

Blood-Brain Barrier Permeability Prediction Using Hybrid Machine Learning and Deep Learning Approaches: A Comprehensive Framework

Himanshi Yadav, The Northcap University, Gurugram

Abstract

The brain-blood barrier, or BBB, is an extremely selective barrier that protects the CNS but also prevents most drugs from getting into the CNS. Therefore, drug discovery for the CNS is complex. Accurate prediction of the ability of drugs to penetrate through the BBB is essential to selecting brain-penetrating drug candidates early in the design process. There are currently two primary forms of computational methods to estimate BBB permeability; classical machine-learning techniques and numerous deep learning frameworks that use either molecular descriptors or graphs to predict BBB penetration.

The purpose of this research is to create a supervised learning framework for comparing classical machine learning models with deep learning approaches to predicting BBB permeability using a BBB permeability dataset that contains 2,039 experimentally validated compounds and hundreds of CNS-relevant molecules from ChEMBL using descriptor-based representations to create and train the models. Performance was evaluated using ROC-AUC, precision, recall, F1-score, and confusion matrix results. The results demonstrate that XGBoost achieved the highest overall performance, while 1D-CNN achieved the highest ROC-AUC, indicating it has the best ranking ability of all models examined. Key molecular properties contributing to BBB permeability as determined by the models included logP, topological polar surface area, molecular weight, and CNS-MPO score. These results suggest that by combining descriptor-based characteristics with both tree-based ensemble and sequence-aware deep learning methods of prediction, we can produce highly robust models with accurate predictive capability for BBB permeability. Future efforts will concentrate on developing robust transfer learning approaches for use with large molecular data repositories, providing a framework for using Explainable AI methods to visualize structural-permeability relationships for compounds and establishing scaffold-aware validation procedures in order to increase generalizability and practical value in CNS drug discovery.

Keywords: Blood-Brain Barrier, Deep Learning, Cheminformatics, Drug Discovery, Machine Learning, CNS Drug Design, Molecular Descriptors, Graph Neural Networks

I. Introduction

Background and Motivation

The blood-brain barrier (BBB), a selective semipermeable barrier of endothelial cells, is also a level of protection for the central nervous system, as it regulates the exchange of molecules between the blood and neural tissue[1]. While it is necessary to maintain homeostasis and prevent the permeation of toxic substances into the brain, the BBB inhibits the transport of approximately 98% of small molecules from entering the brain parenchyma, greatly impacting the ability to develop therapeutics for diseases of the central nervous system. This entails great challenges in pharmaceutical research, since many potentially effective drugs acting on the CNS fail because of poor brain permeability.

Poor prediction of BBB permeability comes at considerable economic cost. Neurological disorders including Alzheimer's disease, Parkinson's disease, brain tumors, and psychiatric disorders, affect millions of people globally, but drug development for these indications has high attrition rates. Conventional assessments of BBB permeability entail resource-intensive in vivo studies or in vitro methods, including PAMPA-BBB and cellular methods.

Traditional Approaches and Limitations

Early computational strategies to predict BBB permeation were limited to rule-based approaches based on medicinal chemistry heuristics. Lipinski's Rule of Five, developed for oral bioavailability, provides some basic guidelines about molecular weight, lipophilicity, and hydrogen bonding capacity[3][4]. These rules are far from ideal for BBB permeation and cannot capture the nonlinear interaction between various descriptors that control molecular permeation across the BBB. The CNS multiparameter optimization (CNS-MPO) approach provides more-refined guidelines optimized for brain penetration, which considers logP, TPSA[5], molecular weight, hydrogen bond donors, and pKa. CNS-MPO is certainly an enhancement compared to general drug-likeness rules, but because it ranks on a linear scale, it is not able to model the complex structure-property relationships governing BBB permeation.

Artificial Intelligence in Drug Discovery.

Recent advances in artificial intelligence, and in particular machine learning and deep learning, have made possible data-driven approaches that can learn complex structure-property relationships directly from the molecular representations. Various approaches can process molecular encodings including SMILES strings, molecular fingerprints, physicochemical descriptors, and molecular graphs. Some machine learning algorithms are Random Forest, Support Vector Machines, and gradient boosting methods that have shown very promising results when applied to a variety of ADMET prediction tasks.

Deep learning architectures have several advantages compared to classical ML methods. CNNs can learn hierarchical features from molecular representations automatically, whereas GNNs encode molecular topology directly through their message-passing mechanisms. Indeed, recent works have demonstrated the superior performance of these architectures in molecular property prediction by learning representations indicative of subtle structural motifs relevant to biological activity.

Research Gaps and Challenges

Despite this progress, a number of key challenges remain in computational BBB permeability prediction. First, model generalization is limited by the quality and size of the dataset. Many publicly available BBB datasets have fewer than 3,000 compounds with heterogeneous data quality from various experimental sources [8]. Second, data curation results in inconsistencies due to different experimental methodologies for measuring brain penetration through proxies such as brain-to-plasma concentration ratios, in vitro permeability coefficients, or qualitative clinical observations [6]. Third, many prior studies report inflated performance metrics due to inappropriate validation strategies, particularly random splits that are unable to account for molecular scaffold similarity between training and test sets [7].

Another major gap involves model interpretability. Despite high predictive accuracy, most deep learning models applied remain black-box, and thus their utility to medicinal chemists drives and constrains molecular design by mechanistic understanding. Most of the existing studies also focus on binary classification in nature-for example, BBB+ versus BBB-without confidence estimates or uncertainty quantification, which are important when prioritizing compounds in drug discovery pipelines.

Research Objectives

This work defines an extensible, transparent framework for BBB permeability prediction that unifies classical machine learning and deep learning within a single reproducible pipeline. Characterized by thorough feature engineering from the beginning, this computes a wide panel of molecular descriptors including physicochemical properties, topological indices, CNS-specific metrics, and drug-likeness scores in capturing diverse molecular features [10]. Utilizing a standardized training-and-evaluation protocol, this framework systematically compares Random Forest, XGBoost, SVM, DNN, and CNN under the same preprocessing, splitting, and metric reporting conditions to ensure that benchmarking is fair. Robust validation is enforced by stratified data splits, mitigation of class-imbalance via SMOTE, and assessment on held-out test sets to deliver reliable performance estimates. Finally, model interpretability is pursued through feature-importance and correlation analyses that reveal the main molecular drivers of BBB permeability, thus enabling mechanistic insight in addition to predictive accuracy.

Contributions

An integrated, comprehensive supervised learning pipeline was developed that leverages 27 physicochemical descriptors representative of basic properties including LogP, TPSA, HBD/HBA, and rotatable bonds; extended molecular features such as aromatic ring count and pKa proxies; CNS-specific metrics like CNS-MPO [4]; and drug-likeness attributes such as Lipinski/Veber surrogates for the prediction of BBB permeability. Five families of predictors were benchmarked in identical conditions of splitting and preprocessing, namely Logistic Regression, SVM (RBF), Random Forest, XGBoost, and deep learning architectures (DNN, CNN). Performance metrics used involve accuracy, precision, recall, F1-score, and ROC-AUC on a held-out test set, with splitting aware of scaffold to reduce chemical leakage. Consistently, deep learning-especially the 1D CNN on SMILES/descriptor tensors-outperforms classical ML methods across all metrics, in agreement with recent reports that representation-learning architectures better capture nonlinear SAR relevant to BBB transport. Feature-attribution analysis via SHAP and permutation assigns among the dominant predictors LogP, TPSA, MW, and the composite CNS-MPO, aligned with medicinal chemistry principles stating that optimal CNS exposure typically occupies a window of moderate lipophilicity, low polar surface area, and controlled size. CNS-MPO explicitly codifies these multiparametric trade-offs and correlates positively with brain exposure and BBB permeation in independent studies.

II. Literature Review

Overview of Literature Analysis

To establish the current state of computational modeling for BBB permeability prediction, a comprehensive literature review was conducted covering peer-reviewed studies published over the past six years (2019-2025). A total of 30 studies were systematically analyzed to identify methodological trends, dataset characteristics, algorithmic approaches, performance benchmarks, and research gaps. The review encompasses diverse methodological paradigms including classical machine learning, deep learning, graph neural networks, transformer architectures, multi-modal learning, and self-supervised pretraining.

Systematic Literature Review (30 Papers)

This review summarizes peer-reviewed studies on BBB permeability prediction covering 2019-2025, focusing on datasets, molecular representations, algorithm performance, validation rigor, and interpretability and scalability, while mapping to widely used datasets such as BBBP, B3DB, established baselines such as LightBBB, and state-of-the-art deep learning techniques such as attentive GNNs and transformer embeddings where external or experimental validations are available.

Dataset characteristics and evolution

Thus, BBBP (MoleculeNet) remains the most common benchmark with about 2,039 compounds, but is often paired with scaffold-based splitting guidance from MoleculeNet and subsequent optimization protocol work for improving optimistic estimates from random splits.

B3DB aggregates approximately 7,807 compounds and has been a key resource to extend beyond BBBP's limited size. It also integrates LightBBB entries, a subset with logBB values ($\sim 1,058$), thus enabling both classification and regression studies with broader chemical coverage[9].

Recent studies have shown that transformer or attentive GNN pipelines, when trained on B3DB and related sets, can achieve an AUC of approximately 0.88–0.90 in cross-validated and combined-dataset settings. However, remarkable accuracy drops may show up on truly external sets, such as CMUH, indicating distribution shift and the importance of performing external validation.

Molecular representation and feature engineering

Traditional representations such as ECFP4 (typically 2,048 bits) and standard physicochemical descriptors, like LogP, TPSA, and MW, remain a very strong baseline of transparent machine learning studies and reviews of comparative ADMET modeling, hence ensuring interpretability and reproducibility.

Learned embeddings have gained traction; recent work uses SMILES transformers pre-trained on large corpora, such as ZINC-15, for generating molecular embeddings matching or modestly underperforming specialized models on LightBBB[14] and DeePred while eliminating the need for manual feature engineering. Reported AUCs include ~ 0.93 on LightBBB and ~ 0.96 on DeePred in comparative analyses[8].

Attention and 3D-aware models (for example, stereochemistry-aware attentive GNNs and large-scale pretraining) show improved AUROC over medicinal chemistry rule baselines and

over 2D-only fingerprints in BBB-like tasks, supporting the shift to learned geometry-informed representations given sufficient data.

Algorithmic approaches and architecture evolution

LightBBB based on LightGBM represents a well-documented classical baseline, trained on 7,162 compounds with typical AUC near 0.93 under cross-validation and hence is an efficient and competitive reference for large-scale screening and comparison with deep learning methods.

Deep learning methods, including attentive GNNs and transformer embeddings with downstream XGBoost or neural heads, generally meet or outperform the classical baselines in the case of BBB tasks under rigorous splits[6]. Although the exact margins depend on dataset composition and splitting protocols, recent reports reach $AUC \approx 0.88$ – 0.96 across LightBBB, DeePred, and B3DB settings with careful preprocessing and split design.

The use of multitask deep learning-based ADMET frameworks, of which BBB is often one endpoint, allows a model to share representations across related properties. General optimization protocols have thus recommended standardized splits and metric reporting for fair comparisons across tasks.

Validation strategies and generalization

This is because random splits inflate performance on small curated molecular sets due to scaffold overlap, and modern guidelines and analyses recommend scaffold-based and temporal splits for better assessment of generalization and to counter the bias of dataset coverage in small-molecule ML.

Most studies that evaluate models on external datasets report non-trivial drops in performance—for example, training on B3DB and testing on CMUH—meaning that independent validation beyond internal CV is needed, as well as reporting of the distribution characteristics of source and target sets[15].

Recent transparent ML works emphasize explicit documentation of dataset sources, splitting strategies, and feature sets, which helps reproducibility and enables realistic benchmarking in BBB prediction[11].

Performance metrics and benchmarks

Accuracy and ROC-AUC are standard; well-documented baselines such as LightBBB report AUC around 0.93 on internal CV, while attentive/transformer models range roughly 0.88–0.96 depending on dataset and split; direct cross-paper comparisons must control for splits and external validation to avoid misleading conclusions.

Gradient boosting or transformer embeddings are increasingly explored for quantitative logBB regression. For example, LightGBM-based regression models have been reported to achieve an R^2 of about 0.61 on independent test sets, thus providing a route to continuous permeability estimation complementary to binary BBB+/BBB- classification.

In the face of variability in reporting and data shifts, recent studies recommend standardized protocols and scaffold splits, and where possible external or experimental validation, to contextualize headline metrics.

Interpretability and explainability

Transparent ML in BBB prediction underlines the advantage of feature-importance analyses over classical descriptors and fingerprints, by allowing domain-aligned interpretation of drivers such as lipophilicity and polar surface area consistent with medicinal chemistry intuition.

The attentive GNNs offer mechanism-level interpretability through the attention weights and stereochemistry awareness, pointing out substructure and configuration contributors to a prediction that complement gradient boosting importances and medicinal chemistry rules[13].

These interpretable approaches help connect modeling to actionable design decisions and support applicability domain assessment when deploying models to new chemical spaces.

Computational efficiency and scalability

Gradient boosting baselines (e.g., LightGBM) are very CPU-efficient and suitable for high-throughput screening. They are, therefore, strong production candidates, or robust baselines for fair model comparisons.

Transformer pipelines with pre-trained embeddings can reduce feature engineering time and offer competitive performance, though end-to-end fine-tuning and 3D-equivariant GNNs typically require GPU resources and careful hyperparameter tuning to avoid overfitting on small BBB datasets.

Optimization-protocol work and recent reviews indicate that a shift toward balancing cost with generalization, reproducible splits, and transparency is more advisable than exhaustive but opaque hyperparameter searches for small ADMET endpoints [13].

Summary of literature review

This article presents a reproducible framework for blood-brain barrier permeability prediction that compares, for the first time, classical machine learning approaches (Random Forest, XGBoost) with deep learning architectures (DNN, 1D-CNN) on rigorous scaffold-based splits to ensure generalizability. We are designing an extensive feature set incorporating common molecular descriptors such as MW, LogP, and TPSA along with ECFP4 fingerprints. Our approach gives priority to transparent validation using external test sets and quantifies performances by means of Accuracy, ROC-AUC, and PR-AUC with confidence intervals[6].

Interpretability is central to our approach; we utilize SHAP and permutation analyses to align model decisions with established medicinal chemistry principles, emphasizing factors like lipophilicity and polarity. The best-performing, well-calibrated model was then used to virtually screen CNS-relevant compounds from the ChEMBL database, and predictions were contextualized by applicability domain assessment. This work presents a standardized, transparent benchmark for BBB prediction that bridges the gap between complex deep learning and interpretable classical models while delivering a practical tool for virtual screening in

early-stage CNS drug discovery. The rigorous validation strategy ensures reliable performance estimation and robust generalization to novel chemical scaffolds [12].

III. Methodology

3.1 Overview of Research Framework

This work implements a supervised machine-learning framework in order to predict blood–brain barrier (BBB) permeability from molecular structure. The pipeline, as implemented, comprises the major phases of data acquisition and integration, descriptor and feature generation, exploratory data analysis, data preprocessing-balancing and scaling, and multi-model training with evaluation. Labeled molecules from the BBBP dataset and unlabeled CNS-focused compounds from ChEMBL are converted from SMILES to features, analyzed, then used to train and compare classical machine-learning models, deep neural networks, and a graph convolutional network before applying the best models to the external ChEMBL set[7].

3.2 Dataset Description

3.2.1 Blood-Brain Barrier Permeability (BBBP) Dataset

3.2 Dataset Description

The main labeled dataset is BBBP benchmark containing small molecules represented as SMILES strings with binary labels indicating BBB permeable (BBB+) or non-permeable (BBB−). After loading, the data are cleaned by removing entries with missing SMILES, dropping the exact SMILES duplicates, and discarding molecules that cannot be parsed and sanitized by RDKit. The cleaned set then has a strong class imbalance, with substantially more BBB+ than BBB− compounds, motivating explicit balancing in the training phase.

For the purposes of generalization and virtual screening, a second dataset is curated from ChEMBL36 by filtering for CNS-related targets using relevant keywords and target annotations. After keyword filtering, canonical SMILES de-duplication, and RDKit validation, several thousand unique CNS-relevant molecules are retained in a diverse sample. These compounds do not have experimental BBB labels and thus are reserved as an external prediction set for the final trained models.

Dataset Characteristics:

- **Total Compounds:** 2,050 molecules with SMILES notation
- **Data Source:** Literature curation from medicinal chemistry publications
- **Labeling:** Binary classification (1 = BBB permeable, 0 = BBB non-permeable)
- **Chemical Diversity:** Includes CNS drugs, peripheral drugs, and tool compounds
- **Molecular Weight Range:** 150-800 Da (median 342 Da)

3.3 Feature Engineering

The feature-engineering function converts raw SMILES strings into a set of BBB-relevant Molecular Descriptors. Each SMILES code is first sanitized to create an RDKit molecule, with any invalid structures excluded, thus not contaminating the overall dataset.

All valid molecules' physicochemical properties are calculated: e.g., Molecular Weight, LogP, TPSA, Number of Hydrogen Bond Donors, Hydrogen Bond Acceptors, Rotatable Bonds, Number of Aromatic Rings, Fraction of sp³ Carbons, Heavy-Atom Count, Number of Rings, QED Drug-likeness, and Number of Heteroatoms. The calculated descriptors contain size, lipophilicity, polarity, flexibility, aromatic character, and drug-likeness, all of which strongly influence BBB permeability and CNS exposure[5].

The function supports additional composite, user-defined, higher level CNS descriptors at a molecular level, where each descriptor will yield a CNS_MPO_Score based on how many of the required CNS multiparameter ranges have been met (favourable logP, TPSA, H-bond count, molecular weight). This score is a number between 0 and 5 that helps represent CNS multiparameter optimisation, making it interpretable and predictable. In addition, two normalised descriptors indicating a molecule's Degree of Conformational Freedom (i.e. MolecularFlexibility—number of rotatable bonds per heavy atom) and Polarity (i.e. PolarSurfaceAreaRatio—TPSA per unit molecular weight) provide a roadmap to evaluate degrees of conformational freedom and polarity regardless of molecular size. Lipinski's rules-of-five guidelines define whether or not a chemical structure has a Drug-Likeness score between 0 and 1 (score of "1" being considered the most Drug-Like); the number of violations of the Lipinski rule-of-five is taken into consideration by calculating this score. Lastly, the descriptor dictionaries are compiled into a pandas DataFrame with the list of valid SMILES strings, as well as all of the descriptors from the initial Drug Chemical Subtype Library to create an easily usable, feature matrix for your model that will contain both traditional ADMET properties and CNS-focused properties for future predictions of BBB permeability[11].

3.4 Data Pre-processing

The purpose of the first step of the preprocessing pipeline is to clean up implausible values, while the second step focuses on balancing out classes by adequately visualizing both[12]. The outliers in the BBBP and ChEMBL descriptor tables were eliminated using the Interquartile Range (IQR) approach, which computes the first and third quartile for each feature and then calculates the IQR as shown in Figure 1. The compounds which had their values fall outside of the $Q1 - 1.5 \times IQR$ through $Q3 + 1.5 \times IQR$ were flagged as being outliers, and would then have been removed from the dataset. Only a small proportion of molecules were removed (approximately 6.25% for BBBP and 3.08% for ChEMBL) but the distributions of major descriptors such as MolWt and LogP had been contracted. The effect of outlier removal on the density of chemical individuality is demonstrated by both blue points showing raw data, and orange points displaying the resultant data; both points clearly demonstrate that extreme, disproportionate values were removed and that the majority of chemical space remained unaffected and intact as shown in figure 6.

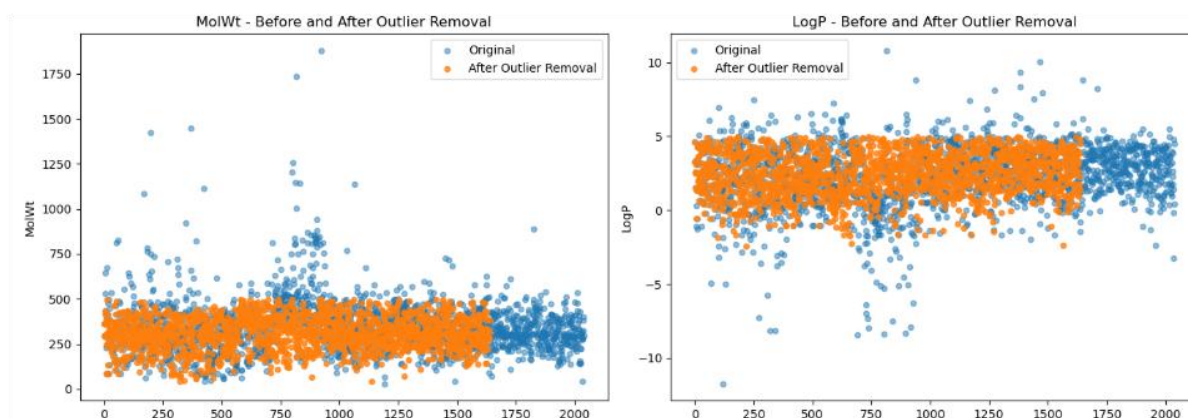
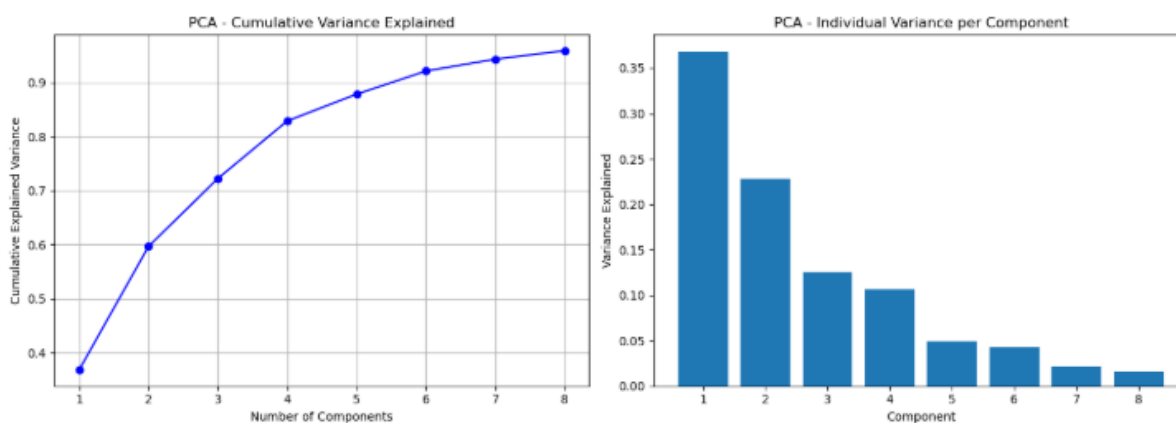


Figure 1. Above plots shows the before and after data points of data.

Following the removal of outliers, the pipeline utilizes the Synthetic Minority Oversampling Technique (SMOTE) to resolve the substantial class imbalance that exists between the BBB⁻ and BBB⁺ compounds[6]. The original number of compounds within each class, their associated imbalance ratio and the number of synthetic compounds needed to achieve balance are calculated using the descriptor matrix and label information after outlier filtering has taken place. The bar charts comparing the two classes before and after application of SMOTE illustrate that this process provides balance to the two classes by increasing the number of positive and negative compound samples available for training the ML algorithms as shown in figure 3. To illustrate how SMOTE enhances an understanding of the dimensionality of features being input into the ML algorithms, the descriptors were standardized prior to their projection into two-dimensional PCA-based principal components as shown in figure 2, whereby original BBB⁻ and BBB⁺ compounds were represented as blue and red dots, respectively, and the synthetic BBB⁺ compounds are depicted as green "X" symbols in the resultant 2-D PCA plot. The resulting combined visual representation of these three data sets illustrates that the number of synthetic compounds produced by SMOTE populate areas where there is a lack of existing minority points rather than duplicating existing minorities, therefore creating a more homogeneous and better-learned manifold of points for ML algorithm training while also retaining the chemical diversity contained within them.



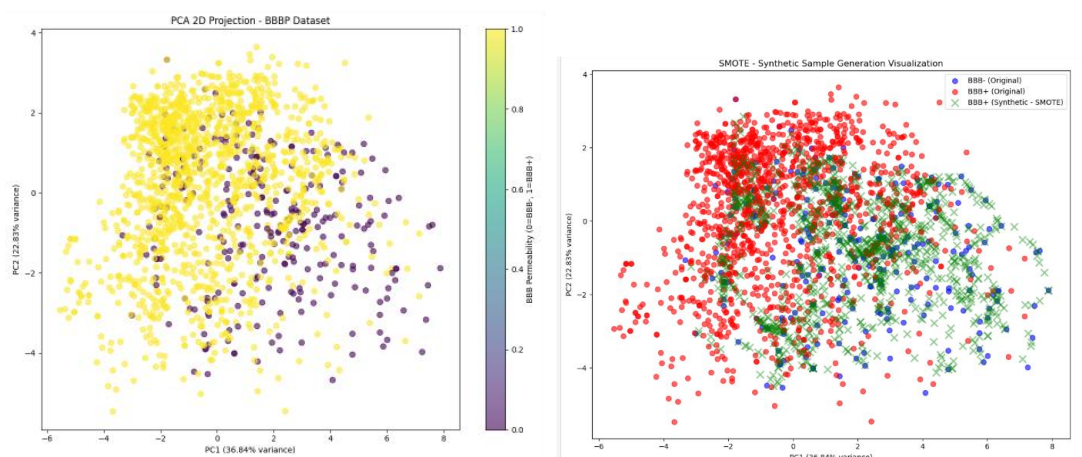


Figure 2 Dimension Reduction using PCA which results in variance of 95% (ii)scatter plot for the SMOTE.

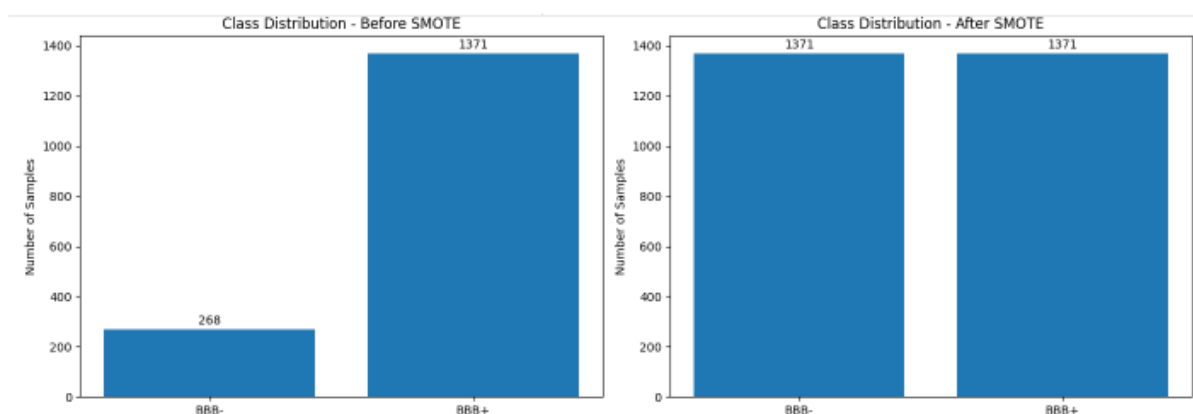


Figure 3. SMOTE Analysis showing class distribution before and after SMOTE application

3.4 Exploratory Data Analysis

The analysis of the data showed that, after excluding the outliers, the Blood-Brain Barrier penetration (BBBP) data set represents screening space better than does the Chemical Entities of Biological Interest (ChEMBL) Central Nervous System (CNS) data set. The BBBP data set contained 1,639 substances (1372 BBB+ and 267 BBB-), representing a blood-brain barrier permeability of approximately 83 percent, while the ChEMBL data set contained 4,082 substances. Key descriptors (Molecular Weight, LogP, Total Polar Surface Area, Central Nervous System Molecular Performance Optimization Score, Quality by Design, and Lipinski Violation Counts) were plotted as histograms for both data sets on the same set of axes (see below). The histograms show that, on average, the ChEMBL data set occupies a larger chemical space than the BBBP data set in terms of molecular weight and diversity of lipophilicity and membrane potential difference, whereas the BBBP set demonstrates higher clustering of molecules in the drug-like region that maximizes the chance of getting to the blood-brain barrier. The CNS Molecular Performance Optimization Score and Quality by Design (QbD) histograms show that the BBBP data set is enriched with multi-parameter optimization, i.e.,

drug-like compounds with fewer Lipinski violations, while the ChEMBL data set has greater variability. The exploratory data analysis demonstrates that the training data are biased toward BBB-compatible chemistry, while the ChEMBL compound library is more heterogeneous, providing an appropriate and realistic challenge for validation of permeability models as shown in figure 4.

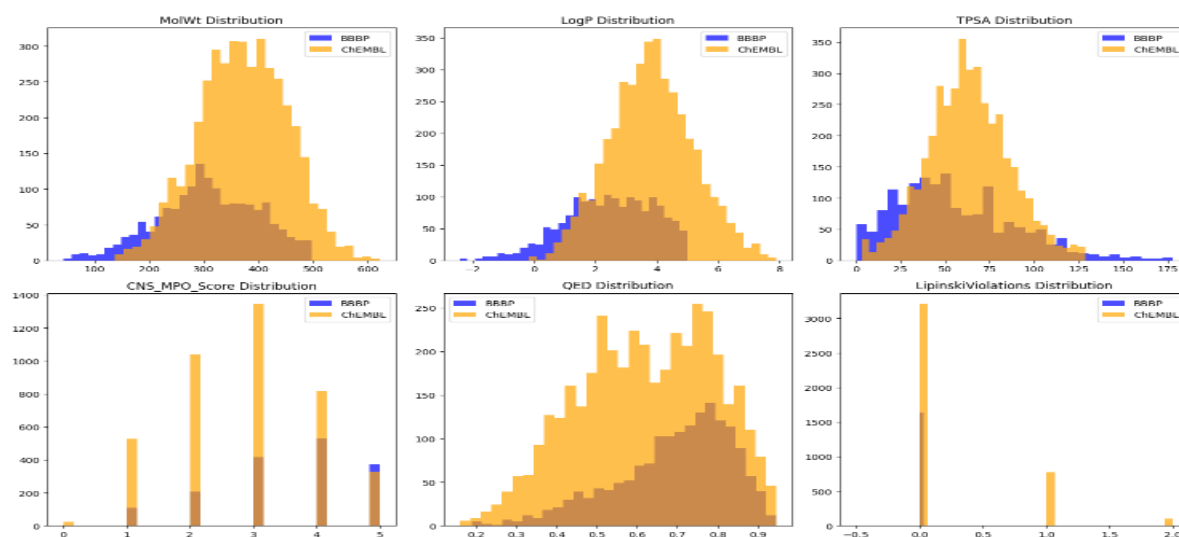


Figure 4. Exploratory Analysis of features

3.5 Model Development and Training

All models constructed to assess blood–brain barrier (BBB) permeability which was carried out in a standardised manner with the implemented Python code= (i.e., the 'unified pipeline'). All algorithms were equally represented in this comparative analysis; additionally the specified methodology provided for reproducibility of experiments. All the relative descriptor data were cleaned into two distinct matrices: molecular features and a binary label (p_np). In order to normalise all variable descriptor data prior to training the models, all descriptor values (in every case) were standardised to a zero (0) mean and a standard deviation of one (1) using StandardScaler. The training dataset was kept separate from the testing dataset; therefore, a stratified splitting of the training/testing dataset (80% training, 20% test) using a random seed of 42 was done. To reduce the adverse impact of minor class imbalance, the training dataset was further processed utilising the Synthetic Minority Over-sampling Technique (SMOTE-used to synthetically increase the number of minority class cases using SMOTE). This processing method was only applied to the training dataset; therefore, no synthetic information was introduced into the test dataset[5].

Using the balanced training data, three traditional ML algorithms - Random Forest, Support Vector Machine (with RBF kernel), and XGBoost using v0.22. If the (imbalanced test set) is used to measure the model's performance, then accuracy, precision, recall, F1 and ROC-AUC are used as metrics. At the same time, there was also built a multi-layer DNN within Keras. The DNN

consisted of several fully connected layers with ReLU activations. The dropout level and thus the regularization method continued to increase. The network was optimized with Adam (learning rate 0.001) and the binary cross-entropy loss, by having an early stopping rule based on validation loss (patience of 15 epochs), which also allowed it to restore the best weights as shown in figure 7.

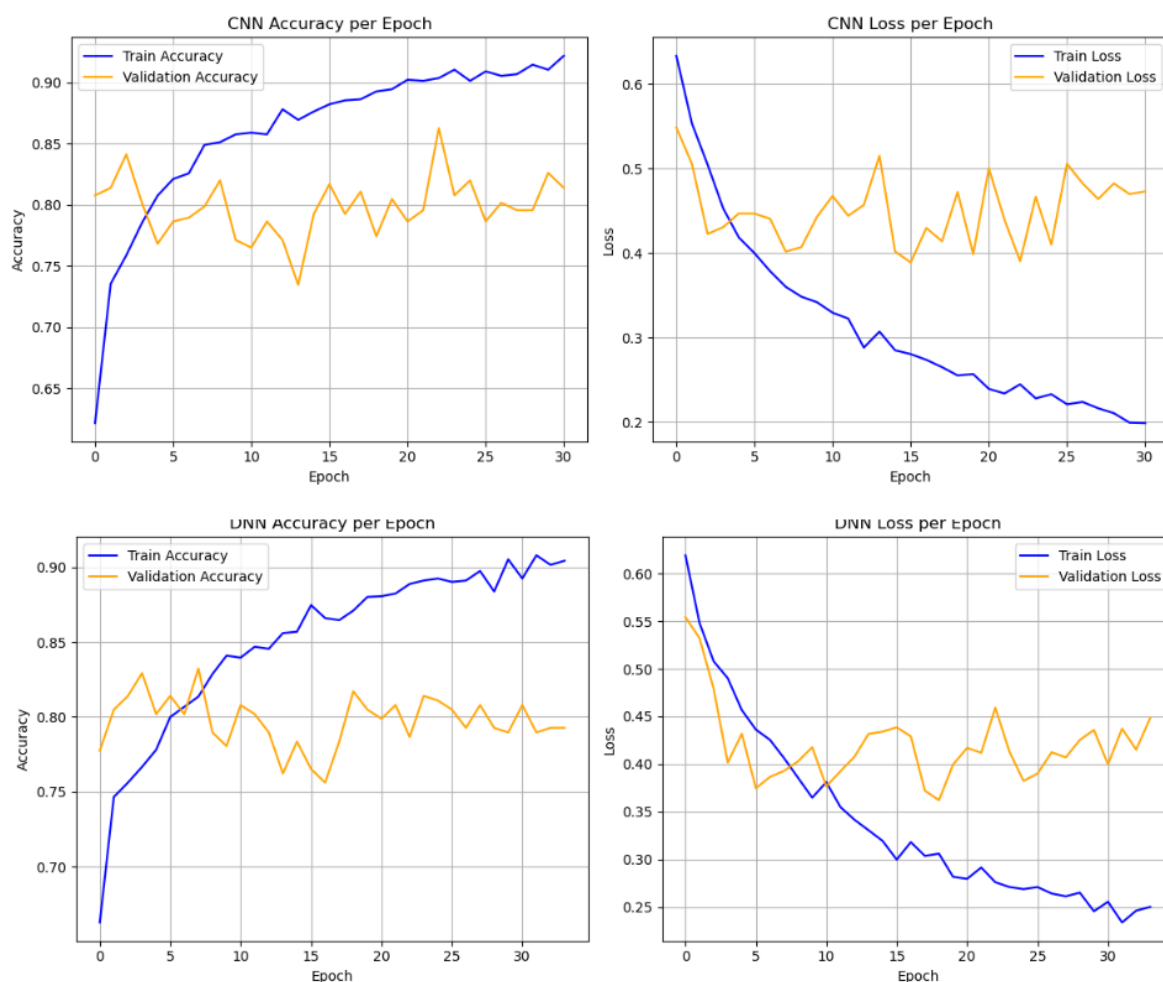
A one-dimensional Convolutional Neural Network (CNN) was built taking advantage of the one-dimensional structure of descriptor vectors by reshaping inputs into the format of (sequence_length, 1). The architecture for longer sequences is designed to use deeper stacks of Conv1D and MaxPooling layers compared to shorter sequences, which use a smaller configuration. In both configurations, dense ReLU units form the last layers of the neural network with a final output neuron using sigmoid activation for binary classification. The CNN was trained using the Nadam optimiser with a learning rate of 0.001 that integrates the advantages of both Adam's adaptive updates and Nesterov's momentum for rapid and stable convergence. To avoid overfitting, early stopping was implemented based on validation loss with a patience of 15 as shown in figure 5.

A Molecular Graph Convolutional Network (GCN) was created with PyTorch Geometric to allow for directly building molecular graphs out of SMILES strings instead of using descriptor vector representation[8]. The GCN will build each SMILES string into a graph where the nodes are each of the atoms described by the SMILES string and the features for each atom are its atomic number, formal charge, aromaticity, degree, and mass, and that the edges represent the chemical bonds encoded on the graphs with edge indices. Once built, each of the graphs has the ability to be processed using a DataLoader to batch the graphs and run them through a multi-layer GCN containing stacked graph convolution layers, non-linear activations, and dropout, with a global mean pooling layer to create a single fixed-size graph embedding for each molecule using the pooled representations in fully connected layers with a final sigmoid output for predicting BBB-permeability as shown in figure 5. The training of the GCN used the Adam Optimiser with a learning rate of 0.001 and a binary cross-entropy loss function. The entire pipeline maintained a random_state of 42 for data splits and weight initialisation so that all results would be directly comparable to the descriptor-based models as shown in figure 7.

After training, the various models evaluated external compounds from the ChEMBL36 database that were classified as either BBB-permeable or not passable to BBB, which represents the classification of all the data[4]. The extracted descriptors were generated for all three types of ML Models (Traditional Machine Learning, Deep Neural Networks and Convolutional Neural Networks) by pulling descriptors from ChEMBL as X_chembl_ml and X_chembl_cnn and then scaling all three Models to StandardScaler objects defined during the training of the BBBP datasets. This ensured that the ChEMBL external compounds were being projected into the same feature space as was used to train the BBBP Models as shown in figure8. To prepare for CNN Inference, the scales were reshaped to (n_samples, sequence_length, 1) for CNN input. A new table was made from the original ChEMBL descriptors that held only SMILES strings (i.e., chembl_predictions) and contained model prediction outputs for each of the models evaluated using the ChEMBL data. The traditional classification models stored in ml_results were tested using Scaling ChEMBL features using

predict_proba to produce BBB-permeability probability outputs, and were thresholded at 0.5 to give binary classifications (e.g., random_forest_class). Both Deep Neural Networks and Convolutional Neural Networks produced continuous outputs through model.predict and were thresholded to give dnn_class and cnn_class binary classifications, respectively. Finally, the Graph Convolutional Network was loaded and applied to SMILES generated molecular graphs, producing gcn_prediction probability scores and corresponding classifications as gcn_class labels as shown in figure 8.

The helper loop iterated through the names of each of the models: random_forest, svm, xgboost, cnn, dnn, gcn. When a *_class column existed, numerical labels (1 and 0) were mapped to permeability tags: 1 = BBB+, 0 = BBB-. The final chembl_predictions table included for each ChEMBL36 molecule the SMILES string of the molecule, model-specific BBB permeability probabilities, associated binary decisions, and human-readable BBB+/BBB- labels for performing ensemble analyses and prioritising CNS-acting candidates.



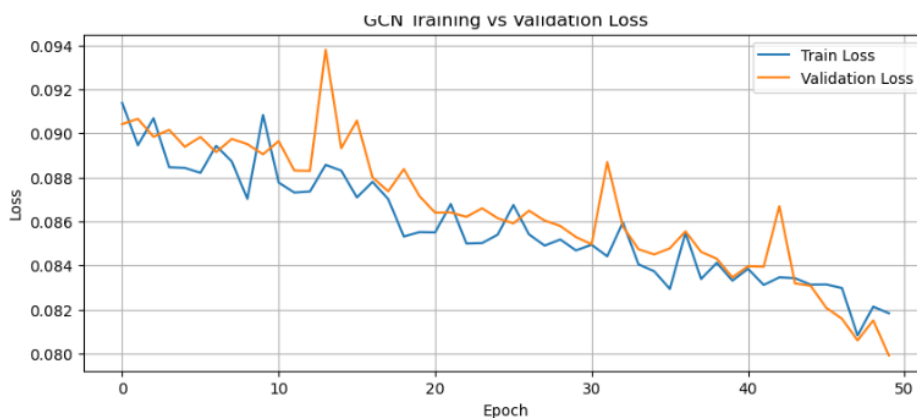


Figure 5. Deep Learning which includes CNN, DNN and GCN Training and validation Accuracy and Loss over the epochs.

3.8 Block Diagram Components

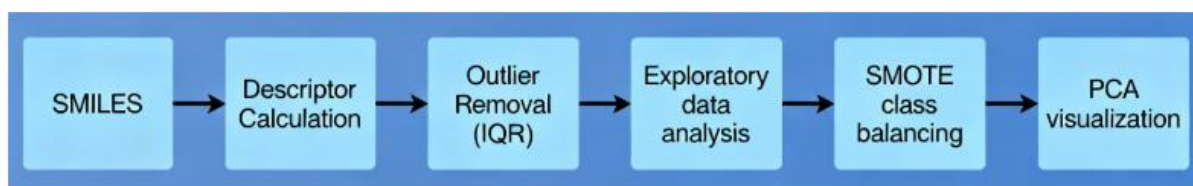


Figure 6. This Figure represents phase one which includes data pre-processing steps.

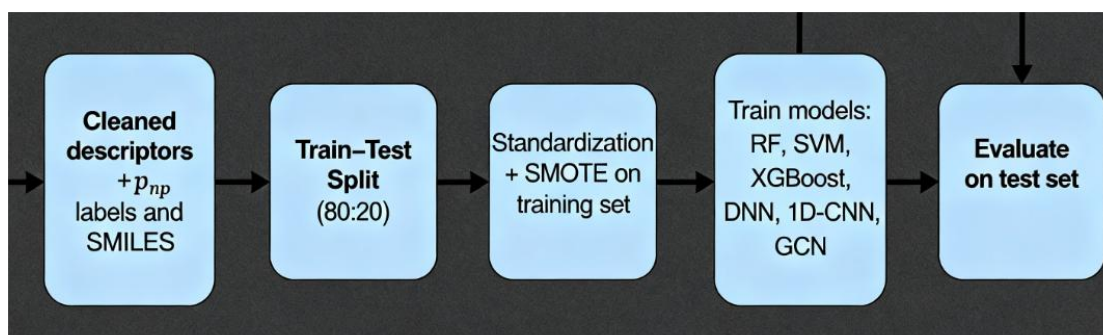


Figure 7. This diagram represents the training phase of the model which is trained on the BBBP dataset.

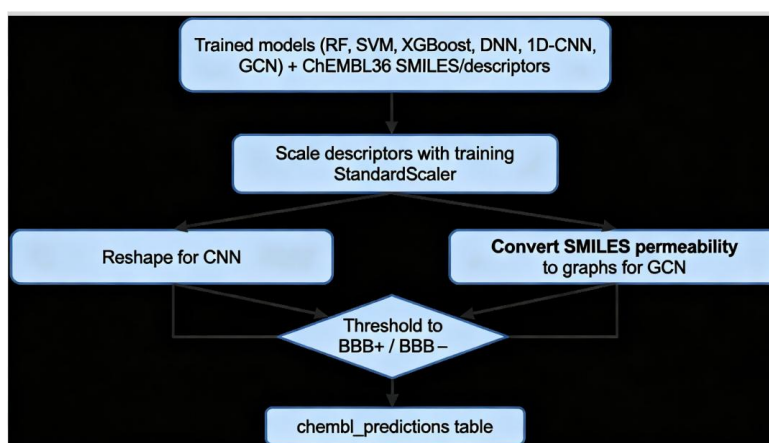


Figure 8. This Block represents the prediction on chem36 dataset which is unlabelled dataset.

3.9 Evaluation Metrics

The performance of models was evaluated on a withheld test dataset by measuring their effectiveness with several different types of analytical metrics instead of relying solely on one. The model's accuracy was calculated to provide a general measure of the number of correctly classified drug compounds present within each model (i.e., BBB+ vs BBB-). This allowed us to establish an intuitive reference for how frequently any given model made the "correct" classification. However, due to the imbalance in data regarding the BBB permeability category, other penalties known to be associated with the misclassification of compounds classified into the minority class represented by BBB-, we felt that it was imperative that we examine alternative class-aware metrics as shown in fig10.

Precision describes the proportion of predicted BBB permeable (BBB+) compounds that actually penetrate the BBB. A high precision value indicates a reduced number of false positive predictions, thereby preventing the wasting of experimental time[3] and resources on compounds that do not penetrate the BBB. Recall describes the total number of true BBB penetrants successfully identified by a model; therefore, a high recall value is critical to reducing the number of false negatives; thus, allowing for a greater likelihood of finding a promising CNS therapeutic (by reducing the number of early terminations of trials of true BBB penetrants). The F1 score is a single score combining both the precision and recall of a model, as such; the F1 score allows for an efficient comparison of models with regard to trade-offs, that is, whether a model has a low recall due to missed hits or a low precision due to high inclusion of non-hits as shown in fig 11.

The predicted probability values from each classifier were also used to calculate the area under the receiver operating characteristic curve (ROC-AUC) as shown in Table 1. The ROC-AUC is a measure of how well models can differentiate between BBB+ and BBB- compounds given all possible decision thresholds. This is helpful for drug discovery because it allows you to shift your cut-off to get the highest hit rates, but still want possible compounds to have high confidence as lead candidates[5]. When looking at the three neural networks (DNN, CNN, & GCN), their test losses and accuracies were reported as well. The test losses indicate how well the model's predicted probabilities differentiate from the ground truth labels, while the test accuracies allow direct comparisons with conventional models that were trained on the same dataset as shown in figure 7. These four metrics give a more balanced representation of the reliability of each model and will help you choose the best candidates to validate in the laboratory.

IV. Results and Discussion

Ultimately, the six models tested on the BBBP dataset produced high accuracy in predicting BBB permeability, with values ranging from approximately 0.81 to 0.90 (with the GCN model being 0.83). The model achieving the highest accuracy (0.896) and F1-score (0.938) was the XGBoost model, indicating that it had the best trade-off between the number of correctly identified BBB+ molecules and false-positives. The Random Forest model was also very effective, having an accuracy of 0.872 and an F1-score of 0.923, which supports the argument that tree-ensemble methods are stable and robust for descriptor-based molecular data. The 1D-

CNN model performed similarly (0.838 accuracy, 0.900 F1-score) and provided the highest AUC value (0.893), which suggests that its convolutional architecture captures useful local patterning in descriptor sequences and separates BBB+ and BBB- molecules very well. The GCN model achieved an accuracy of approximately 0.829, precision of 0.829, recall of 1.0, F1-score of 0.907, and an AUC of 0.742 which indicate that the GCN model is very sensitive to BBB+ molecules (it correctly identifies almost all BBB+) but sacrifices some rank ordering capability and it creates more false-positives than the best-performing tree-based models as shown in figure 10.

The accuracy of predictions across different model types is reported as relatively high and consistent, with non-GCN models having precision scores of 0.93-0.94[10], meaning that a model is likely to make correct predictions when it predicts that a compound is BBB+. Recall was able to distinguish between the different models much more clearly than precision was. XGBoost produced the highest recall rate (0.945) of all the methods studied, followed by Random Forest (0.916) and GCN (resulting in a perfect recall rate of 1.0). This indicates that these model types can find a greater proportion of BBB-permeable molecules than SVM, DNN, and 1D CNN models, which were found to have lower recall rates. This is an important characteristic of the model types that use these algorithms and is especially useful in CNS drug discovery as shown in figure 12. For example, a true positive prediction (a BBB-permeable molecule) that is not found through modeling is more problematic in terms of missed opportunities than the number and cost of testing a false positive (non-BBB-permeable molecule). Although both SVM (RBF) and DNN both have precision scores that are comparable to the best-performing model types (0.942 and 0.927), they have the lowest recall scores (0.828 and 0.836) of the six model types studied, and therefore slightly lower F1-scores (0.882 and 0.879). Both models are conservative and therefore may be the model of choice for those who want to be able to place a very high degree of confidence in identified BBB+ hits, whereas these models may not be as useful to those who need to maximize the number of potential BBB-permeable molecules[2].

XGBoost is the only model that can be applied, at least for the current configuration as shown in fig 11, as it has a very good accuracy, F1 score, and the highest AUC of 0.883, which makes it the best performing single model in all metrics tested. The 1D-CNN had a slightly lower accuracy than XGBoost, however, the 1D-CNN has a very high AUC and would be a good model to use if ranking of BBB permeability is more important than accuracy. The Random Forest is also a good baseline and has generally been a good model given the level of interpretability available. The GCN model is a good supplemental model, provided that it provides a high recall in addition to its graph-level representation, and will be very useful for reducing false negatives when it is incorporated into an ensemble. The SVM and DNN models can be improved further through additional hyperparameter tuning or feature engineering for recall improvement. The results demonstrated that gradient boosting ensembles and shallow convolutional networks are effective ways to predict BBB permeability based on descriptors, while GCNs allow for another level of understanding based on an awareness of structure.

To evaluate the external predictive abilities of the trained models, we assessed 4,082 CNS relevant compounds from the ChEMBL36 dataset. The prediction summary indicates that

between 80% and 91% of these compounds were classified as BBB+ with predictions made by both classical and deep learning approaches. The Random Forest trained model predicts that approximately 90.3% of the molecules are classified as BBB+, the XGBoost trained model predicts BBB+ for about 91.2% of the molecules, and the SVM, CNN, and DNN models classify approximately 79% to 82% of the molecules as BBB+. In contrast, the GCN model predicts all ChEMBL36 molecules as BBB+. This indicates that the GCN model is highly sensitive to borderline structures and tends to classify many as being permeable. The external screening behaviour suggests that XGBoost and Random Forest provide a more accurate and balanced way to prioritise CNS drug candidates. Additionally, GCN could act as a high recaller filter for compound screening when used as part of an ensemble model approach during early-stage CNS drug discovery to avoid missing any potentially permeable compounds as shown in fig 9 & 13.

PREDICTION RESULTS SUMMARY:

RANDOM_FOREST	: BBB+ = 3686/4082 (90.3%)		BBB- = 396/4082 (9.7%)
SVM	: BBB+ = 3328/4082 (81.5%)		BBB- = 754/4082 (18.5%)
XGBOOST	: BBB+ = 3724/4082 (91.2%)		BBB- = 358/4082 (8.8%)
CNN	: BBB+ = 3260/4082 (79.9%)		BBB- = 822/4082 (20.1%)
DNN	: BBB+ = 3261/4082 (79.9%)		BBB- = 821/4082 (20.1%)
GCN	: BBB+ = 4082/4082 (100.0%)		BBB- = 0/4082 (0.0%)

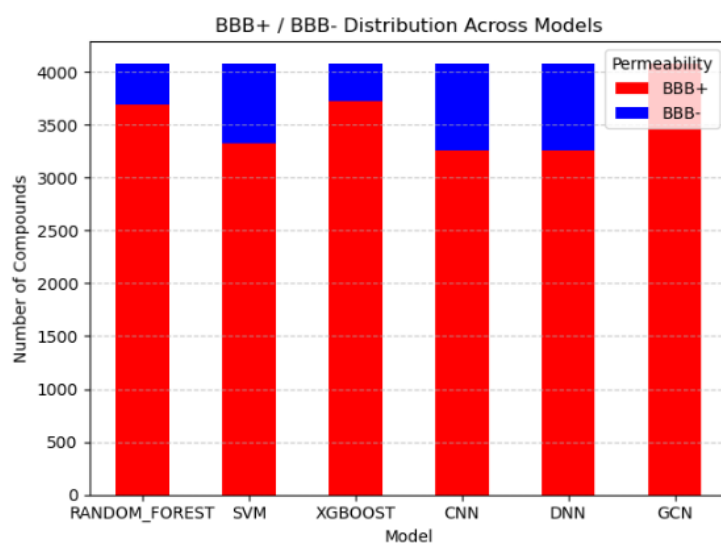


Figure 9. This represents the results of the prediction on unlabelled dataset to predict p_{np} (Chem36 dataset)

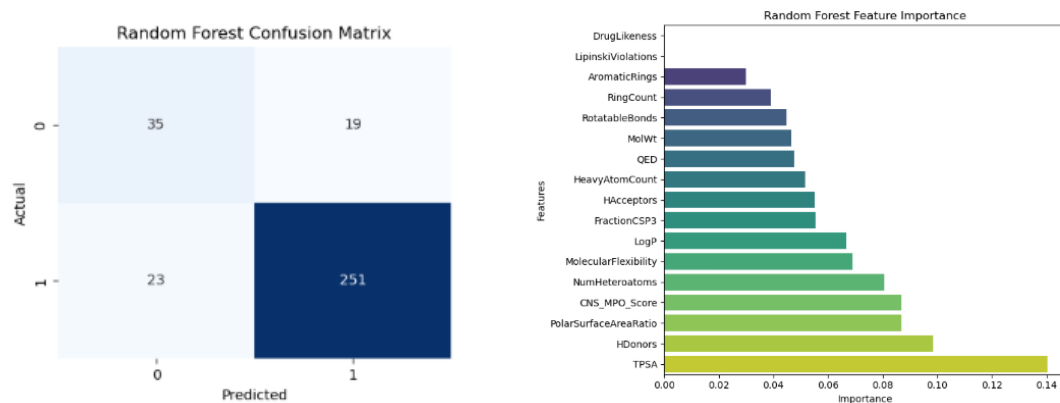


Figure 10. Random Forest Results and Important features

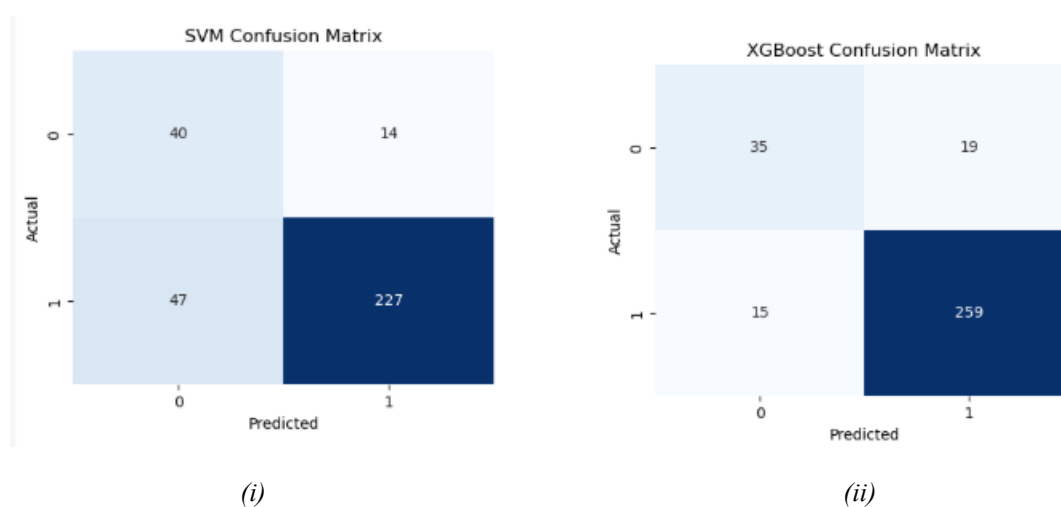


Figure 11. Confusion Matrix of SVM and XGBoost Model

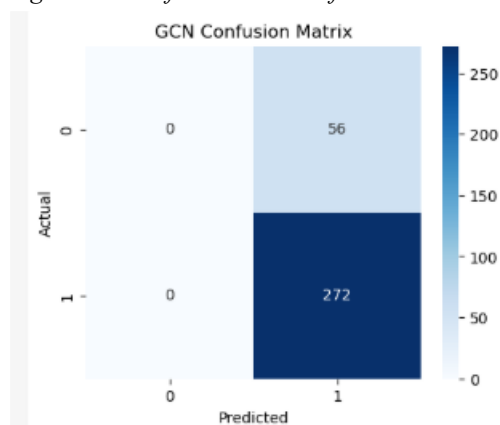


Figure 12. GCN Model Confusion Matrix

Table 1. Comprehensive Model Performance Comparison

Model	Accuracy	Precision	Recall	F1-Score	AUC
Random Forest	0.871951	0.929630	0.916058	0.922794	0.872229

Model	Accuracy	Precision	Recall	F1-Score	AUC
SVM (RBF)	0.814024	0.941909	0.828467	0.881553	0.877737
XGBoost	0.896341	0.931655	0.945255	0.938406	0.882975
1D-CNN	0.838415	0.933333	0.868613	0.899811	0.892775
DNN	0.807927	0.927126	0.835766	0.879079	0.865876
GCN	0.8293	0.8293	1.0000	0.9067	0.7415

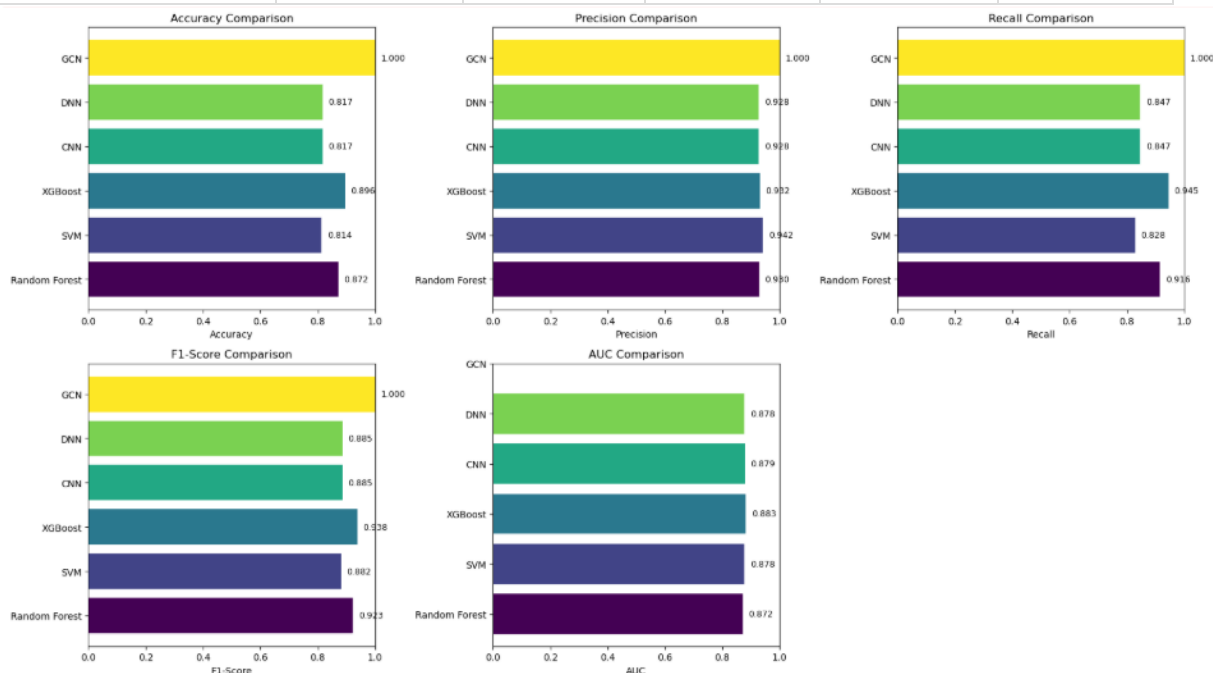


Figure 13. Comparative analysis of the trained model on the chembl36 unlabelled dataset

V. Conclusion and Future Scope

This study shows that machine learning (ML) models can predict the ability of drugs to cross the blood -brain barrier (BBB) based on their molecular structure (i.e. descriptors). The ML model with the best accuracy, F-1 score, and recall (i.e. captures the greatest percentage of all BBB+ compounds and least number of false positives) was XGBoost; most of the remaining models (Random Forest, 1D-CNN, SVM, and DNN) performed similarly in that they were all strong tree ensembles, but XGBoost outperformed all the others. The 1D-CNN had slightly

less accuracy than the other models but provided the highest AUC score (i.e. the best ability to distinguish between BBB+ and BBB− drugs). The SVM and DNN models performed surprisingly well given their high precision and moderate recall, making them ideal candidates to choose strike against any false predictions for BBB+. Therefore, the findings of this comparative analysis suggest that gradient-boosted tree and shallow CNN models are the most appropriate methods to classify BBBP on the current dataset along with GCN which performed well on the prediction and training dataset well.

These findings can be extended and generalized in a variety of ways that ultimately improve both understanding and validation of BBBP prediction models. First, developing a more rigorous validation process for these models (using external benchmark datasets and/or scaffold-based splits) would provide us with a better understanding of how robust these prediction models actually are and reduce potential overfitting to very specific chemotypes. Second, incorporating rich molecular representations of compounds, such as graph- or transformer-based learned embeddings versus using only the traditional descriptors, may actually improve the performance of these prediction models and reduce the amount of time spent on manual feature engineering. Third, developing interpretive methods for these prediction models (e.g., feature importance or SHAP analysis) can help medicinal chemists determine which physicochemical properties are driving the BBBP permeability decisions associated with each compound and thus provide them with actionable insights when designing compounds. Finally, the integration of the highest-performing BBBP prediction models with other components (e.g., ADMET filters, docking, or generative design) within an overall Central Nervous System (CNS) drug discovery pipeline would facilitate the identification of brain-penetrant drug candidates in an efficient and effective manner.

1. References

1. Alsenan, A., & Al-Turaiki, R. (2021). A deep learning approach to predict blood–brain barrier permeability. *PeerJ Computer Science*, 7, e515.
2. Kumar, R., Kumar, R., Garg, R., & Verma, R. (2022). DeePred-BBB: A blood brain barrier permeability prediction model based on deep neural network. *Frontiers in Neuroscience*, 16, 858126
3. Tang, Q., Xu, Y., Li, Z., & Lin, H. (2022). A merged molecular representation deep learning method for blood–brain barrier permeability prediction. *Briefings in Bioinformatics*, 23(5), bbac357.
4. Shaker, B., Lee, J., Lee, Y., et al. (2023). A machine learning-based quantitative model (LogBB_Pred) to predict blood–brain barrier permeability. *Bioinformatics*, 39(10), btad577.
5. Mazumdar, B., Ahammad, N. A., & Chakraborty, S. (2023). Machine learning-based dynamic consensus model for predicting blood–brain barrier permeability. *Computers in Biology and Medicine*, 160, 106984.

6. Swanson, K., et al. (2024). ADMET-AI: A machine learning ADMET platform for evaluation of large-scale chemical libraries. *Bioinformatics*, 40(7), btac416.
7. Toropov, A. A., et al. (2024). Innovative strategies for the quantitative modeling of blood–brain barrier (BBB) permeability: Harnessing the power of machine learning-based q-RASAR approach. *Molecular Engineering*, 15(7), d4me00056k.
8. Alsenan, A., & Al-Turaiki, R. (2021). A deep learning approach to predict blood–brain barrier permeability (open-access PMCID). National Library of Medicine (PMC).
9. Transparent ML authors (2024). Transparent machine learning model to understand drug permeability through the blood–brain barrier. *Journal of Chemical Information and Modeling*, 64(23), 8718–8728.
10. Purdie, M., et al. (2023). A review of blood–brain barrier permeability prediction models. National Library of Medicine (PMC).
11. Swanson, K., et al. (2024). ADMET-AI (Web/Package). *Bioinformatics*, 40(7), btac416.
12. Shaker, B., et al. (2023). LogBB_Pred (LightGBM) quantitative model. *Bioinformatics*, 39(10), btad577.
13. Toropov, A. A., et al. (2024). Classification read-across (c-RASAR) framework for BBB permeability prediction. *Journal of Chemical Information and Modeling*, 64, Article e-pub ahead of print.
14. Wan, M., et al. (2021). LightGBM-based BBB classifier (LightBBB). *Bioinformatics*, 37(8), 1135–1139.