

Learning co-evolution information with natural language processing for protein folding problem*

Egor Zverev¹, Sergei Grudinin², and Ilia Igashov^{1,2}

zverev.eo@phystech.edu; sergei.grudinin@inria.fr; igashov.is@phystech.edu

¹Organization, address; ²Organization, address

Co-evolution information is crucial for building protein sequence descriptors. Its computation is traditionally based on Multiple Sequence Alignments. Applications of MSA are limited, for instance, it fails for sequences with shallow alignment. This work investigates pre-trained language models as potential alternatives to MSA. The main focus is on fold classification problem. The authors examine several state-of-the-art methods and modify them by replacing MSA-based feature-generation parts with pre-trained LMs. They compare language models and MSA in their ability to extract sequence evolution information. The performance of models is compared using TOP1-accuracy score.

Keywords: *protein fold classification; co-evolution, feature generation; transformers; BERT;*

DOI: 10.21469/22233792

1 Introduction

Protein properties are determined by its shape, which is specified by its amino acid sequence [14]. Therefore it is important to learn how to analyze this sequence. The first step is constructing protein sequence descriptors. Ideally the descriptors contain co-evolution information[...]. This information is obtained by searching for homologues of the protein in large databases and computing multiple sequence alignment (MSA) on it. This model requires significant computational effort. Alignment databases are finite. If for a given sequence there are no matches in the database, MSA-based methods fail to produce reliable results. It does not guarantee precise results since it fully relies on finite databases. MSA-based methods fail for sequences with shallow alignment.

Co-evolution-based methods assume[...] that a mutation of one amino acid leads to mutations of others. Therefore, given an amino acid sequence it is natural to look for other similar sequences. That is performed with MSA [5]. This method fails when we are dealing with the sequences that did not evolve a lot. They have shallow alignments. That means there are only a few similar sequences in the database. Therefore, information extracted from these alignments is insignificant.

The authors propose to analyse pre-trained language models (LMs) [9] as potential alternative to traditional MSA-based models. It is assumed that protein sequences are not random [4], There is structure in amino acid sequences. Therefore, a set of all amino acids could be seen as a language with complicated inner rules. BERT [6] learns the structure of any abstract language. We assume that by learning amino acids structure, BERT will be able to implicitly learn co-evolution information.

Our main aim is to study the efficiency of pre-trained LMs in application to the fold classification problem. To start with, we study the state-of-the-art method DeepSF [10]. It solves the protein fold classification problem. By replacing its MSA-based feature-generation part with a pre-trained LM, we study how NLP approach affects learning the fold-related information. The whole framework is schematically represented in Figure 1. Further, we perform the same with other fold classification algorithms [7, 15] as well as with end-to-end protein structure prediction methods [11, 16].

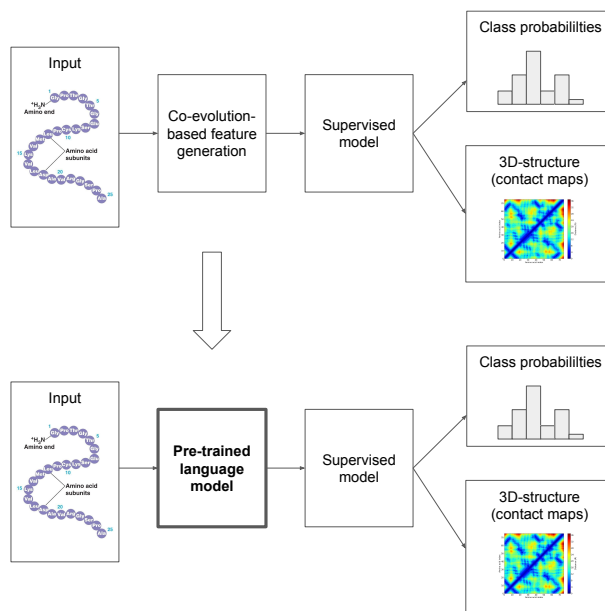


Figure 1 The main framework

2 Problem statement

Proteins that share similarities in their structure are divided into classes – folds [13]. In this work we investigate a fold classification problem.

2.1 General description

We denote Y a set of all known classes, A a set of all proteins with known folds and $f : A \rightarrow Y$ a mapping between proteins and their classes.

Suppose, $x \notin A$ is a protein with unknown fold. We aim to build a model capable of classifying x as an element of one of the present classes, extending f on the set $A \cup \{x\}$.

2.2 Feature generation

In the fold classification problem x_i represents i -th protein description in D -dimensional space. However, initially each protein is described by its amino acid sequence. Let Σ be the 20-letter alphabet of acid sequences.

There is a set $A = \{a_i \in \Sigma^+\}_{i=1}^n$ of protein sequences. Let

$$g_{\theta,D} : \Sigma^+ \rightarrow \mathbb{R}^D$$

be an embedding of space of acid sequences into D -dimensional set of descriptors. These embeddings are parameterized by θ .

We separate A into training and validation set. Training set is used to formulate fold classification problem. Validation set is used to formulate feature generation problem.

$$A = A_{\text{train}} \cup A_{\text{val}}$$

Protein sequence descriptors are generated using $g_{\theta,D}$ as $x_i = g_{\theta,D}(a_i)$

2.3 Fold classification problem

Suppose there is a given set of pairs $S = \{(x_i, Y_i)\}_{i=1}^n$ where $x_i \in \mathbb{R}^D$ denotes the sequence descriptor of i -th protein and

$$Y_i \in \mathbb{Y} = \{(1, 0, \dots, 0), (0, 1, 0, \dots, 0) \dots (0, \dots, 0, 1)\}$$

denotes the i -th protein known class, $|\mathbb{Y}| = m$

We separate S into training and validation set. Training set is used to formulate fold classification problem.

$$S = S_{\text{train}} \cup S_{\text{val}}$$

Let $W = \{(w_1, \dots, w_m) | w_i \in \mathbb{R}^D\}$ be a space of parameters for classification models. Let g_w be a model parameterized by w . We use e_k to denote a vector with 1 on k -th position, with zeros on all other positions

$$e_k = \begin{matrix} [0, \dots, 0, & 1, & 0, \dots, 0]. \\ \uparrow \\ k \end{matrix}$$

55 We work under the following assumption:

$$P_w(Y = e_k|x) = \frac{\exp(x^T w_k)}{\sum_j \exp(x^T w_j)}$$

56 Let $p_{w,x}(y)$ be discrete density of Y , $L_{Y,x}(w)$ be likelihood function of Y .

$$p_{w,x}(y) = \prod_{d=1}^m (P_w(Y = e_d|x))^{y_d}$$

57 ,

$$L_{Y,x}(w) = \prod_{i=1}^n p_{w,x}(Y_i)$$

58 Given S_{train} , classification problem is a maximization of likelihood on training set:

$$L_{train}(w) = L_{Y_{train}, x_{train}}(w) \rightarrow \max_{w \in W}$$

59 2.4 Feature generation problem

For fixed θ, D and descriptors x_i we denote

$$\hat{w} = \arg \max_{w \in W} L_{Y_{train}, x_{train}}(w).$$

60 .

61 Feature generation problem is a maximization of likelihood function on validation set:

$$L_{Y_{val}, x_{val}}(\hat{w}) \rightarrow \max_{\theta, D}$$

62 2.5 Overview

63 The problem of determining proteins' class is separated into two parts: first, we encode
64 amino acid sequences as D -dimensional descriptors, then these descriptors are used to train a
65 prediction model.

66 3 Computational experiment

67 The goal of our experiments is to verify that BERT is capable of extracting important
68 features from the given amino acid sequences. For this we will use pre-trained AlBert [8] model
69 to generate these features and then train our own neural networks to solve fold classification
70 problem.

71 3.1 Data

72 We use SCOP2 [1,2] dataset. It contains information about proteins whose structure is already
73 known. All the proteins in the database have classes assigned to them (folds). At present,

there are 1517 known folds [2]. The number of proteins contained in each fold is less than 50 for most of the folds. Detailed histogram is represented in Figures 2, 3.

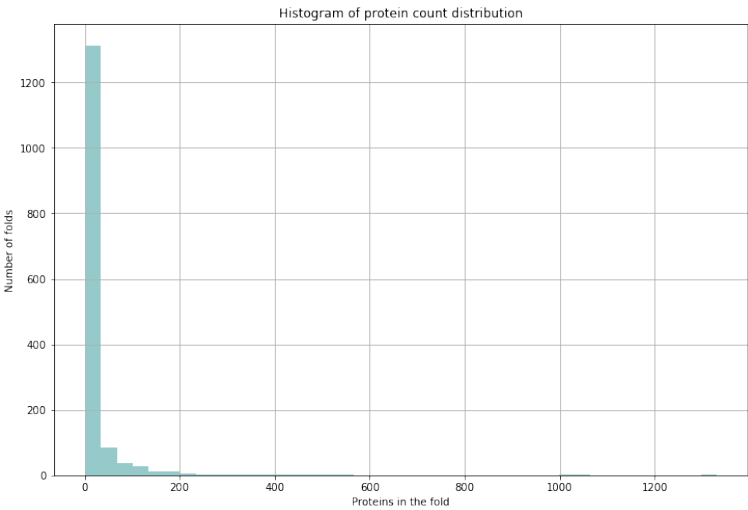


Figure 2 Protein count histogram

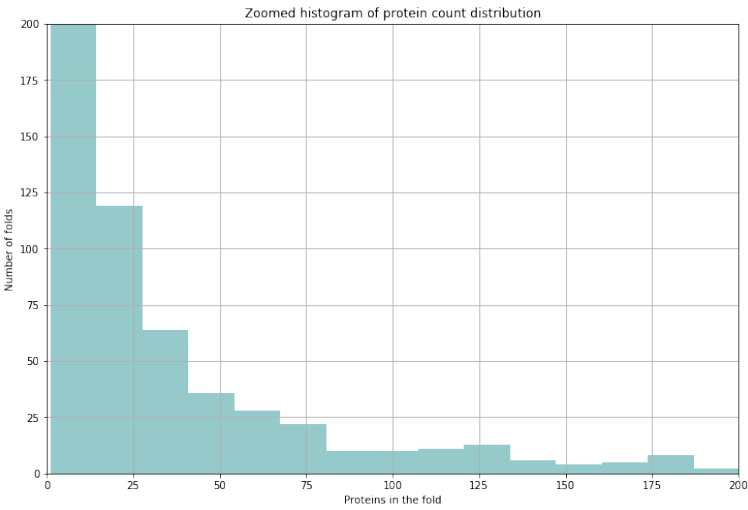


Figure 3 Protein count zoomed histogram

In SCOP2 proteins are represented by their amino acid sequences. Each sequence is stored as a string in the 20 letter alphabet and an ID. The length of the string varies mostly between 25 and 500 from protein to another. The length diagram is represented on Figure 2. Search by ID allows database users to extract information about target proteins, including their fold classes.

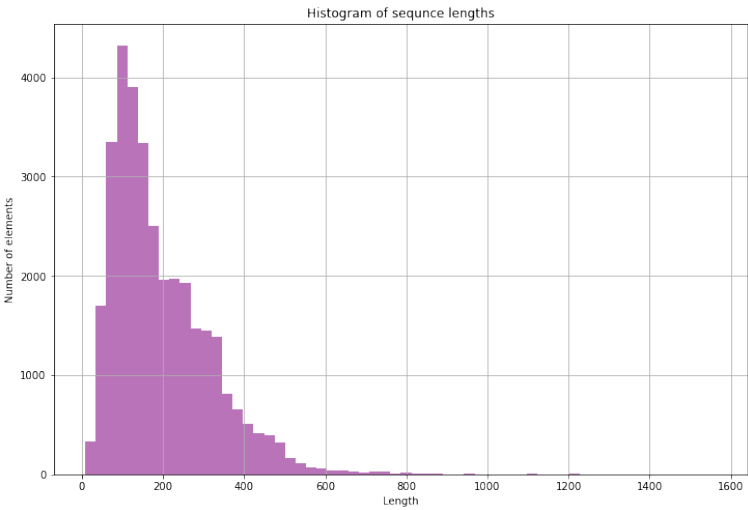


Figure 4 Sequence length histogram

3.2 Baseline solution

Our solution is a modification of the DeepSF [10] method. The designers of DeepSF use PSI-BLAST [3] to generate features from MSA as well as SCRATCH [12] to extract secondary structure (3 classes) and solvent accessibility (2 classes). In the core of this classification method PSI-BLAST is used during the first stage of prediction [12].

Hoe et al. then use these features as an input to the deep convolutional neural network. The depth of CNN used in DeepSF is 10. The output of the model is a vector of fold probabilities. Fold with the maximal probability is accepted as a target protein fold.

3.3 Proposed solution

DeepSF feature generation is entirely based on MSA. The main point of this work is to suggest another approach to feature generation which is independent of MSA.

We apply BERT to the language of amino acids. We take the representation of acid sequences BERT creates as a new set of features. Then we train DeepSF model on these features.

It shall be noted that we do not reject co-evolution information in our research. Instead of directly applying MSA we let BERT learn information about amino acid language. We suppose that the model learns evolution information implicitly.

3.4 Evaluation

We use standard cross-validation procedure to compute accuracy of our model. We then compare it to the accuracy score of DeepSF method based on MSA feature generation as well as several other state-of-the-art techniques.

4 Experiment details

Though the data is available for 1517 classes, during the experiments we approached the problem slowly. First, we solved a classification problem for 2 classes, then for 10 classes.

4.1 First experiment

In this experiment we took two balanced classes. They contain 66 and 80 elements respectively. The distributions of their sequence length are similar. They are displayed on figure 5. We used a model that consists of two essential layers - 30-max-pooling followed by a linear layer. Using this architecture we achieved 93% score on validation data. The architecture and training curves are represented on figures 6 and 7.

The architecture presented here is the simplest possible one - it contains max-pooling layer to transform input to fixed size and FC layer for classification. Motivation for this experiment is to verify whether features extracted by BERT are meaningful.

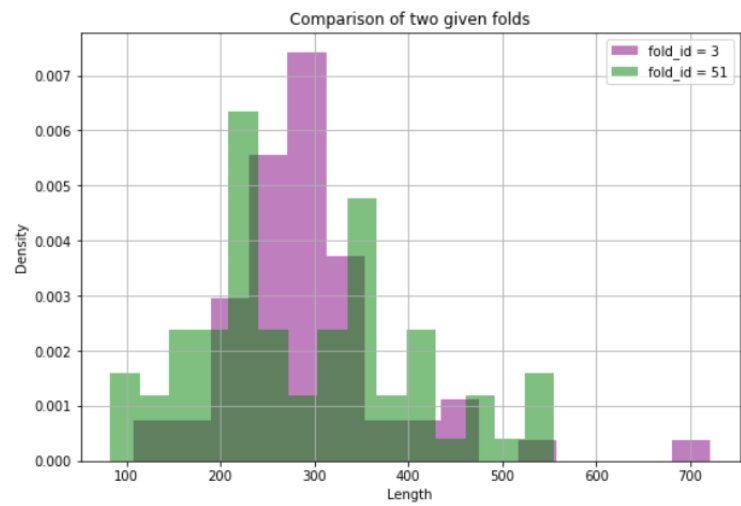


Figure 5 Sequence lengths of two folds

Layer (type:depth-idx)	Output Shape	Param #
K_max_pooling_1d: 1-1	[32, 4096, 30]	--
ReLU: 1-2	[32, 4096, 30]	--
Flatten: 1-3	[32, 122880]	--
Linear: 1-4	[32, 2]	245,762

Total params: 245,762
Trainable params: 245,762
Non-trainable params: 0
Total mult-adds (M): 7.86

Input size (MB): 157.29
Forward/backward pass size (MB): 0.00
Params size (MB): 0.98
Estimated Total Size (MB): 158.27

Figure 6 Experiment 1 architecture

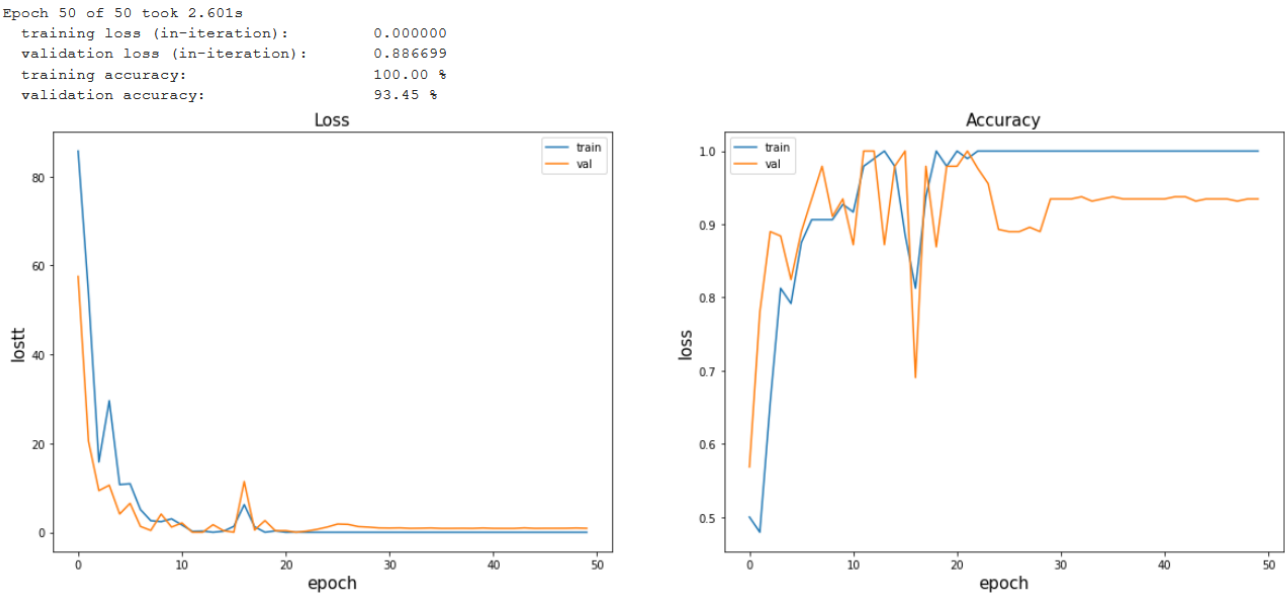


Figure 7 Experiment 1 training curve

4.2 Second experiment

We decided to use exactly the same architecture for 10-fold classification problem. Our experiments show that increasing the k hyper-parameter in k -max-pooling layer improves the quality of the model. We chose $k = 60$. It allowed us to reach 95% accuracy on validation and 93% accuracy on test samples. The training curve is represented on figure 8.

The second experiment is an extension of the first one. Its goal is to validate that simple 2-layer architecture could deal with more complicated classification problem.

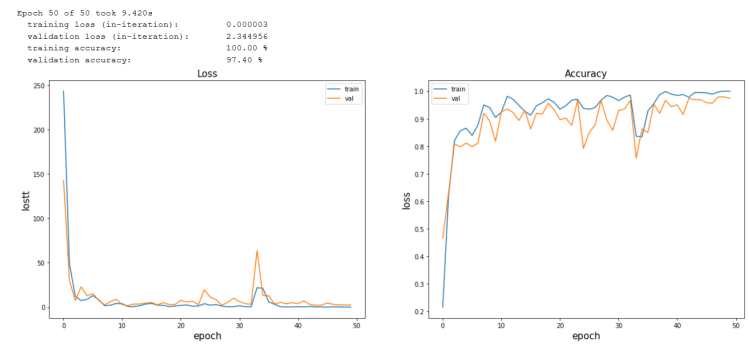


Figure 8 Experiment 2 training curve

5 Feature work

In this paper we investigated only classification problem for at most classes. In the future we would like to extend our experiments for 1517-fold classification problem and compare the results.

Here we used BERT as a "black box". We generated features with BERT and used them as an input for another model. The desired improvement here would be to learn how to extract attention mechanisms from BERT, possibly gaining more co-evolution information in the process.

6 Results

Conducted experiments show that BERT successfully captured the structure of the protein sequence language (or at least its part). We managed to solve 10-fold classification problem with 93% accuracy.

Therefore, BERT is capable of extracting crucial information from the sequences of amino acids.

References

- [1] Antonina Andreeva, Dave Howorth, Cyrus Chothia, Eugene Kulesha, and Alexey G. Murzin. SCOP2 prototype: a new approach to protein structure mining. *Nucleic Acids Research*, 42(D1):D310–D314, 11 2013.
- [2] Antonina Andreeva, Eugene Kulesha, Julian Gough, and Alexey G Murzin. The SCOP database in 2020: expanded classification of representative family and superfamily domains of known protein structures. *Nucleic Acids Research*, 48(D1):D376–D382, 11 2019.
- [3] Medha Bhagwat and L. Aravind. Psi-blast tutorial. *Methods in molecular biology (Clifton, N.J.)*, 395:177–186, 2007. 17993673[pmid].
- [4] Davide De Luca, Debora Slanzi, Irene Poli, Fabio Polticelli, and Giovanni Minervini. Do natural proteins differ from random sequences polypeptides? natural vs. random proteins classification using an evolutionary neural network. *PLOS ONE*, 7(5):1–10, 05 2012.
- [5] S. de Oliveira and C. Deane. Co-evolution techniques are reshaping the way we do structural bioinformatics. *F1000Res*, 6:1224, 2017.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [7] W. Elhefnawy, M. Li, J. Wang, and Y. Li. DeepFrag-k: a fragment-based deep learning approach for protein fold recognition. *BMC Bioinformatics*, 21(Suppl 6):203, Nov 2020.
- [8] Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghaliya Rehawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, DEBSINDHU BHOWMIK, and Burkhard Rost. Prottrans: Towards cracking the language of life’s code through self-supervised deep learning and high performance computing. *bioRxiv*, 2020.
- [9] Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghaliya Rihawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Debsindhu Bhowmik, and Burkhard Rost. Prottrans: Towards cracking the language of life’s code through self-supervised deep learning and high performance computing. *bioRxiv*, 2020.
- [10] Jie Hou, Badri Adhikari, and Jianlin Cheng. DeepSF: deep convolutional neural network for mapping protein sequences to folds. *Bioinformatics*, 34(8):1295–1303, 12 2017.
- [11] Shaun M Kandathil, Joe G Greener, Andy M Lau, and David T Jones. Deep learning-based prediction of protein structure using learned representations of multiple sequence alignments. *bioRxiv*, 2020.
- [12] C. N. Magnan and P. Baldi. SSpro/ACCpro 5: almost perfect prediction of protein secondary structure and relative solvent accessibility using profiles, machine learning and structural similarity. *Bioinformatics*, 30(18):2592–2597, Sep 2014.

- [13] R. Dustin Schaeffer and Valerie Daggett. Protein folds and protein folding. *Protein engineering, design & selection : PEDS*, 24(1-2):11–19, Jan 2011. 21051320[pmid].
- [14] Akif Uzman. Molecular biology of the cell (4th ed.): Alberts, b., johnson, a., lewis, j., raff, m., roberts, k., and walter, p. *Biochemistry and Molecular Biology Education*, 31(4):212–214, 2003.
- [15] A. Villegas-Morcillo, A. M. Gomez, J. A. Morales Cordovilla, and V. E. Sanchez Calle. Protein fold recognition from sequences using convolutional and recurrent neural networks. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, pages 1–1, 2020.
- [16] Jinbo Xu, Matthew Mcpartlon, and Jin Li. Improved protein structure prediction by deep learning irrespective of co-evolution information. *bioRxiv*, 2020.

Received January 01, 2017