
Deep Learning for Protein-Protein Interaction Prediction

— Samuel Sledzieski —

-
- Legend: ▲ TF ● non-TF

The diagram illustrates the central dogma of molecular biology, showing the flow of genetic information from DNA to RNA to Protein.

DNA: A double helix structure is shown with four segments of the template strand (bottom strand) exposed. The sequences of these segments are: A T G C T T C G T A, A T G C T T C G T A, A T G C T T C G T A, and A T G C T T C G T A. The complementary coding strand (top strand) sequences are: T C G A A G C C A, T C G A A G C C A, T C G A A G C C A, and T C G A A G C C A.

Transcription: An arrow labeled "Transcription" points from the DNA to the RNA.

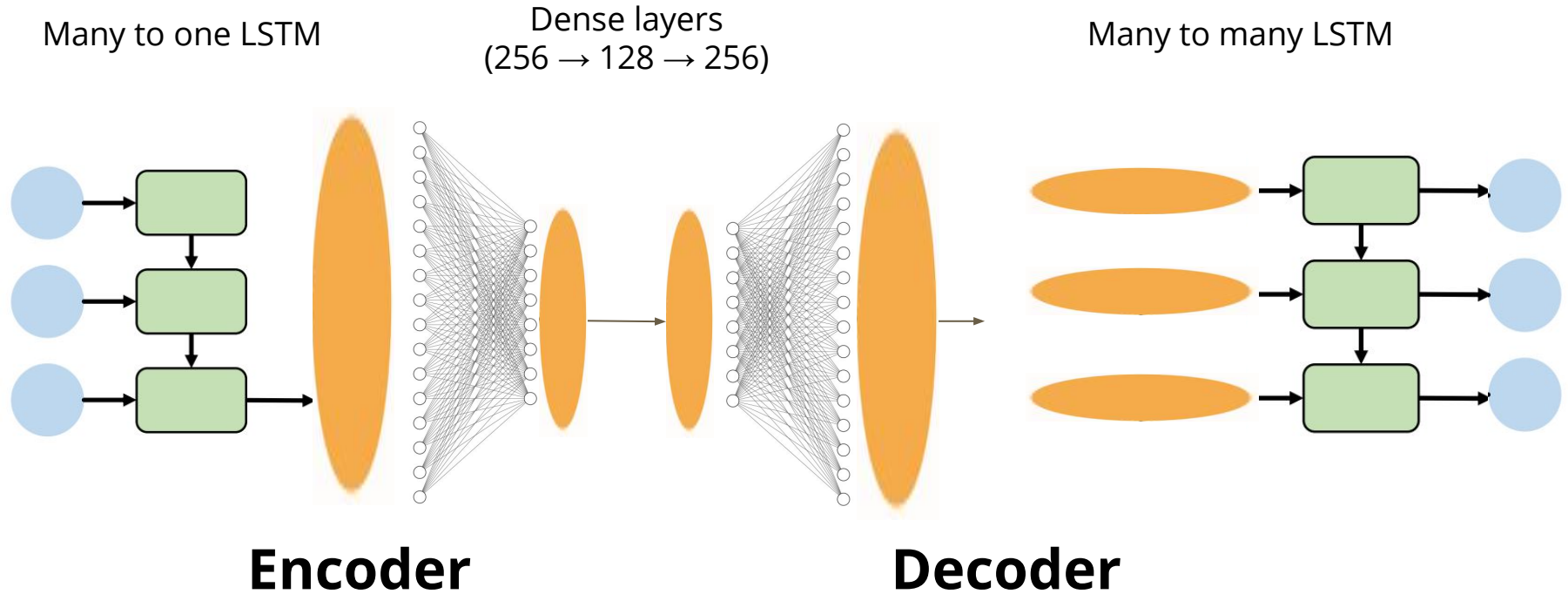
RNA: A single-stranded green line represents the mRNA. It has four segments corresponding to the DNA segments above. The sequences are: A U G C U U U C G U A U, U A C G A A A G C C A U A, A U G C U U U C G U A U, and U A C G A A A G C C A U A.

Translation: An arrow labeled "Translation" points from the RNA to the Protein.

Protein: A black line represents the polypeptide chain. It has four segments corresponding to the RNA segments above. The amino acid sequences are: Met Leu Ser Tyr, Tyr Glu Ser Ile, Met Leu Ser Tyr, and Tyr Glu Ser Ile.

- Protein-Protein Interaction prediction from **sequence only** is a challenging task
- **Our contribution: we present a deep learning framework for PPI prediction**
 - Learn protein sequence embeddings with an **LSTM autoencoder**
 - Introduce the idea of an **interaction fingerprint**
 - Predict protein interaction using a **CNN**
 - Explore PPI **transfer learning** to other species

LSTM Autoencoder Architecture



Protein Interaction Embedding



Subsample data

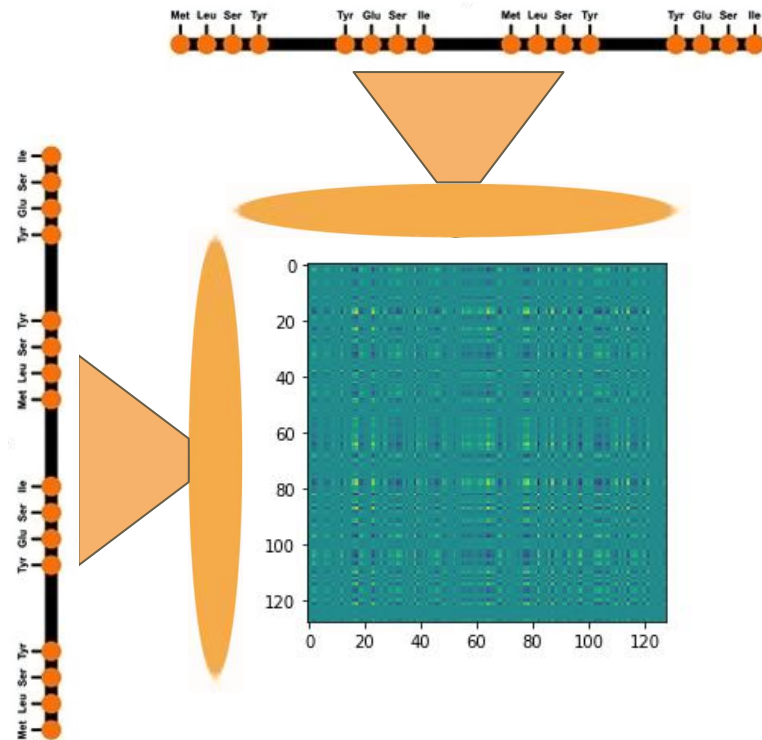
- *Homo sapiens* protein-protein interactions from STRING database
- Sample 2,000 true positive interactions
- Generate 13,000 negative interactions by randomly sampling pairs of proteins
 - Random interaction has ~3% chance of actually being true

Learn latent space of training data

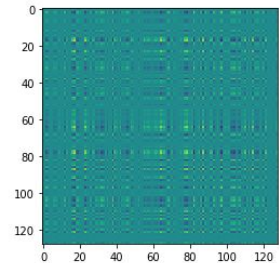
- Use trained LSTM-AE to generate a length 128 vector

Generate interaction “fingerprints”

- Outer product of latent-space embeddings generates a unique image for each protein-protein interaction



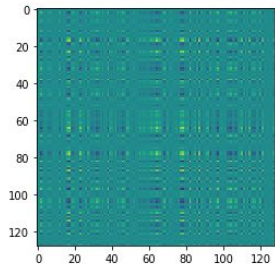
9606.ENSP00000481559 9606.ENSP00000228347
Interaction? False



64@126x126
1@128x128



9606.ENSP00000000233 9606.ENSP00000380185
Interaction? True



CNN Architecture

64@31x31

64@31x31

32@29x29

32@14x14

32@14x14

128
2

Convolution

Max-Pool

Dropout

Max-Pool

Dropout

Convolution

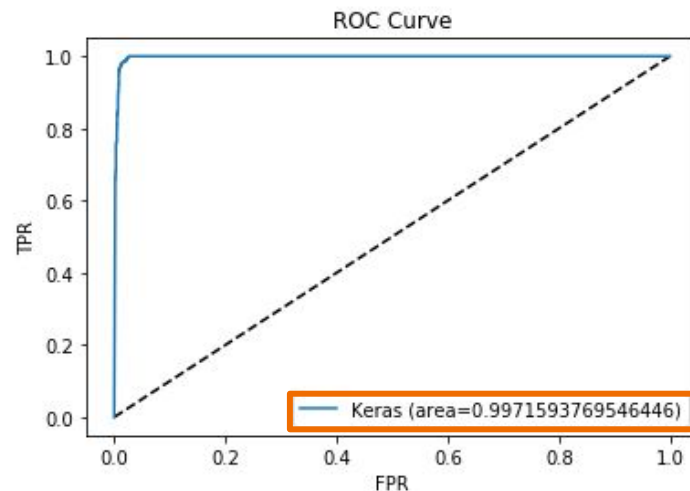
Dense

0	1
0.8991	0.1009

0	1
0.00001	0.99999

Model Training & Evaluation

- 75-25 Train-Test Split
- Trained LSTM-AE on 1,000 samples for 5 epochs, batch size 32
- Trained CNN on 11,251 samples (2,000 positive, 13,000 negative) for 50 epochs, batch size 128
- Evaluated CNN on 3,750 never-seen samples
- Achieved 96.24% training accuracy, 94.13% testing accuracy, 0.997 AUROC



Takeaways:

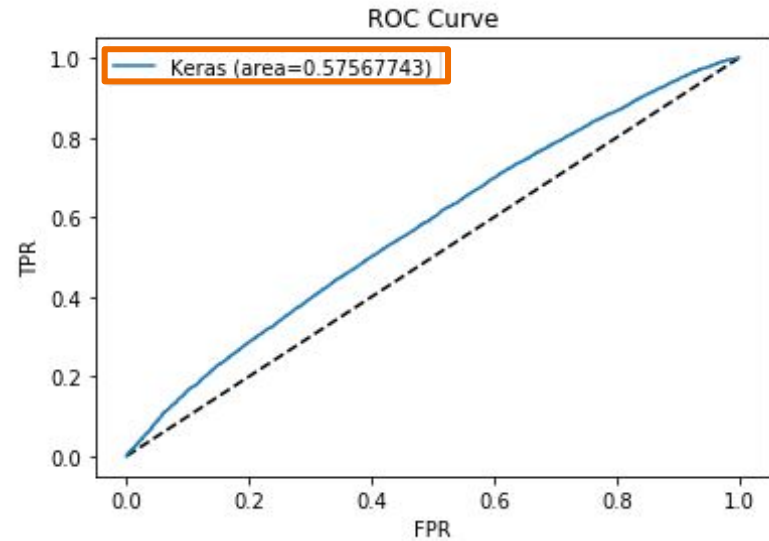
- LSTM-AE training is **slow**
- CNN training is **fast**
- Embedding + outer product aids prediction of protein interaction

```
Epoch 45/50
11251/11251 [=====] - 2s 204us/sample - loss: 0.1224 - accuracy: 0.9525
Epoch 46/50
11251/11251 [=====] - 2s 203us/sample - loss: 0.1196 - accuracy: 0.9548
Epoch 47/50
11251/11251 [=====] - 2s 204us/sample - loss: 0.1135 - accuracy: 0.9556
Epoch 48/50
11251/11251 [=====] - 2s 202us/sample - loss: 0.1152 - accuracy: 0.9539
Epoch 49/50
11251/11251 [=====] - 2s 201us/sample - loss: 0.1095 - accuracy: 0.9571
Epoch 50/50
11251/11251 [=====] - 2s 202us/sample - loss: 0.1022 - accuracy: 0.9624

3750/3750 [=====] - 0s 117us/sample - loss: 0.1440 - accuracy: 0.9413
Test Loss, Test Accuracy: [0.14399235029617946, 0.94133335]
```

Generalization to other species

- Can we apply our trained models to data from different species?
 - *Saccharomyces Cerevisiae* (Yeast)
 - *Drosophila melanogaster* (Fruit fly)
- Sampled 10,000 positive interactions between the two species
- Sampled 10,000 negative interactions with one protein coming from each species
- Embed proteins and predict interaction using pre-trained model
- Use same embeddings but retrain CNN



Takeaways:

- Accuracy is **poor** using same embedding and CNN
- Accuracy **improves** using the same embeddings but retraining the CNN
- Embeddings, and to an extent prediction models, are difficult to transfer across species
- There may be potential for weight sharing / inductive transfer learning to speed training

Download the models and data

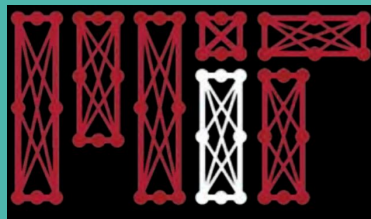
shorturl.at/buyF1

Download the models only

shorturl.at/acCFL



https://github.com/samsledje/Deep_PPI



References in Github README