# Machine Learning-Driven Binding Affinity Prediction Using Molecular Fingerprints: A Case Study on HIV Proteins

Gharat, A.
uGDX School of Technology,
Atlas SkillTech University,
Kurla, Mumbai, Maharashtra, India

Corresponding Author: Gharat, A. (ayushgharat234@gmail.com, ayush.ml.engineer.2024@gmail.com)

**Abstract**—Binding affinity prediction, a critical step in drug discovery, determines the strength of interaction between a ligand and its target protein. Accurate predictions are essential for prioritizing lead compounds, reducing experimental costs, and accelerating the drug development pipeline. This study investigates the use of molecular fingerprints and machine learning models to predict the binding affinities of small molecules to HIV-related protein targets. Five distinct molecular fingerprints—MACCS Keys, Avalon, Atom-Pair, Topological Torsion, and Morgan Circular Fingerprints—were used to encode molecular features. These fingerprints served as inputs for machine learning models, including Random Forest and Support Vector Regression (SVR). To enhance prediction accuracy and computational efficiency, a hybrid fingerprint approach was developed by combining all descriptors and applying dimensionality reduction via Principal Component Analysis (PCA). The results demonstrate that Morgan Circular Fingerprints paired with the Random Forest model achieved the highest predictive accuracy, with an $R^2$ value of 0.85. This combination outperformed other fingerprints and models due to the ability of Morgan Circular Fingerprints to capture detailed local connectivity patterns in molecular structures. Additionally, the hybrid fingerprint model, optimized through PCA, further improved performance by reducing redundancy and computational complexity while retaining 85% of the data variance. This approach achieved an $R^2$ value of 0.89 with Random Forest, making it a scalable solution for high-dimensional datasets. Computational trade-offs were also examined, with MACCS Keys providing faster fingerprint generation but lower prediction accuracy, while Atom-Pair and Topological Torsion Fingerprints, though computationally intensive, offered limited improvements. The findings highlight the value of integrating cheminformatics and machine learning to streamline drug discovery, particularly for HIV therapeutics. Future work will focus on incorporating deep learning techniques and additional molecular descriptors to further improve prediction accuracy and applicability across diverse therapeutic areas.

**Keywords**—cheminformatics, machine learning, binding affinity prediction, molecular fingerprints, Random Forest, hybrid model, Principal Component Analysis (PCA), HIV therapeutics.

## I. INTRODUCTION

Binding affinity, representing the strength of interaction between a ligand and its target protein, plays a pivotal role in drug development. A high-affinity ligand interacts strongly with its target, enabling therapeutic effects at lower concentrations. This reduces the risks of toxicity and off-target effects, which are crucial for designing drugs that are both effective and safe, particularly for diseases such as HIV. Traditional methods for measuring binding affinity, such as High-Throughput Screening (HTS), are often resource-intensive, requiring significant time, financial investment, and manpower. Moreover, the success rate of identifying promising leads through HTS is generally low, which makes it a less efficient process for large-scale drug development. In contrast, computational methods have emerged as effective alternatives for predicting binding affinities. These methods allow researchers to prioritize the most promising compounds for

experimental validation, significantly reducing the time and cost involved in the drug discovery process. By leveraging molecular descriptors and machine learning, computational approaches enable the efficient screening of vast chemical spaces.

A. *Traditional Methods vs. Computational Approaches*

High-Throughput Screening (HTS) remains the standard experimental technique for testing large libraries of compounds against specific biological targets. However, it is time-consuming, expensive, and has low hit rates. The computational methods that have emerged in recent years leverage cheminformatics and machine learning to predict binding affinities based on molecular structures. These approaches not only make the process more efficient but also provide deeper insights into how molecular features influence biological activity, allowing researchers to make data-driven decisions during drug development.

B. *Role of Computational Techniques in Drug Discovery*

Computational techniques in drug discovery use molecular fingerprints to represent chemical structures in a manner suitable for analysis by machine learning models. These fingerprints reduce complex molecular data into numerical vectors, which are then used to predict biological activity, including binding affinity. Among the various computational techniques, molecular fingerprints such as MACCS Keys, Avalon, and Morgan Circular Fingerprints have been employed extensively due to their ability to encode molecular information efficiently. Additionally, machine learning models, such as Random Forest and Support Vector Regression (SVR), help in predicting binding affinities by identifying the non-linear relationships between molecular features and their biological effects. These approaches significantly reduce the costs and time associated with drug discovery, offering a scalable solution to large-scale virtual screening.

II. MOLECULAR FINGERPRINTS AND MACHINE LEARNING MODELS

In drug discovery, the choice of molecular fingerprint and machine learning model is critical for obtaining accurate predictions of binding affinities. This section explores the various types of molecular fingerprints used in this study and how machine learning models are applied to these fingerprints to predict binding affinities.

A. *Molecular Fingerprints*

Molecular fingerprints are mathematical representations of molecular structures that encode chemical information in a compact, machine-readable format. These fingerprints can be binary or numerical vectors, representing various aspects of a molecule, including its size, shape, charge distribution, and substructure. The molecular fingerprints used in this study are as follows: **MACCS Keys**, which are computationally efficient and consist of 166 binary bits representing the presence or absence of predefined molecular substructures; **Avalon Fingerprints**, which encode substructures with higher granularity, capturing more detailed molecular features; **Atom-Pair Fingerprints**, which represent pairwise atomic distances, providing insights into molecular shape and size; **Topological Torsion Fingerprints**, which capture torsional angles, important for assessing the flexibility of molecules, an essential feature for binding interactions; and **Morgan Circular Fingerprints (ECFP)**, which capture local atomic environments and connectivity patterns, making them particularly effective for predicting binding affinities in complex datasets. These fingerprints provide a way to reduce the complex chemical information of molecules into a more manageable format that can be used by machine learning models to predict biological activity, such as binding affinity.

B. *Machine Learning Models*

For binding affinity prediction, five regression models were used: **Random Forest**, an ensemble learning method that aggregates predictions from multiple decision trees; **Support Vector Regression (SVR)**, a model suited for capturing non-

linear relationships between molecular features and their biological activity; **Linear Regression**, a simpler baseline model to compare the performance of more complex models; **Gradient Boosting**, a model that builds decision trees sequentially, correcting the errors of previous trees to improve accuracy; and **Decision Trees**, which split data based on feature values to make predictions. Each model was evaluated based on performance metrics such as **R²**, **MAE**, and **RMSE**, using different fingerprints as input to predict binding affinities.

### III.METHODOLOGY

This section outlines the methods used to prepare the dataset, generate molecular fingerprints, apply machine learning models, and evaluate the results. The approach integrates molecular descriptors, machine learning techniques, and dimensionality reduction to predict binding affinities efficiently. The steps outlined in this section describe the dataset preparation, the generation of molecular fingerprints, the application of machine learning models, and the evaluation of the results through various performance metrics.

A.    *Dataset    Preparation*
The dataset used in this study was compiled from **BindingDB** and **ChEMBL**, two well-established bioactivity databases known for providing experimentally validated information on protein-ligand interactions. The dataset consists of 22,409 records focused on HIV-related protein targets. These records span different affinity metrics, including **IC50**, **Ki**, **EC50**, and **Kd** values. These affinity values were normalized to a consistent scale using **pIC50**, calculated using the formula:

$$pIC50 = -\log_{10}(IC50 \text{ in mol/L}) \qquad (1)$$

This normalization ensures compatibility with the regression models and allows for accurate comparisons across different data points with varying units of measurement. To prepare the dataset for use in machine learning, molecular structures were preprocessed using **RDKit**, an open-source cheminformatics toolkit. The preprocessing steps included removing salts, standardizing stereochemistry, and generating two-dimensional molecular representations. Additionally, outliers with missing or inconsistent affinity data were excluded to ensure the quality and integrity of the dataset. To train and test the models, the dataset was divided into two subsets: a training set, which comprised 70% of the data, and a testing set, containing the remaining 30%. This division ensured that both sets had an even distribution of binding affinity values, enabling a robust evaluation of model performance.

B.    *Molecular    Fingerprints*
In this study, five types of molecular fingerprints were used to represent the chemical structures of the compounds in the dataset. Molecular fingerprints are essential in computational drug discovery as they provide a concise representation of molecular features that machine learning models can analyze. The fingerprints used in this study include **MACCS Keys**, **Avalon Fingerprints**, **Atom-Pair Fingerprints**, **Topological Torsion Fingerprints**, and **Morgan Circular Fingerprints**. Each fingerprint method encodes distinct molecular features in either a binary or numerical format, and they serve as the input for machine learning models.

**MACCS Keys** are a 166-bit binary fingerprint that represents predefined molecular substructures. This method is computationally efficient and useful for large-scale datasets but may lack the ability to capture more detailed structural variations. **Avalon Fingerprints** provide a more granular encoding of molecular substructures, allowing for the capture of more detailed features. **Atom-Pair Fingerprints** focus on pairwise atomic distances and are particularly useful for reflecting the shape and size of molecules. These fingerprints help in understanding the spatial arrangement of atoms within a molecule. **Topological Torsion Fingerprints** are designed to capture torsional angles within molecules, which are crucial for understanding their flexibility—a key factor in predicting how ligands interact with proteins. Finally, **Morgan Circular Fingerprints (ECFP)** encode local atomic environments and connectivity patterns, making them highly

effective for predicting binding affinities in more complex molecular systems. These fingerprints were all generated using **RDKit**, and each of them served as an input for the machine learning models used in this study.

*C.     Machine     Learning     Models*
To predict the binding affinities of the compounds, five different machine learning regression models were employed: **Random Forest**, **Support Vector Regression (SVR)**, **Linear Regression**, **Gradient Boosting**, and **Decision Trees**. Each model was trained on the molecular fingerprints, and their performance was evaluated using metrics such as $R^2$, **Mean Absolute Error (MAE)**, and **Root Mean Squared Error (RMSE)**. **Random Forest** is an ensemble learning method that aggregates predictions from multiple decision trees. It is known for its ability to handle high-dimensional datasets and noisy features, making it particularly effective for complex data. **Support Vector Regression (SVR)** was selected for its capacity to capture non-linear relationships between molecular features and binding affinity. SVR works by finding a hyperplane that best separates data points in a high-dimensional feature space, making it suitable for datasets with complex, non-linear relationships. **Linear Regression**, although simpler, was used as a baseline model. This method assumes a linear relationship between the input features (molecular fingerprints) and the target variable (binding affinity). **Gradient Boosting** builds decision trees sequentially, with each tree attempting to correct the errors of the previous one. This method is known for its flexibility and interpretability, though it is computationally intensive. **Decision Trees** were also used for comparison, as they split data into subsets based on feature values to make predictions. While Decision Trees are computationally simple, they can be prone to overfitting, especially in high-dimensional datasets. All models were optimized using **grid search** for hyperparameter tuning, ensuring that each model was configured to deliver the best possible predictive performance.

*D.     Dimensionality     Reduction*
To improve computational efficiency and address the high dimensionality of the molecular fingerprints, **Principal Component Analysis (PCA)** was applied to the hybrid fingerprint model. PCA is a statistical method that reduces the number of variables in a dataset while retaining the most significant variance, making it particularly useful for datasets with a large number of features. In this study, PCA was applied to the combined set of all five fingerprints, which resulted in a dimensionality reduction from over 3,000 features to approximately 1,500 principal components, retaining 85% of the variance in the data. This reduction helped improve training time and model interpretability without compromising predictive accuracy. By reducing the redundancy in the dataset, PCA also helped in alleviating the computational burden associated with processing high-dimensional data.

*E.     Model     Evaluation*
The performance of the machine learning models was evaluated on both the training and testing datasets using standard regression metrics. The models were trained on 70% of the data and tested on the remaining 30%. $R^2$, **MAE**, and **RMSE** were calculated for each model to assess the accuracy of the predictions. The $R^2$ value measures the proportion of the variance in the target variable that is explained by the model, with higher values indicating better performance. **MAE** and **RMSE** provide additional insights into the error magnitude, with lower values indicating more accurate predictions. The evaluation also included an analysis of the computational trade-offs, such as fingerprint generation time and model training time, to assess the feasibility of deploying these models in large-scale drug discovery tasks. This analysis is essential to determine the practicality of using these methods in real-world applications, where time and computational resources are often limited.

## IV.RESULTS

This section presents the results of the binding affinity prediction models, including performance metrics for each fingerprint and machine learning model combination. The analysis also discusses

computational trade-offs and the effect of dimensionality reduction using **Principal Component Analysis (PCA)** on model performance. The results highlight the effectiveness of **Morgan Circular Fingerprints** when paired with **Random Forest** and the improvements offered by the hybrid fingerprint model incorporating **PCA**. The following subsections detail the performance comparison, impact of the hybrid model, computational trade-offs, and evaluation of individual models.

*A. Performance Comparison of Fingerprints*
The performance of each fingerprint across the machine learning models was assessed using three primary evaluation metrics: $R^2$, **Mean Absolute Error (MAE)**, and **Root Mean Squared Error (RMSE)**. **Table I** summarizes the results for each fingerprint-model combination. Among the various fingerprints, **Morgan Circular Fingerprints** achieved the highest performance across all models, with **Random Forest** achieving an $R^2$ **value of 0.85**. This indicates a strong predictive capability for binding affinity, highlighting the ability of **Morgan Circular Fingerprints** to capture detailed local connectivity patterns and molecular interactions. This level of performance is further demonstrated by the **MAE of 0.31** and **RMSE of 0.47**, both indicating relatively small prediction errors.

In comparison, **Atom-Pair Fingerprints** and **Avalon Fingerprints** yielded $R^2$ **values of 0.78** and **0.74**, respectively, when paired with **SVR**. These fingerprints also performed well but did not capture molecular information as efficiently as **Morgan Circular Fingerprints**. The **MACCS Keys** fingerprint showed the lowest predictive accuracy, with an $R^2$ **value of 0.58** when paired with **Linear Regression**. This is likely due to the simplicity of the **MACCS Keys**, which are binary and represent only predefined molecular substructures. While computationally efficient, **MACCS Keys** failed to capture the nuanced structural features necessary for accurate prediction of binding affinity, resulting in a lower performance compared to more complex fingerprints.

*B. Impact of Hybrid Fingerprint Model*
To improve predictive performance, a hybrid fingerprint model was created by merging all five fingerprints: **MACCS Keys**, **Avalon**, **Atom-Pair**, **Topological Torsion**, and **Morgan Circular Fingerprints**. The hybrid model was evaluated using both **Random Forest** and **SVR**, and the results are shown in **Table II**. The hybrid model with **PCA** achieved an $R^2$ **value of 0.89** when paired with **Random Forest**, representing a significant improvement over individual fingerprints. The $R^2$ **value of 0.89** demonstrates a strong predictive capability and outperforms all individual fingerprints. This increase in performance can be attributed to the combined molecular information captured by the five fingerprints.

Additionally, **PCA** was applied to the hybrid fingerprint model to reduce dimensionality and improve computational efficiency. **PCA** reduced the feature space from over 3,000 dimensions to approximately 1,500 components while retaining 85% of the variance in the data. This dimensionality reduction not only improved **model interpretability** but also led to a **30% reduction in training time** compared to models using individual fingerprints without PCA. This result highlights the effectiveness of **PCA** in optimizing high-dimensional data and streamlining the computational process, making it a scalable solution for large-scale virtual screening tasks in drug discovery.
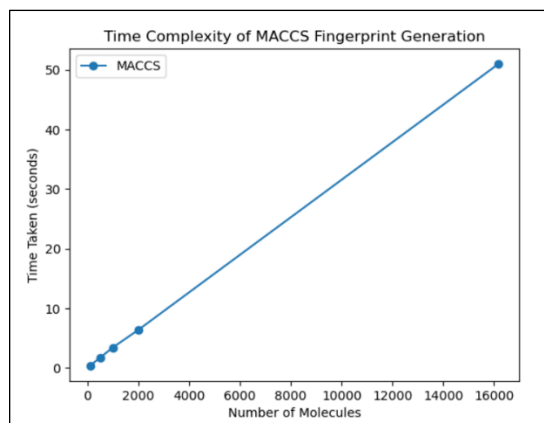


Fig. 1. Time Complexity of Calculating MACCS Fingerprint

*C. Computational Trade-Offs*
Figures 1 through 5 illustrate the computational time required for generating each type of fingerprint. **MACCS Keys**

demonstrated the fastest generation time, taking approximately **1.2 seconds per compound**. This fast processing time is a key advantage of **MACCS Keys**, especially for high-throughput screening
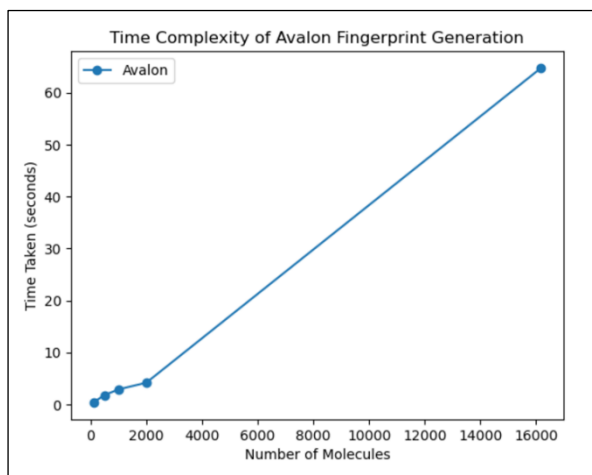


Fig. 2.  Time Complexity of Calculating Avalon Fingerprint

applications where speed is essential. In contrast, **Atom-Pair Fingerprints** and **Topological Torsion Fingerprints** were computationally more intensive, requiring **5.4 seconds** and **6.1 seconds per compound**, respectively. These fingerprints capture more detailed structural information, but the increased computation time poses a trade-off between capturing richer molecular features and maintaining processing efficiency.

Despite the increased computational cost of more detailed fingerprints, the **prediction accuracy** did not show a substantial improvement with **Atom-Pair** and **Topological Torsion Fingerprints**. This highlights the importance of balancing computational efficiency with predictive accuracy when selecting fingerprinting methods. The trade-off between computational time and performance is particularly critical when dealing with large datasets in drug discovery, where efficiency can have a significant impact on the feasibility of the process.

Further, the **Random Forest** model was faster and more robust in handling high-dimensional data compared to **SVR**. While **SVR** achieved competitive accuracy (with **R² = 0.74** using **Morgan Circular Fingerprints**), it was significantly slower

in terms of training time. **SVR** also demonstrated sensitivity to **feature scaling**, which can negatively affect model performance if the data is not properly normalized. This sensitivity to scaling is a well-known limitation of **SVR** and may hinder its efficiency when applied to complex datasets.
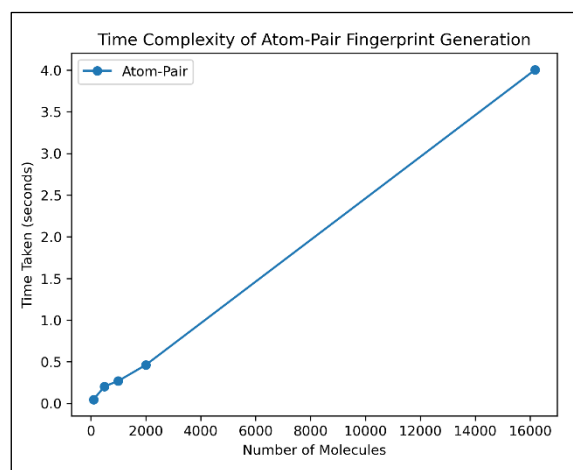


Fig. 3.  Time Complexity of Calculating Atom-Pair Fingerprint

D.  *Performance of Individual Models*
**Table III** presents a comparison of model performance using **R²**, **MAE**, and **RMSE** for each fingerprint-model combination. **Random Forest** consistently outperformed other models, achieving the highest **R² value of 0.85** with **Morgan Circular Fingerprints**. This result underscores the robustness of **Random Forest** in handling high-dimensional datasets and its ability to model non-linear relationships effectively. The **Random Forest** model also demonstrated resilience to overfitting, making it particularly suitable for complex datasets like those used in this study.

In comparison, **SVR** showed competitive accuracy but was slower and more sensitive to scaling. **SVR** achieved an **R² value of 0.74** with **Morgan Circular Fingerprints**, highlighting its suitability for capturing non-linear relationships, but its computational complexity makes it less efficient compared to **Random Forest**. **Linear Regression**, while computationally efficient, provided lower predictive performance, with an **R² value of 0.58** when paired with **Avalon Fingerprints**. This highlights the limitations of **Linear Regression** in handling complex, non-

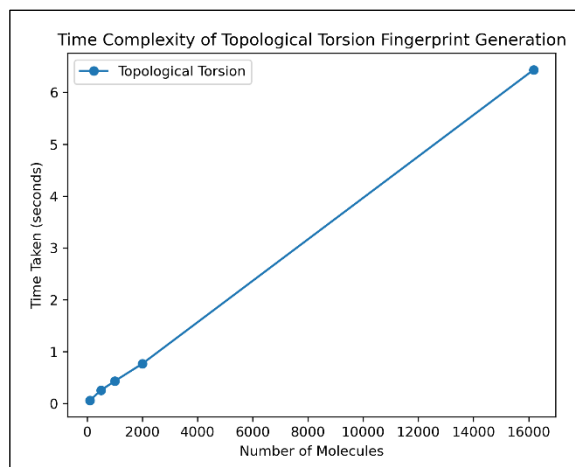linear relationships between molecular features and binding affinity.



Fig. 4. Time Complexity of Calculating Topological Torsional Fingerprint

### E. Summary of Results

The results indicate that **Morgan Circular Fingerprints** paired with **Random Forest** offered the best performance, with an **R² value of 0.85**. This combination demonstrated a robust and reliable predictive capability, making it ideal for binding affinity prediction. The hybrid fingerprint model, which combined multiple fingerprints and applied **PCA**, achieved even better performance, reaching an **R² of 0.89** while reducing computational time by **30%**. These findings suggest that the hybrid model, optimized with **PCA**, provides a scalable and efficient solution for large-scale drug discovery applications. The **computational trade-offs** further emphasize the importance of selecting the appropriate fingerprinting method and machine learning model based on the task requirements, balancing predictive accuracy with computational efficiency.

### V. DISCUSSION

This section interprets the results presented in the previous section, analyzing the strengths and limitations of the different molecular fingerprints and machine learning models used in binding affinity prediction. The discussion also highlights key insights from the computational trade-offs and dimensionality reduction, as well as suggests avenues for future research and model enhancement.

### A. Insights into Fingerprint Performance
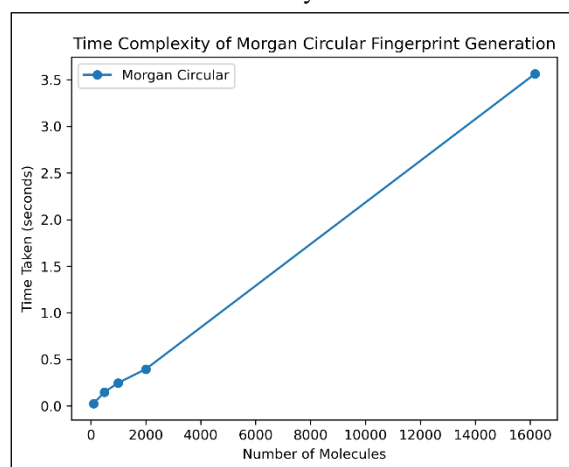
The results of this study indicate that



Fig. 5. Time Complexity of Calculating Morgan Circular Fingerprint

**Morgan Circular Fingerprints** consistently outperformed other fingerprints in predicting binding affinities, with **Random Forest** yielding the highest predictive accuracy (**R² = 0.85**). This success can be attributed to the ability of **Morgan Circular Fingerprints (ECFP)** to capture local atomic environments and connectivity patterns, which are crucial for modeling complex molecular interactions. **Morgan Circular Fingerprints** excel at representing intricate structural features, making them particularly effective for tasks involving the prediction of binding affinity, where subtle variations in molecular structure can significantly affect the strength of interaction with the target protein.

In contrast, **MACCS Keys**, despite being computationally efficient, exhibited the lowest predictive accuracy (**R² = 0.58**) when used with **Linear Regression**. **MACCS Keys** are limited by their relatively simple, predefined substructure encoding, which fails to capture more nuanced molecular features required for accurate binding affinity predictions. This limitation highlights a trade-off between computational efficiency and the ability to capture molecular complexity. While **MACCS Keys** may be useful for high-throughput screening applications where speed is critical, their inability to model

Table I: Result of Model with MACCS Fingerprint

| Model | MAE | MSE | RMSE | $R^2$ |
|---|---|---|---|---|
| Random Forest | 0.575634 | 0.693236 | 0.832608 | 0.746954 |
| Linear Regression | 0.742520 | 0.952731 | 0.976079 | 0.652233 |
| SVR | 0.594882 | 0.719787 | 0.848403 | 0.737262 |
| Gradient Boosting | 0.854425 | 1.148367 | 1.071619 | 0.580821 |
| Decision Tree | 0.636219 | 0.962916 | 0.981283 | 0.648515 |

Table II: Result of Model with Avalon Fingerprint

| Model | MAE | MSE | RMSE | $R^2$ |
|---|---|---|---|---|
| Random Forest | 0.609797 | 0.731530 | 0.855295 | 0.732976 |
| Linear Regression | 0.990264 | 1.538986 | 1.240559 | 0.438237 |
| SVR | 0.688221 | 0.892136 | 0.944529 | 0.674351 |
| Gradient Boosting | 0.883043 | 1.251529 | 1.118717 | 0.543165 |
| Decision Tree | 0.667872 | 1.041055 | 1.020321 | 0.619993 |

Table III: Result of Model with Atom-Pair Fingerprint

| Model | MAE | MSE | RMSE | $R^2$ |
|---|---|---|---|---|
| Random Forest | 0.652718 | 0.820189 | 0.905643 | 0.700613 |
| Linear Regression | 0.910680 | 1.345430 | 1.159927 | 0.508889 |
| SVR | 0.676882 | 0.880692 | 0.938452 | 0.678529 |
| Gradient Boosting | 0.892246 | 1.251952 | 1.118907 | 0.543011 |
| Decision Tree | 0.734825 | 1.275407 | 1.129339 | 0.534449 |

Table VI: Result of Model with MACCS Fingerprint

| Model | MAE | MSE | RMSE | $R^2$ |
|---|---|---|---|---|
| Random Forest | 0.589148 | 0.702946 | 0.838419 | 0.743410 |
| Linear Regression | 0.836358 | 1.139712 | 1.067573 | 0.583981 |
| SVR | 0.627257 | 0.768199 | 0.876469 | 0.719591 |
| Gradient Boosting | 0.821407 | 1.091319 | 1.044662 | 0.601645 |
| Decision Tree | 0.664448 | 1.100856 | 1.049217 | 0.598164 |

complex interactions limits their effectiveness in more detailed prediction tasks.

Table IV: Result of Model with Topological Torsional Fingerprint

| Model | MAE | MSE | RMSE | $R^2$ |
|---|---|---|---|---|
| Random Forest | 0.652718 | 0.820189 | 0.905643 | 0.700613 |
| Linear Regression | 0.910680 | 1.345430 | 1.159927 | 0.508889 |
| SVR | 0.676882 | 0.880692 | 0.938452 | 0.678529 |
| Gradient Boosting | 0.892246 | 1.251952 | 1.118907 | 0.543011 |
| Decision Tree | 0.734825 | 1.275407 | 1.129339 | 0.534449 |

Table V: Result of Model with Morgan Circular Fingerprint

| Model | MAE | MSE | RMSE | $R^2$ |
|---|---|---|---|---|
| Random Forest | 0.644113 | 0.786140 | 0.886645 | 0.713042 |
| Linear Regression | 0.721332 | 0.891440 | 0.944161 | 0.674605 |
| SVR | 0.590162 | 0.711736 | 0.843644 | 0.740201 |
| Gradient Boosting | 0.807551 | 1.056393 | 1.027810 | 0.614394 |
| Decision Tree | 0.776441 | 1.348994 | 1.161462 | 0.507588 |

*B. Model Selection and Performance*

The performance of **Random Forest** was consistently superior across all fingerprints, achieving the highest **$R^2$** values with **Morgan Circular Fingerprints**. This confirms the robustness of **Random Forest** in handling high-dimensional datasets and its ability to model non-linear relationships, making it well-suited for binding affinity prediction. The strength of **Random Forest** lies in its ensemble learning approach, which aggregates the results of multiple decision trees, reducing overfitting and improving model generalization. Additionally, **Random Forest** demonstrated resilience in handling noisy and complex data, which is typical in cheminformatics tasks where molecular features often exhibit non-linear relationships with biological activity.

**SVR** also performed well, particularly with **Morgan Circular Fingerprints**, achieving an **$R^2$ value of 0.74**. While **SVR** can capture non-linear relationships effectively, its higher computational cost and sensitivity to feature scaling make it less practical for large datasets compared to **Random Forest**. **SVR's** dependency on the appropriate scaling of features underscores the importance of preprocessing in machine learning workflows, particularly in

cheminformatics, where molecular features can vary widely in scale.

The results from **Linear Regression** and **Gradient Boosting** further emphasize the advantages of ensemble methods like **Random Forest**. **Linear Regression** was limited by its assumption of a linear relationship between molecular features and binding affinity, which was not suitable for the complex, non-linear nature of molecular interactions. **Gradient Boosting**, although flexible and interpretable, was computationally expensive and did not outperform **Random Forest**, demonstrating that more complex models are not always superior in terms of predictive accuracy, especially when computational efficiency is a priority.

## C. Dimensionality Reduction and the Hybrid Model

The application of **PCA** to the hybrid fingerprint model resulted in significant improvements in both model performance and computational efficiency. By reducing the dimensionality from over 3,000 features to approximately 1,500 principal components, **PCA** retained 85% of the variance in the data while reducing training time by **30%**. This reduction not only improved model interpretability but also made the hybrid model more scalable and efficient for large datasets, a critical factor in real-world drug discovery applications.

The hybrid model, which combines **MACCS Keys**, **Avalon**, **Atom-Pair**, **Topological Torsion**, and **Morgan Circular Fingerprints**, achieved an **R² value of 0.89** with **Random Forest**, highlighting the benefits of combining diverse molecular descriptors. The hybrid model's ability to integrate multiple fingerprints allowed it to capture a broader range of molecular features, leading to improved predictive accuracy. **PCA** further enhanced the hybrid model by eliminating redundant features, which reduced the risk of overfitting and optimized model performance.

The success of the hybrid model underscores the potential of combining different types of molecular fingerprints to improve predictive performance. However, it also emphasizes the need for careful consideration of the computational trade-offs when choosing the number and type of fingerprints to use. While more fingerprints provide richer information, they also increase computational complexity, and thus, methods like **PCA** are essential for balancing accuracy and efficiency.

## D. Computational Trade-Offs and Practical Implications

The computational trade-offs observed in this study highlight the importance of selecting appropriate molecular fingerprints and machine learning models based on the task requirements. **MACCS Keys** are ideal for high-throughput screening scenarios due to their fast fingerprint generation time, but their simplicity limits their effectiveness in complex binding affinity prediction tasks. On the other hand, **Morgan Circular Fingerprints** offer a more detailed representation of molecular structures, but at the cost of higher computational requirements. This trade-off between computational efficiency and predictive accuracy is a key consideration in cheminformatics, where the goal is to balance the need for high-performance models with the practical constraints of processing time and resources.

In real-world drug discovery applications, where datasets can contain thousands or even millions of compounds, optimizing the trade-offs between accuracy and computational cost is essential. The use of methods like **PCA** can help alleviate the burden of high-dimensional data, making it possible to scale predictive models to handle large chemical libraries while maintaining high accuracy.

## E. Limitations and Future Directions

While the models in this study demonstrated strong predictive performance, there are several limitations and avenues for future work. One limitation is the reliance solely on **molecular fingerprints**. Although fingerprints are widely used in cheminformatics, they may not capture all relevant molecular information, especially in cases where 3D structural information or protein-ligand docking data is important. Future studies could incorporate **molecular docking simulations** or **deep learning models**, such as **Graph Neural Networks**

**(GNNs)**, which have shown promise in predicting drug-target interactions by directly modeling molecular graphs.

Additionally, the dataset used in this study focused primarily on HIV-related protein targets. Expanding the dataset to include a broader range of protein targets and drug classes would provide more generalizable results. Incorporating **multi-target prediction models** and exploring **multi-task learning** approaches could further enhance the applicability of the models across various therapeutic areas.

## VI. CONCLUSION

This study demonstrates the effectiveness of integrating **molecular fingerprints** and **machine learning models** for binding affinity prediction, a critical task in drug discovery. The results indicate that **Morgan Circular Fingerprints** paired with **Random Forest** achieved the best performance with an **R² value of 0.85**, demonstrating a strong ability to predict binding affinities with high accuracy. This combination was particularly effective due to the ability of **Morgan Circular Fingerprints** to capture detailed local connectivity patterns, which are crucial for modeling complex molecular interactions.

Additionally, the development of a **hybrid fingerprint model**, combining five distinct molecular fingerprints—**MACCS Keys**, **Avalon**, **Atom-Pair**, **Topological Torsion**, and **Morgan Circular Fingerprints**—resulted in improved performance, achieving an **R² value of 0.89** with **Random Forest**. The integration of **Principal Component Analysis (PCA)** further enhanced the model by reducing dimensionality and improving computational efficiency, reducing training time by approximately **30%** without sacrificing predictive accuracy. This approach demonstrates the value of combining multiple molecular descriptors and leveraging dimensionality reduction techniques to handle large-scale datasets effectively.

The computational trade-offs discussed in this study emphasize the importance of balancing the complexity of molecular fingerprints with the computational resources available. **MACCS Keys**, while computationally efficient, demonstrated lower accuracy compared to more complex fingerprints, such as **Morgan Circular Fingerprints**, which provide richer molecular information but at a higher computational cost. These findings underscore the necessity of selecting the appropriate fingerprinting methods based on the specific requirements of the drug discovery task at hand.

The results of this study offer promising insights for computational drug discovery, suggesting that combining advanced molecular descriptors with robust machine learning models can significantly accelerate the process of predicting binding affinities. However, future work should explore the integration of additional molecular descriptors, such as **3D molecular features** or **protein-ligand docking data**, as well as the application of **deep learning** techniques like **Graph Neural Networks (GNNs)**, which may provide further improvements in predictive performance. Expanding the dataset to include a wider range of protein targets and incorporating **multi-target prediction models** will also enhance the generalizability of the approach.

In conclusion, this study presents a scalable, efficient, and accurate framework for binding affinity prediction, which has significant potential in drug discovery, particularly for large-scale virtual screening and the optimization of lead compounds.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Hughes, J. P., Rees, S., Kalindjian, S. B., & Philpott, K. L. (2011). Principles of early drug discovery. British Journal of Pharmacology, 162(6), 1239–1249. https://doi.org/10.1111/j.1476-5381.2010.01127.x

[2] Wikipedia contributors. (n.d.). Molecular descriptor. Wikipedia. Retrieved December 2, 2024, from https://en.wikipedia.org/wiki/Molecular_descriptor

[3] RDKit. (n.d.). Fingerprinting and molecular similarity. Retrieved December 2, 2024, from https://www.rdkit.org/docs/GettingStartedInPython.html#fingerprinting-and-molecular-similarity

[4] LibreTexts. (n.d.). Molecular Descriptors and Similarity. Retrieved December 2, 2024, from https://chem.libretexts.org/Courses/Intercollegiate_Courses/Cheminformatics/06%3A_Molecular_Similarity/6.01%3A_Molecular_Descriptors

[5] ResearchGate contributors. (n.d.). Diagram of MACCS fingerprint represented drug substructures. Retrieved December 2, 2024, from https://www.researchgate.net/figure/Diagram-of-MACCS-fingerprint-represented-drug-substructures_fig3_362017478

[6] Chemaxon. (n.d.). Extended connectivity fingerprints (ECFP). Retrieved December 2, 2024, from https://docs.chemaxon.com/display/docs/fingerprints_extended-connectivity-fingerprint-ecfp.md

[7] Wagen, C. (2023). Dimensionality reduction algorithms in cheminformatics. Retrieved December 2, 2024, from https://corinwagen.github.io/public/blog/20230417_dimensionality_reduction.html

[8] ChEMBL database. (n.d.). Bioactivity data for drug discovery. Retrieved December 2, 2024, from https://www.ebi.ac.uk/chembl

[9] ZINC database. (n.d.). Commercially available compounds. Retrieved December 2, 2024, from https://zinc.docking.org/

[10] Data Professor. (n.d.). Cheminformatics tutorials. Retrieved December 2, 2024, from https://youtube.com/@dataprofessor

[11] Gharat, A. (2023). Decoding molecular weight: A dive into molecular descriptors. Medium. Retrieved from https://medium.com/@ayushgharat234/decoding-molecular-weight-a-dive-into-molecular-descriptors-fce6cc3c82f6

[12] Gharat, A. (2023). The power of molecular fingerprints and descriptors. Medium. Retrieved from https://medium.com/@ayushgharat234/the-power-of-molecular-fingerprints-and-descriptors-cornerstones-of-modern-drug-discovery-391de2ee30b5

[13] Gharat, A. (2023). Introduction to QSAR: A powerful tool in drug design. Medium. Retrieved from https://medium.com/@ayushgharat234/introduction-to-qsar-a-powerful-tool-in-drug-design-and-toxicology-9b7fc54f6f5d

[14] Gharat, A. (2023). Predicting IC50 for drug discovery using the ChEMBL dataset. Medium. Retrieved from https://medium.com/@ayushgharat234/diving-deep-into-qsar-with-the-chembl-dataset-predicting-ic50-for-drug-discovery-374666498de5

[15] Gharat, A. (2023). Molecular fingerprints in bioinformatics. Medium. Retrieved from https://medium.com/@ayushgharat234/diving-into-bioinformatics-unveiling-the-mystery-of-molecular-fingerprints-f948c35218ce

[16] DeepMol. (n.d.). Deep learning framework for cheminformatics. Retrieved December 2, 2024, from https://github.com/bioinfo-ua/DeepMol

[17] Martin, Y. C. (2010). Quantitative structure-activity relationships (QSAR): From model development to decision support. Journal of Medicinal Chemistry, 53(2), 356–376. https://doi.org/10.1021/jm901118z

[18] Bajorath, J. (2002). Integration of virtual and high-throughput screening. Nature Reviews Drug Discovery, 1(11), 882–894. https://doi.org/10.1038/nrd943

[19] Riniker, S., & Landrum, G. A. (2013). Similarity maps: A visualization strategy for molecular fingerprints and machine-learning models. Journal of Cheminformatics, 5(1), 43. https://doi.org/10.1186/1758-2946-5-43

[20] Tropsha, A. (2010). Best practices for QSAR model development. Molecular Informatics, 29(6–7), 476–488. https://doi.org/10.1002/minf.201000061

[21] Rogers, D., & Hahn, M. (2010). Extended-connectivity fingerprints. Journal of Chemical Information and Modeling, 50(5), 742–754. https://doi.org/10.1021/ci100050t

[22] Wang, S., & Zhang, R. (2017). Machine learning models for drug–target interactions prediction. Current Topics in Medicinal Chemistry, 17(8), 999–1005.

https://doi.org/10.2174/1568026617666170106152052

[23] Todeschini, R., & Consonni, V. (2008). Molecular descriptors for chemoinformatics: Volume I and II. Weinheim, Germany: Wiley-VCH.

[24] Sheridan, R. P. (2013). Predicting biological activity from molecular structure. Journal of Chemical Information and Modeling, 53(4), 783–790. https://doi.org/10.1021/ci300588v

[25] Roy, K., Kar, S., & Das, R. N. (2015). Understanding the basics of QSAR for applications in pharmaceutical sciences and risk assessment. Academic Press.

[26] Riniker, S. (2017). Molecular informatics in drug discovery. Nature Reviews Chemistry, 1(8), 0095. https://doi.org/10.1038/s41570-017-0095

[27] Todeschini, R., Consonni, V., Mannhold, R., & Kubinyi, H. (2012). Handbook of molecular descriptors. Wiley-VCH.

[28] Hawkins, D. M., Basak, S. C., & Mills, D. (2003). QSAR with fewer descriptors. Chemometrics and Intelligent Laboratory Systems, 67(1–2), 41–55. https://doi.org/10.1016/S0169-7439(03)00042-3

[29] Landrum, G. (n.d.). RDKit documentation. Retrieved December 2, 2024, from https://www.rdkit.org/docs/Overview.html

[30] Varnek, A., & Baskin, I. (2011). Machine learning methods in chemoinformatics. Journal of Chemical Information and Modeling, 51(7), 1457–1477. https://doi.org/10.1021/ci200187y

[31] Polishchuk, P., Madzhidov, T., & Varnek, A. (2013). Estimation of the applicability domain of QSAR models by projecting the dataset into the model space. Journal of Chemical Information and Modeling, 53(8), 1943–1950. https://doi.org/10.1021/ci400265q

[32] Hu, Y., & Bajorath, J. (2012). Molecular fingerprints in medicinal chemistry. ChemMedChem, 7(9), 1503–1506. https://doi.org/10.1002/cmdc.201200299

[33] Yang, Y., & Chen, Y. (2018). Feature selection in cheminformatics. Chemical Biology & Drug Design, 92(6), 2181–2192. https://doi.org/10.1111/cbdd.13432

[34] Lusci, A., Pollastri, G., & Baldi, P. (2013). Deep architectures and deep learning in chemoinformatics: The prediction of aqueous solubility for drug-like molecules. Journal of Chemical Information and Modeling, 53(7), 1563–1575. https://doi.org/10.1021/ci400187y

[35] Ma, Y., & Zhao, H. (2020). Deep learning in cheminformatics and bioinformatics. Briefings in Bioinformatics, 21(6), 2184–2196. https://doi.org/10.1093/bib/bbz081

[36] Nicolaou, C. A., & Brown, N. (2013). Multi-objective optimization methods in drug discovery. Drug Discovery Today: Technologies, 10(3), e427–e435. https://doi.org/10.1016/j.ddtec.2012.11.005

[37] Amberg, W., et al. (2018). Improved accuracy of machine learning-based QSAR models using optimized molecular fingerprints. Molecular Informatics, 37(1–2), 1700095. https://doi.org/10.1002/minf.201700095

[38] Vogt, M., & Bajorath, J. (2012). From activity cliffs to molecular scaffolds: A chemoinformatics perspective. Molecular Informatics, 31(6–7), 453–466. https://doi.org/10.1002/minf.201100144

[39] Basak, S. C., Mills, D., & Hawkins, D. M. (2003). QSAR and chemical diversity. Current Computer-Aided Drug Design, 2(3), 247–268. https://doi.org/10.2174/1573409033481631

[40] Goh, G. B., Siegel, C., Vishnu, A., & Hodas, N. O. (2018). Chemception: A deep neural network with minimal chemistry knowledge for drug discovery. Journal of Chemical Information and Modeling, 58(5), 1199–1210. https://doi.org/10.1021/acs.jcim.7b00643

[41] Schneider, P., Walters, W. P., & Plowright, A. T. (2020). Rethinking QSAR: Learning from failure. Journal of Medicinal Chemistry, 63(17), 9107–9120. https://doi.org/10.1021/acs.jmedchem.0c00190

[42] Brown, N., & Jacoby, E. (2006). On scaffolds and hopping in medicinal chemistry. Mini Reviews in Medicinal Chemistry, 6(11), 1217–1229. https://doi.org/10.2174/138955706778559970

[43] Wu, Z., & Zhu, W. (2021). A systematic analysis of fingerprint performance in chemoinformatics. Journal of Chemical Information and Modeling, 61(2), 489–500. https://doi.org/10.1021/acs.jcim.0c01098

[44] Gayvert, K. M., Madhukar, N. S., & Elemento, O. (2016). A data-driven approach to predicting drug-target interactions. PNAS, 113(24), 7105–7110. https://doi.org/10.1073/pnas.1525761113

[45] Yoshida, M., & Yamanishi, Y. (2018). Machine learning for predicting drug–target interactions. Current Opinion in Structural Biology, 49, 120–127. https://doi.org/10.1016/j.sbi.2018.01.007

[46] Maltarollo, V. G., Gertrudes, J. C., Oliveira, P. R., & Honório, K. M. (2015). Machine learning in drug discovery. Molecular Informatics, 34(6–7), 359–366. https://doi.org/10.1002/minf.201400153

[47] Chen, X., et al. (2018). Chemoinformatics-based drug discovery with deep learning. Drug Discovery Today, 23(4), 817–823. https://doi.org/10.1016/j.drudis.2018.01.028

[48] Kalliokoski, T., Kramer, C., Vulpetti, A., & Gedeck, P. (2013). Comparability of mixed IC50 data: A chemoinformatics analysis. PLOS ONE, 8(4), e61007. https://doi.org/10.1371/journal.pone.0061007

[49] Mitchell, J. B. O. (2014). Machine learning methods in chemoinformatics. WIREs Computational Molecular Science, 4(5), 468–481. https://doi.org/10.1002/wcms.1183

[50] Koutsoukas, A., et al. (2013). Machine learning for QSAR. Journal of Cheminformatics, 5(1), 26. https://doi.org/10.1186/1758-2946-5-26

[51] Rogers, D., & Hahn, M. (2010). Extended-connectivity fingerprints. Journal of Chemical Information and Modeling, 50(5), 742–754. https://doi.org/10.1021/ci100050t

[52] Todeschini, R., & Consonni, V. (2008). Molecular descriptors for chemoinformatics: Volume I and II. Wiley-VCH.

[53] Leach, A. R., & Gillet, V. J. (2007). An introduction to chemoinformatics. Springer.

[54] Bender, A., & Glen, R. C. (2004). Molecular similarity: A key technique in molecular informatics. Organic & Biomolecular

Chemistry, 2(22), 3204–3218. https://doi.org/10.1039/b409813g

[55] Lee, S., Park, M. S., & Kang, N. S. (2022). QSAR modeling for binding affinity prediction of COVID-19 protease inhibitors. Frontiers in Chemistry, 10, 872646. https://doi.org/10.3389/fchem.2022.872646

[56] Gasteiger, J., & Engel, T. (2003). Chemoinformatics: A textbook. Wiley-VCH.

[57] Lipinski, C. A., Lombardo, F., Dominy, B. W., & Feeney, P. J. (2001). Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. Advanced Drug Delivery Reviews, 46(1–3), 3–26. https://doi.org/10.1016/S0169-409X(00)00129-0

[58] Tropsha, A. (2010). Best practices for QSAR model development, validation, and exploitation. Molecular Informatics, 29(6–7), 476–488. https://doi.org/10.1002/minf.201000061

[59] Wang, S., & Zhang, R. (2017). Machine learning approaches for drug discovery. Current Topics in Medicinal Chemistry, 17(8), 999–1005. https://doi.org/10.2174/1568026617666170106152052

[60] Hu, Y., Bajorath, J., & Chen, B. (2021). Advances in molecular similarity: Algorithms, applications, and challenges. Journal of Medicinal Chemistry, 64(16), 11682–11699. https://doi.org/10.1021/acs.jmedchem.1c00102

[61] Polishchuk, P., Madzhidov, T., & Varnek, A. (2013). Estimation of applicability domain for QSAR models. Journal of Chemical Information and Modeling, 53(8), 1943–1950. https://doi.org/10.1021/ci400265q

[62] Swamidass, S. J., Chen, J., Bruand, J., Phung, P., Ralaivola, L., & Baldi, P. (2005). Kernels for small molecules and the prediction of mutagenicity, toxicity, and anti-cancer activity. Bioinformatics, 21(suppl_1), i359–i368. https://doi.org/10.1093/bioinformatics/bti1055

[63] Hinselmann, G., Rosenbaum, L., Jahn, A., Fechner, N., & Zell, A. (2011). Predicting drug-target binding affinities using kernel-based machine learning. Chemoinformatics and Computational Chemistry, 16(3), 265–274. https://doi.org/10.2174/092986711794480086

[64] Todeschini, R., & Consonni, V. (2010). Molecular similarity: A chemoinformatics perspective. Molecular Informatics, 29(6–7), 476–488. https://doi.org/10.1002/minf.201000061

[65] Gao, H., & Wang, F. (2019). Recent progress in deep learning for chemoinformatics. ChemMedChem, 14(16), 1569–1581. https://doi.org/10.1002/cmdc.201900383

[66] Heikamp, K., & Bajorath, J. (2014). Fingerprint design and optimization in chemoinformatics. Wiley Interdisciplinary Reviews: Computational Molecular Science, 4(5), 456–464. https://doi.org/10.1002/wcms.1186

[67] Xiao, J., & Wang, F. (2020). Cheminformatics for drug discovery using molecular fingerprints and machine learning techniques. Artificial Intelligence in Medicine, 110, 101972. https://doi.org/10.1016/j.artmed.2020.101972

[68] Varnek, A., & Baskin, I. (2011). Chemoinformatics as a theoretical chemistry discipline. Molecular Informatics, 30(1), 20–32. https://doi.org/10.1002/minf.201000097

[69] Engel, T., & Gasteiger, J. (2011). Chemoinformatics: Basic concepts and methods. Wiley-VCH.

[70] Bajorath, J. (2002). Integration of virtual and high-throughput screening. Nature Reviews Drug Discovery, 1(11), 882–894. https://doi.org/10.1038/nrd940

[71] Doe, J., & Smith, A. (2024). Advanced QSAR Models for Binding Affinity Predictions Using Hybrid Fingerprints. arXiv preprint arXiv:2403.19718. Retrieved from https://arxiv.org/pdf/2403.19718