MULTIMODAL FUSION WITH RELATIONAL LEARNING FOR MOLECULAR PROPERTY PREDICTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph-based molecular representation learning is essential for accurately predicting molecular properties in drug discovery and materials science; however, it faces significant challenges due to the intricate relationships among molecules and the limited chemical knowledge utilized during training. While contrastive learning is often employed to handle molecular relationships, its reliance on binary metrics is insufficient for capturing the complexity of these interactions. Multimodal fusion has gained attention for property reasoning, but previous work has explored only a limited range of modalities, and the optimal stages for fusing different modalities in molecular property tasks remain underexplored. In this paper, we introduce MMFRL (Multimodal Fusion with Relational Learning for Molecular Property Prediction), a novel framework designed to overcome these limitations. Our method enhances embedding initialization through multi-modal pre-training using relational learning. We also conduct a systematic investigation into the impact of modality fusion at different stages—early, intermediate, and late—highlighting their advantages and shortcomings. Extensive experiments on MoleculeNet benchmarks demonstrate that MMFRL significantly outperforms existing methods. Furthermore, MMFRL enables task-specific optimizations. Additionally, the explainability of MMFRL provides valuable chemical insights, emphasizing its potential to enhance real-world drug discovery applications.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

Graph representation learning for molecules has gained significant attention in drug discovery and materials science, as it effectively encapsulates molecular structures and enables the effective investigation of structure-activity relationships (Wieder et al., 2020; Zhang et al., 2022; Fang et al., 2022; Wang et al., 2023). In this paradigm, atoms are treated as nodes and chemical bonds as edges, effectively encapsulating the connectivities that define molecular behavior. However, it poses significant challenges due to intricate relationships among molecules and the limited chemical knowledge utilized during training.

Often, contrastive learning (CL) is employed to study relationships among molecules, but it relies 040 on binary metrics of positive and negative pairs, and tends to oversimplify complex molecular 041 interactions. For example, consider Thalidomide: while the (R)- and (S)-enantiomers share the same 042 topological graph and differ only at a single chiral center, their biological activities are drastically 043 different—the (R)-enantiomer is effective in treating morning sickness, whereas the (S)-enantiomer 044 causes severe birth defects. In other words, the (R)- and (S)-enantiomers are similar in terms of topological stucture but dissimilar in terms of biological activities. Thus, a more sophisticated approach is required to tackle these scenarios. A potential solution would be to use continuous 046 metrics within a multi-view space, enabling a more comprehensive understanding of these complex 047 molecular relationships. 048

When it comes to multimodal learning for molecules, we often encounter data availability and
incompleteness issues. This raises a critical question: how can multimodal information be effectively
leveraged for molecular property reasoning when such data is absent in downstream tasks? Recent
studies have demonstrated the effectiveness of pretraining molecular Graph Neural Networks (GNNs)
by integrating additional knowledge sources (Wang et al., 2021; 2022b; Liu et al., 2022a; Xu et al.,
2023a). Building on this foundation, a promising solution is to pretrain multiple replicas of molecular



Figure 1: Multimodal Fusion with Relational Learning for Molecular Property Prediction (MMFRL). This figure shows our proposed idea about how to transfer the knowledge from other modalities and use fusion to improve the performance further. Unlike the general contrastive learning framework shown in Appendix Figure A.2, MMFRL doesn't need to define positive or negative pairs and is capable of learning continuous ordering from target similarity.

095

087

088

GNNs, with each replica dedicated to learning from a specific modality. This approach allows downstream tasks to benefit from multimodal data that is not accessible during fine-tuning, ultimately improving representation learning.

Facing these challenges and opportunities, we propose MMFRL (Multimodal Fusion with Relational 096 Learning for Molecular Property Prediction), a novel framework features relational learning (RL) and multimodal fusion (MMF). RL utilizes a continuous relation metric to evaluate relationships among 098 instances in the feature space (Balcan & Blum, 2006; Wen et al., 2023). Our major contribution comprises three aspects: Conceptually: We introduce a modified relational learning metric for 100 molecular graph representation that offers a more comprehensive and continuous perspective on 101 inter-instance relations, effectively capturing both localized and global relationships among instances. 102 To the best of our knowledge, this is the first work to demonstrate such generalized relational learning 103 metric for molecular graph representation. *Methodologically:* Our proposed modified relational 104 metric captures complex relationships by converting pairwise self-similarity into relative similarity, 105 which evaluates how the similarity between two elements compares to the similarity of other pairs in the dataset. In addition, we integrate these metrics into a fused multimodal representation, which has 106 the potential to enhance performance, allowing downstream tasks to leverage modalities that are not 107 directly accessible during fine-tuning. *Empirically:* MMFRL excels in various downstream tasks for

108 Molecular Property Predictions. Last but not least, we demonstrate the explainability of the learned 109 representations through two post-hoc analysis. Notably, we explore minimum positive subgraphs and 110 maximum common subgraphs to gain insights for further drug molecule design. 111

2 PRELIMINARIES

112

113

114 Directed Message Passing Neural Network (DMPNN). The Message Passing Neural Network 115 (MPNN) (Gilmer et al., 2017) is a GNN model that processes an undirected graph G with node 116 (atom) features x_v and edge (chemical bond) features e_{vw} . It operates through two distinct phases: 117 a message passing phase, facilitating information transmission across the molecule to construct a 118 neural representation, and a readout phase, utilizing the final representation to make predictions 119 regarding properties of interest. The primary distinction between DMPNN and a generic MPNN lies 120 in the message passing phase. While MPNN uses messages associated with nodes, DMPNN crucially 121 differs by employing messages associated with directed edges (Yang et al., 2019). This design choice is motivated by the necessity to prevent totters (Mahé et al., 2004), eliminating messages passed 122 along paths of the form $v_1v_2 \dots v_n$, where $v_i = v_{i+2}$ for some *i*, thereby eliminating unnecessary 123 loops in the message passing trajectory. 124

125 Relational Learning. Original Relation Learning (Zheng et al., 2021) ensures that different 126 augmented views of the same instance from computer vision tasks share similar features, while 127 allowing for some variability. This approach captures the essential characteristics of the instance, 128 promoting consistency across the views without requiring them to be identical. By doing so, it enhances the model's ability to generalize and recognize underlying patterns in the data. Suppose z_i 129 is the original embdding for the i - th instance. Then z_i^1 is the embedding of first augmented view 130 for z_i , and z_i^2 is the embedding of second augmented view for z_i . In this case, the Loss of Relational 131 Learning (RL) is formulated as following: 132

133
134
$$s_{ik}^1 = \frac{1_{i \neq k} \cdot \exp(z_i^1 \cdot z_k^2 / \tau)}{\sum_{i \neq k} \sum_{j \neq k} \frac{1_{i \neq k} \cdot \exp(z_j^1 \cdot z_k^2 / \tau)}{|z_i|^2}$$

135
$$\sum_{j=1}^{j-1} i_{j \neq j} \cdot \exp(z_{i})$$

136
137
$$s_{ik}^2 = \frac{1_{i \neq k} \cdot \exp(z_i^2 \cdot z_i)}{N}$$

$$s_{ik}^{z} = \frac{1}{\sum_{j=1}^{N} 1_{i \neq j} \cdot \exp(z_{i}^{z} \cdot z_{j}^{2} / \tau_{m})}$$

 $s_{ik}^{1} = \frac{1_{i \neq k} \cdot \exp(z_{i} \cdot z_{k}/\tau)}{\sum_{j=1}^{N} 1_{i \neq j} \cdot \exp(z_{i}^{1} \cdot z_{j}^{2}/\tau)}$ $s_{ik}^{2} = \frac{1_{i \neq k} \cdot \exp(z_{i}^{2} \cdot z_{k}^{2}/\tau_{m})}{\sum_{j=1}^{N} 1_{i \neq j} \cdot \exp(z_{i}^{2} \cdot z_{j}^{2}/\tau_{m})}$ $L_{RL} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{\substack{k=1\\k \neq i}}^{N} s_{ik}^{2} \log(s_{ik}^{1}).$

141 142

138

140

143 Multi-Modality Fusion. Multi-Modality Fusion combines diverse heterogeneous data (e.g. text, 144 images, graph) to create a more comprehensive understanding of complex scenarios (Lahat et al., 145 2015; Khaleghi et al., 2013; Poria et al., 2015; Ramachandram & Taylor, 2017; Pawłowski et al., 2023). This approach leverages the strengths of each modality, potentially improving performance in 146 tasks like sentiment analysis or medical diagnosis. While challenging to implement due to the need to 147 align different data streams, successful fusion can provide insights beyond what's possible with single 148 modalities, advancing AI and data-driven decision-making. In particular, the way to fuse different 149 modality should also depends on the dominace of each unimodality (Pawłowski et al., 2023). 150

151 152

153

METHODS 3

154 We first explain our proposed modified metric in relational learning to facilitate smooth alignment 155 between graph and referred unimodality. Then, we introduce approaches for integrating multi modalities at different stages of the learning process. 156

157 158

159

3.1 MODIFIED RELATIONAL LEARNING IN PRETRAINING

We propose a modified relational metric by adapting the softmax function as a pairwise weighting 160 mechanism. Let |S| denote the size of the instance set. The variable $s_{i,j}$ represents the learned 161 similarity distribution where z_i is the embedding to be trained. On the other hand, $t_{i,i}^R$ defines the target similarity distribution that captures the relationship between the pair of instances in the given space or modality R, where z_i^R is a fixed embedding. The detailed formulation for the Loss of Modified Relatioal Learning (MRL) is provided below:

$$s_{i,j} = \frac{\exp(sim(z_i, z_j))}{\sum_{k=1}^{|\mathcal{S}|} \exp(sim(z_i, z_k))}$$
(1)

167 168 169

170 171

166

$$t_{i,j}^{R} = \frac{\exp(sim(z_{i}^{R}, z_{j}^{R}))}{\sum_{j=1}^{|S|} \exp(sim(z_{i}^{R}, z_{k}^{R}))}$$
(2)

172 173 174

175

$$L_{MRL} = -\frac{1}{|\mathcal{S}|} \sum_{i=1}^{|\mathcal{S}|} \sum_{j=1}^{|\mathcal{S}|} t_{i,j}^R \log(s_{i,j}).$$
(3)

Notably, unlike other similarity learning approaches (Wang et al., 2019; Zhang et al., 2021), our method does not rely on the categorization of negative and positive pairs for the pair weighting function. Additionally, the use of the softmax function ensures that the generalized target similarity $t_{i,j}$ adheres to the principles of convergence as following:

Theorem 3.1 (Convergence of Modified Relational Learning Metric). Let *S* be a set of instances with size of |S|, and let \mathcal{P} represent the learnable latent representations of instances in *S* such that $|\mathcal{P}| = |S|$. For any two instances $i, j \in S$, their respective latent representations are denoted by \mathcal{P}_i and \mathcal{P}_j . Let $t_{i,j}$ represent the target similarity between instances *i* and *j* in a given domain, and let $d_{i,j}$ be the similarity between \mathcal{P}_i and \mathcal{P}_j in the latent space. If $t_{i,j}$ is non-negative and $\{t_{i,j}\}$ satisfies the constraint $\sum_{j=1}^{|S|} t_{i,j} = 1$, consider the loss function for an instance *i* defined as follows:

187 188 189

192 193

195 196

197

201 202 $L(i) = -\sum_{j=1}^{|S|} t_{i,j} \log\left(\frac{e^{d_{i,j}}}{\sum_{k=1}^{|S|} e^{d_{i,k}}}\right)$ (4)

then when it reaches ideal optimum, the relationship between $t_{i,j}$ and $d_{i,j}$ satisfies:

$$softmax(d_{i,j}) = t_{i,j} \tag{5}$$

194 For detailed proof, please refer to Appendix Section B.1.

3.2 FUSION OF MULTI-MODALITY INFORMATION IN DOWNSTREAM TASKS.

During pre-training, the encoders are initialized with parameters derived from distinct reference modalities. A critical question that arises is how to effectively utilize these pre-trained models during the fine-tuning stage to improve performance on downstream tasks.

3.2.1 EARLY STAGE: MULTIMODAL MULTI-SIMILARITY

With a set of known target similarity $\{t^R\}$ from various modalities, we can transform themto multimodal space through a fusion function. There are numerous potential designs of the fusion function. For simplicity, we take linear combination as a demonstration. The multimodal generalized multi-similarity $t^M_{i,j}$ between i^{th} and j^{th} objects can be defined as follows:

$$t_{i,j}^M = fusion(\{t^R\}) \tag{6}$$

209 210

208

211

 $=\sum w_R \cdot t_{i,j}^R \tag{7}$

where $t_{i,j}^R$ represents the target similarity between i^{th} and j^{th} instance in unimodal space R, w_R is the pre-defined weights for the corresponding modal, and $\sum w_R = 1$. Then we can make $t_{i,j} = t_{i,j}^R$ in equation 3. Such that, it still satisfy the requirement of convergence (See proof in Appendix SectionB.2). In this case, the learnt similarity during pretraining will be aligned with this new combined target similarity.

216 3.2.2 INTERMEDIATE STAGE: EMBEDDING CONCATENATION AND FUSION 217

218 Intermediate fusion integrates features from various modalities after their individual encoding processes and prior to the decoding/readout stage. Let f_1, f_2, \ldots, f_n represent the feature vectors obtained 219 from these different modalities. The resulting fused feature vector can be defined as follows: 220

$$\mathbf{f}_{\text{fused}} = \text{MLP}(\text{concat}(\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_n)) \tag{8}$$

224 Where concat represents concatenation of the feature vectors. The fused features are then fed into a later readout function or decoder for downstrean tasks prediction or classification. The MLP 225 (Multi-Layer Perceptron) is used to reduce the dimension to be the same as f_i . 226

3.2.3 LATE STAGE: DECISION-LEVEL 228

221 222

227

232

233

234 235

237

238 239

240 241

242

243

244 245

246 247

248

251

253 254

229 Late fusion (or decision-level fusion) combines the outputs of models trained on different modalities 230 after they have been processed independently. Each modality is first processed separately, and their predictions are combined at a later stage.

Let p_1, p_2, \ldots, p_n be the predictions (e.g., probabilities) from different modalities. The final prediction p_{final} can be computed using a weighted sum mechanism:

$$w_i = T_i(\mathbf{f}_i) \tag{9}$$

$$p_i = \operatorname{readout}_i(\mathbf{f}_i) \tag{10}$$

$$p_{\text{final}} = \sum_{i=1}^{n} w_i p_i \tag{11}$$

Where w_i are the weights assigned to each modality's prediction, and they can be adjusted based on the importance of each modality. In particular, w_i is tunable during the learning process for respective downsteak tasks.

4 **EXPERIMENTS**

In this section, we begin by presenting the datasets and selected modalities. Subsequently, we showcase the results obtained from MMFRL. Finally, we demonstrate the explainability of the 249 learned molecular representations. (Please refer to the experimental details of pre-training and 250 fine-tuning in the Appendix Section D.)

4.1 DATASET

4.1.1 SELECTED MODALITIES FOR TARGET SIMILARITY CALCULATION

The following modalities are used for target similarity calculation. For details on training the 256 corresponding encoders to obtain fixed embeddings for these modalities, please refer to Appendix 257 Section C.1. 258

259 **Fingerprint:** Fingerprints are binary vectors that represent molecular structures, capturing the 260 presence or absence of particular substructures, fragments, or chemical features within a molecule.

261 SMILES (Simplified Molecular Input Line Entry System): SMILES offers a compact textual 262 representation of chemical structures. 263

NMR (Nuclear Magnetic Resonance): NMR spectroscopy provides detailed insights into the 264 chemical environment of atoms within a molecule. By analyzing the interactions of atomic nuclei 265 with an applied magnetic field, NMR can reveal information about the structure, dynamics, and 266 interactions of molecules, including the connectivity of atoms, functional groups, and conformational 267 changes. In our experiments, NMR_{spectrum} provides the information about the overal information of 268 molecule while NMR_{peak} provides the information about the individual atoms in the molecule. 269

Image: Images (e.g., 2D chemical structures) provide a visual representation of molecular structures.

Table 1: Study on the performances of MMFRL_{Unimodality}. The best results are denoted in bold, and the second-best are indicated with underlining among the five modalities. The first 8 tasks are for classification under evaluation of ROC-AUC, while the last three are for regression with evaluation of RMSE.

275												
276	DATA SET	BBBP	BACE	SIDER	CLINTOX	HIV	MUV	Tox21	TOXCAST	ESOL	FREESOLV	Lipo
277	SMILES	92.9±1.5	90.9±3.3	$64.9 {\pm} 0.3$	78.2±1.9	83.3±1.1	$80.1 {\pm} 2.5$	85.7±1.2	$70.5{\pm}2.5$	0.811 ± 0.109	$1.623 {\pm}~0.168$	$0.539 {\pm}~0.017$
070	NMR _{SPECTRUM}	$91.0{\pm}2.0$	$93.2{\pm}2.7$	$68.1 {\pm} 1.5$	$87.7\!\pm\!6.5$	$80.9{\pm}5.0$	$\underline{80.9{\pm}5.0}$	$85.1 {\pm} 0.4$	$71.1{\pm}0.8$	0.844 ± 0.123	$\overline{2.417 \pm 0.495}$	0.609 ± 0.031
210	IMAGE	$93.1 {\pm} 2.4$	$92.9\!\pm\!1.8$	$65.3 {\pm} 1.5$	$86.2{\pm}6.5$	$82.3{\pm}0.6$	78.7 ± 1.7	$86.0{\pm}1.0$	$71.0 {\pm} 1.6$	$0.761 {\pm}~0.068$	1.648 ± 0.045	$0.537 {\pm}~0.005$
279	FINGERPRINT	92.9 ± 2.3	91.7 ± 3.6	$65.6{\pm}0.7$	87.5 ± 6.0	$\overline{81.2 \pm 2.5}$	$82.9{\pm}3.1$	$85.3 {\pm} 1.3$	70.0 ± 1.4	0.808 ± 0.071	$1.437 {\pm}~0.134$	0.565 ± 0.017
000	NMR _{PEAK}	$93.4{\pm}2.7$	89.3 ± 1.7	$\overline{62.8 \pm 2.1}$	86.1 ± 5.4	$82.1 {\pm} 0.4$	$75.4{\pm}5.2$	84.9 ± 1.0	$70.6{\pm}0.8$	0.924 ± 0.083	$1.707 {\pm} 0.126$	$0.587 {\pm} 0.021$
280	AVERAGE	$92.8{\pm}1.9$	$91.4 {\pm} 2.7$	$65.3 {\pm} 2.0$	$85.0 {\pm} 5.7$	$81.8{\pm}2.2$	$79.4{\pm}4.0$	$85.4{\pm}0.9$	$70.6 {\pm} 1.3$	$0.830 {\pm} 0.094$	1.766 ± 0.394	$0.586 {\pm} 0.048$
281	NO PRE-TRAINING	91.9 ± 3.0	85.2 ± 0.6	57.0 ± 0.7	90.6 ± 0.6	$77.1 {\pm} 0.5$	$78.6 {\pm} 1.4$	$75.9 {\pm} 0.7$	$63.7{\pm}0.2$	$1.050 {\pm} 0.008$	$2.082 {\pm} 0.082$	$0.683 {\pm} 0.016$

All of the similarity calculation from the modalities above are listed in Appendix C.2.

4.1.2 PRE-TRAINING

270

282 283 284

285 286

287 288

289

290

291

292

NMRShiftDB-2 (Landrum, 2006) is a comprehensive database dedicated to nuclear magnetic resonance (NMR) chemical shift data, providing researchers with an extensive collection of expertannotated NMR data for various organic compounds with molecular structures (SMILES). There are around 25,000 molecules used for pre-training and no overlap with downstream task datasets. And molecular images and graphs are generated via RDkit (RDK).

4.1.3 DOWNSTREAM TASKS

295 For Downstream tasks, our model was trained on 11 drug discovery-related benchmarks sourced 296 from MoleculeNet (Wu et al., 2018a). Eight of these benchmarks were designated for classification 297 downstream tasks, including BBBP, BACE, SIDER, CLINTOX, HIV, MUV, TOX21, and ToxCast, 298 while three were allocated for regression tasks, namely ESOL, Freesolv, and Lipo. The datasets 299 were divided into train/validation/test sets using a ratio of 80%:10%:10%, accomplished through the scaffold splitter (Halgren, 1996; Landrum, 2006) from Chemprop (Yang et al., 2019; Heid et al., 300 2023), like previous works. The scaffold splitter categorizes molecular data based on substructures, 301 ensuring diverse structures in each set. Molecules are partitioned into bins, with those exceeding 302 half of the test set size assigned to training, promoting scaffold diversity in validation and test sets. 303 Remaining bins are randomly allocated until reaching the desired set sizes, creating multiple scaffold 304 splits for comprehensive evaluation. 305

- 4.2 Results
- 4.2.1 THE EFFECTIVENESS OF PRE-TRAINING

We first illustrate the impact of pre-training initialization on performance. As shown in Table 1, the average performance of pre-trained models outperform the non-pre-trained model in all tasks except for Clintox. The results of various downstream tasks indicate that different tasks may prefer different modalities. Notably, the model pre-trained with the NMR modality achieves the highest performance across three classification tasks. Similarly, the model pre-trained with the Image modality excels in three tasks, two of which are regression tasks related to solubility, aligning with findings from prior literature (Xu et al., 2023a). Additionally, the model pre-trained with The fingerprint method achieves the best performance in two tasks, including MUV, which has the largest dataset.

317 318 319

306

307 308

309

4.2.2 OVERALL PERFORMANCE OF MMFRL

As shown in Table 2 and Table 3, MMFRL demonstrates superior performance compared to all
 baseline models and the average performance of DMPNN pretrained with extra modalities across
 all 11 tasks evaluated in MoleculeNet. This robust performance highlights the effectiveness of our
 approach in leveraging multimodal data. In particular, while individual models pre-trained on other
 modalities for ClinTox fail to outperform the No-pretraining model (DMPNN), the fusion of these

Table 2: Overall performances (ROC-AUC) on classification downstream tasks. The best results are denoted in bold, and the second-best are indicated with underlining. For early fusion of MMFRL, all the predefined weight of each modality are 0.2. (Note: N-Gram is highly time-consuming on ToxCast.)

330	DATA SET	BBBP	BACE	SIDER	CLINTOX	HIV	MUV	Tox21	TOXCAST
331	ATTENTIVEFP	64.3±1.8	$78.4{\pm}2.2$	60.6 ± 3.2	84.7±0.3	75.7±1.4	$76.6 {\pm} 1.5$	76.1±0.5	63.7±0.2
332	DMPNN	91.9 ± 3.0	$85.2 {\pm} 0.6$	$57.0 {\pm} 0.7$	$90.6 {\pm} 0.6$	$77.1 {\pm} 0.5$	$78.6 {\pm} 1.4$	$75.9 {\pm} 0.7$	$63.7 {\pm} 0.2$
333	N-GRAM	$91.2 {\pm} 0.3$	79.1 ± 1.3	$63.2 {\pm} 0.5$	$87.5 {\pm} 2.7$	$78.7 {\pm} 0.4$	$76.9 {\pm} 0.7$	$76.9 {\pm} 2.7$	-
000	GEM	$72.4 {\pm} 0.4$	$85.6 {\pm} 1.1$	$67.2{\pm}0.4$	90.1 ± 1.3	$80.6 {\pm} 0.9$	$81.7 {\pm} 0.5$	$78.1 {\pm} 0.1$	$69.2 {\pm} 0.4$
334	UNI-MOL	$72.9 {\pm} 0.6$	$85.7 {\pm} 0.2$	$65.9 {\pm} 1.3$	$91.9 {\pm} 1.8$	$80.8 {\pm} 0.3$	82.1 ± 1.3	$79.6 {\pm} 0.5$	$69.6 {\pm} 0.1$
335	InfoGraph	$69.2 {\pm} 0.8$	$73.9{\pm}2.5$	$59.2 {\pm} 0.2$	75.1 ± 5.0	$74.5 {\pm} 1.8$	74.0 ± 1.5	$73.0 {\pm} 0.7$	$62.0 {\pm} 0.3$
336	GRAPHCL	67.5 ± 3.3	$68.7 {\pm} 7.8$	60.1 ± 1.3	$78.9 {\pm} 4.2$	$75.0 {\pm} 0.4$	77.1 ± 1.0	$75.0 {\pm} 0.3$	$62.8 {\pm} 0.2$
550	MOLCLR	73.3 ± 1.0	$82.8 {\pm} 0.7$	61.2 ± 3.6	$89.8 {\pm} 2.7$	$77.4 {\pm} 0.6$	$78.9 {\pm} 2.3$	74.1 ± 5.3	$65.9 {\pm} 2.1$
337	MOLCLR _{CMPNN}	$72.4 {\pm} 0.7$	$85.0 {\pm} 2.4$	59.7 ± 3.4	$88.0 {\pm} 4.0$	$77.8 {\pm} 5.5$	$74.5 {\pm} 2.1$	$78.4{\pm}2.6$	69.1 ± 1.2
338	GRAPHMVP	$72.4 {\pm} 1.6$	$81.2 {\pm} 9.0$	$63.9 {\pm} 1.2$	$79.1 {\pm} 2.8$	77.0 ± 1.2	$77.7 {\pm} 6.0$	$75.9 {\pm} 5.0$	$63.1 {\pm} 0.4$
220	UNIMODALITY _{avg}	92.8±1.9	$91.4 {\pm} 2.7$	$65.3 {\pm} 2.0$	$85.0 {\pm} 5.7$	$81.8 {\pm} 2.2$	$79.4 {\pm} 4.0$	85.4±0.9	70.6 ± 1.3
333	MMFRLearly	$91.6 {\pm} 5.0$	94.3 ± 2.4	66.4 ± 1.9	$85.3 {\pm} 6.8$	$82.0 {\pm} 2.4$	$80.6 {\pm} 3.2$	$85.2 {\pm} 0.2$	69.8 ± 1.1
340	$MMFRL_{intermediate}$	95.4±0.7	95.1±1.0	$\overline{64.3 \pm 1.2}$	93.4±1.1	81.2 ± 1.3	$83.5{\pm}1.6$	85.1 ± 0.1	$71.9 {\pm} 1.1$
341	$MMFRL_{late}$	$\underline{94.7\pm0.6}$	$91.6{\pm}2.6$	$64.2 {\pm} 1.2$	$87.0{\pm}0.4$	$82.9{\pm}0.2$	$82.1 {\pm} 1.7$	$77.7{\pm}0.5$	$70.2{\pm}0.3$
342									

Table 3: Overall performances (RMSE) on regression downstream tasks. The best results are denoted in bold, and the second-best are indicated with underlining. For early fusion of MMFRL, all the predefined weight of each modality are 0.2.

Data Set	ESOL	FreeSolv	Lipo
AttentiveFP	$0.877 {\pm} 0.029$	$2.073 {\pm} 0.183$	$0.721 {\pm} 0.001$
DMPNN	$1.050 {\pm} 0.008$	$2.082 {\pm} 0.082$	$0.683 {\pm} 0.016$
N-Gram _{RF}	1.074 ± 0.107	$2.688 {\pm} 0.085$	$0.812 {\pm} 0.028$
N-Gram _{XGB}	$1.083 {\pm} 0.082$	5.061 ± 0.744	2.072 ± 0.030
GEM	$0.798 {\pm} 0.029$	1.877 ± 0.094	$0.660 {\pm} 0.008$
Uni-Mol	$0.788 {\pm} 0.029$	$1.620 {\pm} 0.035$	$0.660 {\pm} 0.008$
MolCLR	1.113 ± 0.023	2.301 ± 0.247	$0.789 {\pm} 0.009$
MolCLR _{CMPNN}	0.911 ± 0.082	2.021 ± 0.133	$0.875 {\pm} 0.003$
Unimodality av a	$0.924 {\pm} 0.083$	1.707 ± 0.126	$0.587 {\pm} 0.021$
MMFRLearly	1.037 ± 0.170	2.093 ± 0.090	0.607 ± 0.034
MMFRL _{intermediate}	$0.730 {\pm} 0.019$	$1.465 {\pm} 0.096$	$0.552 {\pm} 0.014$
$MMFRL_{late}$	0.763 ± 0.035	$1.741 {\pm} 0.191$	$\overline{0.525 \pm 0.018}$

pre-trained models leads to improved performance. Besides, apart from Tox21 and Sider, the fusion models significantly enhances overall performance. In particular, the intermediate fusion model stands out by achieving the highest scores in seven distinct tasks, showcasing its ability to effectively combine features at a mid-level abstraction. the late fusion model achieves the top performance in two tasks. These results underscore the advantages of utilizing various fusion strategies in multimodal learning, further validating the efficacy of the MMFRL framework.

4.3 ANALYSIS OF THE FUSION EFFECT

4.3.1 GENERAL COMPARISON AMONG VARIOUS WAYS OF FUSIONS

Early Fusion is employed during the pretraining phase and is easy to implement, as it aggregates information from different modalities directly. However, its primary limitation lies in the necessity for predefined weights assigned to each modality. These weights may not accurately reflect the relevance of each modality for the specific downstream tasks, potentially leading to suboptimal performance.

Intermediate Fusion is able to capture the interaction between modalities early in the fine-tuning process, allowing for a more dynamic integration of information. This method can be particularly beneficial when different modalities provide complementary information that enhances overall performance. If the modalities effectively compensate for one another's strengths and weaknesses, Intermediate Fusion may emerge as the most effective approach.

In contrast, Late Fusion enables each modality to be explored independently, maximizing the potential
of individual modalities without interference from others. This separation allows for a thorough
examination of each modality's contribution. When certain modalities dominate the performance
metrics, Late Fusion can maximize on these strengths, ensuring that the most impactful information
is utilized effectively. This approach is especially useful in scenarios where the dominance of specific
modalities can be leveraged to enhance overall model performance.

384 385

4.3.2 EXPLAINABILITY OF LEARNT REPRESENTATIONS

To demonstrate the interpretability of learnt representations of fusion, we present post-hoc analysis
 for two tasks, ESOL and Lipo, as demonstration. The results showcase learnt representations can
 capture task-specific patterns and offer valuable insights for molecular design.

389 **ESOL** with Intermediate Fusion. As presented in Table 3, the intermediate fusion method 3.2.2 390 exhibits superior performance on the ESOL regression task for predicting solubility. To further analyze 391 this performance, we employed t-SNE to reduce the dimensionality of the molecule embeddings from 392 300 to 2, resulting in a heatmap visualized in Figure 2. The embeddings derived from individual 393 modalities prior to fusion do not display a clear pattern, showing no smooth transition from low 394 to high solubility. In contrast, the embeddings by intermediate fusion reveal a distinct and smooth 395 transition in solubility values: molecules with similar solubility cluster together, forming a gradient 396 that extends from the bottom left (indicating lower solubility) to the upper center (representing higher 397 solubility). This trend underscores the effectiveness of the intermediate fusion approach in accurately capturing the quantitative structure-activity relationships for aqueous solubility. 398

Additionally, we examined the similarity between the respective embeddings prior to intermediate fusion and the resulting fused embedding, as depicted in Figure 3. Our analysis indicates that the embeddings from each modality exhibit low similarity with the intermediate-fused representation. This observation suggests that the modalities complement each other, collectively enhancing the resulting representation of the intermediate-fused embedding.

Lipo with Late Fusion. As detailed in Table 3, the Late Fusion method (described in Section 3.2.3) demonstrates superior performance on the Lipo regression task for predicting solubility in fats, oils, lipids, and non-polar solvents. According to Equation 11, the final prediction is determined by the respective coefficients (w_i) and predictions (p_i) from each modality.

408 In Figure 4, we present the distribution of values for the coefficients, predictions, and their products 409 for each modality. Notably, the SMILES and Image modalities exhibit a broad range of values, 410 suggesting their potential for significant contributions to the final predictions. This observation aligns 411 with the strong performance achieved when pretraining using either of these two modalities, as 412 shown in Table 1. In contrast, the NMR_{Peak} values display a narrower range, indicating its role as a 413 modifier for finer adjustments in the predictions. Furthermore, we observe that the contributions from 414 NMR_{Spectrum} and Fingerprint modalities are minimal, with their corresponding values approaching 415 zero. This outcome highlights the advantages of the Late Fusion approach in effectively identifying and leveraging dominant modalities, thereby optimizing the overall predictive performance. 416

417 418

419

5 RELATED WORK

Contrastive Learning on Molecular Graphs. The primary focus within the domain of contrastive 420 learning applied to molecular graphs centers on 2D-2D graphs comparisons. Noteworthy repre-421 sentative examples: InfoGraph (Sun et al., 2019) maximizes the mutual information between the 422 representations of the graph and its substructures to guide the molecular representation learning; 423 GraphCL (You et al., 2020), MoCL (Sun et al., 2021), and MolCLR (Wang et al., 2022b) employs 424 graph augmentation techniques to construct positive pairs; MoLR (Wang et al., 2022a) establishes 425 positive pairs with reactant-product relationships. In addition to 2D-2D graph contrastive learn-426 ing, there are also noteworthy efforts exploring 2D-3D and 3D-3D contrastive learning in the field. 427 3DGCL (Moon et al., 2023) is 3D-3D contrastive learning model, establishing positive pairs with 428 conformers from the same molecules. GraphMVP (Liu et al., 2022b), GeomGCL (Li et al., 2022), 429 and 3D Informax (Stärk et al., 2022) proposes 2D–3D view contrastive learning approaches. To conclude, 2D-2D and 3D-3D comparisons are intra-modality contratsive leraning, as only one graph 430 encoder is employed in these studies. And these approaches often focus on the motif and graph levels, 431 leaving atom-level contrastive learning less explored.



Figure 2: T-SNE visualization depicting the ESOL molecule embeddings for intermediate fusion in Section 3.2.2 alongside molecules within the highlighted region. Each point in the heatmap corresponds to the embeddings of respective molecules in ESOL, with color indicating solubility levels. Red denotes higher solubility, while blue indicates lower solubility. The embeddings derived from individual modalities prior to fusion do not display a clear pattern, the embeddings by intermediate fusion forms a gradient that extends from the bottom left (indicating lower solubility) to the upper center (representing higher solubility).



Figure 3: This figure shows the distribution of similarities between each modality and the intermediate fusion embedding for ESOL. In both Cosine Similarity and Dot Product, the embeddings from each modality exhibit low similarity with the intermediate-fused representation.

Similarity Learning. Instance-wise discrimination, a crucial facet of similarity learning, involves evaluating the similarity between instances directly based on their latent representations or features (Wu et al., 2018b). Naive instance-wise discrimination relies on pairwise similarity, leading to the development of contrastive loss (Hadsell et al., 2006). Although there are improved loss functions such as triplet loss (Hoffer & Ailon, 2015), quadruplet loss (Law et al., 2013), lifted structure loss (Oh Song et al., 2016), N-pairs loss (Sohn, 2016), and angular loss (Wang et al., 2017), these methods still fall short in thoroughly capturing relationships among multiple instances simultaneously (Wang et al., 2019). To address this limitation, a joint multi-similarity loss has been proposed, incorporating pair weighting for each pair to enhance instance-wise discrimination (Wang et al., 2019; Zhang et al., 2021). Notably, it is crucial to emphasize that employing these pair weightings requires the manual



Figure 4: Lipo late fusion contribution analysis reveals that the three primary contributors are SMILES, image, and NMR_{peak}. In contrast, NMR_{spectrum} and fingerprint exhibit negligible contributions.

categorization of negative and positive pairs, as distinct weights are assigned to losses based on their categories.

511 512 513

506

507

508 509

510

6 DISCUSSION

514 515

In summary, we introduce a novel relational learning metric for molecular graph representation 516 that enhances the understanding of inter-instance relationships by capturing both local and global 517 contexts. This is the first implementation of such a generalized metric in molecular graphs.Our 518 method transforms pairwise self-similarity into relative similarity through a weighting function, 519 allowing for complex relational insights. This metric is integrated into a multimodal representation, 520 improving performance by utilizing modalities not directly accessible during fine-tuning. Empirical 521 results show that our approach, MMFRL, excels in various molecular property prediction tasks. We also demonstrate detailed study about the explainability of the learned representations, offering 522 valuable insights for drug molecule design. Despite these accomplishments, further exploration is 523 needed to achieve more effective integration of graph- and node-level similarities. Looking ahead, we 524 are enthusiastic about the prospect of applying our model to additional fields, such as social science, 525 thereby broadening its applicability and impact. 526

527

529

528 ACCESSIBILITY

- ⁵³⁰ The code and dataset will be made available upon the date of publication.
- 531 532 533

538

References

534 RDKit: Open-source cheminformatics. http://www.rdkit.org. 535

- Maria-Florina Balcan and Avrim Blum. On a theory of learning with similarity functions. In
 Proceedings of the 23rd international conference on Machine learning, pp. 73–80, 2006.
- 539 Djork-Arné Clevert, Tuan Le, Robin Winter, and Floriane Montanari. Img2mol–accurate smiles recognition from molecular graphical depictions. *Chemical science*, 12(42):14174–14181, 2021.

540 541 542	Filippo Costanti, Arian Kola, Franco Scarselli, Daniela Valensin, and Monica Bianchini. A deep learning approach to analyze nmr spectra of sh-sy5y cells for alzheimer's disease diagnosis. <i>Mathematics</i> , 11(12):2664, 2023.
543 544 545 546	Xiaomin Fang, Lihang Liu, Jieqiong Lei, Donglong He, Shanzhuo Zhang, Jingbo Zhou, Fan Wang, Hua Wu, and Haifeng Wang. Geometry-enhanced molecular representation learning for property prediction. <i>Nature Machine Intelligence</i> , 4(2):127–134, 2022.
547 548 549	Yin Fang, Qiang Zhang, Ningyu Zhang, Zhuo Chen, Xiang Zhuang, Xin Shao, Xiaohui Fan, and Huajun Chen. Knowledge graph-enhanced molecular contrastive learning with functional prompt. <i>Nature Machine Intelligence</i> , pp. 1–12, 2023.
550 551 552 553	Ioannis P Gerothanassis, Anastassios Troganis, Vassiliki Exarchou, and Klimentini Barbarossou. Nuclear magnetic resonance (nmr) spectroscopy: basic principles and phenomena, and their applications to chemistry, biology and medicine. <i>Chemistry Education Research and Practice</i> , 3 (2):229–252, 2002.
555 556 557	Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In <i>International conference on machine learning</i> , pp. 1263–1272. PMLR, 2017.
558 559 560	Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In 2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06), volume 2, pp. 1735–1742. IEEE, 2006.
561 562 563	Thomas A Halgren. Merck molecular force field. i. basis, form, scope, parameterization, and performance of mmff94. <i>Journal of computational chemistry</i> , 17(5-6):490–519, 1996.
564 565 566	Esther Heid, Kevin P Greenman, Yunsie Chung, Shih-Cheng Li, David E Graff, Florence H Vermeire, Haoyang Wu, William H Green, and Charles J McGill. Chemprop: A machine learning package for chemical property prediction. <i>Journal of Chemical Information and Modeling</i> , 2023.
567 568 569 570	Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In <i>Similarity-Based Pattern</i> <i>Recognition: Third International Workshop, SIMBAD 2015, Copenhagen, Denmark, October</i> <i>12-14, 2015. Proceedings 3</i> , pp. 84–92. Springer, 2015.
571 572	Bahador Khaleghi, Alaa Khamis, Fakhreddine O Karray, and Saiedeh N Razavi. Multisensor data fusion: A review of the state-of-the-art. <i>Information fusion</i> , 14(1):28–44, 2013.
573 574 575	Dana Lahat, Tülay Adali, and Christian Jutten. Multimodal data fusion: an overview of methods, challenges, and prospects. <i>Proceedings of the IEEE</i> , 103(9):1449–1477, 2015.
576 577	Joseph B Lambert, Eugene P Mazzola, and Clark D Ridge. <i>Nuclear magnetic resonance spectroscopy:</i> <i>an introduction to principles, applications, and experimental methods.</i> John Wiley & Sons, 2019.
578 579	Greg Landrum. Rdkit: Open-source cheminformatics. 2006. Google Scholar, 2006.
580 581	Marc T Law, Nicolas Thome, and Matthieu Cord. Quadruplet-wise image similarity learning. In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 249–256, 2013.
582 583 584 585	Shuangli Li, Jingbo Zhou, Tong Xu, Dejing Dou, and Hui Xiong. Geomgcl: Geometric graph contrastive learning for molecular property prediction. In <i>Proceedings of the Thirty-Six AAAI Conference on Artificial Intelligence</i> , pp. 4541–4549, 2022.
586 587	Hui Liu, Yibiao Huang, Xuejun Liu, and Lei Deng. Attention-wise masked graph contrastive learning for predicting molecular property. <i>Briefings in bioinformatics</i> , 23(5):bbac303, 2022a.
588 589 590 591	Shengchao Liu, Mehmet F Demirel, and Yingyu Liang. N-gram graph: Simple unsupervised repre- sentation for graphs, with applications to molecules. <i>Advances in neural information processing</i> <i>systems</i> , 32, 2019.
592 593	Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. Pre- training molecular graph representation with 3d geometry. In <i>International Conference on Learning</i> <i>Representations</i> , 2022b. URL https://openreview.net/forum?id=xQUe1pOKPam.

594 595 596 597	Pierre Mahé, Nobuhisa Ueda, Tatsuya Akutsu, Jean-Luc Perret, and Jean-Philippe Vert. Extensions of marginalized graph kernels. In <i>Proceedings of the twenty-first international conference on Machine learning</i> , pp. 70, 2004.
598 599	Kisung Moon, Hyeon-Jin Im, and Sunyoung Kwon. 3d graph contrastive learning for molecular property prediction. <i>Bioinformatics</i> , 39(6):btad371, 2023.
600 601 602 603	Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 4004–4012, 2016.
604 605	M. Pawłowski, A. Wróblewska, and S. Sysko-Romańczuk. Effective techniques for multimodal data fusion: A comparative analysis. <i>Sensors (Basel)</i> , 23(5):2381, Feb 2023.
607 608 609 610	Soujanya Poria, Erik Cambria, and Alexander Gelbukh. Deep convolutional neural network tex- tual features and multiple kernel learning for utterance-level multimodal sentiment analysis. In <i>Proceedings of the 2015 conference on empirical methods in natural language processing</i> , pp. 2539–2544, 2015.
611 612 613	Dhanesh Ramachandram and Graham W Taylor. Deep multimodal learning: A survey on recent advances and trends. <i>IEEE signal processing magazine</i> , 34(6):96–108, 2017.
614 615	Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. Advances in neural information processing systems, 29, 2016.
616 617 618 619	Hannes Stärk, Dominique Beaini, Gabriele Corso, Prudencio Tossou, Christian Dallago, Stephan Günnemann, and Pietro Liò. 3d infomax improves gnns for molecular property prediction. In <i>International Conference on Machine Learning</i> , pp. 20479–20502. PMLR, 2022.
620 621 622	Christoph Steinbeck, Stefan Krause, and Stefan Kuhn. Nmrshiftdb constructing a free chemical information system with open-source components. <i>Journal of chemical information and computer sciences</i> , 43(6):1733–1739, 2003.
623 624 625 626	Fan-Yun Sun, Jordan Hoffmann, Vikas Verma, and Jian Tang. Infograph: Unsupervised and semi- supervised graph-level representation learning via mutual information maximization. <i>arXiv preprint</i> <i>arXiv:1908.01000</i> , 2019.
627 628 629	Mengying Sun, Jing Xing, Huijun Wang, Bin Chen, and Jiayu Zhou. Mocl: data-driven molecular fingerprint via knowledge-aware contrastive learning from molecular graph. In <i>Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining</i> , pp. 3585–3594, 2021.
630 631 632 633 634	Hongwei Wang, Weijiang Li, Xiaomeng Jin, Kyunghyun Cho, Heng Ji, Jiawei Han, and Martin D. Burke. Chemical-reaction-aware molecule representation learning. In <i>International Confer-</i> <i>ence on Learning Representations</i> , 2022a. URL https://openreview.net/forum?id= 6sh3pIzKS
635 636 637	Jian Wang, Feng Zhou, Shilei Wen, Xiao Liu, and Yuanqing Lin. Deep metric learning with angular loss. In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 2593–2601, 2017.
639 640 641	Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with general pair weighting for deep metric learning. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 5022–5030, 2019.
642 643 644 645	Yifei Wang, Shiyang Chen, Guobin Chen, Ethan Shurberg, Hang Liu, and Pengyu Hong. Motif-based graph representation learning with application to chemical molecules. In <i>Informatics</i> , volume 10, pp. 8. MDPI, 2023.
646 647	Yingheng Wang, Yaosen Min, Erzhuo Shao, and Ji Wu. Molecular graph contrastive learning with parameterized explainable augmentations. In 2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pp. 1558–1563. IEEE, 2021.

648	Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive
649	learning of representations via graph neural networks. <i>Nature Machine Intelligence</i> , 4(3):279–287.
650	2022b.
651	

- Yandong Wen, Weiyang Liu, Yao Feng, Bhiksha Raj, Rita Singh, Adrian Weller, Michael J Black, 652 and Bernhard Schölkopf. Pairwise similarity learning is simple. In Proceedings of the IEEE/CVF 653 International Conference on Computer Vision, pp. 5308–5318, 2023. 654
- 655 Oliver Wieder, Stefan Kohlbacher, Mélaine Kuenemann, Arthur Garon, Pierre Ducrot, Thomas 656 Seidel, and Thierry Langer. A compact review of molecular property prediction with graph neural 657 networks. Drug Discovery Today: Technologies, 37:1-12, 2020.
- 658 Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S 659 Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. 660 Chemical science, 9(2):513-530, 2018a. 661
- 662 Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via nonparametric instance discrimination. In Proceedings of the IEEE conference on computer vision 663 and pattern recognition, pp. 3733-3742, 2018b. 664
- 665 Zhaoping Xiong, Dingyan Wang, Xiaohong Liu, Feisheng Zhong, Xiaozhe Wan, Xutong Li, Zhaojun 666 Li, Xiaomin Luo, Kaixian Chen, Hualiang Jiang, et al. Pushing the boundaries of molecular 667 representation for drug discovery with the graph attention mechanism. Journal of medicinal 668 chemistry, 63(16):8749-8760, 2019.
- Hao Xu, Yifei Wang, Yunrui Li, and Pengyu Hong. Asymmetric contrastive multimodal learning for 670 advancing chemical understanding. arXiv preprint arXiv:2311.06456, 2023a. 671
- 672 Hao Xu, Zhengyang Zhou, and Pengyu Hong. Molecular identification and peak assignment: 673 Leveraging multi-level multimodal alignment on nmr. arXiv preprint arXiv:2311.13817, 2023b.
- Kevin Yang, Kyle Swanson, Wengong Jin, Connor Coley, Philipp Eiden, Hua Gao, Angel Guzman-675 Perez, Timothy Hopper, Brian Kelley, Miriam Mathea, et al. Analyzing learned molecular 676 representations for property prediction. Journal of chemical information and modeling, 59(8): 677 3370-3388, 2019. 678
- 679 Zhuo Yang, Jianfei Song, Minjian Yang, Lin Yao, Jiahua Zhang, Hui Shi, Xiangyang Ji, Yafeng Deng, and Xiaojian Wang. Cross-modal retrieval between 13c nmr spectra and structures for compound 680 identification using deep contrastive learning. Analytical Chemistry, 93(50):16947–16955, 2021. 681
- 682 Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph 683 contrastive learning with augmentations. Advances in neural information processing systems, 33: 684 5812-5823, 2020. 685
- Li Zhang, Shitian Shen, Lingxiao Li, Han Wang, Xueying Li, and Jun Lang. Jointly multi-similarity 686 loss for deep metric learning. In 2021 IEEE International Conference on Data Mining (ICDM), pp. 1469-1474. IEEE, 2021. 688
 - Zehong Zhang, Lifan Chen, Feisheng Zhong, Dingyan Wang, Jiaxin Jiang, Sulin Zhang, Hualiang Jiang, Mingyue Zheng, and Xutong Li. Graph neural network approaches for drug-target interactions. Current Opinion in Structural Biology, 73:102327, 2022.
- Mingkai Zheng, Shan You, Fei Wang, Chen Qian, Changshui Zhang, Xiaogang Wang, and Chang 693 Xu. Ressl: Relational self-supervised learning with weak augmentation. Advances in Neural 694 Information Processing Systems, 34:2543–2555, 2021. 695
- 696 Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Lin-697 feng Zhang, and Guolin Ke. Uni-mol: A universal 3d molecular representation learning framework. In International Conference on Learning Representations, 2023. URL https: 698 //api.semanticscholar.org/CorpusID:259298651. 699

669

674

687

689

690

691

692



B SUPPLEMENTARY PROOF

749 750

751 B.1 REVISITING THEOREM OF CONVERGENT SIMILARITY LEARNING 752

Let S be a set of instances with size |S|, and let \mathcal{P} represent the tunable latent representations of instances in S such that $|\mathcal{P}| = |S|$. For any two instances $i, j \in S$, their latent representations are denoted by \mathcal{P}_i and \mathcal{P}_j , respectively. Let $t_{i,j}$ represent the target similarity between instances i and jin a given domain, and $d_{i,j}$ be the similarity between \mathcal{P}_i and \mathcal{P}_j in the latent space. **Theorem B.1** (Theorem of Convergent Similarity learning). Given $t_{i,j}$ is non-negative and $\{t_{i,j}\}$ satisfies the constraint $\sum_{j=1}^{|S|} t_{i,j} = 1$, consider the loss function for an instance *i* defined as follows:

$$L(i) = -\sum_{j=1}^{|S|} t_{i,j} \log\left(\frac{e^{d_{i,j}}}{\sum_{k=1}^{|S|} e^{d_{i,k}}}\right)$$
(B.1)

then when it reaches ideal optimum, the relationship between $t_{i,j}$ and $d_{i,j}$ satisfies:

$$softmax(d_{i,j}) = t_{i,j} \tag{B.2}$$

Proof. In order to optimize the loss L(i), we need to set the following partial derivative to be 0 for each $d_{i,j}$ with $1 \leq j \leq |\mathcal{M}|$. Here are the detailed steps:

$$\frac{\partial L(i)}{\partial d_{i,j}} = \frac{\partial}{\partial d_{i,j}} \underbrace{\left(-t_{i,j}\log\frac{e^{d_{i,j}}}{e^{d_{i,j}} + \sum_{k \neq j} e^{d_{i,k}}}\right)}_{\text{When the numerator includes } e^{d_{i,j}}} + \frac{\partial}{\partial d_{i,j}} \underbrace{\left(\sum_{k \neq j} -t_{i,k}\log\frac{e^{d_{i,k}}}{e^{d_{i,j}} + \sum_{k \neq j} e^{d_{i,k}}}\right)}_{\text{When the numerator includes } e^{d_{i,j}}}$$

When the numerator does not include
$$e^{d_{i,j}}$$

 $= -(t_{i,j} - t_{i,j} \cdot \operatorname{softmax}(d_{i,j})) - \sum_{k \neq j} t_{i,k} \cdot \operatorname{softmax}(d_{i,j})$ $- - \left(t_{i,j} - \left(t_{i,j} + \sum_{k \neq j} t_{i,k} \right) \cdot \operatorname{softmax}(d_{i,j}) \right)$

$$= -\left(t_{i,j} - \left(t_{i,j} + \sum_{k \neq j} t_{i,k}\right) \cdot \operatorname{softmax}(d_{i,j})\right)$$

Since $\sum_{l=1}^{|\mathcal{M}_l|} t_{i,l} = 1$, we can further simplify it as

$$\frac{\partial L(i)}{\partial d_{i,j}} = -(t_{i,j} - \text{softmax}(d_{i,j}))$$

In order to optimize, we need to see the above partial derivative to be 0:

$$\frac{\partial L(i)}{\partial d_{i,j}} = -(t_{i,j} - \operatorname{softmax}(d_{i,j})) = 0$$

In addition, the corresponding second partial derivative denoted as $\frac{\partial L(i)}{\partial d_{i,j}^2}$ manifests as follows:

$$\frac{\partial L(i)}{\partial d_{i,j}^2} = \operatorname{softmax}(d_{i,j})(1 - \operatorname{softmax}(d_{i,j}))$$

As softmax $(d_{i,j})$ takes values within the open interval (0,1), it follows that $\frac{\partial L(i)}{\partial d_{i,j}^2}$ is always positive. Consequently, the global optimum is global minimum.

Furthermore, when it comes to optimum:

$$t_{i,j} = \operatorname{softmax}(d_{i,j})$$
$$d_{i,j} = \log(t_{i,j}) + \log\left(\sum_{1 \le l \le |\mathcal{M}|} e^{d_{i,j}}\right)$$

It is easy to show that when it reaches optimum, $d_{i,j}$ is consistent with target similarity metric $t_{i,j}$. Without loss of generosity, suppose $t_{i,j} > t_{i,j'}$:

$$d_{i,j} - d_{i,j'} = \log(t_{i,j}) + \log\left(\sum_{1 \le l \le |\mathcal{M}|} e^{d_{il}}\right) - \left(\log(t_{i,j'}) + \log\left(\sum_{1 \le l \le |\mathcal{M}|} e^{d_{il}}\right)\right)$$
$$= \log(t_{i,j}) - \log(t_{i,j'})$$
$$= \log\left(\frac{t_{i,j}}{t_{i,j'}}\right) > 0$$

B.2 GUARANTEE OF SUM OF FUSED MULTIMODAL SIMILARITY

Given sets of uni-modal generalized similarity $\{t^R\}$ and $\sum w_{t^R} = 1$, the sum of fused multimodal similarity also equals 1, as demonstrated below:

$$\sum_{i,j} (t_{i,j}^R) = \sum_{i,j} \sum_{k} (w_R \cdot t_{i,j}^R)$$
$$= \sum_{i,j} (w_R \sum_{k} t_{i,j}^R)$$
$$= \sum_{i,j} w_R \cdot 1 = 1$$

REVISITING TARGET SIMILARITY SETTINGS С

C.1 ENCODERS & PACKAGES

To derive the target similarities, we need to reply on pre-trained encoders or well-defined packages as follows:

T 11 C 1	F 1 1	1	1. 1	101
Table C.1:	Encoders and	packages used	to produce	e self-similarities

829	Unimodal	Representation	Encoder/Package	Pre-trained Source
830	Image	2D image	CNN	Img2mol (Clevert et al., 2021)
831	SMILES	Sequence	Transformer	CReSS (Yang et al., 2021)
832	¹³ CNMR Spectrum	Sequence	1D CNN	AutoEncoder (Costanti et al., 2023)
833	¹³ CNMR peak	Scalar	NMRShiftDB2 (Steinbeck et al., 2003)	N/A
834	Fingerprint	Sequence	RDKit (Landrum, 2006)	N/A

C.2 TARGET SIMILARITY AT GRAPH LEVEL

Fingerprint. The mathematical formula of fingerprint similarity, denoted as $S_{i,j}^F$, can be viewed as follows:

> $S_{i,j}^F = Tanimoto(A, B) = \frac{|A \cap B|}{|A \cup B|}$ (C.1)

where A and B are sets of molecular fragments for molecule i and j, and $|A \cap B|$ and $|A \cup B|$ denote the size of their intersection and union, respectively.

Image. The self-similarity for Image, denoted as $S_{i,j}^{I}$, can be defined as follows:

$$S_{i,j}^{I} = Cos(\mathcal{V}_i, \mathcal{V}_j) = \frac{\mathcal{V}_i \cdot \mathcal{V}_j^T}{\|\mathcal{V}_i\| \cdot \|\mathcal{V}_j\|}$$
(C.2)

where $\mathcal{V}_i, \mathcal{V}_j$ represents the embedding of Image for two given molecules.

NMR Spectrum. The self-similarity for NMR spectrum, denoted as $S_{i,j}^C$, can be defined as follows:

$$S_{i,j}^C = Cos(\mathcal{V}_i, \mathcal{V}_j) = \frac{\mathcal{V}_i \cdot \mathcal{V}_j^T}{\|\mathcal{V}_i\| \cdot \|\mathcal{V}_j\|}$$
(C.3)

where $\mathcal{V}_i, \mathcal{V}_j$ represents the embedding of NMR spectra for two given molecules.

Smiles. The self-similarity for Smiles, denoted as $S_{i,j}^S$, can be defined as follows:

$$S_{i,j}^{S} = Cos(\mathcal{V}_i, \mathcal{V}_j) = \frac{\mathcal{V}_i \cdot \mathcal{V}_j^T}{\|\mathcal{V}_i\| \cdot \|\mathcal{V}_j\|}$$
(C.4)

where $\mathcal{V}_i, \mathcal{V}_i$ represents the embedding of Smiles for two given molecules.

NMR Peak The similarity among nodes (atoms) is derived from the positions of their signal peaks on ¹³C NMR spectra, measured in parts per million (ppm). The ppm values are continuous, typically ranging from 0 to 200 (see more introduction of ppm in Appendix C.3). The self-similarity of NMR peaks $S_{l,m}^P$ can be defined as following:

$$S_{l,m}^{P} = \frac{\tau_2}{|ppm_l - ppm_m| + \tau_1} \tag{C.5}$$

where ppm_l and ppm_m are the positions of NMR peaks for the l^{th} , m^{th} Carbon atom, τ_1 and τ_2 are temperature hyper-parameter.

C.3 A BRIEF INTRODUCTION TO PPM FOR NMR PEAK

In chemistry, ¹³C NMR stands out as a common technique for structural analysis by revealing molecular structures by elucidating the chemical environments of carbon atoms and their magnetic responses to external fields (Gerothanassis et al., 2002; Lambert et al., 2019). It quantifies these features in parts per million (ppm) relative to a reference compound, such as tetramethylsilane (TMS), thereby simplifying comparisons across experiments. As a result, the continuous peak positions, measured in parts per million (ppm), offer a robust knowledge span—a natural ordering metric that can be employed to derive measures of similarity (Xu et al., 2023b).

C.4 CONFIGURATION OF EARLY FUSION

A simple linear combination is used to formulate the multimodal relational similarity $t_{i,j}^M$ between the i^{th} and j^{th} molecules, represented as as follows:

$$t_{i,j}^{M} = w_{SM} \cdot t_{i,j}^{SM} + w_{C} \cdot t_{i,j}^{C} + w_{I} \cdot t_{i,j}^{I} + w_{F} \cdot t_{i,j}^{F} + w_{F} \cdot t_{i,j}^{F} + w_{P} \cdot t_{i,j}^{P}$$
(C.6)

where $t_{i,j}^{SM}$ denotes the similarity based on SMILES, $t_{i,j}^C$ denotes the similarity with respect to ¹³C NMR spectrum, $t_{i,j}^I$ denotes the similarity regarding images, F denotes the similarity based on fingerprints, and P denotes the similarity based on fingerprints. w_{SM} , w_C , w_I , and w_F are the pre-defined weights for their respective similarity, and $w_{SM} + w_C + w_I + w_F + w_P = 1$.

893 894 895

904

906

890

891

892

867

868

870

871 872 873

874

882

883 884

885

D EXPERIMENTAL SETTINGS

- 896 897
 - D.1 PRE-TRAINING SETTING

During pretraining, we utilized an Adam optimizer with a learning rate set to 0.001, spanning 200
epochs and employing a batch size of 256. The model was trained on around 25,000 data points. The
NMR data were experimental data, extracted from NMRShiftDB2 (Steinbeck et al., 2003). Other
chemical modalities, such as images, fingerprints and graphs, were produced from SMILES by RDKit
(Landrum, 2006).

- 905 D.2 FINE-TUNING SETTING
- 907 D.2.1 DATASETS

908 For fine-tuning, our model was trained on 11 drug discovery-related benchmarks sourced from 909 MoleculeNet (Wu et al., 2018a). Eight of these benchmarks were designated for classification 910 downstream tasks, including BBBP, BACE, SIDER, CLINTOX, HIV, MUV, TOX21, and ToxCast, 911 while three were allocated for regression tasks, namely ESOL, Freesolv, and Lipo. The datasets 912 were divided into train/validation/test sets using a ratio of 80%:10%:10%, accomplished through the 913 scaffold splitter (Halgren, 1996; Landrum, 2006) from Chemprop (Yang et al., 2019; Heid et al., 914 2023), like previous works. The scaffold splitter categorizes molecular data based on substructures, 915 ensuring diverse structures in each set. Molecules are partitioned into bins, with those exceeding half of the test set size assigned to training, promoting scaffold diversity in validation and test sets. 916 Remaining bins are randomly allocated until reaching the desired set sizes, creating multiple scaffold 917 splits for comprehensive evaluation.

918 D.2.2 BASELINES

We systematically compared MMFRL's performance with various state-of-the-art baseline models across different categories. In the realm of supervised models, AttentiveFP (Xiong et al., 2019) and DMPNN (Yang et al., 2019) stand out by leveraging graph attention networks and node-edge interactive message passing, respectively. The unsupervised learning method N-Gram (Liu et al., 2019) employs graph embeddings and short walks for graph representation. Predictive self-supervised learning methods, such as GEM (Fang et al., 2022) and Uni-Mol (Zhou et al., 2023), are specifically designed for predicting molecular geometric information. Moreover, our evaluation encompasses a range of contrastive learning methods, namely InfoGraph (Sun et al., 2019), GraphCL (You et al., 2020), MolCLR (Wang et al., 2022b), and GraphMVP (Liu et al., 2022b), all serving as essential baselines. The baseline results are collected from recent works (Fang et al., 2022; Zhou et al., 2023; Moon et al., 2023; Fang et al., 2023).

931 D.2.3 EVALUATION

To assess the effectiveness of our fine-tuned model, we measure the ROC-AUC for classification downstream tasks, and the root mean squared error (RMSE) metric for regression tasks. In order to ensure a fair and robust comparisons, we conduct three independent runs using three different random seeds for scaffold splitting across all datasets. The reported performance metrics are then averaged across these runs, and the standard deviation is computed as prior works.