# DOI, FAIR, an MX COVID-19 Case Study

**Author: Frank Von Delft**

**Affiliation: Diamond Light Source**

**26th October 2021**

# Driver

Urgency to contribute to global pandemic response

- XChem experiment on SARS-Cov-2 Main Protease (Mpro): <u>starting data for drug discovery</u>

- Results *(derived data)* made public for non-specialist community: Fragalysis Cloud

- Triggered COVID Moonshot – *crowdsourced open drug discovery effort*



Requests from specialist community to make raw data p

- With necessary metadata

- With links to published structures

# XChem experiment (Animation)



→ 2000 diffraction datasets (8Gb each)

→ 80 protein crystal structures

→ rich metadata (kbs)

# Publishing in Fragalysis
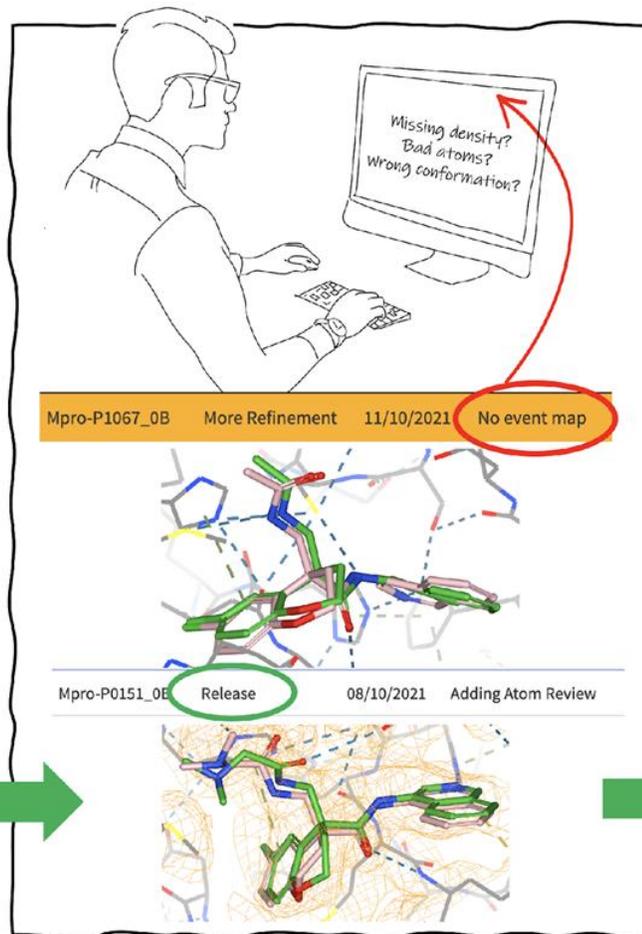


Curate

Review

Publish

# Overall: lengthy experiment, analysis & release

# FAIR data in Macromolecular Crystallography (MX)

- Protein Data Bank (PDB) since 1971

- Structures made public
  - more recently (20 years) with "reduced data" *(=not raw)*
  - allows re-analysis

- Most relevant data and metadata included

- Free to all users

- PDB uniquely defined by 4-character name e.g. 5RE7

# Action for *raw data*

- Use Zenodo: well-defined API, and familiarity
- "Mail merge" spreadsheet with general data
  - *(authors, text descriptions etc.)*
- Zip data for convenience for end user
  - (one file / deposition)
- Programmatic uploaded Python)
  - (unpublished) API
  - good help from developers
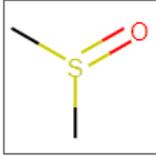
- Main work: *organising the metadata*

# Links

- List of Protein Data Bank (PDB) ID / Zenodo DOI sent to the PDB team who set up the data DOI links

- Links are bidirectional: **Zenodo** ⟵⟶ **PDB**

- **Easier to find for users.**

# (Specialist) User perspective

- The data can be **downloaded** and used with **standard software**
- No "site" file is needed, all metadata needed are present in the dataset
- The results can be reproduced using the model in PDB
- MX users think in terms of PDB IDs ➔ We are making it the main persistent identifier
- There are no entry barriers (registration etc.)
- Open format – users are familiar with Zip files

## The data are  Findable  Accessible  Interoperable  Reusable

# For non-specialist user



- The data can be **viewed, downloaded** and annotated
- No pre-knowledge is needed…
- Trivial to share:  click "share", copy URL

## The data can *also* be Implementation: HARD

# Outcome

- COVID Moonshot consortium receives £8m funding from Wellcome on behalf of the COVID-19 Therapeutics Accelerator

- Funds pre-clinical phase:  develop 5 candidates to clinical trials

- A growing collaboration:



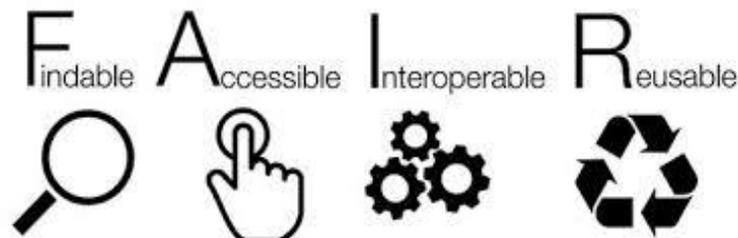COVID Moonshot collaborators (left to right): Ben Perry, Discovery Open Innovation Leader at DND*i*; Frank von Delft, Professor of Structural Chemical Biology at the University of Oxford and Principal Beamline Scientist at Diamond Light Source; John Chodera, Associate Member at the Memorial Sloan Kettering Cancer Center and founding member of the Folding@home Consortium; Annette von Delft, Translational Scientist at the University of Oxford; Ed Griffen, Technical Director and Co-founder of MedChemica; Alpha Lee, Chief Scientific Officer at PostEra and Faculty Member at the University of Cambridge.

**2nd European Photon & Neutron EOSC Symposium**

26 October 2021

# Thank you

**Contact details:** *frank.von-delft@diamond.ac.uk*