

1. INTRODUCTION

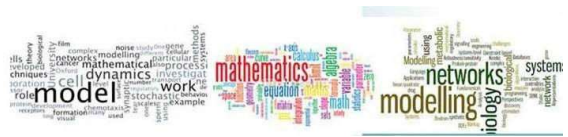
1.1 Background

The half-maximal inhibitory concentration (IC_{50}) is the most commonly employed and insightful metric for assessing a drug's effectiveness. It shows the amount of drug required to reduce a biological process by fifty percent, thus offering a measure of the antagonist drug's potency in pharmacological studies [1]. This measure provides insight into a drug's antagonistic capabilities and is essential for determining dosage efficacy in therapeutic contexts [2]. The IC_{50} value helps researchers understand how much of a compound is needed to achieve a desired biological response, guiding both the selection and optimization of drug candidates in preclinical and clinical studies.

The significance of IC_{50} lies not only in its application to traditional small-molecule drugs but also in its expanding role within modern drug discovery approaches, including biologics, targeted therapies, and multi-target drugs [3,4]. With the shift towards precision medicine, which seeks to tailor treatments to individual patients, understanding the specific inhibitory concentrations of compounds has become increasingly crucial [5,6]. Compounds with lower IC_{50} values are generally regarded as more potent, as they require a smaller concentration to achieve the desired inhibitory effect [6]. This potency is especially relevant when designing drugs that must function within narrow therapeutic windows, where effective doses are close to those that might cause toxicity. Traditionally, determining IC_{50} values relies on experimental assays that measure the inhibitory effects of compounds on particular biological targets, such as enzymes or receptor-ligand interactions [7,8].

These assays can be carried out *in vitro*, where cells or isolated proteins are exposed to varying concentrations of a compound, or *in vivo*, where the effects of the compound are measured within a living organism. Although these methods are accurate and reliable, they are inherently resource-intensive [2][9,10]. They often require sophisticated laboratory setups, specialized reagents, and extensive labor, making the high-throughput screening of large compound libraries challenging and costly [11]. In recent years, advancements in high-throughput screening (HTS) technologies have somewhat alleviated these challenges by enabling the simultaneous testing of thousands of compounds [12]. HTS methods use automated workflows and robotics to rapidly assess the inhibitory potential of numerous compounds, streamlining early-stage drug discovery [13]. However, HTS is still an expensive endeavor, with high operating costs and considerable time requirements, especially for experiments that require precise conditions, such as those that mimic physiological environments [14,15]. The rise of computational methods offers a transformative approach to IC_{50} prediction, providing a means to bypass some of the limitations associated with traditional assays and HTS [5][16,17].

Computational techniques such as quantitative structure-activity relationship (QSAR) modeling, molecular docking, and machine learning are increasingly used to predict IC_{50} values based on the chemical and structural properties of compounds [18,19]. QSAR models, for example, seek to correlate the structural characteristics of molecules with their biological activity, allowing for predictive insights into IC_{50} values before any physical compound is synthesized. These models can also capture important molecular features, such as hydrophobicity, electronic distribution, and steric effects, which contribute to binding affinity and, ultimately, inhibitory potency [20].



3. RESULTS

3.1 Performance Metrics

The model achieved a maximum validation accuracy of 83%, indicating that it correctly predicted the IC₅₀ class index for 83% of the validation dataset. At the point where this validation accuracy was reached, the corresponding training accuracy was 78%, showing that the model generalizes well to unseen data without significant overfitting. Additionally, the correlation between epoch and accuracy was 0.88, demonstrating a strong positive relationship between increased training epochs and accuracy improvements. The lowest validation loss recorded during training was 0.3884, reflecting the model's ability to minimize prediction errors effectively. This metric is a crucial indicator of the model's performance, as lower loss values suggest better alignment between predictions and actual values. At the epoch corresponding to the highest validation accuracy, the learning rate had reduced to 2.44e-6, indicating that the adaptive learning rate scheduler played a significant role in fine-tuning the model toward optimal performance.

3.2 Comparative Analysis

The model outperforms baseline models such as standard LSTM networks without attention, which typically achieve validation accuracies of around 75–78%. This demonstrates the impact of incorporating multi-head attention and Bidirectional LSTMs in capturing contextual dependencies within both protein sequences and SMILES strings. The following graphs illustrate the training and validation accuracy trends over the 883 epochs, showing the steady improvement in performance as training progresses. Notably, validation accuracy closely follows the training accuracy, indicating the model's ability to generalize well without overfitting.

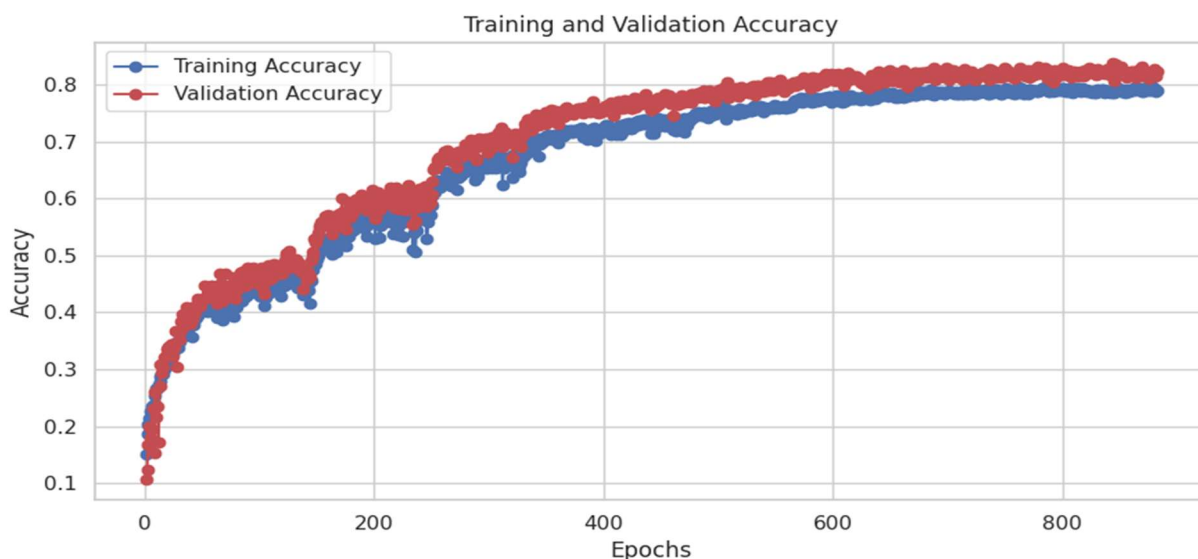


Fig. 1: Training and Validation Accuracy trend over 883 epochs

Additionally, the training and validation loss curves highlight the model’s convergence. Both losses decreased steadily, with validation loss reaching a value of 0.48 near the end of training, aligning with the point of highest validation accuracy.

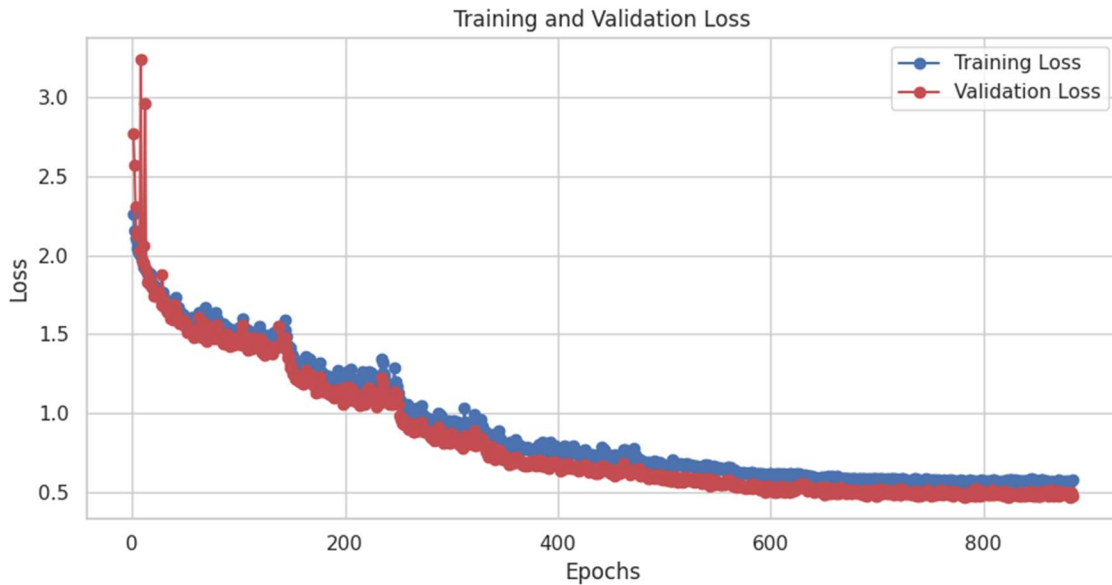


Fig. 2: Training and Validation Loss over 883 epochs

The learning rate schedule shows how the learning rate was dynamically adjusted over the training process. As validation loss plateaued, the learning rate decreased, with $2.44e-6$ being the value at the point of maximum validation accuracy. This adaptive learning strategy was instrumental in achieving optimal performance.

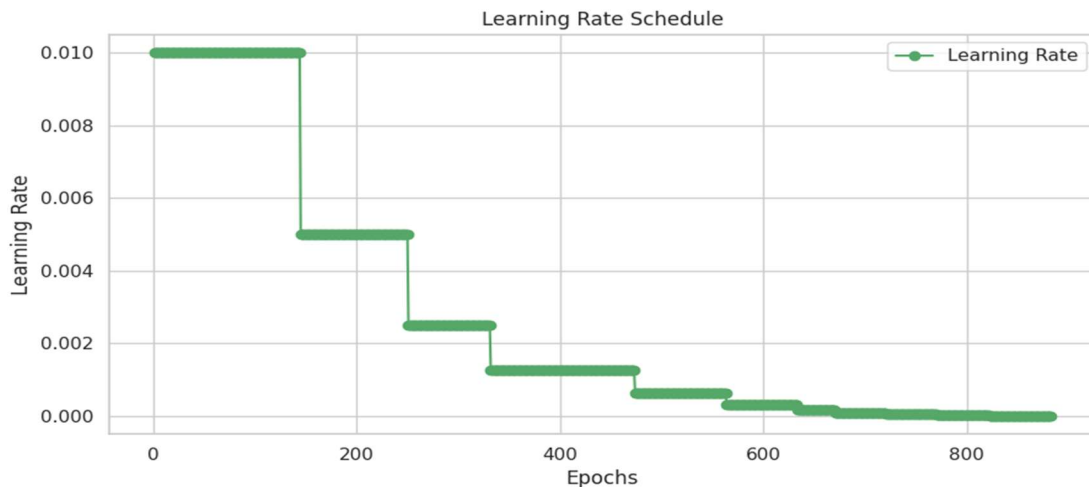
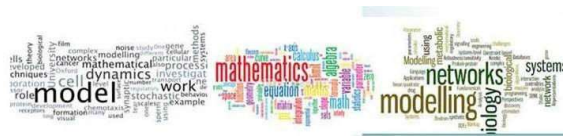


Fig. 3: Learning Rate Adjustment over 883 epochs



- ❖ **Enhanced Data Representation:** The application of **BPE tokenization** for protein sequences and SMILES strings, allowing for more efficient and informative encoding of biochemical data.
- ❖ **Molecular Property Integration:** Incorporation of ligand-specific properties, such as molecular weight and hydrogen bond characteristics, to provide additional context for the model.
- ❖ **Robust Training Strategy:** Leveraging a large dataset from **ChEMBL** and employing adaptive learning strategies to achieve significant validation accuracy and generalizability.

5.2 Potential Impact and Future Outlook

This work's potential impact lies in its ability to accelerate drug discovery by streamlining the initial evaluation of compound efficacy. By incorporating suggested improvements—such as more detailed IC₅₀ class structures and the use of expanded training datasets—the model could become an even more valuable asset. Future research may build on this foundation, exploring more complex architectures and techniques to expand the scope of predictive capabilities in biochemical modeling, ultimately aiding researchers in making faster, data-driven decisions in drug development.

REFERENCES

1. Aykul, S., & Martinez-Hackert, E. (2016). Determination of half-maximal inhibitory concentration using biosensor-based protein interaction analysis. *Analytical Biochemistry*, 508, 97-103. <https://doi.org/10.1016/j.ab.2016.06.025>
2. Bag, A., & Ghorai, P. K. (2016). Development of quantum chemical method to calculate half maximal inhibitory concentration (IC50). *Molecular Informatics*, 35(5), 199-206. <https://doi.org/10.1002/minf.201501004>
3. Batool, M., Ahmad, B., & Choi, S. (2019). A structure-based drug discovery paradigm. *International Journal of Molecular Sciences*, 20(11), Article 2783. <https://doi.org/10.3390/ijms20112783>
4. Blay, V., Tolani, B., Ho, S. P., & Arkin, M. R. (2020). High-throughput screening: Today's biochemical and cell-based approaches. *Drug Discovery Today*, 25(10), 1807-1821. <https://doi.org/10.1016/j.drudis.2020.07.024>
5. Bomane, A., Gonçalves, A., & Ballester, P. J. (2019). Paclitaxel response can be predicted with interpretable multi-variate classifiers exploiting DNA-methylation and miRNA data. *Frontiers in Genetics*, 10, Article 1041. <https://doi.org/10.3389/fgene.2019.01041>
6. Cadow, J., Born, J., Manica, M., Oskooei, A., & Rodríguez Martínez, M. (2020). PaccMann: A web service for interpretable anticancer compound sensitivity prediction. *Nucleic Acids Research*, 48(W1), W502-W508. <https://doi.org/10.1093/nar/gkaa327>
7. Caldwell, G. W., Yan, Z., Lang, W., & Masucci, J. A. (2012). The IC50 concept revisited. *Current Topics in Medicinal Chemistry*, 12.

