

PAPER

Assessing the Reliability of AlphaFold3 Predictions for Protein-Ligand Affinity Prediction via Sfcnn

Guo Yu,¹ Yiming Wu¹ and Yiyang Tan¹¹School of Information Science and Technology, ShanghaiTech University, 393 Middle Huaxia Road, 201210, Shanghai, China

Abstract

This study systematically evaluates the reliability of AlphaFold3 (AF3)-predicted protein structures for protein-ligand affinity (PLA) prediction tasks. The Sfcnn model, a 3D convolutional neural network (CNN) for PLA prediction, was reproduced using PyTorch. Model performance was validated on the PDBbind v2019 refined set for training and the CASF-2016 core set for testing. Subsequently, AF3-derived protein structures from the CASF-2016 core set were assessed and compared to experimentally determined structures using Sfcnn scores to determine the suitability of AF3 predictions in PLA applications.

Key words: AlphaFold3, protein-ligand affinity, CNN scoring function, CASF-2016

Introduction

Background

AlphaFold3 (AF3), DeepMind's latest AlphaFold model, predicts protein and protein-ligand structures with high accuracy. It extends AlphaFold2 by adding explicit ligand modeling, enhanced multimer assembly support, and optimized multiple sequence alignments (MSAs) via deep neural networks trained on extensive sequence and structural data [1].

Sfcnn is a 3D convolutional neural network-based scoring function introduced by Wang et al. [4] in 2022, designed to provide accurate and reliable predictions of binding affinities for protein-ligand complexes.

Objective

The primary objective of this study is to evaluate the reliability of AlphaFold3-predicted protein-ligand complex structures for protein-ligand affinity (PLA) prediction. Using the Chai-1 server for AF3 structure generation, we ensure support for custom ligands and robust MSA construction. The resulting AF3 structures are assessed with the reproduced Sfcnn model, and predicted affinities are compared to those from experimentally determined structures. This enables a direct evaluation of AF3's suitability for PLA prediction and highlights its current strengths and limitations.

Materials and Methods

Datasets

The Sfcnn network was trained using protein-ligand complexes from the PDBbind v2019 refined set[3], which includes experimentally determined binding affinities (pKa values). The

model was evaluated on the CASF-2016 [2] core set, comprising 285 protein-ligand complexes. To prevent data leakage, 266 overlapping protein complexes between the training and test sets were excluded, resulting in 4,852 unique training complexes.

Data Augmentation

To increase the effective size of the training set, each protein-ligand complex was randomly rotated nine times using random rotation matrices, yielding ten variants per complex. All variants share the same PLA score, resulting in a total of 48,520 training samples.

Featurization

Protein-ligand complexes are represented as 3D grids of size $20 \times 20 \times 20$, with each grid cell encoded as a one-hot vector of length 28. This vector comprises 14 protein atom types¹ and 14 ligand atom types. The resulting training tensor has shape (48520, 20, 20, 20, 28).

¹ <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-022-04762-3>

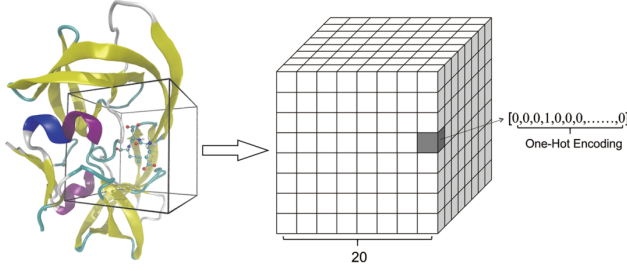


Fig. 1. Featurization of protein-ligand complexes. Example shown: PDB ID 1a30. Default resolution is $20 \times 20 \times 20$ with 28 atomic categories.

Sfcnn Network Architecture and Implementation

Architecture

The original Sfcnn publication describes four network architectures and three featurization strategies. The architecture depicted in Figure 2, combined with the aforementioned featurization, achieved optimal validation performance.

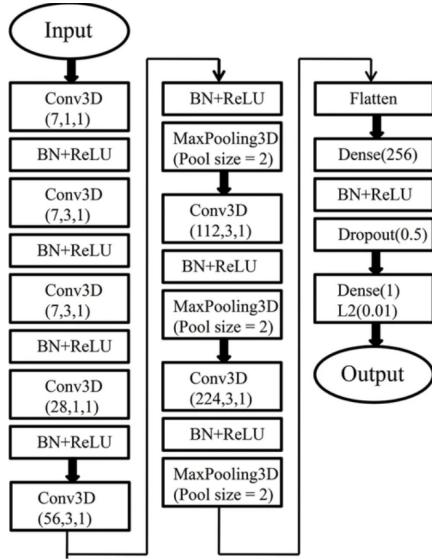


Fig. 2. Final CNN architecture for the Sfcnn network.

This architecture employs 3D convolutional layers with batch normalization and ReLU activation. L2 regularization is applied to the output layer to mitigate overfitting and enhance generalization.

Implementation Details

The PyTorch implementation closely mirrors the original TensorFlow version, with two key differences:

1. Due to PyTorch's `Conv3D` API, input tensors are permuted to shape (batch_size, 28, 20, 20, 20).
2. PyTorch lacks a direct L2 regularization API; instead, weight decay is applied to the final fully connected layer to approximate this effect.

Data Storage

The original Sfcnn implementation stored data as concatenated arrays in a single `.pkl` (pickle) file, requiring all data to reside in memory, which is impractical for extremely large datasets. The HDF5 format (`.h5`) via `h5py` was adopted for incremental writing and efficient storage. The resulting training grid occupies 40.1 GiB.

Training Procedure

Training and validation sets are partitioned as in the original study, with the validation set comprising indices 41,000 to 48,520. The resulting splits are: training (41,000, 20, 20, 20, 28), validation (7,520, 20, 20, 20, 28), and test (285, 20, 20, 20, 28), corresponding to a ratio of 84.00% : 15.42% : 0.58%.

Note that the original hyperparameters did not yield convergence in our PyTorch experiments. Both sets of hyperparameters are summarized below.

Table 1. Original and Reproduced Hyperparameters

Parameter	Original	Reproduced
Learning rate	0.004	0.00068
Batch size	64	32
Dropout rate	0.5	0.15
L2 regularization / FC weight decay	0.01	0.01
Epochs	200	400

Reproduced Results

Evaluation Metrics

Sfcnn performance is evaluated using the following metrics:

$$\begin{aligned}
 \text{RMSE} &= \sqrt{\frac{1}{N} \sum_{i=1}^N (y_{\text{predict}} - y_{\text{true}})^2} \\
 \text{MAE} &= \frac{1}{N} \sum_{i=1}^N |y_{\text{predict}} - y_{\text{true}}| \\
 \text{SD} &= \sqrt{\frac{1}{N-1} \sum_{i=1}^N ((ay_{\text{predict}} + b) - y_{\text{true}})^2} \\
 R &= \frac{\mathbb{E}[(y_{\text{predict}} - \mu_{y_{\text{predict}}})(y_{\text{true}} - \mu_{y_{\text{true}}})]}{\sigma_{y_{\text{predict}}} \sigma_{y_{\text{true}}}}
 \end{aligned}$$

where a and b are the slope and intercept of the linear regression between predicted and measured values, $\mathbb{E}[\cdot]$ denotes expectation, and μ and σ represent means and standard deviations, respectively.

Table 2. Performance Metrics on CASF-2016 Core Set

Metric	Reproduced Sfcnn	Original Sfcnn
Pearson R	0.7678	0.7928
RMSE	1.4647	1.3263
MAE	1.1633	1.0277
SD	1.3928	1.3252

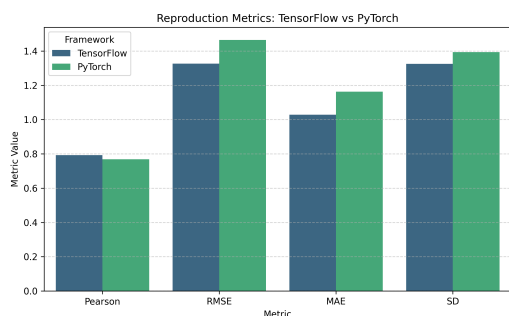


Fig. 3. Training curve for reproduced hyperparameters.

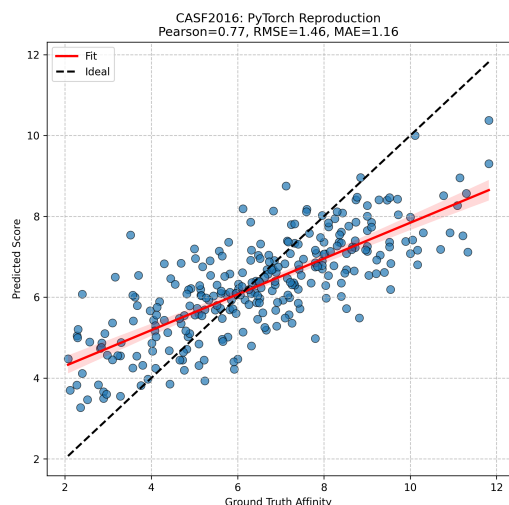


Fig. 4. Training curve for original hyperparameters.

Although the original Sfcnn reports superior metrics, its training process did not converge during reproduction, raising concerns regarding the reproducibility of the reported performance. Due to the lack of access to the original training data and the absence of author responses to data requests on GitHub², the original training process is deemed irreproducible. The training, validation, and testing results for the four metrics are presented in Figure 3 and Figure 4.

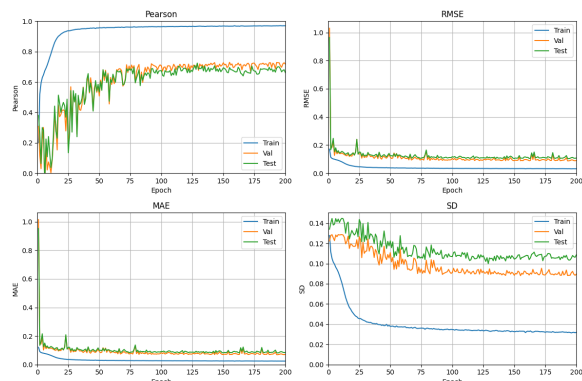


Fig. 5. Training curve for reproduced hyperparameters.

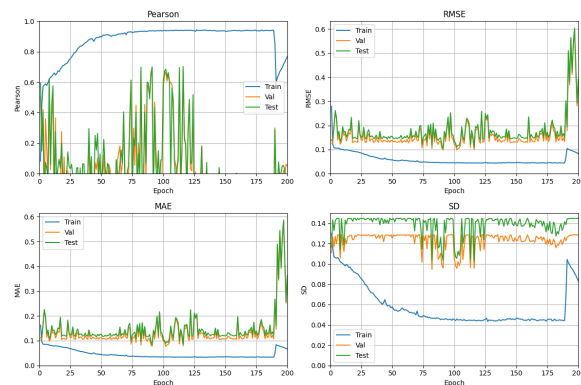


Fig. 6. Training curve for original hyperparameters.

The residuals of the test set (gaps between predicted and true values) are visualized in Figure 7. The histogram approximately follows a normal distribution, indicating that the model's predictions are generally unbiased, with most errors concentrated around zero.

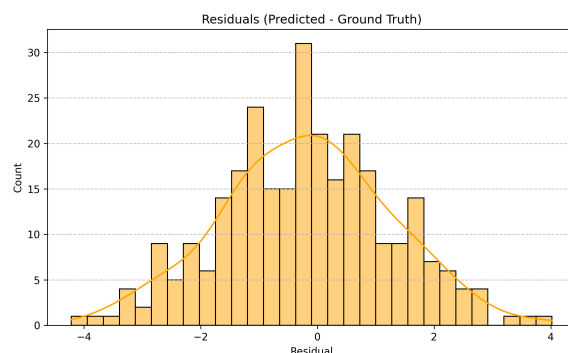


Fig. 7. Test residual histogram

We conjectured that the divergent results are caused by unusually high learning rate and dropout rate; therefore, the reproduced results are considered the baseline for subsequent AF3 result assessment.

AlphaFold3 Result Assessment

Dataset Selection

The assessment utilizes the CASF-2016 core set, excluding six proteins with structural complexity beyond AlphaFold3's predictive capacity, resulting in 279 proteins. Proteins with more than five isomorphous or heterogeneous chains were excluded, as detailed in Table 3.

² <https://github.com/bioinfocqupt/Sfcnn/issues/1>

Table 3. Excluded Complex Protein Structures

PDB ID	Number of Chains
2xb8	12
2ymd	10
3n76	12
3n7a	12
3n86	12
4ciw	12

Structure Generation and Processing

Protein structures were generated using the *Chai-1 online server*³. The AlphaFold3 online server⁴ was not used due to its inability to accept specific ligand SMILES codes. The MSA (Multiple Sequence Alignment) option was enabled with the MMseqs2 algorithm for all generations.

Server outputs were downloaded as zip archives containing multiple ranked structures and associated metrics. The top-ranked structure (`pred.rank.0.cif`) was selected for analysis. To avoid conversion errors between file formats (.cif, .pdb, .mol2), structures were parsed directly using the `MMCIFParser` from the `Bio.PDB` Python package, followed by featurization and grid mapping.

Atoms or isotopes not included in the 14 predefined atom types are categorized as *other*.

Scoring Protocol

To quantitatively assess the reliability of AlphaFold3 (AF3)-predicted structures for protein-ligand affinity (PLA) prediction, we evaluated the predicted complexes using the reproduced Sfcnn network, initialized with pre-trained weights (Pearson $R = 0.768$). The experimentally determined PLA values from the PDBbind v2019 core set served as the ground truth. For benchmarking, the predicted scores for AF3-generated structures are compared against both the ground truth and the Sfcnn scores obtained from experimentally resolved structures, employing the same evaluation metrics as in Section 4.

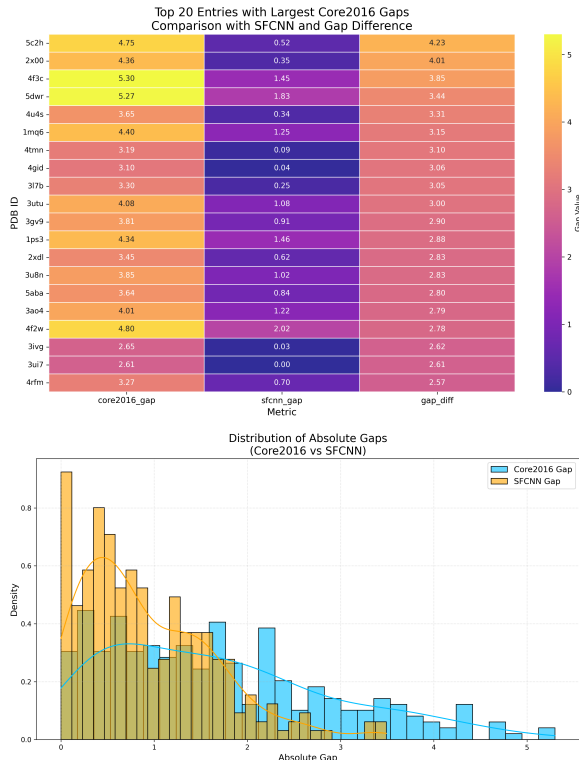
Table 4. Performance metrics for Sfcnn AF3-predicted structures compared to Sfcnn CASF2016 and groundtruth.

Metric	vs. CASF2016 Sfcnn	vs. Groundtruth
Pearson R	0.2930	0.3850
RMSE	2.0836	1.1669
MAE	1.6933	0.9231
SD	2.0825	1.0970

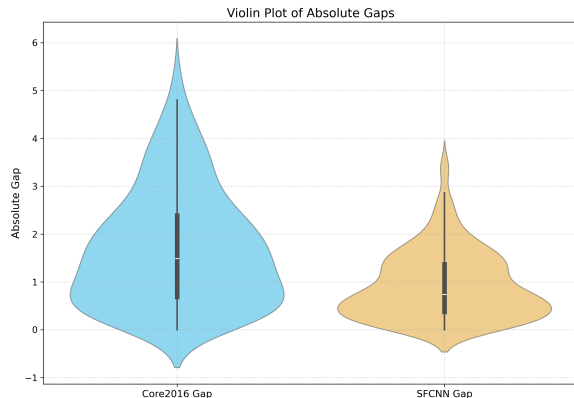
As summarized in Table 4, the correlation between predicted and experimental affinities decreases when using AF3-generated structures (Pearson $R = 0.2930$) compared to experimentally determined structures (Pearson $R = 0.3850$). Both RMSE and MAE are notably higher for AF3 predictions, indicating increased prediction error and reduced model reliability on these structures.

Residual Analysis and Visualization

To further dissect the sources and distribution of prediction errors, a series of visual analyses are presented below:

**Fig. 8.** Heatmap (top) of the top 20 complexes with the largest prediction gaps and histogram (bottom) of absolute prediction gaps for AF3 structures.

Figures 8 illustrate the distribution and magnitude of prediction errors. The heatmap highlights complexes with the largest discrepancies between AF3-predicted and ground truth PLA scores, while the histogram reveals a right-skewed error profile, indicating that a substantial fraction of complexes exhibit large deviations.

**Fig. 9.** violin plot of prediction gaps for AF3 structures.

³ <https://lab.chaidiscovery.com/dashboard>

⁴ <https://alphafoldserver.com/>

Figure 9 provides a complementary perspective: the violin plot visualizes the full error distribution, emphasizing the significantly higher presence of large gaps.

Interpretation and Implications

Collectively, these results demonstrate that while the reproduced Sfcnn model maintains moderate agreement with experimental affinities on the CASF-2016 core set, its predictive performance deteriorates when applied to AF3-predicted structures. The observed increase in RMSE and MAE, coupled with reduced correlation, indicates that current AF3 structural models introduce additional uncertainty into affinity prediction pipelines.

The top 10 error complexes and their corresponding protein structures and Chai-1 overall confidence scores are demonstrated in Table 5.

Table 5. Top 10 complexes with the largest prediction gaps. ‘heter’ denotes heterogeneous, ‘iso’ denotes isomorphous.

PDB ID	Aggregate Score	Prediction Gap	chain structure
5c2h	0.95	4.23	1 chain
2x00	0.93	4.01	5 heter
4f3c	0.97	3.85	2 iso
5dwr	0.90	3.44	1 chain
4u4s	0.89	3.31	2 heter
1mq6	0.97	3.15	2 heter
4tmn	0.63	3.10	2 iso
4gid	0.96	3.06	1 chain
3l7b	0.93	3.05	2 iso
3utu	0.97	3.00	2 heter

Statistics of these complexes show no significant correlation with protein structures, suggesting that AF3’s performance degradation is **not** attributable to the incapability of handling certain particular complex structures. However, the overall high aggregate scores raise concerns about the reliability of AF3’s confidence scores in the context of PLA prediction.

Hypotheses for AF3’s Reduced Performance

The observed decrease in PLA prediction accuracy for AF3-generated structures may be attributed to AF3’s methodological limitations and constraints in the training dataset. The following hypotheses are proposed with reference to the work of Zheng et al. [5]:

1. Diffusion Module Limitations:

AF3’s diffusion-based architecture lacks explicit physics-based energy calculations essential for accurate binding affinity prediction. While it generates structurally plausible conformations, the absence of energy functions prevents reliable ranking of binding interactions. Additionally, the model may inadequately represent conformational ensembles due to its reliance on static crystallographic structures during training.

2. AF3 Dataset Constraints:

The training dataset exhibits temporal limitations (pre-cutoff structures only) and conformational bias toward static X-ray structures, limiting the representation of dynamic binding processes. The underrepresentation of binding affinity data and potential memorization of training structures rather than learning generalizable physical

principles may contribute to reduced performance on novel protein-ligand pairs.

These hypotheses suggest that both target-specific structural challenges and broader methodological limitations of AF3 contribute to the observed performance gap. Future work could involve targeted analysis of outlier complexes and the development of structure prediction models specifically tailored for affinity prediction tasks.

Conclusion

This study presents a systematic evaluation of AlphaFold3-predicted protein structures for protein-ligand affinity prediction using a reproduced Sfcnn model. The findings highlight the importance of reproducibility in deep learning models for structural biology and demonstrate the current capabilities and limitations of AF3 in the context of PLA prediction.

Author Contributions

Guo Yu

- Designed and implemented the data pipeline, including dataset curation, preprocessing, featurization, data augmentation, and storage.
- Managed exclusion of overlapping complexes and ensured data compatibility with the reproduced network.
- Designed and maintained the entire AF3 generation-evaluation workflow, implemented K-fold cross-validation.

Yiming Wu

- Implemented and reproduced the Sfcnn neural network architecture in PyTorch, adapting the original TensorFlow design.
- Designed and executed evaluation process for both experimentally determined and AF3-predicted structures.

Yiyang Tan

- Conducted comparative studies and error analysis of model outputs, finished calculation of evaluation metrics (RMSE, MAE, SD, Pearson R).
- Generated all analysis visualizations, such as heatmaps, histograms etc..

All Authors

- Contributed to the training process and hyperparameter tuning.
- Participated in the AF3 result generation.
- Participated in the interpretation of results and manuscript writing.
- Provided critical review of the paper and process.

The project was conducted collaboratively with regular discussions to refine methodology and analysis.

External Libraries

- PyTorch:** Custom neural network implementation.
- Pandas:** Data storage and analysis.
- Numpy:** Data processing and manipulation.
- Matplotlib:** Visualization of training loss curves.

Acknowledgments

This work was completed as the final project for the CS177: Bioinformatics—Software Development and Applications course at ShanghaiTech University. The authors thank the course instructors and teaching assistants for their guidance and support throughout the project.

References

1. Dunger J. et al. Abramson J., Adler J. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, 2024.
2. Minyi Su, Qifan Yang, Yu Du, Guoqin Feng, Zhihai Liu, Yan Li, and Renxiao Wang. Comparative assessment of scoring functions: The casf-2016 update. *Journal of Chemical Information and Modeling*, 59(2):895–913, 2019.
3. Renxiao Wang, Xueliang Fang, Yipin Lu, Chao-Yie Yang, and Shaomeng Wang. The pdbind database: Methodologies and updates. *Journal of Medicinal Chemistry*, 48(12):4111–4119, 2005.
4. Yu Wang, Zhengxiao Wei, and Lei Xi. Sfcnn: a novel scoring function based on 3d convolutional neural network for accurate and stable protein–ligand affinity prediction. *BMC Bioinformatics*, 23(1):222, 2022.
5. Haiyang Zheng, Hanfeng Lin, Adebawale A. Alade, Jingjing Chen, Erika Y. Monroy, Min Zhang, and Jin Wang. Alphafold3 in drug discovery: A comprehensive assessment of capabilities, limitations, and applications. *bioRxiv*, 2025.