# Final Project

# DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

Piyar Ali , Ali



CSCI E-104 Advanced Deep Learning, 2023

**Harvard University Extension School**

Prof. Zoran B. Djordjević

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

# Contents

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

# Project Summary:

[DeepPurpose](#) is a powerful deep learning library that specializes in drug discovery and drug repurposing tasks. The library employs various deep learning models such as GCN and RNN to learn representations of molecular structures and predict their properties. DeepPurpose is equipped with a variety of features that include drug-target interaction prediction, compound property prediction, protein-protein interaction prediction, and protein function prediction.

However, it does not provide protein to protein amino acid sequence similarity. Protein sequence pairwise similarity is a valuable tool that can aid in the identification and characterization of drug targets, as well as optimizing drug design and minimizing off-target effects. Our project aims to leverage the built-in target encoders in DeepPurpose to encode protein sequence data from uniport to train a model to calculate the similarity pairwise scores.

DeepPurpose protein to protein interaction prediction utilizes affinity scores. Protein affinity scores, or binding affinity scores, measure the attraction between proteins and ligands and help us understand how proteins and ligands interact with each other. In contrast, pairwise similarity is more focused on comparing the amino acid sequence and inferring evolutionary relationships between proteins. While both protein affinity scores and protein pairwise similarity provide information about proteins, they measure different aspects of protein biology.

To verify the protein to protein pairwise scores generated by our model, we can use [NCBI BLAST](#), which is a widely used algorithm developed and maintained by the National Center for Biotechnology Information (NCBI) for comparing and analyzing biological sequences. BLAST can provide alignment scores if the two proteins are aligned. Overall, DeepPurpose is a powerful tool for drug discovery and drug repurposing tasks, and it offers a wide range of functionalities and tools for various aspects of drug discovery.

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

# Problem Statement:

DeepPurpose is a powerful deep learning library for drug discovery, but it lacks protein to protein amino acid sequence similarity, which is an essential tool for identifying and characterizing drug targets, optimizing drug design, and minimizing off-target effects. The objective of the project is to leverage the built-in target encoders in DeepPurpose to encode protein sequence data from uniport and train a model to calculate the similarity pairwise scores. To verify the protein to protein pairwise scores generated by the model, the project will use NCBI BLAST, which is a widely used algorithm for comparing, analyzing, and scoring biological sequences.
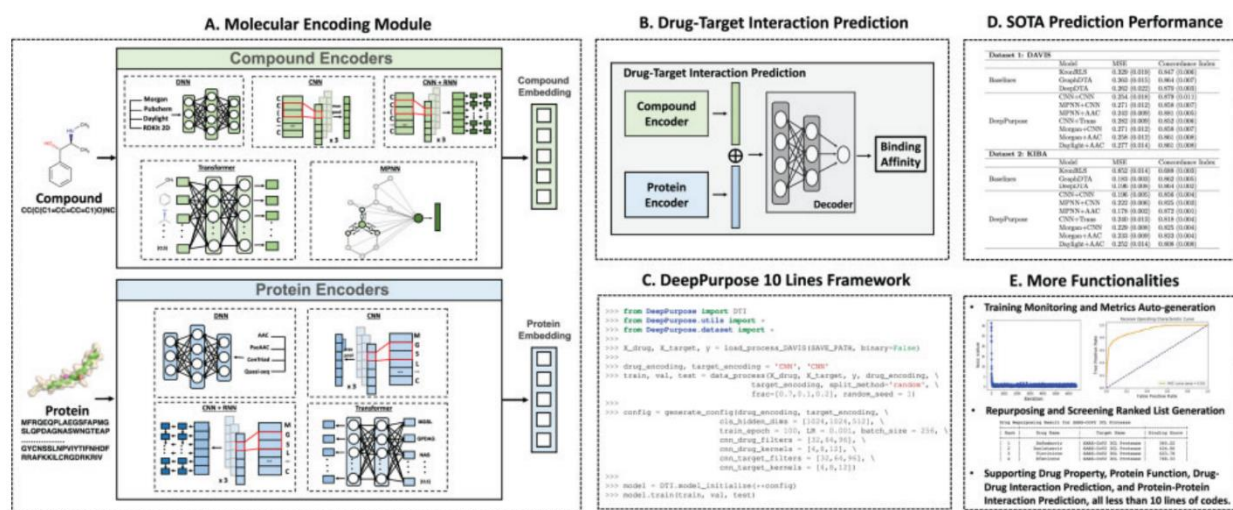
## DeepPurpose Library



Image Source: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8016467/pdf/btaa1005.pdf

Fig. 1. Overview of DeepPurpose library. (A) DeepPurpose takes as input the SMILES of a compound and a protein's amino acid sequence and then generates embeddings for them. (B) The learned embeddings are then concatenated and fed into a decoder to predict DTI binding affinity. (C) DeepPurpose provides a simple but flexible programming framework that implements over 50 state-of-the-art DL models for DTI prediction. (D) DeepPurpose models achieve comparable performance with three other DTI prediction algorithms on two benchmark datasets. (E) Finally, DeepPurpose has many functionalities, including monitoring the training process, debugging and generation ranked lists for repurposing and screening. Further, DeepPurpose supports other downstream prediction tasks (e.g., drug–drug interaction prediction, compound property prediction)

@Ali PiyarAli

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

## Protein Sequence Similarity



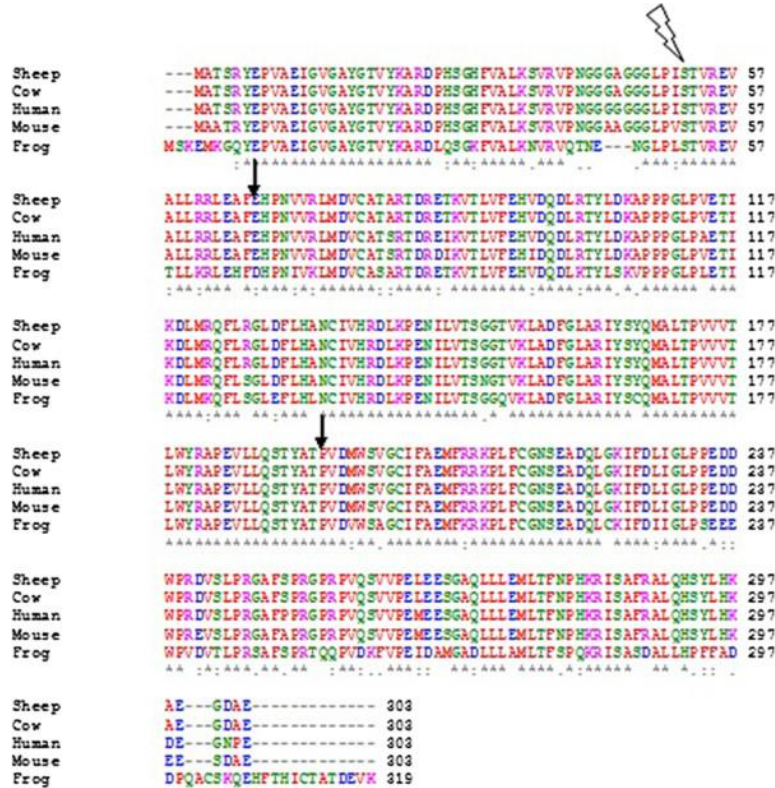Image Source: https://www.wikiwand.com/en/Clustal

## Dataset:

## DeepPurpose
*Public Drug-Target Binding Benchmark Dataset*

| Data | Function |
|------|----------|
| BindingDB | `download_BindingDB()` to download the data and `process_BindingDB()` to process the data |
| DAVIS | `load_process_DAVIS()` to download and process the data |
| KIBA | `load_process_KIBA()` to download and process the data |

*Repurposing Dataset*

@Ali PiyarAli

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

| Data | Function |
|---|---|
| Curated Antiviral Drugs Library | `load_antiviral_drugs()` to load and process the data |
| Broad Repurposing Hub | `load_broad_repurposing_hub()` downloads and process the data |

*Bioassay Data for COVID-19* (Thanks to MIT AI Cures)

| Data | Function |
|---|---|
| AID1706 | `load_AID1706_SARS_CoV_3CL()` to load and process |

Protein Sequence Data

| Data | Function |
|---|---|
| Isoform sequences | `list(SeqIO.parse(seq_path, 'fasta'))` |

# Installation/configuration

The project is setup in Python Notebook using Google Colab

## Install Libraries

```
!pip install rdkit
!pip install git+https://github.com/bp-kelley/descriptastorus
!pip install DeepPurpose
!pip install Bio
!pip install PySpark
!pip install wget
```

## Import Libraries

```
from DeepPurpose import DTI as modelsdti
from DeepPurpose import PPI as modelsppi
from DeepPurpose.utils import *
from DeepPurpose.dataset import *
from DeepPurpose import dataset, CompoundPred, utils, encoders


import Bio
from Bio.Blast import NCBIWWW, NCBIXML
```

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

```python
from Bio import SeqIO
from Bio.SeqRecord import SeqRecord
from Bio.Seq import Seq

import torch
import torch.nn as nn
import torch.optim as optim

import pandas as pd
import numpy as np
import time
import random
import csv
import os
import wget

import multiprocessing
from multiprocessing import pool

import matplotlib.pyplot as plt
```

## Download Dataset

```python
# Download the affinity scores file

!wget -P /content http://staff.cs.utu.fi/~aatapa/data/DrugTarget/drug-
target_interaction_affinities_Kd__Davis_et_al.2011.txt

# Get the fasta (protein sequence data) from
https://www.uniprot.org/help/downloads

!wget -P /content
https://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebas
e/complete/uniprot_sprot_varsplic.fasta.gz

# extract/unzip the file

!gunzip /content/uniprot_sprot_varsplic.fasta.gz

# set the fasta file path

seq_path = '/content/uniprot_sprot_varsplic.fasta'
```

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

## Load Dataset

```
# Read and extract protein records using Bio library


seq_records = list(SeqIO.parse(seq_path, 'fasta'))
```

*Note: the dataset is large total records 40,957, possible Paris 40,957 x 40,957 = 1,677,475,849*

*for model training purpose, use only first 200 records.*

```
seq_records = list(SeqIO.parse(seq_path, 'fasta'))[:200]
```

```
# Extract Sequence Pair
seqs = []
for i in range(len(seq_records)):
  seqs.append(str(seq_records[i].seq))
```

## DeepPurpose Encoders

```
# Encode the protein sequence pair using DeepPurpose protein encoders


X = utils.encode_protein(pd.DataFrame(seq, columns=['Target Sequence']),
'AAC')
```

*Note: Encoding AAC takes time. Time Reference: 24s for ~100 sequences in a CPU*

**DeepPurpose provides following different encoders, you can use any of them to improve the model:**

| Target Encodings | Description |
|------------------|-------------|
| AAC | Amino acid composition up to 3-mers |
| PseudoAAC | Pseudo amino acid composition |
| Conjoint_triad | Conjoint triad features |
| Quasi-seq | Quasi-sequence order descriptor |
| ESPF | Explainable Substructure Partition Fingerprint |

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

| CNN | Convolutional Neural Network on target seq |
|-----|--------------------------------------------|
| CNN_RNN | A GRU/LSTM on top of a CNN on target seq |
| Transformer | Transformer Encoder on ESPF |

```python
# Convert the target encoding to a tensor
X_tensor = torch.tensor(X['target_encoding'].values.tolist(),
dtype=torch.float)
```

## Similarity Model

```python
# Instantiate the similarity scoring model
model = nn.Sequential(
    nn.Flatten(),
    nn.Linear(16840, 256),
    nn.ReLU(),
    nn.Linear(256, 64),
    nn.ReLU(),
    nn.Linear(64, 1),
    nn.Sigmoid()
)
```

```python
# Define the optimization algorithm and loss function
optimizer = optim.Adam(model.parameters(), lr=0.001)
criterion = nn.MSELoss()
```

```python
# Train the model on a pair of protein sequences
x1 = X_tensor[0].unsqueeze(0)
x2 = X_tensor[1].unsqueeze(0)
y = torch.tensor([[0.8]], dtype=torch.float)
```

```python
losses = []

for epoch in range(50):
    optimizer.zero_grad()

    # Concatenate x1 and x2 along the column axis
```

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

```python
    x = torch.cat((x1, x2), dim=1)

    output = model(x)
    loss = criterion(output, y)
    loss.backward()
    optimizer.step()

    losses.append(loss.item())

    print('Epoch {}: Loss = {:.4f}'.format(epoch+1, loss.item()))
```
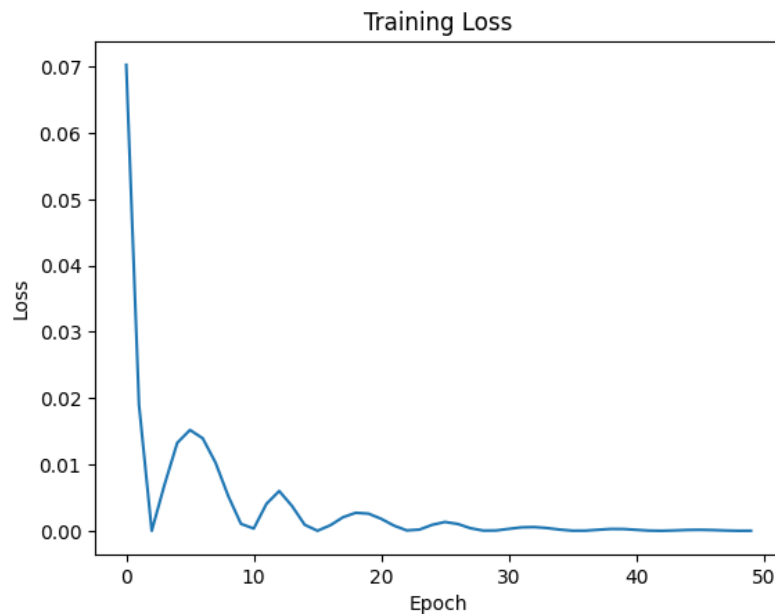
# Results

## Training Loss



## Similarity Score

For the first two pairs

```python
# Calculate the similarity score for the pair of protein sequences
```

@Ali PiyarAli

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

```
score = model(torch.cat((x1, x2), dim=1)).item()
print('Similarity score: {:.4f}'.format(score))
Result: Similarity score: 0.7929
```

x1 =
MENERAKQVYLAKLNEQAERYDEMVEAMKKVAALDVELTIEERNLLSVGYKNVIGARRASWRILSSIEQKEESKGNEQN
AKRIKDYRTKVEEELSKICYDILAVIDKHLVPFATSGESTVFYYKMKGDYFRYLAEFKSGADREEAADLSLKAYEAATSSASTE
LSTTHPIRLGLALNFSVFYYEILNSPERACHLAKRAFDEAIAELDSLNEDSYKDSTLIMQLLRDNLTLWTSDLEEGGK

x2 =
MENERAKQVYLAKLNEQAERYDEMVEAMKKVAALDVELTIEERNLLSVGYKNVIGARRASWRILSSIEQKEESKGNEQN
AKRIKDYRTKVEEELSKICYDILAVIDKHLVPFATSGESTVFYYKMKGDYFRYLAEFKSGADREEAADLSLKAYEAATSSASTE
LSTTHPIRLGLALNFSVFYYEILNSPERACHLAKRAFDEAIAELDSLNEDSYKDSTLIMQLLRDNLTLWTSDLEEGGEQSKG
HNQQDEVNKI

## Verification

The scores can be verified using NCBI Blast: https://blast.ncbi.nlm.nih.gov/Blast.cgi

NCBI BLAST (Basic Local Alignment Search Tool) is a widely used algorithm developed and maintained by the National Center for Biotechnology Information (NCBI) for comparing and analyzing biological sequences, such as DNA, RNA, and protein sequences. Blast can provide alignment scores, if the two proteins are aligned. Blast can be access via web interface or using python library.

Using python library is very slow if verifying large dataset of protein-to-protein pairwise similarity sequence as the API has a limit of 15 seconds wait after each protein-to-protein pair request. Each request response takes approximately 30 seconds. For example, sending an API request for 200 protein pairs can take approximately 1.5 to 2 hours.

Python notebook provides the code for accessing the NCBI Blast using python. For the purposes of this project, we used the web interface.

## Web Interface Verification

Enter the sequence in the NCBI Blast Search Interface

@Ali PiyarAli

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction
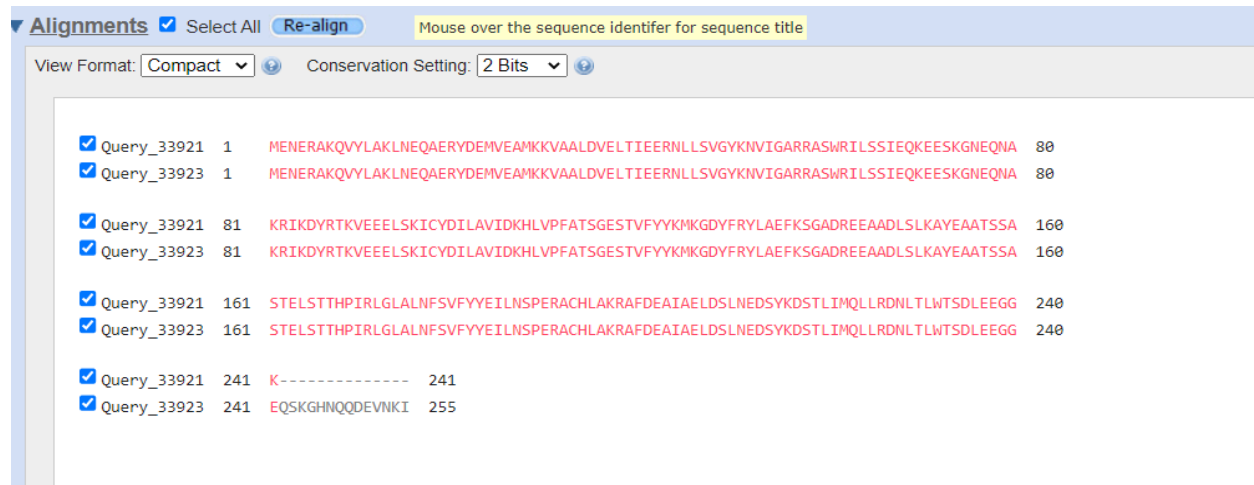


Get the results.



Analyze the Results

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

Our model predicted the similarity score of 0.7929

The NCBI Blast database score is 241/255 = 0.9450

Our model did very good. The model can be improved by experimenting with increasing epochs or changing the model hyper parameters or using different DeepPurpose available encoders.

## Challenges

1. Data format challenge: DeepPurpose requires data in a specific format, but the documentation does not provide clear guidelines on this format. Users must either inspect their dataset or look at the code to determine the required format.

2. Large dataset challenge: The Uniprot protein sequence data is massive, with over 40,000 records, resulting in a possible 1.6 billion protein to protein pairs. Running a model on a good sample size requires a computer with high RAM, even when utilizing big data technologies such as Spark or Multiprocessing. The provided multiprocessing code in Python Notebook can help apply the similarity model on the full dataset, but most computers will run out of memory. This requires either a high-powered machine or chunking the data into batches and running the model on each batch.

**Final Project**
**Name:** Piyar Ali, Ali
CSCI E-104 Advanced Deep Learning, 2023 | Harvard University Extension School | Prof. Zoran B. Djordjević

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction

3.  Verification challenge: Our model similarity scores can be verified using NCBI Blast, either via a web interface or a Python library. However, using the Python library to verify large datasets of protein to protein pairwise similarity is very slow, as the API has a 15-second wait time after each protein-to-protein pair request. Each request response can take approximately 30 seconds. For example, sending an API request for 200 protein pairs can take approximately 1.5 to 2 hours.