

# 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:

Wednesday 12<sup>th</sup> June, 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        |
| Code Listings  | VII       |
| <b>1 Introduction</b>  | <b>1</b>  |
| 1.1 Motivation . . . . .   | 1         |
| 1.2 Problem Statement . . . . .                                    | 1         |
| 1.3 Purpose and Research Question . . . . .                        | 2         |
| 1.4 Structure of the Thesis . . . . .                              | 2         |
| <b>2 Literature Review</b>   | <b>4</b>  |
| 2.1 Enzymatic Mechanisms Involved in Pesticide Breakdown . . . . . | 4         |
| 2.2 Deep Learning Techniques in Environmental Science . . . . .    | 5         |
| <b>3 Theoretical Background</b>                                    | <b>7</b>  |
| 3.1 Principles of Enzymology . . . . .                             | 7         |
| 3.1.1 Enzyme Classification and Function . . . . .                 | 7         |
| 3.1.2 Role of Enzymes in Biodegradation . . . . .                  | 9         |
| 3.2 Fundamentals of Machine Learning and Deep Learning . . . . .   | 10        |
| 3.2.1 Introduction to Deep Learning . . . . .                      | 10        |
| 3.2.2 Fundamentals of Ligand Binding Site Prediction . . . . .     | 12        |
| 3.2.3 Sequence-Based Machine Learning . . . . .                    | 13        |
| 3.2.4 Evaluation of Machine Learning Models . . . . .              | 14        |
| <b>4 Methodology</b>   | <b>17</b> |
| 4.1 Data Collection . . . . .                                      | 17        |
| 4.2 Data Preprocessing . . . . .                                   | 19        |
| 4.3 Feature Engineering . . . . .                                  | 19        |
| 4.4 Model Development . . . . .                                    | 19        |

|          |   |           |
|----------|---|-----------|
| 4.4.1    | Model Architecture . . . . .                        | 19        |
| 4.4.2    | Model Training . . . . .                            | 19        |
| 4.4.3    | Model Evaluation . . . . .                          | 19        |
| <b>5</b> | <b>Results</b>                                      | <b>21</b> |
| 5.1      | Model Performance . . . . .                         | 21        |
| 5.2      | Comparative Analysis with Existing Models . . . . . | 21        |
| 5.3      | Interpretation of Model Predictions . . . . .       | 21        |
| 5.1      | Implications of Findings . . . . .                  | 22        |
| 5.2      | Strengths and Limitations . . . . .                 | 22        |
| 5.1      | Summary of Findings . . . . .                       | 23        |
| 5.2      | Contributions to the Field . . . . .                | 23        |
| 5.3      | Final Remarks and Future Work . . . . .             | 23        |
|          | <b>Appendix</b>                                     | <b>24</b> |
|          | <b>List of references</b>                           | <b>26</b> |
|          | <b>Statement of independent work</b>                | <b>27</b> |

## List of Figures

|  |    |
|--|----|
| Abbildung 1: Organisation of enzyme structure and lysozyme example. . .      | 8  |
| Abbildung 2: Diagram of a multi-layer feedforward artificial neural network. | 11 |

## List of Tables

## Code Listings

- 1 Python script for data retrieval and preprocessing from Uniprot . . . 20

# 1 Introduction

## 1.1 Motivation

In recent years, the prediction of pesticide degradation has gained significant attention due to the environmental and health impacts of pesticide residues. Traditional experimental methods for determining the degradation pathways and rates of pesticides 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.

This need is particularly pronounced at Bayer Crop Science, where the efficient and accurate prediction of pesticide degradation is crucial for developing environmentally friendly and safe agricultural products. The implementation of advanced computational methods, such as deep learning, can significantly enhance Bayer Crop Science's ability to predict and manage the environmental impact of their pesticide products, ensuring compliance with regulatory standards and promoting sustainable agricultural practices.

## 1.2 Problem Statement

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 pesticide degradation. 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.

### **1.3 Purpose and Research Question**

This thesis aims to develop a deep learning model to predict the degradation of pesticides based on enzyme classes. The core research question guiding this study is: "How can deep learning be applied to predict pesticide degradation pathways based on enzyme class data?" By addressing this question, the study seeks to contribute to the fields of computational biology and environmental science, providing a tool that can aid in the rapid assessment of pesticide biodegradation potential. 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.

### **1.4 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 purpose and research question 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 also includes sections on the principles of enzymology, detailing enzyme classification, function, and their role in biodegradation, as well as the fundamentals of deep learning, including introductions to ligand-binding prediction and model evaluation techniques.

The methodology chapter provides a detailed description of the research design and procedures followed in this study. It begins with 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 improve model performance. The chapter then explains the model development process, including the architecture of the deep learning model, the training process, and the techniques used for model evaluation to ensure its reliability and accuracy.

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 breakdown of pesticides in the environment is a complex process involving various mechanisms, primarily driven by microbial enzymes. These enzymes catalyze reactions that convert toxic pesticide compounds into less harmful substances, facilitating their removal from the environment. This section explores the key enzymatic mechanisms involved in pesticide degradation, focusing on hydrolytic, oxidative, and reductive enzymes.

Microbial enzymes play a pivotal role in the biodegradation of soil contaminants, including pesticides. They can be categorized based on the reactions they catalyze:

**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.

**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.

**Reductive Enzymes:** Reductive enzymes, including reductases, catalyze the reduction 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.

Incorporating microbial enzymes into bioremediation strategies can significantly enhance the degradation of pesticides in contaminated soils. This approach leverages the natural capabilities of microbes to detoxify pollutants through enzymatic reactions. According to a review on the function of microbial enzymes in breaking down soil contaminated with pesticides, these enzymes are highly effective in transforming and mineralizing pesticides, thus reducing their environmental impact.<sup>1</sup>

Another study highlights the advancements and applications of microbial enzymes in biodegradation processes. This review emphasizes the critical role of enzymes in the degradation pathways of various pesticides and discusses the potential for engineered enzymes to improve bioremediation efficiency.<sup>2</sup>

## 2.2 Deep Learning Techniques in Environmental Science

Deep learning has emerged as a powerful tool in environmental science, offering advanced methods for predicting and understanding complex biochemical processes. In the context of pesticide degradation, deep learning models can analyze vast amounts of biochemical data to predict enzyme interactions and degradation pathways. There are several deep learning architectures that have been successfully applied to enzyme classification and prediction tasks, providing valuable insights into the mechanisms of pesticide breakdown.

For instance, the DEEPre model uses deep learning to predict enzyme commission (EC) numbers from raw sequence data. This model has shown significant improvements in prediction accuracy over traditional methods by utilizing convolutional and sequential feature extraction techniques. Such models can be crucial for predicting the biodegradation pathways of pesticides, enabling more accurate and efficient environmental risk assessments.<sup>3</sup>

Another example is the DeEPn model, which uses a deep neural network to classify enzymes into their functional classes, including all seven EC classes. This model

---

<sup>1</sup>Singh, Brajesh K. and Walker, Allan (2006).

<sup>2</sup>Chia, Xing Kai et al. (2024).

<sup>3</sup>Li, Yu et al. (2017).

has demonstrated high precision and accuracy, making it a valuable tool for environmental scientists looking to understand and predict enzyme-mediated degradation processes. By accurately classifying enzymes, DeEPn facilitates the identification of potential candidates for bioremediation and other environmental applications.<sup>4</sup>

Despite the advances made by these models, there is still a need for new approaches to further improve the accuracy and applicability of pesticide degradation predictions. Traditional models often rely on pre-defined features and limited datasets, which can restrict their performance and generalizability. By contrast, my 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.<sup>5</sup> This method can potentially provide a more detailed and accurate prediction of enzyme classes responsible for pesticide degradation, enhancing our understanding of the biodegradation pathways and mechanisms involved.

---

<sup>4</sup>*DeEPn: A Deep Neural Network Based Tool for Enzyme Functional Annotation - Consensus* (2024).

<sup>5</sup>Krivák, Radoslav and Hoksza, David (2018).

## 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.<sup>6</sup>

#### 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:

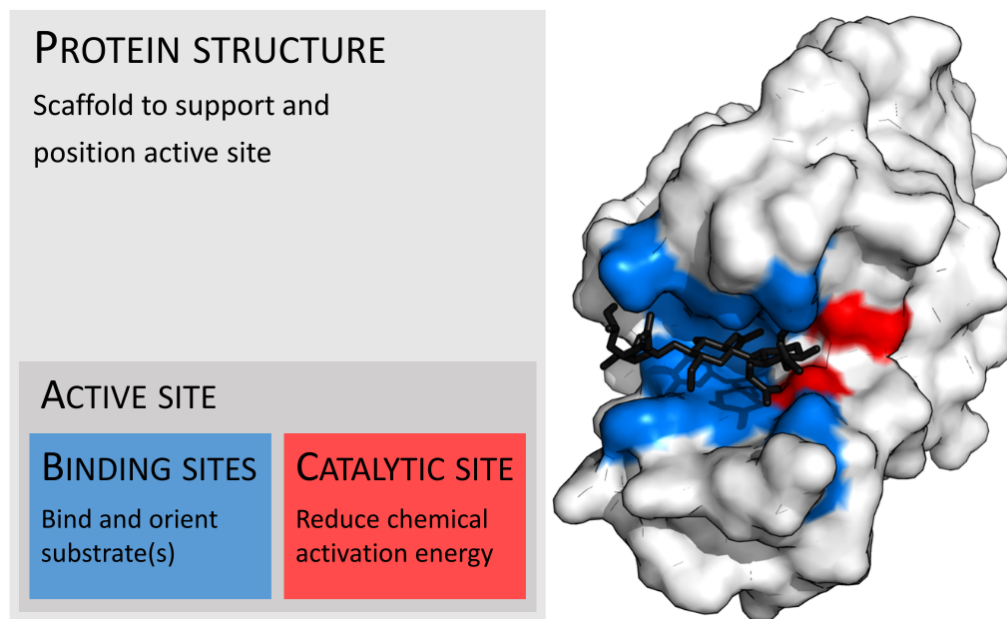
1. **Oxidoreductases:** These enzymes catalyze oxidation-reduction reactions, where the transfer of electrons occurs between molecules. Examples include dehydrogenases and oxidases.
2. **Transferases:** These enzymes transfer functional groups from one molecule to another. Examples include kinases, which 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.

---

<sup>6</sup>Robinson, Peter K. (2015).

- Überleitung finden – 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 1:** 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 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.

A study by Veselovsky et al. (2001) emphasizes the importance of visualizing active site structures, even for enzymes with unknown 3D structures. By analyzing enzyme interactions with reversible competitive inhibitors and molding the substrate-binding region, researchers can predict the shape and dimensions of the active site. This approach has been validated by comparing it with known enzyme-inhibitor complexes, demonstrating its utility in understanding enzyme function and aiding in the search for new ligands.<sup>7</sup>

### 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 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.<sup>8</sup>

**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.<sup>9</sup>

---

<sup>7</sup>Veselovsky, A. et al. (2001).

<sup>8</sup>Munnecke, D. (1976).

<sup>9</sup>Bello, Angelica, Carreon, Yessica, and Nava-Ocampo, Alejandro (2000).



**Reductive Enzymes:** Reductive enzymes, including reductases, catalyze the reduction 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 Machine Learning and Deep Learning**

Machine Learning and Deep Learning are two powerful techniques in the realm of machine learning, each offering unique strengths for different types of data and prediction tasks. While deep learning excels at handling unstructured data and automatically extracting features, Machine Learning e.g. Random Forest is known for its robustness, interpretability, and effectiveness in handling structured data with high-dimensional features. This section focuses on the use of Random Forest for sequence-based predictions, particularly in the context of predicting enzyme classes based on protein sequences.

### **3.2.1 Introduction to Deep Learning**

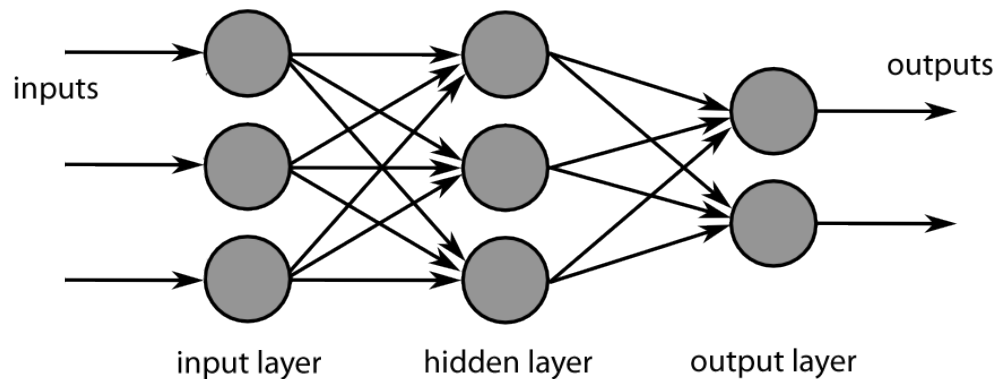
Deep learning has dramatically transformed various fields by enabling the analysis and interpretation of complex datasets. Unlike traditional machine learning methods, which often require manual feature extraction, deep learning models can automatically learn relevant features from raw data. This capability is largely due to the hierarchical structure of neural networks, which can capture multiple levels of abstraction.

A deep neural network is composed of an input layer, multiple hidden layers, and an output layer. Each layer consists of nodes (neurons) that are interconnected

with nodes from the previous and next layers. The strength of these connections is determined by weights, which are adjusted during the training process to minimize prediction error. The training is typically performed using a variant of stochastic gradient descent (SGD) and backpropagation, a method for computing the gradient of the loss function with respect to each weight.<sup>10</sup>

The following image illustrates the basic structure of a deep neural network:

**Figure 2:** Diagram of a multi-layer feedforward artificial neural network.



**Source:** Chrislbderivative work, CC BY-SA 3.0 via Wikimedia Commons

In this structure, the input layer receives raw data. This data is then processed through one or more hidden layers, where the neurons apply weights and activation functions to capture complex patterns and features. Finally, the processed information reaches the output layer, where the final prediction or classification is made.

Deep learning is particularly effective for problems involving high-dimensional and unstructured data, such as images, audio, and text. Its ability to automatically extract and learn complex features from raw data makes it superior to traditional machine learning methods, which often rely on manually engineered features that may not capture the full complexity of the data.

<sup>10</sup>Bishop, Christopher M. (2006).

### 3.2.2 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 enzyme function and substrate specificity. Several computational methods have been developed to predict ligand-binding sites from protein structures, including geometric, physicochemical, and machine learning-based approaches.

P2Rank is 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.

The following image illustrates the workflow of P2Rank, highlighting the process of predicting ligand binding sites from protein structures:

– P2Rank Workflow –

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.<sup>11</sup>

---

<sup>11</sup>Krivák, Radoslav and Hoksza, David (2018).

### 3.2.3 Sequence-Based Machine Learning

In this thesis context, sequence-based machine learning refers to the use of protein sequences as input data for predicting enzyme classes. The sequences are typically represented as strings of amino acids, with each amino acid corresponding to a specific position in the sequence. Machine learning models, such as Random Forest, can analyze these sequences and extract relevant features to predict the enzyme class. This approach uses a Random Forest algorithm, an ensemble learning method that operates by constructing a multitude of decision trees during training and outputting the mode of the classes (classification) or mean prediction (regression) of the individual trees. This method was introduced by Leo Breiman in 2001 and has since become a staple in bioinformatics and other scientific fields due to its versatility and performance.<sup>12</sup>

The workflow for sequence-based machine learning involves the following steps:

1. **Feature Extraction:** The sequence data, such as amino acid sequences, are first transformed into numerical features. This can include k-mer counts, physicochemical properties of amino acids, and other sequence-derived features.
2. **Model Training:** These features are used to train the Random Forest model, which learns to associate specific patterns in the sequence data with enzyme classes.
3. **Prediction:** Once trained, the model can predict the enzyme class of new, unseen sequences.

A significant advantage of using Random Forest in this context is its ability to handle complex interaction structures and high-dimensional feature spaces, which are common in biological sequence data.<sup>13</sup>

In the context of this study, deep learning techniques are utilized for predicting ligand-binding sites, which are crucial for understanding protein function and interaction with other molecules. Tools like p2rank can be used to predict these binding

---

<sup>12</sup>Breiman, Leo (2001).

<sup>13</sup>Díaz-Uriarte, R. and Andrés, S. Á D. (2006).

sites from protein structures, providing the base for a Random Forest model to predict the enzyme class based on the identified ligand-binding sites. This combined approach leverages the strengths of deep learning for feature extraction and Random Forest for classification, enhancing the accuracy and reliability of enzyme class predictions

### 3.2.4 Evaluation of Machine Learning Models

Evaluating deep learning models involves several metrics and techniques to ensure their accuracy and generalizability. This is essential not only for validating the model's performance but also for comparing it against other models. Using independent datasets for benchmarking is crucial to demonstrate the model's robustness and applicability to real-world scenarios. Several key metrics are commonly used to evaluate deep learning models:

1. **Accuracy:** The ratio of correctly predicted instances to the total instances. It provides a straightforward measure of performance but can be misleading if the data is imbalanced.<sup>14</sup>
2. **Precision:** The ratio of true positive predictions to the total predicted positives. Precision is crucial when the cost of false positives is high.
3. **Recall (Sensitivity):** The ratio of true positive predictions to the total actual positives. Recall is important when the cost of false negatives is high.<sup>15</sup>
4. **F1-Score:** The harmonic mean of precision and recall, providing a single metric that balances both. It is useful when there is an uneven class distribution.<sup>16</sup>
5. **ROC-AUC (Receiver Operating Characteristic - Area Under Curve):** A performance measurement for classification problems at various threshold settings. It tells how much the model is capable of distinguishing between classes.<sup>17</sup>

---

<sup>14</sup>“Accuracy and Precision” (2024).

<sup>15</sup>“Precision and Recall” (2024).

<sup>16</sup>“F-Score” (2024).

<sup>17</sup>“Receiver Operating Characteristic” (2024).

Using these metrics, the performance of deep learning models can be tested on independent datasets. This is critical for ensuring that the models are not just overfitting to the training data but can generalize well to new, unseen data.

In the context of comparing different models, these metrics allow for a standardized evaluation process. By applying the models to benchmark datasets, researches can objectively measure and compare their performance. For instance, in the case of predicting enzyme-mediated pesticide degradation, using independent datasets ensures that the model's predictions are reliable and can be generalized across various enzyme and pesticide types.

Model evaluation is essential for several reasons:<sup>18</sup>

1. **Validation of Model Performance:** Ensuring that the model performs well not just on the training data but also on new, unseen data. This is crucial for the model's reliability and applicability in real-world scenarios.
2. **Comparison with Other Models:** By using standardized evaluation metrics and independent datasets, we can objectively compare the performance of different models. This helps in identifying the most effective model for a given task.
3. **Identification of Overfitting:** Evaluating the model on independent datasets helps in identifying overfitting, where the model performs well on training data but poorly on new data. Regularization techniques and cross-validation can be used to mitigate this issue.
4. **Continuous Improvement:** Regular evaluation allows for continuous monitoring and improvement of the model. Feedback from evaluation metrics can guide further tuning and optimization of the model.
5. **Transparency and Trust:** Providing clear and objective evaluation metrics builds trust with stakeholders and users of the model. It ensures that the model's predictions are transparent and can be trusted for decision-making.

---

<sup>18</sup>Emmert-Streib, F. et al. (2020).

Using benchmark datasets and standardized metrics is a best practice in machine learning and deep learning. It ensures that the models are robust, reliable, and ready for deployment in practical applications. In the context of environmental science and bioremediation, such rigorous evaluation is crucial for developing effective and reliable models for predicting pesticide degradation and other complex processes.

## 4 Methodology

### 4.1 Data Collection

Data collection is a critical step in developing predictive models, as the quality and relevance of the data directly impact the model's performance. In this thesis, the data was collected from Uniprot, a comprehensive resource for protein sequence and functional information. The focus was on obtaining high-quality, reviewed entries with 3D structural data and catalytic properties to ensure the reliability and applicability of the data for predicting enzyme functions.<sup>19</sup>

Uniprot, or the Universal Protein Resource, is a central repository of protein sequence and annotation data. It is widely recognized for its comprehensive, high-quality data, making it an essential resource for bioinformatics and computational biology. Uniprot integrates information from various sources, including experimental studies, computational analysis, and literature, providing a rich and reliable dataset for scientific research. Key features of Uniprot include:

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.

---

<sup>19</sup>UniProt Consortium (2021).



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 next sections describe the data preprocessing, feature engineering, and model development steps in detail.

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:<sup>20</sup>

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.

---

<sup>20</sup>Polley, Tobias (2024).

## **4.2 Data Preprocessing**

## **4.3 Feature Engineering**

## **4.4 Model Development**

### **4.4.1 Model Architecture**

### **4.4.2 Model Training**

### **4.4.3 Model Evaluation**

**Code 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')

```

## **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

|             |                  |    |
|-------------|------------------|----|
| Anhang 1:   | Filler . . . . . | 25 |
| Anhang 1.1: | Filler . . . . . | 25 |

## **Appendix 1    Filler**

### **Appendix 1.1    Filler**

- Filler



## List of references

## 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, Wednesday 12<sup>th</sup> June, 2024

---

Tobias Polley