# The DRI Project and forthcoming activities to support metadata-rich depositions

## Genevieve Evans

**Structural Biologist
& wwPDB data curator**

**Email: gle@ebi.ac.uk**

**PDBe**
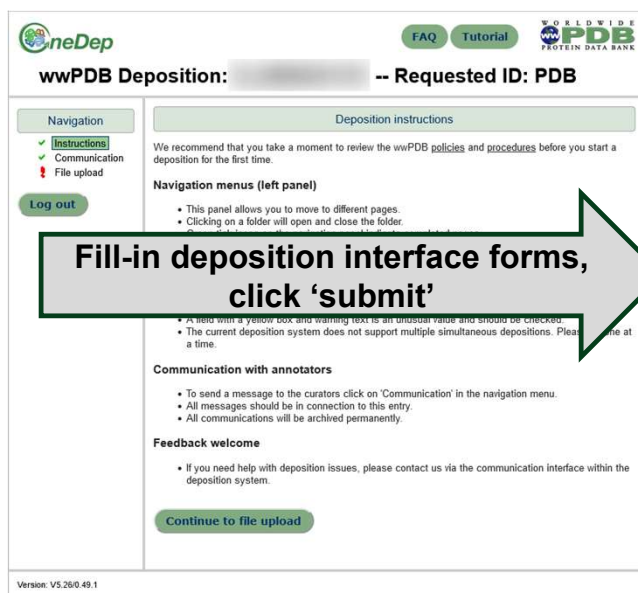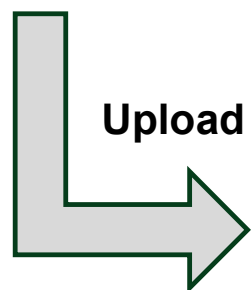Protein Data Bank in Europe

**EMBL-EBI**

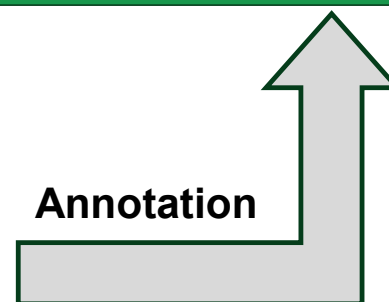# Submitting data to the PDB… is hard work

Model file =
Coordinates (mandatory)
+
Metadata (optional)

Experimental data file =
Experimental data (mandatory)
+
Metadata (tiny bit)

Experimental data file
+
Annotated model file =
Coordinates
+
Metadata ("full")

**Upload**



**Fill-in deposition interface forms, click 'submit'**

OneDep

**Annotation**

The richer in metatada in the uploaded file(s), the less forms to fill-in the deposition interface.

PDBe
Protein Data Bank in Europe

PDBe-KB
Protein Data Bank in Europe - Knowledge Base

AlphaFold
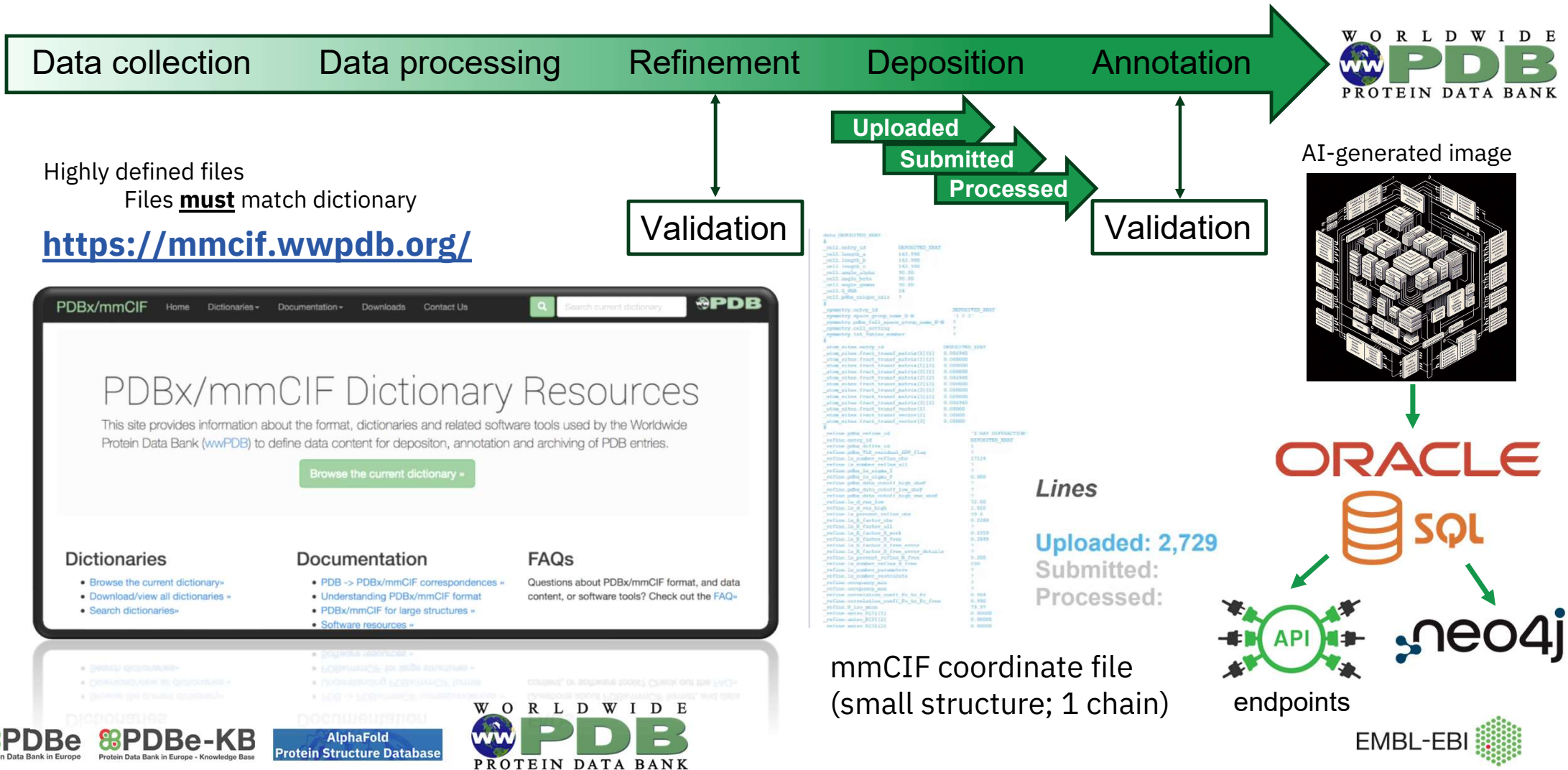Protein Structure Database

WORLDWIDE
wwPDB
PROTEIN DATA BANK

EMBL-EBI

# The road to retain metadata

**Each step must**
- generate its own metadata
- retain the previous step's metadata

Data collection | Data processing | Refinement | Deposition | Annotation

**Currently:**
Metadata stored in variety of forms
(heterogeneity)

Validation

Validation

**Goal:** the file outputted by refinement software is as complete as possible to minimize the need to enter metadata manually in the deposition interface.

# Metadata → Database



Data collection    Data processing    Refinement    Deposition    Annotation

**WORLDWIDE PDB PROTEIN DATA BANK**

Uploaded
Submitted
Processed

Validation

Validation

Highly defined files
Files **must** match dictionary

**https://mmcif.wwpdb.org/**

AI-generated image

Lines

Uploaded: 2,729
Submitted:
Processed:

mmCIF coordinate file
(small structure; 1 chain)

ORACLE
SQL

API

neo4j

endpoints

EMBL-EBI

Data collection — Data processing — Refinement — Deposition — Annotation

UK Research and Innovation

**DRI Project**

ISPyB/SynchWeb

SQL

diamond
**Diamond Light Source**

JSON

script

mmCIF-gen
(using gemmi)

https://pypi.org/project/mmcif-gen/

- Experimental method (X-ray, NMR, EM, EC…)
- Data collection date
- Wavelength
- Synchrotron name
- Beamline name
- Detector type
- Detector name
- Serial crystallography experiment Y/N

- Temperature (not yet in SynchWeb / DLS ISPyB)
- Multiple datasets
- **Oscillation angle per image**
- **Total oscillation angle**
- **Beam transmission**

metadata mmCIF file

GΦL Global Phasing Limited

CCP4

embedded in **mtz file**

PDBe · PDBe-KB · AlphaFold Protein Structure Database · WORLDWIDE PDB PROTEIN DATA BANK · EMBL-EBI

GΦL
Global Phasing Limited

- ❑ The idea is to **append** (additional) **metadata** in some standard format to the end of an MTZ file:

  - ▪ This way metadata travel together with the most popular format for reflection data (MTZ), avoiding risks of de-synchronisation and loss of information along the path from experiment to deposition.

  - ▪ That appendix does not impact on existing MTZ-reading programs as far as we know (tested: CCP4, Phenix, GPhL, Gemmi and other programs/packages using the CCP4 library)

- ❑ Things to consider:

  - ▪ How to append?

  - ▪ Define start and end of appendix.

  - ▪ Define content of appendix.

  - ▪ Associate metadata with rest of MTZ file (i.e. reflection data and existing MTZ header).

  - ▪ How to extract?

Data collection → Data processing → Refinement → Deposition → Annotation

IUCr Journals
CRYSTALLOGRAPHY JOURNALS ONLINE

home    submit    subscribe    open access

Table 3
Data collection and processing

Values given in parentheses are for the highest resolution shell.

Diffraction source
Wavelength (Å)
Temperature (K)
Detector
Crystal to detector distance (mm)
Total rotation range (°)
Exposure time per degree (s) or rotation per image (°)
Exposure time per image (s)
Space group
$a, b, c$ (Å)
$\alpha, \beta, \gamma$ (°)
Mosaicity (°)
Resolution range (Å)

Total no. of reflections
No. of unique reflections
Completeness (%)
Redundancy
$<I/\sigma(I)>$ from merged data
$CC_{1/2}$
$R_{r.i.m.}$ or $R_{meas}$
Overall $B$ factor from Wilson plot (Å²)

Example appendices extraction:

```
gemmi mtz -A truncate-unique.mtz
```

```
#MTZAPPENDIX-START\n
#MTZAPPENDIX-ITEM CIF DatasetID=1\n
…
#MTZAPPENDIX-END\n
```

autoPROC tool:
aP_mtz_appendix –h

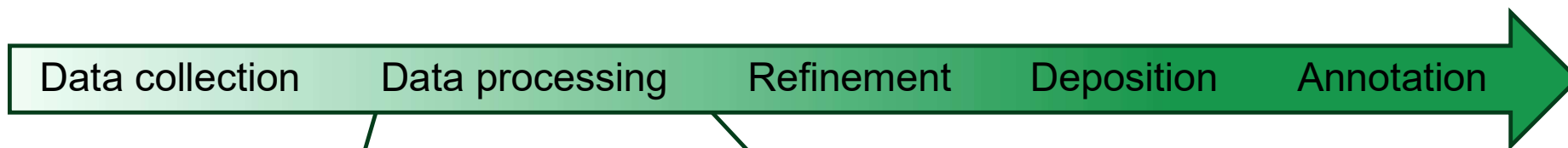https://colab.research.google.com/drive/1L0FhSWeqqp50LBtNJBK7LOmEP0Vz4YnR?usp=sharing

https://github.com/glevans/PDB_Notebooks/tree/main/GemmiRecipes

```
!gemmi mtz -A truncate-unique.mtz
_reflns.d_resolution_low 54.110
_reflns.d_resolution_high 1.517
_reflns.pdbx_Rmerge_I_obs 0.0535
_reflns.pdbx_Rrim_I_all 0.0547
_reflns.pdbx_Rpim_I_all 0.0112
_reflns.pdbx_number_measured_all 818462
_reflns.number_obs 40884
_reflns.pdbx_netI_over_sigmaI 23.41
_reflns.percent_possible_obs 100.0
_reflns.pdbx_redundancy 20.02
_reflns.pdbx_CC_half 0.999
_reflns.pdbx_percent_possible_anomalous 99.6
_reflns.pdbx_redundancy_anomalous 10.68
_reflns.pdbx_CC_half_anomalous -0.054
_reflns.pdbx_absDiff_over_sigma_anomalous 0.636
```

mtz from autoPROC

PDBe — Protein Data Bank in Europe
PDBe-KB — Protein Data Bank in Europe - Knowledge Base
AlphaFold Protein Structure Database
WORLDWIDE PDB PROTEIN DATA BANK
EMBL-EBI

**Data collection** | **Data processing** | **Refinement** | **Deposition** | **Annotation**

**Working together towards metadata capture:**

Metadata values & corresponding wwPDB capturing (mmCIF categories items)

https://github.com/glevans/automating_data_collection_stats (private)

- '_reflns' for Overall **'Collection Statistics'**
- '_reflns_shell' for Inner and Outer Shell **'Collection Statistics'**

These categories have items that correlate to each of the statistical values listed in the form. Below is useful table (Table 1) relating the mmCIF item labels in the '_reflns' category, to the mmCIF items in '_reflns_shell' category and their correlation to items in the wwPDB deposition interface **'Collection Statistics'** form values:

**Table 1**: Items in the '_reflns' category, '_reflns_shell' category and human-readable names for the values:

| | _reflns | _reflns_shell | More information |
|---|---|---|---|
| High resolution limit [Å] | d_resolution_high | d_res_high | The high resolution limit used for data processing, or for the shell. |
| Low resolution limit [Å] | d_resolution_low | d_res_low | The low resolution limit used for data processing, or for the shell. |
| Total number unique reflections | number_obs | number_unique_obs | The number of symmetry-unique observations. |
| Total number of reflections | pdbx_number_measured_all | number_measured_all | The number of measured intensities just before the final merging step. |

**Interested?**

**Email: gle@ebi.ac.uk
with GitHub handle**

**Overall Wilson B factor** is often NOT listed in logs (or mmCIF) with processing stats from autoprocessing pipelines.
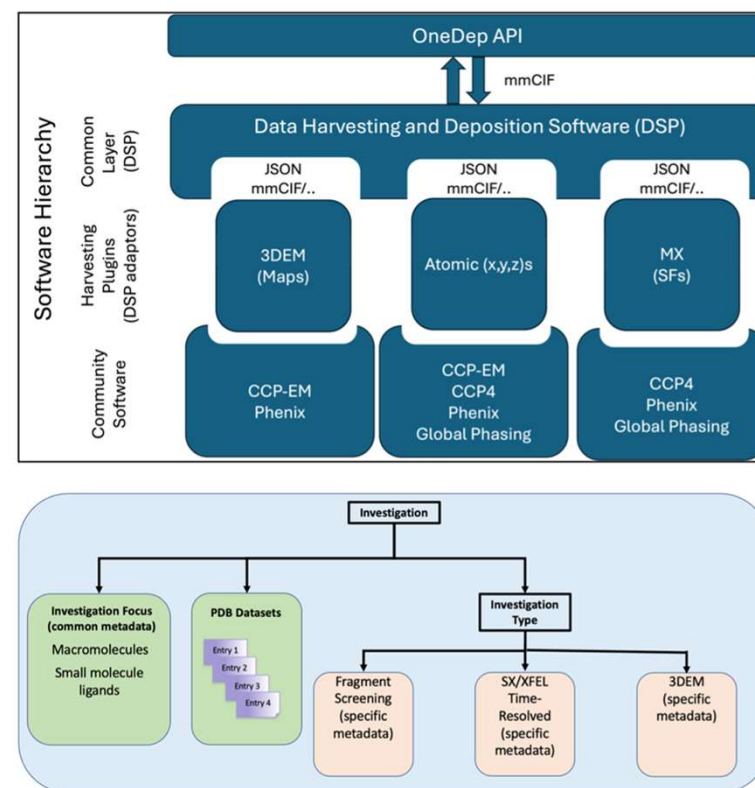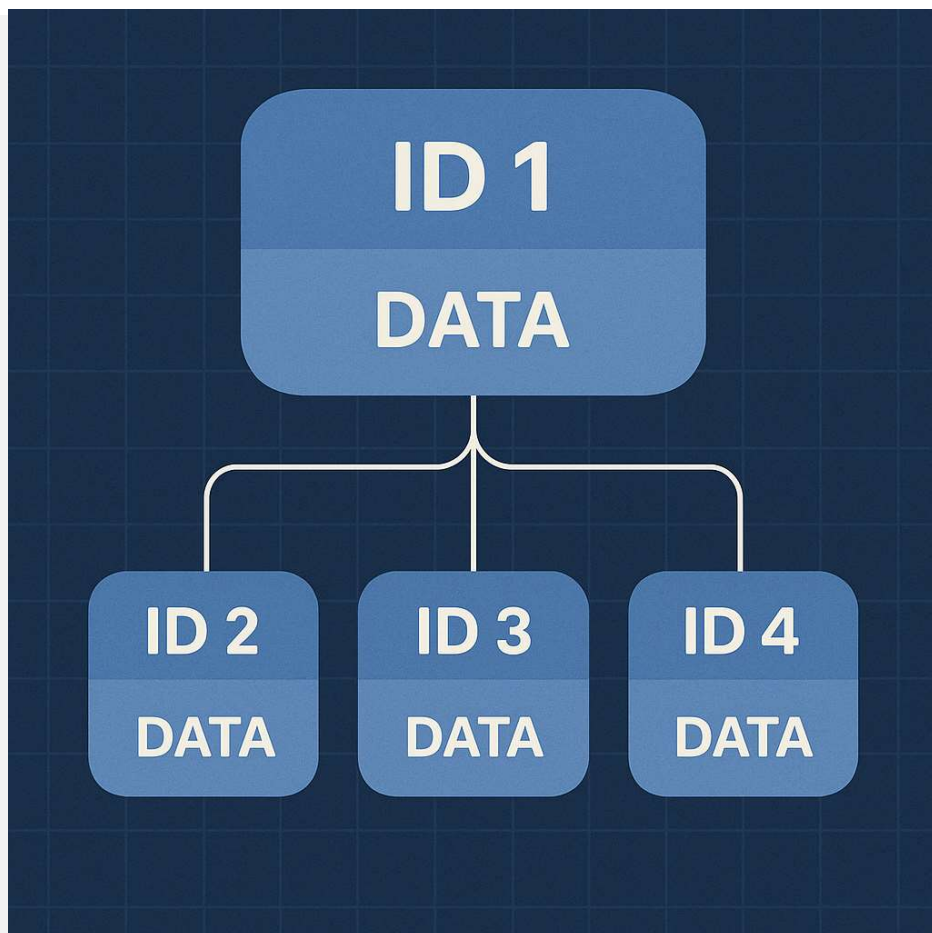
## Project awarded: Streamlining Data Deposition

3 years project: November 2025 to October 2028

- Make Deposition API available for depositors
- Improve metadata copy based on previous depositions
- New OneDep APIs to:
  - Access ligand curation & validation
  - Add sequence annotation
  - Add assembly annotation
  - Access wwPDB validation
- Enable batch depositions and processing
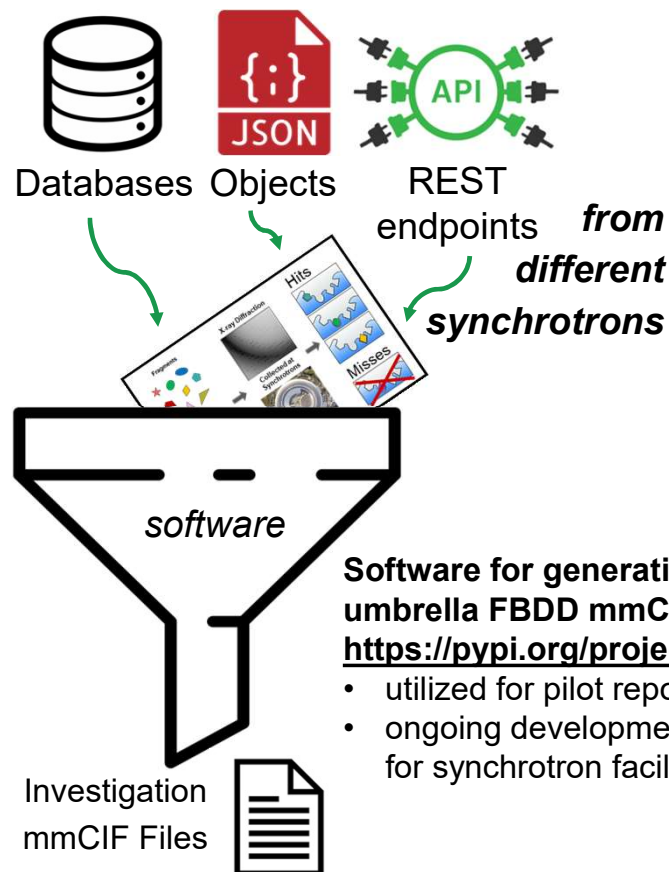
# Batch metadata → INVESTIGATION files



**Concept:**

❖ Umbrella file (ID 1) with defined data model to connect multiple entries (ID 2, ID 3 &ID 4).

**Key features:**

❖ **Value-adding, enriching metadata**

❖ Define types of umbrella files for different use case, *e.g.*:
- Fragment screening
- Time-resolved
- EMDB specific cases – composite maps

❖ Only contains a subset of information found in a coordinate mmCIF file (<u>not</u> coordinates)

❖ Links to non-wwPDB databases
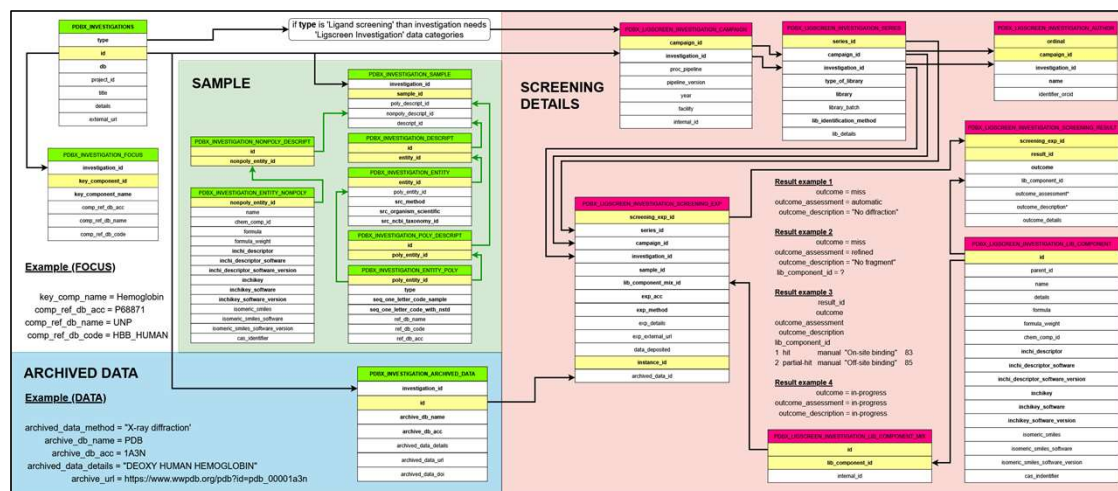
# X-ray-based Fragment-based drug discovery (FBDD) data

Databases  Objects  REST endpoints  *from different synchrotrons*

*software*

Investigation mmCIF Files

**Software for generating umbrella FBDD mmCIF file**
**https://pypi.org/project/mmcif-gen/**
- utilized for pilot repository
- ongoing development for synchrotron facilities

**Standardizing data schema for FBDD data**
**https://github.com/PDBeurope/InvestigationCIF**

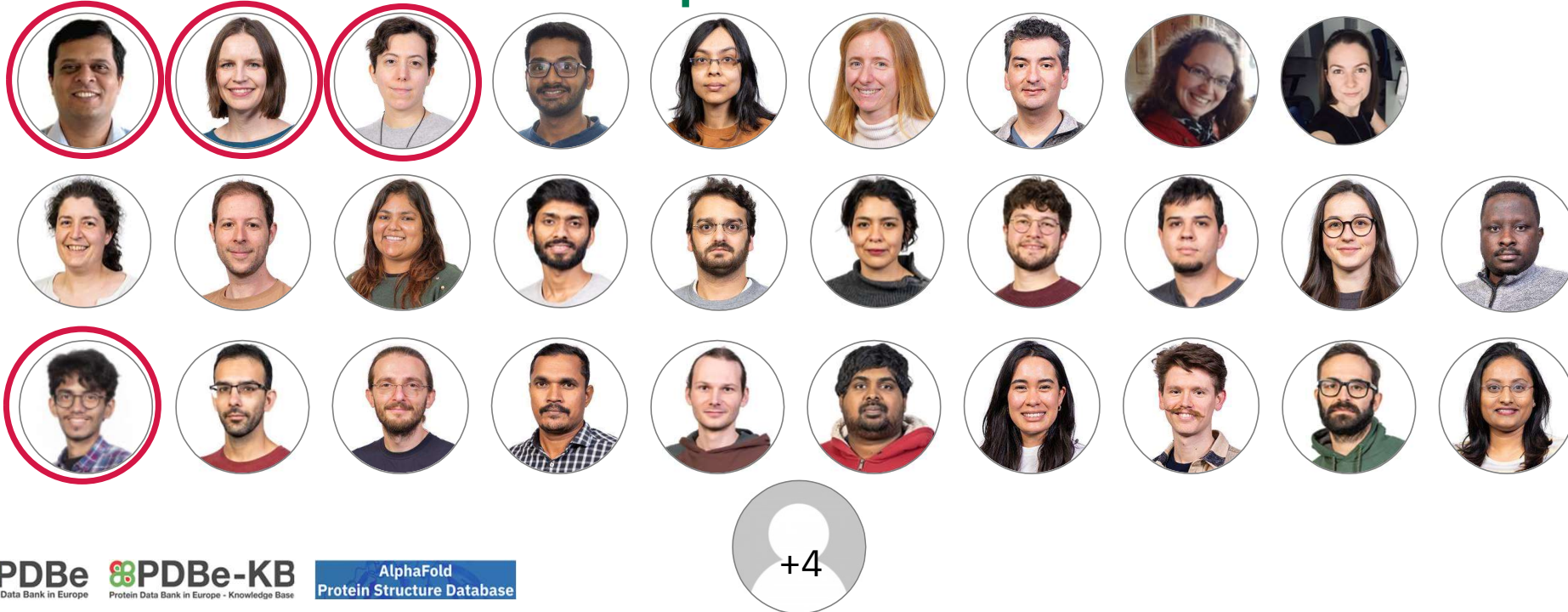Synchrotrons and associated facilities involved in developing this data model:
- **The Crystallisation Facility** at the **European Molecular Biology Laboratory (EMBL) Grenoble** and **European Synchrotron Radiation Facility (ESFR)** in France
- **XChem: Diamond Fragment Screening** at **Diamond Light Source (DLS)** in the United Kingdom
- **Fragment Screening Facility** at **Berlin synchrotron BESSY-MX and Helmholtz-Zentrum Berlin/HZB** in Germany
- **FragMAX** at Swedish synchrotron **MAX IV** in Sweden

# Acknowledgements

## Resources

## The Protein Databank in Europe team

# Acknowledgements

# Thank you, feedback welcome!

## DLS

- ❖ I03 Beamline
  - Dave Hall
  - Mark Williams
  - Neil Paterson
- ❖ ISPYB
  - Irakli Sikharulidze
- ❖ XChem / I04-1 Beamline
  - Frank von Delft
  - Daren Fearon
  - Warren Thompson
  - Jasmin Aschenbrenner

## Global Phasing Ltd

- ❖ Software developers
  - Gerard Bricogne
  - Clemens Vonrhein
  - Marcin Wojdyr

## CCP4

- ❖ Software developers
  - Eugene Krissinel

## MaxIV

- ❖ FragMax
  - Tobias Krojer

## EMBL Grenoble

- ❖ CRIMS at ESRF (Crystallographic Information Management System / CRIMS)
  - José A. Marquez
  - Raphaël Bourgeas

## BESSY

- ❖ MX Team
  - Manfred Weiss

**pdbhelp@ebi.ac.uk**

**PDBe**

PDBe — Protein Data Bank in Europe

EMBL-EBI