
Light Attention Predicts Protein Location from the Language of Life

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

Anonymous Authors¹

Abstract

Although knowing where a protein functions in a cell is important to characterize biological processes, this information remains unavailable for most known proteins. Machine learning narrows the gap through predictions from expertly chosen input features leveraging evolutionary information that is resource expensive to generate. We showcase using embeddings from protein language models for competitive localization predictions not relying on evolutionary information. Our lightweight deep neural network architecture uses a softmax weighted aggregation mechanism with linear complexity in sequence length referred to as light attention (LA). The method significantly outperformed the state-of-the-art for ten localization classes by three to five percentage points (Q10). The novel models are available as a web-service and as a stand-alone application at OMITTED.com.

1. Introduction

Proteins are the machinery of life involved in all essential biological processes (biological background in Appendix). Knowing where in the cell a protein functions, referred to as its *subcellular localization* or *cellular compartment*, is important for unraveling biological function (Nair & Rost, 2005; Yu et al., 2006). Experimental determination of protein function is complex, costly, and selection biased (Ching et al., 2018). In contrast, the costs of determining protein sequences continuously decrease (Consortium, 2021), increasing the sequence-annotation gap (gap between proteins of known sequence and unknown function). Computational methods have been bridging this gap (Rost et al., 2003); one way has been to predict protein subcellular location (Goldberg et al., 2012; 2014; Almagro Armenteros et al., 2017; Savojardo et al., 2018). The standard tool in molecular

biology, namely homology-based inference (HBI), accurately transfers annotations from experimentally annotated to sequence-similar un-annotated proteins. However, HBI is not available or unreliable for most proteins (Goldberg et al., 2014; Mahlich et al., 2018). Machine learning methods perform less well (lower precision) but are available for all proteins (high recall). The best methods use evolutionary information from families of related proteins as input (Goldberg et al., 2012; Almagro Armenteros et al., 2017). Although the marriage of evolutionary information and machine learning has influenced computational biology for decades (Rost & Sander, 1993), due to database growth, this information becomes increasingly costly to generate.

Recently, protein sequence representations (embeddings) have been learned from databases (Steinegger & Söding, 2018; Consortium, 2021) using language models (LMs) (Heinzinger et al., 2019; Rives et al., 2019; Alley et al., 2019; Elnaggar et al., 2020) initially used in natural language processing (NLP) (Radford, 2018; Devlin et al., 2019; Radford et al., 2019). Models trained on protein embeddings via transfer learning tend to be outperformed by approaches using evolutionary information (Rao et al., 2019; Heinzinger et al., 2019). However, embedding-based solutions can even outshine HBI (Littmann et al., 2021) and models predicting aspects of protein structure (Bhattacharya et al., 2020; Rao et al., 2020). Yet, for location prediction, embedding-based models (Heinzinger et al., 2019; Elnaggar et al., 2020; Littmann et al., 2021) remained inferior to the state-of-the-art using evolutionary information, e.g., represented by DeepLoc (Almagro Armenteros et al., 2017).

In this work, we leveraged protein embeddings to predict cellular location without evolutionary information. We proposed a deep neural network architecture using light attention (LA) inspired by previous attention mechanisms (Bahdanau et al., 2015; Vaswani et al., 2017).

2. Related Work

Previous state-of-the-art (SOTA) models for subcellular location prediction combined homology, evolutionary information, and machine learning, often building prior knowledge about biology into model architectures. For instance, LocTree2 (Goldberg et al., 2012) implemented profile-kernel SVMs (Cortes & Vapnik, 1995; Rui Kuang et al.,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

055 which identified k-mers conserved in evolution and
 056 put them into a hierarchy of models inspired by cellular
 057 sorting pathways. BUSCA (Savojardo et al., 2018) com-
 058 bines three compartment-specific prediction methods based
 059 on SVMs using evolutionary information (Pierleoni et al.,
 060 2006; 2011; Savojardo et al., 2017). DeepLoc (Alma-
 061 gro Armenteros et al., 2017) uses convolutions followed
 062 by a bidirectional LSTM (Hochreiter & Schmidhuber, 1997;
 063 Schuster & Paliwal, 1997) that employs Bahdanau-Attention
 064 (Bahdanau et al., 2015). Using evolutionary information,
 065 DeepLoc rose to become the SOTA. Embedding-based meth-
 066 ods (Heinzinger et al., 2019) have not yet outperformed this
 067 SOTA, although ProtTrans (Elnaggar et al., 2020), based on
 068 very large data sets, came close.

069

070 3. Methods

071 3.1. Data

072 **Standard set DeepLoc.** Following previous work
 073 (Heinzinger et al., 2019; Elnaggar et al., 2020), we mainly
 074 used a data set introduced by *DeepLoc* (Almagro Ar-
 075 menteros et al., 2017) for training and testing. The training
 076 set contained 13 858 proteins annotated with experimen-
 077 tal evidence for one of ten location classes (nucleus, cyto-
 078 plasm, extracellular space, mitochondrion, cell membrane,
 079 Endoplasmatic Reticulum, plastid, Golgi apparatus, lyso-
 080 some/vacuole, peroxisome). Another 2 768 proteins made
 081 up the test set (henceforth called *setDeepLoc*), which had
 082 been redundancy reduced to the training set (but not to
 083 itself) at 30% pairwise sequence identity (PIDE) or to an
 084 E-value cutoff of 10^{-6} . To tune model parameters and avoid
 085 overestimating performance, we further split the DeepLoc
 086 training set into a training set containing 9 503 sequences
 087 and a validation set (redundancy reduced to training by 30%
 088 PIDE) containing 1 158 sequences (distribution of classes
 089 in Appendix: Datasets).

090 **Novel setHARD.** To rule out that methods had been op-
 091 timized for the static standard test set (*setDeepLoc*) used
 092 by many developers, we created a new independent test set
 093 from SwissProt (Consortium, 2021). Applying the same
 094 filtering mechanisms as the DeepLoc developers (only eu-
 095 karyotes; only proteins longer than 40 residues; no frag-
 096 ments; only experimental location annotations) gave 5 947
 097 proteins. Using MMseqs2 (Steinegger & Söding, 2017), we
 098 removed all proteins from the new set with more than 20%
 099 PIDE to any protein in DeepLoc (both training and testing
 100 data). Next, we mapped location classes from DeepLoc to
 101 SwissProt, merged duplicates, and removed multi-localized
 102 proteins (protein X both in class Y and Z). Finally, we clus-
 103 tered proteins to representatives at 20% PIDE and obtained
 104 a new and more challenging test set (dubbed *setHARD*) with
 105 490 proteins. Class distributions differed between the two
 106 sets (see Appendix: Datasets).

Table 1. Parameters and implementation details of SeqVec (Heinzinger et al., 2019), ProtBert and ProtT5 (Elnaggar et al., 2020). The time it takes to embed a single sequence (sec per se-
 quence) is averaged over embedding 10 000 proteins taken from the Protein Data Bank (PDB) (Berman et al., 2000). The num-
 ber of sequences used for the pre-training task is detailed in "# sequences".

	SEQVEC	PROTBERT	PROTT5
PARAMETERS	93M	420M	3B
# SEQUENCES	33M	2.1B	2.1B
SEC PER SEQUENCE	0.03	0.06	0.1
ATTENTION HEADS	-	16	32

3.2. Models

Input: protein embeddings. As input to the LA archi-
 tectures, we extracted embeddings from three pre-trained
 protein language models (LMs; Table 1): the bidirectional
 LSTM SeqVec (Heinzinger et al., 2019) based on ELMo (Pe-
 ters et al., 2018), the encoder-only model ProtBert (Elnaggar
 et al., 2020) based on BERT (Devlin et al., 2019), and the
 encoder-only model ProtT5 (Elnaggar et al., 2020) based on
 T5 (Raffel et al., 2020). We obtained embeddings for each
 residue (NLP equivalent: word) in a protein sequence (NLP
 equivalent: document) using the bio-embeddings software
 (Dallago et al., 2020). For SeqVec, the per-residue embed-
 dings were generated by summing the representations of
 each layer. For ProtBert and ProtT5, the per-residue em-
 beddings were extracted from the last hidden layer of the
 models. With a hidden size of 1024 for each LM, inputs to LA
 were of size $1024 \times L$, where L is the length of the
 protein sequence.

Light Attention (LA) architecture. The input to the light
 attention (LA) classifiers (Figure 1) was a protein embed-
 ding $x \in \mathbb{R}^{1024 \times L}$. The input was transformed by two
 separate 1D convolutions with filter sizes s parameterized
 by learned weights $W^{(e)}, W^{(v)} \in \mathbb{R}^{s \times 1024 \times d_{out}}$. The con-
 volutions were applied over the length dimension to produce
 attention coefficients and value features $e, v \in \mathbb{R}^{d_{out} \times L}$.

$$e_{i,j} = b_i + \sum_{k=1}^{1024} \sum_{l=-\lfloor \frac{s}{2} \rfloor}^{\lceil \frac{s}{2} \rceil} W_{l,i}^{(e)} x_{:,j+l} \quad (1)$$

where $b \in \mathbb{R}^{d_{out}}$ is a learned bias and $x_{:,j}$ denotes the j-th
 residue embedding. Null vectors were set to $x_{:,j}$ for j out-
 side the interval $[0, L]$. To use the coefficients as attention
 distribution over all j , we softmax-normalized over protein
 length. The attention weight $\alpha_{i,j} \in \mathbb{R}$ for the j-th residue
 and the i-th feature dimension was calculated as:

$$\alpha_{i,j} = \frac{\exp(e_{i,j})}{\sum_{l=1}^L \exp(e_{i,l})} \quad (2)$$

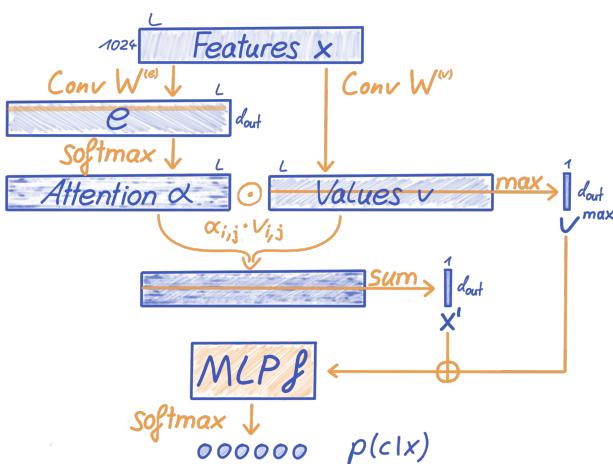


Figure 1. Sketch of LA solution. The LA architecture is parameterized by two weight matrices $W^{(e)}, W^{(v)} \in \mathbb{R}^{s \times 1024 \times d_{out}}$ and the weights of an MLP $f : \mathbb{R}^{2d_{out}} \mapsto \mathbb{R}^{d_{class}}$.

As the weight distributions for each feature dimension i are independent, they might generate different attention patterns. We used the normalized attention distributions to compute weighted sums over the transformed residue embeddings $v_{i,j}$. Thus, we obtained a fixed-size representation $x' \in \mathbb{R}^{d_{out}}$ for the whole protein, independent of its length.

$$x'_i = \sum_{j=1}^L \alpha_{i,j} v_{i,j} \quad (3)$$

We concatenated x'_i with the maximum of the values over the length dimension $v_i^{\max} \in \mathbb{R}^{d_{out}}$, meaning $v_i^{\max} = \max_{1 \leq j \leq L} (v_{i,j})$. This concatenated vector was input into a two layer multi-layer perceptron (MLP) $f : \mathbb{R}^{2d_{out}} \mapsto \mathbb{R}^{d_{class}}$ with d_{class} as the number of classes. The softmax over the MLP output represents the individual class probabilities with classes indexed by c :

$$p(c|x) = \text{softmax}_c(f(x' \oplus m)) \quad (4)$$

where \oplus denotes concatenation.

Training models. For the LA architecture, we trained three models, one for each protein embedding (SeqVec, ProtBert and ProtT5) for subsets of the training set. The models were trained using filter size $s = 9$, $d_{out} = 1024$, the Adam (Kingma & Ba, 2015) optimizer (learning rate 5×10^{-5}) with a batch size of 150 protein embeddings, and early stopping after no improvement in validation loss for 80 epochs. We selected the hyperparameters via random search (Appendix: Hyperparameters). Training was done on either an Nvidia Quadro RTX 8000 with 48GB vRAM or an Nvidia GeForce GTX 1060 with 6GB vRAM.

Methods used for comparison. For comparison, we trained a two layer feed-forward network (FFN) proposed previously (Heinzinger et al., 2019). Instead of per-residue embeddings in $\mathbb{R}^{1024 \times L}$, the FFNs used sequence-embeddings in \mathbb{R}^{1024} , which derived from residue-embeddings averaged over the length dimension (i.e. mean pooling). Furthermore, for these representations, we performed annotation transfer (dubbed AT) based on embedding similarity (Littmann et al., 2021). Following this approach, proteins in *setDeepLoc* and *setHARD* were annotated by transferring the class of the nearest neighbor in the DeepLoc training set (given by L1 distance).

3.3. Evaluation.

Following previous work, we assessed performance mostly through the mean ten-class accuracy (Q10), giving the percentage of correctly predicted proteins in one of ten location classes. Additional measures tested did not provide any additional insights and were, therefore, confined to the Appendix: Additional Results, e.g., Matthew's correlation coefficient (MMC) for multiple classes (Gorodkin measure (Gorodkin, 2004)). Error estimates were calculated over ten random seeds on both test sets. For previous methods (DeepLoc and DeepLoc62 (Almagro Armenteros et al., 2017), LocTree2 (Goldberg et al., 2012), MultiLoc2 (Blum et al., 2009), SherLoc2 (Briesemeister et al., 2009), CELLO (Yu et al., 2006), iLoc-Euk (Chou et al., 2011), YLoc (Briesemeister et al., 2010) and WoLF PSORT (Horton et al., 2007)) published performance values were used (Almagro Armenteros et al., 2017) for *setDeepLoc* and computed for *setHARD*. As a naive baseline, we used the majority classifier. All evaluation scripts to reproduce results are available¹.

4. Results

Embeddings outperformed evolutionary information. The simple AT approach already outperformed some methods using evolutionary information (Figure 2: AT*). The FFN trained on ProtT5 (Elnaggar et al., 2020) outperformed DeepLoc (Almagro Armenteros et al., 2017) (Figure 2: FNN ProtT5 vs. DeepLoc). Methods based on ProtT5 embeddings consistently yielded better results than ProtBert and SeqVec (Heinzinger et al., 2019) (*ProtT5 vs *ProtBert/*SeqVec in Figure 2). Results on Q10 are consistent with MCC (Appendix: Additional Results).

LA architecture best. The light attention (LA) architecture consistently outperformed other embedding-based approaches, irrespective of the protein LM used (LA* vs. AT/FFN* in Figure 2). Using ProtBert embeddings, LA outperformed the state-of-the-art (SOTA) (Almagro Armenteros et al., 2017) by 1 – 2 percentage points (LA Prot-

¹ OMITTED

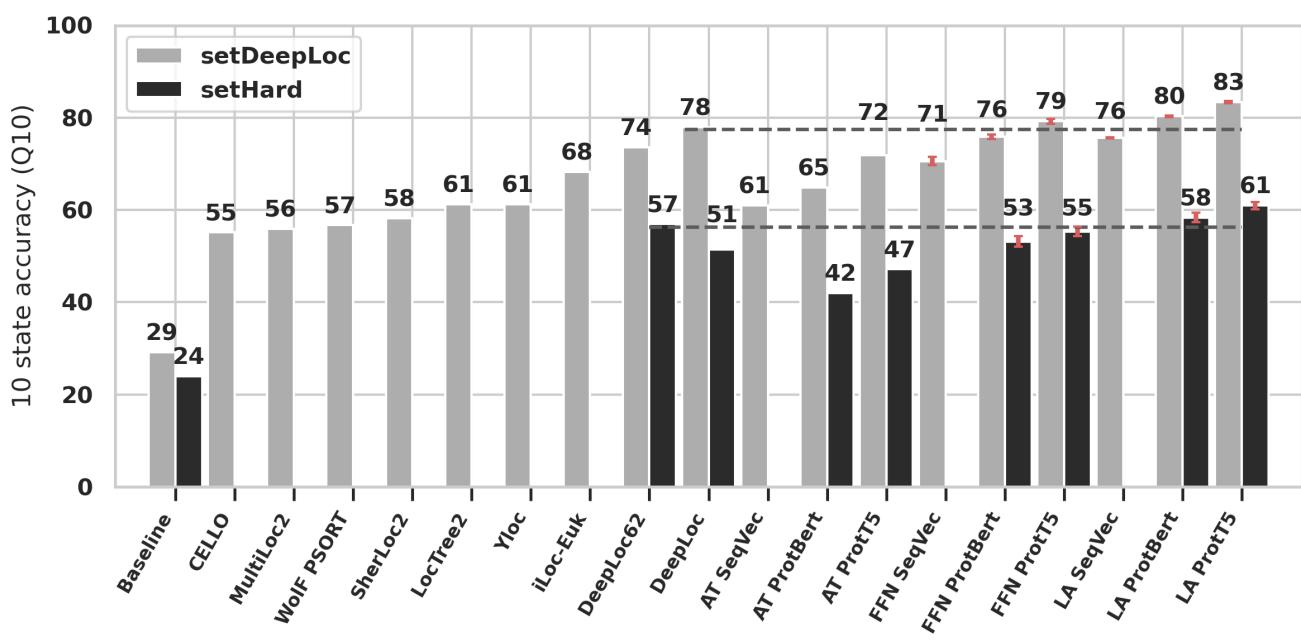


Figure 2. **LA architectures perform best.** Bars give the ten-class accuracy (Q10) for popular location prediction methods on *setDeepLoc* (light-gray bars) and *setHARD* (dark-gray bars). Baseline is the most common class in each set. Horizontal gray dashed lines mark the previous SOTA on either set. Estimates for standard errors are marked in orange for the methods introduced here. *setHARD* results are provided for a subset of methods that yielded the best results on *setDeepLoc* (Methods for detail on the external methods used; tabular data in Appendix: Additional Results). Two results stood out: (i) the LA approaches introduced here outperformed the top methods although not using evolutionary information (highest bars), and (ii) the performance estimates differed completely between the two data sets (difference light/dark gray).

Bert Figure 2) and by 3 – 5 percentage points using ProtT5 (Figure 2) for *setHARD* and *setDeepLoc*, respectively.

Overfitting by using standard data set. The substantial drop in performance (around 22 percentage points) between results for the standard *setDeepLoc* and the new challenging *setHARD* (Figure 2: light-gray vs. dark-gray, respectively) suggests some level of overfitting. Mimicking the distribution of classes found in *setDeepLoc* by sampling with replacement from *setHARD* led to better results (in Q10: *DeepLoc62*=63%; *DeepLoc*=54%; *LA ProtBert*=62%; *LA ProtT5*=66%). *DeepLoc* performed worse on *setHARD* using evolutionary information than using simple sequence information (Figure 2: *DeepLoc* vs. *DeepLoc62*). Otherwise, the relative ranking and difference of models largely remained consistent between *setDeepLoc* and *setHARD*.

Low performance for minority classes. The confusion matrix of predictions for *setDeepLoc* using LA trained on ProtT5 embeddings highlighted how many proteins were incorrectly predicted in the most prevalent class, *cytoplasm*, and that even the two majority classes were often confused with each other (Figure 3: *nucleus* and *cytoplasm*). In line with the previous SOTA (Almagro Armenteros et al., 2017), the performance was particularly low for the most

under-represented classes, namely *Golgi apparatus*, *lysosome/Vacuole*, and *peroxisome* (accounting for 2.6%, 2.3%, and 1.1% of the data, respectively).

Light aggregation (LA) mechanism crucial. To further evaluate the effectiveness of the LA architecture’s aggregation mechanisms, we replaced the light attention that produced x' with averaging the coefficient features e over the length dimension. Performance using ProtT5 embeddings dropped from 83.37% to 81.54 \pm 0.13% (Q10(*setDeepLoc*)). Similarly, we dropped the max-pooled values v^{max} as input to the MLP such that only the aggregated light attention features were used. This reduced performance from 83.37% to 82.23 \pm 0.44% (Q10(*setDeepLoc*)).

Model trainable on consumer hardware. After embeddings for proteins were generated, the final LA architecture, made of 18 940 224 parameters, could be trained on an Nvidia GeForce GTX 1060 with 6GB vRAM in 18 hours or on a Quadro RTX 8000 with 48GB vRAM in 2.5 hours.

5. Discussion

Light attention beats pooling. The central challenge for the improvement introduced here was to convert the

220	Mem	0.86	0.08	0.02		0.01				
221	Cyt		0.82							
222	End	0.17	0.09	0.65	0.01	0.02				
223	Gol	0.20	0.24	0.07	0.39	0.01				
224	Lys	0.31	0.11	0.11	0.03	0.22				
225	Mit				0.84		0.03			
226	Nuc					0.90				
227	Per	0.03	0.30	0.13		0.03	0.07	0.03	0.40	
228	Pla					0.03	0.01		0.94	
229	Ext	0.02	0.02	0.01						0.95
230	Mem									
231	Cyt									
232	End									
233	Gol									
234	Lys									
235	Mit									
236	Nuc									
237	Per									
238	Pla									
239	Ext									

Figure 3. **Mostly capturing majority classes.** Confusion matrix of LA predictions on ProtT5 (Elnaggar et al., 2020) embeddings for *setDeepLoc* (Almagro Armenteros et al., 2017) annotated with the fraction of the true class. Y-axis (vertical): true class, X-axis (horizontal): predicted class. Labels: Mem=cell Membrane; Cyt=Cytoplasm; End=Endoplasmic Reticulum; Gol=Golgi apparatus; Lys=Lysosome/vacuole; Mit=Mitochondrion; Nuc=Nucleus; Per=Peroxisome; Pla=Plastid; Ext=Extracellular

residue-embeddings (NLP equivalent: word embeddings) from protein language models such as SeqVec (Heinzinger et al., 2019), ProtBert, or ProtT5 (Elnaggar et al., 2020) to meaningful sequence-embeddings (NLP equivalent: document). A qualitative evaluation of the influence of the attention mechanism (Figure 4) highlighted its ability to steer predictions. Although averaging surpassed evolutionary-information-based methods using simple similarity-based annotation transfer (Figure 2: AT*) and in one instance even SOTA using a simple feed-forward network (Figure 2: DeepLoc vs. FNN ProtT5), LA was able to consistently distill more information from embeddings. Most likely, the improvement can be attributed to LA's ability to regulate the immense difference in lengths of proteins (varying from 30 to 30 000 residues (Consortium, 2021)) by learning attention distributions over the sequence positions. LA models appeared to have captured relevant long-range dependencies while retaining the ability to focus on specific sequence regions such as beginning and end, which play a particularly important role in determining protein location for some proteins (Lange et al., 2007; Almagro Armenteros et al., 2017).

First win over evolutionary information. Effectively, LA trained on protein LM embeddings from ProtT5 (Elnaggar et al., 2020) was at the heart of the first method that clearly appeared to outperform the best existing method (*DeepLoc*,

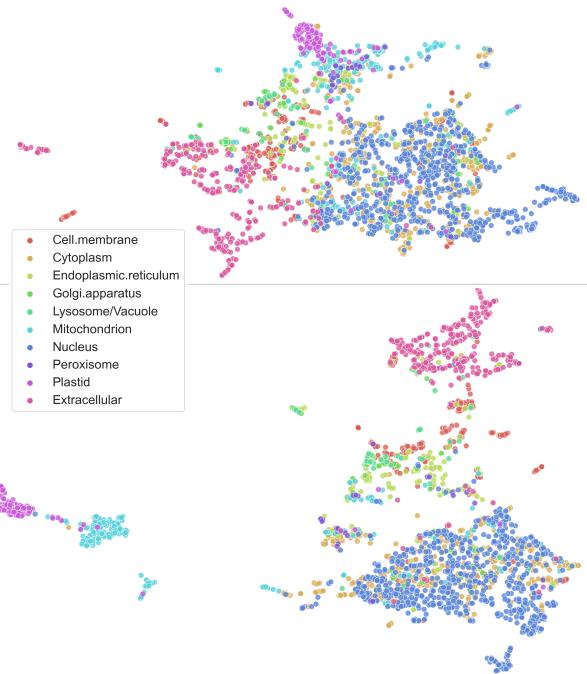


Figure 4. **Qualitative analysis confirms: attention effective.** UMAP (McInnes et al., 2018) projections of per-protein embeddings colored according to subcellular location (*setDeepLoc*). Top: ProtT5 embeddings (LA input; x) mean-pooled over protein length (as for FFN/AT input). Bottom: ProtT5 embeddings (LA input; x) weighted according to the attention distribution produced by LA (this is not x' as we sum the input features x and not the values v after the convolution).

(Almagro Armenteros et al., 2017; Heinzinger et al., 2019)) in a statistically significant manner on two test sets (Figure 2). To the best of our knowledge, this improvement was the first instance ever that embedding-based transfer learning substantially outperformed AI/ML methods using evolutionary information for function prediction. Even if embeddings are extracted from LMs trained on large sequence data originating from evolution, the vast majority of data learned originates from much more generic constraints informative of protein structure and function.

Better and faster. The embeddings needed as input for the LA models come with three advantages over evolutionary-information-based input essential for methods such as *DeepLoc* (Almagro Armenteros et al., 2017). Chiefly, embeddings can be obtained in far less time than is needed to generate evolutionary information and require fewer compute resources. Even considering the lightning-fast MMseqs2 (Steinegger & Söding, 2017), which is not the standard in bioinformatics (other methods 10-100x slower), in our experience required about 0.3 seconds to generate evolutionary information input for a large set of 10 000 pro-

teins. The slowest but most informative embedder (ProtT5) is 3x faster, while the second most informative (ProtBert) is 5x faster (Table 1). Additionally, these MMseqs2 stats derive from runs on a machine with > 300GB of RAM and 2x40cores/80threads CPUs, while generating LM embeddings required only a moderate machine (8 cores, 16GB RAM) equipped with a modern GPU with >10GB of vRAM. Lastly, extracting evolutionary information relies on the use of tools (e.g., MMseqs2) that are sensitive to parameter changes, ultimately an extra complication for users. In contrast, generating embeddings doesn't require a parameter choice beyond which trained model to use (e.g., ProtBert vs. ProtT5).

Overfitting through standard data set? For protein subcellular location prediction, the data sets from *DeepLoc* (Almagro Armenteros et al., 2017) have become a standard in the field. Such static standards facilitate method comparisons. To further probe results, we created a new test set (*setHARD*), which was redundancy-reduced both with respect to itself and all proteins in the *DeepLoc* set (comprised of training data and *setDeepLoc*, used for testing). For this set, the 10-state accuracy (Q10) dropped, on average, 22 percentage points with respect to the static standard (Figure 2). We argue that this large margin may be attributed to some combination of the following coupled effects.

(1) All new methods may simply have been substantially overfitted to the static data set, e.g., by misusing the test set for hyperparameter optimization. This could partially explain the increase in performance on *setHARD* when mimicking the class distributions in the training set and *setDeepLoc*.

(2) The static standard set allowed for some level of sequence-redundancy (information leakage) at various levels: certainly within the test set, which had not been redundancy reduced to itself, maybe also between the training and test set. Methods with many free parameters might more easily zoom into exploiting such residual sequence similarity for prediction because proteins with similar sequence locate in similar compartments. In fact, this may explain the somewhat surprising observation that *DeepLoc* appeared to perform worse on *setHARD* using evolutionary information instead of a generic BLOSUM metric (Figure 2: *DeepLoc62* vs. *DeepLoc*). Residual redundancy is much easier to capture via evolutionary information than by BLOSUM (Urban et al., 2020) (for computational biologists: the same way in which PSI-BLAST (Altschul et al., 1997) outperforms pairwise BLAST).

(3) The confusion matrix (Figure 3) demonstrated how classes with more experimental data tended to be predicted more accurately. As *setDeepLoc* and *setHARD* differed in their class composition, even without overfitting and redundancy, prediction methods would perform differently on

the two. In fact, this can be investigated by recomputing the performance on a similar class-distributed superset of *setHARD*, on which performance dropped only by 11, 24, 18, and 17 percentage points for *DeepLoc62*, *DeepLoc*, *LA ProtT5*, and *LA ProtBert*, respectively.

Overall, several overlaying effects caused the performance to drop between the two data sets. Interestingly, different approaches behaved alike: both for alternative inputs from protein language models (ProtVec, ProtBERT, ProtT5) and for alternative methods (AT, FFN, LA), of which one (AT) refrained from weight optimization.

What can users expect from subcellular location predictions? If the top accuracy for one data set was Q10 ~ 60% and Q10 ~ 80% for the other, what can users expect for their next ten queries: six correct or eight, or 6-8? The answer depends on the query: if those proteins were sequence similar to proteins with known location (case: redundant): the answer is eight. Conversely, for new proteins (without homologs of known location), six in ten will be correctly predicted, on average. In turn, this implies that for novel proteins, there seems to be significant room for pushing performance to further heights, possibly by combining *LA ProtBert/LA ProtT5* with evolutionary information.

6. Conclusion

We presented a light attention mechanism (LA) in an architecture operating on language model embeddings of protein sequences, namely those from SeqVec (Heinzinger et al., 2019), ProtBert, and ProtT5 (Elnaggar et al., 2020). LA efficiently aggregated information and coped with arbitrary sequence lengths, thereby mastering the enormous range of proteins spanning from 30-30 000 residues. By implicitly assigning a different importance score for each sequence position, the method succeeded in predicting protein subcellular location much better than methods based on simple pooling. More importantly, for two protein LMs, LA succeeded in outperforming the state-of-the-art without using evolutionary-based inputs, i.e., the single most important input feature for previous methods. This constituted an important breakthrough: although many methods had come close to the state-of-the-art using embeddings instead of evolutionary information, none had ever overtaken as the methods presented here. Our best method was based on the largest protein LM, namely on ProtT5 (*LA ProtT5* in Figure 2). Many location prediction methods have been assessed on a standard data set (here: *setDeepLoc*) introduced a few years ago (Almagro Armenteros et al., 2017). Using a new, more challenging data set (*setHARD*), the performance of all methods appeared to drop by around 22 percentage points. While class distributions and data set redundancy (or homology) may explain some of this drop, over-fitting might have also contributed. Overall, the drop underlined

330 that many challenges remain to be addressed by future methods.
 331 For the time being, the best methods *LA ProtBert* and
 332 *LA ProtT5*, are freely available via a web-server² and as part
 333 of a high-throughput pipeline (Dallago et al., 2020).

Acknowledgements

References

339 Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi,
 340 M., and Church, G. M. Unified rational protein
 341 engineering with sequence-based deep representation
 342 learning. *Nature Methods*, 16(12):1315–1322, De-
 343 cember 2019. ISSN 1548-7105. doi: 10.1038/
 344 s41592-019-0598-1. URL <https://www.nature.com/articles/s41592-019-0598-1>. Number:
 345 12 Publisher: Nature Publishing Group.

347 Almagro Armenteros, J. J., Sønderby, C. K., Sønderby, S. K.,
 348 Nielsen, H., and Winther, O. DeepLoc: prediction of
 349 protein subcellular localization using deep learning. *Bioinfor-*
 350 *matics*, 33(21):3387–3395, November 2017. ISSN 1367-
 351 4803. doi: 10.1093/bioinformatics/btx431. URL <https://academic.oup.com/bioinformatics/article/33/21/3387/3931857>. tex.ids: al-
 352 magroarmenterosDeepLocPredictionProtein2017a
 353 publisher: Oxford Academic.

354 Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang,
 355 J., Zhang, Z., Miller, W., and Lipman, D. J. Gapped
 356 BLAST and PSI-BLAST: a new generation of protein
 357 database search programs. *Nucleic Acids Research*, 25
 358 (17):3389–3402, September 1997. ISSN 0305-1048. doi:
 359 10.1093/nar/25.17.3389. URL <https://doi.org/10.1093/nar/25.17.3389>.

360 Bahdanau, D., Cho, K., and Bengio, Y. Neural Machine
 361 Translation by Jointly Learning to Align and Translate.
 362 In Bengio, Y. and LeCun, Y. (eds.), *3rd International
 363 Conference on Learning Representations, ICLR 2015,
 364 San Diego, CA, USA, May 7-9, 2015, Conference Track
 365 Proceedings*, 2015. URL <http://arxiv.org/abs/1409.0473>.

366 Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G., Bhat,
 367 T. N., Weissig, H., Shindyalov, I. N., and Bourne, P. E.
 368 The Protein Data Bank. *Nucleic Acids Research*, 28(1):
 369 235–242, January 2000. ISSN 0305-1048. doi: 10.1093/
 370 nar/28.1.235. URL <https://doi.org/10.1093/nar/28.1.235>.

371 Bhattacharya, N., Thomas, N., Rao, R., Dauparas, J.,
 372 Koo, P. K., Baker, D., Song, Y. S., and Ovchinnikov,
 373 S. Single Layers of Attention Suffice to Predict
 374 Protein Contacts. *bioRxiv*, pp. 2020.12.21.423882,

375 December 2020. doi: 10.1101/2020.12.21.423882.
 376 URL <https://www.biorxiv.org/content/10.1101/2020.12.21.423882v2>. Publisher:
 377 Cold Spring Harbor Laboratory Section: New Results.

378 Blum, T., Briesemeister, S., and Kohlbacher, O. MultiLoc2:
 379 integrating phylogeny and Gene Ontology terms improves
 380 subcellular protein localization prediction. *BMC bioin-*
 381 *formatics*, 10(1):274, 2009. Publisher: Springer.

382 Briesemeister, S., Blum, T., Brady, S., Lam, Y., Kohlbacher,
 383 O., and Shatkay, H. SherLoc2: a high-accuracy hybrid
 384 method for predicting subcellular localization of proteins.
 385 *Journal of proteome research*, 8(11):5363–5366, 2009.
 386 Publisher: ACS Publications.

387 Briesemeister, S., Rahnenführer, J., and Kohlbacher, O.
 388 YLoc—an interpretable web server for predicting sub-
 389 cellular localization. *Nucleic acids research*, 38(suppl_2):
 390 W497–W502, 2010. Publisher: Oxford University Press.

391 Ching, T., Himmelstein, D. S., Beaulieu-Jones, B. K.,
 392 Kalinin, A. A., Do, B. T., Way, G. P., Ferrero, E.,
 393 Agapow, P.-M., Zietz, M., Hoffman, M. M., Xie, W.,
 394 Rosen, G. L., Lengerich, B. J., Israeli, J., Lanchantin,
 395 J., Woloszynek, S., Carpenter, A. E., Shrikumar, A.,
 396 Xu, J., Cofer, E. M., Lavender, C. A., Turaga, S. C.,
 397 Alexandari, A. M., Lu, Z., Harris, D. J., DeCaprio, D.,
 398 Qi, Y., Kundaje, A., Peng, Y., Wiley, L. K., Segler,
 399 M. H. S., Boca, S. M., Swamidass, S. J., Huang,
 400 A., Gitter, A., and Greene, C. S. Opportunities and
 401 obstacles for deep learning in biology and medicine.
 402 *Journal of The Royal Society Interface*, 15(141):
 403 20170387, April 2018. doi: 10.1098/rsif.2017.0387.
 404 URL <https://royalsocietypublishing.org/doi/10.1098/rsif.2017.0387>. Publisher:
 405 Royal Society.

406 Chou, K.-C., Wu, Z.-C., and Xiao, X. iLoc-Euk: a multi-
 407 label classifier for predicting the subcellular localization
 408 of singleplex and multiplex eukaryotic proteins. *PloS one*,
 409 6(3):e18258, 2011. Publisher: Public Library of Science.

410 Consortium, T. U. UniProt: the universal protein knowl-
 411 edgebase in 2021. *Nucleic Acids Research*, 2021. doi:
 412 10.1093/nar/gkaa1100. URL <https://academic.oup.com/nar/advance-article/doi/10.1093/nar/gkaa1100/6006196>.

413 Cortes, C. and Vapnik, V. Support-vector networks. *Ma-*
 414 *chine Learning*, 20(3):273–297, September 1995. ISSN
 415 0885-6125, 1573-0565. doi: 10.1007/BF00994018.
 416 URL <http://link.springer.com/10.1007/BF00994018>.

417 Dallago, C., Schütze, K., Heinzinger, M., Olenyi, T.,
 418 and Rost, B. bio_embeddings: python pipeline for

²OMITTED

- 385 fast visualization of protein features extracted by lan-
 386 guage models. *F1000Research*, 9, August 2020. doi:
 387 10.7490/f1000research.11181631. URL <https://f1000research.com/posters/9-876>.
- 388
- 389 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K.
 390 BERT: Pre-training of Deep Bidirectional Transfor-
 391 mers for Language Understanding. In Burstein, J., Do-
 392 ran, C., and Solorio, T. (eds.), *Proceedings of the 2019*
 393 *Conference of the North American Chapter of the As-*
 394 *sociation for Computational Linguistics: Human Lan-*
 395 *guage Technologies, NAACL-HLT 2019, Minneapolis,*
 396 *MN, USA, June 2-7, 2019, Volume 1 (Long and Short*
 397 *Papers)*, pp. 4171–4186. Association for Computational
 398 Linguistics, 2019. doi: 10.18653/v1/n19-1423. URL
 399 <https://doi.org/10.18653/v1/n19-1423>.
- 400
- 401 Elnaggar, A., Heinzinger, M., Dallago, C., Rehawi, G.,
 402 Wang, Y., Jones, L., Gibbs, T., Feher, T., Angerer, C.,
 403 Steinegger, M., Bhowmik, D., and Rost, B. ProtTrans:
 404 Towards Cracking the Language of Life’s Code Through
 405 Self-Supervised Deep Learning and High Perfor-
 406 mance Computing. *bioRxiv*, pp. 2020.07.12.199554,
 407 July 2020. doi: 10.1101/2020.07.12.199554. URL
 408 <https://www.biorxiv.org/content/10.1101/2020.07.12.199554v2>. tex.ids: elnaggarProtTransCrackingLanguage2020a publisher: Cold
 409 Spring Harbor Laboratory section: New Results.
- 410
- 411
- 412 Goldberg, T., Hamp, T., and Rost, B. LocTree2 predicts
 413 localization for all domains of life. *Bioinformatics*, 28
 414 (18):i458–i465, September 2012.
- 415
- 416 Goldberg, T., Hecht, M., Hamp, T., Karl, T., Yachdav, G.,
 417 Ahmed, N., Altermann, U., Angerer, P., Ansorge, S.,
 418 Balasz, K., Bernhofer, M., Betz, A., Cizmadija, L., Do,
 419 K. T., Gerke, J., Greil, R., Joerdens, V., Hastreiter, M.,
 420 Hembach, K., Herzog, M., Kalemanov, M., Kluge, M.,
 421 Meier, A., Nasir, H., Neumaier, U., Prade, V., Reeb, J.,
 422 Sorokoumov, A., Troshani, I., Vorberg, S., Waldraff, S.,
 423 Zierer, J., Nielsen, H., and Rost, B. LocTree3 predic-
 424 tion of localization. *Nucleic Acids Research*, 42(W1):
 425 W350–W355, 2014. ISSN 0305-1048. doi: 10.1093/nar/
 426 gku396. URL <https://doi.org/10.1093/nar/gku396>. eprint: <https://academic.oup.com/nar/article-pdf/42/W1/W350/17423232/gku396.pdf>.
- 427
- 428
- 429 Gorodkin, J. Comparing two K-category assign-
 430 ments by a K-category correlation coefficient.
Computational Biology and Chemistry, 28(5):
 431 367 – 374, 2004. ISSN 1476-9271. doi:
 432 <https://doi.org/10.1016/j.combiolchem.2004.09.006>.
 433 URL <http://www.sciencedirect.com/science/article/pii/S1476927104000799>.
- 434
- 435 Heinzinger, M., Elnaggar, A., Wang, Y., Dallago, C.,
 436 Nechaev, D., Matthes, F., and Rost, B. Modeling as-
 437 pects of the language of life through transfer-learning
 438 protein sequences. *BMC Bioinformatics*, 20(1):723,
 439 December 2019. ISSN 1471-2105. doi: 10.1186/s12859-019-3220-8. URL <https://doi.org/10.1186/s12859-019-3220-8>. tex.ids: heinzinger-
 440 ModelingAspectsLanguage2019a.
- 441
- Hochreiter, S. and Schmidhuber, J. Long Short-Term
 442 Memory. *Neural Computation*, 9(8):1735–1780, 1997.
 443 doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>. eprint:
<https://doi.org/10.1162/neco.1997.9.8.1735>.
- Horton, P., Park, K.-J., Obayashi, T., Fujita, N., Harada,
 444 H., Adams-Collier, C., and Nakai, K. WoLF PSORT:
 445 protein localization predictor. *Nucleic acids research*, 35
 446 (suppl_2):W585–W587, 2007. Publisher: Oxford University Press.
- Kingma, D. P. and Ba, J. Adam: A Method for Stochastic Optimization. In Bengio, Y. and LeCun, Y. (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- Lange, A., Mills, R. E., Lange, C. J., Stewart, M.,
 447 Devine, S. E., and Corbett, A. H. Classical Nuclear
 448 Localization Signals: Definition, Function, and Inter-
 449 action with Importin alpha,. *Journal of Biological
 450 Chemistry*, 282(8):5101–5105, February 2007. ISSN
 451 0021-9258. doi: 10.1074/jbc.R600026200. URL
<http://www.sciencedirect.com/science/article/pii/S0021925820688019>.
- Littmann, M., Heinzinger, M., Dallago, C., Olenyi, T., and
 452 Rost, B. Embeddings from deep learning transfer GO
 453 annotations beyond homology. *Scientific Reports*, 11(1):
 454 1160, January 2021. ISSN 2045-2322. doi: 10.1038/s41598-020-80786-0. URL <https://www.nature.com/articles/s41598-020-80786-0>. Number:
 455 1 Publisher: Nature Publishing Group.
- Mahlich, Y., Steinegger, M., Rost, B., and Bromberg,
 456 Y. HFSP: high speed homology-driven function an-
 457 notation of proteins. *Bioinformatics*, 34(13):i304–
 458 i312, July 2018. ISSN 1367-4803. doi: 10.1093/
 459 bioinformatics/bty262. URL <https://doi.org/10.1093/bioinformatics/bty262>.
- McInnes, L., Healy, J., Saul, N., and Großberger, L.
 460 UMAP: Uniform Manifold Approximation and Projec-
 461 tion. *J. Open Source Softw.*, 3(29):861, 2018. doi:
 462 10.21105/joss.00861. URL <https://doi.org/10.21105/joss.00861>.

- 440 Nair, R. and Rost, B. Mimicking Cellular Sorting
 441 Improves Prediction of Subcellular Localization.
 442 *Journal of Molecular Biology*, 348(1):85–100, April
 443 2005. ISSN 0022-2836. doi: 10.1016/j.jmb.2005.02.
 444 025. URL <http://www.sciencedirect.com/science/article/pii/S0022283605001774>.
- 445
- 446 Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark,
 447 C., Lee, K., and Zettlemoyer, L. Deep Contextualized
 448 Word Representations. In *Proceedings of the 2018*
 449 *Conference of the North American Chapter of the Asso-*
 450 *ciation for Computational Linguistics: Human Lan-*
 451 *guage Technologies, Volume 1 (Long Papers)*, pp. 2227–
 452 2237, New Orleans, Louisiana, June 2018. Association
 453 for Computational Linguistics. doi: 10.18653/v1/N18-1202. URL <https://www.aclweb.org/anthology/N18-1202>.
- 454
- 455 Pierleoni, A., Martelli, P. L., Fariselli, P., and Casadio, R.
 456 BaCeLo: a balanced subcellular localization predictor.
 457 *Bioinformatics*, 22(14):e408–416, July 2006.
- 458
- 459 Pierleoni, A., Martelli, P. L., and Casadio, R. MemLoci:
 460 predicting subcellular localization of membrane proteins
 461 in eukaryotes. *Bioinformatics*, 27(9):1224–1230, May
 462 2011.
- 463
- 464 Radford, A. Improving Language Understanding by Gener-
 465 ative Pre-Training. 2018.
- 466
- 467 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and
 468 Sutskever, I. Language Models are Unsupervised Multi-
 469 task Learners. 2019.
- 470
- 471 Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S.,
 472 Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring
 473 the Limits of Transfer Learning with a Unified Text-
 474 to-Text Transformer. *J. Mach. Learn. Res.*, 21:140:1–
 475 140:67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- 476
- 477 Rao, R., Bhattacharya, N., Thomas, N., Duan, Y., Chen, X.,
 478 Canny, J., Abbeel, P., and Song, Y. S. Evaluating Pro-
 479 tein Transfer Learning with TAPE. *Advances in neural*
 480 *information processing systems*, 32:9689–9701, Decem-
 481 ber 2019. ISSN 1049-5258. URL <https://pubmed.ncbi.nlm.nih.gov/33390682>.
- 482
- 483 Rao, R., Ovchinnikov, S., Meier, J., Rives, A., and Sercu, T.
 484 Transformer protein language models are unsupervised
 485 structure learners. *bioRxiv*, pp. 2020.12.15.422761,
 486 December 2020. doi: 10.1101/2020.12.15.422761.
 487 URL <https://www.biorxiv.org/content/10.1101/2020.12.15.422761v1>. tex.ids: rao-
 488 TransformerProteinLanguage2020a publisher: Cold
 489 Spring Harbor Laboratory section: New Results.
- 490
- 491 Rives, A., Goyal, S., Meier, J., Guo, D., Ott, M., Zitnick,
 492 C. L., Ma, J., and Fergus, R. Biological Structure and
 493 Function Emerge from Scaling Unsupervised Learning to
 494 250 Million Protein Sequences. *bioRxiv*, 2019. doi: 10.
 495 1101/622803. URL <https://www.biorxiv.org/content/early/2019/04/29/622803>.
- 496
- 497 Rost, B. and Sander, C. Prediction of protein secondary
 498 structure at better than 70% accuracy. *Journal of molec-*
 499 *ular biology*, 232(2):584–599, 1993. tex.ids: RN1 pub-
 500 lisher: Elsevier Science type: Journal article.
- 501
- 502 Rost, B., Liu, J., Nair, R., Wrzeszczynski, K. O., and
 503 Ofran, Y. Automatic prediction of protein function. *Cel-*
 504 *lular and Molecular Life Sciences*, 60(12):2637–2650,
 505 2003. URL http://www.rostlab.org/papers/2003_rev_func/. Type: Journal article.
- 506
- 507 Rui Kuang, Ie, E., Ke Wang, Kai Wang, Siddiqi, M., Fre-
 508 und, Y., and Leslie, C. Profile-based string kernels for
 509 remote homology detection and motif extraction. In *Pro-*
 510 *ceedings. 2004 IEEE Computational Systems Bioinfor-*
 511 *matics Conference, 2004. CSB 2004.*, pp. 146–154, Stan-
 512 ford, CA, USA, 2004. IEEE. ISBN 978-0-7695-2194-
 513 7. doi: 10.1109/CSB.2004.1332428. URL <http://ieeexplore.ieee.org/document/1332428/>.
- 514
- 515 Savojardo, C., Martelli, P. L., Fariselli, P., and Casadio, R.
 516 SChloro: directing Viridiplantae proteins to six chloro-
 517 plastic sub-compartments. *Bioinformatics*, 33(3):347–
 518 353, 2017.
- 519
- 520 Savojardo, C., Martelli, P. L., Fariselli, P., Profi-
 521 ti, G., and Casadio, R. BUSCA: an integrative web server
 522 to predict subcellular localization of proteins. *Nucleic Acids Research*, 46(W1):W459–
 523 W466, 2018. ISSN 0305-1048. doi: 10.1093/nar/gky320. URL <https://doi.org/10.1093/nar/gky320>. _eprint: <https://academic.oup.com/nar/article-pdf/46/W1/W459/25110557/gky320.pdf>.
- 524
- 525 Schuster, M. and Paliwal, K. K. Bidirectional recurrent
 526 neural networks. *IEEE Trans. Signal Process.*, 45(11):
 527 2673–2681, 1997. doi: 10.1109/78.650093. URL
 528 <https://doi.org/10.1109/78.650093>.
- 529
- 530 Steinegger, M. and Söding, J. MMseqs2 enables sen-
 531 sitive protein sequence searching for the analysis of
 532 massive data sets. *Nature Biotechnology*, 35(11):1026–
 533 1028, November 2017. ISSN 1546-1696. doi: 10.1038/nbt.3988. URL <https://doi.org/10.1038/nbt.3988>.
- 534
- 535 Steinegger, M. and Söding, J. Clustering huge protein se-
 536 quence sets in linear time. *Nature Communications*, 9
 537 (1):2542, June 2018. ISSN 2041-1723. doi: 10.1038/s41467-018-04964-5. URL <https://doi.org/10.1038/s41467-018-04964-5>.

- 495 Urban, G., Torrisi, M., Magnan, C. N., Pollastri, G.,
496 and Baldi, P. Protein profiles: Biases and proto-
497 cols. *Computational and Structural Biotechnology*
498 *Journal*, 18:2281 – 2289, 2020. ISSN 2001-0370.
499 doi: <https://doi.org/10.1016/j.csbj.2020.08.015>. URL
500 <http://www.sciencedirect.com/science/article/pii/S2001037020303688>. tex.ids:
501 urbanProteinProfilesBiases2020a.
502
- 503 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
504 L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention
505 is All you Need. In Guyon, I., Luxburg, U. V., Bengio, S.,
506 Wallach, H., Fergus, R., Vishwanathan, S., and Garnett,
507 R. (eds.), *Advances in Neural Information Processing
508 Systems*, volume 30, pp. 5998–6008. Curran Asso-
509 ciates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fb053c1c4a845aa-Paper.pdf>.
510
- 511
- 512
- 513
- 514 Yu, C.-S., Chen, Y.-C., Lu, C.-H., and Hwang, J.-K. Predic-
515 tion of protein subcellular localization. *Proteins: Struc-
516 ture, Function, and Bioinformatics*, 64(3):643–651, 2006.
517 tex.ids: yuPredictionProteinSubcellular2006a publisher:
518 Wiley Online Library.
519
- 520
- 521
- 522
- 523
- 524
- 525
- 526
- 527
- 528
- 529
- 530
- 531
- 532
- 533
- 534
- 535
- 536
- 537
- 538
- 539
- 540
- 541
- 542
- 543
- 544
- 545
- 546
- 547
- 548
- 549

Appendix: Light Attention Predicts Protein Location from the Language of Life

Anonymous Authors¹

1. Protein Preliminaries

Protein Sequences. Proteins are built by chaining and arbitrary number of one of 20 amino acids in a particular order. When amino acids come together to form protein sequences, they are dubbed residues. During the assembly in the cell, constrained by physiochemical forces, the one-dimensional chains of residues fold into unique 3D shapes based solely on their sequence that largely determine protein function. The ideal machine learning model would predict a protein’s 3D shape and thus function from just protein sequence (the ordered chain of residues).

Protein Subcellular Location. Eukaryotic cells contain different organelles/compartments. Each organelle serves a purpose, e.g., ribosomes chain together new proteins while mitochondria synthesize ATP. Proteins are the machinery used to perform these functions, including transport in and out and communication between different organelles and a cell’s environment. For some compartments, e.g., the nucleus, special stretches of amino acids, e.g., nuclear localization signals (NLS), help identifying a protein’s location via simple string matching. However, for many others, the localization signal is diluted within the whole sequence, requiring sequence-level predictions. Furthermore, some organelles (and the cell itself) feature membranes with different biochemical properties than the inside or outside, requiring protein gateways.

Homology-inference. Two highly similar protein sequences will most likely fold in similar 3D structures and more likely to perform similar functions. Homology based inference (Nair & Rost, 2002; Mahlich et al., 2018), which transfers annotations of experimentally validated proteins to query protein sequences, is based on this assumption (Sander & Schneider, 1991). Practically this means searching a database of annotated protein sequences for sequences that meet both an identity threshold and a length-of-match threshold to some query protein sequence. Sequence homology delivers good results, but its stringent requirements

render it applicable to only a fraction of proteins (Rost, 1999).

Machine learning Function Prediction. When moving into territory where sequence similarity is less conserved for shorter stretches of matching sequences (Mahlich et al., 2018; Rost, 2002), one can try predicting function using evolutionary information and machine learning (Goldberg et al., 2012; Almagro Armenteros et al., 2017). Evolutionary information from protein profiles, encoding a protein’s evolutionary path, is obtained by aligning sequences from a protein database to a query protein sequence and computing conservation metrics at the residue level. Using profiles leads to impressively more accurate predictions for sequences with no close homologs and has been the standard for most protein prediction tasks (Urban et al., 2020), including subcellular localization (Goldberg et al., 2012; Almagro Armenteros et al., 2017; Savojardo et al., 2018). While profiles provide a strong and useful inductive bias, their information content heavily depends on a balance of the number of similar proteins (depth), the overall length of the matches (sequence coverage), the diversity of the matches (column coverage), and their generation is parameter sensitive.

2. Hyperparameters

The following describes the search space used to find hyperparameters of our final LA and FFN models. We performed random search over these parameters. The evaluated learning rates were in the range of $[5 \times 10^{-6} - 5 \times 10^{-3}]$. For the light attention architecture, we tried filter sizes $[3, 5, 7, 9, 11, 13, 15, 21]$ and hidden sizes $d_{out} \in [32, 128, 256, 512, 1024, 1500, 2048]$, as well as concatenating outputs of convolutions with different filter sizes. For the FFN, we searched over the hidden layer sizes $[16, 32, 64, 512, 1024]$, where 32 was the optimum. We maximized batch size to fit a Quadro RTX 8000 with 48GB vRAM, resulting in the batch size of 150. Note that the memory requirement is dependent on the size of the longest sequence in a batch. In the DeepLoc dataset, the longest sequence had 13 100 residues.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

055
056 *Table 1.* Accuracy and Matthew's correlation coefficient (MCC)
057 on *setDeepLoc*.

METHOD	ACCURACY	MCC
LOCTREE2	61.20	0.525
MULTILOC2	55.92	0.487
SHERLOC2	58.15	0.511
YLOC	61.22	0.533
CELLO	55.21	0.454
iLOC-EUK	68.20	0.641
WOLF PSORT	56.71	0.479
DEEPLOC62	73.60	0.683
DEEPLOC	77.97	0.735
AT SEQVEC	60.97	0.508
AT PROTBERT	64.85	0.567
AT PROT5	71.89	0.661
FFN SEQVEC	70.57± 0.93	0.636± 0.011
FFN PROTBERT	75.88± 0.45	0.702± 0.006
FFN PROT5	79.20± 0.55	0.749± 0.007
LA SEQVEC	75.63± 0.11	0.705± 0.002
LA PROTBERT	80.29± 0.21	0.762± 0.002
LA PROT5	83.37± 0.24	0.800± 0.003

077 *Table 2.* Accuracy and Matthew's correlation coefficient (MCC)
078 on *setHARD*.

METHOD	ACCURACY	MCC
DEEPLOC62	56.94	0.476
DEEPLOC	51.36	0.410
AT PROTBERT	42.04	0.306
AT PROT5	47.14	0.368
FFN PROTBERT	53.16± 1.19	0.429± 0.014
FFN PROT5	55.31± 1.04	0.457± 0.012
LA PROTBERT	58.36± 1.02	0.490± 0.012
LA PROT5	60.92± 0.82	0.522± 0.010

3. Additional Results

We provide results for both *setDeepLoc* (Table 1) and *setHARD* (Table 2) in tabular form, including the Matthew's Correlation Coefficients (MCC).

4. Datasets

Table 3 shows the distribution of subcellular localization classes in the *setDeepLoc* and our new *setHARD*.

4.1. New test set creation

In the following, we lay out the steps taken to produce the new test set (*setHARD*). The starting point is a filtered UniProt search with options as selected in Figure 1. Python code used is available here: *OMITTED*.

- Download data as FASTA & XML:

```
 wget "https://www.uniprot.org/
```

058
059 *Table 3.* Number of proteins and percentage of dataset for each
060 class for the DeepLoc dataset and our *setHARD*. ER abbreviates
061 Endoplasmatic Reticulum

LOCATION	DEEPLOC		SETHARD	
	#	%	#	%
NUCLEUS	4043	28.9	99	20.2
CYTOPLASM	2542	19.3	117	23.8
EXTRACELLULAR	1973	14.0	92	18.8
MITOCHONDRION	1510	11.8	10	2.0
CELL MEMBRANE	1340	9.5	98	20.0
ER	862	6.2	34	6.9
PLASTID	757	5.4	11	2.6
GOLGI APPARATUS	356	2.6	13	2.6
LYSOSOME/VACUOLE	321	2.3	13	2.2
PEROXISOME	154	1.1	3	0.6

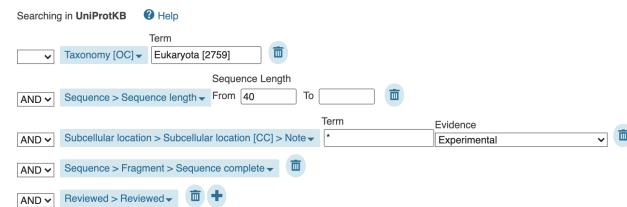


Figure 1. Screenshot of the filtering options applied to the advanced UniProt search (uniprot.org/uniprot).

```
uniprot/?query=taxonomy:%
22Eukaryota%20[2759]%22%
20length:[40%20TO%20*]%
20locations:(note:/*%20evidence:%
22Inferred%20from%20experiment%
20[ECO:0000269]%22)%20fragment:no%
20AND%20reviewed:yesformat=
xmlforce=true sort(scorecompress=
yes"
```

```
wget "https://www.uniprot.org/
uniprot/?query=taxonamy:%
22Eukaryota%20[2759]%22%
20length:[40%20TO%20*]%
20locations:(note:/*%20evidence:%
22Inferred%20from%20experiment%
20[ECO:000026%22)%20fragment:no%
20AND%20reviewed:yesformat=
fastaforce=true sort(scorecompress=
yes"
```

- Download deeploc data:

```
 wget http://www.cbs.dtu.dk/services/
DeepLoc-1.0/deeploc_data.fasta
```

- Align sequences in swissprot to deeploc that have more than 20% PIDE:

```

110 mmseqs easy-search swissprot.fasta
111 deeploc_data.fasta -s 7.5
112 --min-seq-id 0.2 --format-output
113 query,target,fident,alnlen,mismatch,
114 gapopen,qstart,qend,tstart,tend,
115 evalue,bits,pident,nident,qlen,tlen,
116 qcov,tcov alignment.m8 tmp
117
118 • Extract localizations from SwissProt XML:
119 python extract_localizaiotns_from_
120 swissprot.py
121
122 • Map deeploc compartments on swissprot localiza-
123 tions & remove duplicates ([P123, Nucleus] appear-
124 ing twice), remove multilocated ([P123, Nucelus] and
125 [P123, Cytoplasm] → remove P123) empty or not
126 experimental annotations:
127 python map_and_filter_swissprot-
128 annotations.py
129
130 • Create FASTA like deeploc from sequences not in
131 alignment:
132 python extract_unaligned-
133 sequences.py
134
135 • Redundancy reduce new set to 20%:
136 mmseqs easy-cluster --min-seq-id
137 0.2 new_test_set_not_redundancy-
138 reduced.fasta new_hard_test_set-
139 PIDE20.fasta tmp
140
141
142
143
```

References

- 144 Almagro Armenteros, J. J., Sønderby, C. K., Sønderby, S. K.,
145 Nielsen, H., and Winther, O. DeepLoc: prediction of pro-
146 tein subcellular localization using deep learning. *Bioinfor-
147 matics*, 33(21):3387–3395, November 2017. ISSN 1367-
148 4803. doi: 10.1093/bioinformatics/btx431. URL <https://academic.oup.com/bioinformatics/article/33/21/3387/3931857>. tex.ids: al-
149 magroarmenterosDeepLocPredictionProtein2017a
150 publisher: Oxford Academic.
- 151
- 152
- 153
- 154 Goldberg, T., Hamp, T., and Rost, B. LocTree2 predicts
155 localization for all domains of life. *Bioinformatics*, 28
156 (18):i458–i465, September 2012.
- 157
- 158 Mahlich, Y., Steinegger, M., Rost, B., and Bromberg,
159 Y. HFSP: high speed homology-driven function an-
160 notation of proteins. *Bioinformatics*, 34(13):i304–
161 i312, July 2018. ISSN 1367-4803. doi: 10.1093/
162 bioinformatics/bty262. URL <https://doi.org/10.1093/bioinformatics/bty262>.
- 163
- 164

Nair, R. and Rost, B. Sequence conserved for subcellu-
lar localization. *Protein Science*, 11(12):2836–2847,
2002. ISSN 1469-896X. doi: <https://doi.org/10.1110/ps.0207402>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1110/ps.0207402>.

Rost, B. Twilight zone of protein sequence alignments.
Protein Engineering, Design and Selection, 12(2):85–
94, February 1999. ISSN 1741-0126. doi: 10.
1093/protein/12.2.85. URL <https://doi.org/10.1093/protein/12.2.85>.

Rost, B. Enzyme Function Less Conserved than
Anticipated. *Journal of Molecular Biology*, 318
(2):595–608, April 2002. ISSN 0022-2836.
doi: 10.1016/S0022-2836(02)00016-5. URL
<http://www.sciencedirect.com/science/article/pii/S0022283602000165>.

Sander, C. and Schneider, R. Database of homology-
derived protein structures and the structural meaning
of sequence alignment. *Proteins: Structure, Function,
and Bioinformatics*, 9(1):56–68, 1991. ISSN 1097-0134.
doi: <https://doi.org/10.1002/prot.340090107>. URL
<https://onlinelibrary.wiley.com/doi/abs/10.1002/prot.340090107>.

Savojardo, C., Martelli, P. L., Fariselli, P., Profi-
ti, G., and Casadio, R. BUSCA: an integrative web
server to predict subcellular localization of pro-
teins. *Nucleic Acids Research*, 46(W1):W459–
W466, 2018. ISSN 0305-1048. doi: 10.1093/nar/gky320. URL <https://doi.org/10.1093/nar/gky320>. eprint: <https://academic.oup.com/nar/article-pdf/46/W1/W459/25110557/gky320.pdf>.

Urban, G., Torrisi, M., Magnan, C. N., Pollastri, G.,
and Baldi, P. Protein profiles: Biases and proto-
cols. *Computational and Structural Biotechnology
Journal*, 18:2281 – 2289, 2020. ISSN 2001-0370.
doi: <https://doi.org/10.1016/j.csbj.2020.08.015>. URL
<http://www.sciencedirect.com/science/article/pii/S2001037020303688>. tex.ids:
urbanProteinProfilesBiases2020a.