



HAL
open science

Graph Neural Network Based on Molecular and Pharmacophoric Features for Drug Design Applications

Mariana Brito Azevedo, Luc Brun, Pierre Héroux, Jean-Luc Lamotte, Ronan Bureau, Alban Lepailler

► To cite this version:

Mariana Brito Azevedo, Luc Brun, Pierre Héroux, Jean-Luc Lamotte, Ronan Bureau, et al.. Graph Neural Network Based on Molecular and Pharmacophoric Features for Drug Design Applications. Workshop on Graph-based Representations in Pattern Recognition, Jun 2025, Caen, France. pp.47-57, <10.1007/978-3-031-94139-9_5>. <hal-05110964>

HAL Id: hal-05110964

<https://hal.science/hal-05110964v1>

Submitted on 13 Jun 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Graph Neural Network Based on Molecular and Pharmacophoric Features for Drug Design Applications

Mariana Brito Azevedo¹[0000-0003-4609-8800], Luc Brun¹[0000-0002-1658-0527], Pierre Héroux²[0000-0002-3509-2609], Jean-Luc Lamotte¹[0000-0001-6493-1769], Ronan Bureau³[0000-0001-9404-8117], and Alban Lepailleur³[0000-0003-0202-1588]

¹ Normandie Univ, ENSICAEN, CNRS, UNICAEN, GREYC UMR 6072, 14000 Caen, France {mariana.brito-azevedo, luc.brun, jean-luc.lamotte}@unicaen.fr

² Normandie Univ, UNIROUEN, UNIHAVRE, INSA Rouen, LITIS, 76000 Rouen, France {pierre.heroux}@univ-rouen.fr

³ Normandie Univ, UNICAEN, CERMN UR4258, F-14000 Caen, France {ronan.bureau,alban.lepailleur}@unicaen.fr

Abstract. Research fields that leverage relational data, like many others, have been significantly impacted by Deep Learning (DL) techniques, particularly Graph Neural Networks (GNNs). Among these fields, drug design, which aims to create new molecules with optimal affinities for specific targets, is a crucial step in the development of new medicinal drugs. In silico approaches in this area often rely on molecular graphs that encode the atoms and bonds of a molecule, without prior knowledge of the biological properties to be predicted. To address this limitation, pharmacophoric features are essential, as they contain structural information that captures important biological properties. These features have proven effective in tasks involving protein-ligand interactions. In this context, we propose the MCP-GNN model, which combines molecular representations with complete graphs of pharmacophoric features, both based on 2D information, to classify biological activity. Our experimental results demonstrate that this approach, using simple yet efficient techniques, achieves better performance than more complex architectures.

Keywords: Graph Neural Networks · Drug Design · Pharmacophores.

1 Introduction

Deep Learning (DL) methods have revolutionized some fields of research, such as text analysis and computer vision, and are increasingly impacting various domains. Among these domains, relational data forms a unique category where the embedding of input data is not in Euclidean space but corresponds to a set of entities and the relationships between them. Graph Neural Networks (GNNs) have quickly established themselves as a type of DL particularly suited for processing this data. The applications of GNNs include modern recommender systems,

computer vision, natural language processing, program analysis and urban intelligence, among others [7].

Across these applications, the task of predicting the biological activity of a molecule is particularly promising in drug design [3]. The success of GNNs in this field is attributed to the fact that chemical molecules can be intuitively represented as graphs. These graphs correspond to the fundamental entities manipulated by GNNs, which are capable of extracting meaningful features to provide accurate predictions. This approach is currently recognized for delivering better performance than other techniques, including traditional DL methods, in tasks related to drug design, such as structure-based drug-target interaction (DTI) prediction [19] and kinase-profiling prediction [9].

The prediction of molecular properties by GNNs may be based on molecular graphs where nodes represent atoms and edges represent chemical bonds. However, this type of representation is application-agnostic, and different GNNs based on the same molecular representation can predict physical (e.g., boiling point), chemical, or biological properties of a molecule [5].

In the context of drug design applications, an alternative representation is based on graphs with pharmacophoric features. As described by [14], a pharmacophore represents the essential structural information of a molecule that enables it to bind to a specific biological target [15]. From this perspective, a graph with pharmacophoric features contains meaningful and essential characteristics for identifying active molecules and has proven effective in tasks related to protein-ligand interactions [10].

In this paper, we introduce a novel approach in drug design called Molecular and Complete Pharmacophoric Features Graph Neural Network (MCP-GNN). This model addresses the limitations of previous methods [11, 18] by offering a simpler yet more effective solution. Our goal is to classify the biological activity of molecules independently of their 3D conformation, i.e. based on their 2D representation, categorizing them as active or inactive, with active molecules holding potential for developing new drugs. This method achieves state-of-the-art results in classifying BCR-ABL data and nine popular kinase datasets.

This article is organized as follows: Section 2 introduces the concepts of pharmacophores, pharmacophoric features and GNNs, along with other approaches that have developed models for similar purposes. Section 3 describes the methodology for developing the proposed model, including data preprocessing and model architecture design. Section 4 presents the datasets and configurations used to conduct the experiments and details the results obtained from our proposal, comparing them with other techniques in the literature.

2 Related Works

2.1 Graph representations for predicting biological activity

The biological action of a molecule is generally not induced by the entire molecule but rather through the interactions of some of its parts with biological targets,

such as proteins. From this perspective, it is possible to use different levels of graph representation to predict biological activity, making the concepts of pharmacophores and pharmacophoric features essential.

A pharmacophore can be defined as "the ensemble of steric and electronic features that are necessary to ensure optimal supramolecular interactions with a specific biological target structure and to trigger (or block) its biological response" [17]. Hence, a molecule, often referred to as a ligand, whose characteristics match those of a pharmacophore, should interact with the biological target on which the pharmacophore has been defined. In this sense, a pharmacophore may be considered a molecular pattern that characterizes biological activity.

However, the pharmacophores corresponding to a given biological target are usually unknown. An alternative solution involves using pharmacophoric features, which correspond to connected groups of atoms in the molecule that contribute to favorable interactions with target molecules [12]. In this sense, it is possible to construct a pharmacophore graph, where the vertices encode these features, and the edges encode the types of relationships between these functional groups. Thus, a pharmacophore graph can be understood as a super-graph that should contain all pharmacophores associated with a given biological target.

Different types of graph representations have been used to predict biological activity. In [11], the authors propose the model RG-MPNN. In this architecture, two levels of graph representation are used. The first level uses molecular graphs while the second level uses pharmacophore reduced graphs. These reduced graphs present 18 types of pharmacophoric features, and the edges represent the total number of chemical bonds shared between two pharmacophore nodes plus one.

The AttentiveFP model [18] is based on molecular representation and employs graph attention mechanisms. This model utilizes 9 types of atomic features and 4 types of bond features at the molecular level. Non-local effects at the intramolecular level, which refer to the interactions and relationships within a single molecule, are captured through attention mechanisms. This strategy enables accurate predictions of chemical, biological, and physiological properties.

However, it should be noted that these methods have some limitations. In AttentiveFP, only the molecular and intramolecular representation levels are used and learned, without incorporating pharmacophoric features, which may limit its performance in predicting biological activity. Moreover, in RG-MPNN, the model utilizes pharmacophoric features, but the edges in the pharmacophoric graph only represent adjacency between the corresponding features in the molecular graph. Since biological activity may involve non-adjacent features, the pharmacophoric graph fails to act as a true super-graph of all relevant pharmacophores.

2.2 GNNs

A GNN is a neural network that performs a sequence of transformations on an input graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w)$, where \mathcal{V} represents the set of nodes, \mathcal{E} represents the set of edges, and $w : \mathcal{E} \rightarrow \mathbb{R}$ assigns a real value to each edge. Alternatively, a graph may be represented as $\mathcal{G} = (A, X)$ where A and X encode the adjacency

matrix and the vertex attributes of \mathcal{G} , respectively. In this case, the weight function w is encoded by the entries of the matrix A .

In GNNs, convolution operations are used to aggregate the value of each node with its k -hop neighborhood, where k is a parameter of the convolution. Spectral convolutions [4] based on the graph’s Laplacian are primarily designed for varying signals on a fixed graph topology. Spatial convolutions do not have this limitation and are typically described through the message-passing neural network (MPNN) framework [7], where nodes iteratively exchange and update information based on their local neighborhoods.

GNNs usually apply multiple one hop convolutions in order to encode within each vertex information about its l -hop neighborhood, where l corresponds to the number of convolution layers. However, for graph classification, it is necessary to map each input graph to a single output value. This operation is usually performed by mapping the graph obtained by the last convolutional layer onto a vector of fixed size using an operation called global pooling [16] defined by Equation 1. Here, \mathbf{h}_i^L corresponds to the feature vector of node i and \mathcal{V}^L corresponds to the vertex set of the graph defined at the last convolutional layer.

$$\mathbf{h}_{\mathcal{G}} = POOL(\mathbf{h}_i^L | i \in \mathcal{V}^L) \quad (1)$$

The vector $\mathbf{h}_{\mathcal{G}}$ having a fixed size, the global pooling must be invariant to the number of nodes in a graph. Additionally, since the set of nodes in a graph is unordered, the pooling operation must also be permutation-invariant. Common operations that satisfy these conditions include sum, mean, max/min, or their combinations computed on the vertex set.

The pooled vector $\mathbf{h}_{\mathcal{G}}$ is then processed through one or more Multi-Layer Perceptrons (MLPs), followed by a softmax function for classification. Additionally, $\mathbf{h}_{\mathcal{G}}$ can be constructed by pooling features from multiple GNN layers using residual connections, creating a multi-scale representation that reduces smoothing effects. However, it is important to note that by reducing node features to a single vector, this approach loses structural and local information, which may be critical for large, complex graphs.

3 Methodology

3.1 Data preprocessing

Our preprocessing generates two graph representations for each molecule: a molecular graph and a pharmacophore graph.

Each molecular graph’s nodes and edges are enriched as follows: nodes include atom type, degree, charge, hybridization, bonded hydrogens, and aromaticity. Edges encode bond type, conjugation, and ring participation.

Our pharmacophore graph is defined as a Complete Graph of Pharmacophoric Features (CGPFs). The nodes are derived from the 2D structure of a molecule using the extraction method for pharmacophoric features proposed in [13], the workflow tool Norns [1] and the library OpenBabel to manage the chemical data.

Figure 1 illustrates the transformation of a molecule into a CGPF. Let us note that a ligand should not be too large in order to interact with a protein. Consequently, the associated CGPFs usually corresponds to a small complete graph. The mean size of CGPFs in all datasets considered in Section 4 does not exceed 14 nodes.

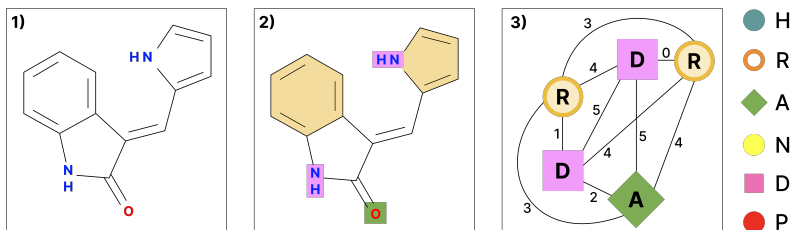


Fig. 1. Transformation of a molecule into a CGPF. Starting from the initial molecule (1), the pharmacophore features are identified (2), and then the topological distances between each of the features are computed (3) to obtain the CGPF.

In this graph representation, a node indicates one of six possible pharmacophoric features: hydrogen-bond acceptors (A) or donors (D), negatively (N) and positively (P) charged ionizable groups, hydrophobic regions (H) and aromatic rings (R). The type of each node is thus encoded by a one-hot vector of dimension $d = 6$. In addition, to represent the edges, following previous studies in the drug design field [13, 6], the topological distances between pairs of pharmacophoric features in the chemical molecule were used. Such distances are defined as the number of bonds in the shortest path connecting any pair of pharmacophoric feature.

It is worth noting that our pharmacophore graph is based on a selection of traditional pharmacophoric features that are considered as the most important based on pharmacophoric modelling theory [8], ensuring that the representation is both chemically and biologically meaningful. Moreover, the use of a complete graph allows to obtain the interactions between non adjacent pharmacophoric features in order to better capture the phenomenon of ligand/target interaction.

3.2 Model architecture

The second stage consists in defining the GNN model, whose architecture is illustrated in Figure 2.

After obtaining the two graph representations, a linear layer is applied to the edge features of each graph, converting the one-hot encoding vector into a single scalar that captures their importance.

For CGPFs, the layer learns the relevance of distances between pharmacophoric features, with each distance associated with a one-hot vector of maximum size $d_{max} = 118$. For molecular graphs, the layer focuses on learning the

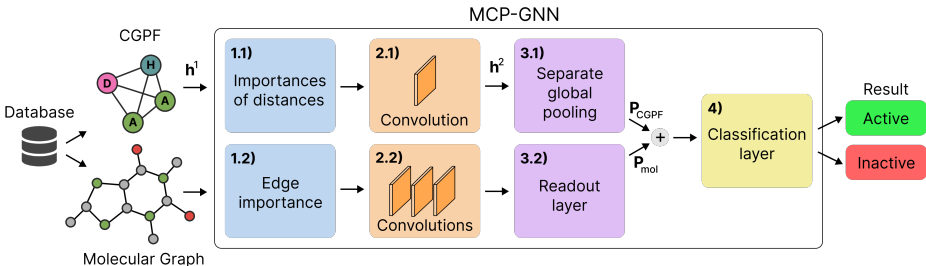


Fig. 2. Architecture of MCP-GNN model. The stages include learning the importance of edge features (1.1 and 1.2), graph convolution layers (2.1 and 2.2), implementation of a separate global pooling for each type of node for CGPFs (3.1), implementation of a readout layer for molecular graphs (3.2) and a classification layer (4).

importance of features related to atomic bonds, which is a vector of size $d = 6$ that encodes bond type, conjugation, and ring participation. This edge weighting step has the advantage of requiring few parameters, which is well-suited to biological datasets that are typically not extensive.

The relevance of edges is then incorporated into adjacency matrices to perform convolutions weighted by the edges’ relevance. For the representation with pharmacophoric features, since complete graphs are used, the computation of a single graph convolution layer is sufficient for each node to aggregate information from all its neighbors. For the molecular graph representation, multiple layers of graph convolution are applied.

The third stage of our model consists of performing global pooling for each graph representation to obtain a global representation of each graph as a vector of fixed dimension. As mentioned in Section 2.2, global pooling is typically performed using max, mean, or sum operations applied to the entire set of vertices of the graph. Since this operation can lead to a significant loss of information, for CGPFs, which have a limited number of vertex labels, a separate global pooling strategy was applied to mitigate this loss. Let $\mathcal{V}_1, \dots, \mathcal{V}_6$ define a partition of the initial vertex set \mathcal{V} , where \mathcal{V}_j contains all vertices of \mathcal{V} of type j . We first apply a global mean pooling on each type of vertex:

$$\mathbf{P}_j = \frac{1}{|\mathcal{V}_j|} \sum_{i \in \mathcal{V}_j} \mathbf{h}_i^2 \quad (2)$$

where \mathbf{h}_i^2 defines the features of node i after our convolution step (Figure 2). Then, the global pooling for CGPFs is defined as the concatenation of all partial global poolings:

$$\mathbf{P}_{CGPF} = \text{concat}(\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_6) \quad (3)$$

For the molecular representation, it was not feasible to apply the same strategy due to the great diversity of atom types. Therefore, a readout layer was

implemented, where global mean pooling and global max pooling were concatenated to obtain a global representation at the atomic level, defined as P_{mol} in Figure 2.

Finally, the results obtained from each pooling operation (P_{CGPF} and P_{mol}) are concatenated, and an MLP is used to classify molecules as active or inactive.

4 Experiments

4.1 Datasets, training and evaluation

We consider the BCR-ABL dataset and nine kinase family datasets (Table 2). All datasets consist of biological data where the action of the input molecule on a target must be classified as either active (label 0) or inactive (label 1). The BCR-ABL dataset was obtained from the study proposed by [6], while the kinase datasets were the same as those used in the study by [11].

To compare the performance of the proposed model MCP-GNN with RG-MPNN and AttentiveFP, the same rules were followed for training and evaluation. The data was split into 80% for training, 10% for validation, and 10% for testing using a stratified division. The training step was executed for 200 epochs with an early stopping strategy, and the following metrics were computed: accuracy (ACC), Area Under the ROC Curve (AUC), Matthews Correlation Coefficient (MCC), and mean F1 Score (MEAN F1). Additionally, a 10-fold cross-validation strategy was employed to assess the model’s ability to generalize.

A hyperparameter optimization for the AttentiveFP and MCP-GNN models was conducted using the Optuna library [2]. The hyperparameters of the RG-MPNN model were set to those proposed by the authors for each dataset. The hyperparameters optimized in this work included: learning rate and weight decay, pertaining to the optimizer; batch size; and hidden size, referring to the number of neurons in the model’s intermediate layers. For MCP-GNN, the number of graph convolutions for the molecular representation was also considered.

4.2 Ablation studies

We conducted an ablation study on three balanced datasets, with statistics provided in Table 2 (first column). To evaluate our graph representations, we assessed our model’s performance using each representation individually. For each, we tested three variations: the full architecture (Original), a version without edge importance calculation (using original distance values for CGPFs and unweighted edges in the molecular representation), and a version without pooling strategies (using a single global mean pooling for both representations).

Table 1 shows the AUC values for each dataset and ablation scenario. Non-ablation results are on the third line for each dataset. Notably, introducing edge relevance and our enhanced pooling strategy significantly improves the model with CGPFs by more than one standard deviation, though improvements are less pronounced with the molecular graph representation.

Table 1. AUC values obtained during ablation studies

Dataset	Graph representation	Original	Without edge importance	Without pooling strategy
AKT1	Only molecules	0.944±0.010	0.944±0.007	0.935±0.016
	Only CGPFs	0.924±0.012	0.905±0.014	0.903±0.015
	Complete	0.949±0.009	-	-
BCR-ABL	Only molecules	0.943±0.015	0.945±0.019	0.938±0.030
	Only CGPFs	0.925±0.023	0.891±0.026	0.896±0.030
	Complete	0.949±0.017	-	-
CDK2	Only molecules	0.904±0.019	0.905±0.018	0.895±0.021
	Only CGPFs	0.827±0.032	0.762±0.040	0.773±0.029
	Complete	0.905±0.013	-	-

We observed that the molecular representation slightly outperforms the pharmacophoric one, likely due to its larger number of node features and parameters. For instance, in the BRAF dataset (Table 2), our model had 225,000 parameters, with 79% related to molecular graphs and 21% to CGPFs. Combining both representations yields a better AUC than either alone, indicating complementary information.

Additionally, regarding model complexity, AttentiveFP and RG-MPNN have approximately 534,000 and 1,950,000 parameters for the same dataset, highlighting our model’s efficiency despite fewer parameters (see Table 2).

4.3 Evaluation

Table 2 presents the performance of the three models in classifying biological activity across kinase datasets and the BCR-ABL dataset. The results highlight the superior performance of the proposed model in most scenarios.

Notably, the MCP-GNN model shows significant improvements for smaller datasets such as CK1, MAP4K2, and CDK2, achieving performance gains of between 2% and 8% across multiple metrics. For larger datasets with a higher number of molecules, the performance gains are less pronounced. Nevertheless, the proposed model consistently delivers better results, demonstrating its robustness and generalization capability across different datasets.

5 Conclusion

The proposed model integrates molecular and pharmacophoric features using specific architectural techniques. Our experiments show that this approach effectively captures complex molecular relationships, making it a promising tool for drug design. Future work will involve linking pharmacophoric features to atoms to provide more information, adding 3D positional data to each node, and introducing explainability techniques to enhance understanding of its decision-making process.

Table 2. Results with BCR-ABL and kinase datasets. The **Dataset** column displays the database name in the first row, the total number of molecules in the second row, and the distribution of active/inactive molecules in the third row. In addition, bold values represent the best performance, while blue values indicate the second best.

Dataset	Model	ACC	AUC	MCC	MEAN F1
AKT1	RG-MPNN	0.875±0.016	0.935±0.012	0.748±0.032	0.873±0.016
3967 mol.	AttentiveFP	0.877±0.010	0.941±0.006	0.753±0.021	0.872±0.018
2175/1792	MCP-GNN	0.892±0.012	0.949±0.009	0.783±0.023	0.890±0.012
AURKA	RG-MPNN	0.829±0.030	0.889±0.034	0.644±0.063	0.821±0.034
3886 mol.	AttentiveFP	0.846±0.016	0.916±0.013	0.681±0.033	0.840±0.017
2290/1596	MCP-GNN	0.860±0.016	0.928±0.014	0.710±0.034	0.854±0.017
BCR-ABL	RG-MPNN	0.868±0.025	0.936±0.015	0.739±0.049	0.868±0.028
1479 mol.	AttentiveFP	0.858±0.052	0.926±0.028	0.721±0.104	0.877±0.027
773/706	MCP-GNN	0.883±0.024	0.949±0.017	0.761±0.050	0.878±0.025
BRAF	RG-MPNN	0.895±0.020	0.934±0.033	0.709±0.060	0.877±0.023
4824 mol.	AttentiveFP	0.922±0.013	0.955±0.014	0.783±0.038	0.894±0.024
3629/1132	MCP-GNN	0.932±0.011	0.960±0.013	0.808±0.033	0.901±0.016
BTK	RG-MPNN	0.898±0.023	0.944±0.018	0.749±0.051	0.870±0.023
3640 mol.	AttentiveFP	0.902±0.018	0.947±0.015	0.759±0.043	0.876±0.024
1892/1748	MCP-GNN	0.915±0.023	0.955±0.013	0.788±0.053	0.886±0.021
CDK2	RG-MPNN	0.774±0.040	0.853±0.036	0.552±0.081	0.763±0.035
2603 mol.	AttentiveFP	0.796±0.021	0.873±0.015	0.595±0.041	0.791±0.023
1303/1300	MCP-GNN	0.830±0.017	0.905±0.013	0.662±0.035	0.827±0.032
CK1	RG-MPNN	0.832±0.049	0.780±0.064	0.397±0.190	0.626±0.120
807 mol.	AttentiveFP	0.827±0.039	0.833±0.055	0.439±0.111	0.699±0.055
155/632	MCP-GNN	0.825±0.036	0.852±0.029	0.517±0.062	0.718±0.044
MAP4K2	RG-MPNN	0.768±0.044	0.770±0.083	0.413±0.165	0.679±0.097
897 mol.	AttentiveFP	0.777±0.045	0.834±0.049	0.470±0.113	0.742±0.059
280/617	MCP-GNN	0.791±0.063	0.846±0.069	0.511±0.148	0.748±0.049
EGFR	RG-MPNN	0.854±0.014	0.923±0.013	0.709±0.028	0.854±0.015
8788 mol.	AttentiveFP	0.860±0.014	0.938±0.009	0.721±0.028	0.860±0.014
4513/4275	MCP-GNN	0.874±0.011	0.946±0.006	0.749±0.023	0.874±0.011
PIM1	RG-MPNN	0.881±0.018	0.942±0.012	0.708±0.042	0.869±0.017
4507 mol.	AttentiveFP	0.908±0.014	0.956±0.011	0.763±0.039	0.883±0.019
3321/1186	MCP-GNN	0.919±0.012	0.970±0.009	0.786±0.031	0.892±0.016

References

1. Norns packages. <https://projects.greyc.fr/chemoinfo/software/>
2. Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M.: Optuna: A next-generation hyperparameter optimization framework. CoRR **abs/1907.10902** (2019)
3. Bechelli, S., Delhommelle, J.: Ai’s role in pharmaceuticals: Assisting drug design from protein interactions to drug development. Artificial Intelligence Chemistry **2**(1), 100038 (2024)
4. Bo, D., Wang, X., Liu, Y., Fang, Y., Li, Y., Shi, C.: A survey on spectral graph neural networks (2023)

5. Chen, B., Pan, Z., Mou, M., Zhou, Y., Fu, W.: Is fragment-based graph a better graph-based molecular representation for drug design? a comparison study of graph-based models. *Computers in Biology and Medicine* **169**, 107811 (2024)
6. Geslin, D., Lepailleur, A., Manguin, J.L., Vo, N.V., Lamotte, J.L., Cuissart, B., Bureau, R.: Deciphering a pharmacophore network: A case study using bcr-abl data. *Journal of Chemical Information and Modeling* **62**(3), 678–691 (2022)
7. Gilmer, J., Schoenholz, S.S., Riley, P.F., Vinyals, O., Dahl, G.E.: Neural message passing for quantum chemistry. In: *Proceedings of the 34th International Conference on Machine Learning - Volume 70*. p. 1263–1272. ICML'17, JMLR.org (2017)
8. Giordano, D., Biancaniello, C., Argenio, M.A., Facchiano, A.: Drug design by pharmacophore and virtual screening approach. *Pharmaceuticals (Basel, Switzerland)* **15**(5), 646 (2022)
9. Gu, S., Liu, H., Liu, L., Hou, T., Kang, Y.: Artificial intelligence methods in kinase target profiling: Advances and challenges. *Drug Discovery Today* **28**(11), 103796 (2023)
10. Kengkanna, A., Ohue, M.: Enhancing property and activity prediction and interpretation using multiple molecular graph representations with mmgx. *Communications Chemistry* **7**, 74 (2024)
11. Kong, Y., Zhao, X., Liu, R., Yang, Z., Yin, H., Zhao, B., Wang, J., Qin, B., Yan, A.: Integrating concept of pharmacophore with graph neural networks for chemical property prediction and interpretation. *Journal of Cheminformatics* **14**(52), 103796 (2022)
12. Lehembre, E., Giovannini, J., Geslin, D., Lepailleur, A., Lamotte, J.L., Auber, D., Ouali, A., Cremilleux, B., Zimmermann, A., Cuissart, B., Bureau, R.: Towards a partial order graph for interactive pharmacophore exploration: extraction of pharmacophores activity delta. *Journal of Cheminformatics* **15**(1), 116 (Nov 2023)
13. Métivier, J.P., Cuissart, B., Bureau, R., Lepailleur, A.: The pharmacophore network: A computational method for exploring structure–activity relationships from a large chemical data set. *Journal of Medicinal Chemistry* **61**(8), 3551–3564 (2018)
14. Muhammed, M.T., Aki-yalcin, E.: Pharmacophore modeling in drug discovery: Methodology and current status. *Journal of the Turkish Chemical Society Section A: Chemistry* **8**(3), 749–762 (2021)
15. Oselusi, S.O., Dube, P., Odugbemi, A.I., Akinyede, K.A., Ilori, T.L., Egieyeh, E., Sibuyi, N.R., Meyer, M., Madiehe, A.M., Wyckoff, G.J., Egieyeh, S.A.: The role and potential of computer-aided drug discovery strategies in the discovery of novel antimicrobials. *Computers in Biology and Medicine* **169**, 107927 (2024)
16. Stanovic, S., Gaüzère, B., Brun, L.: Maximal independent vertex set applied to graph pooling. In: Krzyzak, A., Suen, C.Y., Torsello, A., Nobile, N. (eds.) *Structural, Syntactic, and Statistical Pattern Recognition*. pp. 11–21. Springer International Publishing, Cham (2022)
17. Wermuth, C.G., Ganellin, C.R., Lindberg, P., Mitscher, L.A.: Glossary of terms used in medicinal chemistry (iupac recommendations 1998). *Pure and Applied Chemistry* **70**(5), 1129–1143 (1998)
18. Xiong, Z., Wang, D., Liy, X., Zhong, F., Wan, X., Li, X., Li, Z., Luo, X., Chen, K., Jiang, H., Zheng, M.: Pushing the boundaries of molecular representation for drug discovery with the graph attention mechanism. *J Med Chem* **63**(16), 8749–8760 (2020)
19. Zhang, Z., Chen, L., Zhong, F., Wang, D., Jiang, J., Zhang, S., Jiang, H., Zheng, M., Li, X.: Graph neural network approaches for drug-target interactions. *Current Opinion in Structural Biology* **73**, 102327 (2022)