

## Beginner's Guide to the PDBbind Database (v.2018)

The PDBbind database provides a comprehensive collection of experimental binding affinity data for the biomolecular complexes in the Protein Data Bank (PDB). This type of information is the much needed basis for various computational and statistical studies on molecular recognition. PDBbind was first released to the public in May 2004. Over 6,500 users from over 70 countries have already registered to use this database. The PDBbind database is now updated annually to keep up with the growth of PDB. The current release is **version 2018**.

### What information does PDBbind provide?

- ❑ **Binding affinity data:** Originally, PDBbind only considered the complexes formed between proteins and small-molecule ligands. Other types of biomolecular complexes in PDB have been covered by PDBbind as well since 2008. This release contains binding data ( $K_d$ ,  $K_i$  &  $IC_{50}$  values) for protein-ligand (16,151), protein-protein (2,416), protein-nucleic acid (896), and nucleic acid-ligand (125) complexes. All binding data are curated by ourselves from over 34,700 original references.
- ❑ **Processed structural files:** PDBbind also provides processed, “clean” structural files for the protein-ligand complexes (16,100+) included in this release. In brief, the biological unit of each complex is downloaded from PDB and then split into a protein molecule (in PDB format) and a ligand molecule (in Mol2 and SDF format). Atom/bond types on the ligand molecule are assigned by a special computer program and then examined manually. Such structural files can be readily utilized by most molecular modeling software. They are wrapped in a data package for download from the PDBbind-CN web site.
- ❑ **Web-based display and analysis tools:** The user can access PDBbind through a web portal at <http://www.pdbbind-cn.org/>. Registration is free for academic and commercial users. On the PDBbind-CN web site, basic information of each complex is summarized on a single page. Text-based and structure-based search among the contents of PDBbind is also enabled. This web site actually provides structural information for all valid protein-ligand complexes in the Protein Data Bank, not limited to those with known binding data.

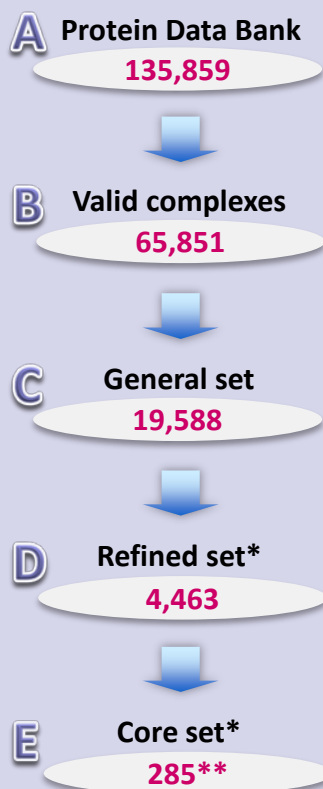
### Basic Information of the PDBbind Database\*

Version	Entries In PDB	All complex with binding data	Protein-ligand complex	Protein-protein complex	Protein-nucleic acid complex	Nucleic acid-ligand complex
2004	28,991	2,276	2,276	N.A.	N.A.	N.A.
...	...	...	...	...	...	...
2014	96,952	12,995	10,656	1,592	660	87
2015	105,183	14,620	11,987	1,807	717	109
2016	114,344	16,179	13,308	1,976	777	118
2017	124,962	17,900	14,761	2,181	837	121
<b>2018</b>	<b>135,859</b>	<b>19,588</b>	<b>16,151</b>	<b>2416</b>	<b>896</b>	<b>125</b>

\*: Information of some earlier versions (v.2005 – v.2013) are not included in this table due to space limit.

## Basic structure of the PDBbind data set

The data sets in PDBbind are compiled through a stepwise process as follows.



\* Only complexed formed between proteins and small-molecule ligands are considered in this data set.  
\*\* This number is for core set v.2016.

(A) The PDBbind v.2018 is based on the contents of PDB officially released at the first week of 2018, which contained a total of **135,859** experimentally determined structures. Theoretical models are not considered by us.

(B) The entire PDB was screened by a set of computer programs to identify four major categories of molecular complexes, including protein-small ligand, nucleic acid-small ligand, protein-nucleic acid and protein-protein complexes. This step identified a total of **65,851** entries as valid complexes.

(C) The primary reference of each complex was examined to collect experimentally determined binding affinity data ( $K_d$ ,  $K_i$  and  $IC_{50}$ ) of the given complex. Binding data for **19,588** complexes were collected in this way. They are the main body of the PDBbind database, which is referred to as the “**general set**”.

(D) As an additional feature, a “**refined set**” was compiled to select the protein-ligand complexes with better quality out of the general set. A number of filters regarding binding data, crystal structures, as well as the nature of the complexes were applied to selection (see ref.3 below for details). The refined set in this release consists of **4,463** protein-ligand complexes.

(E) A “**core set**” was also included in previous releases of PDBbind. Compilation of the core set aims at providing a relatively small set of high-quality protein-ligand complexes for validating docking/scoring methods. In particular, the core set has served as the primary test set in the Comparative Assessment of Scoring Functions (CASF) benchmark developed by our group. The core set is not included in the PDBbind data package any more because it is not updated annually as PDBbind itself. Besides, the core set is more than a list of protein-ligand complexes but with a huge amount of derivative data. Researchers can obtain the core set by downloading the CASF data package at <http://www.pdbbind-cn.org/casf.asp>.

## References and notes

The PDBbind database is currently maintained by Prof. Renxiao Wang’s group at the Shanghai Institute of Organic Chemistry, Chinese Academy of Sciences. To cite the PDBbind database, please refer to the following references:

- (1) Liu, Z.H. et al. *Acc. Chem. Res.* 2017, *50*, 302-309. (PDBbind v.2016)
- (2) Liu, Z.H. et al. *Bioinformatics*, 2015, *31*, 405-412. (PDBbind v.2014)
- (3) Yan, L.; et al. *J. Chem. Inf. Model.*, 2014, *54*, 1700-1716. (PDBbind v.2013 & CASF-2013)
- (4) Cheng, T. J.; et al. *J. Chem. Inf. Model.*, 2009, *49*, 1079-1093. (PDBbind v.2007 & CASF-2007)
- (5) Wang, R. X.; et al. *J. Med. Chem.* 2005, *48*, 4111-4119; *J. Med. Chem.* 2004, *47*, 2977-2980. (early versions)