

1 We thank the reviewers for their detailed and thoughtful feedback. We respond to each reviewer individually below.

## 2 **Reviewer 1**

- 3 -Our model's computational requirements scale linearly with the size of the dataset. We have recently trained the model  
4 on much larger datasets using distributed data-level parallelism on a GPU cluster. Given a sufficient number of GPUs, it  
5 can be trained quickly on an arbitrarily large dataset.
- 6 -As suggested, we will revise the drawings in Fig. 2E to clarify that our representation is truly 3D, not a 2D projection.
- 7 -Regarding Table 3 labels: DB4 is a subset of DB5, and DB4 is the training set in Table 2 as well (hence certain  
8 numbers are the same). We will clarify this point, perhaps by dropping references to DB4 and simply referring to  
9 train/validation/test subsets of DB5.
- 10 -We only use CAUROC for the test set. This is computed per training run, and we show the mean and standard deviation  
11 of CAUROC across training seeds. So it is a mean (across training seeds) of the median (across complexes) of AUROC.  
12 We will clarify this in the Table 2 and Fig. 3 captions.
- 13 -Regarding the plot not reaching 0.9 in Figure 3: there are two factors at play here. First, we are plotting mean of  
14 CAUROC across 5 training seeds without showing the individual values. When we select the best one by validation  
15 loss, we can consistently select a higher CAUROC than the mean value. Second, our final reported results come from a  
16 dataset size of 163840, which increases performance over this plot (though we do not currently have experiments for all  
17 grid sizes at that dataset size). We will clarify this by plotting all the individual training seed CAUROC as opposed to  
18 the standard deviation, and we will generate this plot for dataset size 163840.

## 19 **Reviewer 2**

- 20 -Machine learning applied to molecular structures (in drug design, materials science, and structural biology) is a rapidly  
21 growing field of interest to machine learning practitioners—in terms of techniques, datasets, and properties of atomic  
22 data. ImageNet has had a major impact in machine learning as an effective dataset against which to pretrain for vision  
23 tasks. Likewise, our work shows how one can achieve high performance on molecular learning tasks using large sets of  
24 related data and provides a dataset for doing so.
- 25 -The use of learned features for protein interface prediction is novel. [23] feeds high-level, hand-defined features into  
26 their graph convolution architecture.
- 27 -Regarding learned flexibility: We show that unlike previous methods, our method performs well on cases involving  
28 conformational change even when the training dataset includes no conformational change at all. The training datasets  
29 for previous methods included cases that exhibit substantial conformational change. We will clarify this point.
- 30 -Regarding formatting/citations: We will make the suggested changes.
- 31 -We use a 3D CNN because graph convolutions do not directly model complex three-dimensional geometric arrange-  
32 ments, and proteins are complex three-dimensional structures.
- 33 -We will clarify that Fig. 3B shows that our model—unlike competing methods—is capable of leveraging additional  
34 data to increase its performance and achieve state-of-the-art results. We will avoid the term “scalability” in this context.

## 35 **Reviewer 3**

- 36 -Regarding dataset availability: We have already uploaded a preliminary version of the dataset on Harvard Dataverse  
37 (the subset we trained on) and will upload the full dataset.
- 38 -Regarding pruning only against the test set and smoothing across voxels: These are good points and should only  
39 improve performance.
- 40 -Regarding analyzing why other methods fail to improve performance using the larger dataset: We agree this would be  
41 informative. For example, we can determine which high-level features don't translate well to DIPS.
- 42 -Regarding information on voxel size and number: Some of this is stated in the Methods section (voxels are 1 Å in size  
43 and we use 35x35x35 voxels total per surfacelet), but we will add to and highlight this information.
- 44 -Our model could indeed be used to predict regions within the same protein that interact with one another, and we agree  
45 that this would result in significantly more training data. The key question is how to exclude parts of one surfacelet  
46 from the other so signal does not leak about each interacting component. We will point this out.
- 47 -We will better explain the experimental setup of Fig. 3A. Grid size refers to the entire window, and voxel size is kept  
48 constant. If the voxel size were allowed to vary, then the learned filters would have to encode information at different  
49 scales, which might tax the model's capacity.
- 50 -Our model scales linearly with the amount of data provided, and we have recently trained it on much larger datasets  
51 using distributed data-level parallelism on a GPU cluster. The model should indeed be useful in many molecular  
52 prediction tasks, as the reviewer notes.