# RCSB Protein Data Bank: A Resource for Chemical, Biochemical, and Structural Explorations of Large and Small Biomolecules

Christine Zardecki,*,[†,§] Shuchismita Dutta,[†,§] David S. Goodsell,[‡,§] Maria Voigt,[†,§] and Stephen K. Burley[†,§,∥,⊥,#]

[†]Center for Integrative Proteomics Research and Department of Chemistry and Chemical Biology, Rutgers, The State University of New Jersey, 174 Frelinghuysen Road, Piscataway, New Jersey 08854, United States

[‡]Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, California, 92037, United States

[§]RCSB Protein Data Bank

[∥]Institute for Quantitative Biomedicine, Rutgers, The State University of New Jersey, 174 Frelinghuysen Road, Piscataway, New Jersey, 08854, United States

[⊥]Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San Diego, 9500 Gilman Drive, La Jolla, California 92093-0657, United States

[#]San Diego Supercomputer Center, University of California, San Diego, 10100 Hopkins Dr, La Jolla, California, 92093

**ABSTRACT:** The Research Collaboratory for Structural Bioinformatics (RCSB) Protein Data Bank (PDB) supports scientific research and education worldwide by providing access to annotated information about three-dimensional (3D) structures of macromolecules (e.g., nucleic acids, proteins), and associated small molecules (e.g., drugs, cofactors, inhibitors) in the PDB archive. Researchers, educators, and students use RCSB PDB resources to study the shape and interactions of biological molecules and their implications in molecular biology, medicine, biotechnology, and beyond. RCSB PDB supports development of standards for data deposition, representation, annotation, and validation of atomic structural data obtained from various experimental methods. Uniform representation of PDB data is essential for providing consistent search and analysis capabilities for all PDB users, from beginning students to domain experts. The RCSB PDB Web site provides tools for searching, visualizing, and analyzing PDB data, including easy exploration of chemical interactions that stabilize macromolecules and play important roles in their interactions and functions. In addition, educational resources are available for free and unrestricted use in the classroom for exploring chemistry and biology at the molecular level.

**KEYWORDS:** *General Public, High School/Introductory Chemistry, First-Year Undergraduate/General, Graduate Education/Research, Biochemistry, Interdisciplinary/Multidisciplinary, Internet/Web-Based Learning, Nucleic Acids/DNA/RNA, Proteins/Peptides, X-ray Crystallography*

The Protein Data Bank (PDB) is the first open access digital resource in biology for sharing three-dimensional (3D) protein structures.[1] The PDB was established in 1971 with 7 structures, and has grown exponentially to provide access to more than 113,000 entries of natural and designed macromolecules (proteins, nucleic acids and carbohydrates), more than 84,000 of which are complexed with small chemical components (solvent molecules, ions, cofactors, inhibitors, and drugs). Originally, PDB was a resource designed for the structural biology community, but through the years, its utility has grown and the PDB users now include biologists, software developers, computational and other scientists, bioinformaticians, students, educators and the general public.

The PDB archive of data files is one of the most heavily used biological data resources worldwide. In 2014, more than 505,000,000 atomic coordinate and experimental data files were downloaded for research and education. These downloads also include routine downloads by pharmaceutical and biotechnology companies for use in proprietary drug discovery efforts. A huge number of free resources and tools utilize PDB data to serve their users. These range from educational resources such as the NIH 3D Print Exchange[2] and Proteopedia;[3] molecular viewers including Jmol/JSmol,[4] Chimera,[5] Pymol;[6] and many scientific research tools and databases.[7]

The PDB archive is managed by a collection of regional data centers, called the Worldwide Protein Data Bank (wwPDB),[8] spread across the United States, Europe, and Japan. wwPDB centers collaborate on data deposition and annotation/validation practices. Each member hosts a distribution center, and provides tools for access and usage. Research Collaboratory for Structural Bioinformatics (RCSB) PDB, based at Rutgers, The State University of New Jersey and the University of California, San Diego, is focused on providing resources for research and education.[9] As part of the wwPDB, RCSB PDB members curate PDB data and develop data standards and software for the deposition and annotation pipeline. RCSB PDB also aims to enable breakthroughs in scientific inquiry, medicine, drug discovery, and technology by offering tools that provide rich structural views

**Table 1. Data Dictionaries Used in PDB Biocuration**

| Resource | Description |
|---|---|
| PDB Exchange (PDBx)/macromolecular Crystallographic Information File (mmCIF) Dictionary http://mmcif.wwpdb.org (accessed 6 Nov 2015) | Crystallographic data dictionaries with extensions describing NMR, 3DEM, and protein production, contains >4300 data items |
| Chemical Component Dictionary (CCD) http://www.wwpdb.org/data/ccd (accessed 6 Nov 2015) | Each chemical definition includes descriptions of chemical properties such as stereochemical assignments, chemical descriptors (SMILES[16] and InChI[21]), systematic chemical names, and idealized coordinates (generated using Molecular Networks' Corina, and if there are issues, OpenEye's OMEGA). Contains >18,000 small molecules |
| Biologically Interesting molecule Reference Dictionary (BIRD) http://www.wwpdb.org/data/bird (accessed 6 Nov 2015) | Molecular weight and formula, polymer sequence and connectivity, descriptions of structural features and functional classification, natural source (if any), and external references to corresponding UniProt or Norine entries. Contains ~750 small molecules |

of biological molecules and systems. In addition to supporting biological and chemical learning, RCSB PDB is an exemplar of the new discipline of data science; it provides a glimpse into science history, and serves as a resource for developing database query and analysis skills.

In this paper, we describe data annotation practices, and highlight RCSB PDB resources available for query and analysis, and education.

## ARCHIVING PDB DATA: DEPOSITION, ANNOTATION, AND VALIDATION

The PDB archive includes 3D structures of macromolecules (primarily proteins, DNA and RNA) as determined by experiments using X-ray crystallography, nuclear magnetic resonance (NMR), and/or 3D electron microscopy (3DEM). For each new structure, researchers submit atomic coordinates, experimental data, and molecular information using specialized tools, and wwPDB biocurators then review, annotate and validate the entry. Atomic coordinates are checked for consistency with the known sequence of the macromolecule and chemical structure of small molecules, and biological assemblies are defined and annotated. The entry is also extensively annotated with experimental information and cross-referenced to related entries and external resources. The wwPDB collaborates closely with archives that maintain related data, including the Cambridge Structural Database of small molecule crystal structures,[10] and the EMDataBank for 3D Electron Microscopy maps and models.[11]

Uniform annotation of PDB data enables consistent searching and analysis across the archive. To facilitate uniformity, PDB curation relies upon standard data dictionaries that define the representation of all components in the entry (Table 1). While the PDB Exchange (PDBx)[12] and macromolecular Crystallographic Information File (mmCIF)[13] data dictionaries provide the bases for internal data cross-referencing, processing, annotation, validation and database management operations, the Chemical Component Dictionary[14] describes all standard and modified amino acids/nucleotides, small molecule ligands, and solvent molecules. All chemical components are checked against this dictionary during annotation.[17] Protein and nucleic acid polymers can be built by linking together individual chemical components in a specified order denoted by the polymer sequence. Specialized molecules with unusual chemistries and interesting biological and pharmaceutical functions, such as peptide-like inhibitors and many antibiotics, are included in the Biologically Interesting molecule Reference Dictionary.[15] Use of these dictionaries enables specialized query and access to small molecule information, from specialized information pages for all ligands in the PDB to tools for visualizing ligand−protein interactions.

wwPDB uses community-accepted standards to "validate" deposited data, and produces reports that provide an assessment of structure quality based upon geometric and experimental data validation. wwPDB has convened method-specific Validation Task Forces[18] to develop recommendations for validation standards and software for use in annotation.[19] During annotation, validation reports are provided to the depositor to highlight any areas of concern. Journal editors may request these reports from authors to inform manuscript review.

The PDB archive includes structural information with a wide range of quality, due to the many challenges inherent in the experimental methods, and the nature of the molecule(s) or complex(es) being studied. Validation reports for all PDB structures determined by X-ray crystallography include a "slider" graphic (Figure 1) to summarize the quality of the
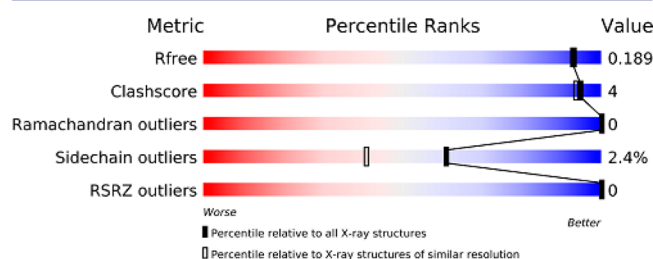


**Figure 1.** Validation report slider graphic indicates the quality of an entry as compared with other PDB entries. Shown is the slider image for an entry with better overall quality relative to all X-ray structures (PDB entry 1cbs, a small protein with a ligand at 1.8 Å resolution).[20]
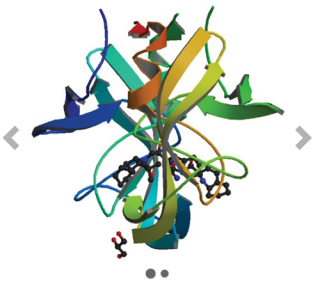
determined structure as compared with other structures in the archive. These graphics are displayed on the RCSB PDB Web site to help users find the structures of highest quality and to provide a warning to be critical when using structures with less experimental support.

## QUERY, REPORTING, AND ANALYSIS

The RCSB PDB Web site[21] integrates PDB data, related information about the structure from external scientific resources, and precalculated comparative and statistical information for query, analysis, and visualization.[22] On average, the Web site is accessed by ~325,000 unique users every month from ~190 countries. The top search bar supports simple keyword searches (ID, author, molecule name, chemical name), and suggests results options organized by different categories, including organism, molecule name, or experimental technique. Advanced searching allows users to combine searches for many specific data items, such as molecule name, authors, experimental techniques, and resolution. Browsers are available to find PDB structures organized using data annotations from external resources (e.g., Gene Ontology terms describing biological process, cellular component, and molecular function;[23] Enzyme Classification;[24] the World Health Organization Collaborating Centre's Anatom-

**Figure 2.** Structure Summary page highlights for PDB entry 4qgi, an HIV protease complexed with the drug saquinavir.[29] (A) PDB ID, title, 2D image with links to interactive 3D viewers; (B) information about the publication describing the entry, with links to PubMed and reference information; (C) small molecule information with links to summary information and 3D views. Structure Summary pages also include links to the atomic coordinates, sequence information, experimental data, and validation information.

ical Therapeutic Chemical (ATC) Classification System); or by exploring drill-down distributions of standard characteristics (e.g., polymer type, organism, resolution). Since general searches, such as for "hemoglobin", can return hundreds of matching structures, a variety of tools are available to help narrow the focus of the inquiry. For example, search autosuggestions, query refinement options, and sorting results by most recently solved entry, molecular weight, and resolution can help guide users to structures of higher quality and relevance. When available, searches will also return corresponding educational *Molecule of the Month* features, which are described in more detail below.

Every entry has a Structure Summary page that provides access to many aspects of the structure (Figure 2). Interactive 3D viewers, including Jmol/JSmol and Ligand Explorer,[25] can be used to rotate the molecule, select specific residues, and highlight ligand–protein interactions.[4] Many of the data items shown can be used to query for other entries with the same data (e.g., sequence database reference, specific chemical component). The entry's corresponding validation slider described above

(Figure 1) is also displayed. Additional links provide related information from external scientific sources, such as functional annotations from CATH[26] and SCOP,[27] and sequence information from UniProt.[28]

The RCSB PDB Web site also supports small molecule searching by ID, name, formula, or chemical drawing. Summary pages are available for each chemical component to provide 2D and 3D visual representations, any subcomponent information, corresponding DrugBank[7d] information, and access to atomic coordinates (Figure 3).

Additionally, the RCSB PDB *Mobile* app supports access to molecular data on mobile devices (Android and iOS).[30] App users can search the entire PDB database by ID, molecule name, or author name, and view a summary report for corresponding structures. The 3D visualization program NDKMol[31] allows app users to interactively view structures and save views as images. The app can be used to search and explore structures during lectures, symposia, and poster sessions.
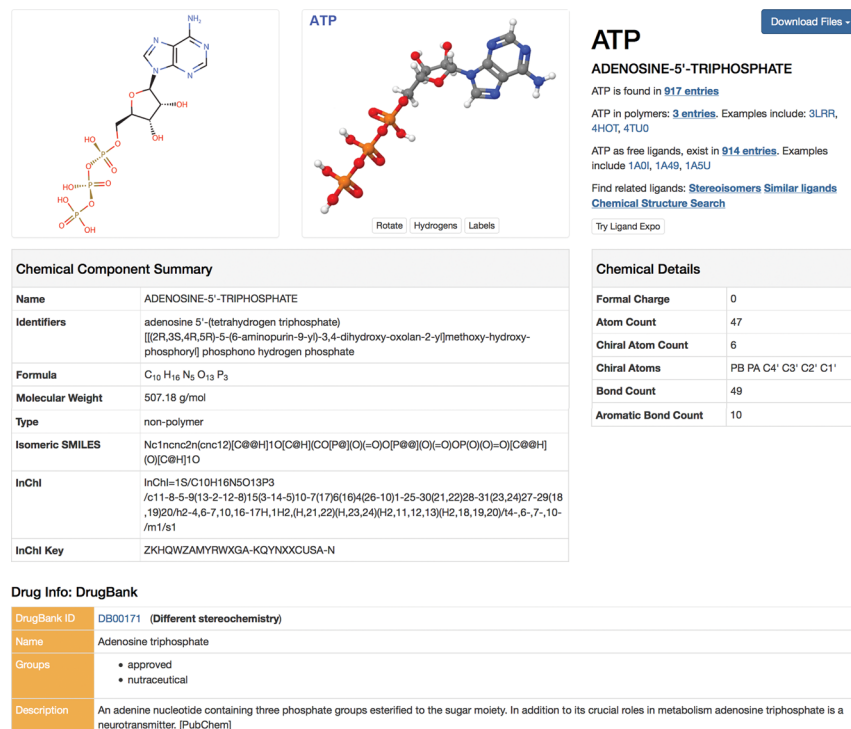
**Figure 3.** Ligand summary for ATP. The information shown on this page is built using the corresponding entry for ATP contained in the Chemical Component Dictionary. Blue text links to "query-by-example" searches of the archive, and may be used to find entries that include the ligand. Data highlighted in orange are integrated from external resources such as DrugBank to provide any pharmaceutical context to the structure not included in the deposited PDB entry.[7d] The top ligand image can be downloaded, and the 3D view of the molecule rotated using Jmol.[4]

## ■ EDUCATIONAL RESOURCES

A recent survey revealed that the RCSB PDB Web site is used by a wide range of communities, including educators and students at the high school and university levels. For the educational community, RCSB PDB tools take a subject-based approach, allowing chemistry students to visualize the chemical and structural basis of biological processes, such as *how is oxygen stored and transported to different cells in the human body* or *how do specific drugs act on their target proteins*. Tools are available to find and visualize PDB molecules of interest and explore their interactions. In addition, a number of resources provide nonexperts with information and examples of how to interpret the molecule's functions in the context of chemical interactions.

To enable broader access by educators and students, RCSB PDB established an education-focused portal to PDB data. The PDB-101 Web site[32] hosts regularly published articles, educational materials, and introductions to information specific to PDB data and their representations within the archive.

The *Molecule of the Month* column[33] serves as the foundation of many PDB-101 resources. Since 2000, this feature has highlighted selected biological structures with text, images, and interactive views (Figure 4). The column provides a curated set of example structures from the archive to illustrate key points about a molecule. The collection of Molecule of the Month articles has been organized by biological concept into a browser to enable top-down searching by functional category. Rather than searching by a particular molecule, users can browse articles about specific topics (such as viruses or the immune system).

Beyond *Molecule of the Month*, PDB-101 offers a wide range of educational materials to explore biomolecular structure and



**Figure 4.** Images from the Molecule of the Month column on HIV Capsid.[34] (A) The feature begins with a description of the general biology, with an illustration of a single capsid protein (left; PDB entry 1e6j)[35] and the full capsid (right; 3j3q).[36] (B) A description of molecules that interact with the HIV capsid follows, showing TRIM5 (2lm3)[37] and cyclophilin (1ak4).[38] (C) The final section presents a JSmol[39] view to allow interactive exploration of the HIV capsid hexamer (3mge)[40] and pentamer (3p05).[41]

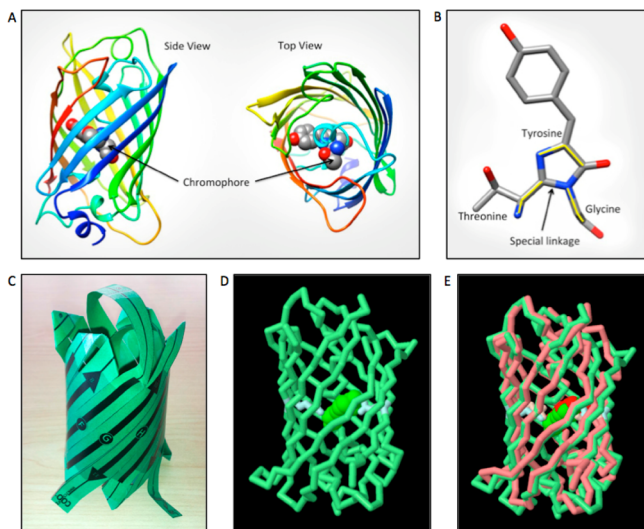**Figure 5.** Exploring the structure of fluorescent proteins. Images from the GFP activity at PDB-101. (A) Ribbon view of GFP with the chromophore highlighted in its core, based on entry 1ema.[42] (B) Close-up of specific residues chemically linked to form the chromophore. (C) A downloadable PDF can be used to create a paper model of GFP. Curated JSmol view of GFP highlights (D) conserved residues in the protein core that play a role in the chromophore formation and (E) the same structure superimposed with a distant relative DsRed (PDB entry 1g7k),[43] showing high structure conservation despite limited sequence similarity.

function. Videos and animations explore specific topics, from the biology of HIV to protein folding. *Understanding PDB Data* offers a general introduction to structural biology and PDB data files, with topics such as crystallographic resolution and biological assemblies. Hands-on model activities can be used to explore the folding of proteins and nucleic acids.

Resources have been compiled to provide activities, lesson plans, and curricula. For example, the Green Fluorescent Protein (GFP) activity (Figure 5) references the GFP Molecule of the Month, and uses a paper model to provides a hands-on understanding of the polymer nature of the protein, overall shape and folding of the protein, and the assembly of the GFP chromophore from chemical interactions between three specific amino acids in the core of the protein. An interactive JSmol view demonstrates the chemistry involved in the creation of chromophore.

In December 2014, high school curricula were launched at PDB-101 to combine a variety of PDB-101 and external resources (videos, animations, games, activities and exercises) for comprehensive learning about the biology of HIV/AIDS at introductory and advanced high school levels. Using these materials, classes studied the HIV lifecycle, interactions with the immune system, and the basis of current infection treatments. Web site materials were accessed more than 8000 times during this pilot session. Based on feedback from high school instructors who participated in the pilot, the curricula have been reorganized into individual modules (Biomolecular Structures and Models, Molecular Immunology, Molecular View of HIV/AIDS) that can be implemented in a variety of other lesson plans. Development of similar modules focusing on the structural basis of diabetes is underway.

PDB-101 is also used as a resource in college education. In a 2014 survey of PDB-101 usage, 28% of respondents were undergraduate students, and 33% were graduate students. The most popular areas of research selected were the life sciences (66%), chemistry (34%), and computational sciences (20%). Examples of how these PDB-related materials have been incorporated at the collegiate level are frequently highlighted in our *Education Corner*, a guest column published in our quarterly newsletter. Examples have included drug discovery projects,[44] interesting molecular visualizations,[45] and cell biology.[46] The usage of PDB data and the utility of accurate molecular visualizations in undergraduate education has been a topic of much study.[47] Other published examples of undergraduate classroom usage include 3D printing and models,[48] molecular modeling,[49] pharmaceutical and medicinal chemistry,[50] and beyond.[51]

## ■ CONCLUSION

The RCSB PDB offers free access to a broad range of primary research data and educational materials for all users. Structural entries in the PDB are extensively annotated and validated according to current community standards, providing a rich resource for chemical exploration. With the RCSB PDB tools and resources, users may explore detailed information about small and large biomolecules, their chemical interactions, as well as study broader structural and functional concepts in biology.

## ■ AUTHOR INFORMATION

**Corresponding Author**

*E-mail: Zardecki@rcsb.rutgers.edu.

**Notes**

The authors declare no competing financial interest.

## ■ ACKNOWLEDGMENTS

## ■ REFERENCES

(1) (a) Berman, H. M.; Kleywegt, G. J.; Nakamura, H.; Markley, J. L. How community has shaped the Protein Data Bank. *Structure* **2013**, *21* (9), 1485−91. (b) Protein Data Bank. Crystallography: Protein Data Bank. *Nat. New Biol.* **1971**, *233* (42), 223.

(2) NIH 3D Print Exchange. http://3dprint.nih.gov/ (accessed 6 Nov 2015).

(3) Prilusky, J.; Hodis, E.; Canner, D.; Decatur, W. A.; Oberholser, K.; Martz, E.; Berchanski, A.; Harel, M.; Sussman, J. L. Proteopedia: a status report on the collaborative, 3D web-encyclopedia of proteins and other biomolecules. *J. Struct. Biol.* **2011**, *175* (2), 244−52.

(4) Jmol: an open-source Java viewer for chemical structures in 3D. http://www.jmol.org/ (accessed 6 Nov 2015).

(5) Pettersen, E. F.; Goddard, T. D.; Huang, C. C.; Couch, G. S.; Greenblatt, D. M.; Meng, E. C.; Ferrin, T. E. UCSF Chimera−a visualization system for exploratory research and analysis. *J. Comput. Chem.* **2004**, *25* (13), 1605−12.

(6) DeLano, W. The PyMOL molecular graphics system, http://www.pymol.org (accessed 6 Nov 2015).

(7) (a) Kirchmair, J.; Markt, P.; Distinto, S.; Schuster, D.; Spitzer, G. M.; Liedl, K. R.; Langer, T.; Wolber, G. The Protein Data Bank (PDB), its related services and software tools as key components for in silico guided drug discovery. *J. Med. Chem.* **2008**, *51* (22), 7021−40. (b) Berman, H. M.; Kleywegt, G. J.; Nakamura, H.; Markley, J. L. The Protein Data Bank archive as an open data resource. *J. Comput.-Aided Mol. Des.* **2014**, *28*, 1009−1014. (c) Fernandez-Suarez, X. M.; Rigden, D. J.; Galperin, M. Y. The 2014 Nucleic Acids Research Database Issue and an updated NAR online Molecular Biology Database Collection. *Nucleic Acids Res.* **2014**, *42* (Database issue), D1−6. (d) Law, V.; Knox, C.; Dioumbou, Y.; Jewison, T.; Guo, A. C.; Liu, Y.; Maciejewski,

# DOCKET ALARM

# Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

## Real-Time Litigation Alerts

Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

## Advanced Docket Research

With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

## Analytics At Your Fingertips

Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

## API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

### LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

### FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

### E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.

fastcase®
Smarter legal research.