

# Contrastive learning on protein embeddings enlightens midnight zone

Michael Heinzinger<sup>1,2</sup>, Maria Littmann<sup>1</sup>, Ian Sillitoe<sup>3</sup>, Nicola Bordin<sup>3</sup>, Christine Orengo<sup>3</sup> & Burkhard Rost<sup>1,4</sup>

1 TUM (Technical University of Munich) Dept Informatics, Bioinformatics & Computational Biology - i12, Boltzmannstr. 3, 85748 Garching/Munich, Germany

2 TUM Graduate School, Center of Doctoral Studies in Informatics and its Applications (CeDoSIA), Boltzmannstr. 11, 85748 Garching, Germany

3 Institute of Structural and Molecular Biology, University College London, London WC1E 6BT, UK

4 Institute for Advanced Study (TUM-IAS), Lichtenbergstr. 2a, 85748 Garching, Germany & TUM School of Life Sciences Weihenstephan (WZW), Alte Akademie 8, Freising

\* Corresponding author: mheinzinger@rostlab.org, <http://www.rostlab.org/> Tel: +49-289-17-811 (email rost: assistant@rostlab.org)

## Abstract

Experimental structures are leveraged through multiple sequence alignments, or more generally through homology-based inference (HBI), facilitating the transfer of information from a protein with known annotation to a query without any annotation. A recent alternative expands the concept of HBI from sequence-distance lookup to embedding-based annotation transfer (EAT). These embeddings are derived from protein Language Models (pLMs). Here, we introduce using single protein representations from pLMs for contrastive learning. This learning procedure creates a new set of embeddings that optimizes constraints captured by hierarchical classifications of protein 3D structures defined by the CATH resource. The approach, dubbed *ProtTucker*, has an improved ability to recognize distant homologous relationships than more traditional techniques such as threading or fold recognition. Thus, these embeddings have allowed sequence comparison to step into the “midnight zone” of protein similarity, i.e., the region in which distantly related sequences have a seemingly random pairwise sequence similarity. The novelty of this work is in the particular combination of tools and sampling techniques that ascertained good performance comparable or better to existing state-of-the-art sequence comparison methods. Additionally, since this method does not need to generate alignments it is also orders of magnitudes faster. The code is available at <https://github.com/Rostlab/EAT>.

**Key words:** protein structure classification, protein language model, self-supervised learning, contrastive learning, CATH.

**Abbreviations used:** **3D**, three-dimensional; **BFD**, Big Fantastic Database (11); **CATH**, hierarchical classification of protein 3D structures in Class, Architecture, Topology and Homologous superfamily (1,2); **DL**, Deep Learning; **EAT**, Embedding-based Annotation Transfer; **EI**, evolutionary information; **embeddings** fixed-size vectors derived from pre-trained pLMs; **ESM-1b**, pLM from Facebook dubbed Evolutionary Scale Modeling (12); **FNN**, Feed-forward Neural Network; **FunFams**, functional families as sub-classification of the most fine-grained H level in CATH (13); **HBI**, Homology Based Inference; **HMM**, Hidden Markov Model; **HMMer**, particular method for HMM-profile alignments (6); **HSSP**, homology-derived secondary structure of proteins (14); **HVAL**, distance from empirical curve separating proteins with similar structure recognizable from pairwise alignments (15); **LM**, Language Model; **MMseqs2**, fast database search and multiple sequence alignment method (10); **MSA**, Multiple Sequence Alignment; **NLP**, Natural Language Processing; **PDB**, Protein Data Bank; **PIDE**, percentage pairwise sequence identity; **pLM**, protein Language Model; **ProSE**, pLM based on long short-term memory (LSTM) cells dubbed Protein Sequence Embeddings (16); **ProtBERT**, pLM based on the LM BERT (17); **ProtT5**, pLM based on the LM T5 (18).

## Introduction

**Phase-transition from daylight through twilight into midnight zone.** Protein sequence determines structure which determines function. This simple chain underlies the success of grouping proteins into families from sequence (3,19-21). Information from experimental high-resolution three-dimensional (3D) structures expands the perspective from families to super-families (1,22) that often reveal evolutionary and functional connections not recognizable from sequence alone (23,24). Thus, 3D information helps us to penetrate through the twilight zone of sequence alignments (25,26) into the midnight zone of distant evolutionary relationships (27).

The transition from daylight, through twilight and into the midnight zone is characterized by a phase-transition, i.e., a sigmoid function describing an order of magnitude increase in recall (relations identified) at the expense of a decrease in precision (relations identified correctly) over a narrow range of sequence similarity. Measuring sequence similarity by the HSSP-value (**HVAL**) (26,28) for the *daylight zone* at  $HVAL > 5$  (>25% PIDE - pairwise sequence identity over >250 aligned residues) over 90% of all protein pairs have similar 3D structures, while at the beginning of the *midnight zone* for  $HVAL < -5$  (<15% PIDE for >250 aligned residues), over 90% have different 3D structures. Thus, the transition from daylight to midnight zone is described by a phase-transition in which over about ten percentage points in PIDE precision drops from 90% to 10%, i.e., from almost *all correct* to almost *all incorrect* within  $\pm 5$  points PIDE. The particular point at which the twilight zone begins and how extreme the transition is, depends on the

62 phenotype: steeper at lower PIDE for structure (26) and flatter  
63 at higher PIDE for function (29,30).

64 If two proteins have highly similar structures, it is still possible  
65 for their sequences to be found in this *midnight zone*, i.e., have  
66 seemingly random sequence similarity (27). Thus, if we could  
67 safely lower the threshold just a little, we would gain many  
68 annotations of structural and functional similarity. In fact, any  
69 push a little lower reveals many proteins with similar pheno-  
70 type, e.g., structure or function. Unfortunately, without impro-  
71 ving the search method, such a lowering usually comes at the  
72 expense of even more proteins with dissimilar phenotype.

73 This simple reality has been driving the advance of methods  
74 using sequence similarity to establish relations: from advanced  
75 pairwise comparisons (14,31) over sequence-profile (32-35) to  
76 profile-profile comparisons (24,36-41) or efficient shortcuts to  
77 the latter (10,42). All those methods share one simple idea,  
78 namely, to use evolutionary information (**EI**) to create families of  
79 related proteins. These are summarized in multiple sequence  
80 alignments (MSAs). Using such information as input to machine  
81 learning methods has been generating essentially all state-of-  
82 the-art (SOTA) prediction methods for almost three decades  
83 (43-45). Using MSAs has also been one major key behind the  
84 breakthrough in protein structure prediction through  
85 *AlphaFold2* (46), and subsequently of *RoseTTAFold* (47) which  
86 builds on ideas introduced by *AlphaFold2*, i.e., allowing for  
87 communication between different sequence- and structure  
88 modules within the network. Transfer- or representation-  
89 learning offer a novel route toward comparisons of and  
90 predictions for single sequences without MSAs.

## 1 Embeddings capture language of life written in proteins.

2 The introduction of LSTM- or attention-based Language  
3 Models (LMs) such as ELMo (48) or BERT (17) enabled a better  
4 use of large, unlabeled text corpora which arguably improved  
5 all tasks in natural language processing (NLP) (49). These  
6 advances have been transferred to proteins through protein  
7 Language Models (pLMs) equating amino acids with words in  
8 NLP and the sequence of entire proteins with sentences. Such  
9 pLMs learn to predict masked or missing amino acids using  
10 large databases of raw protein sequences as input (9,12,50-  
11 54), or by refining the pLM through another supervised task  
12 (16,55). Processing the information learned by the pLM, e.g.,  
13 by using the output of the last hidden layers of the networks  
14 forming the pLMs, yields a representation of protein sequences  
15 referred to as embeddings (Fig. 1 in (9)). Embeddings have  
16 been used successfully as exclusive input to predicting  
17 secondary structure and subcellular localization at  
18 performance levels almost reaching (12,50,51) or even  
19 exceeding (9,56,57) the SOTA using evolutionary information  
20 from MSAs as input. Embeddings can even substitute  
21 sequence similarity for homology-based annotation transfer  
22 (58,59). The power of such embeddings has been increasing  
23 with the advance of algorithms and the growth of data (9). The  
24 recent advances have shown that a limit to such improvements  
25 has not nearly been reached when writing this (22.02.2022).

26 Embeddings from pLMs capture a diversity of higher-level  
27 features of proteins, including various aspects of protein  
28 function and structure (9,12,51,58-61). In fact, pLMs such as  
29 ProtT5 (9) or ESM-1b (12) capture aspects about protein  
30 structure so impressively that inter-residue distances – and  
31 consequently 3D structure – can be predicted without using  
32 MSAs, even with relatively small (few free parameters) Deep  
33 Learning (DL) architectures (62).

34 Supervised learning directly maps the input to the class  
35 output. Instead, contrastive learning (63), optimizes a new  
36 embedding space in which similar samples are pushed closer,  
37 dissimilar samples farther apart. Contrastive learning relies only  
38 on the similarity between pairs (or triplets) of samples instead  
39 of on class label. The definition of similarity in embedding rather  
40 than sequence space, combined with contrastive learning,  
41 offered an alternative to sequence-based protein comparisons.  
42 This led us to hypothesize that we might find structurally and  
43 functionally consistent sub-groups within protein families from  
44 raw sequences. As a proof-of-principle, a rudimentary  
45 precursor of this work helped to cluster FunFams (2,59). The  
46 benefit of optimizing embeddings specifically for SCOPe fold  
47 recognition (64) has recently been shown (16,60,65). Other  
48 approaches toward fold recognition deep learn fold-specific  
49 motifs (66), pairwise similarity scores (67) or sequence  
50 alignments (68). However, most of the top-performing solutions  
51 rely on information extracted from MSAs (69) and do not utilize  
52 the transfer-learning capabilities offered by recent pLMs.

53 Here, we expand on the hypothesis that replacing supervised  
54 learning by contrastive learning intrinsically fits the hierarchy of  
55 CATH (1,2). We propose an approach that marries both, self-  
56 supervised pretraining and contrastive learning, by  
57 representing protein sequences as embeddings, and using  
58 increasing overlap in the CATH hierarchy as a notion of  
59 increasing structural similarity to contrastively learn a new  
60 embedding space. We used the pLM ProtT5 (9) as static feature  
61 encoder (no fine-tuning of the pLM) to retrieve initial  
62 embeddings that were then mapped by a feed-forward neural  
63 network (FNN) to a new, learned embedding space optimized  
64 on CATH through contrastive learning. More specifically, the  
65 Soft Margin Loss was used with triplets of proteins (anchor,  
66 positive, and negative) to optimize the new embedding space

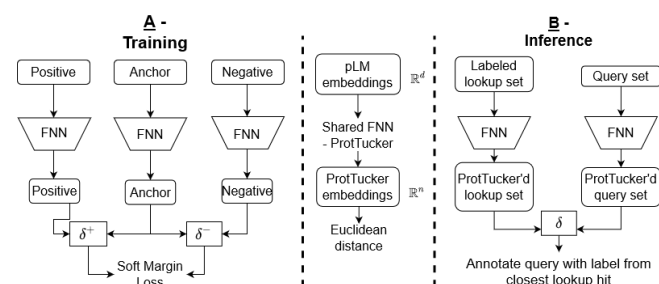


Fig. 1: Sketch of ProtTucker. Panel A illustrates how protein triplets were used to contrastively learn the CATH hierarchy (1,2). First, protein Language Models (pLMs) were used as static feature encoders to derive embeddings for protein sequences (anchor, positive, negative). The embedding of each protein was processed separately by the same, shared FNN with hard parameter sharing, called *ProtTucker*. During optimization, the Soft Margin Loss was used to maximize the distance between proteins from different CATH classes (anchor-negative pairs) while minimizing the distance between proteins in the same CATH class (anchor-positive pairs). All four CATH-levels were simultaneously learned by the same FNN. This resulted in a newly, learned CATH-optimized embedding for each protein. Panel B sketches how the contrastive learning FNN is used for prediction of new proteins (inference). For all proteins in a lookup set with experimental annotations (labeled proteins; here the CATH lookup set), as well as for a query protein without experimental annotations (unlabeled proteins) all embeddings are extracted in two steps: (1) extract per-residue embeddings from original pLM and create per-protein embeddings by averaging over protein length. (2) Input those embeddings into the pre-trained FNNs, i.e., *ProtTucker*. Similar to homology-based inference (HBI), predictions are generated by transferring the annotation of the closest hit from the lookup set to the query protein. The embedding-based annotation transfer (EAT) transferred annotations to the hit with the smallest Euclidean distance in *ProtTucker* embedding space.

67 toward maximizing the distance between proteins from  
68 different CATH classes (anchor-negative pairs) while  
69 minimizing the distance between proteins in the same CATH  
70 class (anchor-positive pairs). Triplets of varying structural  
71 similarity were used simultaneously to optimize a single, shared  
72 network: all four CATH-levels were simultaneously learned by  
73 one FNN. The resulting model was called *ProtTucker* and its  
74 embeddings were established to identify more distant relations  
75 than is possible from sequence alone. One important objective  
76 of *ProtTucker* is to study entire functional modules through  
77 identifying more distant relations, as found to be crucial for  
78 capturing mimicry and hijacking of SARS-CoV-2 (70).

## 79 Methods

80 **CATH hierarchy.** The CATH (2,22) hierarchy (v4.3) classifies three-  
81 dimensional (3D) protein structures from the PDB (Protein Data Bank  
82 (71)) at the four levels Class, Architecture, Topology and Homologous  
83 superfamily. On average, higher levels (further away from root: H>T>A  
84 >C) are more similar in their 3D structure or have more residues for  
85 which the same level of 3D similarity is reached. We used increasing  
86 overlap in this hierarchical classification as a proxy to define increasing  
87 structural similarity between protein pairs. For example, we assumed  
88 that any two proteins with the same topology (T) are structurally more  
89 similar than any two proteins with identical architecture (A) but different  
90 topology (T). In more formal terms: SIM3D(P1,P2)>SIM3D(P3,P4),  
91 where T(P1)=T(P2) & T(P3)≠T(P4) & A(P3)=A(P4). This notion of similarity  
92 was applied on all four levels of CATH.

93  
94 **Data set.** The sequence-unique datasets provided by CATH (1,2) v4.3  
95 (123k proteins, CATH-S100) provided training and evaluation data for  
96 *ProtTucker*. A test set (300 proteins, dubbed test300 in the following)  
97 for final evaluation and a validation set (200 proteins, dubbed *val200*)

1 for early stopping were randomly split off from CATH-S100 while  
2 ensuring that (1) every homologous superfamily appeared maximally  
3 once in test300 n val200 and (2) each protein in test300 & val200 has a  
4 so called Structural Sub-group (SSG) annotation, i.e., clusters of  
5 domain structure relatives that superpose within 5Å (0.5nm), in CATH.  
6 To create the training set, we removed any protein from CATH-S100  
7 that shared more than 20% pairwise sequence identity (PIDE) to any  
8 validation or test protein according to MMSeqs2 (10) applying its  
9 iterative profile-search (--num-iterations 3) with highest sensitivity (-s  
10 7.5) and bidirectional coverage (--cov-mode 0). Additionally, large  
11 families (>100 members) within CATH-S100 were clustered at 95%  
12 PIDE and length coverage of 95% of both proteins using MMSeqs2  
13 (bidirectional coverage; --cov-mode 0). The cluster representatives  
14 were used for training (66k proteins, dubbed **train66k**) and as lookup  
15 set during early stopping on set val200. We needed a lookup set from  
16 which to transfer annotations because contrastive learning outputs  
17 embeddings instead of class predictions. For the final evaluation on  
18 test300, we created another lookup set but ignored val200 proteins  
19 during redundancy reduction (69k proteins, dubbed **lookup69k**). This  
20 provided a set of proteins for validation which had similar sequence  
21 properties to those during the final evaluation while “hiding” them  
22 during training and not using them for any other optimization. To ensure  
23 strict non-redundancy between lookup69k and test300, we further  
24 removed any protein from test300 with HVAL>0 (26) to any protein in  
25 lookup69k (219 proteins, dubbed **test219** in the following). All  
26 performance measures were computed using test219.

27 Data augmentation can be crucial for contrastive learning to reach  
28 performance in other fields (72). However, no straightforward way  
29 exists to augment protein sequences as randomly changing sequences  
30 very likely alters or even destroys protein structure and function.  
31 Therefore, we decided to use homology-based inference (HBI) for data  
32 augmentation during training, i.e., we created a new training set based  
33 on Gene3D (5) (v21.0.1) which uses Hidden Markov Models (HMMs)  
34 derived from CATH domain structures to transfer annotations to  
35 labeled CATH to unlabeled UniProt. We first clustered the 61M protein  
36 sequences in Gene3D at 50% PIDE and 80% coverage of both proteins  
37 (bidirectional coverage; --cov-mode 0) and then applied the same  
38 MMSeqs2 profile-search (--num-iterations 3 -s 7.5) as outlined above  
39 to remove cluster representatives with  $\geq 20\%$  PIDE to any protein in  
40 test300 or val200 ( $P_{\text{train}}, P_{\text{test300|val200}} \leq 20\%$ ). This filtering yielded  
41 11M sequences for an alternative training set (dubbed **train11M**).

42 The CATH detection of ProtTucker was further analyzed using a  
43 strictly non-redundant, high-quality dataset. This set was created by  
44 first clustering CATH v4.3 at 30% using HMM profiles from HMMER  
45 and additionally discarding all proteins without equivalent entry in  
46 SCOPe, i.e., the domain boundaries and the domain-superfamily  
47 assignment had to be nearly identical (3186 proteins, CATH-S30). We  
48 used the highly sensitive structural alignment scoring tool SSAP (7,8)  
49 to compute the structural similarity between all protein pairs in this set.

50 We probed whether ProtTucker embeddings might also help in  
51 solving tasks unrelated to protein structure/CATH, using as proxy a  
52 dataset assessing subcellular location prediction in ten states (56,73).  
53 We embedding-transferred annotations (EAT) from the standard  
54 **DeepLocTrain** set to 490 proteins in a recently proposed test set  
55 (**setHard**) that was strictly non-redundant to **DeepLocTrain**. Datasets  
56 described elsewhere in more detail (56,73). Finally, we showcased  
57 predictions for entire organisms using three UniProt reference  
58 proteome: *Escherichia coli* (E. Coli; reviewed, Swiss-Prot (74)),  
59 *Armillaria ostoyae* (A. ostoyae; unreviewed, TrEMBL (74)) and  
60 *Megavirus Chilensis* (M. Chilensis; unreviewed TrEMBL (74)).

61  
62 **Data representation.** Protein sequences were encoded through  
63 distributed vector representations (embeddings) derived from four  
64 different pre-trained protein language models (pLM): (1) **ProtBERT** (9)  
65 based on the NLP (Natural Language Processing) algorithm BERT (17)  
66 but trained on BFD (Big Fantastic Database) with over 2.3 billion protein  
67 sequences (11). (2) **ESM-1b** (12) is similar to (Prot)BERT but trained on  
68 UniRef50 (74). (3) **ProtT5-XL-U50** (9) (*ProtT5* for simplicity) based on  
69 the NLP sequence-to-sequence model T5 (18) trained on BFD and fine-  
70 tuned on Uniref50. (4) **ProSE** (16) trained long short-term memory cells  
71 (LSTMs) either solely on 76M unlabeled sequences from UniRef90  
72 (**ProSE-DLM**) or on additionally predicting intra-residue contacts and  
73 structural similarity from 28k SCOPe proteins (64)(multi-task: **ProSE-**  
74 **MT**). While ProSE, ProtBert and ESM-1b were trained on

75 reconstructing corrupted tokens/amino acids from non-corrupted  
76 (protein) sequence context (masked language modeling), ProtT5 was  
77 trained by *teacher forcing*, i.e., input and targets were fed to the model  
78 with inputs being corrupted protein sequences and targets being  
79 identical to inputs but shifted to the right (span generation with span  
80 size of 1 for ProtT5). Except for ProSE-MT, all pLMs were optimized  
81 only through self-supervised learning exclusively using unlabeled  
82 sequences for pre-training.

83 pLMs output a single vector for each residue yielding an  $L \times N$ -  
84 dimensional matrix (L: protein length, N: embedding dimension;  
85 N=1024 for ProtBERT/ProtT5; N=1280 for ESM-1b; N=6165 for ProSE).  
86 From this  $L \times N$  embedding matrix, we derived a fixed-size N-  
87 dimensional vector representing each protein by averaging over protein  
88 length (Fig. 1 (9)). The pLMs were used as static feature encoder only,  
89 i.e., no gradient was backpropagated for fine-tuning. As recommended  
90 in the original publication (9), for ProtT5, we only used the encoder part  
91 of ProtT5 in half-precision to embed protein sequences. Similarly,  
92 ProtBERT embeddings were derived in half-precision.

93  
94 **Contrastive learning: architecture.** A two-layer feedforward neural  
95 network (FNN) projected fixed-size per-protein (sentence-level)  
96 embeddings from 1024-d (1280-d/6165-d for ESM-1b and ProSE  
97 respectively) to 256 and further to 128 dimensions with the standard  
98 hyperbolic tangent (*tanh*) as non-linearity between layers. We also  
99 experimented with deeper/more sophisticated networks without any  
100 gain from more free parameters (data not shown). This confirmed  
101 previous findings that simple networks suffice when inputting  
102 advanced embeddings (9,12,62,75). As the network was trained using  
103 contrastive learning, no final classification layer was needed. Instead,  
104 the 128-dimensional output space was optimized directly.

105  
106 **Contrastive learning: training.** During training, the new embedding  
107 space spanned by the output of the FNN was optimized to capture  
108 structural similarity using triplets of protein embeddings. Each triplet  
109 had an anchor, a positive and a negative. In each epoch, all **train66k**  
110 proteins were used once as anchor, while positives and negatives were  
111 sampled randomly from **train66k** using the following **hierarchy-**  
112 **sampling**. First, a random level  $\alpha$  ( $\alpha=[1,2,3,4]$ ) describing the increasing  
113 structural overlap between triplets was picked. This defined a positive  
114 (same CATH-number as anchor up to level  $\alpha \leq \alpha$ ) and a negative label  
115 (same CATH-number as anchor up to level  $\alpha' < \alpha$ ). For instance, assume  
116 the anchor has the CATH-label 1.25.10.60 (Rad61, Wapl domain) and  
117 we randomly picked  $\alpha=3$  (topology-level), only proteins with the  
118 anchor's topology (1.25.10.x; Leucine-rich Repeat Variant) qualify as  
119 positives while all negatives share the anchor's architecture (1.25.y.z;  
120 Alpha Horseshoe) with different topology ( $y \neq 10$ ). Self-hits of the anchor  
121 were excluded. From the training proteins fulfilling those constraints,  
122 one positive and one negative were picked at random. If no triplets  
123 could be formed (e.g.,  $\alpha=4$  with a single-member homologous  
124 superfamily had no positive for this anchor/ $\alpha$  combination),  $\alpha$  was  
125 changed at random until a valid triplet could be formed (eventually, all  
126 proteins found a class-level partner). This flexibility in sampling enabled  
127 training on all proteins, independent of family size.

128 Unlike randomly sampling negatives, the hierarchical sampling could  
129 be described as semi-hard sampling as it zoomed into triplets that were  
130 neither too easy (little signal) nor too hard (outliers) to classify by  
131 ensuring a minimal overlap between the anchor and the chosen  
132 negative/positive pair. Thereby, trivial triplets are under-sampled  
133 (avoided), i.e., those with 3D structures so different that the separation  
134 becomes trivial (*daylight zone*). As the final triplet selection was still  
135 random, anchor-positive pairs could still be too easy/similar which was  
136 shown to hinder the success of contrastive learning (76). To solve this  
137 issue, we paired hierarchy-sampling with so called **batch-hard**  
138 **sampling** (76) which offers a computationally efficient solution for  
139 sampling *semi-hard* triplets within one mini-batch. More specifically,  
140 we combined **batch-hard sampling** with the triplets created using  
141 **hierarchy-sampling** by re-wiring all proteins, irrespective of anchor,  
142 positive or negative, within one mini-batch such that they satisfied the  
143 hierarchy-sampling criterion and had maximum/minimal Euclidean  
144 distance for anchor-positive/anchor-negative pairs. Sampling hard  
145 triplets only within each mini-batch instead of across the entire data set  
146 avoided extreme outliers (potentially too hard/noisy) while increasing  
147 the rate of semi-hard anchor-positive/anchor-negative pairs. Assume  
148 multiple proteins of the topology Leucine-rich Repeat Variant were



1 within one mini-batch, the hardest positive for each anchor would be  
2 picked by choosing the anchor-positive pair with the largest Euclidean  
3 distance. Accordingly, anchor-negative pairs would be picked based  
4 on the smallest Euclidean distance. For each mini-batch, this sampling  
5 was applied to all four levels of the CATH-hierarchy, so triplets were  
6 wired on all four CATH levels resulting in a total batch-size of about:  
7  $batch\_size * 3 * 4$ . This was an “about” instead of “equal” because for  
8 some mini-batches, not all proteins had valid triplets for all four levels.

9 Finally, the same two-layer FNN was used (hard parameter sharing)  
10 to project the 1024-d (or 1280-d/6165-d for ESM-1b or ProSE  
11 respectively) embeddings of all proteins, irrespective of anchor,  
12 positive or negative, to a new 128-d vector space. The *Soft Margin Loss*  
13 was used to optimize this new embedding space such that anchor-  
14 positive pairs were pulled together (reduction of Euclidean distance)  
15 while pushing apart anchor-negative pairs (increase of Euclidean  
16 distance). The efficiency of combining the Soft Margin Loss with batch-  
17 hard sampling was established before (76), although without prior  
18 hierarchical triplet sampling. Here, we used Soft Margin Loss as  
19 implemented in PyTorch:

$$20 \quad Loss(d, y) = \sum_t \frac{\log(1+e^{-(y[t]+d[t])})}{|d|} \quad (\text{Eqn. 1})$$

$$21 \quad d = \left\{ \min_{\substack{n \in B \wedge a \neq n \\ C_i(a) \neq C_i(n)}} D(f(a), f(n)) - \max_{\substack{n \in B \wedge p \neq a \\ C_i(a) \neq C_i(p)}} D(f(a), f(p)) \right\} \quad (\text{Eqn. 2})$$

$$22 \quad \forall a \in \{B\} \wedge \forall i \in 1,2,3,4 \wedge \forall C_i \in \{CATH \text{ labels}\}$$

23  $B$  represents one mini-batch created through hierarchical sampling,  
24  $f(a)$ ,  $f(p)$  and  $f(n)$  represent the *ProtTucker* embeddings of proteins  $a$   
25 (anchor),  $n$  (negative), and  $p$  (positive) represented as pLM  
26 embeddings;  $C_i$  represents the CATH annotation of a protein on the  $i$ 'th  
27 hierarchy level of CATH; finally,  $D(f(a), f(x))$  represents the Euclidean  
28 distance between the embeddings for proteins  $a$  and  $x$ . We created the  
29 mini-batch  $B$  used for training by choosing for each protein or anchor  
30  $a$  in  $B$  the hardest negative  $n$  and the hardest positive  $p$  by picking those  
31 proteins in  $B$  that have the smallest | largest Euclidean distance  $D$  to  $a$   
32 *ProtTucker* embedding space while not sharing | sharing  $C_i$ ,  
33 respectively. Those semi-hard triplets are indexed by  $t$  and  $d[t]$  referring  
34 to the difference between  $D$  of anchor-negative and  $D$  of anchor-  
35 positive. In our case, the label for the  $t$ 'th triplet  $y[t]$  is always 1 as the  
36 sign of  $x$  indicates training success, i.e., whether the distance anchor-  
37 positive is smaller than that between anchor-negative.

38 Consequently, triplets of varying structural similarity were used  
39 simultaneously to optimize a single, shared network, i.e., all four CATH-  
40 levels were learned by the same network at the same time (Fig. 1A). We  
41 used the *Adam optimizer* (77) with a learning rate of 0.001, and a batch-  
42 size of 256 to optimize the network. The effective batch-size increased  
43 due to batch-hard sampling to a maximum of 3072, depending on the  
44 number of valid triplets that could be formed within the current mini-  
45 batch. Training terminated (*early stopping*) at the highest accuracy in  
46 predicting the correct homologous superfamily for set *val200*.

47  
48  
49 **Evaluation and prediction (inference).** While supervised training  
50 directly outputs class predictions, contrastive learning, outputs a new  
51 embedding space. Thus, predictions (inferences) were generated as for  
52 homology-based inference (HBI), i.e., given protein  $X$  with experimental  
53 annotation (CATH assignment) and a query protein  $Q$  without, then HBI  
54 transfers the annotation from  $X$  to  $Q$  if sequence similarity  $SIM(X, Q) >$   
55 threshold. For contrastive learning, we replaced  $SIM(X, Q)$  by  $D(f(X), f(Q))$   
56 with  $D$  as the shortest Euclidean distance in embedding space (Fig.  
57 1B). In previous studies (9,58,59), we found the Euclidean distance  
58 performed better than the cosine distance which is more common in  
59 AI/NLP. The Euclidean distance also optimized the *ProtTucker*  
60 embeddings. Set *test219* with the *lookup69k* as lookup set (set of all  $X$ )  
61 served as final evaluation. If no protein in the lookup set shared the  
62 annotation of the query protein at a certain CATH-level (more likely for  
63 H than for C), the sample was excluded from the evaluation of this  
64 CATH-level as no correct prediction was possible (Table S1).

65 During evaluation, we compared the performance of our embedding-  
66 based annotation transfer (EAT) to HBI using the sequence  
67 comparisons from MMSegs2 (10). While transferring only the HBI hit  
68 with the lowest E-value, we searched for hits up to an E-value of 10 to  
69 ensure that most proteins had at least one hit while using the highest  
70 sensitivity setting (-s 7.5). Additionally, we used publicly available

71 CATH-Gene3D (2) Hidden Markov Models (HMMs) along with HMMER  
72 (6) to detect remote homologs up to an E-value of 10.

73 For both approaches, EAT and HBI, we computed the accuracy as  
74 the fraction of correct hits for each CATH-level. A hit at lower CATH-  
75 levels could be correct if and only if all previous levels were correctly  
76 predicted. Due to varying number of samples at different CATH-levels  
77 (Table S1), performance measures not normalized by background  
78 numbers could be higher for lower levels. Predictions were counted as  
79 incorrect if a query did not have a hit in the lookup set but a lookup  
80 protein of the same CATH annotation existed. This not only affected  
81 the number of proteins available at different CATH-levels (Table S2) but  
82 also the number of classes (Table S3). A random baseline was  
83 computed by transferring annotations from a randomly picked protein  
84 in *lookup69k* to *test219*.

85  
86 **Performance measures.** The four coarse-grained classes at the top  
87 CATH level (“C”) are defined by their secondary structure content.  
88 These four branch into 5481 different superfamilies with distinct  
89 structural and functional aspects (CATH v4.3.0). However, most  
90 standard metrics are defined for binary cases which requires some  
91 grouping of predictions into four cases: 1) TP (true positives): correctly  
92 predicted to be in the *positive* class, 2) TN (true negatives): correctly  
93 predicted to be in the *negative* class, 3) FP (false positives): incorrectly  
94 predicted to be positives, and 4) FN (false negatives): incorrectly  
95 predicted to be in the negative class. Here, we focused on  
96 performance measures applicable for multiclass problems and are  
97 implemented in scikit (78). These were in particular: **accuracy** (Acc,  
98 Eqn. 3) as the fraction of correct predictions

$$99 \quad Accuracy(y, \hat{y}) = \frac{1}{n_{samples}} \sum_{i=0}^{n_{samples}-1} 1(\hat{y}_i = y_i) \quad (\text{Eqn. 3})$$

100 with  $y_i$  being the ground truth (experimental annotation) and  $\hat{y}_i$  the  
101 prediction for protein  $i$ . In analogy, we defined **coverage** as the  
102 proportion of the *test219* proteins for which a classifier made a  
103 prediction at a given prediction reliability  $\hat{y}_i^r$  and reliability threshold  $\theta$ :

$$104 \quad Coverage(y, \hat{y}) = \frac{1}{n_{samples}} \sum_{i=0}^{n_{samples}-1} 1(\hat{y}_i^r < \theta) \quad (\text{Eqn. 4})$$

105 In these definitions accuracy corresponds to precision, and coverage  
106 to recall binarizing a multiclass problem through micro-averaging, i.e.,  
107 by counting the total TPs, FPs and FNs globally, irrespective of the  
108 class. The multi-class extension of **Matthew's correlation coefficient**  
109 (**MCC**, (45)) was defined as:

$$110 \quad MCC = \frac{c \times s - \sum_k^K p_k \times t_k}{\sqrt{(s^2 - \sum_k^K p_k^2) \times (s^2 - \sum_k^K t_k^2)}} \quad (\text{Eqn. 5})$$

111 with  $t_k = \sum_i^K C_{ik}$  as the number of times class  $k$  truly occurred,  $p_k =$   
112  $\sum_i^K C_{ki}$  as the number of times class  $k$  was predicted,  $c = \sum_i^K C_{kk}$ , the  
113 total number of samples correctly predicted, and  $s = \sum_i^K \sum_j^K C_{ij}$ , the  
114 total number of samples.

115  
116 **95% confidence intervals** for accuracy and MCC were estimated over  
117  $n=1000$  bootstrap sets; for each bootstrap set we randomly sampled  
118 predictions from the original test set with replacement. Standard  
119 deviation (or in the case of bootstrapping: bootstrap standard error)  
120 was calculated as the difference of each test set ( $x_i$ ) from the average  
121 performance ( $\bar{X}$ ) (Eqn. 6). 95% confidence intervals were estimated  
122 by multiplying the bootstrap standard error by 1.96.

$$123 \quad StdDev = \sqrt{\frac{\sum_i (x_i - \bar{X})^2}{n}} \quad (\text{Eqn. 6})$$

## 124 Results

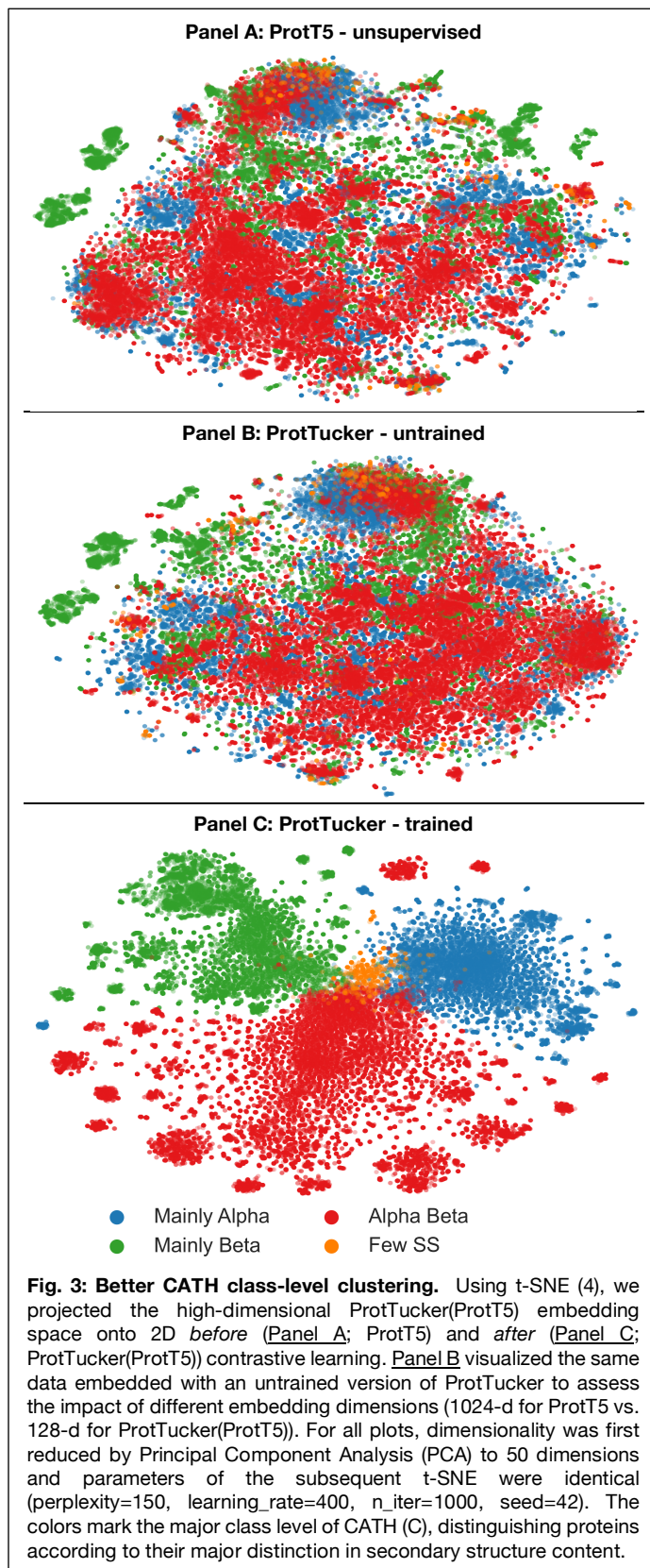
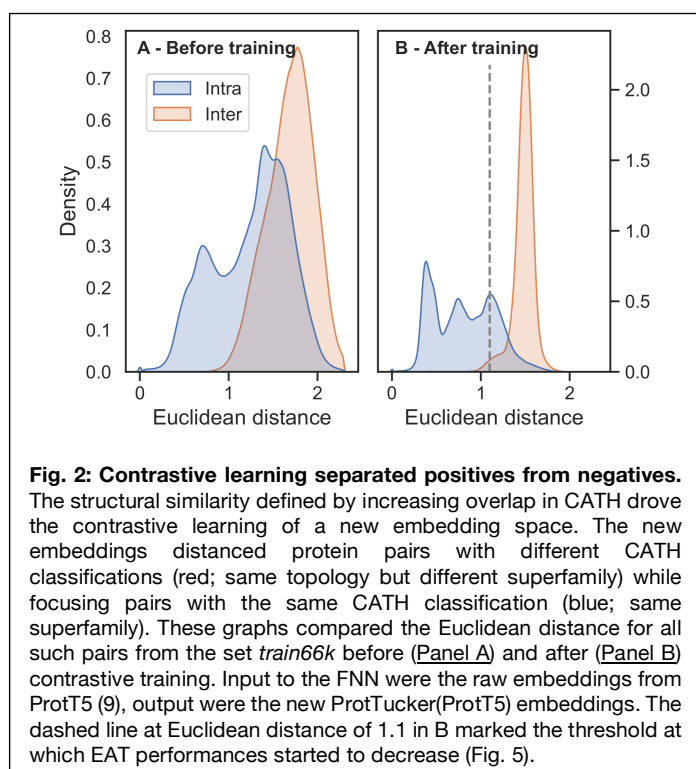
125 **Generalization of HBI to EAT.** Homology-based inference  
126 (HBI) uses sequence similarity to transfer annotations from  
127 experimentally characterized (labelled) to uncharacterized  
128 (unlabelled) proteins. More specifically, an unlabeled query  
129 protein  $Q$  is aligned against a set of proteins  $X$  with  
130 experimental annotations (dubbed *lookup set*) and the  
131 annotation of the best hit, e.g., measured as lowest E-value, is  
132 transferred if it is below a certain threshold (e.g., E-  
133 value( $Q, X$ ) < 10<sup>-3</sup>). This relates to inferring the annotation of a  
134 query protein from the  $k$ -Nearest Neighbors ( $k$ -NN) in sequence

1 space. More recently, similar approaches expanded from  
 2 distance in sequence to distance in embedding space (Fig. 1B)  
 3 as means for k-NN based annotation transfer (58,60). Here, we  
 4 refer to such methods as embedding-based annotation  
 5 transfer (EAT). We used EAT as a proxy for the comparison of  
 6 embeddings from five different protein Language Models  
 7 (pLMs): ProSE-DLM & ProSE-MT (16), ProtBERT & ProtT5 (9),  
 8 and ESM-1b (12). Next, we used triplets of proteins (anchor,  
 9 positive, negative) to learn a new embedding space by pulling  
 10 protein pairs from the same CATH class (anchor-positive)  
 11 closer together while pushing apart pairs from different CATH  
 12 classes (anchor-negative; Fig. 1A). We referred to this method  
 13 as *ProtTucker*. At this stage, we did not fine-tune the pre-  
 14 trained pLMs. Instead, we created a new embedding space  
 15 using a two-layer feed-forward neural network (FNN).

16  
 17 **EAT with raw embeddings level with HBI.** First, we  
 18 transferred annotations from all proteins in *lookup69k* to any  
 19 protein in *test219*. All pLMs significantly (at 95% CI -  
 20 confidence interval) outperformed random annotation transfer  
 21 (Table 1). Performance differed between pLMs (Table 1), with  
 22 ProtBERT (9) being consistently worse than LSTM-based  
 23 ProSE-DLM or more advanced transformers (ESM-1b, ProtT5).  
 24 ESM-1b and ProtT5 also numerically outperformed ProSE-  
 25 DLM and HBI using MMseqs2 (10), especially on the most fine-  
 26 grained and thus hardest level of superfamilies. However,  
 27 MMseqs2 had been used for redundancy-reduction, i.e., the  
 28 data set had been optimized for minimal performance of  
 29 MMseqs2. HBI using publicly available HMM-profiles from  
 30 CATH-Gene3D (2) along with the profile-based advanced  
 31 HMMER (6) designed for more remote homology detection,  
 32 outperformed all raw embeddings for homologous  
 33 superfamilies, while embeddings from ESM-1b and ProtT5  
 34 appeared superior on the class- and architecture-level (Table  
 35 1). In fact, all HBI values, for MMseqs2 and HMMER, were  
 36 highest for the H-level, and 2- highest for the C-level. In  
 37 contrast, raw pLM embeddings mirrored the random baseline

38 trend, with numbers being inversely proportional to the rank in  
 39 C, A, T, H (highest for C, lowest for H, Table 1).

40  
 41 **EAT improved through supervised embeddings.** ProSE-MT  
 42 expands ProSE-DLM by additionally training on intra-residue  
 43 contacts and structural similarity using labeled data from  
 44 SCOPe (16). This additional effort was reflected by the higher





**Table 1: Accuracy for annotation transfer to queries in test219 \***

	Method/Input	C	A	T	H	Mean
Baseline	Random	29 ±6	9 ±4	1 ±2	0 ±0	10 ±3
HBI	MMSeqs2 (sequence)	52 ±7	36 ±6	29 ±6	35 ±6	38 ±6
	HMMER (profile)	70 ±6	60 ±6	59 ±7	77 ±7	67 ±6
EAT Unsupervised	ProSE-DLM	74 ±6	48 ±7	28 ±6	25 ±7	44 ±6
	ESM-1b	79 ±5	61 ±6	50 ±7	57 ±8	62 ±7
	ProtBERT	67 ±6	38 ±6	22 ±6	18 ±6	36 ±6
	ProtT5	84 ±5	67 ±6	57 ±6	64 ±8	68 ±6
EAT - Supervised	ProSE-MT	82 ±5	65 ±6	52 ±7	56 ±8	64 ±7
EAT - Contrastive learning – ProtTucker	ProSE-DLM	78 ±4	53 ±6	32 ±6	29 ±7	48 ±6
	ProSE-MT	87 ±4	68 ±6	53 ±7	55 ±8	66 ±6
	ESM-11b	87 ±4	68 ±6	59 ±7	70 ±7	71 ±6
	ProtBERT	81 ±5	52 ±7	37 ±6	39 ±8	52 ±7
	ProtT5	<b>89 ±4</b>	75 ±6	64 ±6	76 ±6	76 ±6
	ProtT5 (train11M)	88 ±4	<b>77 ±5</b>	<b>68 ±5</b>	<b>79 ±7</b>	<b>78 ±6</b>

\* Accuracy (Eqn. 3) for predicting CATH (2,3) levels (C, A, T, H) by transferring annotations from *Lookup69k* (lookup set) to *test219* (queries); shown for each of the four levels from the most coarse-grained level class C to the most fine-grained level of homology H. The column *Mean* marked the simple arithmetic average over the four performance values. Queries with at least one lookup protein of the same CATH classification but without any hit at E-value<10 for MMSeqs/HMMER were counted as incorrect predictions. Errors indicate bootstrapped 95% confidence intervals, i.e., 1.96 bootstrap standard errors (Eqn. 6). Queries with at least one lookup protein of the same CATH annotation but without any hit (no hit with E-value<10 for MMSeqs/HMMER; irrelevant for EAT) were counted as wrong predictions. **Bold** letters mark the numerically highest values (averages over all *test219* proteins) in each column irrespective of the confidence interval.

Methods: Baseline: Random transferred the label of a randomly picked protein; HBI: MMSeqs2 (10) used single sequence search to transfer the annotation of the hit with the lowest E-value; HBI: HMMER used HMM-profiles (6); EAT-unsupervised: embedding-based transfer of annotations using the smallest Euclidean distance measured in embedding space of unsupervised pLMs ProSE-DLM, ESM-1b (12), ProtBERT and ProtT5 (9); EAT-supervised: annotation transfer using ProSE-MT trained on structural data in SCOPe; EAT: contrastive learning ProtTucker: contrastive learning trained on CATH classifications in *train66k* using as input embeddings from ProSE-DLM, ProSE-MT, ESM-1b, ProtBERT, and ProtT5; ProtTucker-ProtT5 (train11M) trained on additional data from Gene3D (*train11M*).

1 performance for all CATH levels (Table 1, ProSE-MT>ProSE-DLM). The supervision pushed the LSTM-based ProSE-MT to reach performance levels close to the unsupervised, raw embeddings from transformer-based ProtT5. The performance gap increased with classification difficulty (Table 1, ProtT5>ProSE-MT, especially at the H-level).

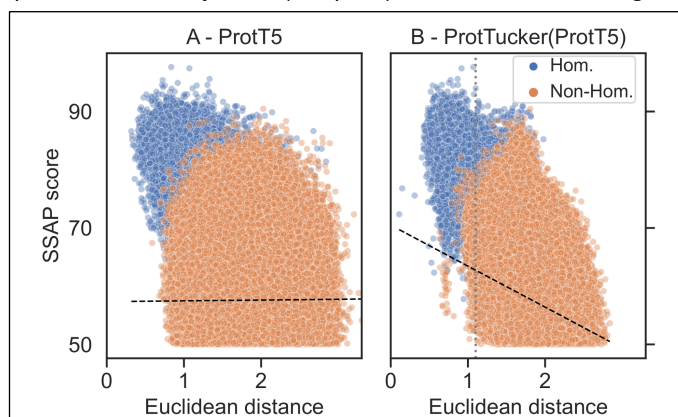
#### EAT improved by contrastively learning embeddings.

2 Contrastive learning tries to bring members from the same class/CATH-level closer while pushing those from different classes further apart. One success is the degree to which these two distributions (same vs. different) were separated through training: the distribution of all pairwise Euclidean distances within (intra/same) and between (inter/different) superfamilies in *train66k* changed substantially through contrastive learning (Fig. 2). Before applying contrastive learning, the distributions between (inter, Fig. 2: red) and within (intra, Fig. 2: blue) overlapped much more than after.

3 Displaying the information learned by the embeddings, we compared t-SNE projections colored by the four main CATH classes before (Fig. 3A) and after (Fig. 3C) contrastive learning. These two projections compared 1024 dimensions from ProtT5 (Fig. 3A) with 128 dimensions from ProtTucker (Fig. 3C). To rule out visual effects from higher dimensionality, we also compared the untrained, randomly initialized version of ProtTucker using pre-trained ProtT5 embeddings as input (Fig. 3B). For all cases, the data set (*train66k*) and the parameters for dimensionality reduction were identical. T-SNE projections of raw ProtT5 embeddings qualitatively suggested some class separation (clustering). The information underlying this separation was preserved when projecting ProtT5 embeddings through an untrained ProtTucker (Fig. 3B). Embeddings from ProtTucker(ProtT5), i.e., those resulting through contrastive learning, separated the classes much more clearly.

4 To further probe the extent to which contrastive learning captured remote homologs, we compared the Euclidean

37 distance between all protein pairs in a 30% non-redundant dataset (CATH-S30) with the structural similarity of those pairs computed via SSAP (7,8) (Fig. 4). From the ~10M possible pairs between the 3,186 proteins in CATH-S30 (problem not fully symmetric, therefore  $N*(N-1)$ : 10.1M), 7.1M had to be discarded due to low quality (SSAP-score <50), leaving 2.9M pairs of which only 1.8% (53k pairs) had the same homologous



**Fig. 4: ProtTucker captured fine-grained structural similarity.**

3186 non-redundant proteins (CATH-S30) probed the remote homology detection of embeddings *before* (Panel A) and *after* contrastive learning (Panel B). The Euclidean distance between ProtTucker embeddings (Panel B) correlated better with structural similarity computed by SSAP (7,8) than unsupervised embeddings (Panel A): Spearman  $\rho=0.22$  and  $p=0.05$  (black dashed lines). This correlation increased to  $\rho=0.37$  and  $p=0.26$  for structurally more similar protein pairs (SSAP-score>70). Only 1.8% (53k) of all structurally similar pairs were in the same homologous superfamily (blue). The unsupervised ProtT5 already separated homologous pairs from others, but ProtTucker(ProtT5) improved, especially, for hard cases with low structural similarity. The gray dashed line at Euclidean distance=1.1 in Panel B marked the threshold at which EAT performances started to decrease (Fig. 5).

**Table 2: Ablation study. \***

	C	A	T	H	Mean
Baseline	89 ±4	75 ±6	64 ±6	76 ±6	76 ±6
-batch-hard	88 ±4	73 ±6	62 ±7	69 ±7	73 ±6
-hierarchy	83 ±5	69 ±6	62 ±7	71 ±7	71 ±6
-both	83 ±5	63 ±6	51 ±7	57 ±8	64 ±7

\* **Accuracy** (Eqn. 3) and 95% CI (Eqn. 6) for predicting CATH-levels (2,3) through EAT from *Lookup69k* (lookup set) to *test219* (queries). We investigate the effect on performance when dropping either *batch-hard* sampling, *hierarchy*-sampling or *both* from the *Baseline* model (ProtTucker(ProtT5)).

1 superfamily (Fig. 4: blue). Despite this imbalance, unsupervised  
 2 ProtT5 (Fig. 4A) already to some extent separated the same  
 3 from different superfamilies. Still, ProtTucker(ProtT5) improved  
 4 this separation, especially, for pairs with low structural similarity  
 5 (Fig. 4B). This was supported by the Spearman correlation  
 6 coefficient between the structural similarity and the Euclidean  
 7 distance increasing from 0.05 to 0.22 after contrastive learning.  
 8 When considering only the subset of pairs that likely have  
 9 similar folds (SSAP-score>70), this correlation increased to  
 10 0.26 and 0.37 for ProtT5 and ProtTucker(ProtT5), respectively.

11 The trend captured by the better separation of distributions  
 12 (Fig. 2) and structural features (Fig. 3, Fig. 4) translated directly  
 13 into performance increases: all embeddings optimized on the  
 14 CATH hierarchy through contrastive learning yielded better EAT  
 15 classifications than the raw embeddings from pre-trained pLMs  
 16 (Table 1). As ProtTucker described the process of refining raw  
 17 embeddings through contrastive learning, we used the  
 18 annotation  $\text{ProtTucker}(X)$  – in this section also shortened to  
 19  $\text{PT}(X)$  – to refer to the embeddings output by inputting the pre-  
 20 trained pLM  $X$  into the contrastive learning. The improvements  
 21 were larger for more fine-grained CATH levels: all models  
 22 improved significantly for the H-level, while only PT(ProtBERT)  
 23 and PT(ESM-1b) improved from 4 to 14 or from 0 to 21

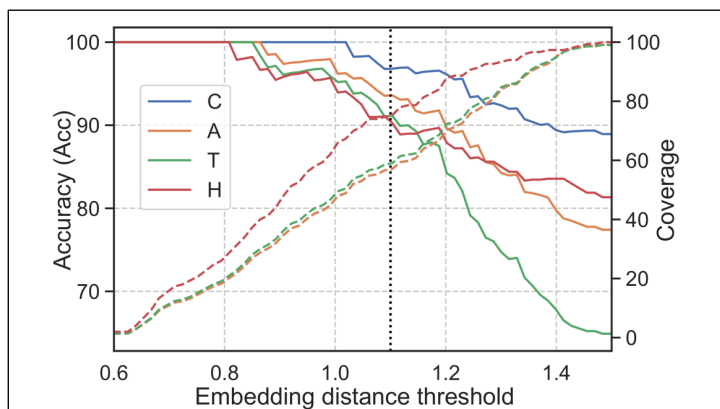
24 percentage points for the C-, and the H-level, respectively.  
 25 PT(ProtT5) consistently outperformed all other pLMs on all four  
 26 CATH-levels, with an increasing performance gap toward the  
 27 more fine-grained H-level at which all pLMs except for  
 28 PT(ESM-1b) performed significantly worse. The improvements  
 29 from contrastive learning for PT(ProSE-DLM) and PT(ProSE-  
 30 MT) were mostly consistent but largely insignificant. Especially,  
 31 the model already optimized using labeled data (ProSE-MT)  
 32 hardly improved through another round of supervision by  
 33 contrastive learning and even worsened slightly at the H-level.  
 34 We augmented the training set for PT(ProtT5) by adding HBI-  
 35 hits from HMM-profiles provided by CATH-Gene3D (if  
 36 sequence dissimilar to test300/val200). This increased the  
 37 training set from 66k (66\*10) to 11m (11\*10) proteins (15-fold  
 38 increase) and raised performance, although the higher values  
 39 were neither statistically significant nor consistent (Table 1:  
 40 values in last row not always higher than those in 2- to last row).

41  
 42 **Ablation study.** We studied the effect of *batch-hard* and  
 43 *hierarchical* sampling on performance by removing each  
 44 component when training PT(ProtT5) (Table 2). Benchmarking  
 45 on EAT from *lookup69k* to *test219* established that removing  
 46 either component lowered performance for all CATH levels.  
 47 Dropping both sampling methods substantially lowered  
 48 performance. While dropping *batch-hard* sampling still reached  
 49 high performance for the coarse-grained C- and A-level,  
 50 dropping hierarchy-sampling dropped performance for both.  
 51 Dropping both sampling technique, performed worse for all  
 52 CATH levels but the decrease for the more fine-grained  
 53 superfamily level was much larger than for the C-level.

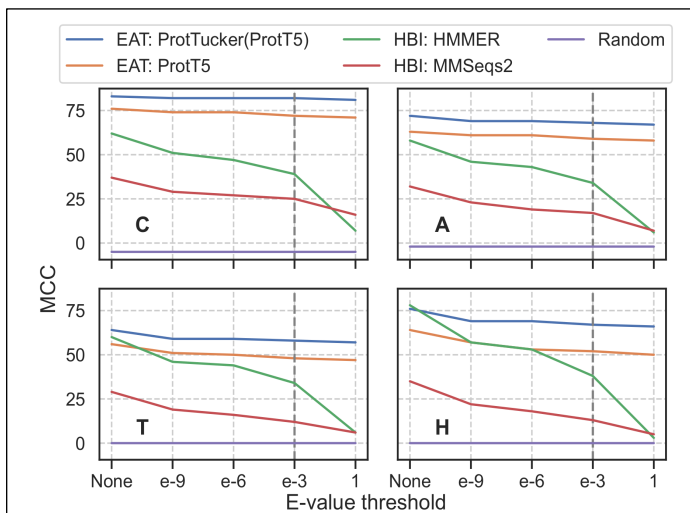
54  
 55 **Embedding distance correlated with accuracy.** The MCC,  
 56 (Eqn. 5) of HBI inversely correlated with E-value (Fig. 6, HBI-  
 57 methods): more significant hits (lower E-values) more often  
 58 shared the same CATH level than less significant hits (higher  
 59 E-values). In analogy, we explored the corresponding relation for  
 60 EAT, namely the correlation between accuracy (Eqn. 3) and  
 61 embedding distance for ProtTucker(ProtT5). Indeed, accuracy  
 62 correlated with embedding distance (Fig. 5: solid lines) while  
 63 recall inversely correlated (Fig. 5: dashed lines) for all four  
 64 classes. For instance, when transferring only annotations for  
 65 closest hits with Euclidean distances of 1.1 or less, predictions  
 66 were made for 57%, 57%, 59% or 75% of the test set  
 67 (coverage, Eqn. 4) of these 96%, 93%, 91% or 90% were  
 68 correct for levels C, A, T, H, respectively.

69  
 70 **ProtTucker reached into the midnight zone.** Annotation  
 71 transfer by HBI crucially depends on the sequence similarity  
 72 between query (unknown annotation) and template  
 73 (experimental annotation). Usually, the significance of an  
 74 inference is measured as the chance of finding a hit at random  
 75 for a given database size (E-value; the lower the better). Here,  
 76 we compared the effect of gradually removing hits depending  
 77 on their E-values. Essentially, this approach measured how  
 78 sensitive performance was to the degree of redundancy  
 79 reduction between query and lookup set. For instance, at a  
 80 value of  $10^{-6}$  (dashed vertical lines in Fig. 6), all pairs with  
 81  $E\text{-values} \leq 10^{-6}$  were removed. HBI based on sequence alone  
 82 performed much better with than without residual redundancy  
 83 (Fig. 6). The trend was similar for EAT, but much less  
 84 pronounced: EAT succeeded for pairs with very different  
 85 sequences (Fig. 6 toward right) almost as well as for proteins  
 86 with more sequence similar matches in the set (Fig. 6 toward  
 87 left: EAT almost as high as toward right).

88 **ProtTucker not a generalist.** We evaluated the generality of  
 89 ProtTucker embeddings by (mis)-using them as exclusive input



**Fig. 5: Embedding distance correlated with reliability.** Similar to varying E-value cut-offs for HBI, we examined whether the fraction of correct predictions (accuracy; left axis; Eqn. 3) depended on embedding distance (x-axis) for EAT. This was shown by transferring annotations for all four CATH levels (Class: blue; Architecture: orange; Topology: green; Homologous superfamily: red) from *lookup69k* to the queries in set *test219* (Panel B in Fig. 1) using the hit with lowest Euclidean distance. The fraction of *test219* proteins having a hit below a certain distance threshold (coverage; right axis, dashed lines; Eqn. 4) was evaluated separately for each CATH level. For example, at an Euclidean distance of 1.1 (vertical dotted line), 75% of the proteins found a hit at the H-level (Cov(H)=75%) and 90% were correctly predicted (Acc(H)=90%; SOM Tables S3 and S4 for more details). Thus, decreasing embedding distance correlated with EAT performance. This correlation enables users to select only the, e.g., 10% top hits, or as many hits to a certain CATH level as possible, depending on the objectives.



**Fig. 6: Performance decreasing with lower residual sequence similarity.** We analyzed the change of performance in MCC (Eqn. 3) through removing proteins from *lookup69k* based on their E-value with respect to *test219* for two HBI-based (green: HMMER (6); red: MMSeqs2 (10)) and two EAT-based methods (orange: raw ProtT5 (9); blue: contrastive learning optimized ProtTucker(ProtT5)). The E-values were derived by searching sequences in *test219* against *lookup69k* using (i) HMM-profiles from CATH-Gene3D (2) through HMMer and (ii) MMSeqs2 sequence search with highest sensitivity ( $-s\ 7.5, -cov\ 0$ ). “None” referred to the performance without applying any threshold, i.e., all proteins in *lookup69k* were used for annotation transfer; all other thresholds referred to removing proteins below this E-value from *lookup69k*. Predictions were considered as false positives when no hit was found; pairs without CATH class matches were ignored. While the performance of EAT using raw ProtT5 and refined ProtTucker(ProtT5) embeddings decreased upon removing sequence similar pairs (toward right), HBI-based methods dropped significantly more. The default threshold for most sequence searches (E-value < 10<sup>-3</sup>) was highlighted by vertical, gray, dashed lines.

1 to predict subcellular location in ten states. To this end, we EAT  
 2 transferred annotations from an established data set (Table S5)  
 3 to a strictly non-redundant test set (*setHard*, Table S5).  
 4 ProtTucker(ProtT5) embeddings outperformed the raw ProtT5  
 5 embeddings in the CATH classification for which they were  
 6 optimized (structural similarity; Table 1), there appeared no  
 7 performance gain in predicting location. Conversely,  
 8 performance also did not decrease significantly, indicating that  
 9 the new embeddings retained some of the information available  
 10 in ProtT5 embeddings.

11  
 12 **Family size mattered.** By clustering very large protein families  
 13 (>100 members after redundancy reduction) at 95% PIDE, we  
 14 constrained the redundancy in set *train66k*. Nevertheless, when  
 15 splitting *test219* into three bins of varying family sizes, we still  
 16 observed a trend towards higher accuracy (Eqn. 3) for larger  
 17 families at the H-level (Fig. S1). We chose the three bins such  
 18 that they contained about the same number of samples (small  
 19 families: ≤10 members, medium: 11-70, and large: ≥70  
 20 members). Especially, unsupervised EAT using the raw ProtT5  
 21 embeddings exhibited a clear trend towards higher accuracy  
 22 with increasing family size. In contrast, the two HBI-methods  
 23 (MMseqs2, HMMER), as well as EAT using the optimized  
 24 ProtTucker(ProtT5) embeddings performed similarly for small  
 25 and medium-sized families and much better for large families.

26 **EAT complements HBI.** As previously shown (59), ProtTucker  
 27 can improve clustering functional families (79). Here, we  
 28 showed how EAT can be used to detect outliers. Firstly, we  
 29 computed pairwise Euclidean distances between the

30 embeddings of all protein pairs in set *train66k* and analyzed the  
 31 five pairs (10 proteins) with the highest Euclidean distance in  
 32 the same homologous superfamily (Table S6). High distance  
 33 within the same homologous superfamily indicates potentially  
 34 wrong annotations. Secondly, we computed the nearest  
 35 neighbors of those ten proteins to find an alternative, potentially  
 36 more suitable annotation. For instance, the proteins in the  
 37 *Phosphorylase Kinase* superfamily with the largest embedding  
 38 distance (4pdyA01, bacterial aminoglycoside phosphotrans-  
 39 ferase) to any other protein within this family (3skjF00, human  
 40 *Galactose-binding domain-like* (80)) linked to different UniProt  
 41 entries (C8WS74\_ALIAD and EPHA2\_HUMAN). In contrast, the  
 42 nearest neighbor (3heiA00, human *phosphorylase kinase* (81))  
 43 of 3skjF00 linked to the same UniProt entry (EPHA2\_HUMAN  
 44 (74)) with the same enzymatic activity (EC number 2.7.10.1  
 45 (82)). Such analyses may indicate impure homologous  
 46 superfamilies along with suggesting alternative labels to be  
 47 confirmed or rejected through manual curation.

48  
 49 **EAT predicts entire proteomes in minutes.** Training  
 50 ProtTucker(ProtT5) required generating ProtT5 embeddings for  
 51 *train66k*. This took 23m and 11m, respectively. Embeddings  
 52 were generated using ProtT5 in half-precision with batch  
 53 processing. All times were measured on a single Nvidia RTX  
 54 A6000 with 48GB of vRAM and an AMD EPYC ROME 7352.

55 When predicting for new queries, ProtTucker requires *labeled*  
 56 *lookup* proteins from which annotations can be transferred to  
 57 unlabeled query proteins. Embeddings for this lookup set are  
 58 pre-computed for the first query and can be re-used for all  
 59 subsequent queries at any future time. The time required to  
 60 labeled *lookup* proteins from which annotations can be  
 61 transferred to unlabeled query proteins. Embeddings for this  
 62 lookup set are pre-computed for the first query and can be re-  
 63 used for all subsequent queries at any future time. The time  
 64 required to generate ProtTucker embeddings from the  
 65 embeddings of pLMs was negligible as its generation required  
 66 only a single forward pass through a two-layer FNN. This  
 67 implied that the total time for EAT with ProtTucker was largely  
 68 determined by the embedding generation speed. For instance,  
 69 creating per-protein embeddings from ProtT5 for the 123k  
 70 proteins in CATH-S100 required 23 minutes (m). The total time  
 71 for creating ProtTucker(ProtT5) embeddings from ProtT5  
 72 embeddings for the same set on the same machine was 0.5  
 73 seconds (s), i.e., ProtTucker added about 0.3%. Creating HMM  
 74 profiles for the same set using either MSAs from MMSeqs2 ( $--$   
 75  $num\_iterations\ 3, -s\ 7.5$ ) or jackhammer took 15m or 30h,  
 76 respectively (Table 3).

77 To predict using EAT, users have to embed only single query  
 78 proteins requiring, on average, 0.01s per protein for the CATH-  
 79 S100 set. Using either single protein sequence search  
 80 (MMSeqs2), pre-computed HMM profiles (HMMER) or pre-

**Table 3: Runtime. \***

Methods	Pre-processing [s]	Inference [s]
MMSeqs2 (sequence)	0.2 * 10 <sup>3</sup>	2.5 10 <sup>-2</sup>
HMMER (HMM)	114 * 10 <sup>3</sup>	150 10 <sup>-2</sup>
ProtTucker(ProtT5)	1.4 * 10 <sup>3</sup>	1.4 10 <sup>-2</sup>

\* Runtime to transfer annotations from CATH-S100 (123k proteins) to a single query. All times measured in seconds [s] on a single Nvidia RTX A6000 with 48GB of vRAM and an AMD EPYC ROME. *Methods*: two HBI-based methods (MMSeqs2 and HMMER) and one EAT-based method (ProtTucker(ProtT5)). *Pre-processing* measured the time required for building datasets (indexed database for MMSeqs2; MSA for jackhammer plus HMM profiles (HMMER) or ProtT5 embeddings (ProtTucker)); *Inference* the time for each new protein with a transfer.



**Table 4: CATH predictions for three model proteomes. \***

Proteome	Size	Gene3D	ProtTucker(ProtT5)@0.9	Agreement	Gene3D-multi	Agreement-multi	Inference time [s]
<i>E. Coli</i> (K12)	2,033	59% (1,193)	97% ( 1,982)	81% (963)	23% ( 275)	86% (235)	113 (0.06)
<i>A. Ostoyae</i>	22,192	31% (6,902)	79% (17,416)	75% (4,707)	18% (1,223)	65% (684)	1,384 (0.06)
<i>M. Chilensis</i>	1,120	35% ( 392)	87% ( 974)	84% (320)	10% ( 40)	73% ( 29)	95 (0.09)

\* Comparison of the annotation-transfer from 123k CATH-S100 proteins (1,2) through HBI (*Gene3D* (5)) and through EAT as introduced here (ProtTucker(ProtT5), or PT(ProtT5)) for three entire reference proteomes: *Escherichia coli* (*E. Coli*), *Armillaria ostoyae* (*A. ostoyae*) and *Megavirus Chilensis* (*M. Chilensis*). In other words, all proteins in the three organisms were mapped to proteins of known structure using the CATH hierarchy. Gene3D predictions were taken from UniProt; PT(ProtT5) predictions were derived from the single nearest neighbor in Euclidean space. Coverage-related numbers refer to the percentage of proteins in the entire proteome (Size; in brackets: actual number of proteins), while those pertaining to the agreement are percentages of the set with annotations. Size: number of proteins; Gene3D: fraction of proteins with Gene3D annotation (coverage); PT(ProtT5)@0.9: coverage of PT(ProtT5) at Euclidean distance  $\leq 0.9$  (5% error rate; Table S3); Agreement: fraction of proteins for which Gene3D and PT(ProtT5) had a prediction and reported the same homologous CATH superfamily (for multi-domain proteins with multiple Gene3D annotations, matching any domain by PT(ProtT5) was considered as correct); Gene3D Multi: fraction of Gene3D proteins with multi-domain annotation; Agreement Multi: fraction of multi-domain proteins for which the homologous CATH superfamily predicted by PT5 agreed with one of the Gene3D domain annotations; Inference time: the total time needed for proteome-wide embedding-based annotation transfer (EAT) measured in seconds [s] on a single Nvidia RTX A6000 with 48GB of vRAM and an AMD EPYC ROME (in brackets the average time per protein).

1 computed embeddings (ProtTucker) to transfer annotations  
2 from CATH-S100 to a single query protein took on average  
3 0.025s, 1.5s or 0.0008s, respectively. Proteins in the PDB and  
4 CATH are, on average, roughly half as long (173 residues) as  
5 those from UniProt (343 residues). This is relevant for runtime,  
6 because embedding generation scales quadratically with  
7 sequence length (Fig. 13 in SOM of (9)).

8 This increase was also reflected for the proteome-wide  
9 annotation transfer (Table 4), although these values included  
10 computations required for all aspects of EAT (1: load ProtT5  
11 embeddings for pre-computed CATH-S100 lookup set; 2: load  
12 ProtT5 and embedding for query proteome; 3: generate  
13 ProtTucker(ProtT5) embeddings for queries and lookup; 4:  
14 compute pairwise Euclidean distances between query/lookup).  
15 We compared EAT using ProtTucker(ProtT5) embeddings to  
16 HBI proxied by existing Gene3D annotations taken from  
17 UniProt for three different proteomes (Table 4). At an expected  
18 error rate of 5% (Euclidean distance  $\leq 0.9$ , Table S3), EAT  
19 predicted substantially more proteins than Gene3D at HMMER  
20 E-value $<10^{-3}$ . For the subset of proteins for which both  
21 methods transferred annotations, those largely agreed  
22 (*Agreement*, Table 4; Table S7 for other thresholds). All values  
23 for coverage decreased for multi-domain proteins, as proxied  
24 by “multiple Gene3D annotations”, while the agreement  
25 between Gene3D and ProtTucker(ProtT5) increased for 2 of 3  
26 proteomes (Table 4: *multi*).

## Discussion

28 **Prototype for representation learning of hierarchies.** We  
29 have presented a new solution for combining the information  
30 implicitly contained in the embeddings from protein Language  
31 Models (pLMs) and contrastive learning to learn directly from  
32 hierarchically sorted data. As proof-of-concept, we applied the  
33 concept to the CATH hierarchy of protein structures (2,3,22,83).  
34 Hierarchies are difficult to handle by traditional supervised  
35 learning solutions. One *shortcut* is to learn each level in the  
36 hierarchy independently (84-86) at the price of having less  
37 information for other levels and of not explicitly benefiting from  
38 the hierarchy. Instead, our solution of contrastively learning  
39 protein triplets (anchor, positive, negative) to extract a new  
40 embedding space by condensing positives and moving  
41 negatives apart benefits from CATH’s hierarchical structure.  
42 Simultaneously training a single, shared feed-forward neural  
43 network (FNN) on triplets from all four CATH classification  
44 levels allowed the network to directly capture the hierarchy.

45 Encoding protein sequences through previously trained pLMs  
46 enabled ready information transfer from large but unlabeled  
47 sequence databases such as BFD (11) to 10,000-times smaller  
48 but experimentally annotated (labeled) proteins of known 3D  
49 structure classified by CATH. In turn, this allowed us to readily  
50 leverage aspects of protein structure captured by pLMs that  
51 are informative enough to predict structure from embeddings  
52 alone (62). Although the raw, pre-trained, unoptimized  
53 embeddings captured some aspects of the classification (Fig.  
54 2A, Fig. 3A, Fig. 4A, Table 1), contrastive learning boosted this  
55 signal significantly (Fig. 2B, Fig. 3C, Fig. 4B, Table 1).

56 Crucial for this success was the novel combination of  
57 hierarchy- and batch-hard sampling (Table 2). Presumably,  
58 because those techniques enforce so-called *semi-hard* triplets  
59 that are neither too simple nor too hard to learn (76). This  
60 training setup learned different classifications for the same  
61 protein pair, depending on the third protein forming the triplet,  
62 thereby forcing the network to learn the complex hierarchy. The  
63 ambivalence in the notion of *positive/negative pair* facilitated  
64 training by allowing to include superfamilies with few members  
65 (otherwise to be skipped) and it increased the number of  
66 possible triplets manifold compared to only sampling on the  
67 level of superfamilies. These advantages might partially explain  
68 the synergy of both sampling techniques (Table 2).  
69

70 **Raw embedding EAT matched profile alignments in hit  
71 detection.** In technical analogy to homology-based inference  
72 (HBI), we used embedding based annotation transfer (EAT, Fig.  
73 1B) to transfer annotations from labeled lookup proteins  
74 (proteins with a known CATH classification) to unlabeled query  
75 proteins (any protein of known sequence without known  
76 structure). Instead of transferring annotations from the closest  
77 hit in sequence space, EAT transferred annotations to the hit  
78 with smallest Euclidean distance in embedding space. This  
79 relatively simple approach was shown previously to predict  
80 protein function as defined by Gene Ontology (GO) better than  
81 hand-crafted features (60) even to levels competitive to much  
82 more complex approaches (58).

83 The concept of EAT was so successful that raw embeddings  
84 from two different pre-trained pLMs (ESM-1b (12), and ProtT5  
85 (9)) already set the bar high for predicting CATH levels. The raw,  
86 general-purpose ESM-1b and ProtT5 outperformed HBI based  
87 on advanced HMM-profiles from HMMER (6) on the C- and A-  
88 level while falling short on the H-level (Table 1). Furthermore,  
89 we showed that ProtT5 already separated protein pairs with the  
90 same from those with different homologous superfamilies even

1 when using a lookup set that consisted only of proteins with  
2 maximally 30% pairwise sequence identity (Fig 4A).  
3 Importantly, this competitive performance was achieved at a  
4 much smaller cost in terms of runtime (Table 3, Table 4).

5 As the lookup embeddings or HMM profiles are computed  
6 only once, we neglected this additional step. Such preparations  
7 cost much more than single queries: pre-computing HMM  
8 profiles using MMseqs2 took 15m, pre-computing embeddings  
9 about 23m (Table 3) using the same set and machine but  
10 utilizing CPUs in one (MMseqs2) and GPUs in the other  
11 (ProtT5). Only MMSeqs2 generated and indexed its database  
12 rapidly (19.5s). However, pre-processing is required only once,  
13 rapidly amortizing when running many queries. The ability to  
14 pre-compute such representations is also a crucial difference  
15 between ProtTucker and other learned methods (16,68). For  
16 pairwise protein comparisons, those methods typically require  
17  $N$  comparisons/forward-passes to search with a single query  
18 against  $N$  proteins. Instead, ProtTucker only needs a single  
19 forward pass to embed the new query; subsequent similarity  
20 scoring simply and quickly computes an Euclidean distance.

21 This makes ProtTucker search speed scale well with  
22 database growth suggesting the tool as a fast but sensitive pre-  
23 filter for other methods that in turn provide residue-level  
24 information as showcased on three model organisms (Table 4),  
25 including one of the largest organisms on earth (fungus *A.*  
26 *Ostoyae*, 22,192 proteins) and one of the largest viruses (*M.*  
27 *Chilensis*, 1,120 proteins). In less than 27 min on a single  
28 machine (Table 4), ProtTucker transferred substantially more  
29 CATH annotations mapping proteins from their sequence to 3D  
30 structures through the CATH resource than Gene3D (5) at a  
31 similar level of expected error (Table 4).

32 For the virus and the bacterium (*E. coli*) fungus the  
33 annotations agreed to over 80% with Gene3D, while this value  
34 dropped to 75% for the fungus (Table 4). Although high, the  
35 agreement was lower than expected: if ProtTucker and Gene3D  
36 each had fewer than 5% errors, then both should agree for over  
37 90% of the proteins for which both transfer annotations. Most  
38 likely, this discrepancy ( $\Delta(90-80)$ ) arose partially from multi-  
39 domain proteins. Despite carefully cross-validating ProtTucker,  
40 an alternative explanation for the discrepancy is underestima-  
41 ting the expected error a distances  $\leq 0.9$  as 5% instead of up to  
42 15%. The “functional shape” of the agreement between  
43 ProtTucker and Gene3D at different distance thresholds (Table  
44 S7) suggested that the “errors” (lack of agreement) did not only  
45 originate from ProtTucker. Carefully annotating the five proteins  
46 with the lowest distance and a different CATH annotation  
47 (Table S4) supported this perspective.

48 The agreement for multi-domain proteins dropped less than  
49 expected (11 percentage points drop for *M. Chilensis*, 5  
50 percentage points increase for *E. coli*), possibly suggesting that  
51 ProtTucker using averages over an entire protein for  
52 comparison did not trip substantially more over the multi-  
53 domain challenge than the local alignment-based Gene3D  
54 using HMMER (6). This might suggest ProtTucker to have  
55 added correct annotations over Gene3D in multi-domain  
56 proteins, although developed exclusively on single domain  
57 proteins. The substantial increase in coverage from the level  
58 expected at distances  $\geq 0.9$  (Fig. 5, Table S4) for the proteomes  
59 (Table 4) might be misleading: to establish performance  
60 coverage (Fig. 5, Table S4), we used a highly non-redundant  
61 lookup set, presumably removing many easy hits. In contrast,  
62 analyzing proteomes, we transferred annotations for all CATH-  
63 S100 proteins, leveraging “redundant annotation transfers” to  
64 increase coverage.

65 As for HBI, the accuracy of EAT also increased for larger  
66 families (Fig. S1). One explanation is that the larger the family,

67 the higher the random hit rate, simply because there are more  
68 possible hits. Another, more subtle (and given the enormous  
69 compute time needed to train ProtT5, more difficult to test)  
70 explanation is that the largest CATH families represent most of  
71 the largest protein families (2). In fact, a few hundred of the  
72 largest superfamilies cover half of the entire sequence space  
73 (2,88). Simply due to their immense size, these large families  
74 have been sampled more during the pre-training of ProtT5.

75  
76 **ProtTucker embeddings intruded into midnight zone.** The  
77 embedding space resulting from contrastive learning,  
78 introduced here, improved performance consistently for all four  
79 pLMs (Table 1). This was revealed through several ways of  
80 looking at the results from embeddings with and without  
81 contrastive learning: (1) the increased separation of protein  
82 pairs within the same protein superfamily and between different  
83 superfamilies (Fig. 2), (2) the qualitative improvement in the  
84 clustering of t-SNE projections (Fig. 3), the better correlation of  
85 embedding distance and structural similarity (Fig. 4) and (3) the  
86 quantitative improvement in the EAT benchmark (Table 1). On  
87 top, the Euclidean distance correlated with accuracy (Fig. 5,  
88 Table S3). Similar to an E-value in HBI, this lets users gauge the  
89 reliability of a hit between query and annotated protein.

90 While the accuracy of the best performing pLM  
91 (ProtTucker(ProtT5)) was similar to HBI using HMM-profiles on  
92 the most fine-grained level of homologous superfamilies (CATH  
93 level H, Table 1), the relative advantage of EAT increased, the  
94 more diverged the level of inference, i.e., EAT outperformed  
95 HBI for more distant relations from the midnight zone (CATH  
96 level C, Table 1). When further reducing data redundancy, i.e.,  
97 removing more similar sequences, this trend became clearer  
98 (Fig. 6). Despite increasing difficulty, the performance of EAT  
99 decreased almost insignificantly where HBI approached  
100 random for insignificant E-values. This trend was supported by  
101 the correlation of structural similarity as defined by SSAP (7,8)  
102 and the Euclidean distance between protein pairs in a 30%  
103 non-redundant data set (Fig. 4).

104 ProtTucker and tools such as HMMer have very different  
105 resolution: ProtTucker considers only per-protein averages to  
106 match query to template. In contrast, HMMer – or similar  
107 methods – align each residue between both proteins. The  
108 coarse-grain yields the speedup (Table 4), and pitches  
109 ProtTucker as a fast pre-filter. Once the hit is found by scanning  
110 large data sets, the slower, fine-grained methods for per-  
111 residue alignments and 3D prediction can be employed.  
112 However, the per-protein average also implies limitations, e.g.,  
113 when Q and T have very different numbers of domains or the  
114 number of domains for Q is not known (Table 3).

115 Ultimately, the coarse-grained ProtTucker can compete at all  
116 because embeddings intrinsically abstract the constraints  
117 under which protein sequences evolve, including constraints  
118 upon structure, function, and the environment. The same  
119 constraints coin the evolutionary information contained in  
120 profiles of protein families. Apparently, pLMs such as ESM-1b  
121 (12), ProtBERT (18), or ProtT5 (18) are successfully condense  
122 these constraints. In fact, pLMs are arguably more successful  
123 than profile-based methods because a simple length-average  
124 over the position-specific scoring metrics (PSSM) would not  
125 suffice to predict CATH numbers very accurately.

126 ProtTucker builds upon this success to explicitly capture the  
127 constraints relevant for the CATH hierarchy. Thus, the less a  
128 particular aspect of function depends on structure, the less  
129 likely the new ProtTucker embeddings will reflect this aspect.  
130 On the other hand, an approach similar to ProtTucker focused  
131 on particular functional hierarchies, e.g., EC numbers appears  
132 to work well (SM Akmes & M Heinzinger, unpublished).



1 Taken together, these results indicated that contrastive  
2 learning captured structural hierarchies and provides a novel,  
3 powerful tool to uncover structural similarities clearly beyond  
4 what has been achievable with 50 years of optimizing  
5 sequence-based alignment techniques. Using EAT to  
6 complement HBI could become crucial for a variety of  
7 applications, ranging from finding remote structural templates  
8 for protein 3D structure predictions over prioritizing new  
9 proteins without any similarity to an existing structure to  
10 filtering potentially wrong annotations. One particular example  
11 has recently been shown for the proteome of SARS-CoV-2 to  
12 unravel entire functional components possibly relevant for  
13 fighting COVID-19 (70).

#### 14 **ProtTucker embeddings improved FunFams clustering.**

15 Previously (59), we showed that a simplistic predecessor of  
16 ProtTucker helped to refine the clustering of FunFams (79). By  
17 adding an additional, more fine-grained hierarchy level in  
18 CATH, FunFams link the structure-function continuum of  
19 proteins. The functional consistency within FunFams was  
20 proxied through the enzymatic activity as defined by the EC  
21 (Enzyme Commission (82)) number. Even the preliminary  
22 ProtTucker improved the annotation transfer of ligand binding  
23 and EC numbers (59) by removing outliers from existing  
24 FunFams and by creating new, more functionally coherent  
25 FunFams. As for CATH, the contrastively trained  
26 ProtTucker(ProtBERT) also improved over its unsupervised  
27 counterpart, ProtBERT, for FunFams. It improved functional  
28 consistency especially for proteins in the twilight zone (<35%  
29 PIDE, Fig. 5 in (59)). Thus, ProtTucker embeddings improved  
30 functional (FunFams) and structural (CATH) consistency  
31 beyond sequence similarity. Here, we expanded upon this  
32 analysis by showing how EAT can be improved even more  
33 through contrastively learning hierarchies. Using the proposed  
34 method, we could spot potential outliers, i.e., samples with the  
35 same annotation but large embedding distance. This might  
36 become essential to clean up databases. Aside from outlier-  
37 spotting, we could also obtain labels from the nearest  
38 neighbors of outliers (Table S6). Although we could not  
39 reproduce the same level of success when applying EAT to  
40 inferring subcellular location in ten states (Table 3), the CATH-  
41 optimized ProtTucker embeddings also did not perform worse.  
42  
43

#### 44 **Generic advantages of contrastive learning.**

45 Contrastive learning benefits from hierarchies as opposed to supervised  
46 training which usually flattens the hierarchy thereby losing its  
47 intrinsic advantage. Other possible advantages of contrastive  
48 learning include the following three. (1) Dynamic data update  
49 (online learning): While supervised networks require re-training  
50 to benefit from new data, contrastively trained networks can  
51 benefit from new data by simply updating the lookup set. This  
52 could even add completely new classes, such as proteins for  
53 which the classification will become available only in the future.  
54 HBI shares this advantage that originates from the difference  
55 between classifying proteins into existing families versus  
56 classifying by identifying the most similar proteins in that family.  
57 (2) Learn the access, not the data: Instead of forcing the  
58 supervised network to memorize the training data, contrastive  
59 learning teaches how to access the data stored in an external  
60 lookup set. (3) Compression: As many other learning  
61 techniques, contrastive learning can act as a compression  
62 technique. For instance, we reduced the disk space required to  
63 store protein embeddings threefold by projecting 1024-  
64 dimensional vectors onto 128 dimensions while improving  
65 performance (Table 1). This renders new queries (inference)  
66 more efficient and enables scaling up to very large lookup sets.

67 (4) Interpretability: knowing from which protein an annotation  
68 was transferred might help users benefit more from a certain  
69 prediction than just the prediction itself. For instance, knowing  
70 that an unnamed query protein shares all CATH levels with a  
71 particular glucocorticoid receptor might suggest some  
72 functional implications helping to design future experiments.  
73

74 Gaining from embeddings for sequence comparisons. The  
75 solution described here established the power for contrastive  
76 learning. Other, more straightforward solutions

## 77 **Conclusions**

78 Embeddings from protein Language Models (pLMs) extract the  
79 information learned by these models from unlabeled protein  
80 sequences. Embedding-based Annotation Transfer (EAT)  
81 replacing the proximity in sequence space used by homolog-  
82 based inference (HBI) through proximity in embedding space  
83 already reaches traditional alignment methods in transferring  
84 CATH annotations from a template protein with experimental  
85 annotations to an unlabeled query protein. Although not quite  
86 reaching the performance of advanced profile-profile searches  
87 by HMMer for all four CATH levels, the best embeddings  
88 surpassed HMMer for two of the four levels (C and A). When  
89 optimizing embeddings through contrastive learning for the  
90 goal of transferring CATH annotations, EAT using these new  
91 embeddings consistently outperformed all sequence  
92 comparison techniques tested. This higher performance was  
93 reached at a fraction (three orders of magnitude) of the  
94 computational time. Although the new embeddings optimized  
95 through contrastive learning for CATH did not improve  
96 performance for a completely different task, namely the  
97 prediction of subcellular location in ten classes, the CATH-  
98 optimized solution did also not perform significantly worse.  
99 Remarkably, just like HBI, the performance of EAT using the  
100 optimized ProtTucker embeddings was proportional to family  
101 size with increased accuracy for larger families.  
102

103 **Availability:** Building on top of bio\_embeddings package (89)  
104 we have made a script available that simplifies EAT  
105 <https://github.com/Rostlab/EAT>.  
106

## 107 **Acknowledgements**

108 Thanks primarily to Tim Karl, but also to Guy Yachdav & Laszlo  
109 Kajan (all TUM) for invaluable help with hardware and software;  
110 to Inga Weise (TUM) for support with many other aspects of  
111 this work. Many thanks for both anonymous reviewers, deep  
112 kudos for the detailed help from No. 1! This work was  
113 supported by the Bavarian Ministry of Education through  
114 funding to the TUM and by a grant from the Alexander von  
115 Humboldt foundation through the German Ministry for  
116 Research and Education (BMBF: Bundesministerium für  
117 Bildung und Forschung), by two grants from BMBF (031L0168  
118 and program "Software Campus 2.0 (TUM) 2.0" 01IS17049) as  
119 well as by a grant from Deutsche Forschungsgemeinschaft  
120 (DFG-GZ: RO1320/4-1). Last, but not least, thanks to all those  
121 who maintain public databases in particular Steven Burley  
122 (PDB, Rutgers), Alan Bridge (Swiss-Prot, SIB, Lausanne), Alex  
123 Bateman (UniProt, EBI Hinxton) and their crews, and to all  
124 experimentalists who enabled this analysis by making their  
125 data publicly available.

## References

- 1
- 2 1. Orengo, C.A., Flores, T.P., Taylor, W.R. and Thornton, J.M. (1993)
- 3 Identification and classification of protein fold families. *Protein Eng*,
- 4 **6**, 485-500.
- 5 2. Sillitoe, I., Bordin, N., Dawson, N., Waman, V.P., Ashford, P.,
- 6 Scholes, H.M., Pang, C.S.M., Woodridge, L., Rauer, C., Sen, N. *et*
- 7 *al.* (2021) CATH: increased structural coverage of functional space.
- 8 *Nucleic Acids Res*, **49**, D266–D273.
- 9 3. Das, S., Sillitoe, I., Lee, D., Lees, J.G., Dawson, N.L., Ward, J. and
- 10 Orengo, C.A. (2015) CATH FunFHMmer web server: protein
- 11 functional annotations using functional family assignments. *Nucleic*
- 12 *Acids Res*, **43**, W148-153.
- 13 4. Van der Maaten, L. and Hinton, G. (2008) Visualizing data using t-
- 14 SNE. *J Mach Learning Res*, **9**, 2579-2605.
- 15 5. Lewis, T.E., Sillitoe, I., Dawson, N., Lam, S.D., Clarke, T., Lee, D.,
- 16 Orengo, C. and Lees, J. (2018) Gene3D: extensive prediction of
- 17 globular domains in proteins. *Nucleic Acids Res*, **46**, D435-D439.
- 18 6. Finn, R.D., Clements, J. and Eddy, S.R. (2011) HMMER web server:
- 19 interactive sequence similarity searching. *Nucleic Acids Res*, **39**,
- 20 W29-37.
- 21 7. Taylor, W.R. and Orengo, C.A. (1989) A holistic approach to protein
- 22 structure alignment. *Protein Eng*, **2**, 505-519.
- 23 8. Orengo, C.A. and Taylor, W.R. (1996) SSAP: Sequential structure
- 24 alignment program for protein structure comparison. *Meth*
- 25 *Enzymol*, **266**, 617-635.
- 26 9. Elnaggar, A., Heinzinger, M., Dallago, C., Rehawi, G., Wang, Y.,
- 27 Jones, L., Gibbs, T., Feher, T., Angerer, C., Steinegger, M. *et al.*
- 28 (2021) ProtTrans: Towards Cracking the Language of Life's Code
- 29 Through Self-Supervised Learning. *IEEE TPAMI*, **14**, 30.
- 30 10. Steinegger, M. and Soding, J. (2017) MMseqs2 enables sensitive
- 31 protein sequence searching for the analysis of massive data sets.
- 32 *Nat Biotechnol*, **35**, 1026-1028.
- 33 11. Steinegger, M., Mirdita, M. and Soding, J. (2019) Protein-level
- 34 assembly increases protein sequence recovery from metagenomic
- 35 samples manyfold. *Nat Methods*, **16**, 603-606.
- 36 12. Rives, A., Meier, J., Sercu, T., Goyal, S., Lin, Z., Liu, J., Guo, D., Ott,
- 37 M., Zitnick, C.L., Ma, J. *et al.* (2021) Biological structure and
- 38 function emerge from scaling unsupervised learning to 250 million
- 39 protein sequences. *PNAS*, **118**.
- 40 13. Das, S., Lee, D., Sillitoe, I., Dawson, N.L., Lees, J.G. and Orengo,
- 41 C.A. (2015) Functional classification of CATH superfamilies: a
- 42 domain-based approach for protein function annotation.
- 43 *Bioinformatics*, **31**, 3460-3467.
- 44 14. Sander, C. and Schneider, R. (1991) Database of homology-derived
- 45 structures and the structural meaning of sequence alignment.
- 46 *Proteins*, **9**, 56-68.
- 47 15. Mika, S. and Rost, B. (2003) UniqueProt: creating representative
- 48 protein sequence sets. *Nucleic Acids Res*, **31**, 3789-3791.
- 49 16. Bepler, T. and Berger, B. (2021) Learning the protein language:
- 50 Evolution, structure, and function. *Cell Syst*, **12**, 654-669 e653.
- 51 17. Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2019),
- 52 *Proceedings of the 2019 Conference of the North American Chapter*
- 53 *of the Association for Computational Linguistics: Human Language*
- 54 *Technologies, Volume 1 (Long and Short Papers)*. Association for
- 55 Computational Linguistics, pp. 4171–4186.
- 56 18. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M.,
- 57 Zhou, Y., Li, W. and Liu, P.J. (2020) Exploring the Limits of Transfer
- 58 Learning with a Unified Text-to-Text Transformer. *J Mach Learning*
- 59 *Res*, **21**, 1-67.
- 60 19. Sonnhammer, E.L. and Kahn, D. (1994) Modular arrangement of
- 61 proteins as inferred from analysis of homology. *Protein Sci*, **3**, 482-
- 62 492.
- 63 20. Bateman, A., Birney, E., Durbin, R., Eddy, S.R., Finn, R.D. and
- 64 Sonnhammer, E.L. (1999) Pfam 3.1: 1313 multiple alignments and
- 65 profile HMMs match the majority of proteins. *Nucleic Acids Res*, **27**,
- 66 260-262.
- 67 21. Gough, J. and Chothia, C. (2002) SUPERFAMILY: HMMs
- 68 representing all proteins of known structure. SCOP sequence
- 69 searches, alignments and genome assignments. *Nucleic Acids Res*,
- 70 **30**, 268-272.
- 71 22. Orengo, C.A., Michie, A.D., Jones, S., Jones, D.T., Swindells, M.B.
- 72 and Thornton, J.M. (1997) CATH - a hierarchic classification of
- 73 protein domain structures. *Structures*, **5**, 1093-1108.
- 74 23. Todd, A.E., Orengo, C.A. and Thornton, J.M. (2001) Evolution of
- 75 function in protein superfamilies, from a structural perspective. *J*
- 76 *Mol Biol*, **307**, 1113-1143.
- 77 24. Yona, G. and Levitt, M. (2002) Within the twilight zone: a sensitive
- 78 profile-profile comparison tool based on information theory. *J Mol*
- 79 *Biol*, **315**, 1257-1275.
- 80 25. Doolittle, R.F., Feng, D.-F., Johnson, M.S. and McClure, M.A. (1986)
- 81 Origins and evolutionary relationships of retroviruses. *Q Rev Biol*,
- 82 **64**, 1-30.
- 83 26. Rost, B. (1999) Twilight zone of protein sequence alignments.
- 84 *Protein Eng*, **12**, 85-94.
- 85 27. Rost, B. (1997) Protein structures sustain evolutionary drift. *Folding*
- 86 *& Design*, **2**, S19-S24.
- 87 28. Mika, S. and Rost, B. (2003) UniqueProt: creating representative
- 88 protein sequence sets. *Nucleic Acids Res*, **31**, 3789-3791.
- 89 29. Rost, B. (2002) Enzyme function less conserved than anticipated. *J*
- 90 *Mol Biol*, **318**, 595-608.
- 91 30. Nehrt, N.L., Clark, W.T., Radivojac, P. and Hahn, M.W. (2011)
- 92 Testing the ortholog conjecture with comparative functional
- 93 genomic data from mammals. *PLoS Comput Biol*, **7**, e1002073.
- 94 31. Higgins, D.G., Bleasby, A.J. and Fuchs, R. (1992) CLUSTAL V:
- 95 improved software for multiple sequence alignment. *CABIOS*, **8**,
- 96 189-191.
- 97 32. Thompson, J., Higgins, D. and Gibson, T. (1994) CLUSTAL W:
- 98 improving the sensitivity of progressive multiple sequence
- 99 alignment through sequence weighting, position-specific gap
- 100 penalties and weight matrix choice. *Nucleic Acids Res*, **22**, 4673-
- 101 4690.
- 102 33. Sjölander, K., Karplus, K., Brown, M.P., Hughey, R., Krogh, A.,
- 103 Mian, I.S. and Haussler, D. (1996) Dirichlet mixtures: a method for
- 104 improving detection of weak but significant protein sequence
- 105 homology. *CABIOS*, **12**, 327-345.
- 106 34. Altschul, S.F., Madden, T.L., Schaeffer, A.A., Zhang, J., Zhang, Z.,
- 107 Miller, W. and Lipman, D.J. (1997) Gapped Blast and PSI-Blast: a
- 108 new generation of protein database search programs. *Nucleic*
- 109 *Acids Res*, **25**, 3389-3402.
- 110 35. Eddy, S.R. (1998) Profile hidden Markov models. *Bioinformatics*, **14**,
- 111 755-763.
- 112 36. Jaroszewski, L., Rychlewski, L. and Godzik, A. (2000) Improving the
- 113 quality of twilight-zone alignments. *Protein Sci*, **9**, 1487-1496.
- 114 37. Sadreyev, R. and Grishin, N. (2003) COMPASS: a tool for
- 115 comparison of multiple protein alignments with assessment of
- 116 statistical significance. *J Mol Biol*, **326**, 317-336.
- 117 38. Edgar, R.C. and Sjolander, K. (2004) COACH: profile-profile
- 118 alignment of protein families using hidden Markov models.
- 119 *Bioinformatics*, **20**, 1309-1318.
- 120 39. Wang, G. and Dunbrack, R.L., Jr. (2004) Scoring profile-to-profile
- 121 sequence alignments. *Protein Sci*, **13**, 1612-1626.
- 122 40. Soding, J. (2005) Protein homology detection by HMM-HMM
- 123 comparison. *Bioinformatics*, **21**, 951-960.
- 124 41. Sievers, F., Wilm, A., Dineen, D., Gibson, T.J., Karplus, K., Li, W.,
- 125 Lopez, R., McWilliam, H., Remmert, M., Soding, J. *et al.* (2011) Fast,
- 126 scalable generation of high-quality protein multiple sequence
- 127 alignments using Clustal Omega. *Mol Syst Biol*, **7**, 539.
- 128 42. Przybylski, D. and Rost, B. (2007) Consensus sequences improve
- 129 PSI-BLAST through mimicking profile-profile alignments. *Nucleic*
- 130 *Acids Res*, **35**, 2238-2246.
- 131 43. Rost, B., Liu, J., Nair, R., Wrzeszczynski, K.O. and Ofra, Y. (2003)
- 132 Automatic prediction of protein function. *Cellular and Molecular Life*
- 133 *Sciences*, **60**, 2637-2650.
- 134 44. Rost, B. (1996) PHD: predicting one-dimensional protein structure
- 135 by profile based neural networks. *Meth Enzymol*, **266**, 525-539.
- 136 45. Rost, B. and Sander, C. (1993) Prediction of protein secondary
- 137 structure at better than 70% accuracy. *J Mol Biol*, **232**, 584-599.
- 138 46. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M.,
- 139 Ronneberger, O., Tunyasuvunakool, K., Bates, R., Zidek, A.,
- 140 Potapenko, A. *et al.* (2021) Highly accurate protein structure
- 141 prediction with AlphaFold. *Nature*, **599**, 583-589.
- 142 47. Baek, M., Dimaio, F., Anishchenko, I., Dauparas, J., Ovchinnikov,
- 143 S., Lee, G.R., Wang, J., Cong, Q., Kinch, L.N., Schaeffer, R.D. *et al.*



- 1 (2021) Accurate prediction of protein structures and interactions  
2 using a three-track neural network. *Science*, **373**, 871-876.
- 3 48. Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee,  
4 K. and Zettlemoyer, L. (2018) Deep contextualized word  
5 representations. *arXiv*, 10.48550/arXiv.2005.14165.
- 6 49. Brown, T.B., Mann, B., Ryder, N., Subbiah, N., Kaplan, J., Dhariwal,  
7 P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A. et al. (2020)  
8 Language Models are Few-Shot Learners. *arXiv*, 1802.05365.
- 9 50. Alley, E.C., Khimulya, G., Biswas, S., AlQuraishi, M. and Church,  
10 G.M. (2019) Unified rational protein engineering with sequence-  
11 based deep representation learning. *Nature Meth*, **16**, 1315-1322.
- 12 51. Heinzinger, M., Elnaggar, A., Wang, Y., Dallago, C., Nechaev, D.,  
13 Matthes, F. and Rost, B. (2019) Modeling aspects of the language  
14 of life through transfer-learning protein sequences. *BMC  
15 Bioinformatics*, **20**, 723.
- 16 52. Rao, R., Meier, J., Sercu, T., Ovchinnikov, S. and Rives, A. (2020)  
17 Transformer protein language models are unsupervised structure  
18 learners. *bioRxiv*, 2020.2012.2015.422761.
- 19 53. Madani, A., McCann, B., Naik, N., Shirish Keskar, N., Anand, N.,  
20 Eguchi, R.R., Huang, P. and Socher, R. (2020) ProGen: Language  
21 Modeling for Protein Generation. *arXiv*, 2004.03497.
- 22 54. Ofer, D., Brandes, N. and Linial, M. (2021) The language of proteins:  
23 NLP, machine learning & protein sequences. *Comp Structural  
24 Biotechn J*, **19**, 1750-1758.
- 25 55. Bepler, T. and Berger, B. (2019), *Seventh International Conference  
26 on Learning Representations*. 2019/02/22 ed.
- 27 56. Stärk, H., Dallago, C., Heinzinger, M. and Rost, B. (2021) Light  
28 Attention Predicts Protein Location from the Language of Life.  
29 *bioRxiv*, 2021.2004.2025.441334.
- 30 57. Littmann, M., Heinzinger, M., Dallago, C., Weissenow, K. and Rost,  
31 B. (2021) Protein embeddings and deep learning predict binding  
32 residues for various ligand classes. *Sci Rep*, **11**, 23916.
- 33 58. Littmann, M., Heinzinger, M., Dallago, C., Olenyi, T. and Rost, B.  
34 (2021) Embeddings from deep learning transfer GO annotations  
35 beyond homology. *Sci Rep*, **11**, 1160.
- 36 59. Littmann, M., Bordin, N., Heinzinger, M., Schütze, K., Dallago, C.,  
37 Orengo, C. and Rost, B. (2021) Clustering FunFams using sequence  
38 embeddings improves EC purity *Bioinformatics*, **37**, 3449-3455.
- 39 60. Villegas-Morcillo, A., Makrodimitris, S., van Ham, R.C.H.J., Gomez,  
40 A.M., Sanchez, V. and Reinders, M.J.T. (2021) Unsupervised  
41 protein embeddings outperform hand-crafted sequence and  
42 structure features at predicting molecular function. *Bioinformatics*,  
43 **37**, 162-170.
- 44 61. Hamid, M.-N. and Friedberg, I. (2019) Identifying antimicrobial  
45 peptides using word embedding with deep recurrent neural  
46 networks. *Bioinformatics*, **35**, 2009-2016.
- 47 62. Weissenow, K., Heinzinger, M. and Rost, B. (2021) Protein language  
48 model embeddings for fast, accurate, alignment-free protein  
49 structure prediction. *bioRxiv*, 2021.2007.2031.454572.
- 50 63. Le-Khac, P.H., Healy, G. and Smeaton, A.F. (2020) Contrastive  
51 representation learning: A framework and review. *IEEE Access*.
- 52 64. Fox, N.K., Brenner, S.E. and Chandonia, J.-M. (2014) SCOPe:  
53 Structural Classification of Proteins—extended, integrating SCOP  
54 and ASTRAL data and classification of new structures. *Nucleic  
55 Acids Res*, **42**, D304-D309.
- 56 65. Nallapareddy, V., Bordin, N., Sillitoe, I., Heinzinger, M., Littmann,  
57 M., Waman, V., Sen, N., Rost, B. and Orengo, C. (2022) CATHe:  
58 Detection of remote homologues for CATH superfamilies using  
59 embeddings from protein language models. *bioRxiv* doi  
60 0.1101/2022.03.10.483805.
- 61 66. Li, C.-C. and Liu, B. (2019) MotifCNN-fold: protein fold recognition  
62 based on fold-specific features extracted by motif-based  
63 convolutional neural networks. *Brief Bioinform*, **21**, 2133-2141.
- 64 67. Liu, B., Li, C.-C. and Yan, K. (2020) DeepSVM-fold: protein fold  
65 recognition by combining support vector machines and pairwise  
66 sequence similarity scores generated by deep learning networks.  
67 *Brief Bioinform*, **21**, 1733-1741.
- 68 68. Gao, M. and Skolnick, J. (2021) A novel sequence alignment  
69 algorithm based on deep learning of the protein folding code.  
70 *Bioinformatics*, **37**, 490-496.
- 71 69. Chen, J., Guo, M., Wang, X. and Liu, B. (2018) A comprehensive  
72 review and comparison of different computational methods for  
73 protein remote homology detection. *Brief Bioinform*, **19**, 231-244.
- 74 70. O'Donoghue, S.I., Schafferhans, A., Sikta, N., Stolte, C., Kaur, S.,  
75 Ho, B.K., Anderson, S., Procter, J., Dallago, C., Bordin, N. et al.  
76 (2021) SARS-CoV-2 structural coverage map reveals viral protein  
77 assembly, mimicry, and hijacking mechanisms. *Molecular Systems  
78 Biology* **12**, e10079.
- 79 71. Burley, S.K., Berman, H.M., Bhikadiya, C., Bi, C., Chen, L., Di  
80 Costanzo, L., Christie, C., Dalenberg, K., Duarte, J.M., Dutta, S. et  
81 al. (2019) RCSB Protein Data Bank: biological macromolecular  
82 structures enabling research and education in fundamental biology,  
83 biomedicine, biotechnology and energy. *Nucleic Acids Res*, **47**,  
84 D464-D474.
- 85 72. Chen, T., Kornblith, S., Norouzi, M. and Hinton, G. (2020),  
86 *International conference on machine learning*. PMLR, pp. 1597-  
87 1607.
- 88 73. Almagro Armenteros, J.J., Sonderby, C.K., Sonderby, S.K., Nielsen,  
89 H. and Winther, O. (2017) DeepLoc: prediction of protein subcellular  
90 localization using deep learning. *Bioinformatics*, **33**, 3387-3395.
- 91 74. The UniProt Consortium. (2021) UniProt: the universal protein  
92 knowledgebase in 2021. *Nucleic Acids Res*, **49**, D480-D489.
- 93 75. Marquet, C., Heinzinger, M., Olenyi, T., Dallago, C., Erckert, K.,  
94 Bernhofer, M., Nechaev, D. and Rost, B. (2021) Embeddings from  
95 protein language models predict conservation and variant effects.  
96 *Human Genetics*, 10.21203/rs.3.rs-584804/v1.
- 97 76. Hermans, A., Beyer, L. and Leibe, B. (2017) In defense of the triplet  
98 loss for person re-identification. *arXiv*, 170307737.
- 99 77. Kingma, D.P. and Ba, J. (2014) Adam: A method for stochastic  
100 optimization. *arXiv*, 1412.6980.
- 101 78. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B.,  
102 Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R. and Dubourg, V.  
103 (2011) Scikit-learn: Machine learning in Python. *J Mach Learning  
104 Res*, **12**, 2825-2830.
- 105 79. Sillitoe, I., Cuff, A.L., Dessailly, B.H., Dawson, N.L., Furnham, N.,  
106 Lee, D., Lees, J.G., Lewis, T.E., Studer, R.A., Rentzsch, R. et al.  
107 (2013) New functional families (FunFams) in CATH to improve the  
108 mapping of conserved functional sites to 3D structures. *Nucleic  
109 Acids Res*, **41**, D490-498.
- 110 80. Peng, L., Oganessian, V., Damschroder, M.M., Wu, H. and  
111 Dall'Acqua, W.F. (2011) Structural and functional characterization  
112 of an agonistic anti-human EphA2 monoclonal antibody. *J Mol Biol*,  
113 **413**, 390-405.
- 114 81. Himanen, J.P., Goldgur, Y., Miao, H., Myshkin, E., Guo, H., Buck,  
115 M., Nguyen, M., Rajashankar, K.R., Wang, B. and Nikolov, D.B.  
116 (2009) Ligand recognition by A-class Eph receptors: crystal  
117 structures of the EphA2 ligand-binding domain and the  
118 EphA2/ephrin-A1 complex. *EMBO Rep*, **10**, 722-728.
- 119 82. Webb, E.C. (1992) *Enzyme Nomenclature 1992. Recommendations  
120 of the Nomenclature committee of the International Union of  
121 Biochemistry and Molecular Biology*. 1992 ed. Academic Press,  
122 New York.
- 123 83. Sillitoe, I., Dawson, N., Lewis, T.E., Das, S., Lees, J.G., Ashford, P.,  
124 Touloupe, A., Scholes, H.M., Senatorov, I., Bujan, A. et al. (2019)  
125 CATH: expanding the horizons of structure-based functional  
126 annotations for genome sequences. *Nucleic Acids Res*, **47**, D280-4.
- 127 84. Jensen, L.J., Gupta, R., Blom, N., Devos, D., Tamames, J., Kesmir,  
128 C., Nielsen, H., Staerfeldt, H.H., Rapacki, K., Workman, C. et al.  
129 (2002) Prediction of human protein function from post-translational  
130 modifications and localization features. *J Mol Biol*, **319**, 1257-1265.
- 131 85. Nair, R. and Rost, B. (2005) Mimicking cellular sorting improves  
132 prediction of subcellular localization. *J Mol Biol*, **348**, 85-100.
- 133 86. Kernytsky, A. and Rost, B. (2009) Using genetic algorithms to select  
134 most predictive protein features. *Proteins*, **75**, 75-88.
- 135 87. Liu, J. and Rost, B. (2004) CHOP proteins into structural domain-  
136 like fragments. *Proteins*, **55**, 678-688.
- 137 88. Dessailly, B.H., Nair, R., Jaroszewski, L., Fajardo, J.E., Kouranov,  
138 A., Lee, D., Fiser, A., Godzik, A., Rost, B. and Orengo, C. (2009)  
139 PSI-2: Structural genomics to cover protein domain family space.  
140 *Structure*, **17**, 869-881.
- 141 89. Dallago, C., Schuetze, K., Heinzinger, M., Olenyi, T., Littmann, M.,  
142 Lu, A.X., Yang, K.K., Min, S., Yoon, S., Morton, J.T. et al. (2021)  
143 Learned embeddings from deep learning to visualize and predict  
144 protein sets. *Curr Protoc*, **1**, e113.