



CARP

bioRxiv

22

Representation

- Convolutional masked language model
- ByteNet encoder
- Dilated convolution

Autoencoder

CNN

Convolutions are competitive with transformers for protein sequence pretraining

Protein design aims to build new proteins from scratch thereby holding the potential to tackle many environmental and biomedical problems. Recent progress in the field of natural language processing (NLP) has enabled the implementation of ever-growing language models capable of understanding and generating text with human-like capabilities. Given the many similarities between human languages and protein sequences, the use of NLP models offers itself for predictive tasks in protein research. Motivated by the evident success of generative Transformer-based language models such as the GPT-x series, we developed ProtGPT2, a language model trained on protein space that generates *de novo* protein sequences that follow the principles of natural ones. In particular, the generated proteins display amino acid propensities which resemble natural proteins. Disorder and secondary structure prediction indicate that 88% of ProtGPT2-generated proteins are globular, in line with natural sequences. Sensitive sequence searches in protein databases show that ProtGPT2 sequences are distantly related to natural ones, and similarity networks further demonstrate that ProtGPT2 is sampling unexplored regions of protein space. AlphaFold prediction of ProtGPT2-sequences yielded well-folded non-idealized structures with embodiments as well as large loops and revealed new topologies not captured in current structure databases. ProtGPT2 has learned to speak the protein language. It has the potential to generate *de novo* proteins in a high throughput fashion in a matter of seconds. The model is easy-to-use and freely available.

Design and Engineering

Deep learning and architecture

Annotation

ReLeaSE
#RL #SMILES

GEGL
#RL #GA #SMILES

SAC
#RL

PPO
#RL #PO

Decision Transformer
#RL #Transformer

Color box
Reviewed!
Click to see the summary

Grey box
Not yet reviewed

DyNA-PPO
#RL #Binding #PO

ProGPT2
#GPT #De-novo

3D-MolGNN
#Graph #Actor-critic

eUniRep
#LSTM

MSA2Prot
#Transformer

UniRep
#LSTM

CARP
#CNN

GVP
#Graph

ProteinBERT
#BERT

LM-GVP
#Graph #LM

DeepFRI
#GCN

Hallucination
#MCMC

Augmentation

ProteinGAN
#GAN #Sequence

AR-VAE
#VAE #Sequence

AlphaFold2
#Graph #Attention

ProteinQA
#Graph #QA

trRosetta
#ResNet

ProteinSolver
#Graph #Inverse

ProteinGCN
#GCN #QA

GVP+AF
#Inverse #Structure

Folding / Inverse folding

DeepGO
#CNN

DeepGOPlus
#CNN

DeepGOZero
#CNN #Zero-shot

TAPE
#Benchmark #Transformer

GO
#Benchmark