

# Light Attention Predicts Protein Location from the Language of Life

Hannes Stärk, Christian Dallago,  
Michael Heinzinger, Burkhard Rost

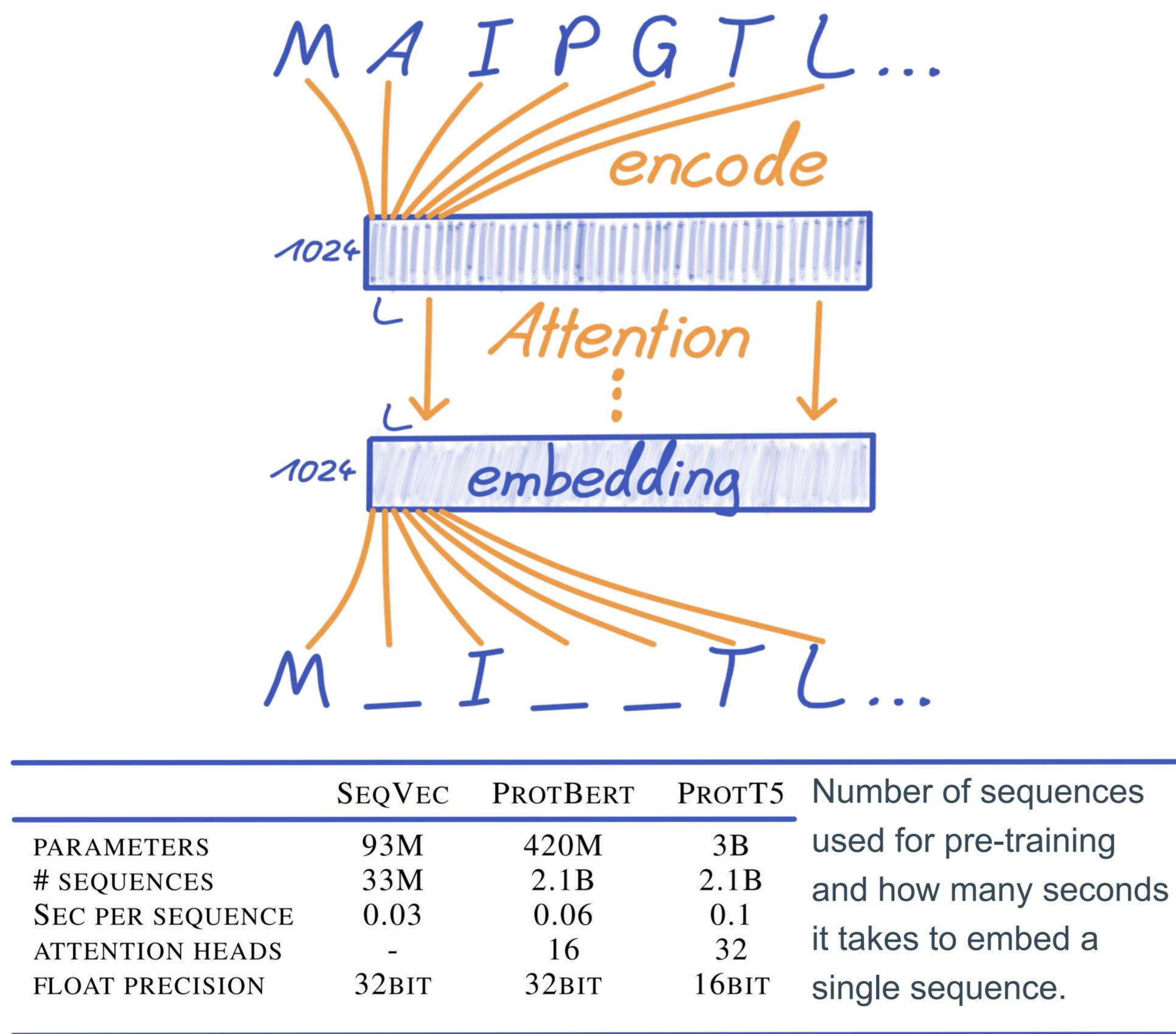
## Subcellular localization predictions with 8 percentage points higher accuracy than before

Proteins are the machinery of life involved in all essential biological processes. A part of finding out how they function is knowing their location in a cell. To predict this subcellular localization, we use language models that were pre-trained on protein sequences to generate protein embeddings. These are processed by our very simple and low-size Light Attention (LA) architecture to beat the previous methods. Advantages:

- Much more accurate
- Faster
- Only sequence input required
- No database search
- Cheap: training and inference on consumer hardware
- Simple!

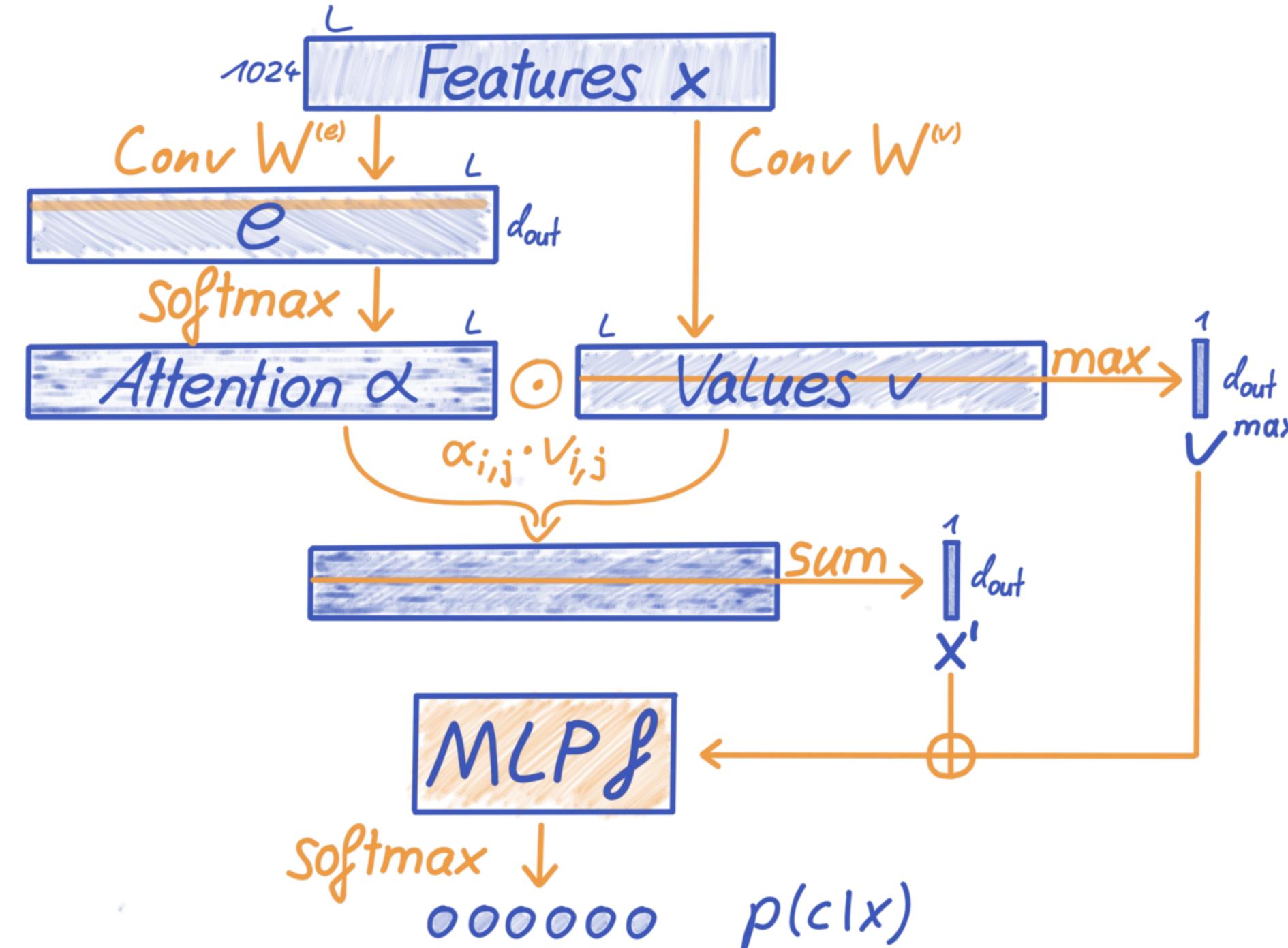
## Protein Language Models

We pre-train three language models either with masked language modeling or next-token-prediction. Each amino acid in a sequence is seen as a token. We try the LSTM SeqVec, and two transformers: ProtBert and ProtT5.



## Light Attention Architecture

Takes protein language model embeddings as input and classifies their subcellular localization.



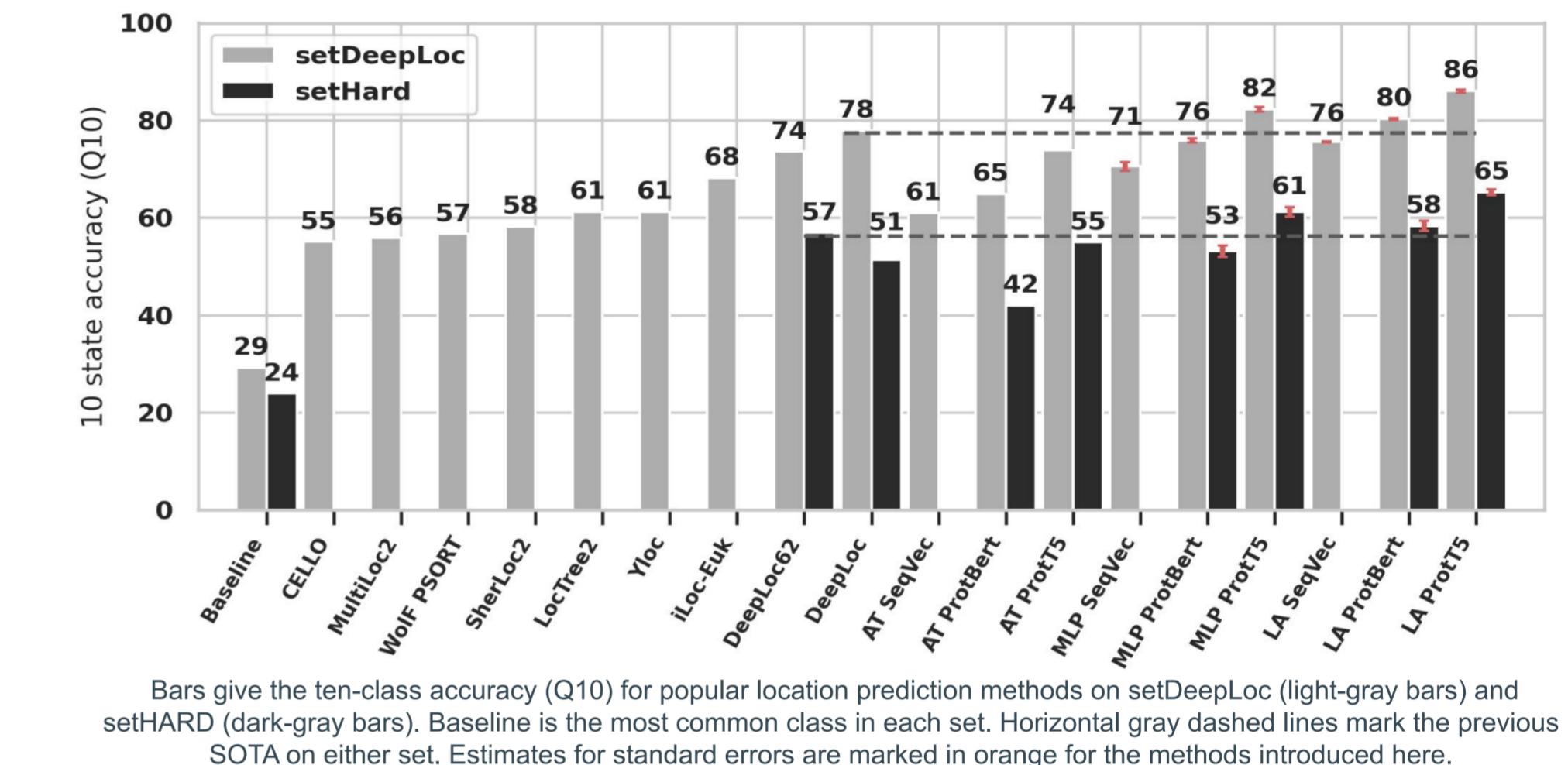
The only learnable weights are in the two convolutions and the final MLP. The “Light Attention” part is simply a weighted sum where the weights are given by a learned function of the inputs.

## Localization prediction Datasets

We train on the standard localization dataset DeepLoc. It has 10 different subcellular localization classes. We evaluate on the DeepLoc [1] test set, which is redundancy reduced to the training data to **30% PIDE**. We also evaluate on a new, harder test set with at most **20% PIDE** to the train set.

LOCATION	DEEPLOC #	DEEPLOC %	SETHARD #	SETHARD %	The number of sequences in each class and the percentage the class makes up of the total set. Shown for the DeepLoc test set and our new harder test set.
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	

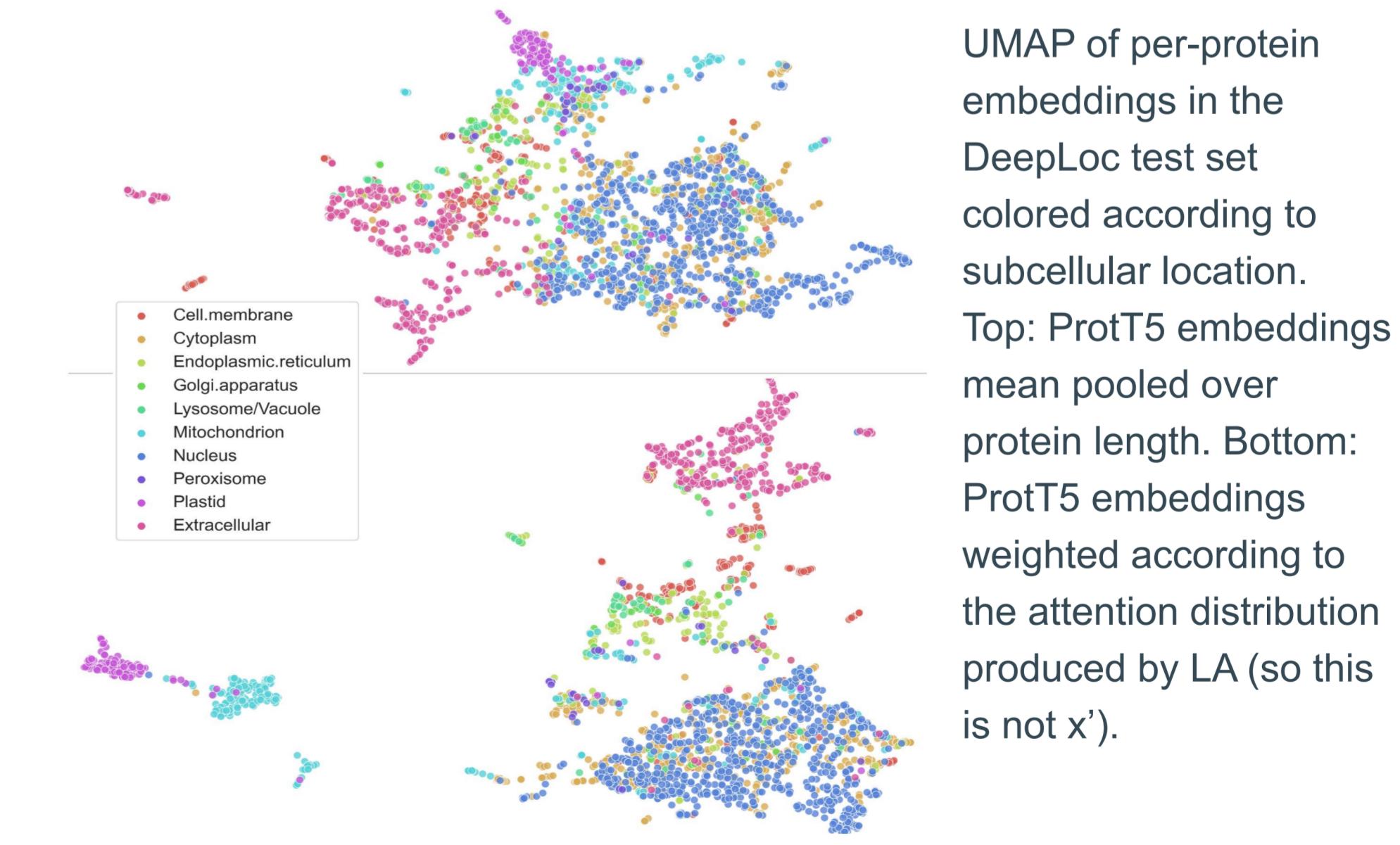
## Comparison with previous Methods



The LA architecture outperforms the previous SOTA by **8** percentage points on both test sets. Also, DeepLoc is worse on setHARD with profile input than with BLOSUM. Maybe profiles just leverage sequence redundancy?

## Other methods on ProtT5 embeddings

METHOD	SETDEEPLOC	SETHARD	Accuracy of additional baselines and ablations:
LA	<b>86.01</b> ± 0.34	<b>65.21</b> ± 0.61	drop softmax, drop maxpool, calculate attention coefficients from v, DeepLoc LSTM, stacked convolutions
LA - SOFTMAX	85.30 ± 0.32	64.72 ± 0.70	
LA - MAXPOOL	84.79 ± 0.19	63.84 ± 0.67	
ATTENTION FROM v	85.41 ± 0.27	64.77 ± 0.93	
DEEPLOC LSTM	79.40 ± 0.88	59.36 ± 0.84	
CONV + ADAPOOL	82.09 ± 0.92	60.79 ± 2.01	
MEANPOOL + MLP	82.27 ± 0.51	61.27 ± 0.97	
LA ON ONEHOT	43.53 ± 1.48	32.57 ± 2.38	
LA ON PROFILES	43.78 ± 1.25	33.35 ± 1.82	



[1] Almagro Armenteros et al. (2017). “DeepLoc: prediction of protein subcellular localization using deep learning.” In: Bioinformatics, 33(21):3387–3395.