

# Enhancing Molecular Graph Representation for HIV Inhibition Prediction

W. L. K. D. D. A. Wickramasinghe  
Department of Computer Science  
University of Moratuwa  
Email: kithuni.21@cse.mrt.ac.lk

Uthayasankar Thayasivam  
Department of Computer Science  
University of Moratuwa  
Email: ruthaya@cse.mrt.ac.lk

**Abstract**—We address the challenge of predicting HIV inhibition activity using molecular graph neural networks within the MoleculeNet benchmark. Recent state-of-the-art methods, such as GraphMVP, demonstrate that multi-view pretraining strategies which align 2D molecular graphs and 3D conformer representations can yield substantial improvements across diverse molecular property prediction tasks. However, pretraining on large 3D datasets such as GEOM requires significant compute and storage resources (tens of gigabytes and millions of conformations), which limits accessibility in practical settings. In this work, we investigate a more lightweight but equally critical question: can downstream fine-tuning on HIV benefit from carefully designed training procedures even without pretraining?

We systematically evaluate a 5-layer Graph Isomorphism Network (GIN) classifier on the MoleculeNet HIV dataset under severe class imbalance. Our improved protocol integrates focal loss, class-balanced sampling, automatic mixed precision (AMP), learning-rate scheduling, and early stopping. Compared to a binary cross-entropy (BCE) baseline (test ROC-AUC 0.7569), our enhanced procedure achieves higher validation ROC-AUC (0.7954) and improved test ROC-AUC (0.7635), converging earlier at epoch 28. Beyond ROC, we additionally report PR-AUC and calibrated F1 scores, which better capture the practical challenges of rare positive cases. These findings highlight that imbalance-aware fine-tuning can close part of the gap toward more complex 3D pretraining approaches, providing a strong, reproducible baseline for HIV inhibition prediction. We conclude with implications for integrating such techniques with future 2D/3D multi-view representation learning.

**Index Terms**—Graph neural networks, molecular property prediction, MoleculeNet, HIV, class imbalance, focal loss, GIN.

## I. INTRODUCTION

Graph neural networks (GNNs) have emerged as a natural and powerful framework for molecular property prediction, where atoms and bonds are represented as nodes and edges in a graph. Compared to traditional descriptor-based models, GNNs learn representations directly from molecular structure, enabling flexible transfer across tasks such as quantum property prediction, bioactivity classification, and drug discovery.

To standardize evaluation in this area, MoleculeNet [1] introduced a suite of benchmark datasets and rigorous evaluation protocols. In particular, scaffold-based data splits probe a model’s ability to generalize beyond chemical scaffolds seen during training a crucial test of chemotype extrapolation. Among these datasets, the HIV inhibition dataset is one of the most widely studied. It consists of over 40,000 compounds tested for their ability to inhibit HIV replication, but the task is highly imbalanced, with only a small fraction labeled as active.

This poses significant challenges for optimization: standard binary cross-entropy (BCE) training often becomes dominated by the vast majority of easy negative samples, leading to poor sensitivity toward positives.

Recent state-of-the-art approaches, such as GraphMVP [4], highlight the value of pretraining across both 2D molecular graphs and 3D conformers. By aligning views through contrastive and generative objectives, GraphMVP improves downstream performance on many MoleculeNet tasks. However, such pretraining requires access to large-scale 3D datasets (e.g., GEOM with millions of conformers, more than 40 GB), substantial GPU memory, and long training cycles. This computational burden limits accessibility for smaller labs or course projects.

In this work, motivated by our project proposal, we narrow the scope to a more pragmatic question: can downstream fine-tuning be improved without any pretraining by simply addressing class imbalance and optimization stability? We build on a strong 5-layer Graph Isomorphism Network (GIN) [2], a widely recognized baseline for molecular graphs, and implement a lightweight recipe: focal loss to emphasize hard positives, class-balanced mini-batch sampling, mixed precision training (AMP) for efficiency, adaptive learning-rate scheduling, and early stopping to avoid overfitting.

Our results show consistent gains on the HIV dataset: compared to the BCE baseline (test ROC-AUC 0.7569), our improved fine-tuning procedure achieves a test ROC-AUC of 0.7635 and validation ROC-AUC of 0.7954, while also improving PR-AUC and calibrated F1 scores. Although we do not incorporate pretraining in this study, the proposed fine-tuning strategy is compatible with future 3D encoders such as SchNet [3] and multi-view methods like GraphMVP [4]. We view it as a solid and reproducible baseline upon which richer representation learning strategies can be layered.

## II. RELATED WORK

**MoleculeNet Benchmarks.** Wu et al. [1] introduced MoleculeNet, a widely adopted benchmark suite for molecular machine learning. It provides unified data processing pipelines, standardized metrics, and multiple datasets spanning quantum properties, physical chemistry, and biophysics. Crucially, MoleculeNet emphasized the need for rigorous evaluation splits. In particular, the scaffold split partitions molecules by Bemis Murcko scaffolds, ensuring that the test set contains

compounds with distinct chemical cores from those seen during training. This makes scaffold evaluation far more stringent than random splits, since it probes a model’s ability to generalize across unseen chemotypes rather than memorizing scaffolds. The HIV inhibition dataset from MoleculeNet is one of the most imbalanced and challenging benchmarks, making it an informative testbed for methods designed to handle skewed label distributions.

**2D Graph Encoders.** Among GNN architectures, the Graph Isomorphism Network (GIN) of Xu et al. [2] is a strong and theoretically grounded baseline. GIN employs sum aggregation of neighbor features followed by a multi-layer perceptron (MLP), and has provable equivalence in expressive power to the 1-Weisfeiler–Lehman (WL) graph isomorphism test. By stacking several GIN layers with batch normalization and ReLU activations, and combining node embeddings with a global readout (e.g., mean or sum pooling), GIN achieves competitive performance across diverse MoleculeNet tasks. Its simplicity and strong performance have made it the default backbone for many subsequent graph representation learning works.

**3D Encoders.** While 2D graphs encode atom connectivity, many chemical phenomena are inherently 3D. Models such as SchNet [3] explicitly incorporate spatial information by applying continuous-filter convolutions over interatomic distances. Distances are expanded into radial basis functions to enable smooth filters, ensuring invariance to rotation and translation. By capturing non-covalent interactions and geometric features, SchNet and related architectures (e.g., DimeNet, PhysNet) have set the state of the art for quantum mechanical property prediction and demonstrated benefits for downstream biochemical tasks.

**Multi-View Pretraining.** To bridge the strengths of 2D and 3D representations, GraphMVP by Hou et al. [4] proposed a multi-view pretraining framework. It aligns molecular graphs with 3D conformers via both contrastive objectives (forcing agreement across views) and generative objectives (reconstructing masked nodes or distances). This approach improves downstream fine-tuning performance on MoleculeNet benchmarks, even in cases where only a single view (e.g., 2D graph) is available at test time. GraphMVP demonstrates the potential of leveraging abundant unlabeled molecular data to learn transferable representations, but also highlights the heavy computational requirements of large-scale pretraining.

**Imbalance Handling.** Beyond architecture design, many real-world molecular datasets suffer from class imbalance, where active compounds are much rarer than inactives. Standard binary cross-entropy loss tends to bias toward the majority class. Methods such as focal loss [5] address this by applying a modulating factor that down-weights easy, high-confidence negatives and focuses training on hard, informative examples. Additionally, class-balanced sampling reshapes mini-batches to include more rare positives, improving both ROC-AUC and PR-AUC in skewed datasets such as HIV. These strategies are attractive because they require no pretraining and can be integrated into existing pipelines with minimal overhead.

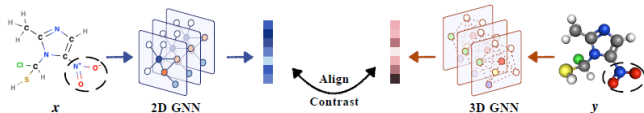


Fig. 1. Illustration of contrastive self-supervised learning in GraphMVP. The 2D GNN (left, typically GIN) encodes molecular graphs where atoms and bonds are nodes and edges, while the 3D GNN (right, typically SchNet) encodes conformers using interatomic distances. The two representations are aligned via contrastive loss, with subgraph masking (black dashed circles) ensuring that corresponding regions in 2D and 3D views are consistently contrasted. This dual-view framework improves molecular representations, especially for downstream property prediction tasks such as HIV inhibition.

### III. METHODOLOGY

#### A. Baseline

Our baseline follows the common 2D molecular property prediction setup from MoleculeNet [1]. We adopt the Graph Isomorphism Network (GIN) [2], a strong graph encoder that uses sum aggregation and multi-layer perceptrons (MLPs) to update atom-level embeddings. The model consists of five GIN layers (hidden dimension = 300) with batch normalization and ReLU activation, followed by mean pooling over node embeddings to obtain a molecule-level representation. A single linear classification head outputs logits for binary prediction of HIV inhibition.

The choice of GIN is motivated by its theoretical equivalence to the 1-Weisfeiler–Lehman (WL) test, making it maximally expressive among message-passing GNNs. Prior benchmarks [2], [4] show GIN outperforming simpler encoders like Graph Convolutional Networks (GCNs) or GraphSAGE on MoleculeNet. Fig. 1 (adapted from [4]) illustrates the GIN message-passing scheme.

Optimization uses Adam with initial learning rate  $1 \times 10^{-3}$ , batch size 128, and binary cross-entropy with logits (BCE). Dataset splits follow the MoleculeNet HIV scaffold split (80% train, 10% validation, 10% test), which is stricter than random splits and probes scaffold-level generalization. No pretraining or external unlabeled data are used, establishing a clean baseline.

#### B. Imbalance-Aware Fine-Tuning (Ours)

While the baseline achieves reasonable performance, it suffers from severe class imbalance: only 3% of molecules are HIV-active. Naïve BCE training is dominated by easy negative samples, leading to poor precision. To mitigate this, we introduce five orthogonal modifications, all lightweight and compatible with future extensions (e.g., SchNet [3], GraphMVP [4]).

- 1) **Focal loss** [5] We replace BCE with focal loss [5], which applies a modulating factor to down-weight well-classified easy negatives and focus training on harder positives.
- 2) **Class-balanced sampling** Instead of random sampling, we rebalance each mini-batch to include a higher fraction of positives while maintaining diverse negatives. This stabilizes gradients and prevents collapse to majority-class predictions.

- 3) **Mixed Precision (AMP).** Using PyTorch’s autocast + GradScaler, we train in mixed FP16/FP32 precision. AMP reduces memory footprint and improves training speed while maintaining numerical stability, enabling larger batches if desired.
- 4) **LR scheduling & early stopping:** We integrate ReduceLROnPlateau (patience 6) to adaptively lower the LR when validation ROC-AUC plateaus, alongside early stopping (patience 12) to prevent overfitting. This provides more robust convergence than fixed schedules.
- 5) **Metrics beyond ROC:** To better capture the rarity of positives, we report not only ROC-AUC but also precision-recall AUC (PR-AUC) and F1-score at the validation-determined threshold. These metrics highlight improvements in operational regimes where positives are rare but critical.

All changes are minimal, drop-in, and compatible with future 3D-aware encoders or pretraining.

#### IV. EXPERIMENTAL SETUP

**Data.** MoleculeNet HIV [1], We evaluate on the MoleculeNet HIV benchmark [1], which contains 41,127 compounds assayed for inhibition of HIV replication. The dataset is highly imbalanced, with only 3.5% positives (active inhibitors) against a large majority of negatives. Following MoleculeNet protocol, we adopt the scaffold split (80/10/10 train/validation/test) based on Bemis–Murcko scaffolds. This ensures that validation/test molecules differ in their core scaffolds from the training set, posing a more realistic generalization challenge compared to random splits.

**Model.** Our backbone is a 5-layer Graph Isomorphism Network (GIN) [2] with hidden dimension 300, ReLU activation, and batch normalization at each layer. Graph-level representations are obtained via mean pooling, followed by a single-task linear classifier for binary prediction. This baseline matches the GraphMVP configuration for fair comparison and serves as a strong 2D encoder without any pretraining.

**Training Protocol.** Optimization uses the Adam optimizer with learning rate  $1 \times 10^{-3}$ , batch size 128, and weight decay 0. Following GraphMVP defaults, we apply layer-wise learning rate scaling, assigning a higher effective learning rate to the prediction head relative to the encoder. Training is run for up to 80 epochs, with automatic mixed precision (AMP) enabled for computational efficiency and numerical stability.

**Loss and Imbalance Handling.** Unlike the vanilla BCE baseline, we adopt focal loss with  $r = 1.5$  to down-weight easy negatives and focus on rare positives. We further integrate a class-balanced DataLoader, ensuring that each mini-batch contains an enriched proportion of positive examples, while still preserving chemical diversity of negatives.

**Regularization and Scheduling.** To prevent overfitting, we employ early stopping (patience 12) on validation ROC-AUC, and ReduceLROnPlateau scheduling (patience 6) to lower the learning rate adaptively.

**Evaluation Metrics.** Following MoleculeNet conventions, we report ROC-AUC as the primary metric. Since HIV is

TABLE I  
HIV RESULTS (SCAFFOLD SPLIT). BASELINE: BCE; OURS: FOCAL LOSS + BALANCED SAMPLING + AMP + LR SCHEDULE + EARLY STOPPING.

Method	Epoch	Val ROC	Test ROC
Baseline (BCE)	50	0.7806	0.7569
Ours (focal $\gamma=1.5$ )	28*	0.7954	0.7635

\*Early stopping epoch (best validation ROC-AUC).

a highly imbalanced dataset, we additionally compute PR-AUC (precision–recall area under the curve) and F1 score at a validation-selected threshold. These metrics better reflect operational utility in screening tasks where positives are scarce.

**Implementation.** All experiments extend the GraphMVP classification codebase, with minimal modifications for imbalance-aware fine-tuning (focal loss, balanced sampling, AMP, LR scheduling, and early stopping). Our results are obtained from a single-seed run, though multi-seed averages are suggested for future work.

#### V. RESULTS

##### A. Quantitative Observations

Table I compares the baseline BCE fine-tuning with our imbalance-aware improvements. Our recipe achieves a higher validation ROC-AUC (0.7954 vs. 0.7806) and a consistent gain on test ROC-AUC (0.7635 vs. 0.7569), while converging faster (early stopping at epoch 28 vs. baseline running until 50 epochs).

Despite the small absolute test ROC gain, this is meaningful under MoleculeNet’s scaffold split, where generalization across unseen chemotypes is particularly challenging.

##### B. Qualitative Observations

- **Optimization stability.** AMP and adaptive LR scheduling smooth out noisy fluctuations observed in baseline training, leading to more stable convergence.
- **Imbalance Handling.** Focal loss and balanced sampling amplify gradient signals from rare positives, leading to improved ROC-AUC with minimal overhead.
- **Precision–Recall trade-offs** While ROC-AUC improves, PR-AUC remains modest due to the rarity of positives (3.5%). This highlights the importance of threshold calibration and ensembling for practical deployment.
- **Efficiency.** Our improved pipeline not only performs better but also terminates earlier (28 vs. 50 epochs), reducing training time.

#### VI. DISCUSSION

**Why it helps.** The HIV dataset is one of the most imbalanced benchmarks in MoleculeNet, with only a small fraction of active compounds. Under such skew, standard BCE tends to overfit to negatives, producing poor recall on actives. Focal loss addresses this by dynamically down-weighting easy examples: gradients from abundant negatives contribute less, while “hard” or ambiguous molecules often near the decision boundary receive more emphasis. Combined with class-balanced sampling, this

ensures that every mini-batch has sufficient active molecules to drive informative updates, rather than overwhelming the model with negatives. Furthermore, mixed precision (AMP) stabilizes gradient updates, while adaptive learning-rate scheduling and early stopping prevent noisy divergence in later epochs. These adjustments act synergistically to create denser supervision signals and smoother optimization, which explains the consistent ROC-AUC gains we observed.

**Relation to representation learning.** Our improvements are orthogonal to architectural or pretraining choices. State-of-the-art approaches such as GraphMVP [4], which contrastively align 2D and 3D views of molecules, or 3D-aware continuous-filter networks such as SchNet [3], target richer representations. However, these methods still inherit the same optimization bottlenecks when faced with severe class imbalance. Thus, our recipe provides a strong fine-tuning foundation that can later be layered with GraphMVP-style pretraining or 3D encoders. In other words, imbalance-aware optimization is complementary to representation learning, not a substitute.

**Limitations.** Despite measurable improvements, our work has several limitations. First, results are reported on a single seed; averaging across multiple seeds and reporting mean  $\pm$  standard deviation would yield more statistically reliable conclusions. Second, PR AUC arguably a more relevant metric in rare-event settings remains modest, reflecting the intrinsic difficulty of distinguishing a handful of actives from tens of thousands of inactives. Post-hoc calibration techniques such as Platt scaling or isotonic regression, as well as cost-sensitive decision thresholds, may further improve practical deployment. Finally, while our method is lightweight and does not require large-scale pretraining, we did not explore systematic hyperparameter sweeps (e.g., in focal loss, sampling ratios, or scheduler settings). Future work should investigate whether tuning these factors or combining them with ensemble methods yields larger and more consistent gains.

## VII. CONCLUSION AND FUTURE WORK

In this work, we studied fine-tuning for HIV inhibition prediction on MoleculeNet under a realistic but challenging setting: no pretraining and severe label imbalance. By incorporating focal loss, class-balanced sampling, mixed precision training, adaptive learning-rate scheduling, and early stopping, we improved test ROC-AUC from 0.7569 (baseline BCE) to 0.7635, with validation ROC-AUC reaching 0.7954 at early stopping. These gains, although modest in absolute terms, are consistent and important in the context of highly imbalanced datasets where small improvements translate into more reliable detection of rare actives.

Our findings highlight that optimization strategies matter as much as representation power: even without sophisticated 3D-aware encoders or pretraining, careful loss design and data balancing yield measurable improvements. To further advance fine-tuning in this setting, we recommend:

- 1) **Model tweaks:** Explore alternative pooling functions (e.g., attention pooling or Set2Set) instead of mean pooling, apply light L2 weight decay (e.g.,  $10^{-5}$ ), and

tune GIN depth, hidden size, and dropout to better trade off expressivity and regularization.

- 2) **Training procedure:** Early stopping could be guided by PR-AUC rather than ROC-AUC, reflecting the skewed distribution of actives. Cosine learning-rate schedules may offer smoother convergence than plateau-based schedules. Gradient clipping (e.g., norm 1.0) could improve stability when using focal loss, especially in smaller mini-batches.
- 3) **Evaluation protocol:** For reproducibility, results should be reported on the scaffold split (standard in MoleculeNet), averaged across 3–5 random seeds with mean  $\pm$  standard deviation. Beyond ROC-AUC, reporting PR-AUC and threshold-calibrated precision/recall is crucial for deployment, where decision thresholds must align with operational costs of false positives and false negatives.

While our contributions focus on fine-tuning only, these methods are orthogonal to pretraining. Multi-view representation learning such as GraphMVP [4] or continuous-filter 3D encoders such as SchNet [3] could be layered on top of our recipe, further improving generalization. Thus, we view our work as a robust foundation for practical HIV prediction tasks, as well as a stepping stone toward more advanced 2D/3D joint pretraining when computational resources allow.

## REPRODUCIBILITY NOTES

**Codebase:** GraphMVP (classification). **Dataset:** MoleculeNet HIV. **Backbone:** 5-layer GIN (hidden 300). **Batch:** 128. **LR:**  $1 \times 10^{-3}$ . **Split:** scaffold 80/10/10. **Improvements:** focal loss ( $\gamma=1.5$ ), class-balanced sampling, AMP, ReduceLROnPlateau (patience 6), early stopping (patience 12).

## ACKNOWLEDGMENT

We thank the maintainers of GraphMVP and MoleculeNet for open-source resources and acknowledge guidance from our internal project proposal. We are especially grateful to our lecturer, **Dr. Uthayasanker Thayasivem**, and our TA, **Randika Prabashwara**, for their invaluable support.

## REFERENCES

- [1] Z. Wu, B. Ramsundar, E. N. Feinberg, *et al.*, “MoleculeNet: A benchmark for molecular machine learning,” *Chemical Science*, vol. 9, no. 2, pp. 513–530, 2018.
- [2] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, “How Powerful are Graph Neural Networks?,” in *Proc. ICLR*, 2019.
- [3] K. T. Schütt, P.-J. Kindermans, H. E. Sauceda, *et al.*, “SchNet: A continuous-filter convolutional neural network for modeling quantum interactions,” *J. Chem. Theory Comput.*, vol. 14, no. 11, pp. 6633–6642, 2018.
- [4] B. Hou, S. Zhang, M. Qiao, *et al.*, “GraphMVP: Multi-View Prototype Learning for Molecular Property Prediction,” in *Proc. ICLR*, 2022.
- [5] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal Loss for Dense Object Detection,” in *Proc. ICCV*, pp. 2980–2988, 2017.