How can PDB submission be made more automatic

Sameer Velankar Protein Data Bank in Europe



wwpdb.org

Worldwide Protein Data Bank (wwPDB)

Sustain freely accessible, interoperating Core Archives of structure data and metadata for biological macromolecules as an enduring public good to promote basic and applied research and education across the sciences.

- Manage the wwPDB Core Archives as a public good according to the FAIR Principles.
- Provide expert deposition, validation, biocuration, and remediation services at no charge to Data Depositors worldwide.
- Ensure universal open access to public domain structural biology data with no limitations on usage.
- Develop and promote community-endorsed data standards for archiving and exchange of global structural biology data.

PDBx/mmCIF Working Group

- PDBx/mmCIF is the archival data standard for the repository
- wwPDB together with the PDBx/mmCIF Working Group of community experts and methods developers oversee the evolution of the standard
- Working Group ensures that the standard is well supported by key community software tools.
- PDB hosts community workshops and maintains mmcif.wwpdb.org serving PDBx/mmCIF data dictionaries, schema and software tools



PDBx/mmCIF Workshop Participants, July 2017

Open Access to single global archive Weekly release

- Regional sharing of data deposition efforts
 - RCSB PDB (US)
 - PDB Japan (PDBj)
 - PDB in Europe (PDBe)
- Partner-hosted websites offer complementary services and views of data



WORLDWIDE

Protein Data Bank Deposition

- Deposition
 - deposit.wwpdb.org



- Validation validate your structures before deposition
 - validate.wwpdb.org
- Validation API



 <u>https://www.wwpdb.org/validation/onedep-validation-</u> web-service-interface

Growth/Complexity

- 13,377 depositions in 2019
- Rapid growth in 3DEM
- More structures per Depositor
- Increasing data complexity
 - Increase in large complexes resolved by EM





Protein Data Bank Deposition



How does deposition work now?



Manual curation review within OneDep

What data is captured?

Deposition type	Database accession codes issued	Mandatory file uploads	Optional file uploads
X-ray and neutron crystallography	PDB	Atomic coordinates Structure factor data Unmerged intensity data	Ligand definition file or image Auxiliary files
Solution and solid-state NMR	PDB BMRB	Atomic coordinates Assigned chemical shifts Restraints used in refinement Auxiliary sequence file from AMBER	Spectral peak lists Ligand definition file or image Auxiliary files
Electron crystallography	PDB EMDB	Atomic coordinates Structure factor data or Mass density map volume	Ligand definition file or image Auxiliary files
3DEM (map and model)	PDB EMDB	Atomic coordinates Mass density map volume Entry image for public display (EMDB)	Any number of additional maps Any number of masks Two half maps Fourier shell correlation (FSC) curve Ligand definition file or image
3DEM (map only)	EMDB	Mass density map volume Entry image for public display	As above

What data is captured?

- Administrative information (e.g., author release instructions).
- Description of each distinct macromolecule present in the sample.
- Description of the experimental setup (e.g., sample preparation and data collection).
- Description of experimental data, refinement (e.g., crystallographic refinement statistics), and software used.
- Description and matching of ligands and modified polymer residues to the PDB Chemical Component Dictionary (CCD).
- Information on the quaternary structure and, whenever possible, experimental support for the biologically relevant assembly.

Deposition of raw data



VALIDATION - DEPOSITION - DICTIONARIES - DOCUMENTATION - TASK FORCES - FTP - STATISTICS - ABOUT -

T - WwPDB Foundation

Starting a deposition session

Can I deposit raw data to the PDB?

Depositors are encouraged to deposit raw data of their structure and then provide the corresponding DOI in the "Related Entries" user interface during structure deposition in wwPDB OneDep.

Depositors are strongly encouraged to deposit their raw data in a curated archive.

- X-ray diffraction
 - ProteinDiffraction.org https://proteindiffraction.org
 - SBGrid https://sbgrid.org
 - CXIDB http://www.cxidb.org
- Electron microscopy
 - EMPIAR https://www.ebi.ac.uk/pdbe/emdb/empiar/deposition
- NMR e.g., FIDs, in particular those relating to NOESY type spectra
 - please contact BMRB bmrbhelp@bmrb.wisc.edu

•	Related experimental data sets			
Please provide the DOI for a related raw data set (e.g. diffraction image data with a DOI, not related PDB entries or related citations)				
	Type of experimental data	DOI for related data set	DOI for metadata	
	~			

Your deposition has been submitted

Thank you for your deposition at the wwPDB.
The PDB code for this deposition is 6Y47.
If you have additional raw X-ray diffraction data, please deposit at one of the raw image archives and inform us of the DOI:
proteindiffraction.org
SBGrid
Zenodo
CXIDB

Group deposition

https://deposit-group.rcsb.rutgers.edu/groupdeposit/

- Provided support for D3R Blind Challenges
- Early Adopters: Roche, Merck Serono, U. Marburg, U. Essex
- Requires a complete PDBx/mmCIF file with all mandatory data
- Each group needs to implement the necessary infrastructure to generate complete PDBx/mmCIF file



Mandatory submission of PDBx/mmCIF – June '19

- Increase in the amount of meta data collected in the entries during deposition
 - Work with Software community to develop automated process for deposition
- Improved consistency in biocuration between related entries
- Increased biocuration efficiency
 - Curating an increasing number of entries with the same resources



Received 21 February 2019 Accepted 3 April 2019

Edited by R. J. Read, University of Cambridge, England

Keywords: PDB; mmClF; OneDep; wwPDB; data dictionary: data archiving: biocuration; validation; macromolecular crystallography; data standards; PDBs/mmClF format; Protein Data Bank; Worldwide Protein Data Bank. Announcing mandatory submission of PDBx/mmCIF format files for crystallographic depositions to the Protein Data Bank (PDB)

Paul D. Adams,^{a,b} Pavel V. Afonine,^a Kumaran Baskaran,^c Helen M. Berrana,^d John Berrisford,^a Gerard Bricogne,^l David G. Brown,^a Stephen K. Burley,^{d.M.a} Minyu Chen,¹ Zukang Teng,² Claus Elensburg,^l Aleksandras Gutmanas,² Jeffrey C. Hoch,¹a Yasuyo Ikegawa,¹ Yumiko Kengaku,¹ Eugene Krissinel,¹ Genji Kurisu,¹ Yuhe Liang,^d Dorothee Liebschner,² Lora Mak,² John L. Markley,² Nigel W. Moriarty,² Garib N. Murshudov,^m Martin Noble,ⁿ Ezra Peisach,⁴ Irina Persikova,^d Billy K. Poon,^a Oleg V. Sobolev,^a Eldon L. Ulrich,² Samer Velankar,⁴ Clemens Vonrhein,¹ John Westbrock,⁴ Marcin Wojdyr,¹ Masashi Yokochi⁴ and Jasmine V. Young⁴

"Molecular Biophysics and Integrated Bioimaging Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA. ^bDenartment of Bioengineering, University of California, Berkeley, CA 94720, USA. ^cBioMagResBank (BMRB) University of Wisconsin-Madison, Madison, WI 53706, USA, dResearch Collaboratory for Structural Bioinformatics Protein Data Bank (RCSB PDB). Institute for Quantitative Biomedicine. Rutrers. The State University of New Jersey. Piscataway, NL08854, LISA, "Protein Data Bank in Europe (PDRe), European Molecular Biology Laboratory, European Bioinformatics Institute (EMBI-EBI). Wellcome Genome Campus. Hinxton. Cambridgeshire CB10.1SD. UK. (Global Phasing Limited, Sheraton House, Castle Park, Cambridge, CB3 0AX, UK, School of Biosciences, University of Kent, Canterbury, Kent CT2 7NL UK, ^bRuteers Cancer Institute of New Jersey, Robert Wood Johnson Medical School, New Brunswick, NI 08903, USA. Research Collaboratory for Structural Bioinformatics Protein Data Bank (RCSB PDB). Sar Diego Supercomputer Center, University of California, San Diego, La Jolla, CA 92093, USA, Protein Data Bank Japan (PDBi) Institute for Protein Research, Osaka University, Osaka 565-0871, Japan ⁸BioMaeResBank (BMRB), UConn Health, 263 Farmington Avenue, Farmington, CT 06030, USA, ¹CCP4, Research Complex at Harwell (RCaH), Rutherford Appleton Laboratory, Didcot, Oxon OX11 0FA, UK, "MRC Laboratory of Molecular Biology, Francis Crick Avenue, Cambridge Biomedical Campus: Cambridge CB2 0OH LIK and [®]Newcastle University. Framlington Place: Newcastle Upon Tyne, NE2 4HH, UK. *Correspondence e-mail: stephen.burley@rcsb.org, hoch@uchc.edu, gkurisu@protein.osaka-u.ac.jp, markley@biochem.wisc.edu, sameer@ebi.ac.uk

The Protein Data Bank (PDB) (wwPDB consortium, 2019) is the single global archive of experimentally determined three-dimensional (30) structure data of biological macromolecules. The continuing growth in the numbers, size and complexity of macromolecular structures in the PDB archive, coupled with the rapid growth of evolving experimental methods such as 3D cryo-electron microscopy (3DEM) has made the traditional PDB format (legacy PDB format) inadequate for fully representing these data. As described below, this format was based on a punched-card format that became obsolete long ago. In the following letter, we describe the changes necessary to address the challenges coming from the extraordinary success of structural biologists.

Since 2003, the PDB has been managed by the Worldwide Protein Data Bank (wwPDB; https://www.wwpdb.org/) (Berman *et al.*, 2003), an international partnership that collaboratively oversees deposition, validation, biocuration and open-access dissemination of 3D macromolecular structure data, adhering to the FAIR principles of Findability, Accessibility, Interoperability and Reusability (Wilkinson *et al.*, 2016). In 2007, the master file format for the archive was officially changed to PDB Exchange/ Macromolecular Crystallographic Information File (PDBs/mmCIF), supported by the PDBs/mmCIF data dictionary, to address new challenges in structure archiving. Later, in 2012, the wwPDB terminated its support of the legacy PDB file format and froze its threthe development (https://wwplo.teg/documentation/file-formats-and-the-pdb).

We now announce that as of 1 July 2019, PDBs/mmCTF will be the only format allowed for deposition of the atomic coordinates for PDB structures resulting from macromolecular crystallography (MX), including X-ray, neutron, fiber and electron diffraction methods, via OneDep (Voung et al., 2017). This requirement will be extended to PDB structures resulting from nuclear magnetic resonance (NRR) spectroscopy and 3DEM methods at a later date to be determined. Elimination of the legacy PDB format will improve the difficiency of the deposition process and enhance validation through capture



© 2019 International Union of Crystallography

Deposition of associated experimental data

- Experience of API's within OneDep
- Interaction with SASBDB during deposition
- API services to interact with federated deposition systems





SFX/XFEL are Revolutionizing MX

- SFX: Serial Femtosecond X-ray crystallography
- XFEL: X-ray Free Electron Laser
- Time-resolved studies
 - Conformational changes in enzyme active sites
 - Response to photoexcitation
 - Membrane protein dynamics
- Rapid advances in data collection and sample handling technologies
- Linking of multiple depositions ("investigations")









SACI

EΞ

Automated deposition

- Improved fidelity and completeness of 3D structure data deposited into the PDB
- Streamlining the wwPDB data deposition, validation, and biocuration system.
- Provide easy mechanism for deposition of multiple structures (SF/XFEL, Fragment screening)
- Improved efficiency for depositors and biocuration
- Consistent biocuration

How can PDB deposition be made more automatic



APIs to expose wwPDB deposition/annotation pipeline

New three-year project (PDBe and RCSB PDB)

- Funded by BBSRC and NSF
- Start date July 2020 (RCSB PDB); Jan 2021 (PDBe)
- Gather requirements on Deposition API requirements
 - Automated deposition
 - Multi-structure deposition/annotation
 - "investigation" based deposition
 - "investigation" based annotation Annotators to provide requirements

New three-year project (PDBe and RCSB PDB)

- Data Content Extension:
 - "Investigation" level semantics
 - Integrate validation data in mmCIF
- New Deposition Data Preparation APIs:
 - pre-deposition registration service for chemical reference data
 - Better validation of ligands before and during deposition

Questions and discussion