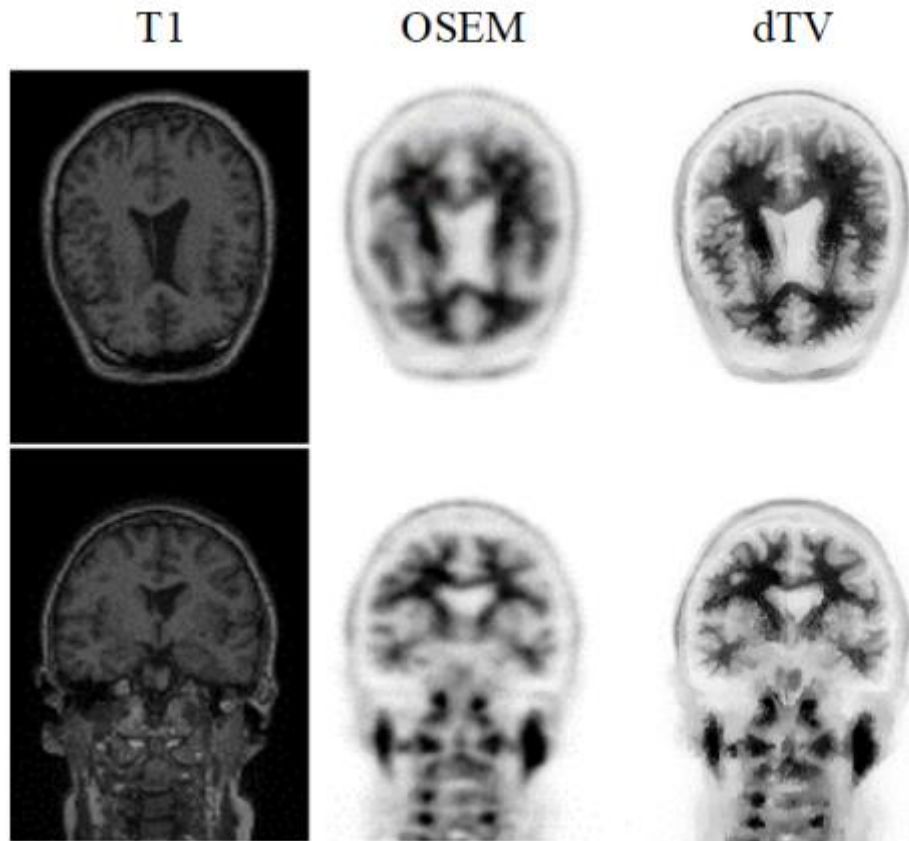


RyR1 receptor proteins embedded in sarcoplasmic reticulum vesicles extracted from mammalian cells.

EMPIAR-10349



# Anatomically guided PET super resolution



Cardio-respiratory MRI motion compensated image reconstruction. →

In preparation and <https://royalsocietypublishing.org/doi/10.1098/rsta.2020.0208>

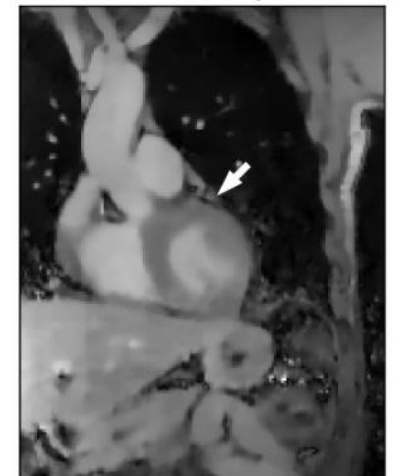
Synergistic Image Reconstruction Framework (uses CIL)

**SIRF**

PDHG 20 epochs



SPDHG 20 epochs



<https://www.medrxiv.org/content/10.1101/2023.04.23.23289004v1>

Super-resolution PET reconstructions, guided by MRI. dTV-regularised, two stages.

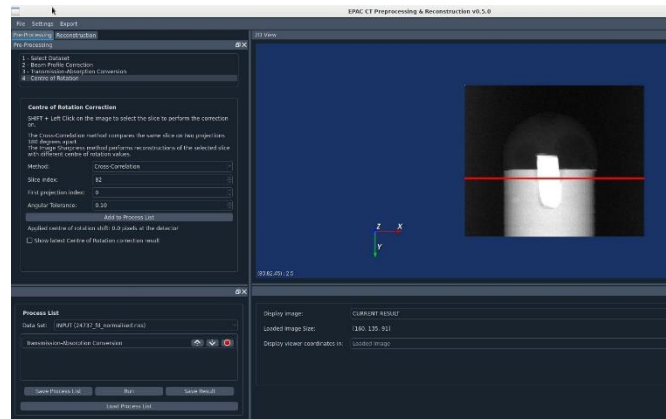
# EPAC

## GUI + Pipeline for CT recon with CIL

- Handles Nexus & lab-based XCT data files
- Pre Processing
- Reconstruction



Danica

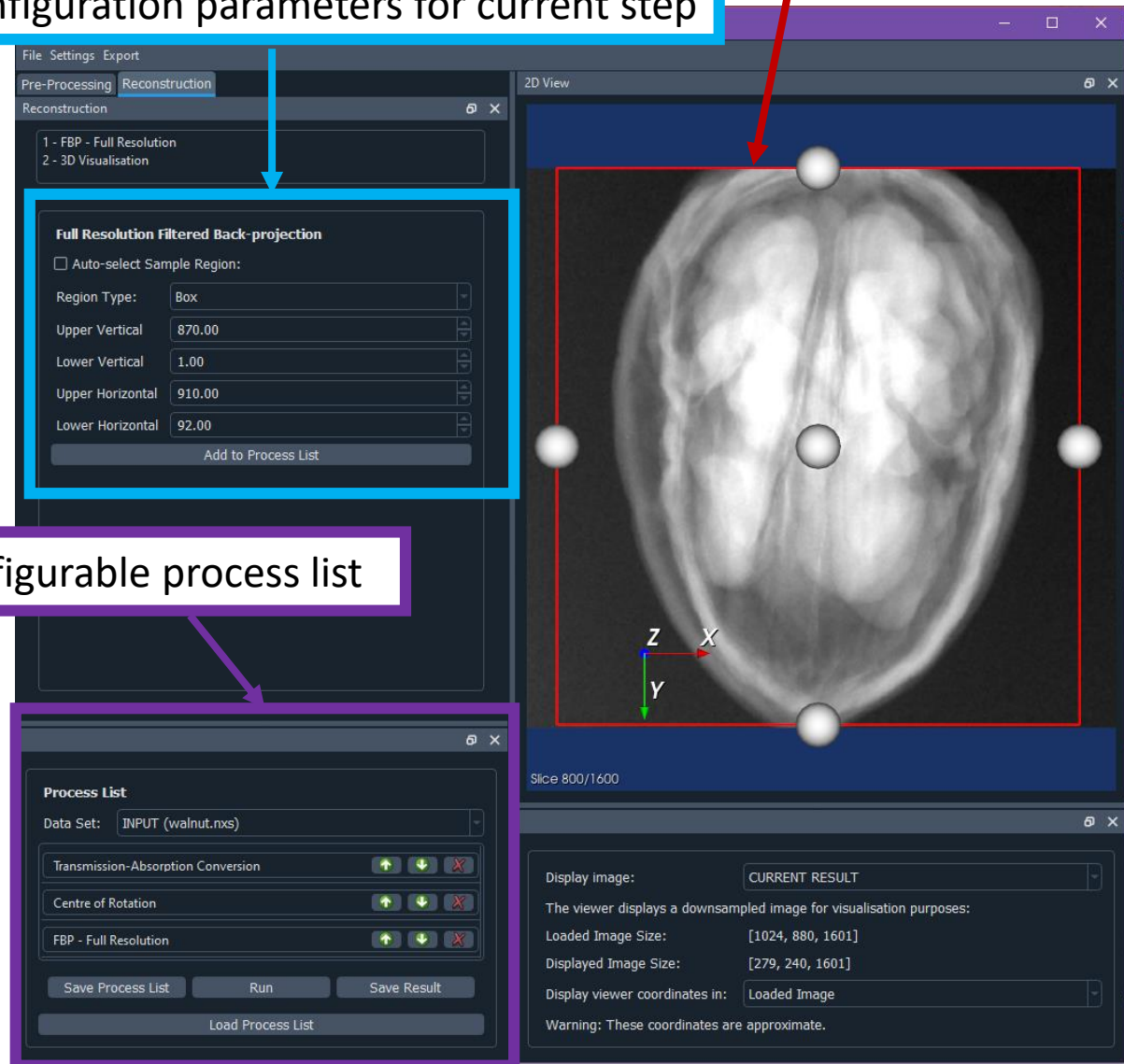


CIL-GUI

Configuration parameters for current step

Visualisation + interaction with data at each step

Configurable process list





# Digital Twinning

A widely used and inconsistently used term ...

Here we consider simulated experiments:

- Generate large amounts of synthetic data

- Full control over sample and instrument parameters, allowing *in silico* trials

Doesn't replace real data.

But useful for development and understanding.

# Parakeet: Digital Twin for cryoEM/cryoET



James

**OPEN BIOLOGY**  
royalsocietypublishing.org/journal/rsob

Methods & Techniques

**Parakeet: a digital twin software pipeline to assess the impact of experimental parameters on tomographic reconstructions for cryo-electron tomography**

Cite this article: Parkhurst JM, Dumoux M, Basham M, Clare D, Siebert CA, Varslot T, Kirkland A, Naismith JH, Evans G. 2021 Parakeet: a digital twin software pipeline to assess the impact of experimental parameters on tomographic reconstructions for cryo-electron tomography. *Open Biol.* 11: 210160. <https://doi.org/10.1098/rsob.210160>

James M. Parkhurst<sup>1,2</sup>, Maud Dumoux<sup>1</sup>, Mark Basham<sup>1,2</sup>, Daniel Clare<sup>2</sup>, C. Alistair Siebert<sup>2</sup>, Trond Varslot<sup>4</sup>, Angus Kirkland<sup>1,3,5</sup>, James H. Naismith<sup>1,6</sup> and Gwyndaf Evans<sup>1,2</sup>

<sup>1</sup>Rosalind Franklin Institute, Harwell Science and Innovation Campus, Didcot OX11 0FA, UK  
<sup>2</sup>Diamond Light Source, and <sup>3</sup>Electron Physical Science Imaging Centre, Diamond Light Source, Harwell Science and Innovation Campus, Didcot OX11 0DE, UK  
<sup>4</sup>Thermo Fisher Scientific, Masimila Pecha, Brno, Czech Republic  
<sup>5</sup>Department of Materials, University of Oxford, Parks Road, Oxford OX1 3PH, UK  
<sup>6</sup>Division of Structural Biology, University of Oxford, Roosevelt Drive, Oxford OX3 7BN, UK

JMP, 0000-0002-9120-8354; MD, 0000-0002-1732-1041; MB, 0000-0002-8438-1415; CAS, 0000-0002-3177-8178; CL, 0000-0002-8176-1979; TV, 0000-0002-7301-6404; DC, 0000-0002-3177-8178

**Roodmus: A toolkit for benchmarking heterogeneous electron cryo-microscopy reconstructions**

Maarten Joosten,<sup>a†</sup> Joel Greer,<sup>b†</sup> James Parkhurst,<sup>c,d</sup> Tom Burnley<sup>b,\*</sup> and Arjen J. Jakobi<sup>a,\*</sup>

<sup>a</sup>Department of Bionanoscience, Kavli Institute of Nanoscience, Delft University of Technology Delft, 2629 HZ Delft, The Netherlands, <sup>b</sup>Science & Technology Facilities Council, Research Complex at Harwell, Oxon, OX11 0FA, United Kingdom, <sup>c</sup>Rosalind Franklin Institute, Harwell Science and Innovation Campus, Oxon, OX11 0QS, United Kingdom, and <sup>d</sup>Diamond Light Source, Harwell Science and Innovation Campus, Oxon, OX11 0DE, United Kingdom. Correspondence e-mail: tom.burnley@stfc.ac.uk, a.jakobi@tudelft.nl

Conformational heterogeneity of biological macromolecules is a challenge in single particle averaging (SPA). Current standard practice is to employ classifica-

**Parakeet:** initiative of the RFI.

Digital Twin of electron tomography experiment. Includes geometry, microscope parameters, acquisition strategy, noise models, etc.

As part of **Roodmus**, we adapted Parakeet for single particle experiments (CCP-EM, TU Delft).

We generated simulated cryoEM datasets starting with atomic models generated in a Molecular Dynamics simulation. Important for development (know ground truth).

Rosalind Franklin Institute

The Rosalind Franklin Institute

CONTACT

Open Datasets at the Franklin

Home → Artificial Intelligence and Informatics → Open Datasets at the Franklin

RFI is hosting synthetic cryoEM datasets generated by RFI, CCP-EM and ATI. Available as Globus collection.

<https://www.rfi.ac.uk/projects/open-datasets/>

# Computed Tomography Digital Twin

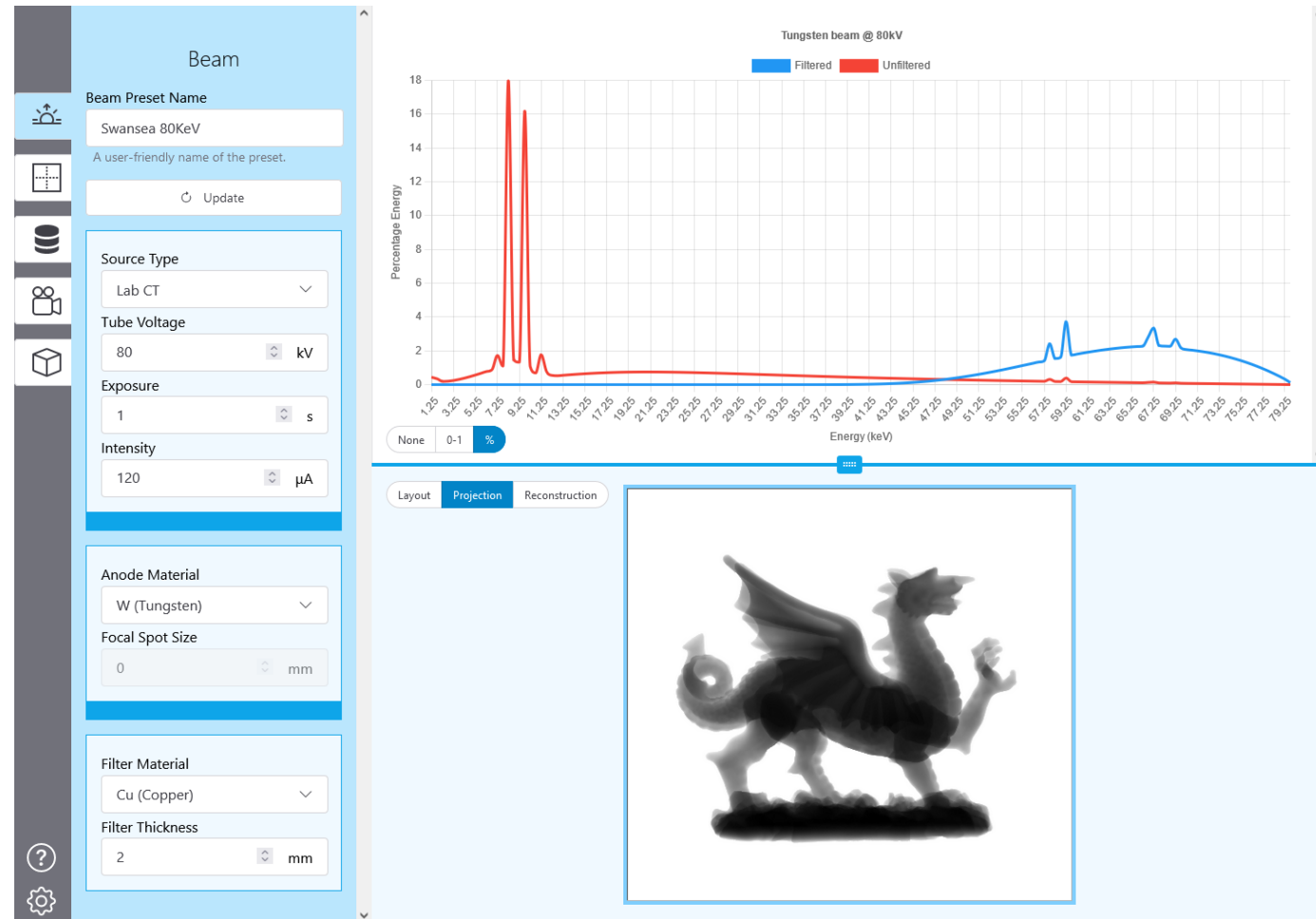
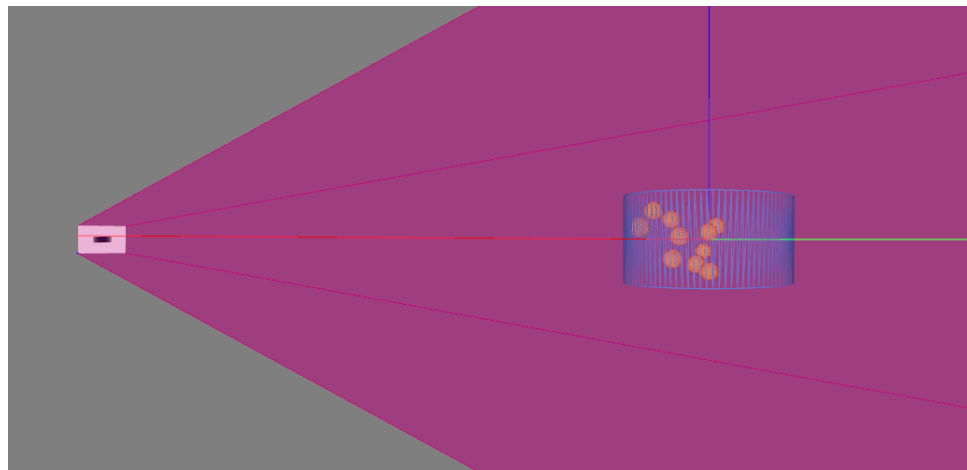
Franck



Developing a Digital Twin for CT beamlines.

Based on gVirtualXRay, a C++ library implemented on GPUs to simulate X-ray imaging.

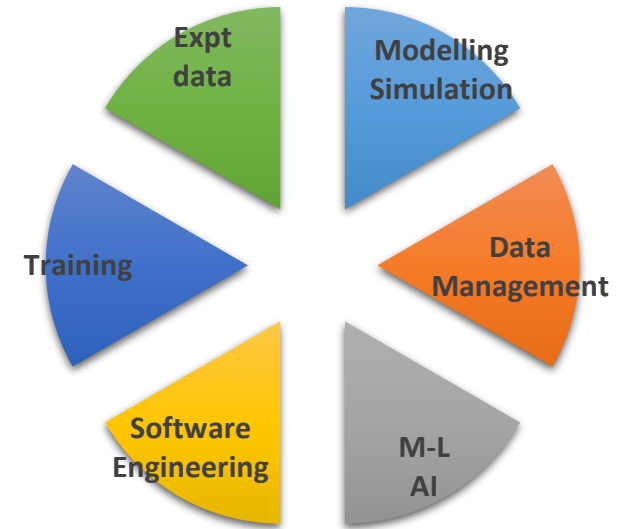
Uses Beer-Lambert law to compute the absorption of light (i.e. photons) by 3D objects (here polygon meshes).



<https://gvirtualxray.fpvidal.net/>

# Future engagement

- We want to build on our scientific strengths
  - but expand into new areas
- Software engineering skills
  - but we are not generic RSEs
- Increased use of machine learning
  - optimise data processing (efficiency)
- Automation via software pipelines
  - but bespoke application work for end users useful for engagement and experience
- Interested in studentships and placements



Software across facilities  
MD simulation for interpretation  
Imaging across lengthscales  
Metadata capture and archiving

# Thanks and contacts

Please come talk to us!



Edoardo  
Pasca



Eugene  
Krissinel



James  
Gebbie-Rayet



Martyn  
Winn



Tom  
Burnley



Science and  
Technology  
Facilities Council

# Questions?