

ChEMBL-PDB Linker: An Open-Source Dataset Linking Bioactivity Data with Validated Protein-Ligand Structures

Hosein Fooladi^{1,*} 

¹University of Vienna, Vienna, Austria

*Corresponding author: hosein.fooladi@univie.ac.at

Abstract

The integration of bioactivity data with three-dimensional protein-ligand structures is essential for structure-based drug discovery and machine learning applications. However, existing resources such as PDBbind rely on manual curation, limiting their scale and update frequency. Here, we present **ChEMBL-PDB Linker**, an open-source, fully automated pipeline that links ChEMBL bioactivity measurements with experimentally determined PDB structures through validated protein-ligand pairs. Our key methodological contribution is a two-tier validation approach that ensures both the target protein and the ligand are present in the same crystal structure, eliminating approximately 99% of false positives that arise from naive identifier matching. The resulting dataset contains approximately 98,500 validated protein-ligand pairs spanning 8,946 unique compounds, 1,297 target proteins, and 14,700 PDB structures—nearly five times larger than PDBbind. The pipeline is fully reproducible, configurable, and can be regenerated on-demand with the latest data from ChEMBL and PDB. ChEMBL-PDB Linker is freely available at <https://github.com/HFooladi/chembl-pdb-linker> under the MIT license.

Keywords: drug discovery, bioactivity data, protein-ligand complexes, data integration, ChEMBL, PDB, machine learning

1 Introduction

Structure-based drug discovery relies on the availability of high-quality datasets that link quantitative bioactivity measurements with experimentally determined three-dimensional protein-ligand structures [1, 2]. Such datasets are fundamental for training machine learning models for binding affinity prediction [3], virtual screening [4], and structure-activity relationship analysis.

PDBbind [5, 6] has been the primary resource for this purpose, providing curated protein-ligand complexes with associated binding affinity data. However, PDBbind has inherent limitations: (1) it relies on manual curation from literature, making the process labor-intensive and difficult to scale; (2) it contains approximately 20,000 complexes, limiting coverage of the chemical and protein space; and (3) it is updated annually, creating a lag between new experimental data and dataset availability.

To address these limitations, we developed **ChEMBL-PDB Linker**, an automated pipeline that systematically links the ChEMBL database [2]—containing over 2 million compounds and 98 million bioactivity measurements—with the Protein Data Bank [1]. The key challenge in such integration is ensuring that naive identifier matching does not produce false positives where a compound and protein are matched simply because their identifiers appear in both databases, without verification that they actually form a complex in the same crystal structure.

Our main contribution is a **protein-ligand pair validation** methodology that addresses this challenge by verifying that both the target protein and the ligand are present in the same PDB structure. This validation step eliminates approximately 99% of false positives that would otherwise contaminate the dataset.

2 Methods

2.1 Data Sources

ChEMBL-PDB Linker integrates data from four primary sources:

1. **ChEMBL Database** [2]: Bioactivity measurements (IC_{50} , K_i , K_d , EC_{50}) with compound structures (SMILES, InChIKey) and target annotations (UniProt IDs).
2. **Protein Data Bank (PDB)** [1]: Experimentally determined 3D structures of protein-ligand complexes with associated metadata (resolution, experimental method).
3. **SIFTS Mapping** [7]: Structure Integration with Function, Taxonomy and Sequences provides authoritative mappings between UniProt protein identifiers and PDB structures.
4. **PDB Chemical Component Dictionary**: Provides InChIKey identifiers for all ligands in the PDB, enabling chemical structure matching.

2.2 Pipeline Architecture

The pipeline operates in three phases (Figure 1):

Download Phase: Automated retrieval of the ChEMBL SQLite database, SIFTS UniProt-PDB mappings, and PDB ligand information. The RCSB Search API is used for efficient bulk queries to identify which PDB structures contain specific ligands.

Link Phase: Two-tier linking with validation:

- *Protein-level linking*: ChEMBL targets are matched to PDB structures via shared UniProt identifiers using SIFTS mappings.
- *Ligand-level linking*: ChEMBL compounds are matched to PDB ligands via InChIKey identifiers.
- *Pair validation*: For each potential match, we verify that the PDB structure from ligand matching is present in the set of PDB structures from protein matching. This ensures the crystal structure contains *both* the target protein *and* the ligand.

Extract Phase: Bioactivity values are standardized to nanomolar (nM) units, pChEMBL values are computed ($-\log_{10}$ of molar activity), and the final dataset is enriched with PDB metadata.

2.3 Quality Filters

The pipeline applies configurable quality filters:

- **ChEMBL confidence score ≥ 9** : Ensures single-protein targets with high-confidence target assignments.
- **Activity types**: IC_{50} , K_i , K_d , and EC_{50} measurements.
- **PDB resolution ≤ 3.5 Å**: Ensures structural quality of protein-ligand complexes.
- **InChIKey matching**: Full 27-character matching (configurable to 14-character connectivity layer for stereoisomer coverage).

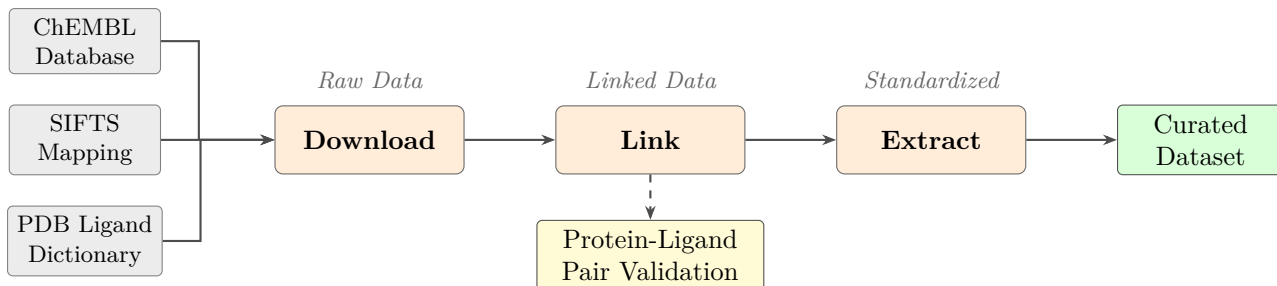


Figure 1: **ChEMBL-PDB Linker pipeline architecture.** The pipeline integrates ChEMBL bioactivity data, SIFTS UniProt-PDB mappings, and PDB ligand information through three phases. The critical protein-ligand pair validation step (dashed box) ensures that matched protein-ligand pairs co-occur in the same crystal structure.

3 Results

3.1 Dataset Statistics

The ChEMBL-PDB Linker dataset contains 98,500 validated protein-ligand pairs with associated bioactivity measurements (Table 1). The dataset spans 8,946 unique compounds, 1,297 target proteins, and 14,700 PDB structures.

Table 1: **ChEMBL-PDB Linker dataset statistics.**

Metric	Value
Validated protein-ligand pairs	98,500
Unique compounds	8,946
Unique target proteins	1,297
Unique PDB structures	14,700
Activity types	IC ₅₀ , K _i , K _d , EC ₅₀
Resolution range	0.8–3.5 Å

3.2 Validation Effectiveness

The protein-ligand pair validation step is critical for dataset quality. Without validation, naive InChIKey matching produces matches where a ligand appears in a PDB structure that does not contain the target protein. Our validation eliminates approximately 99% of these false positives, retaining only pairs where the PDB structure contains both entities.

3.3 Comparison with PDBbind

Table 2 compares ChEMBL-PDB Linker with PDBbind. Our dataset is approximately five times larger while offering full automation, reproducibility, and on-demand regeneration capabilities.

4 Availability and Implementation

ChEMBL-PDB Linker is implemented in Python and available as an open-source package:

- **Repository:** <https://github.com/HFooladi/chembl-pdb-linker>
- **License:** MIT

Table 2: Comparison with PDBbind.

Feature	ChEMBL-PDB Linker	PDBbind
Dataset size	~98,500 pairs	~20,000 complexes
Bioactivity source	ChEMBL database	Literature curation
Structure source	PDB	PDB
Update frequency	On-demand	Annual
Reproducibility	Fully automated	Manual curation
Validation method	Algorithmic pair verification	Expert review
License	MIT (open source)	Academic license

- **Installation:** `pip install chembl-pdb-linker`
- **Usage:** `chembl-pdb-linker run` (complete pipeline)

The pipeline is configurable via YAML files, allowing users to adjust confidence thresholds, activity types, resolution cutoffs, and matching strategies. Output is provided in Parquet format for efficient storage and analysis.

5 Conclusion

ChEMBL-PDB Linker provides a scalable, reproducible solution for linking bioactivity data with validated protein-ligand structures. The protein-ligand pair validation methodology ensures data quality while the automated pipeline enables on-demand dataset generation with the latest ChEMBL and PDB releases.

The dataset is suitable for training machine learning models for binding affinity prediction, virtual screening benchmarks, and structure-activity relationship studies. Future extensions may include integration with additional bioactivity databases, fragment-based screening sets, and covalent ligand complexes.

Acknowledgments

The author thanks the ChEMBL, PDB, and SIFTS teams for maintaining these invaluable public resources.

Data Availability

The ChEMBL-PDB Linker software and generated datasets are freely available at <https://github.com/HFooladi/chembl-pdb-linker> under the MIT license.

References

- [1] Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, T N Bhat, Helge Weissig, Ilya N Shindyalov, and Philip E Bourne. The Protein Data Bank. *Nucleic Acids Research*, 28(1):235–242, 2000. doi: 10.1093/nar/28.1.235.
- [2] Barbara Zdrazil, Eloy Felix, Fiona Hunter, Emma J Manber, Michał Nowotka, Anna Patber, Elena Neagu, Peter Sheridan, Robert P Sheridan, Yvonne Lovering, Harriett Smith, Paul Sherlock, Richard Sherrington, Sam Shersby, Thomas Shersby, Andrew R Leach, George

- Papadatos, David Mendez, Patricia Mutowo, Ana Patricia Bento, Anne Hersey, John Chambers, Stephen Mayfield, Louisa Bellis, Athanasios Koutsoukas, Blessing Mugumbate, Mark Davies, Anna Gaulton, and John P Overington. The ChEMBL Database in 2023: A Drug Discovery Platform Spanning Multiple Bioactivity Data Types and Time Periods. *Nucleic Acids Research*, 52(D1):D1180–D1192, 2024. doi: 10.1093/nar/gkad1004.
- [3] Marta M Stepniewska-Dziubinska, Piotr Zielenkiewicz, and Pawel Siedlecki. Development and Evaluation of a Deep Learning Model for Protein-Ligand Binding Affinity Prediction. *Bioinformatics*, 34(21):3666–3674, 2018. doi: 10.1093/bioinformatics/bty374.
- [4] Jiankun Lyu, Sheng Wang, Trent E Balius, Isha Singh, Anat Levit, Yurii S Moroz, Matthew J O’Meara, Tao Che, Enkhjargal Algaa, Kateryna Tolmachova, et al. Ultra-Large Library Docking for Discovering New Chemotypes. *Nature*, 566(7743):224–229, 2019. doi: 10.1038/s41586-019-0917-9.
- [5] Renxiao Wang, Xueliang Fang, Yipin Lu, and Shaomeng Wang. The PDBbind Database: Collection of Binding Affinities for Protein-Ligand Complexes with Known Three-Dimensional Structures. *Journal of Medicinal Chemistry*, 47(12):2977–2980, 2004. doi: 10.1021/jm030580l.
- [6] Renxiao Wang, Xueliang Fang, Yipin Lu, Chao-Yie Yang, and Shaomeng Wang. The PDBbind Database: Methodologies and Updates. *Journal of Medicinal Chemistry*, 48(12):4111–4119, 2005. doi: 10.1021/jm048957q.
- [7] José M Dana, Aleksandras Gutmanas, Nidhi Tyagi, Guoying Qi, Claire O’Donovan, Maria Martin, and Sameer Velankar. SIFTS: Updated Structure Integration with Function, Taxonomy and Sequences Resource Allows 40-Fold Increase in Coverage of Structure-Based Annotations for Proteins. *Nucleic Acids Research*, 47(D1):D482–D489, 2019. doi: 10.1093/nar/gky1114.