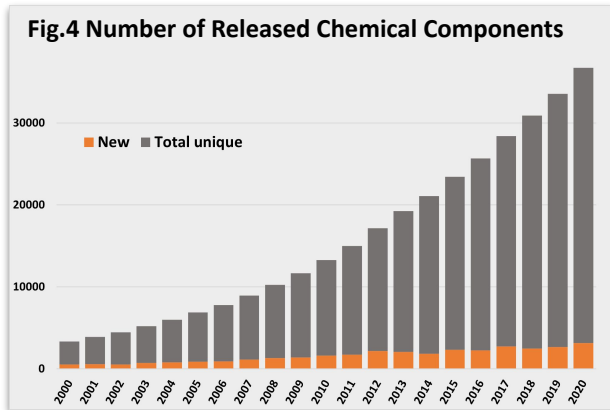
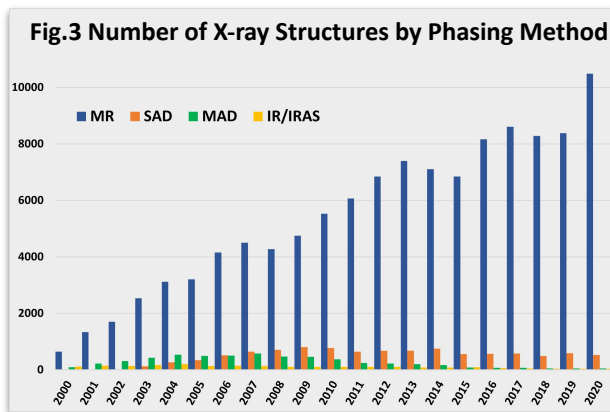
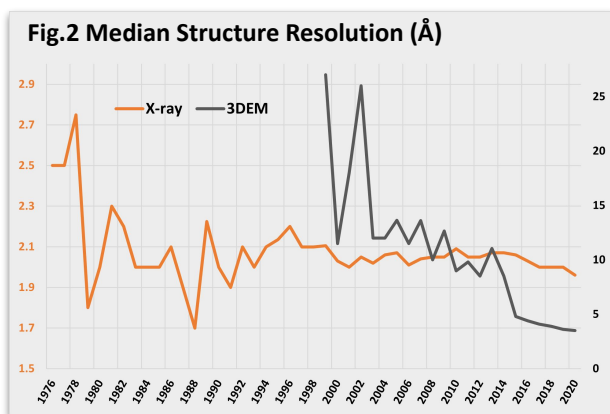
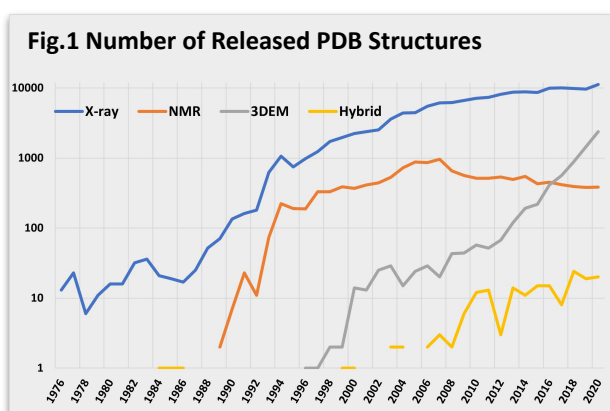


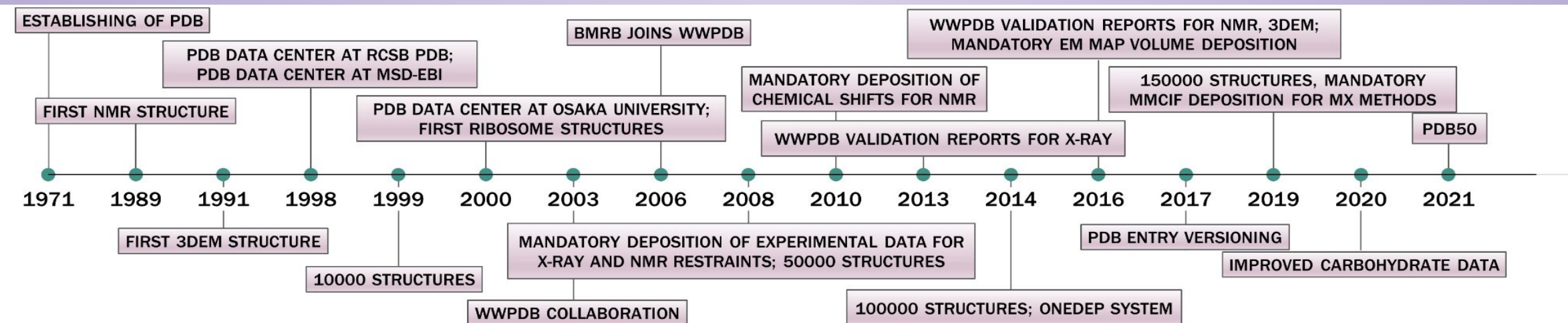
Trends in macromolecular structure data across 50 years of the PDB

Irina Persikova¹, David Armstrong², John Berrisford², Minyu Chen³, Gregg Crichlow¹, Genevieve Evans², Justin Flatt¹, Romana Gaborova², Sutapa Ghosh¹, Deepti Gupta², Deborah Harrus², Brian P. Hudson¹, Reiko Igarashi³, Yumiko Kengaku³, Ju Yaen Kim³, Yuhe Liang¹, Ezra Peisach¹, Osman Salih³, Junko Sato³, Monica Sekharan¹, Chenghua Shao¹, James Tolchard², Jack Turner², Jasmine Young¹

¹RCSB Protein Data Bank, Rutgers, The State University of New Jersey, Piscataway, USA ²Protein Data Bank in Europe, EMBL-European Bioinformatics Institute, Hinxton, UK ³Protein Data Bank Japan, Institute for Protein Research, Osaka University, Japan



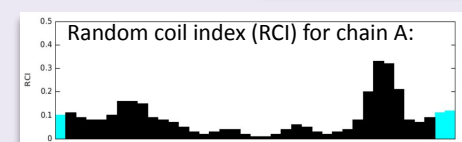
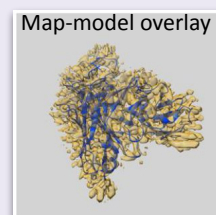
In 2021, we celebrate 50 years of the Protein Data Bank (PDB) archive - one of the longest-running open access scientific databases - managed collaboratively by the wwPDB ([RCSB PDB](#), [PDBe](#), [PDBj](#), and [BMRB](#)). Throughout these 50 years, there has been significant evolution of the data in the PDB archive. This evolution has been driven by a number of factors including development in structural determination techniques, adaptation of biocuration practices, and increase in data capture via updated file formats. Here we present some of the key trends in data across the PDB archive, highlighting how structural biology data has changed over time and how wwPDB biocuration practices have adapted to handle these changes.



Supporting multiple experimental techniques

Experimental method specific:

- Deposition interfaces
- Dictionaries (NMR: 244; EM: 613; X-FEL: 49 dictionary definitions)
- Extended data - anisotropic, unmerged
- Enumerations
- Validation (X-ray, NMR, EM)



Diverse sequences/sources

- Sequence cross-reference with UniProt /GB
- Curation of sequence discrepancies
- Taxonomy ids since 2007 (NCBI)
- Taxonomy enumerations at deposition

ALIGNMENT POSITION	AUTH ENTITY-1	ALIGNMENT POSITION	AUTH ENTITY-2	RESIDUE	ANNOTATION
77	TRP	UNP:A3LT82	77	TRP	engineered mutation
112	CSD	UNP:A3LT82	112	TRP	engineered mutation

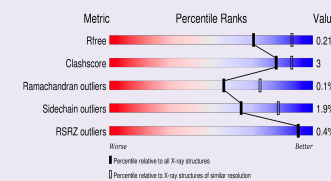
AUTH Entry: 1 V:3
XYZ Chain: A V:1
UNP:A3LT82

TLILTEATFVSPQASGEGAAE
TLILTEATFVSPQASGEGAAE
TLILTEATFVSPQASGEGAAE
61

Data quality and reproducibility

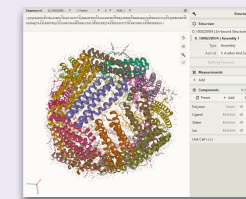
PDF validation reports since 2013:

- Required by the journals
- Overall quality at a glance
- Geometry and chemistry reports
- Model vs experimental data validation



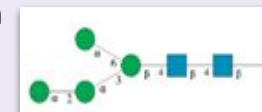
Large structures and assembly annotation

- Large structures are handled as single files at deposition, curation and distribution since 2014
- Biological assembly is defined for each entry
- PISA is used to assess stable interfaces
- Annotated assembly verification in 3D
- Assembly files are available for download



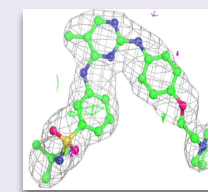
Archival data consistency and remediation

- 2007 - data uniformity across the archive
- 2008 - data enhancement, chemistry, assemblies, binding sites
- 2011 - biological assemblies, residual B factors, peptide inhibitors, nonstandard crystal frames
- 2013 - remediated structure factor files
- 2014 - integration of large structures
- 2016 - remediation of 3DEM entries
- 2017 - PDBx/mmCIF files to support OneDep
- 2020 - Carbohydrate remediation



Small molecules

- Chemical component definitions: stereochemical assignments, chemical descriptors (SMILES & InChI), systematic chemical names, idealized coordinates
- Validation reports: geometrical quality assessment, 2D graphical views of electron density fit



Supporting data uniformity, biocuration efficiency

- OneDep common deposition portal since 2014
- Curation efficiency and consistency (15436 structures curated in 2020)
- Provides enumerations and value limits
- Sequence, taxonomy, chemistry, model vs data, assembly annotation and validation



Looking towards the future

- Group depositions: multiple structures in one session
- PDB and CCD code extension
- Deposition APIs: Data harvesting pipeline from structure determination software packages
- PTM remediation
- Assembly curation improvements
- Validation improvements based on VTF recommendations

