

Title

Author

Date

1 Introduction

2 Related Work

3 Methods

In this section we describe the pipeline used to analyze the embeddings. As shown in Table 1, the input length is different between the models as well as the output produced. We want to address the following problems: 1) compare different methods to join together the aminoacid-specific contextual representations in order to have a representation for the whole chunk; 2) compare different methods to join together the representations of the chunks in order to have a representation for the whole protein; 3) find out if these representations reflect known properties of the proteins.

3.1 Single aminoacid contextual representation

We tried four methods to join together the aminoacid embeddings in order to produce a fixed size embedding for the chunk: aver-

age, maximum, sum and principal component analysis (PCA). Note that even if these operators are commutative, the overall process does take into account the order of the aminoacids precisely because the embeddings are contextual.

The same operator used to combine the aminoacid embeddings is also used to combine the embeddings of the chunks of the sequence.

Name	input length (chunk)	embedding dimension
embedding reproduction (rep)[5]	64	64 per chunk
dnabert [2]	512	768 per chunk
prose [1]	512	100 per amino acid
alphafold [3]	1024	384 per ammino acid
evolutionary scale modeling (esm2) [4]	1024	1280 per ammino acid

Table 1: Embedders used in the experiments, their maximum input length and the dimension of the embedding produced.

4 Results

References

- [1] Tristan Bepler and Bonnie Berger. Learning the protein language: Evolution, structure, and function. *Cell systems*, 12(6):654–669, 2021.
- [2] Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davuluri. Dnabert: pre-trained bidirectional encoder representations from transformers model for dna-language in genome. *Bioinformatics*, 37(15):2112–2120, 2021.
- [3] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [4] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022.
- [5] Kevin K Yang, Zachary Wu, Claire N Bedbrook, and Frances H Arnold. Learned protein embeddings for machine learning. *Bioinformatics*, 34(15):2642–2648, 2018.