

Harvesting Data Collection Information for the RCSB PDB Depositions

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Harvesting Data Collection Information for RCSB PDB Depositions

John Westbrook
Rutgers, The State University of New Jersey

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Data Management at the RCSB PDB

- Dictionary-driven architecture
- Pipeline of data collection spanning experimental structural biology
- Tools to automate data harvesting
- Continuous focus of data standardization and data quality
- Distribution and delivery deeply integrated with biological and medical resources

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Dictionary Resources

<http://mmcif.pdb.org/>

An Information Portal to Biological Macromolecular Structures

Dictionary Resources

The Protein Data Bank (PDB) uses macromolecular Crystallographic Information File (mmcif) data dictionaries to describe the information content of PDB entries. The PDB Exchange data dictionary consolidates content from a variety of crystallographic dictionaries including the IUCr Core, mmcif, Image and Symmetry dictionaries. The PDB Exchange Dictionary also includes extensions describing NMR, Cryo-EM, and protein production data. PDB data processing, data exchange, annotation, and database management operations all make heavy use of the data format and the content of the PDB Exchange Dictionary. Software tools are used to convert mmcif data files to the older PDB format and to PDBML/XML.

- Data files in mmcif format can be downloaded from the RCSB PDB website or by ftp.
- Software tools are available for processing and editing dictionaries.
- Software tools are available for converting mmcif data files to PDB and PDBML formats.
- A complete list of PDB software tools for managing PDB data in mmcif format can be found [here](#).

Dictionary Content and Representation

- Background and Introduction about mmcif
- The Macromolecular Crystallographic Information File (mmcif): *Math. Enzymol.* (1997) 277, 571-590.
- *Elucidation of an Extensive Challenge for Macromolecular Structure and beyond (PDB)*, *Bioinformatics* (2000) 16(2), 159-166.
- mmcif Software Developers Workshop 1997
- mmcif Dictionary Templates
- mmcif Examples
- References

Data Dictionaries

- PDB mmcif Exchange Dictionary (ASCII) | (compressed) | (HTML) | XML Schema
Data dictionary developed as a collaboration between MSD-EBI and RCSB and used by wwPDB members for data exchange.
- mmcif Dictionary (ASCII) | (compressed) | (HTML)
IUCr Standard Dictionary
Original working group members: Paula M. Fitzgerald, Helen Berman, Phil Bourne, Brian McMahon, Keith Wilson, and John Westbrook

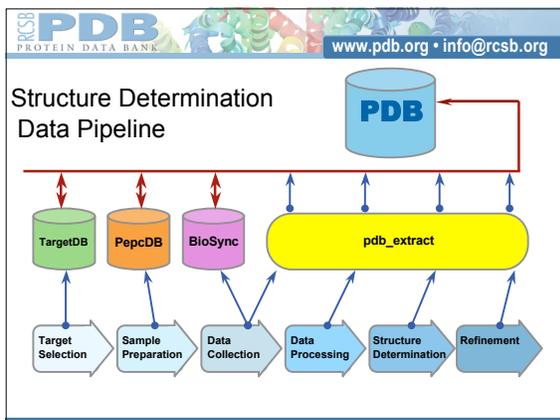
RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

PDBML Schemas

<http://pdbml.pdb.org/>

- New PDBML website
- Schemas provided for the PDB Exchange dictionary and component dictionaries
- Schemas and data dictionaries are updated synchronously

PDBML: the representation of archival macromolecular structure data in XML. John Westbrook, Nobutoshi Ito, Haruki Nakamura, Kim Henrick and Helen M. Berman, *Bioinformatics*, 21(7), 988-992, 2005.



RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Target Registration Database

TargetDB • <http://targetdb.pdb.org/>

- Targets downloadable in XML (~120K targets)
- Targets downloaded from 21 worldwide centers weekly
- Target search by:
 - Sequence (FASTA), project target ID, project site, status (selected, cloned, expressed, ... in PDB), update date, protein name, source organism
- Report output in HTML, FASTA, and XML
- Integrates PDB entry sequences (~90K sequences)
- Includes PDB pre-release sequence data
- Provides links to related project sequence databases
- Summary reports of target and project progress
- 3-Tier architecture (Apache/Tomcat/MySQL)

Harvesting Data Collection Information for the RCSB PDB Depositions

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

PDB EXTRACT  **CCP4**

pdb_extract
<http://pdb-extract.pdb.org/>
<http://sw-tools.pdb.org/>

- Data collection and reduction
 - HKL, SCALEPACK, d'TREK, SAINT, SCALA
- Molecular replacement
 - CNS, CNX, Amore, Morep, EPMR
- Heavy atom phasing
 - CNS, CNX, SOLVE, MLPHARE, SHARP/autoSHARP, SHELXD, SHELX, PHASES, SnB, BnP, Phaser
- Density modification
 - CNS, CNX, DM, Solomon, RESOLVE, SHELXE
- Structure refinement
 - CNS, CNX, REFMAC, RESTRAIN, SHELXL, TNT, WARP, PHENIX

NMR - CNS, CYANNA, DYANNA

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Current Data Item Coverage

Data Item	PDB_EXTRACT		ALL	
	# pop	# total %	# pop	# total %

_reflns.number_all	1049	2239 46.8	16540	37906 43.6
_reflns.number_obs	2229	2239 99.5	34154	37906 90.1
_reflns.pdbx_Rsym_value	693	2239 30.9	11225	37906 29.6
_reflns.percent_possible_obs	2166	2239 96.7	32683	37906 86.2
_reflns.pdbx_redundancy	1694	2239 75.6	26547	37906 70.0
_reflns.pdbx_Rmerge_I_obs	1848	2239 82.5	23562	37906 62.1
_reflns.pdbx_Rsym_value	693	2239 30.9	11225	37906 29.6
_reflns.pdbx_netI_over_av_sigmaI	1782	2239 79.5	24282	37906 64.0

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Current Data Item Coverage

Data Item	PDB_EXTRACT		ALL	
	# pop	# total %	# pop	# total %

_reflns_shell.meanI_over_sigI_obs	1348	2239 60.2	19285	37906 50.8
_reflns_shell.number_measured_all	265	2239 11.8	283	37906 0.7
_reflns_shell.number_measured_obs	450	2239 20.1	490	37906 1.2
_reflns_shell.number_unique_all	1136	2239 50.7	9646	37906 25.4
_reflns_shell.number_unique_obs	286	2239 12.7	296	37906 0.7
_reflns_shell.percent_possible_all	1893	2239 84.5	28230	37906 74.4
_reflns_shell.percent_possible_obs	553	2239 24.7	700	37906 1.8
_reflns_shell.Rmerge_I_obs	1812	2239 80.9	19318	37906 50.9
_reflns_shell.pdbx_redundancy	1506	2239 67.2	17988	37906 47.4
_reflns_shell.pdbx_Rsym_value	654	2239 29.2	9664	37906 25.4

_exptl_crystal.density_Matthews	2201	2239 98.3	19680	37906 51.9
_exptl_crystal.density_percent_sol	2200	2239 98.2	36750	37906 96.9

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Current Data Item Coverage

Data Item	PDB_EXTRACT		ALL	
	# pop	# total %	# pop	# total %

_diffn.ambient_temp	2139	2239 95.5	32538	37906 85.8
_diffn.detector.details	680	2239 30.3	10345	37906 27.2
_diffn.detector.detector	2142	2239 95.6	33018	37906 87.1
_diffn.detector.type	2112	2239 94.3	32289	37906 85.1
_diffn.detector.pdbx_collection_date	2187	2239 97.6	34362	37906 90.6
_diffn.radiation.monochromator	1206	2239 53.8	16397	37906 43.2
_diffn.radiation.rcsb_diffn_protocol	2239	2239 100.0	33669	37906 88.8
_diffn.radiation.wavelength.wavelength	2164	2239 96.6	33233	37906 87.6
_diffn.reflns.number	88	2239 3.9	103	37906 0.2
_diffn.reflns.av_sigmaI_over_netI	11	2239 0.4	14	37906 0.0
_diffn.source.type	2218	2239 99.0	32563	37906 85.9
_diffn.source.pdbx_synchrotron_beamline	1835	2239 81.9	22849	37906 60.2
_diffn.source.pdbx_synchrotron_site	1835	2239 81.9	23005	37906 60.6
_diffn.source.source	2228	2239 99.5	32333	37906 85.3

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

pdb_extract - Online

A MEMBER OF THE RCSB PDB
An Information Portal to Biological Macromolecular Structures

pdb_extract

The **pdb_extract** online tool extracts key information from X-ray crystallographic and NMR software applications in preparation for PDB deposition to:

- Reduce the human effort required to assemble complete and validated protein structure entries ready for PDB deposition.
- Prepare an mmCIF file for deposition quickly.

HOW TO RUN:

- Select experimental method (X-Ray or NMR).
- Provide your workstation for the name of the coordinate file obtained from final structure refinement.
- Select the file type and program name.
- Press the RUN button to start this operation.

X-Ray NMR
 Select Program for Structure Refinement: /
 Coordinate File: File type:

NOTE:

- If the file size is large, it is recommended to upload gzipped (*.gz) or compressed (*.zip) file for faster loading.
- After assembling your data, you will get two mmCIF files for X-Ray or one mmCIF file for NMR. Please load them to ADIT for a complete deposition.

Questions, comments, and suggestions? [Contact Us](#)

RCSB PDB PROTEIN DATA BANK www.pdb.org • info@rcsb.org

Deposition Services

A MEMBER OF THE RCSB PDB
An Information Portal to Biological Macromolecular Structures

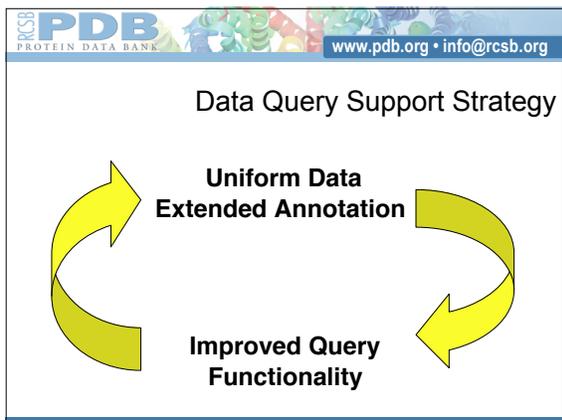
RCSB PDB Data Validation and Deposition Services

- NEW** Beta-ADIT at RCSB | ADIT at RCSB or PDB | [pdb_extract](#)
- NEW** Beta-ADIT NMR | Validation Server at RCSB or PDB | [Ligand Depot](#)

The following **data deposition** tools and instructions can make your structure deposition quick, easy, complete and accurate:

- Prepare** your structural data for deposition using **pdb_extract** and/or the desktop version of **ADIT**
- Validate** your structure using the **Validation Server** at **RCSB PDB** or the **Validation Server** at **PDB**
- Deposit** your structure using the structure deposition tool **beta-ADIT** at **RCSB PDB** or **ADIT** at **RCSB PDB** or **ADIT** at **PDB**, or using **AutoDep** at **MSD-EBI**:
 - The **RCSB PDB**, **PDB**, and **MSD-EBI** are members of the **wwPDB**.
 - Instructions for X-ray crystallography structure depositions.
 - Instructions for NMR structure depositions.
 - Instructions for EM structure depositions.
 - More information on deposition of structures determined by other methods (including Electron diffraction, Fiber-diffraction, Theoretical modeling).
- Search** for
 - Your ligand using **Ligand Depot**
 - Appropriate sequence database references for proteins or nucleic acids in your structure (e.g. using **BLAST**)

Harvesting Data Collection Information for the RCSB PDB Depositions



wwPDB Remediation Project

The wwPDB has collaborated to remediate the PDB archive and create a more consistent set of files.

E-MSD is supported by grants from the Wellcome Trust, the EU (TEMBLOR, NMRQUAL and IIMS), CCP4, the BBSRC, the MRC and EMBL.

PDBj is supported by grant-in-aid from the Institute for Bioinformatics Research and Development, Japan Science and Technology Agency (BIRD-JST), and the Ministry of Education, Culture, Sports, Science and Technology (MEXT).

The BMRB is supported by NIH grant LM05799 from the National Library of Medicine.

The RCSB PDB is supported by grants from the National Science Foundation, National Institute of General Medical Sciences, the Office of Science-Department of Energy, the National Library of Medicine, the National Cancer Institute, the National Center for Research Resources, the National Institute of Biomedical Imaging and Bioengineering, the National Institute of Neurological Disorders and Stroke, and the National Institute of Diabetes & Digestive & Kidney Diseases.

Data Remediation
<http://remediation.wwpdb.org>

- Sequence and taxonomy
 - Resolved anomalies relative to UniProt (61K sequences)
 - Resolved anomalies between chemical and macromolecular sequence
- Atom nomenclature
 - Conforms to IUPAC for standard amino acids and nucleotides
- Ligands and monomers
 - New dictionary with full chemical description (7800+ definitions)
 - All files annotated relative to this dictionary (163K non-polymers + 101K polymer residues)
- Biological assembly defined
 - Viruses now uniformly annotated (250 entries)

Nomenclature Standardization

- IUPAC-compliant H-atom names for standard amino acids and nucleotides (exceptions: OXT and HXT)
- DNA and RNA now differentiated (Adenine is DA for DNA; A for RNA)
- Modified nucleotides expressed as 3-letter codes
 - (replacing +T, +C ... etc.)
- PDB asterisks replaced by single quotes in atom labels (O2* becomes O2')
- Unusual names made more conventional (AC8 becomes C8A)
- White-space characters removed (N 2 becomes N2)

Chemical Dictionary

- Stereochemical assignments
- Aromatic bond assignments
- Nomenclature
 - IUPAC nomenclature for standard amino acids and nucleotides (Pure & Applied Chem., 70, 117-142, 1998)
 - All atom labels begin with element type symbol
 - Retention of all prior names as an alternate identifier
- Model and idealized coordinates
- Chemical descriptors (SMILES & InChI)
- Systematic chemical names
- Redundant chemical components obsoleted
- Additional definitions for protonated forms

Protonation Variants

- Companion dictionary of 220 definitions
- Additional definitions provide nomenclature for protonation states of:
 - N, C-terminal and free amino f
 - common side chain variants
- IUPAC compliant names
- Covers BMRB & CCPN variants

The diagram shows three chemical structures of a histidine side chain. The first shows the neutral form with a neutral imidazole ring. The second shows the protonated form with a positive charge on one of the nitrogen atoms of the imidazole ring. The third shows a different protonation state or tautomer.

Harvesting Data Collection Information for the RCSB PDB Depositions

Worldwide Protein Data Bank
www.wwpdb.org

Virus Structures

Worldwide Protein Data Bank
www.wwpdb.org

<http://remediation.wwpdb.org>
for details and dictionaries

Remediation Test Site

Component definition ATP

wwPDB Remediation Test Site

The full remediated archive is now available at <http://remediation.wwpdb.org> for testing. This site will production site at <http://www.rcsb.org>. Both sites have the same organization structure. The access code <http://www.wwpdb.org/remediation-downloads.html>.

The Chemical Component Dictionary and smaller test sets of data files are also available for testing.

PDB users are encouraged to test the remediated data files between April and July 2007. As of July, 07 site containing remediated data files. The final transition date will be announced on this website.

Comments about the files should be sent to info@wwpdb.org. Major announcements will be made at 5 websites.

About the wwPDB Remediation Project

The evolution of experimental methods, functional knowledge of proteins, and methods used to process archive. The wwPDB has remediated these data to create a more uniform archive.

The results of these efforts can be found in the remediated coordinate files. Highlights include:

Sequence	Updated references to databases and taxonomies	Chemical features
Classon	Resolved differences between chemical and macromolecular sequences	Atom count
Assembly and virus information	Verified and updated primary citation assignments	Disulfide count
Nuclear acid labeling	Improved representation of deposited and experimental coordinate frame	Chiral atoms
Beamline data	Density and ribose nucleotides assigned separate chemical definitions. The RNA forms remain labeled as A, G, C, U, I.	Head count
Atom nomenclature	Standardized to reflect changes in Chemical Component Dictionary (see below)	Aromatic bond count

Worldwide Protein Data Bank
www.wwpdb.org

Testing and Roll-out

- Remediated data files and chemical dictionary previewed by 60 software developers and database maintainers (November 2006)
- Example files files and chemical components dictionary released for public review (January 2007)
- Full remediated archive released for public review in April 2007
- Planned review period continues through July 2007
- Further details of the review updated at <http://www.pdb.org>, <http://remediation.wwpdb.org> and wwPDB sites.

Images of remediated PDB entries 1tk0, 2bvv & 407d obtained using Jmol, Chimera & OpenRasmol

RCSB PDB
www.pdb.org • info@rcsb.org

Access

RCSB PDB

- <http://www.pdb.org/>

wwPDB

- <http://www.wwpdb.org/>

wwPDB Remediation

- <http://remediation.wwpdb.org/>

Dictionary & Schema Resources

- <http://mmcif.pdb.org/> & <http://pdml.pdb.org>

TargetDB & PepcDB

- <http://targetdb.pdb.org> & <http://pepcdb.pdb.org>

RCSB Software Download Site

- <http://sw-tools.pdb.org>
- CVS server rcsb-cvs.rcsb.org - anonymous access

RCSB PDB
www.pdb.org • info@rcsb.org

The RCSB PDB Team

RCSB PDB
www.pdb.org • info@rcsb.org

Acknowledgements

Operated by two members of the RCSB:

The RCSB PDB is a member of the

Supported by: