
Formatting Instructions For NeurIPS 2025

Anonymous Author(s)

Affiliation, Address

anon.email@example.org

Abstract

1 The abstract paragraph should be indented ½ inch (3 picas) on both the left- and
2 right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11
3 points. The word **Abstract** must be centered, bold, and in point size 12. Two line
4 spaces precede the abstract. The abstract must be limited to one paragraph.

5 Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut
6 labore et dolore magna aliqua quaerat voluptatem. Ut enim aeque doleamus animo, cum corpore
7 dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum
8 nobis.

9 1 Introduction

10 1.1 Sfenn Background

11 Sfenn is a 3D convolutional neural network based scoring function model proposed in 2022, which
12 aims to provide accurate and reliable scores for binding affinities of protein-ligand complexes.

13 1.2 Data Methods

14 1.2.1 Dataset

15 The Sfenn network was trained with protein–ligand complexes from the refined set of the PDBbind
16 database version 2019, which contains protein–ligand complexes and their corresponding binding
17 affinities expressed with pKa values, the trained network is later tested on the CASF-2016 core set,
18 which has 285 protein–ligand complexes.

19 Note that the overlaps between train set and test set (266 protein complexes) are excluded, leaving
20 4852 train complexes in total.

21 1.2.2 Augmentation

22 To scale up the training set, each protein-ligand complex is rotated randomly for 9 times using
23 random rotation matrices, those 10 complexes should bear the same PLA (protein-ligand affinity)
24 score, resulting in total 48520 complexes for training

25 1.2.3 Featurization

26 To capture the features of a protein-ligand complex, Sfenn uses the method of grid mapping and one-
27 hot encoding. Each complex is mapped to a 3D grid with resolution $20 \times 20 \times 20$, which is later
28 transformed into a 4D tensor. Each cell within the grid is a formed by an encoding list of length 28,
29 consists of 14 protein atom(isotope)¹ and 14 corresponding ligands, mapped with one-hot encoding
30 method. The final training tensor size is therefore **(48520, 20, 20, 20, 28)**.

¹Please refer to the original Sfenn paper for those atom types: <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-022-04762-3#availability-of-data-and-materials>

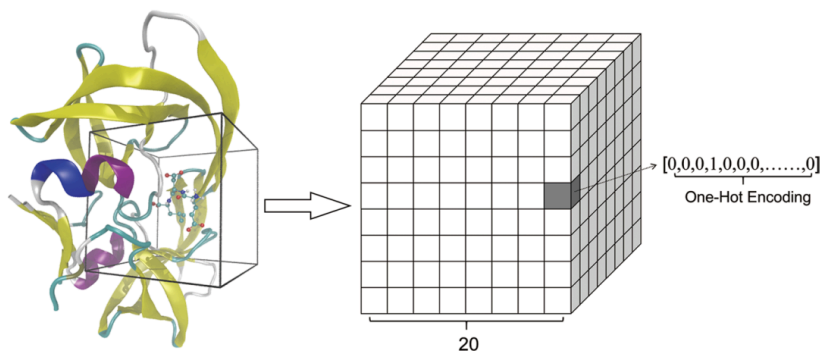


Figure 1: Featurization of the protein-ligand complexes. PDB ID 1a30 is shown as an example. In the default case, the resolution of $20 \times 20 \times 20$ and 28 categories of atomic types were used

1.3 Network

The original paper presents 4 different network structures along with 3 ways of featurization. The network shown in the figure, combining with the featurization method above achieved best performance on validation set.

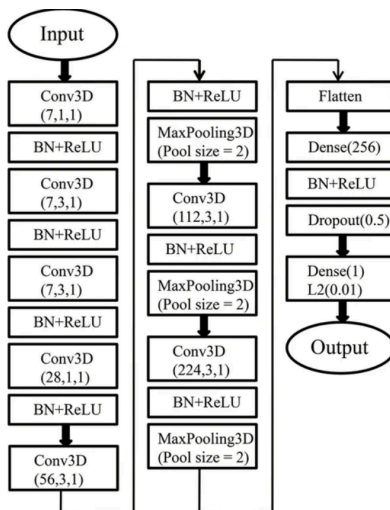


Figure 2: Final CNN structure for Sfcnn network

This network features 3D convolution layers with batch normalization and ReLU activation. L2 regularization was applied on the output layer to reduce the probability of overfitting and improve generalization.

2 Reproduction

2.1 Data Method

2.1.1 Dataset and Featurization

44 The reproduction pipeline uses the same dataset and featurization method, results in training 4D
 45 tensor, shaped **(48520, 20, 20, 20, 28)**, testing 4D tensor, shaped **(285, 20, 20, 20, 28)**.

46 2.1.2 Data Storage

47 It is worth noting that the original Sfcnn data storage uses the format of .pkl (pickle file), which
 48 features concatenate the full arrays first, then dump into the file at once. This approach requires
 49 to store and dispatch all the complexes' information within local memory, which would cause an
 50 extremely high memory consumption due to the high training data volume and is unfeasible on
 51 normal computers.

52 In alternative, our team switched to the format **.h5 (h5py file)**, which supports instant writing and
 53 solves the issue, resulting in 40.1 GiB training grid.

54 2.2 Network

55 2.2.1 Structure

56 The pytorch network structure is similar to the original tensorflow version except for two main
 57 difference:

- 58 • 1. Due to the Conv3D API requirement in pytorch, the input 4D tensor shape is permuted
 59 to (batch_size, 28, 20, 20, 20).
- 60 • 2. Pytorch lacks direct L2 regularization API, the final linear layer in the fully connected
 61 part is therefore set a weight decay to imitate the effect.

62 2.2.2 Training

63 The training process is performed on training set and validation set, the validation set is partitioned
 64 from the training 4D tensor, indexed from 41000 to 48520, same as the original network. The final
 65 training set shape: **(41000, 20, 20, 20, 28)**, validation set shape: **(7520, 20, 20, 20, 28)**, the final
 66 dataset ratio is train : validation : test = **84.00% : 15.42% : 0.58%**

67 Notice that the original training hyperparameters failed to converge in our experiments on the pytorch
 68 network, both the original hyperparameter and our current hyperparameter choice are presented in
 69 the following table:

70 Table 1: Original/Reproduced hyperparams

71 Param	Original	Reproduced
72 lr(learning rate)	0.004	0.0015
73 batch size	64	32
74 dropout rate	0.5	0.15
75 L2 regularization/FC weight decay	0.01	0.01
76 epochs	200	200

77 2.3 Results

78 2.3.1 Metrics

79 The performance of Sfcnn is measured by the following four main metrics:

$$\text{RMSE} = \sqrt{\frac{1}{N} \cdot \sum_{i=1}^N (y_{\text{predict}} - y_{\text{true}})^2} \quad (1)$$

$$\text{MAE} = \frac{1}{N} \cdot \sum_{i=1}^N |y_{\text{predict}} - y_{\text{true}}| \quad (2)$$

$$SD = \sqrt{\frac{1}{N-1} \cdot \sum_{i=1}^N ((a \cdot y_{\text{predict}} + b) - y_{\text{true}})^2} \quad (3)$$

$$R = \frac{E[(y_{\text{predict}} - \mu_{y_{\text{predict}}})(y_{\text{true}} - \mu_{y_{\text{true}}})]}{\sigma_{y_{\text{predict}}} \sigma_{y_{\text{true}}}} \quad (4)$$

where a and b represent the slope and interception of the linear regression line of the predicted and measured values. $E[c \cdot]$ denotes the expectation. $\mu_{y_{\text{predict}}}$ is the expectation of the predicted values. $\mu_{y_{\text{true}}}$ is the expectation of the true values. $\sigma_{y_{\text{predict}}}$ is the standard deviation of the predicted values. $\sigma_{y_{\text{true}}}$ is the standard deviation of the true values.

The result is shown in the following table:

Table 2: Performance Metrics Comparison on CASF-2016 Core Set

Metrics	Reproduced Sfcnn	Original Sfcnn
Pearson R	0.7286	0.7928
RMSE	1.5481	1.3263
MAE	1.2579	1.0277
SD	1.4892	1.3252

Notice that despite the original sfcnn presents better score in all the metrics, its performance is doubtful since its reproduced training process did not reach convergence.

Due to the data storage mentioned above and the author's failure to respond the request raised by another individual of providing the original (.pkl) training set on github², **the original training process is irreproducible.**

The training curves are presented below as comparison.

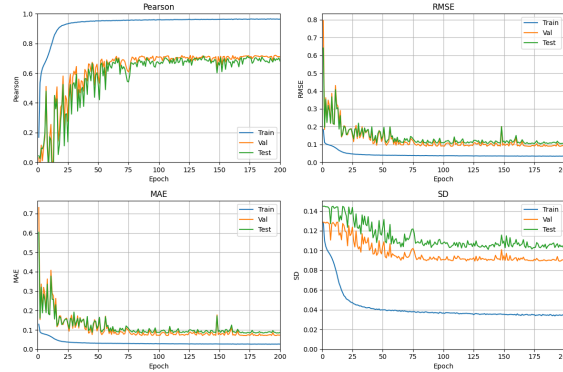


Figure 3: Training process for Reproduced Parameters

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

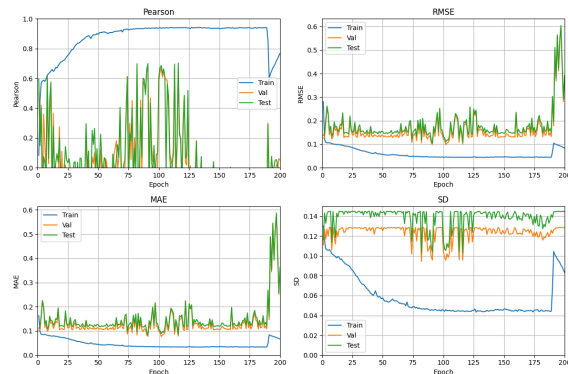


Figure 4: Training process for Original Parameters

After discussion, our team will take the **reproduced(convergent) result** as the normal performance of the Sfenn network and will apply it to the next part of AlphaFold3 result assessing.

3 AlphaFold3 Predictions

3.1 Dataset

The assessment dataset used is the CASF-2016 core set mentioned above except for 6 specific proteins with overly complex structure for AlphaFold3 to make useful predictions, resulting in total **279** proteins.

Proteins with more than 5 Isomorphic/Heterogeneous Chains are deemed too complex and excluded, those proteins are listed below:

Table 3: 6 complex protein structures

PDB ID	Isomorphic/Heterogeneous Chains
2xb8	12
2ymd	10
3n76	12
3n7a	12
3n86	12
4ciw	12

3.2 Pipeline

3.2.1 Online Server

Each protein structure is generated manually on the **Chai-1 online server**³ instead of the AlphaFold3 online server⁴ because it doesn't allow specific ligand SMILES(Simplified Molecular Input Line Entry System) code. The MSA(Multiple Sequence Alignment) option is selected with algorithm **MMseqs2** for each generation.

3.2.2 Docking

Results generated on the server will be downloaded as zip files, each contains multiple scoring ranks and detailed metrics. Structure file with the highest rank (pred.rank_0.cif) will be used as the model result for assessment.

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

⁴<https://alphafoldserver.com/>

126 To avoid the potential issues in converting .cif files⁵ to .pdb⁶ and .mol2⁷ files, the structure files
127 are parsed using the **MMCIFParser** provided in python library Bio.PDB, then go through the
128 featurization and grid mapping process directly.

129 Notice that in the results of AF3, certain atoms or isotopes are not included in the 14 pre-set atom
130 types, those atoms will be included in the **other** part of the pre-set types.

131 **3.2.3 Scoring**

132 The testing grid for predicted structures are scored using the reproduced network, loaded with the
133 pre-trained weight which shows the above performance(pearson 0.728). Detailed analysis of the
134 PLA result and metrics will be analyzed in the following section.

135 **References**

⁵Crystallographic Information File

⁶Protein Data Bank

⁷Tripos molecule structure format