

1 We thank the reviewers for their positive feedback and helpful suggestions for improvement.

2 **Response to Reviewer 1** Both Bepler et al. [11] and Alley et al. [12] were pretrained on the PFAM cor-
3 pus. We have updated our table to make this clearer. We have expanded the discussion of alignment-based
4 features in the Background section to make it clearer what they are. Since alignment-based features are fed
5 directly into the downstream task models, this is separate from pretraining on the language modeling tasks.

6 The reviewer points out the unexpected result that pretraining decreased contact
7 prediction performance with the ResNet. This was caused by an issue with our
8 training/evaluation setup for contact map prediction, which we discovered after
9 the submission deadline. When we re-ran training/evaluation for all models, we
10 found that pretraining improved accuracy across the board, as shown in Table 1.

Method	Contact
Transformer (No PT)	0.32
LSTM (No PT)	0.19
ResNet (No PT)	0.20
Transformer (PT)	0.36
LSTM (PT)	0.39
ResNet (PT)	0.29
Bepler (PT)	0.40
UniRep (PT)	0.34

Table 1: Updated Contact Results

11 **Response to Reviewer 2** When splitting PFAM, we split at both the family and
12 clan level. We will present these results separately in Table 1 of the paper. We
13 ran BLAST to obtain sequence identity between these sets. We find that that on
14 average the maximum sequence identity between the training set and proteins in
15 the random-split set is 73%, in the family-holdout is 37%, and in the clan-holdout
16 is 22%. Hamming distance is measured at the amino acid level, which we chose
17 because we are trying to highlight the protein engineering setting where the goal is to accurately predict the function of
18 an input molecule specified by a protein designer. We have added this reasoning to the task description.

19 The reviewer is correct that the LSTM language model is not making a Markovian assumption, and we have removed
20 this line. The reviewer is correct in their characterization of the LSTM, which is the concatenation of the two vectors
21 of dimension 1024. For the single-prediction tasks, we use a learned attention-weighted sum over the amino-acid
22 representations. These points are in the Appendix, and we have added pointers to these details in the main text.

23 We agree that difference in representation size could make a difference on downstream tasks. Due to cost of pretraining
24 large models, we cannot perform a grid search over number of parameters and representation size for pretrained
25 models. As such, we instead chose to base models on common hyperparameter choices in the literature. In addition,
26 controlling for number of parameters and representation size proved difficult without resorting to extremely odd choices
27 or significant alterations to the network design. For example, lowering the LSTM hidden dimension to 512 requires 10
28 layers to match the parameters in the Transformer. To address these concerns, we are pretraining three smaller models
29 with a fixed representation size to compare architectures with representation size held constant (see Table 2).

30 **Response to Reviewer 3** We agree with the reviewer that the language models we trained are quite large, which
31 may present a problem for fine-tuning on small datasets. Our tasks are of similar size to most of those in the GLUE
32 benchmark, which has been instrumental in demonstrating the success of self-supervision in NLP. Since the models that
33 were applied to GLUE have tens to hundreds of millions of parameters, we chose to make our models roughly the same
34 size (~40M parameters). We also believe that showing results of larger models is important to counter the claim that
35 simply scaling these models would allow them to outperform gold-standard features from bioinformatics. We have
36 added this discussion in the manuscript. To further address questions about size, we are pretraining smaller language
37 models with ~3M parameters for a week on Pfam. Some preliminary results on the Transformer and LSTM after two
38 days of pretraining are reported in Table 2.

Method	SS	Stab	Fluor	RH
LSTM	0.73	0.67	0.66	0.18
Transformer	0.70	0.68	0.68	0.13
ResNet	0.73	0.68	0.43	0.11

Table 2: Small Language Model Results

We agree with the reviewer that usability of the benchmark is of
paramount concern. The larger models require around 11 GB of
memory on a GPU to train, which is on the higher end. To alleviate
this, we will upload weights for the smaller models once they are
trained.

44 In Table 2 of the paper we chose to report simple global metrics.
45 However, we report many detailed metrics in the appendix, including some of the metrics the reviewer suggested:
46 precision-recall, AUPRC, and accuracy for long-range contacts. For the other tasks we provided label/task-specific
47 breakdowns in the appendix, and have added more pointers to the appendix to improve the visibility of these results.
48 We have also revised the discussion in the appendix for ease of reading. We also plan to create an online leaderboard,
49 where these other metrics will be more prominent. We agree with the reviewer that CB513 is built from an old dataset
50 and will look into improving the test set for future iterations of TAPE.

51 **Clarity** We thank all reviewers for positive comments on clarity. All small edits and suggested improvements have
52 been incorporated. In particular we are grateful for Reviewer 2’s catch of the missing beta strand contacts (the contact
53 map was cropped incorrectly, and has been updated). Additionally, as per Reviewer 3’s suggestion, we have changed
54 ‘Protein Modeling’ to ‘Machine learning applied to protein sequences’ in the abstract and throughout the manuscript.