
This is an electronic reprint of the original article.
This reprint may differ from the original in pagination and typographic detail.

Li, Anchen; Casiraghi, Elena; Rousu, Juho

Chemical reaction enhanced graph learning for molecule representation

Published in:
Bioinformatics (Oxford, England)

DOI:
[10.1093/bioinformatics/btae558](https://doi.org/10.1093/bioinformatics/btae558)

Published: 01/10/2024

Document Version
Publisher's PDF, also known as Version of record

Published under the following license:
CC BY

Please cite the original version:
Li, A., Casiraghi, E., & Rousu, J. (2024). Chemical reaction enhanced graph learning for molecule representation. *Bioinformatics (Oxford, England)*, 40(10), 1-9. Article btae558.
<https://doi.org/10.1093/bioinformatics/btae558>

This material is protected by copyright and other intellectual property rights, and duplication or sale of all or part of any of the repository collections is not permitted, except that material may be duplicated by you for your research use or educational purposes in electronic or print form. You must obtain permission for any other use. Electronic or print copies may not be offered, whether for sale or otherwise to anyone who is not an authorised user.

Data and text mining

Chemical reaction enhanced graph learning for molecule representation

Anchen Li ^{1,*}, Elena Casiraghi ^{1,2,3,4}, Juho Rousu ¹

¹Department of Computer Science, Aalto University, Espoo, 02150, Finland

²AnacletoLab, Dipartimento di Informatica "Giovanni degli Antoni", University of Milan, Milan, 20133, Italy

³Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, 94720, United States

⁴ELLIS, European Laboratory for Learning and Intelligent Systems, Milan Unit (University of Milan), Milan, 20133, Italy

*Corresponding author. Department of Computer Science, Aalto University, Espoo, 02150, Finland. E-mail: anchen.li@aalto.fi (A.L.)

Associate Editor: Jonathan Wren

Abstract

Motivation: Molecular representation learning (MRL) models molecules with low-dimensional vectors to support biological and chemical applications. Current methods primarily rely on intrinsic molecular information to learn molecular representations, but they often overlook effectively integrating domain knowledge into MRL.

Results: In this article, we develop a reaction-enhanced graph learning (RXGL) framework for MRL, utilizing chemical reactions as domain knowledge. RXGL introduces dual graph learning modules to model molecule representation. One module employs graph convolutions on molecular graphs to capture molecule structures. The other module constructs a reaction-aware graph from chemical reactions and designs a novel graph attention network on this graph to integrate reaction-level relations into molecular modeling. To refine molecule representations, we design a reaction-based relation learning task, which considers the relations between the reactant and product sides in reactions. In addition, we introduce a cross-view contrastive task to strengthen the cooperative associations between molecular and reaction-aware graph learning. Experiment results show that our RXGL achieves strong performance in various downstream tasks, including product prediction, reaction classification, and molecular property prediction.

Availability and implementation: The code is publicly available at <https://github.com/coder-ACAC/RLM>.

1 Introduction

Molecule representation learning (MRL) techniques are crucial for combining machine learning with biological and chemical sciences (Yi *et al.* 2022). MRL encodes molecules as low-dimensional vectors. These vectors retain molecule information, facilitating their use as features in downstream applications (e.g. product prediction, reaction classification, and molecular property prediction). A variety of MRL methods have been proposed, which roughly fall into the following two categories.

One school is SMILES-based methods (Fabian *et al.* 2020), which utilize SMILES strings as input and employ natural language models as their base architectures. However, they struggle with capturing molecule structures. The other school treats molecule topology as a graph, and models molecules with graph neural networks (GNNs) (Xu *et al.* 2021). Although GNN-based methods generally outperform SMILES-based ones, they typically focus on designing GNN architectures, neglecting the efficient integration of domain knowledge.

Recent studies (Wang *et al.* 2022a) use chemical reactions as domain knowledge for MRL. Typically, reactions are represented by equations, with reactants on the left side and products on the right (cf Definition 2 in Section Preliminaries). These methods first learn molecule embeddings from molecular graphs and then optimize embeddings by equating the sum of reactant embeddings with the sum of product embeddings for each

reaction. Despite effectiveness, we argue that they face at least one of the following issues.

Firstly, these reaction-based methods treat molecules as isolated data instances and rely solely on molecule structures for representation, which ignores the insights from molecule relations inherent in chemical reactions. For example, molecules involved in the same reaction (as reactants/products) may exhibit greater similarities and correlations with each other than with molecules from different reactions. To illustrate this reaction-related relation, we construct a reaction-aware graph (cf Definition 4 in Section Preliminaries) based on a reaction set, as shown in Fig. 1a. In this graph, nodes are molecules and edges denote molecule relations driven by reactions. For molecule *A*, its first-order neighbors (molecules *F*, *G*, and *H*) represent products that can be derived from *A* through reactions. Molecule *A*'s second-order neighbors (molecules *B*, *C*, and *D*) suggest a property/structure similarity with *A*, inferred from shared reaction products. Moreover, molecule *B* is likely more similar to *A* than *C* or *D*, as evidenced by a greater overlap in the reaction products. These analyses inspire us to consider the potential benefits of incorporating molecule relations from the reaction-aware graph into MRL.

Secondly, these methods ignore the transformation relation learning between reactants and products. Their assumption that the summed embeddings of reactants and products should be

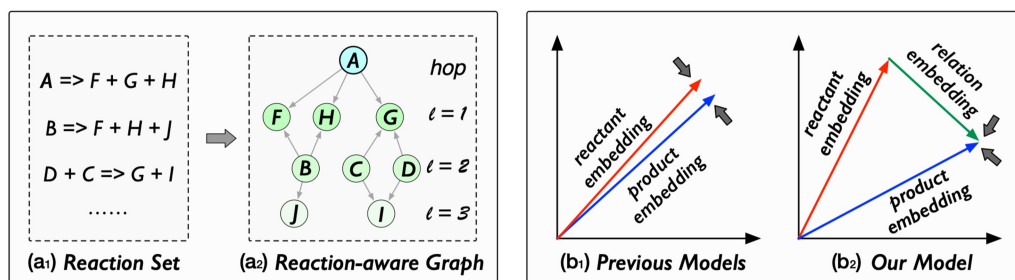


Figure 1. (a₁) and (a₂): An illustration of the reaction-aware graph. (b₁) and (b₂): Comparisons of previous models and our model in considering reactant-product relations of the chemical reaction.

equal (as shown in Fig. 1b1) essentially reduces all reactions to an identity transformation, which oversimplifies the complexity of chemical processes. In reality, reactions involve various changes, such as the number of bonds (e.g. breaking old bonds and forming new ones) and energy variations (e.g. endotherms and exotherms) before and after the reaction. The assumption in current studies fails to model these changes.

Motivated by these gaps, we introduce a reaction-enhanced graph learning framework (RXGL) for MRL. In the molecule modeling stage, we design dual graph learning modules. The first module utilizes graph convolutions on molecular graphs to capture the structural information of molecules. The second module first involves a reaction-aware graph and then creates a GNN to extract reaction-level molecular relations for molecule feature learning. In the optimization stage, we introduce a reaction-based relation learning method that considers the relation between reactants and products in chemical reactions. Specifically, we employ a memory network (Miller et al. 2016) to learn a latent relation vector that connects reactant and product embeddings (as shown in Fig. 1b2). Through the delicate key and memory components in this network, the learned relation vectors could capture the hidden semantic correlations between the reactant and product in each reaction. Furthermore, we incorporate a cross-view contrastive learning task to enhance molecule representations in dual graph modeling. This task treats the molecular and reaction-aware graphs as distinct yet correlated views. By employing contrastive learning, we capture cooperative associations between views, thereby integrating their agreements into molecule representations.

Our contributions are summarized as follows:

- We propose a GNN-based MRL framework RXGL, which introduces the chemical reaction-aware graph to assist in learning molecule representations.
- We devise a reaction-based relation modeling approach, which learns the relations between reactant and product sides in chemical reactions to guide MRL.
- We design a cross-view contrastive learning task to enhance the cooperative association between the molecular graph and reaction-aware graph modeling views.
- Experiment results show that molecule embeddings learned by our RXGL benefit various tasks, i.e. product prediction, reaction classification, and molecular property prediction.

2 Preliminaries

This section introduces notations and formulates the problem.

2.1 Notations

Our method RXGL uses the molecular graph and reaction-aware graph for molecule representation learning. We first introduce the definition of the molecular graph:

Definition 1. Molecular Graph. Given a molecule set \mathcal{M} , each molecule $m \in \mathcal{M}$ has a graph structure $\mathcal{G}_m = (\mathcal{V}_m, \mathcal{E}_m)$, which includes an atom node set \mathcal{V}_m and a bond edge set \mathcal{E}_m . Each atom $a_i \in \mathcal{V}_m$ is represented by a vector encoding its features, and bond $b_i \in \mathcal{E}_m$ is represented by its bond type.

The reaction-aware graph is constructed based on chemical reactions. The chemical reaction is defined as follows:

Definition 2. Chemical Reaction. Consider a molecule set \mathcal{M} and a reaction set \mathcal{X} , a reaction $x_i \in \mathcal{X}$ defines a transformation from a reactant set $R_i \subset \mathcal{M}$ to a product set $P_i \subset \mathcal{M}$, denoted as $x_i: R_i \rightarrow P_i$, where sets $R_i = \{r_{i_1}, r_{i_2}, \dots\}$ and $P_i = \{p_{i_1}, p_{i_2}, \dots\}$ represent the reactants and products involved in reaction x_i , respectively.

According to the definition of the chemical reaction, we define the reaction-level relation between molecules:

Definition 3. Reaction-level Molecule Relation. For a given reaction $x_i: R_i \rightarrow P_i$, a reaction-level molecule relation is defined between any reactant $r_{i_a} \in R_i$ and product $p_{i_b} \in P_i$.

Based on the reaction and reaction-level molecule relation, the reaction-aware graph is defined as follows:

Definition 4. Reaction-aware Graph. A reaction-aware graph is defined as $\mathcal{G}_x = (\mathcal{V}_x, \mathcal{E}_x)$, where $\mathcal{V}_x = R \cup P$ is the node set, and $\mathcal{E}_x \subset R \times P$ is the edge set including reaction-level molecule relations between reactants R and products P .

Figure 1a shows an example of the reaction-aware graph. It is worth noticing that the construction of the reaction-aware graph draws inspiration from the metabolic networks (Wagner and Fell 2001) and chemical reaction networks (Wen et al. 2023) in synthetic biology. Its structural design facilitates the modeling of reaction-level molecule relations.

2.2 Problem formulation

Given the molecule m_i 's graph structure \mathcal{G}_{m_i} , reaction-aware graph \mathcal{G}_x and a chemical reaction set \mathcal{X} , we aim to learn m_i 's

representation \mathbf{h}_i . Then, learned molecule representation \mathbf{h} , can be applied to various downstream tasks.

3 Reaction-enhanced graph learning (RXGL)

This section presents our RXGL method, as shown in Fig. 2. We first introduce dual graph learning modules to model molecule representations. Then, we propose a reaction-based relation learning task and a cross-view contrastive learning task to optimize molecule representations.

3.1 Molecule representation modeling

We introduce dual modules: a molecular graph learning module and a reaction-aware graph learning module. These modules are designed to integrate atom-level and reaction-level features into molecule representations, respectively.

3.1.1 Molecular graph learning module

This module first uses GNNs to model atoms and then utilizes pooling operations on the molecular graph to inject atom-level features into molecule representations, as shown in Fig. 3a. For molecule m 's graph structure $\mathcal{G}_m = (\mathcal{V}_m, \mathcal{E}_m)$, we model representation \mathbf{a}_i^k of atom $a_i \in \mathcal{V}_m$ in the k th GNN layer:

$$\mathbf{a}_i^k = \text{Aggregate}(\mathbf{a}_j^{k-1} | a_j \in \mathcal{N}_i \cup a_i), \quad (1)$$

where \mathcal{N}_i is the atoms connected to atom a_i in graph \mathcal{G}_m . The choice of Aggregate function is the key to designing GNN, leading to the proposal of various structures (Li *et al.* 2022, 2024). We utilize four common GNNs [i.e. GCN (Kipf and Welling 2017), GAT (Veličković *et al.* 2018), SAGE (Hamilton *et al.* 2017), and TAG (Du *et al.* 2017)] for this module. We introduce these GNNs in the Supplementary Materials.

To describe initial feature \mathbf{a}_i^0 of atom a_i in \mathcal{G}_m , we use four types of atom properties: element type, charge, presence in an aromatic ring, and the number of attached hydrogen atoms. Each type of property is represented as a one-hot embedding, and four embeddings are concatenated as the initial atom feature. Note that we do not represent bond types with explicit feature vectors because the bond type could be inferred by the connected atom features (Wang *et al.* 2022a).

By stacking K GNN layers, a pooling function is used to generate atom-level molecule representation \mathbf{h}^A , as follows:

$$\mathbf{h}^A = \text{Pool}(\mathbf{a}_i^K | a_i \in \mathcal{V}_m) = \sum_{a_i \in \mathcal{V}_m} \mathbf{a}_i^K. \quad (2)$$

3.1.2 Reaction-aware graph learning module

This module distills reaction-level relations from the reaction-aware graph to model molecule representations, as

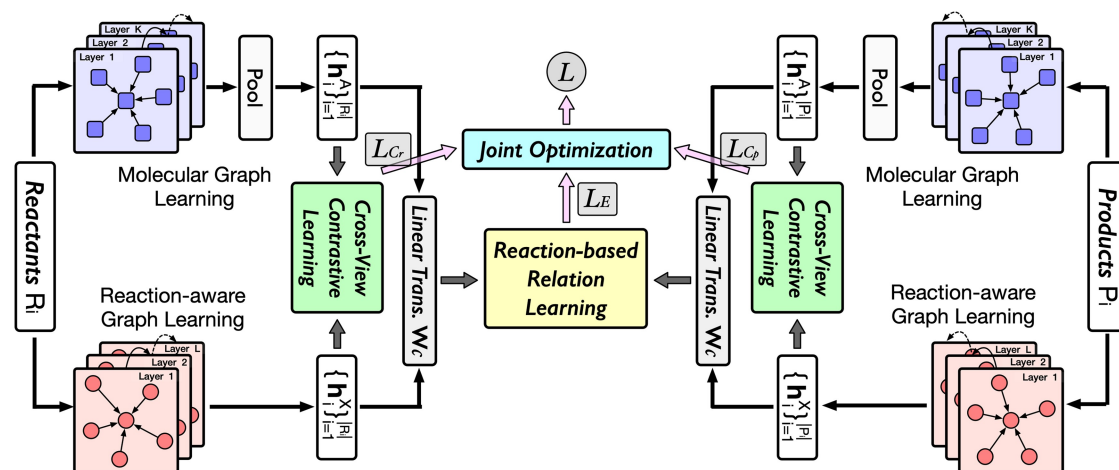


Figure 2. The framework of our RXGL method.

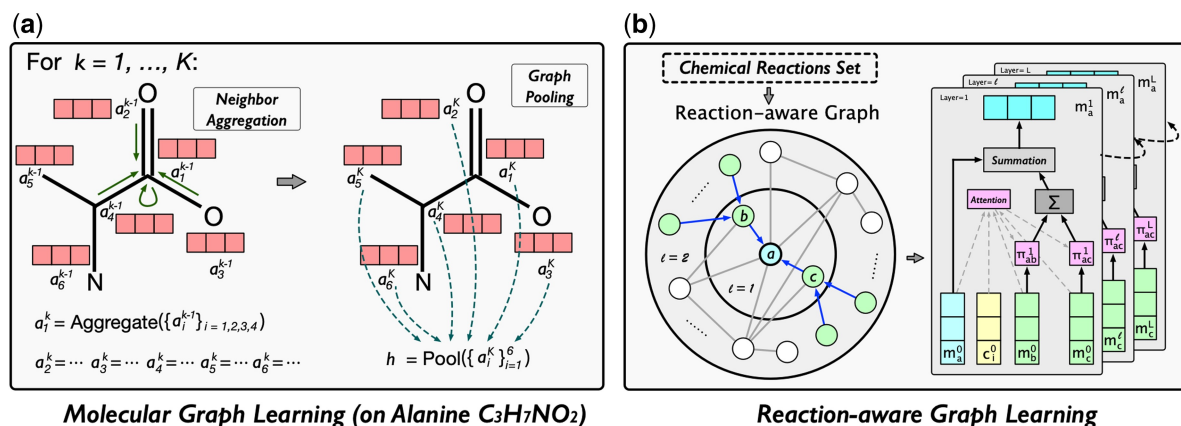


Figure 3. Molecular and reaction-aware graph learning.

shown in Fig. 3b. We first construct a reaction-aware graph $\mathcal{G}_x = (\mathcal{V}_x, \mathcal{E}_x)$. To describe initial feature \mathbf{m}_i of molecule $m_i \in \mathcal{V}_x$ in \mathcal{G}_x , we use the molecule functional group information. We first use a one-hot embedding to denote each functional group, and then the molecule is represented by the summation of its functional group embeddings. We consider 39 functional groups, which can be found in the [Supplementary Materials](#).

We design a GNN to model molecule representations in \mathcal{G}_x . Guided by chemical reactions, it employs an attention mechanism for aggregating neighbor information. Specifically, for molecule m_i in reaction $x: R \rightarrow P$ (i.e. $m_i \in R \cup P$), its neighbor aggregation in the l th layer is defined as:

$$\mathbf{m}_i^l = \mathbf{m}_i^{l-1} + \sum_{j \in \mathcal{N}_i} \pi_{ij}^l \mathbf{m}_j^{l-1}, \quad (3)$$

where \mathbf{m}_i^l is m_i 's representation ($\mathbf{m}_i^0 = \mathbf{m}_i$), π_{ij} denotes the attention score between m_i and its neighbor $m_j \in \mathcal{N}_i$ in \mathcal{G}_x . To enhance computational efficiency, we sample a fixed-size neighbor subset from \mathcal{N}_i for each molecule. The size of the subset is a predefined constant Q .

The attention score π_{ij} is defined as follows:

$$\pi_{ij}^l = \frac{\exp(\mathbf{W}_a^l (\mathbf{m}_i^{l-1} \odot \mathbf{m}_j^{l-1} \odot \mathbf{c}^{l-1}))}{\sum_{k \in \mathcal{N}_i} \exp(\mathbf{W}_a^l (\mathbf{m}_i^{l-1} \odot \mathbf{m}_k^{l-1} \odot \mathbf{c}^{l-1}))}, \quad (4)$$

where \odot denotes the element-wise product between vectors, \mathbf{W}_a is the weight matrix, and \mathbf{c} is the contextual information for molecule m_i in chemical reaction $x: R \rightarrow P$. The reaction context \mathbf{c} is defined as $\mathbf{c}^l = \sum_{m_u \in R \cup P} \mathbf{m}_u^l$. Our attention mechanism considers the reaction context, allowing for the control in neighbor message passing.

After stacking L layers of GNN, we utilize representation \mathbf{m}_i^L as molecule m_i 's reaction-level representation \mathbf{h}_i^X .

3.2 Molecule representation optimization

For each molecule m_i , we derive its two representations (i.e. \mathbf{h}_i^A and \mathbf{h}_i^X) from dual graph learning. To model m_i 's final representation \mathbf{h}_i , we use a linear transformation with weight \mathbf{W}_c and concatenation operation \parallel , as $\mathbf{h}_i = \mathbf{W}_c(\mathbf{h}_i^A \parallel \mathbf{h}_i^X)$.

To further optimize these molecular representations, we introduce a reaction-based relation learning task and a cross-view contrastive learning task. The former models the relation between reactants and products in each reaction. The latter is designed to enhance the cooperative association between molecular and reaction-aware graph learning.

3.2.1 Reaction-based relation learning task

This task models reactant and product pairs using the relation vector, as shown in Fig. 1b2. Given a chemical reaction $x_i: R_i \rightarrow P_i$, we assume that equation $\mathbf{x}_{R_i} + \mathbf{e}_{R_i \rightarrow P_i} = \mathbf{x}_{P_i}$ holds, where \mathbf{x}_{R_i} , \mathbf{x}_{P_i} , and $\mathbf{e}_{R_i \rightarrow P_i}$ are the representations of reactant R_i , product P_i , and relation from R_i to P_i , respectively. First, we use the summation operation to define \mathbf{x}_{R_i} and \mathbf{x}_{P_i} , as:

$$\mathbf{x}_{R_i} = \sum_{m_j \in R_i} \mathbf{h}_j, \quad \mathbf{x}_{P_i} = \sum_{m_k \in P_i} \mathbf{h}_k. \quad (5)$$

To model the relation vector $\mathbf{e}_{R_i \rightarrow P_i}$ of reaction x_i , we introduce a memory network (as shown in Fig. 4) which takes a reactant-product pair $(\mathbf{x}_{R_i}, \mathbf{x}_{P_i})$ as input, and outputs the relation vector $\mathbf{e}_{R_i \rightarrow P_i}$ of equal dimension as \mathbf{x}_{R_i} and \mathbf{x}_{P_i} .

Let \mathbb{R}^d be the dimension of \mathbf{x}_{R_i} and \mathbf{x}_{P_i} . The network first introduces a key matrix $\mathbf{K} \in \mathbb{R}^{N \times d}$ to calculate an attention vector $\mathbf{att} \in \mathbb{R}^N$. The element $att_n \in \mathbf{att}$ is defined as:

$$att_n = \frac{\exp((\mathbf{x}_{R_i} \odot \mathbf{x}_{P_i})^\top \mathbf{k}_n)}{\sum_{m=1}^N \exp((\mathbf{x}_{R_i} \odot \mathbf{x}_{P_i})^\top \mathbf{k}_m)}, \quad (6)$$

where $\mathbf{k} \in \mathbb{R}^d$ is the vector element of matrix \mathbf{K} .

Next, we introduce a memory matrix $\mathbf{M} \in \mathbb{R}^{N \times d}$ to generate the relation vector $\mathbf{e}_{R_i \rightarrow P_i}$, as follows:

$$\mathbf{e}_{R_i \rightarrow P_i} = \sum_{n=1}^N att_n \cdot \mathbf{m}_n, \quad (7)$$

where $\mathbf{m} \in \mathbb{R}^d$ is the vector element of matrix \mathbf{M} .

The relation vector $\mathbf{e}_{R_i \rightarrow P_i}$ is a weighted representation of \mathbf{M} . Intuitively, the memory matrix \mathbf{M} can be interpreted as a store of conceptual building blocks that can be used to describe the relations between reactants and products.

After the relation modeling, we define the following score function for reactant-product pair R_i and P_i of reaction x_i :

$$S(R_i, P_i) = \|\mathbf{x}_{R_i} + \mathbf{e}_{R_i \rightarrow P_i} - \mathbf{x}_{P_i}\|_2. \quad (8)$$

For optimization, we utilize a contrastive learning method similar to (Radford et al. 2021). In a minibatch of reaction data \mathcal{X}_B , we identify matched reactant-product pairs as positive pairs, aiming to minimize their embedding differences, while unmatched pairs are considered negative pairs, with an objective to maximize their embedding discrepancies. We use the margin-based loss function (Bordes et al. 2013), as follows:

$$\mathcal{L}_E = \frac{1}{|\mathcal{X}_B|} \sum_{x_i \in \mathcal{X}_B} S(R_i, P_i) + \frac{1}{|\mathcal{X}_B|^2 - |\mathcal{X}_B|} \sum_{x_i \in \mathcal{X}_B} \sum_{x_j \in \mathcal{X}_B} \max(\gamma - S(R_i, P_j), 0), \quad (9)$$

where $x_i \neq x_j$ and $\gamma > 0$ is a margin hyperparameter.

3.2.2 Cross-view contrastive learning task

This task enhances representations of molecules by aligning outputs from molecular and reaction-aware graph learning modules. We consider different views of the same reactant or product as positive pairs, while views of different reactants or products form negative pairs. This method enables our model to learn distinct representations by contrasting these positive and negative instances. Since contrastive tasks of reactants and products are similar, we illustrate the reactant side task.

Specifically, for a reaction x_i , we first generate its two reactant representations $\mathbf{x}_{R_i}^A$ and $\mathbf{x}_{R_i}^X$ from the molecular and reaction-aware graph learning modules, respectively, as:

$$\mathbf{x}_{R_i}^A = \sum_{m_j \in R_i} \mathbf{h}_j^A, \quad \mathbf{x}_{R_i}^X = \sum_{m_k \in R_i} \mathbf{h}_k^X. \quad (10)$$

Then, we enforce the separation of different reactant representations and align those that are identical. To achieve this, we employ InfoNCE (Oord et al. 2018) to define the cross-view contrastive loss for reactants as follows:

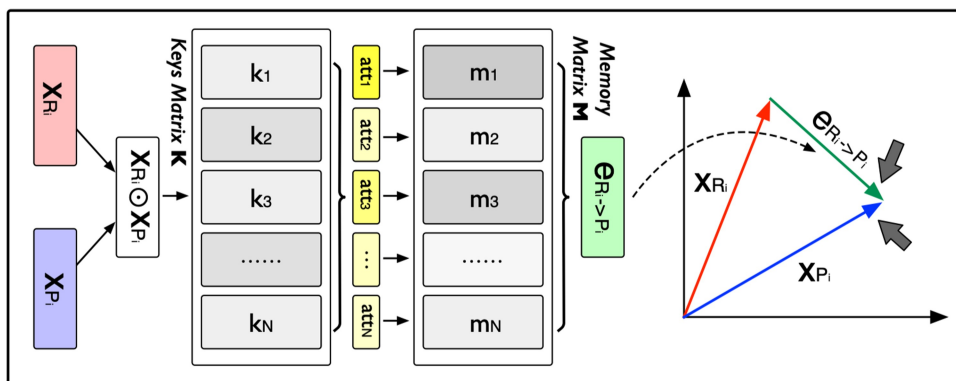


Figure 4. Reaction-based relation learning.

$$\mathcal{L}_C = \sum_{x_i \in \mathcal{X}_B} -\log \frac{\exp(c(\mathbf{x}_{R_i}^A, \mathbf{x}_{R_i}^X)/\tau)}{\sum_{x_j \in \mathcal{X}_B} \exp(c(\mathbf{x}_{R_i}^A, \mathbf{x}_{R_j}^X)/\tau)}, \quad (11)$$

where \mathcal{X}_B denotes a minibatch of reaction data, $c(\cdot, \cdot)$ is the cosine similarity, and τ is a hyper-parameter for the softmax function. Here we use a minibatch-based contrastive strategy for its efficiency in terms of both time and memory.

Similarly, we calculate contrastive loss \mathcal{L}_{C_p} for the product side. Combining two losses, we formulate the objective function of our cross-view contrastive task as $\mathcal{L}_C = \mathcal{L}_C + \mathcal{L}_{C_p}$.

3.2.3 Optimization function

Our RXGL is trained using a weighted sum of losses from the reaction relation learning and cross-view contrastive tasks.

$$\mathcal{L} = \mathcal{L}_E + \alpha \mathcal{L}_C + \lambda \|\Theta\|_F^2, \quad (12)$$

where Θ is model parameters, and α and λ adjust the cross-view contrastive task and L2-regularization, respectively.

We use the minibatch-based negative sampling strategy for two tasks, offering advantages over the traditional approach (Mikolov *et al.* 2013). First, it requires no extra memory to store negative samples. Second, negative samples are refreshed each epoch due to the shuffling of training instances, saving the time for manual re-sampling.

4 Experiments

In this section, we conduct extensive experiments to show the effectiveness of our approach RXGL in several downstream tasks, including product prediction, reaction classification, and molecular property prediction.

4.1 Product prediction

Reaction product prediction is a fundamental problem in the biology and chemistry field, which aims to predict the products of a reaction based on the reactants.

Dataset. We conduct experiments on two biochemistry benchmark datasets: USPTO-15K (Jin *et al.* 2017), containing 15 000 reactions, and USPTO-50K (Liu *et al.* 2024), comprising 50 000 reactions. Each dataset was divided into training, validation, and test sets using an 8:1:1 ratio.

Evaluation protocol. Following (Wang *et al.* 2022a), we formulate the product prediction as a ranking task. Let \mathcal{X}_{test} and P_{test} be the sets of reactions and products in the test set. For each reaction $x_i: R_i \rightarrow P_i$ in \mathcal{X}_{test} , we rank all products

Table 1. Results of the product prediction task on the two datasets.^a

Methods	USPTO-15K		USPTO-50K	
	MRR	Hit@1	MRR	Hit@1
Mol2vec	0.519	0.468	0.835	0.801
MolBERT	0.790	0.734	0.913	0.874
MolR-GCN	0.883 _(0.005)	0.847 _(0.007)	0.958 _(0.003)	0.944 _(0.006)
MolR-GAT	0.881 _(0.003)	0.846 _(0.002)	0.952 _(0.002)	0.931 _(0.004)
MolR-SAGE	0.932 _(0.006)	0.905 _(0.003)	0.972 _(0.005)	0.960 _(0.005)
MolR-TAG	0.925 _(0.005)	0.898 _(0.004)	0.974 _(0.010)	0.965 _(0.009)
RXGL-GCN	0.927 _(0.003)	0.899 _(0.005)	0.965 _(0.007)	0.954 _(0.002)
RXGL-GAT	0.925 _(0.007)	0.894 _(0.006)	0.967 _(0.009)	0.958 _(0.011)
RXGL-SAGE	0.956 _(0.004)	0.936 _(0.007)	0.982 _(0.005)	0.973 _(0.003)
RXGL-TAG	0.941 _(0.006)	0.919 _(0.009)	0.979 _(0.004)	0.974 _(0.009)

^a The numbers in brackets are the standard deviations. Bold values denote the best values of all methods.

$P_j \in P_{test}$ based on the score function $S(R_i, P_j)$ [i.e. Equation (8)]. Then, the ranking of the ground-truth candidate can be used for the evaluation. Our evaluation metrics are MRR (mean reciprocal rank) and Hit@1 (hit ratio at a cut-off value of 1). We conduct each experiment five times, reporting both the mean and standard deviation on the test set, with the results selected based on the best MRR in the validation set.

Baselines. Our RXGL is compared with Mol2vec (Jaeger *et al.* 2018), MolBERT (Fabian *et al.* 2020), and MolR (Wang *et al.* 2022a). For Mol2vec and MolBERT, we employ their pre-trained models to generate embeddings for reactants and products. Then, aligning with MolR, the scoring function is defined through an inner product. For MolR and our RXGL, we use four GNNs (i.e. GCN, GAT, SAGE, and TAG) as encoders for the molecular graph. We present the details of these encoders in the [Supplementary Materials](#).

Hyperparameter settings. We train our model RXGL for 50 epochs with a batch size of 1024, using Adam optimizer with a learning rate of 10^{-4} . The number of GNN layers in molecular and reaction-aware graph learning is set to 2 and 1, respectively, and the output dimension of all layers is 256. In addition, neighbor sampling size Q , memory slice size N , margin γ , temperature τ , and balancing factor α are set to 10, 10, 4, 0.1, and 10^{-3} , respectively. The result of the key hyperparameter sensitivity is reported in Section 4.4.

Performance comparison. The results are shown in Table 1. We find that our RXGL performs best. To be specific, RXGL-SAGE achieves 14.1% average MRR gain and 16.3% average Hit@1 gain over baselines on the USPTO-15K dataset, showing the effectiveness of RXGL in the product prediction.

Table 2. Results of the reaction classification task.^a

Methods	Schneider		USPTO-MTL	
	Accuracy	Recall	Accuracy	Recall
Mol2vec	0.856 _(0.012)	0.850 _(0.009)	0.751 _(0.016)	0.629 _(0.014)
MolBERT	0.849 _(0.014)	0.847 _(0.012)	0.738 _(0.009)	0.583 _(0.008)
MolR-GCN	0.879 _(0.013)	0.882 _(0.013)	0.853 _(0.020)	0.813 _(0.015)
MolR-GAT	0.870 _(0.024)	0.873 _(0.021)	0.862 _(0.017)	0.834 _(0.014)
MolR-SAGE	0.882 _(0.030)	0.881 _(0.027)	0.874 _(0.023)	0.839 _(0.026)
MolR-TAG	0.891 _(0.025)	0.895 _(0.026)	0.888 _(0.024)	0.852 _(0.024)
RXGL-GCN	0.887 _(0.016)	0.889 _(0.017)	0.875 _(0.018)	0.849 _(0.015)
RXGL-GAT	0.872 _(0.026)	0.874 _(0.020)	0.873 _(0.019)	0.836 _(0.016)
RXGL-SAGE	0.907 _(0.029)	0.906 _(0.025)	0.889 _(0.014)	0.858 _(0.018)
RXGL-TAG	0.899 _(0.033)	0.901 _(0.031)	0.895 _(0.019)	0.867 _(0.018)

^a The numbers in brackets are the standard deviations. Bold values denote the best values of all methods.

4.2 Reaction classification

Reaction classification aims to predict the class of reactions.

Dataset. Our RXGL is evaluated using two datasets, i.e. Schneider and USPTO-MTL. Specifically, Schneider (Schneider *et al.* 2015) comprises 38 800 reactions across 46 classes, and USPTO-MTL (Lu and Zhang 2022) includes 143 535 reactions across 1000 classes. Each dataset is randomly split into training, validation, and test sets in an 8:1:1 ratio.

Baselines and experiment settings. Similar to product prediction, we compare our RXGL with Mol2vec, MolBERT, and MolR. We employ our model pre-trained on the USPTO-50k dataset to process all datasets, which includes generating embeddings for reactants and products in each reaction and concatenating them to create a unified reaction feature. For prediction, we use an MLP as the decoder. Accuracy and Recall are used as evaluation metrics. For each experiment, we perform five times and report the mean and standard deviation of the results on the test set.

Performance comparison. Table 2 shows the results of the reaction classification task. We find that our RXGL outperforms baselines. For example, RXGL-SAGE achieves an average Accuracy increase of 4.0% on Schneider dataset. These results highlight RXGL’s efficacy in the reaction classification task, and suggest that the molecule representations it learns are effectively transferable to the downstream task.

4.3 Molecular property prediction

This task is to predict labels of given molecules, which is a classical task to evaluate learned molecule representations.

Dataset. We evaluate our RXGL on four datasets (Wu *et al.* 2018): BBBP, BACE, Tox21, and ClinTox. Each dataset contains molecule SMILES as well as labels indicating the property. Readers can refer to (Wu *et al.* 2018) for a detailed introduction to these four public datasets.

Baselines. We compare our RXGL (i.e. variant RXGL-GCN) with the following four class methods: (i) SMILES-based methods: ChemBERTa (Chithrananda *et al.* 2020) and MolBERT (Fabian *et al.* 2020); (ii) Fingerprint-based methods: Mol2vec (Jaeger *et al.* 2018), ECFP4 (Rogers and Hahn 2010), GraphConv (Duvenaud *et al.* 2015), Weave (Kearnes *et al.*, 2016), D-MPNN (Yang *et al.* 2019), CDDD (Winter *et al.* 2019); (iii) GNN-based methods: GraphCL (You *et al.* 2020), GraphLoG (Xu *et al.* 2021), EdgePred (Hamilton *et al.* 2017), AttrMask (Hu *et al.* 2019), GPT-GNN (Hu *et al.* 2020), InfoGraph (Sun *et al.* 2020),

Table 3. Molecular property prediction results (split type: *random split*).^a

Methods	BBBP	BACE	Tox21	ClinTox
ChemBERTa [★]	0.643	–	0.728	0.733
MolBERT [★]	0.762 _(0.000)	0.866 _(0.000)	–	–
Mol2vec [★]	0.872 _(0.021)	0.862 _(0.027)	0.803 _(0.041)	0.841 _(0.062)
ECFP4 [★]	0.729	0.867	0.822	0.799
GraphConv [★]	0.690	0.783	0.829	0.807
Weave [★]	0.671	0.806	0.820	0.832
D-MPNN [★]	0.708	–	0.688	0.906
CDDD [★]	0.761 _(0.000)	0.833 _(0.000)	–	–
GraphCL [★]	0.695 _(0.005)	0.782 _(0.012)	0.754 _(0.009)	0.701 _(0.019)
GraphLoG [★]	0.725 _(0.008)	0.835 _(0.012)	0.757 _(0.005)	0.767 _(0.033)
MolR [★]	0.895 _(0.031)	0.875 _(0.023)	0.820 _(0.028)	0.913 _(0.043)
ReaKE [♠]	–	0.898	0.824	0.874
RXGL	0.901 _(0.024)	0.906 _(0.017)	0.852 _(0.026)	0.921 _(0.030)

^a The numbers in brackets are the standard deviations. The results of symbols [★] and [♠] are taken from MolR and ReaKE. Bold values denote the best values of all methods.

ContextPred (Hu *et al.* 2019), G-Motif (Rong *et al.* 2020), JOAO (You *et al.* 2021), and GraphMVP (Liu *et al.* 2022); (iv) Reaction-enhanced methods: MolR (Wang *et al.* 2022a) and ReaKE (Wang *et al.* 2022b).

Experiment settings. We split all datasets into training, validation, and test sets in an 8:1:1 ratio, using two split types: Random and scaffold. Our model pre-trained on USPTO-50K datasets is used to generate molecule embeddings, which are formed by concatenating outputs from the molecular graph learning module with the sum of molecule functional group embeddings. Then, we input the molecule embeddings along with their labels into a logistic regression model. We use the AUC (area under the curve) as our evaluation metric. Each experiment is conducted five times, and we report the mean and standard deviation of the results on the test set.

Performance comparison. Different baselines leverage different strategies (i.e. random or scaffold) to split datasets. We compare our model with baselines by adopting their selected dataset splitting. The AUC results for molecular property prediction are presented in Tables 3 and 4. The baseline results are taken from the literature. Our RXGL shows strong performance across four datasets, e.g. RXGL exhibits average AUC gains of 17.8%, 7.2%, 7.9%, and 11.3% on four datasets under random splitting, showing that learned molecule representations effectively transfer to molecule-related tasks.

4.4 Hyper-parameter study

We study three hyper-parameters on RXGL-GCN: Embedding size (d), balancing factor (α), and neighbor sampling size (Q), as shown in Fig. 5. We find: (i) A larger d typically enhances performance by encoding more information, yet an excessively large d increases the parameter size. We suggest setting d at 256 to balance performance and storage requirements. (ii) Nice performance is achieved when α is set to 0.001. We hypothesize that a small α might inadequately reinforce the cooperative associations between molecular and reaction-aware graph learning. (iii) The rank accuracy improves with increasing Q but eventually decreases. This decline may be due to the introduction of irrelevant information from too many neighbors.

4.5 Embedding analysis

This subsection studies the learned molecule embeddings and reaction-based relation embeddings [i.e. Equation (7)].

4.5.1 Molecule embedding analysis

We first use RXGL-GCN to generate embeddings of molecules in the BBBP dataset and visualize them utilizing the t-SNE method (Van der Maaten and Hinton 2008). Specifically, we select the molecules' permeability property, molecule size (i.e. the number of atoms), and 39 functional group properties. These 39 functional groups are provided by RDKit (Landrum 2013). Due to the space limitation, we present visualizations in the main text for the permeability, the molecule size, and the hydroxyl functional group properties. More detailed information about these functional groups and remaining experimental results can be found in [Supplementary Materials](#). From visualization results, we find

Table 4. Property prediction results (split type: *scaffold split*).^a

Methods	BBBP	BACE	Tox21	ClinTox
EdgePred [★]	0.645 _(0.031)	0.646 _(0.047)	0.745 _(0.004)	0.558 _(0.062)
AttrMask [★]	0.702 _(0.005)	0.772 _(0.014)	0.742 _(0.008)	0.686 _(0.096)
GPT-GNN [★]	0.645 _(0.011)	0.776 _(0.005)	0.753 _(0.005)	0.578 _(0.031)
InfoGraph [★]	0.692 _(0.008)	0.739 _(0.025)	0.730 _(0.007)	0.751 _(0.050)
ContextPred [★]	0.712 _(0.009)	0.786 _(0.014)	0.733 _(0.005)	0.737 _(0.040)
G-Contextual [★]	0.703 _(0.016)	0.792 _(0.003)	0.752 _(0.003)	0.599 _(0.082)
G-Motif [★]	0.664 _(0.034)	0.734 _(0.040)	0.732 _(0.008)	0.778 _(0.020)
JOAO [★]	0.660 _(0.006)	0.729 _(0.020)	0.744 _(0.007)	0.663 _(0.039)
GraphMVP [★]	0.724 _(0.016)	0.812 _(0.009)	0.744 _(0.002)	0.775 _(0.042)
MolR [♠]	–	0.774	0.670	0.830
ReaKE [♠]	–	0.781	0.713	0.862
RXGL	0.729 _(0.015)	0.825 _(0.006)	0.736 _(0.003)	0.912 _(0.011)

^a The numbers in brackets are the standard deviations. The results of symbols [★] and [♠] are taken from GraphMVP and ReaKE. Bold values denote the best values of all methods.

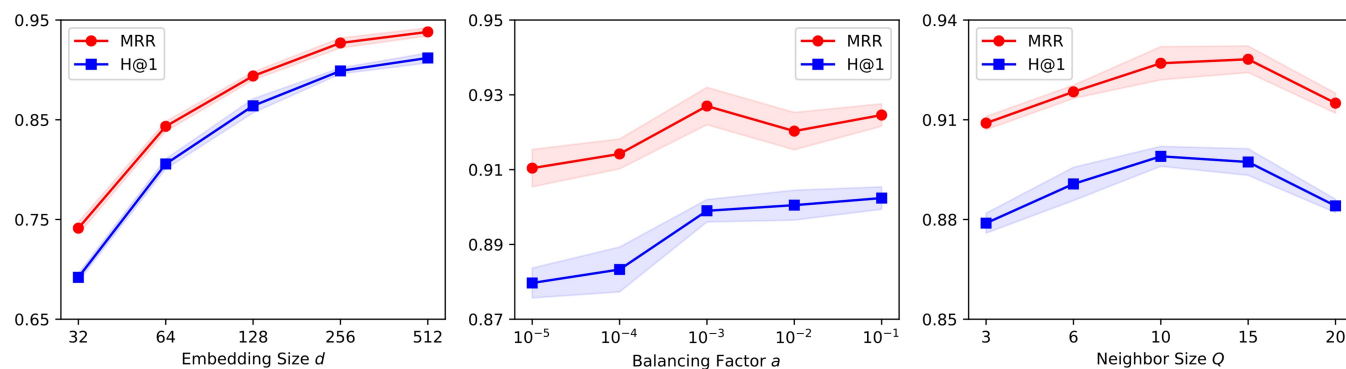


Figure 5. Hyper-parameter study on the USPTO-15K dataset.

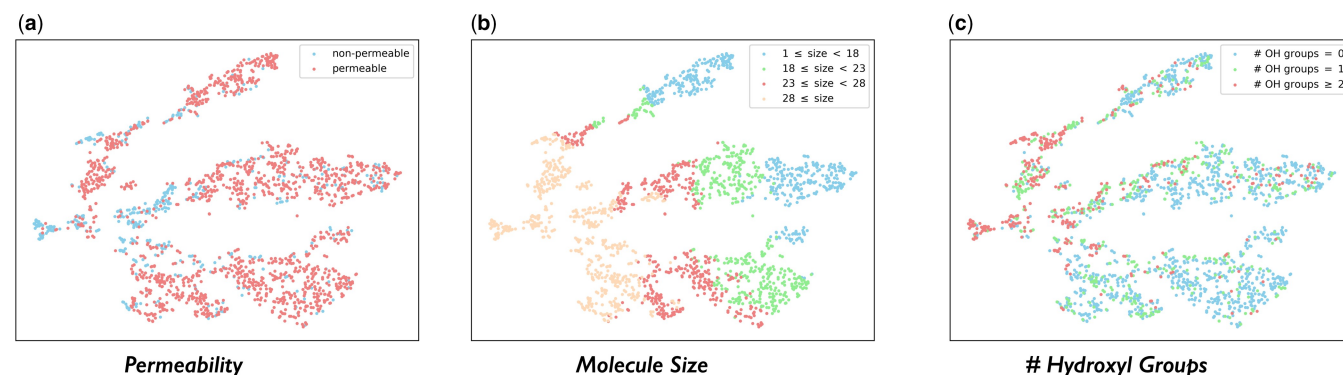


Figure 6. Visualized molecule embeddings on the BBBP dataset.

that: [Figure 6a](#) illustrates molecules colored according to their permeability properties. Several distinct clusters of non-permeable molecules are observed. [Figure 6b](#) shows molecules colored by size. The embedding space distinctly segregates small molecules (located in the right region) from large molecules (located in the left region). [Figure 6c](#) shows molecules colored based on the number of the hydroxyl functional groups (i.e. ‘-OH’). Molecules with a higher number of hydroxyl groups are mainly on the left side of the embedding space, while those without this group are primarily found on the right side. In summary, these results indicate that the molecular embeddings generated by our RXGL framework exhibit a certain degree of correlation with the aforementioned three molecular properties.

4.5.2 Reaction-based relation embedding analysis

In our RXGL, we learn relation embeddings between reactants and products in chemical reactions. This subsection examines the meaningfulness of these relation embeddings.

A key feature of chemical reactions is the change of reactants into products. For example, this process involves breaking old bonds and forming new ones, resulting in a change in the number of bonds before and after the reaction. Intuitively, reactions with similar changes yield similar relation embeddings. To analyze this, we randomly select five reactions (labeled a , b , c , d , and e) from the test set of USPTO-15K dataset and compute the cosine similarity of their relation embeddings, as shown in the left subfigure of [Fig. 7](#). Take reaction a as an example, we observe a high similarity between a and c and a low similarity between a and d .

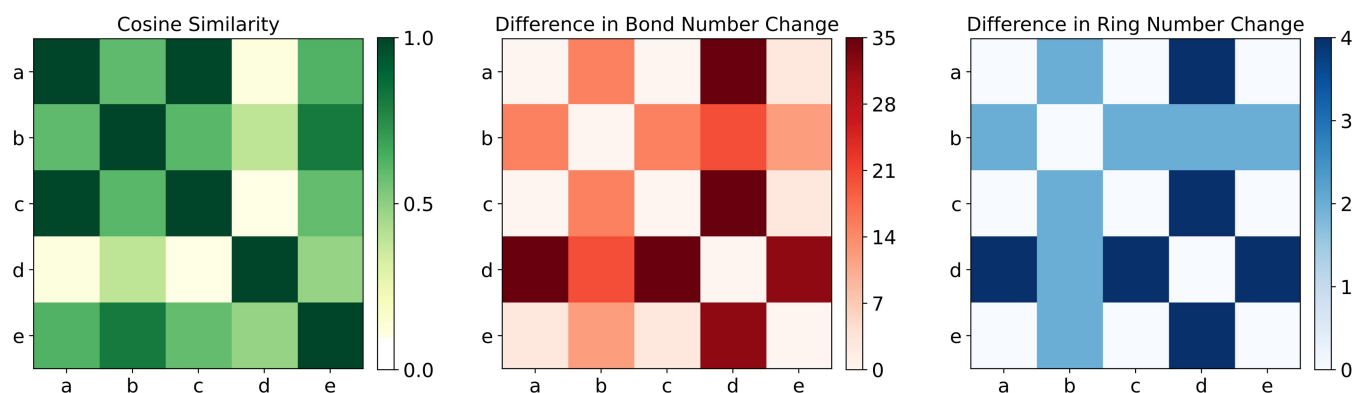


Figure 7. Examples of relation embedding comparison.

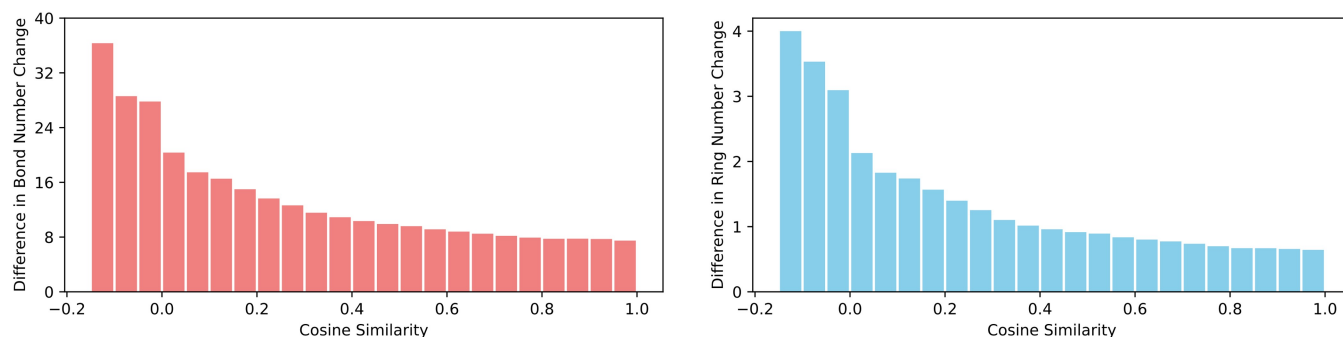


Figure 8. Macro-analysis on the USPTO-15K dataset.

We characterize chemical reaction changes by the difference in bond number (N_B) and ring number (N_R) in reactants and products. The middle and right subfigures of Fig. 7 show differences between these five chemical reactions in terms of bond number changes and ring number changes. For example, a small square in the middle subfigure represents $|N_B(i) - N_B(j)|$, which means the differences between reactions i and j in terms of bond number changes. From the results, we find that the high cosine similarity between reactions a and c is likely due to their similar bond and ring number changes. Conversely, the low similarity between reactions a and d can be attributed to their significant differences in bond and ring number changes.

We further conduct the macro-analysis. We calculate the cosine similarity between every reaction pair in the test set and examine the differences in bond number and ring number changes between each pair. We show the results on the USPTO-15K dataset in Fig. 8, where the y-axis represents the average change within a certain similarity interval. We find that reaction pairs with high cosine similarities tend to exhibit small differences in bond and ring number changes, and vice versa. The above observation underscores the significance of our learned relation embeddings, suggesting their capability to model the changes occurring in chemical reactions.

5 Conclusion and future work

This article proposes RXGL that uses reactions as domain knowledge in MRL. RXGL integrates molecular and reaction-aware graph learning modules to model molecule representations. Also, we enhance molecule representations using a reaction-based relation learning task and a cross-view contrastive task. Experiment results show that RXGL

achieves strong performance across a range of downstream tasks. We believe that the insights of this study will provide valuable guidance for future research exploring the utilization of reactions in MRL. For future work, we plan to consider the incorporation of stereochemical information into RXGL, which will enhance the accuracy and applicability of molecule representations.

Acknowledgements

The authors thank the anonymous reviewers for their valuable suggestions.

Supplementary data

Supplementary data are available at *Bioinformatics* online.

Conflict of interest

None declared.

Funding

This work was supported by Research Council of Finland [339421 and 345802] as well as Jane and Aatos Erkko Foundation (BIODESIGN project). Elena Casiraghi acknowledges travel support from the European Union's Horizon 2020 research and innovation programme under grant agreement No 951847. This article is supported by FAIR (Future Artificial Intelligence Research) project, funded by the NextGenerationEU program within the PNRR-PE-AI scheme (M4C2, Investment 1.3, Line on Artificial Intelligence).

Data availability

No new data were generated or analyzed in support of this research.

References

- Bordes A, Usunier N, Garcia-Duran A *et al.* Translating embeddings for modeling multi-relational data. In: *Advances in Neural Information Processing Systems*, USA, Vol. 26, 2013.
- Chithrananda S, Grand G, Ramsundar B. Chemberta: large-scale self-supervised pretraining for molecular property prediction. arXiv:2010.09885, 2020, preprint: not peer reviewed.
- Du J, Zhang S, Wu G *et al.* Topology adaptive graph convolutional networks. arXiv, arXiv:1710.10370, 2017, preprint: not peer reviewed.
- Duvenaud D, Maclaurin D, Aguilera-Iparraguirre J *et al.* Convolutional networks on graphs for learning molecular fingerprints. In: *Advances in Neural Information Processing Systems*, Canada, Vol. 28, 2015.
- Fabian B, Edlich T, Gaspar H *et al.* Molecular representation learning with language models and domain-relevant auxiliary tasks. arXiv, arXiv:2011.13230, 2020, preprint: not peer reviewed.
- Hamilton W, Ying Z, Leskovec J. Inductive representation learning on large graphs. In: *Advances in Neural Information Processing Systems*, USA, Vol. 30, 2017.
- Hu W, Liu B, Gomes J *et al.* Strategies for pre-training graph neural networks. In: *International Conference on Learning Representations*, Ethiopia, 2019.
- Hu Z, Dong Y, Wang K *et al.* GPT-GNN: generative pre-training of graph neural networks. In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, USA, 2020, 1857–67.
- Jaeger S, Fulle S, Turk S. Mol2vec: unsupervised machine learning approach with chemical intuition. *J Chem Inf Model* 2018;58:27–35.
- Jin W, Coley C, Barzilay R *et al.* Predicting organic reaction outcomes with Weisfeiler–Lehman network. In: *Advances in Neural Information Processing Systems*, USA, Vol. 20, 2017.
- Kearnes S, McCloskey K, Berndl M *et al.* Molecular graph convolutions: moving beyond fingerprints. *J Comput Aided Mol Des* 2016; 30:595–608.
- Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks. In: *International Conference on Learning Representations*, France, 2017.
- Landrum G. Rdkit documentation. *Release* 2013;1:4.
- Li A, Yang B, Huo H *et al.* Hypercomplex graph collaborative filtering. In: *The Web Conference*, France, 2022, 1914–22.
- Li A, Yang B, Huo H *et al.* Structure- and logic-aware heterogeneous graph learning for recommendation. In: *International Conference on Data Engineering*, Netherlands, IEEE, 2024, 544–556.
- Liu J, Yan C, Yu Y *et al.* Mars: a motif-based autoregressive model for retrosynthesis prediction. *Bioinformatics* 2024;40:btac115.
- Liu S, Wang H, Lasenby J *et al.* Pre-training molecular graph representation with 3D geometry. arXiv, arXiv:2110.07728, 2022, preprint: not peer reviewed.
- Lu J, Zhang Y. Unified model for multitask reaction predictions with explanation. *J Chem Inf Model* 2022;62:1376–87.
- Mikolov T, Sutskever I, Chen K *et al.* Distributed representations of words and phrases and their compositionality. In: *Advances in Neural Information Processing Systems*, USA, Vol. 26, 2013.
- Miller A, Fisch A, Dodge J *et al.* Key-value memory networks for directly reading documents. arXiv, arXiv:1606.03126, 2016, preprint: not peer reviewed.
- Oord A.v.d, Li Y, Vinyals O. Representation learning with contrastive predictive coding. arXiv arXiv:1807.03748. 2018, preprint: not peer reviewed.
- Radford A, Kim JW, Hallacy C *et al.* Learning transferable visual models from natural language supervision. In: *International Conference on Machine Learning*, Virtual Event, 2021, 8748–63.
- Rogers D, Hahn M. Extended-connectivity fingerprints. *J Chem Inf Model* 2010;50:742–54.
- Rong Y, Bian Y, Xu T *et al.* Self-supervised graph transformer on large-scale molecular data. *NIPS* 2020;33:12559–71.
- Schneider N, Lowe DM, Sayle RA *et al.* Development of a novel fingerprint for chemical reactions and its application to large-scale reaction classification and similarity. *J Chem Inf Model* 2015;55:39–53.
- Sun F-Y, Hoffmann J, Verma V *et al.* Infograph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In: *International Conference on Learning Representations*, Ethiopia, 2020.
- Van der Maaten L, Hinton G. Visualizing data using t-sne. *J Mach Learn Res* 2008;9:11.
- Veličković P, Cucurull G, Casanova A *et al.* Graph attention networks. In: *International Conference on Learning Representations*, Canada, 2018.
- Wagner A, Fell DA. The small world inside large metabolic networks. *Proc Biol Sci* 2001;268:1803–10.
- Wang H, Li W, Jin X *et al.* Chemical-reaction-aware molecule representation learning. In: *International Conference on Learning Representations*, Virtual Event, 2022a.
- Wang Y, Zheng S, Rao J *et al.* Reake: contrastive molecular representation learning with chemical synthetic knowledge graph. 2022b, preprint: not peer reviewed.
- Wen M, Spotte-Smith EWC, Blau SM *et al.* Chemical reaction networks and opportunities for machine learning. *Nat Comput Sci* 2023;3:12–24.
- Winter R, Montanari F, Noé F *et al.* Learning continuous and data-driven molecular descriptors by translating equivalent chemical representations. *Chem Sci* 2019;10:1692–701.
- Wu Z, Ramsundar B, Feinberg EN *et al.* Moleculenet: a benchmark for molecular machine learning. *Chem Sci* 2018;9:513–30.
- Xu M, Wang H, Ni B *et al.* Self-supervised graph-level representation learning with local and global structure. In: *International Conference on Machine Learning*, Virtual Event, PMLR, 2021, 11548–58.
- Yang K, Swanson K, Jin W *et al.* Are learned molecular representations ready for prime time? arXiv, arXiv:1904.01561, 2019, preprint: not peer reviewed.
- Yi H-C, You Z-H, Huang D-S *et al.* Graph representation learning in bioinformatics: trends, methods and applications. *BIB* 2022; 23:bbab340.
- You Y, Chen T, Sui Y *et al.* Graph contrastive learning with augmentations. *NIPS* 2020;33:5812–23.
- You Y, Chen T, Shen Y *et al.* Graph contrastive learning automated. In: *International Conference on Machine Learning*, Virtual Event, 2021.