

# Protein-Protein Interaction Prediction using a Deep Learning Approach based on Autoencoders

**Albu Alexandra**  
**Babeş-Bolyai University**

**WeADL 2022 Workshop**

The workshop is organized under the umbrella of WeaMyL, project funded by the EEA and Norway Grants under the number RO-NO-2019-0133.

Contract: No 26/2020.



Working together for a **green**, **competitive** and **inclusive** Europe

# Autoencoders

- neural networks formed of an encoder and a decoder
- trained to reconstruct their input

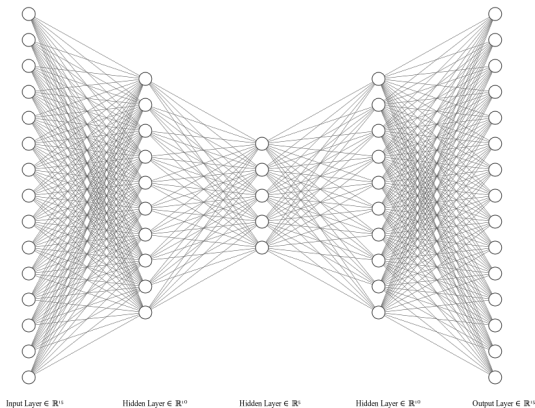
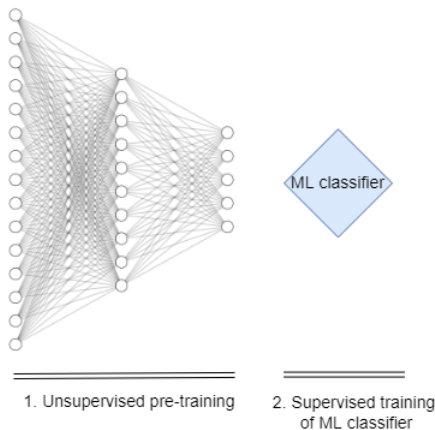


Figure: Autoencoder architecture. Created using <https://alexlenail.me/NN-SVG/>.

# Autoencoders

Autoencoders in classification → three main directions

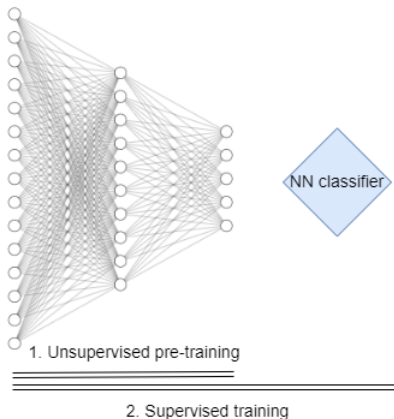
- **Feature extraction:** train a classifier on the learned autoencoder representations → challenge: embedding is performed independently from the classification stage



# Autoencoders

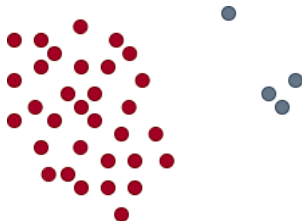
Autoencoders in classification → three main directions

- **Fine-tuning:** fine-tuning the encoder weights together with a neural network classifier



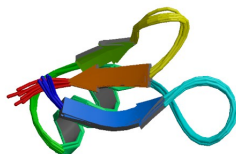
Autoencoders in classification → three main directions

- **Anomaly detection:** train an autoencoder on the majority class, detect outliers



# Protein Data Analysis

- Proteins: complex macromolecules
  - important role in vital biological processes in living organisms.
- Central problem: determining protein functions



**Figure:** 3D view of protein 1I2U. Image obtained from the RCSB PDB<sup>1</sup> representing the protein with PDB ID 1I2U.

---

<sup>1</sup>Berman et al. 2000, *The Protein Data Bank*. <https://www.rcsb.org/>

# Protein-Protein Interaction Prediction

- The majority of proteins perform their roles in complexes
- Experimental determination of PPIs - expensive and prone to false positives → computational approaches try to overcome these limitations
- Challenges: labels are noisy, small number of known interacting pairs compared to non-interacting

# Literature Review. Sequence-based Protein-protein Interaction Prediction

- Machine Learning methods: SVMs<sup>2</sup>, RFs<sup>3</sup>, LightGBM<sup>4</sup>
- Ensembles of machine learning classifiers<sup>5</sup> and ensembles of neural networks<sup>6</sup>

---

<sup>2</sup>Guo et al., 2008, *Using support vector machine combined with auto covariance to predict protein-protein interactions from protein sequences*. Nucleic acids research.

<sup>3</sup>Pan et al., 2010, *Large-Scale prediction of human protein-protein interactions from amino acid sequence based on latent topic features*. Journal of proteome research.

<sup>4</sup>Chen et al., 2019, *Predicting protein-protein interactions through LightGBM with multi-information fusion*. Chemometrics and Intelligent Laboratory Systems.

<sup>5</sup>Chen et al., 2018, *Protein-protein interaction prediction using a hybrid feature representation and a stacked generalization scheme*. BMC Bioinformatics.

<sup>6</sup>Li et al., 2020, *Protein Interaction Network Reconstruction Through Ensemble Deep Learning With Attention Mechanism*. Frontiers in Bioengineering and Biotechnology.



# Literature Review. Sequence-based Protein-Protein Interaction Prediction

- Siamese Architectures: → capture common characteristics of the two proteins in a pair
  - convolutional architecture<sup>7</sup>
  - residual convolutional recurrent architecture<sup>8</sup>
  - Inception convolutional branch and a bidirectional GRU branch<sup>9</sup>

---

<sup>7</sup>Hashemifar et al., 2018, *Predicting protein-protein interactions through sequence-based deep learning*. Bioinformatics

<sup>8</sup>Chen et al., 2019, *Multifaceted protein-protein interaction prediction based on siamese residual RCNN*. BMC Bioinformatics

<sup>9</sup>Zhao et al., 2020, *Conjoint feature representation of go and protein sequence for ppi prediction based on an inception rnn attention network*. Molecular Therapy-Nucleic Acids

# Autoencoder-based Methods in Protein-Protein Interaction Prediction

- Autoencoders as feature extractors + probabilistic SVMs<sup>10,11</sup>
- Autoencoder pretraining + fine-tuning neural network classifier<sup>12</sup>
- Variational graph autoencoder<sup>13</sup> → learn nodes embeddings using the neighbours in the PPI graph

---

<sup>10</sup>Wang et al., 2017, *Predicting protein-protein interactions from protein sequences by a stacked sparse autoencoder deep neural network*. Molecular BioSystems.

<sup>11</sup>Wang et al., 2018, *Predicting protein interactions using a deep learning method-stacked sparse autoencoder combined with a probabilistic classification vector machine*. Complexity.

<sup>12</sup>Sun et al., 2017, *Sequence-based prediction of protein protein interaction using a deep-learning algorithm*. BMC Bioinformatics.

<sup>13</sup>Yang et al., 2020, *Graph-based prediction of protein-protein interactions with attributed signed graph embedding*. BMC bioinformatics.

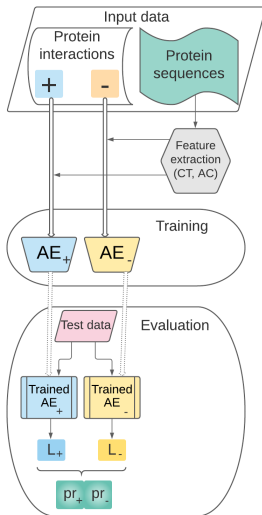
# AutoPPI: Binary classification using a pair of autoencoders

## Approach:

- two autoencoders trained to reconstruct instances belonging to one class
- classification stage: evaluating which of the two autoencoders is able to better reconstruct the testing data point

Czibula, G., Albu, A.I., Bocicor, M.I. and Chira, C., 2021.  
*AutoPPI: An Ensemble of Deep Autoencoders for Protein–Protein Interaction Prediction*. *Entropy*, 23(6), p.643.

# AutoPPI: Binary classification using a pair of autoencoders



# AutoPPI: Binary classification using a pair of autoencoders

- data samples: pairs of proteins  $\rightarrow$  proposed two siamese architectures

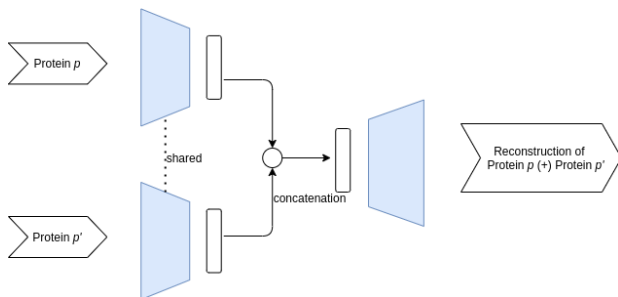


Figure: Siamese-Joint architecture.

# AutoPPI: Binary classification model using a pair of autoencoders

- data samples: pairs of proteins  $\rightarrow$  proposed two siamese architectures

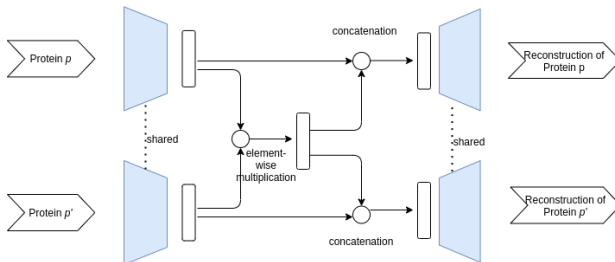


Figure: Siamese-Siamese architecture.

# AutoPPI: Binary classification using a pair of autoencoders

- baseline architecture: simple concatenation of protein features (early fusion)

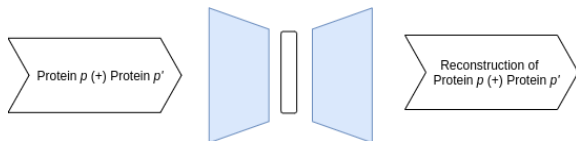


Figure: Joint-Joint architecture.

- **Conjoint Triad (CT)** descriptors
  - group amino acids into seven classes based on their physico-chemical properties
  - compute the frequencies of possible triples of amino acid classes
- **Autocovariance (AC)** descriptors
  - define a *lag* variable
  - compute correlations between amino acids situated in the sequence at at most *lag* positions apart

→ combined representation



## Evaluation Methodology

- $k$ -fold cross-validation - same number of folds as the related work on that data set

## Evaluation metrics:

- Accuracy
- Precision
- Recall
- F1-score
- Specificity
- Area under the ROC curve

# Data sets

- 4 public data sets: one human data set (HPRD) and three multi-species data sets
- Multi-species data sets: obtained by merging three data sets (Caenorhabditis elegans, Escherichia coli and Drosophila melanogaster) - all interactions, proteins filtered using 25% and 1% similarity thresholds

Data set	Number of positive interactions	Number of negative interactions
HPRD	36,630	36,480
Multi-species	32,959	32,959
Multi-species <0.25	19,458	15,827
Multi-species <0.01	10,747	8,065

Table: Data sets used in the experiments.

# Results

Data set	Arch.	Accuracy	$F_1$ -score	Precision	Recall	Specificity	AUC
HPRD	1	$0.977 \pm 0.0006$	$0.977 \pm 0.0007$	$0.986 \pm 0.0009$	$0.968 \pm 0.001$	$0.986 \pm 0.0009$	$0.977 \pm 0.0006$
	2	<b><math>0.979 \pm 0.0007</math></b>	<b><math>0.979 \pm 0.0007</math></b>	$0.973 \pm 0.0015$	<b><math>0.985 \pm 0.009</math></b>	$0.973 \pm 0.0015$	<b><math>0.979 \pm 0.0007</math></b>
	3	$0.96 \pm 0.0014$	$0.959 \pm 0.0015$	<b><math>0.992 \pm 0.006</math></b>	$0.928 \pm 0.0024$	<b><math>0.992 \pm 0.006</math></b>	$0.960 \pm 0.0014$
Multi-species	1	$0.97 \pm 0.0007$	$0.969 \pm 0.0006$	$0.995 \pm 0.0007$	$0.944 \pm 0.0015$	$0.995 \pm 0.0006$	$0.97 \pm 0.0005$
	2	$0.969 \pm 0.0008$	$0.97 \pm 0.0009$	$0.965 \pm 0.0028$	<b><math>0.974 \pm 0.002</math></b>	$0.964 \pm 0.0025$	$0.97 \pm 0.008$
	3	<b><math>0.982 \pm 0.0008</math></b>	<b><math>0.982 \pm 0.0008</math></b>	<b><math>1 \pm 0</math></b>	$0.964 \pm 0.0016$	<b><math>1 \pm 0</math></b>	<b><math>0.982 \pm 0.008</math></b>
Multi-species <0.25	1	$0.973 \pm 0.0011$	$0.975 \pm 0.0009$	$0.995 \pm 0.0011$	$0.956 \pm 0.0017$	$0.995 \pm 0.0012$	$0.975 \pm 0.001$
	2	$0.976 \pm 0.0007$	$0.978 \pm 0.0008$	$0.974 \pm 0.0011$	<b><math>0.983 \pm 0.0008</math></b>	$0.968 \pm 0.0013$	$0.975 \pm 0.0008$
	3	<b><math>0.983 \pm 0.0015</math></b>	<b><math>0.984 \pm 0.0014</math></b>	<b><math>1 \pm 0</math></b>	$0.969 \pm 0.0027$	<b><math>1 \pm 0</math></b>	<b><math>0.985 \pm 0.0013</math></b>
Multi-species <0.01	1	$0.972 \pm 0.0023$	$0.975 \pm 0.0019$	$0.993 \pm 0.001$	$0.958 \pm 0.0035$	$0.991 \pm 0.0015$	$0.975 \pm 0.002$
	2	$0.978 \pm 0.0015$	$0.981 \pm 0.0013$	$0.975 \pm 0.0024$	<b><math>0.987 \pm 0.0027</math></b>	$0.966 \pm 0.0031$	$0.976 \pm 0.0015$
	3	<b><math>0.981 \pm 0.0016</math></b>	<b><math>0.983 \pm 0.0014</math></b>	<b><math>1 \pm 0</math></b>	$0.966 \pm 0.0027$	<b><math>1 \pm 0</math></b>	<b><math>0.983 \pm 0.0014</math></b>

**Table:** Experimental results. 95% CIs are used for the results. 1 - denotes the Joint-Joint architecture, 2 - the Siamese-Joint architecture, 3 - the Siamese-Siamese architecture

- On each data set one of the siamese architectures provides the best results: Siamese-Joint architecture on HPRD and the Siamese-Siamese for the multi-species data sets

# Results. Comparison with related work

Method	Accuracy	F1
<i>AutoPPI</i>	<b>0.979 ± 0.0007</b>	<b>0.979 ± 0.0007</b>
SAE (Sun et al., 2017)	0.9719	-
PIPR (Chen et al., 2019)	0.9811	0.9803
LDA-RF (Pan et al., 2010)	0.979 ± 0.005	-
CT-SVM (Shen et al., 2007) reported by Sun et al., 2017	0.83	-
AC-SVM (Guo et al., 2010) reported by Sun et al., 2017	0.9037	-
Parallel SVM (You et al., 2014) reported by Sun et al., 2017	0.9200–0.9740	-
ELM (You et al., 2014) reported by Sun et al., 2017	0.8480	0.8477
CS-SVM (Zhang et al., 2011)	0.941	0.937
SVM (Nanni et al., 2013)	0.942	-
DNN (Gui et al., 2020)	0.9443 ± 0.0036	-
DNN-PPI (Gui et al., 2019)	0.9726 ± 0.0018	-
DNN-CTAC (Wang et al., 2019)	<b>0.9837</b>	-
S-VGAE (Yang et al., 2020)	<b>0.9915 ± 0.0011</b>	<b>0.9915 ± 0.0012</b>

**Table:** Comparison between our method and related work on the HPRD data set.

# Results. Comparison with related work

Data set	Method	Accuracy	F1
Multi-species	<i>AutoPPI</i>	<b>0.9821 ± 0.0008</b>	<b>0.9818 ± 0.0008</b>
	PIPR (Chen et al., 2019)	0.9819	0.9817
Multi-species <0.25	<i>AutoPPI</i>	<b>0.9829 ± 0.0015</b>	<b>0.9842 ± 0.0014</b>
	PIPR (Chen et al., 2019)	0.9791	0.9808
Multi-species <0.01	<i>AutoPPI</i>	<b>0.9808 ± 0.0016</b>	<b>0.9829 ± 0.0014</b>
	PIPR (Chen et al., 2019)	0.9751	0.9780

**Table:** Comparison between our method and related work on the Multi-species data sets.

# Conclusions and future directions of research

- Introduced a procedure for binary classification of protein–protein interactions
- Proposed two new siamese architectures for the autoencoders
- Evaluated our approach on four data sets including proteins from different species
- Our approach surpassed the majority of related work approaches

# Conclusions and future directions of research

**Challenge:** Random sampling: does not take into consideration whether the testing proteins are included in the training set → drop in performance when testing on unseen proteins<sup>14,15</sup>

Future directions:

- Improve generalization
- Improve performance on imbalanced data sets
- Provide interpretability

---

<sup>14</sup>Dunham and Ganapathiraju, *Benchmark Evaluation of Protein-Protein Interaction Prediction Algorithms.*, Molecules, 2022.

<sup>15</sup>Park and Marcotte, *Flaws in evaluation schemes for pair-input computational predictions.*, Nature methods, 2012.

Thank you!