

# Graph-based Virtual Screening with Heterogeneous Message Passing Networks

**Authors:** Elena Petrov<sup>1</sup>, James Thompson<sup>2</sup>, Maria Santos<sup>1</sup>, Robert Liu<sup>3</sup>

**Affiliations:**

<sup>1</sup>Department of Computational Chemistry, UC Berkeley

<sup>2</sup>Pharmaceutical Sciences Division, Cambridge University

<sup>3</sup>Center for Drug Design, University of Michigan

## Abstract

We propose HeteroMPNN, a heterogeneous message passing neural network for large-scale virtual screening in drug discovery. Our approach models molecular diversity through heterogeneous graph representations and employs specialized message passing for different atom and bond types. Evaluated on the ChEMBL database with over 2 million compounds against 12 protein targets, HeteroMPNN achieves 15% improvement in hit rate compared to traditional fingerprint-based methods while maintaining sub-second screening times. The model demonstrates exceptional performance in identifying structurally diverse active compounds, with enrichment factors up to 3.2x at 1% screening threshold.

**Keywords:** Virtual Screening, Heterogeneous Graphs, Message Passing, Drug Discovery, ChEMBL

## 1. Introduction

Virtual screening is a computational technique used to identify promising drug candidates from large chemical databases before costly experimental validation. Traditional approaches rely on molecular similarity searches using fingerprints such as ECFP or pharmacophore models, which may miss structurally diverse compounds with similar biological activities due to activity cliffs—situations where structurally similar molecules exhibit dramatically different biological activities.

Graph neural networks offer a promising alternative by learning task-specific molecular representations that can capture subtle structural differences crucial for biological activity. However, existing GNN approaches treat all atoms and bonds uniformly, potentially missing important chemical distinctions that influence molecular behavior.

In this work, we introduce HeteroMPNN, which explicitly models chemical diversity through heterogeneous graph representations. Our key contributions include: (1) a heterogeneous molecular graph formulation that distinguishes atom and bond types, (2) specialized message passing functions for different chemical environments, and (3) meta-path aggregation to capture higher-order chemical motifs.

## 2. Related Work

### 2.1 Virtual Screening Methods

Classical virtual screening approaches include ligand-based methods (similarity searching, pharmacophore modeling) and structure-based methods (molecular docking, de novo design). Machine learning approaches have gained popularity, with methods ranging from support vector machines with molecular descriptors to deep neural networks with learned representations.

## 2.2 Heterogeneous Graph Neural Networks

Heterogeneous graphs, containing multiple node and edge types, have been successfully applied in recommendation systems and knowledge graphs. Recent work has explored their application to molecular property prediction, though primarily focusing on single-target scenarios.

## 2.3 Large-Scale Molecular Databases

The ChEMBL database contains over 2 million bioactive compounds with associated target information, making it an ideal testbed for virtual screening methods. Previous studies have used subsets of ChEMBL for method validation, but few have tackled the full-scale screening challenge.

# 3. Methodology

## 3.1 Heterogeneous Molecular Graph Representation

We extend standard molecular graphs to heterogeneous representations that capture chemical diversity:

### 3.1.1 Node Types

We define eight atom types based on chemical properties:

- **Carbon:**  $sp^3$ ,  $sp^2$ ,  $sp$  hybridization subtypes
- **Nitrogen:** Primary, secondary, tertiary, aromatic subtypes
- **Oxygen:** Ether, alcohol, carbonyl, nitro subtypes
- **Sulfur:** Thiol, sulfide, sulfoxide, sulfone subtypes
- **Phosphorus:** Phosphate, phosphonate subtypes
- **Halogens:** F, Cl, Br, I as separate types
- **Metals:** Transition metals, alkaline earth metals
- **Others:** B, Si, Se and rare atoms

### 3.1.2 Edge Types

We categorize bonds into 12 types:

- **Single bonds:**  $sp^3-sp^3$ ,  $sp^3-sp^2$ ,  $sp^2-sp^2$ , aromatic-aromatic
- **Double bonds:** C=C, C=O, C=N, N=O
- **Triple bonds:**  $C\equiv C$ ,  $C\equiv N$
- **Coordination bonds:** Metal-ligand interactions

- **Hydrogen bonds:** Inferred from geometry

### 3.1.3 Meta-paths

We define pharmacologically relevant meta-paths:

- **Pharmacophore patterns:** Donor-Acceptor-Hydrophobic
- **Aromatic systems:** Ring-Ring connections
- **Scaffold hops:** Core-to-core relationships

## 3.2 HeteroMPNN Architecture

### 3.2.1 Type-specific Embeddings

Each atom type  $\tau$  has dedicated embedding parameters:

$$h^0_v = W^{\tau(v)} \cdot x_v + b^{\tau(v)}$$

### 3.2.2 Heterogeneous Message Passing

For each message passing step  $l$ , we compute type-specific messages:

$$m^l_{u \rightarrow v} = \varphi^{\{\tau(u), \tau(v), \tau(e_{uv})\}}(h^l_u, e_{uv})$$

The aggregation function considers edge types:

$$h^{l+1}_v = \psi^{\tau(v)}(h^l_v, \bigoplus_{u \in N(v)} m^l_{u \rightarrow v})$$

### 3.2.3 Meta-path Aggregation

We implement attention-based meta-path aggregation:

$$\begin{aligned} \alpha_p &= \text{softmax}(W_{\text{att}} \cdot \tanh(W_p \cdot h_{\text{path}} + b_p)) \\ h_{\text{meta}} &= \sum_p \alpha_p \cdot h_{\text{path}} \end{aligned}$$

### 3.2.4 Graph-level Representation

The final molecular representation combines local and meta-path features:

$$h_{\text{mol}} = \text{CONCAT}([\text{MEAN\_POOL}(H), \text{MAX\_POOL}(H), h_{\text{meta}}])$$

## 3.3 Multi-target Training Strategy

We employ a multi-target training approach to leverage shared chemical knowledge:

$$L_{\text{total}} = \sum_{t=1}^T \lambda_t \cdot L_t + \lambda_{\text{reg}} \cdot \|\theta\|^2$$

where  $\lambda_t$  weights each target's contribution and  $\lambda_{\text{reg}}$  controls regularization.

### 3.4 Efficient Screening Pipeline

For large-scale screening, we implement:

- 1. **Batch processing:** Process up to 10,000 molecules simultaneously
- 2. **GPU acceleration:** Optimized CUDA kernels for graph operations
- 3. **Caching:** Store computed embeddings for repeated queries
- 4. **Early stopping:** Skip molecules below confidence threshold

## 4. Experimental Setup

### 4.1 Dataset Construction

We construct a comprehensive virtual screening benchmark from ChEMBL v29:

#### 4.1.1 Target Selection

We select 12 diverse protein targets with sufficient bioactivity data:

Target	PDB ID	Protein Family	Active Compounds	Total Compounds
EGFR	1M17	Kinase	15,432	89,234
HIV-1 RT	1RT2	Polymerase	8,967	45,123
AChE	1EVE	Hydrolase	12,543	67,891
COX-2	1CX2	Oxidoreductase	9,876	52,467
BACE-1	1FKN	Protease	7,234	38,901
hERG	1BYW	Ion Channel	6,891	41,567
CYP3A4	1TQN	Cytochrome P450	11,234	59,678
DPP-4	1X70	Peptidase	8,456	43,789
PPAR- $\gamma$	1FM9	Nuclear Receptor	5,678	34,567
CDK2	1AQ1	Kinase	9,123	48,234
5-HT2A	2VT4	GPCR	4,567	29,876
GSK-3 $\beta$	1J1B	Kinase	7,890	42,345

#### 4.1.2 Activity Threshold

We define active compounds as those with  $pIC_{50} \geq 6.0$  ( $IC_{50} \leq 1 \mu M$ ) and inactive compounds as those with  $pIC_{50} < 5.0$  ( $IC_{50} > 10 \mu M$ ).

### 4.1.3 Data Splitting

We use temporal splitting based on first publication date to avoid data leakage:

- **Training:** Publications before 2018 (70%)
- **Validation:** Publications in 2018 (15%)
- **Test:** Publications after 2018 (15%)

## 4.2 Baseline Methods

We compare against established virtual screening approaches:

### 4.2.1 Fingerprint-based Methods

- **ECFP4:** 2048-bit Extended Connectivity Fingerprints
- **MACCS:** 166-bit structural keys
- **AtomPairs:** Topological atom pairs
- **TopologicalTorsions:** 4-atom fragment descriptors

### 4.2.2 Machine Learning Methods

- **Random Forest:** With Morgan fingerprints ( $n_{\text{estimators}}=1000$ )
- **SVM:** With Tanimoto kernel ( $C=1.0$ ,  $\gamma=0.1$ )
- **XGBoost:** Gradient boosting ( $\text{max\_depth}=6$ ,  $n_{\text{estimators}}=1000$ )
- **Neural Fingerprints:** Duvenaud et al. architecture

### 4.2.3 Graph Neural Networks

- **GCN:** Graph Convolutional Network (4 layers)
- **GAT:** Graph Attention Network (4 heads, 4 layers)
- **MPNN:** Message Passing Neural Network
- **AttentiveFP:** Attentive Fingerprint

## 4.3 Implementation Details

### 4.3.1 Model Architecture

- **Hidden dimensions:** 256 for all layers
- **Message passing layers:** 6 layers
- **Meta-path types:** 15 pre-defined patterns
- **Attention heads:** 8 heads for meta-path aggregation
- **Dropout:** 0.3 applied to all layers

### 4.3.2 Training Configuration

- **Framework:** PyTorch 1.9 with DGL 0.7
- **Optimizer:** AdamW with learning rate 0.001
- **Batch size:** 256 molecules per batch
- **Epochs:** 200 with early stopping (patience=25)
- **Hardware:** 8x NVIDIA A100 GPUs (40GB each)
- **Training time:** 72 hours for full dataset

### 4.3.3 Evaluation Metrics

- **ROC-AUC:** Area under receiver operating characteristic curve
- **PRC-AUC:** Area under precision-recall curve
- **Enrichment Factor:**  $EF = (\text{Hits\_selected} / \text{Total\_selected}) / (\text{Hits\_database} / \text{Total\_database})$
- **BEDROC:** Bounded Enrichment at Different Ranks with Cutoff
- **Screening Speed:** Molecules processed per second

## 5. Results and Analysis

### 5.1 Overall Performance

Target	Method	ROC-AUC	PRC-AUC	EF@1%	EF@5%	BEDROC
EGFR	HeteroMPNN	<b>0.891</b>	<b>0.654</b>	<b>28.4</b>	<b>12.7</b>	<b>0.623</b>
	ECFP4+RF	0.847	0.587	23.1	10.8	0.556
	AttentiveFP	0.864	0.612	25.3	11.4	0.581
HIV-1 RT	HeteroMPNN	<b>0.863</b>	<b>0.598</b>	<b>24.7</b>	<b>11.3</b>	<b>0.595</b>
	ECFP4+RF	0.821	0.534	19.8	9.6	0.517
	AttentiveFP	0.839	0.556	21.2	10.1	0.542
AChE	HeteroMPNN	<b>0.876</b>	<b>0.621</b>	<b>26.8</b>	<b>12.1</b>	<b>0.608</b>
	ECFP4+RF	0.838	0.559	22.4	10.5	0.543
	AttentiveFP	0.851	0.578	23.9	11.0	0.567

Results shown for top 3 targets; full results available in supplementary material

### 5.2 Ablation Study

We analyze the contribution of each HeteroMPNN component:

Component Removed	Avg ROC-AUC Drop	Avg EF@1% Drop
Heterogeneous nodes	-4.2%	-18.3%
Heterogeneous edges	-2.8%	-12.7%
Meta-path aggregation	-3.1%	-15.4%
Multi-target training	-2.5%	-9.8%
Type-specific parameters	-3.7%	-16.2%

## 5.3 Molecular Diversity Analysis

We analyze the structural diversity of compounds identified by different methods:

### 5.3.1 Scaffold Diversity

- **HeteroMPNN**: 2,847 unique Bemis-Murcko scaffolds in top 1%
- **ECFP4+RF**: 2,234 unique scaffolds in top 1%
- **AttentiveFP**: 2,456 unique scaffolds in top 1%

### 5.3.2 Chemical Space Coverage

Using t-SNE visualization of molecular descriptors, HeteroMPNN identifies active compounds across broader chemical space regions compared to baselines.

### 5.3.3 Activity Cliff Navigation

HeteroMPNN successfully identifies 73% of activity cliff cases where structurally similar molecules have different activities, compared to 52% for fingerprint methods.

## 5.4 Computational Performance

Method	Training Time	Screening Speed	Memory Usage
HeteroMPNN	72h	<b>892 mol/s</b>	24.3 GB
AttentiveFP	48h	634 mol/s	18.7 GB
ECFP4+RF	2.5h	15,432 mol/s	8.9 GB
Neural FP	36h	745 mol/s	16.2 GB

## 5.5 Target-specific Analysis

Different targets benefit variably from heterogeneous modeling:

### 5.5.1 Kinases (EGFR, CDK2, GSK-3 $\beta$ )

- **Average improvement**: 5.2% ROC-AUC, 22.1% EF@1%
- **Key insight**: Heterogeneous modeling captures ATP-binding site variations

### 5.5.2 GPCRs (5-HT2A)

- **Average improvement:** 3.8% ROC-AUC, 16.7% EF@1%
- **Key insight:** Meta-paths capture pharmacophore patterns

### 5.5.3 Enzymes (AChE, BACE-1, DPP-4)

- **Average improvement:** 4.6% ROC-AUC, 19.3% EF@1%
- **Key insight:** Active site complementarity better modeled

## 5.6 Case Studies

### 5.6.1 EGFR Inhibitor Discovery

HeteroMPNN identified 47 novel EGFR inhibitors in top 0.1% that were missed by fingerprint methods. Subsequent experimental validation confirmed 34 compounds with  $IC_{50} < 1 \mu M$  (72.3% success rate).

### 5.6.2 Multi-target Compounds

The model identified 156 compounds with activity against multiple targets, suggesting potential for polypharmacology applications.

## 6. Discussion

### 6.1 Advantages of Heterogeneous Modeling

The superior performance of HeteroMPNN demonstrates several key advantages:

1. **Chemical Specificity:** Different atom and bond types require specialized treatment
2. **Biological Relevance:** Meta-paths capture pharmacologically important patterns
3. **Transferability:** Multi-target training improves generalization

### 6.2 Interpretability Analysis

#### 6.2.1 Atom Type Importance

Analysis of learned parameters reveals that nitrogen and oxygen types contribute most to binding predictions, consistent with their role in hydrogen bonding.

#### 6.2.2 Meta-path Significance

Donor-acceptor-hydrophobic patterns show highest attention weights, aligning with pharmacophore modeling principles.

### 6.3 Limitations and Challenges

#### 6.3.1 Computational Complexity



Heterogeneous modeling increases computational requirements by 2.8x compared to homogeneous GNNs.

### **6.3.2 Parameter Scaling**

The number of parameters scales quadratically with node/edge types, requiring careful regularization.

### **6.3.3 Type Definition**

Optimal atom/bond type definitions may be target-specific, suggesting need for adaptive typing strategies.

## **7. Future Directions**

### **7.1 Dynamic Type Learning**

Develop methods to automatically learn optimal atom/bond type definitions for specific targets or protein families.

### **7.2 3D Structure Integration**

Incorporate 3D molecular conformations and protein structure information into heterogeneous graph representations.

### **7.3 Active Learning**

Implement active learning strategies to iteratively improve models with minimal experimental validation.

### **7.4 Reaction Prediction**

Extend heterogeneous modeling to chemical reaction prediction and metabolite identification.

## **8. Conclusion**

We introduced HeteroMPNN, a heterogeneous message passing neural network for large-scale virtual screening. Our approach achieves significant improvements over traditional methods by explicitly modeling chemical diversity through heterogeneous graph representations and specialized message passing functions.

Key findings include:

- 15% average improvement in hit rate across 12 diverse protein targets
- Superior identification of structurally diverse active compounds
- Effective navigation of activity cliffs through learned representations
- Scalable performance suitable for large-scale database screening

The success of HeteroMPNN demonstrates the importance of incorporating chemical knowledge into graph neural network architectures for drug discovery applications. Our multi-target training strategy and

meta-path aggregation provide a foundation for future developments in graph-based virtual screening.

## Acknowledgments

We thank ChEMBL for providing comprehensive bioactivity data and the DGL team for their graph neural network framework. This work was supported by NIH grants R01-GM234567 and R01-CA345678, and NSF grant DBI-456789.

## References

1. Gaulton, A., et al. (2017). The ChEMBL database in 2017. *Nucleic Acids Research*, 45(D1), D945-D954.
2. Rogers, D., & Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5), 742-754.
3. Wang, X., et al. (2019). Heterogeneous graph attention network. *WWW Conference*.
4. Duvenaud, D., et al. (2015). Convolutional networks on graphs for learning molecular fingerprints. *NIPS*.
5. Xiong, Z., et al. (2020). Pushing the boundaries of molecular representation for drug discovery. *Journal of Medicinal Chemistry*.
6. Bajorath, J. (2002). Integration of virtual and high-throughput screening. *Nature Reviews Drug Discovery*, 1(11), 882-894.
7. Stumpfe, D., & Bajorath, J. (2012). Exploring activity cliffs in medicinal chemistry. *Journal of Medicinal Chemistry*, 55(7), 2932-2942.
8. Gilmer, J., et al. (2017). Neural message passing for quantum chemistry. *ICML*.
9. Bemis, G. W., & Murcko, M. A. (1996). The properties of known drugs. 1. Molecular frameworks. *Journal of Medicinal Chemistry*, 39(15), 2887-2893.
10. Truchon, J. F., & Bayly, C. I. (2007). Evaluating virtual screening methods: good and bad metrics for the "early recognition" problem. *Journal of Chemical Information and Modeling*, 47(2), 488-508.