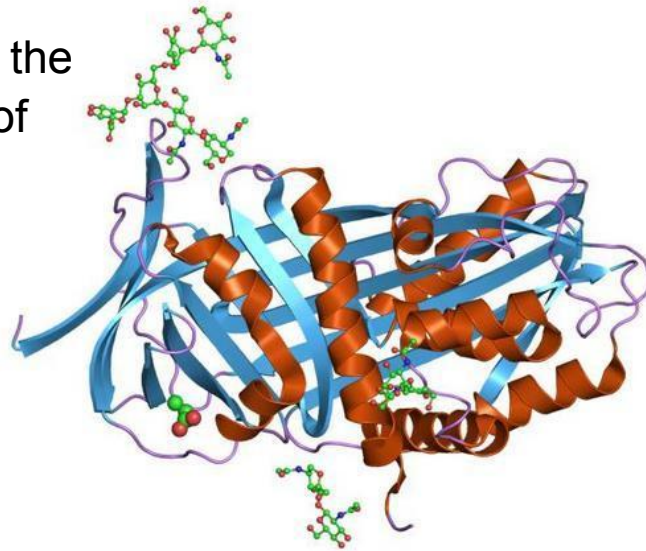# Protein Similarity

**Cecilia Huang, Evelina Danielsson, Joel Bäcker**

**Supervisor: Daniel Varela**

# Andre Lab

Combination of computational and experimental methods to understand the structure, interactions and evolution of proteins.
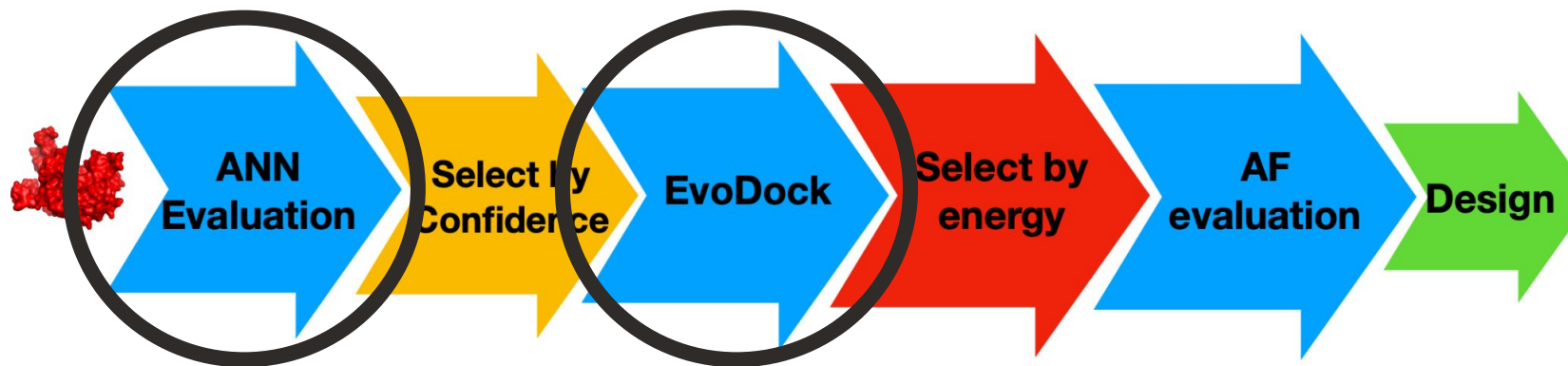
# Protein Similarity in Life Science

The shape of a protein is critical to its function. By finding similarities between the shapes of proteins, functions and relationships can be inferred.
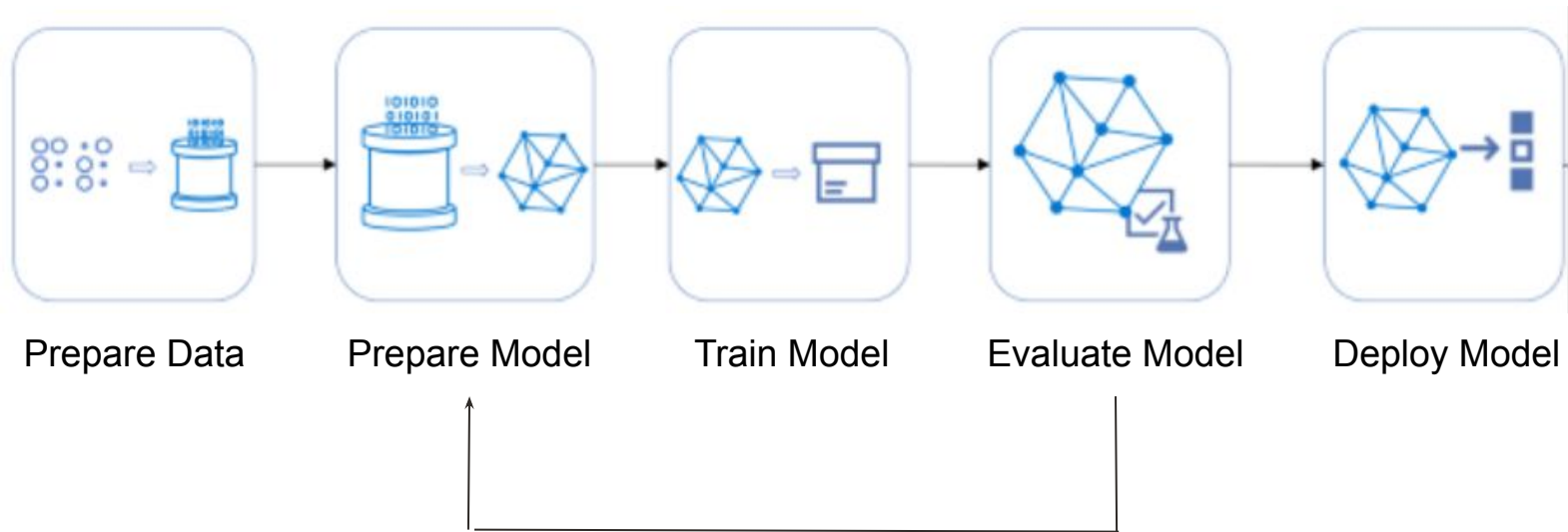
# Predicting protein structure: Challenges

- Large dataset of proteins (500.000)

- Computationally expensive (300 seconds comparing proteins!)

# AndreLab protein design approach

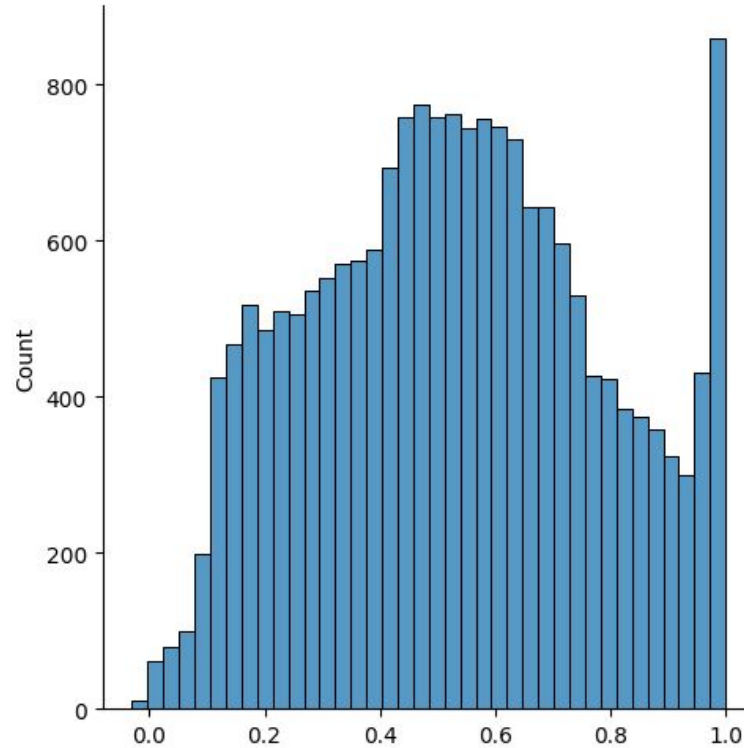# Standard ML pipelining



Prepare Data     Prepare Model     Train Model     Evaluate Model     Deploy Model

# Data Generation

- Random selection

- Zernike-Canterakis shape descriptor

- ZEAL score

# How to create one sample

# ZEAL score distribution

# Features and targets
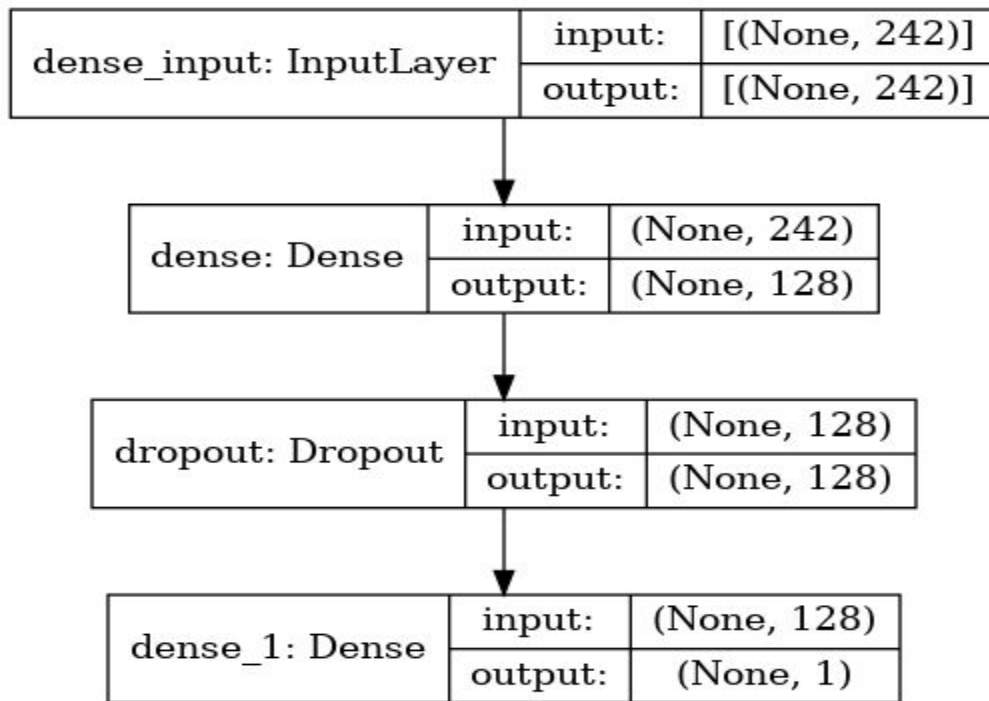
$$X = \begin{bmatrix} zd^1_{T0} & zd^1_{T1} & \cdots & zd^1_{T120} & zd^1_{R0} & \cdots & zd^1_{R120} \\ zd^2_{T0} & zd^2_{T1} & \cdots & zd^2_{T120} & zd^2_{R0} & \cdots & zd^2_{R120} \\ zd^3_{T0} & & \cdots & & \ddots & & \\ \vdots & & & & & \ddots & \\ zd^n_{T0} & & \cdots & & & & zd^n_{R120} \end{bmatrix}$$

$$y = \begin{bmatrix} zeal_1 \\ zeal_2 \\ \vdots \\ zeal_n \end{bmatrix}$$
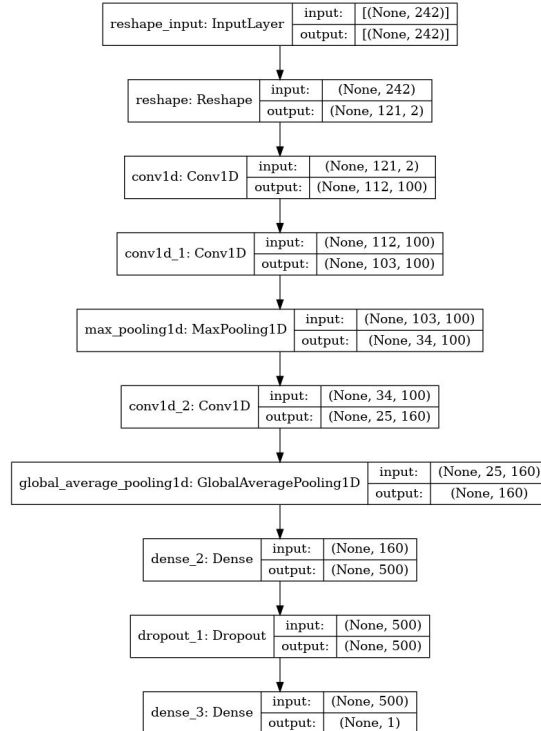
# Model Selection

- Regression
    - Linear Regression
    - Random Forest Regressor
- Neural Networks
    - Feed-Forward Network
    - 1D Convolutional Neural Network
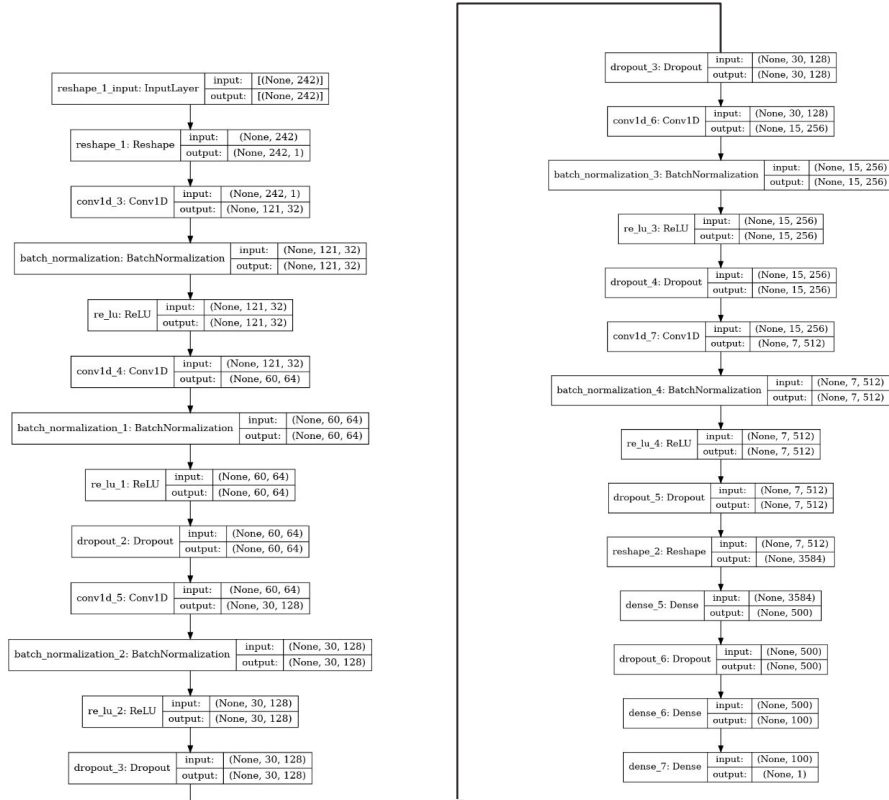    - 1D Fully-Convolutional Network
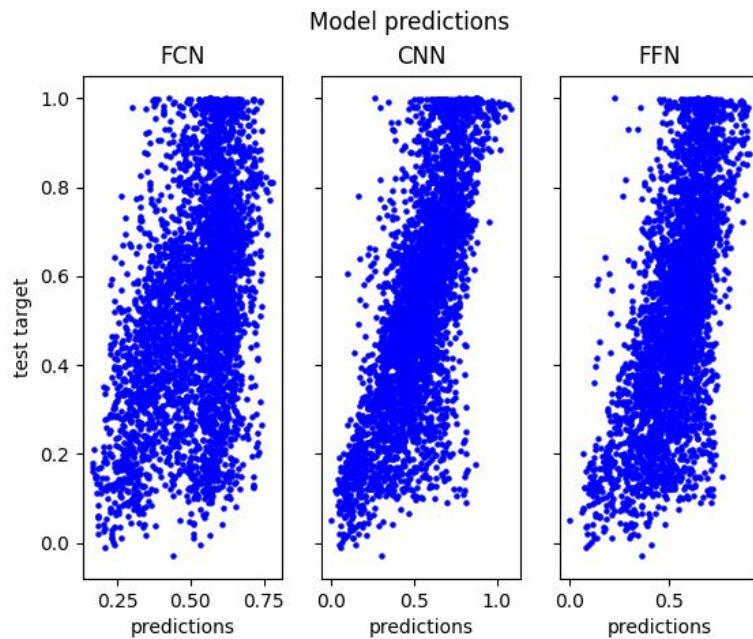
# Feed Forward Network
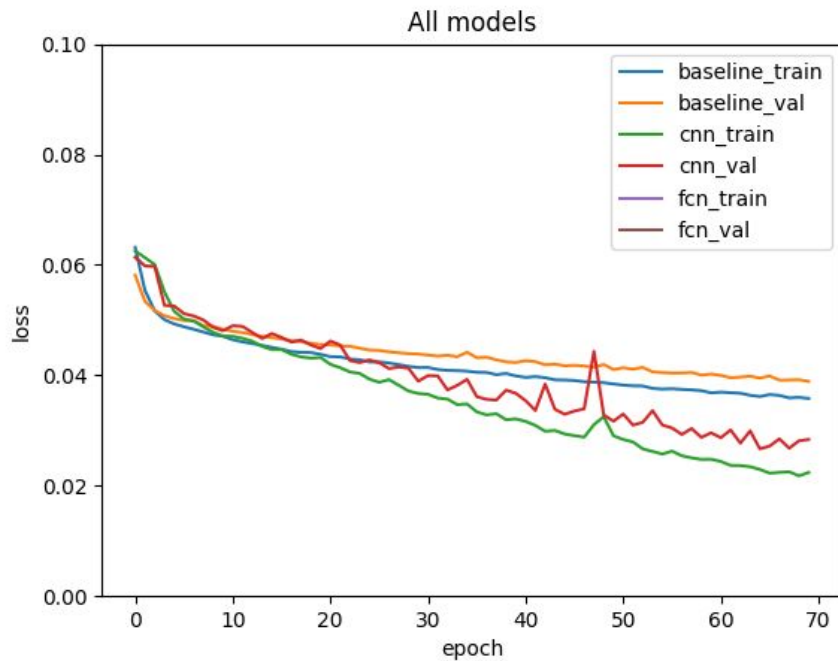
# Convolutional Neural Network



| reshape_input: InputLayer | input: | [(None, 242)] |
|---|---|---|
| | output: | [(None, 242)] |

| reshape: Reshape | input: | (None, 242) |
|---|---|---|
| | output: | (None, 121, 2) |

| conv1d: Conv1D | input: | (None, 121, 2) |
|---|---|---|
| | output: | (None, 112, 100) |

| conv1d_1: Conv1D | input: | (None, 112, 100) |
|---|---|---|
| | output: | (None, 103, 100) |

| max_pooling1d: MaxPooling1D | input: | (None, 103, 100) |
|---|---|---|
| | output: | (None, 34, 100) |

| conv1d_2: Conv1D | input: | (None, 34, 100) |
|---|---|---|
| | output: | (None, 25, 160) |

| global_average_pooling1d: GlobalAveragePooling1D | input: | (None, 25, 160) |
|---|---|---|
| | output: | (None, 160) |

| dense_2: Dense | input: | (None, 160) |
|---|---|---|
| | output: | (None, 500) |

| dropout_1: Dropout | input: | (None, 500) |
|---|---|---|
| | output: | (None, 500) |

| dense_3: Dense | input: | (None, 500) |
|---|---|---|
| | output: | (None, 1) |

# Fully Convolutional Network

# Evaluation



Model predictions

# Evaluation

# Conclusions and further work

- Applicating the classical steps of ML pipelines on biology research

- Data generation was very tedious (surprise!)
    - No previous data
    - Data distribution (ZEAL)

- We are working on testing the model on a much larger data set that we will include in our report.

- Models results