

# Bachelorthesis

## A Deep Learning Approach for Predicting Pesticide Degradation Based on Enzyme Classes

Prüfer(in):

Prof. Dr. Thomas Ströder

Fethi Temiz

Verfasser(in):

Tobias Polley

100853

Gutenbergstr. 5

51469 Bergisch Gladbach

BFWC321B

Cyber Security

Eingereicht am:

June 25, 2024

## Sperrvermerk

Diese Arbeit enthält vertrauliche Informationen über die Firma Bayer AG. Die Weitergabe des Inhalts dieser Arbeit (auch in Auszügen) ist untersagt. Es dürfen keinerlei Kopien oder Abschriften - auch nicht in digitaler Form - angefertigt werden. Auch darf diese Arbeit nicht veröffentlicht werden und ist ausschließlich den Prüfern, Mitarbeitern der Verwaltung und Mitgliedern des Prüfungsausschusses sowie auf Nachfrage einer Evaluierungskommission zugänglich zu machen. Personen, die Einsicht in diese Arbeit erhalten, verpflichten sich, über die Inhalte dieser Arbeit und all ihren Anhängen keine Informationen, die die Firma Bayer AG betreffen, gegenüber Dritten preiszugeben. Ausnahmen bedürfen der schriftlichen Genehmigung der Firma Bayer AG und des Verfassers.

Die Arbeit oder Teile davon dürfen von der FHDW einer Plagiatsprüfung durch einen Plagiatsoftware-Anbieter unterzogen werden. Der Sperrvermerk ist somit im Fall einer Plagiatsprüfung nicht wirksam.

Contents

Sperrvermerk II

List of Figures V

List of Tables VI

Listing Directory VII

1 Introduction 1

1.1 Motivation . . . . . 1

1.2 Structure of the Thesis . . . . . 2

2 Literature Review 4

2.1 Enzymatic Mechanisms Involved in Pesticide Breakdown . . . . . 4

2.2 Deep Learning Techniques in Environmental Science . . . . . 5

3 Theoretical Background 8

3.1 Principles of Enzymology . . . . . 8

3.1.1 Enzyme Classification and Function . . . . . 8

3.1.2 Role of Enzymes in Biodegradation . . . . . 10

3.2 Fundamentals of Ligand Binding Site Prediction . . . . . 12

3.3 Introduction to Recurrent Neural Networks . . . . . 14

3.4 Evaluation of Deep Learning Models . . . . . 17

4 Methodology 20

4.1 Data Collection . . . . . 20

4.2 Data Preprocessing . . . . . 21

4.3 Feature Engineering . . . . . 24

4.4 Model Development . . . . . 28

5 Results 31

5.1 Model Performance . . . . . 31

5.2 Comparative Analysis with Existing Models . . . . . 31

5.3	Interpretation of Model Predictions . . . . .	31
5.1	Implications of Findings . . . . .	32
5.2	Strenths and Limitations . . . . .	32
5.1	Summary of Findings . . . . .	33
5.2	Contributions to the Field . . . . .	33
5.3	Final Remarks and Future Work . . . . .	33
	<b>Appendix</b>	<b>34</b>
	<b>List of References</b>	<b>36</b>
	<b>Statement of independent work</b>	<b>40</b>

## List of Figures

Figure 1: Macro F1 score for different models and EnzymeNet . . . . .	7
Figure 2: Organisation of enzyme structure and lysozyme example. . . . .	9
Figure 3: Lock-and-key model that explains the selectivity of enzymes . . . . .	11
Figure 4: Calculation of feature vectors for neighboring atoms in P2Rank. . . . .	14
Figure 5: A diagram for a one-unit recurrent neural network (RNN). . . . .	16
Figure 6: Cross-validation workflow for evaluating deep learning models. . . . .	18
Figure 7: Source: [17] . . . . .	25
Figure 8: Source: [17] . . . . .	26

## List of Tables

Table 1: Distribution of EC classes in the dataset . . . . .	24
--	----

## Listing Directory

1	Python script for data retrieval and preprocessing from Uniprot . . .	22
2	Python script for downloading the pdb structure . . . . .	23
3	Command Line for running p2rank on a given directory . . . . .	24
4	Python code for creating the model architecture . . . . .	28

# 1 Introduction

## 1.1 Motivation

In recent years, the prediction of pesticide degradation has gained a significant attention due to the environmental and health impacts of pesticides. Traditional methods for determining the degradation of enzymes and the underlying EC-class are labor-intensive and time-consuming. Consequently, there is a growing need for computational methods that can efficiently and accurately predict the degradation behavior of pesticides. One promising approach involves leveraging the capabilities of deep learning to predict enzyme classes responsible for pesticide degradation based on their interaction with specific enzyme binding sites. By combining the prediction of the active binding sites of enzymes with their corresponding protein sequences, it is possible to develop a model that can accurately predict the enzyme class.

The use of advanced computational methodologies, such as Deep Learning for enzyme classification, would enormously increase the development of environmentally friendly and safe agricultural products at Bayer Crop Science. By accurately predicting the functions of enzymes in pesticide degradations, would allow the development of new products that are more sustainable and environmentally friendly. This would also reduce the time and cost of testing existing and new products, as well as the risk of developing products that are harmful to the environment.

Despite advancements in bioinformatics and computational biology, predicting enzyme classes is still fraught with uncertainties. Traditional methods rely heavily on experimental data, which can be resource-intensive and time-consuming. Moreover, the vast diversity of enzyme functions and their complex interactions with various substrates add layers of difficulty to accurate predictions. Thus, there is a pressing need for computational tools that leverage modern machine learning techniques to enhance the prediction accuracy of enzyme-related models. Moreover, there are several models available for predicting enzyme classes based on sequence data, but there is still room for improvement in terms of performance. Therefore, this study



aims to address this gap by developing a deep learning model that can predict enzyme classes responsible for pesticide degradation based on their interaction with specific enzyme binding sites. By combining the prediction of the active binding sites of enzymes with their corresponding protein sequences, it is possible to develop a model that can accurately predict the enzyme class.

By addressing this research question, the study seeks to contribute to the fields of computational biology and environmental science, providing a tool that can accurately predict enzymatic functions and their behavior in pesticide degradations. In addition, this research aims to outperform existing models in predicting enzyme classes responsible for pesticide degradation, thereby enhancing the accuracy of enzyme classification predictions. The findings of this study could have significant implications for the development of environmentally friendly and sustainable agricultural products, as well as the reduction of harmful pesticides in the environment.

## **1.2 Structure of the Thesis**

This thesis is structured into five chapters, each addressing different aspects of the research and providing a comprehensive overview of the study. The first chapter sets the stage for the entire thesis. It begins by outlining the motivation behind the research, highlighting the environmental concerns related to pesticide use and the need for effective degradation prediction methods. The problem statement section identifies the challenges associated with predicting enzyme-mediated pesticide degradation. The introduction section defines the main objective of the study, which is to develop a Deep Learning model to predict pesticide degradation based on enzyme classes. Finally, this chapter provides an overview of the structure of the thesis.

The literature review chapter delves into existing research and foundational theories relevant to the study. It covers enzymatic mechanisms involved in pesticide breakdown, offering insights into how enzymes facilitate the degradation process.

Additionally, it explores the application of Deep Learning techniques in environmental science, emphasizing their potential to enhance predictive accuracy. The chapter concludes with a discussion of the limitations of current models and the need for more advanced approaches to enzyme classification.

The methodology chapter provides a detailed description of the research design and procedures followed in this study. It begins with the Data Collection, specifying the sources and preprocessing steps to prepare the dataset for analysis. The feature engineering section discusses how relevant features were extracted from the data to calculate accurate predictions. The chapter then explains the model development process, including the architecture of the Deep Learning model and the final training process.

The results chapter presents the outcomes of the research. It begins with an evaluation of the model's performance, highlighting key metrics and the effectiveness of the model in predicting pesticide degradation. A comparative analysis with existing models is included to demonstrate the improvements and advantages of the developed model. The chapter also interprets the model predictions, offering insights into the practical implications of the findings and how they can be applied in real-world scenarios.

The discussion chapter summarizes the key findings of the research, reflecting on the significance and impact of the results. It discusses the strengths and limitations of the study, acknowledging areas where the model performed well and identifying potential areas for improvement. The chapter concludes with an overview of the contributions to the field, highlighting the novelty and practical applications of the research. Additionally, it provides recommendations for future work, suggesting directions for further research to build on the findings of this study.

## 2 Literature Review

### 2.1 Enzymatic Mechanisms Involved in Pesticide Breakdown

The degradation of pesticides in the environment is a complex process and occurs by various mechanisms, but mainly through microbial enzymatic activities. The enzymes catalyze reactions in which toxic pesticide compounds are transformed into less harmful compounds facilitating their removal from the environment. In this part, the most critical enzymatic mechanisms applied in pesticide degradation are hydrolytic, oxidative, and reductive enzymes.

Microbial enzymes play critical roles in the biodegradation of soil contaminants, one of which is pesticides. They can further be classified based on the reactions they catalyze:

**Hydrolytic Enzymes:** Hydrolytic enzymes are divided into two groups, esterases and amidases, which hydrolyze ester and amide bonds in pesticide molecules. Their hydrolysis changes them into compounds of much smaller size and more water solubility, which can easily be eliminated through further degradation. For example, microbial esterases act upon organophosphate insecticides in their rapid degradation to hydrolyze them.

**Oxidative Enzymes:** This group of oxidative enzymes includes cytochrome P450 monooxygenases, which introduce oxygen atoms into molecules of pesticides, thus increasing their solubility and reactivity. Commonly, this oxidation makes the pesticides less hazardous, or other enzymes can further degrade the intermediates formed. In particular, the cytochrome P450 enzymes are very versatile, able to metabolize a vast range of xenobiotics, including pesticides.

For example, reductive enzymes catalyze the reduction of pesticides, in most cases by donating electrons and hydrogen atoms on the molecules. Such a reduction may well break complex structures that facilitate the conversion of pesticides

into simpler forms that are much less toxic. For instance, reductive dehalogenases are believed to be important in breaking down halogenated organic compounds.

Degradation of pesticides in contaminated soils by the use of microbial enzymes in bioremediation strategies significantly improves the degradation of the pesticides. The approach is based on the natural capability of microbes to detoxify pollutants via enzymatic reactions. A review found that the efficiency of microbial enzymes in the degradation of pesticides contaminated in soil was well established. [1] Further study focused on the developments and applications of microbial enzymes for enhancing the process of biodegradation. This review will focus on the critical role of enzymes in pesticide-degradation pathways and discuss the potential for engineered enzymes to increase bioremediation efficiency. [2]

Understanding the enzymatic mechanisms is crucial for predicting the enzymatic classes responsible for pesticide degradation. By analyzing the interactions between enzymes and pesticides, it is possible to identify the specific enzyme classes involved in the degradation process. This knowledge can inform the development of more accurate predictive models for pesticide degradation, enabling better risk assessments and environmental management strategies.

## 2.2 Deep Learning Techniques in Environmental Science

Deep Learning has been an essential tool in environmental science, enabling advanced prediction and understanding complex biochemical processes. There are several Deep Learning architectures such as the protein-transformer ESM model, which has made a significant impact on predicting biological properties from sequence data. [3]

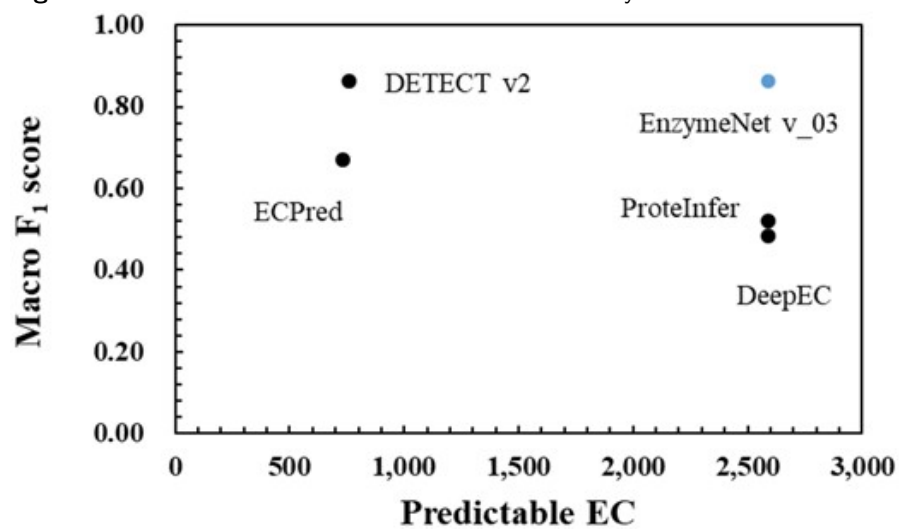
In the context of pesticide degradation and enzyme classification, such models can analyze large quantities of available biochemical data to make predictions about enzyme interactions and functions. Several deep learning architectures have been applied in enzyme classification and prediction tasks, from which valuable insights into the mechanism of pesticide degradation can be obtained.

For instance, the DEEPRe model applies deep learning to predict enzyme commission (EC) numbers based on raw sequence data. Such models apply convolutional and sequential feature extraction techniques, leading to significant improvements in prediction accuracy over methods in current use. In this respect, such models may play a key role in predicting the pesticide biodegradation pathways and help to make environmental risk assessment more precise and fast. [4]

The DeEPn model is one of the examples when EC classification has been done using a deep neural network for enzymes being classified into their functional classes, including all seven EC classes. This model has shown high precision and accuracy and, hence could become an essential tool for environmental scientists interested in understanding and predicting enzyme-mediated degradation processes. The proper classification of enzymes through DeEPn can help predict potential candidates for bioremediation, among other applications related to the environment. [5]

Despite the advances made by these models, there is still a need for new approaches to further improve the accuracy of sequence based predictions. Traditional models often rely on pre-defined features and limited datasets, which can restrict their performance and generalizability. In addition to this, the existing methods only focus on the prediction to the 3rd level of the EC classification, which may not provide sufficient detail for predicting pesticide degradations. For example the accuracy of EnzymeNet, a residual neural network model, across all the sub-subclasses is 0.398. In addition to that there is no score for the 4th level. The following picture shows the macro F1 score for different models and EnzymeNet, which is the best one, but still not good enough. Therefore, there is a need for more advanced deep learning models that can predict enzyme classes with higher accuracy and resolution, enabling more precise predictions of pesticide degradation pathways. [6]

By contrast, the proposed approach leverages the deep learning tool p2rank to analyze the interactive parts of enzymes, focusing on the ligand-binding sites and the specific amino acids involved. This method can potentially provide a more detailed and accurate prediction of enzyme classes responsible for pesticide degradation,

**Figure 1:** Macro F1 score for different models and EnzymeNet

**Source:** Watanabe et al. (2023)

enhancing our understanding of the biodegradation pathways and mechanisms involved. [7]

## 3 Theoretical Background

### 3.1 Principles of Enzymology

Enzymology is the study of enzymes, which are biological catalysts that accelerate biochemical reactions in living organisms. These macromolecules are essential for various cellular processes, including metabolism, DNA replication, and signal transduction. The understanding of enzyme structure, function, and kinetics is crucial for developing applications in biotechnology, medicine, and environmental science. [8]

#### 3.1.1 Enzyme Classification and Function

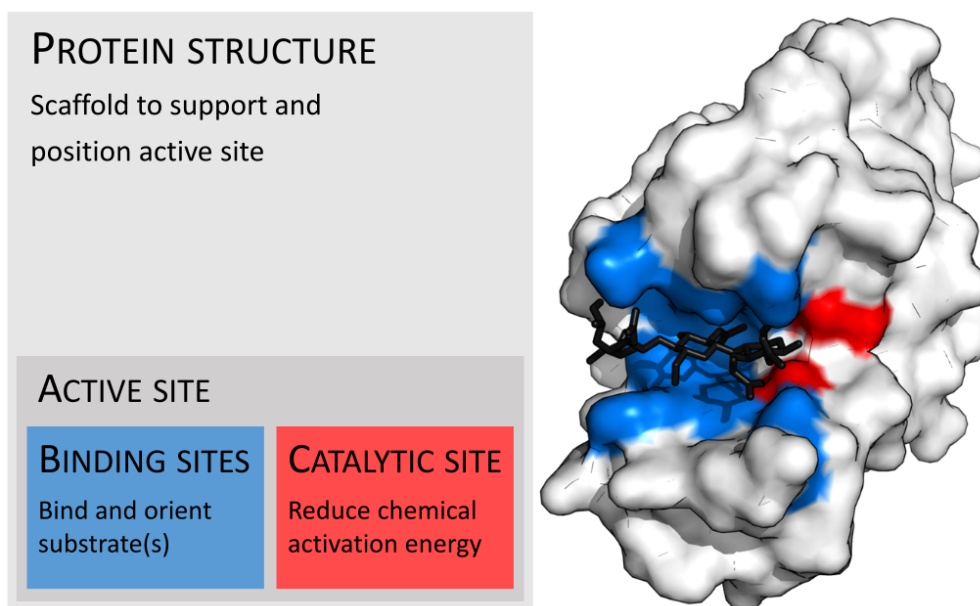
Enzymes are classified based on the reactions they catalyze, following a system established by the Enzyme Commission (EC). This classification system groups enzymes into six main classes, each with specific types of reactions they facilitate:

1. **Oxidoreductases:** These enzymes catalyze oxidation-reduction reactions, where the transfer of electrons occurs between molecules. Examples include dehydrogenases and oxidases.
2. **Transferases:** These enzymes transfer functional groups from one molecule to another. Examples include kinases, which transfer phosphate groups.
3. **Hydrolases:** These enzymes catalyze the hydrolysis of various bonds, including ester, glycosidic, peptide, and others. Examples include proteases and lipases.
4. **Lyases:** These enzymes add or remove groups to form double bonds, without hydrolysis or oxidation. Examples include decarboxylases and dehydratases.
5. **Isomerases:** These enzymes catalyze the rearrangement of atoms within a molecule, leading to isomerization. Examples include racemases and epimerases.
6. **Ligases:** These enzymes catalyze the joining of two molecules with the simultaneous hydrolysis of a diphosphate bond in ATP or a similar triphosphate. Examples include synthetases and carboxylases.

For example, the enzyme tripeptide aminopeptidase has the EC number "3.4.11.4", where the first digit (3) represents the class (Hydrolases in this case), the second digit (4) represents the subclass (hydrolases that act on peptide bonds), the third digit (11) represents the sub-subclass (Hydrolases that cleave off the amino-terminal amino acid from a polypeptide), and the fourth digit (4) represents the serial number of the enzyme within the sub-subclass (Hydrolases that cleave off the amino-terminal end from a tripeptide). This systematic classification allows researchers to identify enzymes based on their catalytic activities and biochemical properties.

Enzymes are not only classified based on their catalytic activities but also based on their biological functions. The three-dimensional (3D) structure of enzymes is fundamental to their function. Enzymes are composed of one or more polypeptide chains that fold into specific shapes to form the active site. The active site is where substrate molecules bind and undergo a chemical reaction. The enzyme structure serves as a scaffold to support and correctly position the active site for optimal catalytic activity.

**Figure 2:** Organisation of enzyme structure and lysozyme example.



**Source:** Thomas Shafee, CC BY 4.0 via Wikimedia Commons



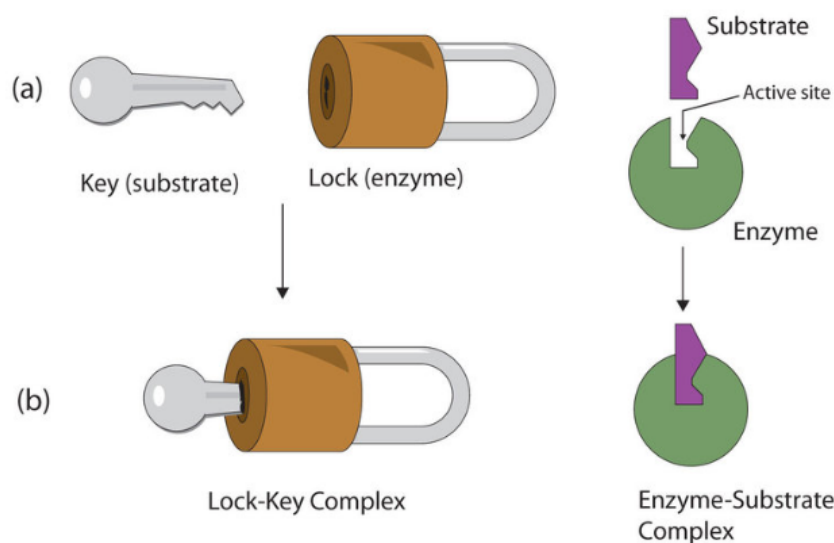
- **Protein Structure:** The overall structure of the enzyme provides the framework that supports and positions the active site. This structure is critical for the enzyme's stability and functionality. The enzyme's polypeptide chains fold into a unique 3D shape, creating a specific environment for the active site.
- **Active Site:** The active site includes two critical regions: binding sites and the catalytic site. The binding sites (highlighted in blue) are regions where substrates bind to the enzyme. These sites ensure that the substrates are properly oriented for the reaction. The catalytic site (highlighted in red) is the region where the chemical reaction occurs. The catalytic site often contains amino acids with specific functional groups that participate directly in the reaction, reducing the activation energy required for the reaction to proceed.

The Key-Lock Principle, first proposed by Emil Fischer in 1894, is a model for understanding the specificity of enzyme-substrate interactions. According to this principle, the enzyme (lock) has a specific active site shape that only fits a particular substrate (key). This model emphasizes the specificity of enzyme-substrate interactions and how enzymes are highly selective for their substrates. This principle is fundamental to understanding enzyme function and the mechanisms of catalysis. The specificity of these interactions is crucial for predicting enzyme activities because it determines the substrates that can bind to the enzyme and undergo catalysis.

The precise arrangement of amino acids in the active site allows enzymes to be highly specific for their substrates, facilitating efficient catalysis. This specificity is a key feature that enables enzymes to perform their roles in various biochemical pathways with high precision. Understanding the structure-function relationship of enzymes is essential for predicting their activities.

### 3.1.2 Role of Enzymes in Biodegradation

Enzymes play a crucial role in the biodegradation of pollutants, including pesticides. The process involves the breakdown of complex organic molecules into simpler, less

**Figure 3:** Lock-and-key model that explains the selectivity of enzymes

**Source:** [poshyvailo-strubeModellingSimulationsEnzymecatalyzed2015]

toxic forms. This degradation is essential for reducing environmental pollution and mitigating the adverse effects of hazardous chemicals.

**Hydrolytic Enzymes:** Hydrolytic enzymes, such as esterases and amidases, catalyze the cleavage of ester and amide bonds in pesticide molecules. This hydrolysis results in the formation of smaller, more water-soluble compounds that are easier to further degrade and eliminate. For example, microbial esterases can hydrolyze organophosphate insecticides, significantly accelerating their breakdown. [munneckeEnzymaticHydrolysisOrganophosphate1976]

**Oxidative Enzymes:** Oxidative enzymes, such as cytochrome P450 monooxygenases, introduce oxygen atoms into the pesticide molecules, increasing their solubility and reactivity. This oxidation process often converts the pesticides into less harmful substances or intermediates that can be further degraded by other enzymes. The cytochrome P450 enzymes are particularly versatile, capable of metabolizing a wide range of xenobiotics, including pesticides. [9]

**Reductive Enzymes:** Reductive enzymes, including reductases, catalyze the reduc-

tion of pesticides, often by adding electrons and hydrogen atoms to the molecules. This reduction can break down complex structures and facilitate the conversion of pesticides into simpler, less toxic forms. Reductive dehalogenases, for instance, play a significant role in the degradation of halogenated organic compounds.

The integration of enzymatic biodegradation with deep learning models can enhance the prediction and analysis of these processes. By using deep learning to analyze enzyme-substrate interactions and their corresponding (EC) classification, we can develop more accurate and efficient bioremediation strategies.

### 3.2 Fundamentals of Ligand Binding Site Prediction

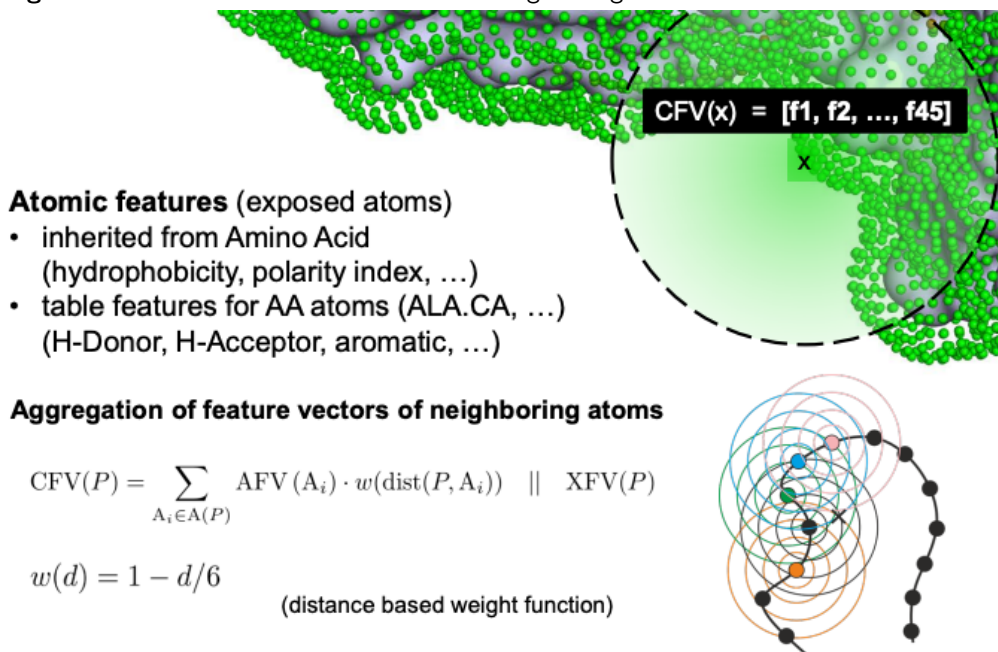
As mentioned earlier, enzymes interact with substrates at specific binding sites, where the catalytic reactions occur. Predicting these ligand-binding sites is crucial for understanding the enzyme function. Several computational methods have been developed to predict ligand-binding sites from protein structures, including geometric, physicochemical, and machine learning-based approaches.

One approach is P2Rank, a machine learning-based tool designed for the rapid and accurate prediction of ligand binding sites from protein structures. It employs a combination of geometric and physicochemical descriptors to analyze protein structures and predict the locations of potential binding sites. P2Rank uses a random forest algorithm, an ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes (classification) or mean prediction (regression) of the individual trees.

The tool focuses on the interactive parts of enzymes, particularly the ligand-binding sites and the specific amino acids involved. This detailed analysis allows for accurate predictions of enzyme classes and their associated degradation pathways. P2Rank's ability to quickly and accurately predict binding sites makes it a valuable tool for drug discovery and environmental bioremediation applications.

P2Rank leverages local chemical neighborhood features near the protein surface to infer potential binding sites for ligands. Here is an overview of the key steps involved in the P2Rank prediction process:

1. **Generation of Connolly Points:** Connolly Points are regularly spaced points generated on the protein’s Connolly surface, representing the solvent-accessible surface area of the protein. These points are generated using a numerical algorithm that ensures even spacing, typically with a solvent radius of 1.6 Å.
2. **Calculation of Feature Descriptors:** Atomic Feature Vectors (AFVs) are calculated for each solvent-exposed heavy atom in the protein, describing various physico-chemical properties such as hydrophobicity, aromaticity, and more. These properties are projected onto the Connolly points using a distance-weighted approach, creating Connolly Feature Vectors (CFVs) for each point. The image shows Connolly Points (green dots) on the protein’s surface, where each point is associated with a Connolly Feature Vector (CFV).
  - a) Atomic Features: Features are inherited from the amino acid, including properties like hydrophobicity and polarity index. Additional features for AA atoms include H-Donor, H-Acceptor, and aromaticity.
  - b) Aggregation of Feature Vectors: The CFV for each Connolly point is calculated by aggregating the AFVs of neighboring atoms using a distance-based weight function  $w(d) = 1 - d/6$ .
3. **Ligandability Prediction:** A Random Forest classifier is used to predict the ligandability score for each Connolly point, indicating the likelihood that a point is part of a ligand-binding site.
4. **Clustering:** Connolly points with high ligandability scores are clustered using a single-linkage clustering method, representing potential binding pockets on the protein surface.

**Figure 4:** Calculation of feature vectors for neighboring atoms in P2Rank.

**Source:** Radoslav Krivák and David Hoksza

5. **Ranking:** Each predicted pocket is assigned a score based on the cumulative ligandability scores of its constituent points, helping prioritize the most likely binding sites for further analysis or docking studies.

P2Rank's approach can significantly enhance the accuracy of predicting enzyme-mediated degradation of pesticides by providing detailed insights into the binding interactions at the molecular level. This integration of deep learning and enzyme analysis forms a robust framework for developing bioremediation strategies and understanding the environmental fate of various pollutants. [7]

### 3.3 Introduction to Recurrent Neural Networks

Recurrent Neural Networks (RNNs) are a class of artificial neural networks designed to recognize patterns in sequences of data such as text, genomes, handwriting, and spoken words. Unlike traditional feedforward neural networks, RNNs have

connections that form directed cycles, allowing information to persist. This makes them particularly powerful for tasks that involve sequential data, where the order of the data points matters.

Recurrent Neural Networks (RNNs) are designed to process sequences of data by maintaining a memory of previous inputs. This memory allows RNNs to make use of information from earlier in the sequence to influence the current processing step, which is essential for understanding context in sequential data. The fundamental difference between RNNs and traditional neural networks is the presence of loops in the network that enable the persistence of information across time steps.

Recurrent Neural Networks (RNNs) are designed to process sequences of data by maintaining a memory of previous inputs. This memory allows RNNs to make use of information from earlier in the sequence to influence the current processing step, which is essential for understanding context in sequential data. The fundamental difference between RNNs and traditional neural networks is the presence of loops in the network that enable the persistence of information across time steps.

The basic structure of an RNN includes an input layer, a hidden layer with recurrent connections, and an output layer. At each time step, the hidden layer receives the input data and its own previous state, allowing it to retain and process information from previous steps in the sequence.

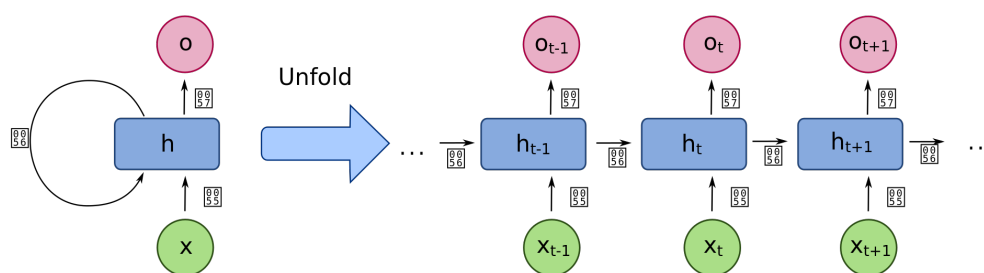
One of the key advancements in RNNs is the development of Long Short-Term Memory (LSTM) networks and Gated Recurrent Units (GRUs), which are designed to overcome the limitations of traditional RNNs, such as the vanishing gradient problem. These architectures use gating mechanisms to control the flow of information, making it easier to capture long-term dependencies in data.

In the context of bioinformatics, RNNs, particularly LSTMs and GRUs, are extensively used for sequence analysis tasks such as protein secondary structure prediction, gene expression analysis, and more. They are effective because they can handle the sequential nature of biological data and capture dependencies that span

over long sequences. LSTM networks are a type of RNN that can learn long-term dependencies. They incorporate memory cells that can maintain their state over long periods. LSTMs have three main gates (input gate, forget gate, and output gate) that regulate the flow of information into and out of the memory cell, thus enabling the network to remember important information for longer durations. [10]

The following image illustrates the basic structure of an RNN:

**Figure 5:** A diagram for a one-unit recurrent neural network (RNN).



**Source:** fdeloche, CC BY-SA 4.0 via Wikimedia Commons

1. Input Sequence ( $x$ ): The green circles represent the input data at different time steps ( $x_{t-1}, x_t, x_{t+1}$ ).
2. Hidden State ( $h$ ): The blue rectangles represent the hidden state of the network. At each time step, the hidden state ( $h$ ) is updated based on the current input and the previous hidden state ( $h_{t-1}, h_t, h_{t+1}$ ).
3. Output Sequence ( $o$ ): The pink circles represent the output of the network at each time step ( $o_{t-1}, o_t, o_{t+1}$ ).

The recurrent connection (arrow looping back) in the hidden state allows information to persist across time steps, enabling the network to maintain context and capture dependencies in the sequence data.

In this study, RNNs are employed for predicting the enzyme class based on the amino acid sequences of a ligand binding site. The sequential nature of the amino

acid sequences makes RNNs well-suited for this task, as they can capture the dependencies and patterns in the data that are crucial for predicting enzyme classes accurately. Especially for complex and long sequences, RNNs, particularly LSTMs, are effective in learning the underlying structure and relationships in the data.

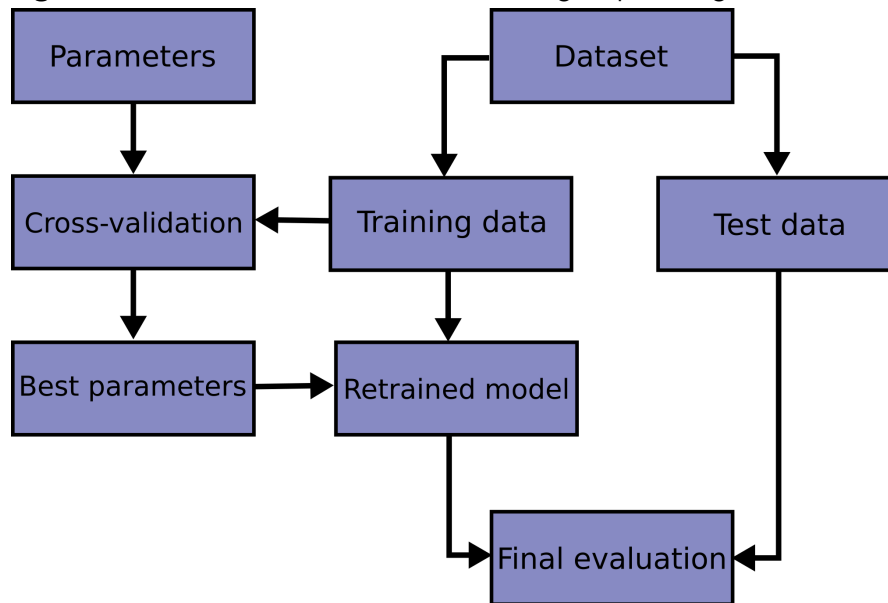
### 3.4 Evaluation of Deep Learning Models

Evaluating deep learning models is essential for understanding their performance and making necessary adjustments to improve accuracy and generalizability. This chapter focuses on model comparison, adjustments, and specific evaluation techniques, including Cross Validation and Grid Search, along with key evaluation metrics.

When developing a deep learning model, it is crucial to compare its performance with other models and make necessary adjustments. This process ensures that the selected model is not only the best fit for the current dataset but also generalizes well to new data. By evaluating multiple models, it is possible to identify which model architecture and hyperparameters yield the best performance. This can be done through adjusting hyperparameters such as learning rate, batch size, and network depth to optimize model performance. Techniques like dropout, L2 regularization, and batch normalization can also be used to prevent overfitting and improve model generalizability. To assess the generalizability of the model, it is essential to evaluate its performance on unseen data, typically through a validation set or cross-validation. Cross-Validation is a widely used technique for evaluating model performance by splitting the dataset into multiple subsets, training the model on different subsets, and testing it on the remaining data. This process helps assess the model's performance across different data partitions and provides a more robust estimate of its generalization capabilities.

The best-known cross-validation method is k-fold. The data is divided into  $k$  equally sized folds. The model is trained  $k$  times, each time using a different fold as the validation set and the remaining  $k - 1$  folds as the training set. Cross-validation



**Figure 6:** Cross-validation workflow for evaluating deep learning models.

**Source:** scikit-learn Documentation

helps in understanding the variability of model performance and reduces the risk of overfitting. Repeating cross-validation multiple times, known as repeated cross-validation, can further enhance reliability. [11] Grid Search is a systematic way of tuning hyperparameters to find the optimal set that maximizes model performance. It involves defining a grid of possible hyperparameter values and exhaustively training and evaluating the model for each combination. For each combination of hyperparameters, the model is trained using the training data. The performance is then evaluated with a cross-validation for each hyperparameter combination. The hyperparameters that yield the best performance are selected as the optimal set. Grid Search is a powerful tool for fine-tuning deep learning models and optimizing their performance. [12]

To effectively evaluate and compare models, various metrics are used: [13] [14] [15]

**Accuracy:** Accuracy measures the proportion of correctly predicted instances out of the total instances. It is a straightforward metric but may not be reliable for

imbalanced datasets.

$$Accuracy = \frac{TotalNumberofPredictions}{NumberofCorrectPredictions} \quad (1)$$

**Precision:** Precision indicates the accuracy of positive predictions, reflecting how many predicted positive instances are actually positive.

$$Precision = \frac{TruePositives}{TruePositives + FalsePositives} \quad (2)$$

**Recall:** Recall measures the model's ability to identify all relevant instances, showing the proportion of actual positives correctly identified.

$$Recall = \frac{TruePositives}{TruePositives + FalseNegatives} \quad (3)$$

**F1 Score:** The F1 Score is the harmonic mean of precision and recall, providing a balanced measure when there is an uneven class distribution.

$$F1Score = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (4)$$

## 4 Methodology

### 4.1 Data Collection

Data collection is the most crucial step in model building since the quality and relevance of the data will directly affect the model's performance. The data for this thesis was collected from UniProt, a very useful resource as it provides protein sequence and functional information. UniProt or Universal Protein Resource is a central protein sequence, and annotation database. It is widely accepted as comprehensive and provides high-quality data, which makes it a must to perform bioinformatics and computational biology. Uniprot pools in much valuable information: experimental findings, different kinds of analyses, and literature information and hence provides rich and reliable sources of further research. Researchers can receive high-quality reviewed entries with 3D structural data and catalytic properties in place, which makes the data reliable and applicable for predicting enzyme functions. [16]

Some of the main features of Uniprot are:

1. **Comprehensive Protein Data:** Uniprot contains a vast collection of protein sequences, functional annotations, and cross-references to other databases, making it a valuable resource for protein research.
2. **Reviewed Entries:** Uniprot contains both reviewed (Swiss-Prot) and unreviewed (TrEMBL) entries. Reviewed entries are manually curated by experts, ensuring high accuracy and reliability.
3. **Functional Annotations:** Each protein entry includes detailed functional annotations, such as catalytic activity, biological processes, and involvement in pathways.
4. **3D Structural Data:** Uniprot links to structural databases like PDB (Protein Data Bank), providing access to 3D structures of proteins, which are crucial for understanding enzyme mechanisms.

5. Cross-references: Extensive cross-references to other databases (e.g., PDB, BRENDA, Reactome) enhance the richness of the data.

For this study, Uniprot was chosen due to its high-quality data, extensive coverage of protein information, and user-friendly interface. The data collection process involved querying Uniprot for enzyme entries with 3D structural data and catalytic activity annotations, extracting relevant information, and preprocessing the data for model development. The data retrieval process involved using the Uniprot REST API to download protein data that matched specific criteria. The criteria included reviewed entries with both 3D structural data and catalytic properties. The Python script below was used to automate the data retrieval and preprocessing steps: [17]

1. API Request: The script constructs a query to the Uniprot REST API to retrieve reviewed protein entries with specified fields and criteria.
2. Data Retrieval: Data is retrieved in compressed format and decompressed using gzip.
3. Data Parsing: The decompressed data is read into a Pandas DataFrame.
4. Data Filtering: The DataFrame is filtered to retain entries with non-null EC numbers and PDB codes.
5. Data Splitting: Entries with multiple PDB codes are split into separate rows for each PDB code.
6. Data Saving: The processed data is saved as a TSV file for further analysis.

## 4.2 Data Preprocessing

Data preprocessing is a crucial step in preparing the dataset for the prediction model. It involves the protein structure download, p2rank workflow, sequence extraction and combination the data into a single dataset. The first step in data preprocessing is to download the 3D protein structures from the Protein Data Bank (PDB) using the PDB ID obtained from Uniprot. The PDB is a repository of experimentally determined protein structures, providing valuable insights into the 3D organization

**Listing 1:** Python script for data retrieval and preprocessing from Uniprot

```

import requests
from tqdm import tqdm
import gzip
from io import BytesIO
import pandas as pd

base_url = "https://rest.uniprot.org/uniprotkb/search?compressed=true&
fields=accession%2Creviewed%2Cid%2Cprotein_name%2Cgene_names%2Corganism_name%
2Cec%2Corganism_id%2Crhea%2Cxref_alphafolddb%2Cxref_pdb%2Cxref_brenda%
2Cxref_biocyc%2Cxref_pathwaycommons%2Cxref_sabio-rk%2Cxref_reactome%
2Cxref_plantreactome%2Cxref_signor%2Cxref_signalink%2Cxref_unipathway&
format=tsv&query=%28*%29+AND+%28reviewed%3Atrue%29+AND+%28proteins_with%
3A1%29+AND+%28proteins_with%3A13%29"

size = 500
offset = 0
all_data = []

response = requests.get(f"{base_url}&size=1")
if response.status_code == 200:
    total_results = int(response.headers.get("x-total-results", 0))
else:
    print(f"Fehler beim Abrufen der Daten: {response.status_code}")
    total_results = 0

with tqdm(total=total_results, desc="Abrufen der Daten", unit=" Eintrag") as pbar:
    while offset < total_results:
        url = f"{base_url}&size={size}&offset={offset}"
        response = requests.get(url)

        if response.status_code == 200:
            with gzip.GzipFile(fileobj=BytesIO(response.content)) as f:
                data = f.read().decode('utf-8')
                if not data.strip():
                    break
                all_data.append(data)
                offset += size
                pbar.update(size)
        else:
            print(f"Fehler beim Abrufen der Daten: {response.status_code}")
            break

combined_data = "\n".join(all_data)

df = pd.read_csv(BytesIO(combined_data.encode('utf-8')), sep='\t')

```

of proteins. The structures are needed for the p2rank workflow, which predicts the interactive site residues in the protein structures.

**Listing 2:** Python script for downloading the pdb structure

```
def download_pdb(pdb_id):
    pdb_url = f"https://files.rcsb.org/download/{pdb_id}.pdb"
    retries = 3
    for attempt in range(retries):
        try:
            response = requests.get(pdb_url, timeout=10)
            if response.status_code == 200:
                with open(f'../data/data_preparation/raw_pdbs/{pdb_id}.pdb', 'w') as file:
                    file.write(response.text)
                return f"Download of {pdb_id} successful."
            else:
                return f"Fehler beim Herunterladen der PDB-Datei {pdb_id}: {response}"
        except requests.exceptions.RequestException as e:
            if attempt < retries - 1:
                continue
            else:
                return f"Error while downloading {pdb_id}: {e}"

with ThreadPoolExecutor(max_workers=10) as executor:
    results = list(tqdm(executor.map(download_pdb, pdb_ids), total=len(pdb_ids), de
```

This script performs the following steps:

1. `requests.get`: Sends an HTTP GET request to the PDB URL to download the protein structure file.
2. `Retry Logic`: Attempts to download the file up to three times in case of failures.
3. `File Writing`: Saves the downloaded PDB file to the specified directory if the download is successful.

The `ThreadPoolExecutor` is used to parallelize the download process and speed up the data retrieval.

The next step in data preprocessing is to run the `p2rank` workflow on the downloaded protein structures to predict the interactive site residues. With the following lines of code, the `p2rank` workflow is executed on the directory containing the PDB files:

**Listing 3:** Command Line for running p2rank on a given directory

```
!./prank.sh predict /Users/tobias.polley/Repositories/DeepZyme/data/data_preparati
```

EC Class	Count
Oxidoreductases	23544
Transferases	59022
Hydrolases	41283
Lyases	2376
Isomerases	4617
Ligases	1323
Translocases	1836

**Table 1:** Distribution of EC classes in the dataset

### 4.3 Feature Engineering

Feature engineering is a critical step in preparing data for prediction models. This process involves transforming raw data into meaningful features that can improve the performance of the model. In this section, the author describes the feature engineering techniques used in this study, focusing on the processing of protein sequences and the calculation of additional features to enhance the predictive power of the Deep Learning model.

To capture meaningful information from protein sequences, this study used several features derived from the sequences, including amino acid composition, molecular weight, isoelectric point, hydrophobicity, and sequence length. These features provide valuable insights into the physicochemical properties of the proteins, enabling the model to learn patterns that correlate with enzyme functions. The ProteinAnalysis class from the Biopython library was used to calculate these features. The following Python code snippet demonstrates the calculation of additional features from the protein sequences:

The first step is to clean the protein sequence shown in chapter 4.2. The sequence itself is used as a feature, and additional features are calculated using the ProteinAnalysis class from Biopython. To convert the cleaned sequences into a format

```

def calculate_features(sequence):
    sequence = clean_sequence(sequence)
    analysis = ProteinAnalysis(sequence)
    amino_acid_composition = list(analysis.get_amino_acids_percent().values)
    molecular_weight = analysis.molecular_weight()
    isoelectric_point = analysis.isoelectric_point()
    hydrophobicity = analysis.gravy()
    sequence_length = len(sequence)
    return amino_acid_composition + [molecular_weight, isoelectric_point, hydrophobicity, sequence_length]

additional_features = df["sequence"].apply(calculate_features)
additional_features = np.array(additional_features.tolist())

```

**Figure 7:** Source: [17]

suitable for the model, a tokenizer is used to encode the sequences into numerical data. In the context of protein sequences, each amino acid is mapped to a unique integer. For example, the sequence "ACDEFGHIKLMNPQRSTVWY" is tokenized into a list of integers.

Tokenization involves converting each amino acid into an integer based on its position in a predefined list of valid amino acids. This process can be mathematically represented as:

token  $(x) = i$  where  $x \in \{A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y\}$  where  $i$  is the index of the amino acid  $x$  in the list.

The tokenized sequence is then passed through an embedding layer that transforms these integers into dense vectors. This embedding process is essential for capturing the contextual meaning of each amino acid within the sequence:  $embedding(i) = v_i$  where  $v_i$  is the embedding vector for the token  $i$ . These embeddings are fed into the RNN, which processes the sequence and updates its hidden states accordingly, allowing the model to capture complex dependencies and interactions between amino acids. The sequences are then padded to ensure they all have the same length, which is necessary for batch processing in Deep Learning models.

Recent advancements have demonstrated that sequence-based models, including language models like ESM-1b, can achieve high accuracy in predicting protein functions and properties. For instance, the study by Hu et al. (2022) highlights



```
tokenizer = Tokenizer()
tokenizer.fit_on_texts(sequences)
encoded_sequences = tokenizer.texts_to_sequences(sequences)

max_sequence_length = max([len(seq) for seq in encoded_sequences])
padded_sequences = pad_sequences(encoded_sequences, maxlen=max_sequence_length,
```

**Figure 8:** Source: [17]

the potential of protein-sequence based models like ESM-1b in predicting protein function from sequences. [18]

In addition to tokenizing the protein sequences, several biochemical features are calculated to provide a comprehensive representation of the proteins. These features include amino acid composition, molecular weight, isoelectric point, hydrophobicity, and sequence length. The Python code for calculating these features is as follows:

1. **Amino Acid Composition:** The amino acid composition represents the relative frequency of each of the 20 standard amino acids in a protein sequence.
  - a) Calculation: It is calculated as the percentage of each amino acid in the sequence.
  - b) Relevance: Different proteins have characteristic amino acid compositions that can provide clues about their function and stability. For example, membrane proteins often have higher hydrophobic amino acid content.
  - c) Example: A protein with a high proportion of hydrophobic amino acids might be involved in membrane-related processes.
2. **Molecular Weight:** Molecular weight is the total mass of all amino acids in the protein sequence.
  - a) Calculation: It is calculated by summing the average atomic masses of the amino acids in the sequence.

- b) **Relevance:** The molecular weight of a protein can influence its physical and chemical properties, such as solubility and interaction with other molecules.
  - c) **Example:** Enzymes with larger molecular weights may have multiple domains or subunits.
3. **Isoelectric Point:** The isoelectric point is the pH at which the protein carries no net electrical charge.
- a) **Calculation:** It is determined by calculating the pH at which the positive and negative charges on the amino acids balance out.
  - b) **Relevance:** The pI affects protein solubility and interaction with other molecules. Proteins are least soluble at their pI and more likely to precipitate.
  - c) **Example:** Proteins with a low pI are often found in acidic environments, such as lysosomal enzymes.
4. **Hydrophobicity (GRAVY Score):** The GRAVY (Grand Average of Hydropathicity) score is a measure of the overall hydrophobic or hydrophilic nature of a protein.
- a) **Calculation:** It is calculated by averaging the hydropathy values of all amino acids in the sequence.
  - b) **Relevance:** Hydrophobicity influences protein folding, stability, and interaction with membranes.
  - c) **Example:** Transmembrane proteins typically have a high GRAVY score due to their hydrophobic transmembrane regions.
5. **Sequence length:** The sequence length is the total number of amino acids in the protein sequence.
- a) **Calculation:** It is simply the count of amino acids in the sequence.

- b) Relevance: The length of a protein can indicate its complexity and the number of functional domains.
- c) Example: Longer proteins may have multiple functional domains or be involved in complex regulatory mechanisms.

These biochemical features provide a multi-dimensional representation of protein sequences, capturing both sequence-specific information and physicochemical properties. This feature-set is essential for analyzing the enzymes and predicting their functions accurately. A study by Gainza et al. (2020) demonstrates the importance of incorporating physicochemical features in protein function prediction models, showing that these features enhance the model's performance. [19]

## 4.4 Model Development

In this section, this study describes the model development process, including data splitting, model architecture, and the rationale behind the chosen methods. The goal is to predict enzyme commission (EC) classes based on protein sequences using a deep learning approach enhanced by additional biochemical features.

To use different models for each hierarchy level, the data is split into four levels of the EC hierarchy. The first level represents the broadest classification, while the fourth level provides the most specific classification. This ensures that models are trained and evaluated on appropriately structured data, allowing for predictions at varying levels of specificity.

The model architecture combines sequence-based features with additional biochemical features to enhance prediction accuracy. The architecture consists of an embedding layer, two LSTM layers, and dense layers that integrate additional features.

**Listing 4:** Python code for creating the model architecture

```
def create_model(num_classes):  
    sequence_input = Input(shape=(max_sequence_length,), dtype='int32')
```

```
embedded_sequences = Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=EMBEDDING_DIM)
x = LSTM(units=32, return_sequences=True)(embedded_sequences)
x = LSTM(units=16)(x)
additional_input = Input(shape=(len(additional_features[0]),))
combined = Concatenate()([x, additional_input])
combined = Dense(units=32, activation='relu')(combined)
output = Dense(units=num_classes, activation='softmax')(combined)
model = Model(inputs=[sequence_input, additional_input], outputs=output)
model.compile(optimizer=Adam(learning_rate=0.001), loss="sparse_categorical_crossentropy")
return model
```

The embedding layer converts amino acid sequences into dense vector representations, capturing semantic similarities between amino acids. This layer allows the model to handle varying sequence lengths and to learn useful representations of amino acids in the context of their sequence. After that Long Short-Term Memory (LSTM) layers are used to capture long-range dependencies in the sequence data, which is crucial for understanding the functional context of amino acids within the sequence. LSTMs are particularly effective in modeling sequential data due to their ability to remember information for long periods and manage the vanishing gradient problem. LSTMs are particularly well-suited for tasks involving sequence data due to their ability to manage long-term dependencies and their robustness against the vanishing gradient problem. Studies have demonstrated the effectiveness of LSTMs in various sequence analysis tasks, including protein function prediction and other bioinformatics applications [20] [21]. Biochemical properties such as molecular weight, isoelectric point, hydrophobicity, and sequence length are included to provide additional context that can enhance the prediction accuracy. These features help the model understand the physical and chemical characteristics of the proteins, which are critical for predicting enzyme functions. The concatenation layer combines the output of the LSTM layers with the additional biochemical features, allowing the model to leverage both sequence-based and property-based information. This integration ensures that the model considers both the sequence context and the biochemical properties of the proteins. Finally the dense layers are used to integrate the combined features and produce the final classification output. These layers apply non-linear transformations to the combined features, enabling the model to learn complex patterns and relationships. This model uses an Adam optimizer with a learning rate of 0.001 and sparse categorical cross-entropy loss

function, which is suitable for multi-class classification tasks. The model is compiled with the specified optimizer, loss function, and evaluation metrics to prepare it for training.

## **5 Results**

### **5.1 Model Performance**

### **5.2 Comparative Analysis with Existing Models**

### **5.3 Interpretation of Model Predictions**

## **Discussion**

### **5.1 Implications of Findings**

### **5.2 Strenths and Limitations**

## **Conclusion**

### **5.1 Summary of Findings**

### **5.2 Contributions to the Field**

### **5.3 Final Remarks and Future Work**



**Appendix**

**Appendix**

Appendix 0.1: Filler . . . . . 35

Filler

## **Appendix 0.1   Filler**

- **Filler**

## List of References

### References

- Singh, Brajesh K. and Allan Walker (May 2006). “Microbial Degradation of Organophosphorus Compounds.” In: *FEMS Microbiology Reviews* 30.3, pp. 428–471. ISSN: 0168-6445. DOI: 10.1111/j.1574-6976.2006.00018.x. (Visited on June 6, 2024).
- Chia, Xing Kai, Tony Hadibarata, Risky Ayu Kristanti, Muhammad Noor Hazwan Jusoh, Inn Shi Tan, and Henry Chee Yew Foo (May 2024). “The Function of Microbial Enzymes in Breaking down Soil Contaminated with Pesticides: A Review.” In: *Bioprocess and Biosystems Engineering* 47.5, pp. 597–620. ISSN: 1615-7605. DOI: 10.1007/s00449-024-02978-6. (Visited on June 6, 2024).
- Rives, Alexander, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus (Apr. 2021). “Biological Structure and Function Emerge from Scaling Unsupervised Learning to 250 Million Protein Sequences.” In: *Proceedings of the National Academy of Sciences* 118.15, e2016239118. DOI: 10.1073/pnas.2016239118. (Visited on June 24, 2024).
- Li, Yu, Sheng Wang, Ramzan Umarov, Bingqing Xie, M. Fan, Lihua Li, and Xin Gao (2017). “DEEPre: Sequence-Based Enzyme EC Number Prediction by Deep Learning.” In: *Bioinformatics* 34, pp. 760–769. DOI: 10.1093/bioinformatics/btx680. (Visited on May 29, 2024).

*DeEPn: A Deep Neural Network Based Tool for Enzyme Functional Annotation - Consensus* (2024).

[https://consensus.app/papers/deepn-network-based-tool-annotation-semwal/c62a6c023f2d5b2f9348b36f6329daae/?utm\\_source=chatgpt](https://consensus.app/papers/deepn-network-based-tool-annotation-semwal/c62a6c023f2d5b2f9348b36f6329daae/?utm_source=chatgpt). (Visited on May 29, 2024).

Watanabe, Naoki, Masaki Yamamoto, Masahiro Murata, Yuki Kuriya, and Michihiro Araki (Jan. 2023). “EnzymeNet: Residual Neural Networks Model for Enzyme Commission Number Prediction.” In: *Bioinformatics Advances* 3.1, vbad173. ISSN: 2635-0041. DOI: 10.1093/bioadv/vbad173. (Visited on June 25, 2024).

Krivák, Radoslav and David Hoksza (Aug. 2018). “P2Rank: Machine Learning Based Tool for Rapid and Accurate Prediction of Ligand Binding Sites from Protein Structure.” In: *Journal of Cheminformatics* 10.1, p. 39. ISSN: 1758-2946. DOI: 10.1186/s13321-018-0285-8. (Visited on May 29, 2024).

Robinson, Peter K. (Nov. 2015). “Enzymes: Principles and Biotechnological Applications.” In: *Essays in Biochemistry* 59, pp. 1–41. ISSN: 0071-1365. DOI: 10.1042/bse0590001. (Visited on June 6, 2024).

Bello, Angelica, Yessica Carreon, and Alejandro Nava-Ocampo (Aug. 2000). “A Theoretical Approach to the Mechanism of Biological Oxidation of Organophosphorus Pesticides.” In: *Toxicology* 149, pp. 63–8. DOI: 10.1016/S0300-483X(00)00222-5.

Hochreiter, Sepp and Jürgen Schmidhuber (Nov. 1997). “Long Short-Term Memory.” In: *Neural Computation* 9.8, pp. 1735–1780. ISSN: 0899-7667. DOI: 10.1162/neco.1997.9.8.1735. (Visited on June 20, 2024).

- Krstajic, D., L. Buturovic, D. Leahy, and Simon Thomas (2014). “Cross-Validation Pitfalls When Selecting and Assessing Regression and Classification Models.” In: *Journal of Cheminformatics* 6. DOI: 10.1186/1758-2946-6-10. (Visited on June 25, 2024).
- Bergstra, James and Yoshua Bengio (n.d.). “Random Search for Hyper-Parameter Optimization.” In: ().
- “Precision and Recall” (Apr. 2024). In: *Wikipedia*. (Visited on May 31, 2024).
- “F-Score” (May 2024). In: *Wikipedia*. (Visited on May 31, 2024).
- “Accuracy and Precision” (Apr. 2024). In: *Wikipedia*. (Visited on May 31, 2024).
- UniProt Consortium (Jan. 2021). “UniProt: The Universal Protein Knowledgebase in 2021.” In: *Nucleic Acids Research* 49.D1, pp. D480–D489. ISSN: 1362-4962. DOI: 10.1093/nar/gkaa1100.
- Polley, Tobias (May 2024). *TobiasPol/DeepZyme*. (Visited on June 5, 2024).
- Hu, Mingyang, Fajie Yuan, Kevin K. Yang, Fusong Ju, Jin Su, Hui Wang, Fei Yang, and Qiuyang Ding (2022). “Exploring Evolution-Based & -Free Protein Language Models as Protein Function Predictors.” In: *ArXiv* abs/2206.06583. DOI: 10.48550/arXiv.2206.06583. (Visited on June 22, 2024).
- Gainza, P., F. Sverrisson, F. Monti, E. Rodolà, D. Boscaini, M. M. Bronstein, and B. E. Correia (Feb. 2020). “Deciphering Interaction Fingerprints from Protein Molecular Surfaces Using Geometric Deep Learning.” In: *Nature Methods* 17.2, pp. 184–192. ISSN: 1548-7105. DOI: 10.1038/s41592-019-0666-6. (Visited on June 22, 2024).
- Liu, Jiale and Xinqi Gong (2019). “Attention Mechanism Enhanced LSTM with Residual Architecture and Its Application for Protein-Protein Interaction

Residue Pairs Prediction.” In: *BMC Bioinformatics* 20. DOI: 10.1186/s12859-019-3199-1. (Visited on June 23, 2024).

Zhang, Yunong (2023). “Encoder-Decoder Models in Sequence-to-Sequence Learning: A Survey of RNN and LSTM Approaches.” In: *Applied and Computational Engineering*. DOI: 10.54254/2755-2721/22/20231220. (Visited on June 23, 2024).

## Statement of independent work

Hiermit erkläre ich, dass ich die vorliegende Bachelorthesis selbständig angefertigt habe. Es wurden nur die in der Arbeit ausdrücklich benannten Quellen und Hilfsmittel benutzt. Wörtlich oder sinngemäß übernommenes Gedankengut habe ich als solches kenntlich gemacht. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Bergisch Gladbach, June 25, 2024

---

Tobias Polley