

---

# Learning Substructure Invariance for Out-of-Distribution Molecular Representations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

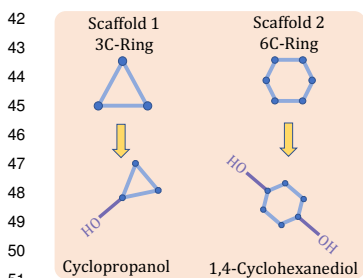
1 Molecule representation learning (MRL) has been extensively studied and current  
2 methods have shown promising power for various tasks, e.g., molecular property  
3 prediction and target identification. However, a common hypothesis of existing  
4 methods is that either the model development or experimental evaluation is mostly  
5 based on i.i.d. data across training and testing. Such a hypothesis can be violated  
6 in real-world applications where testing molecules could come from new environ-  
7 ments, bringing about serious performance degradation or unexpected prediction.  
8 We propose a new representation learning framework entitled MoleOOD to enhance  
9 the robustness of MRL models against such distribution shifts, motivated by an  
10 observation that the (bio)chemical properties of molecules are usually invariantly  
11 associated with certain privileged molecular substructures across different environ-  
12 ments (e.g., scaffolds, sizes, etc.). Specifically, We introduce an environment infer-  
13 ence model to identify the latent factors that impact data generation from different  
14 distributions in a fully data-driven manner. We also propose a new learning objec-  
15 tive to guide the molecule encoder to leverage environment-invariant substructures  
16 that more stably relate with the labels across environments. Extensive experiments  
17 on ten real-world datasets demonstrate that our model has a stronger generalization  
18 ability than existing methods under various out-of-distribution (OOD) settings,  
19 despite the absence of manual specifications of environments. Particularly, our  
20 method achieves up to 5.9% and 3.9% improvement over the strongest baselines  
21 on OGB and DrugOOD benchmarks in terms of ROC-AUC, respectively.

## 22 1 Introduction

23 Predicting molecular properties plays an important role in many related applications like drug  
24 discovery [11] and material design [44]. These professional tasks conventionally take great efforts by  
25 experts e.g. in chemistry and pharmacology. Recent years have witnessed inspiring breakthroughs on  
26 building effective machine learning models for scientific discovery, and solid progress has been made  
27 along the avenue of ML-based molecule representation learning (MRL). In general, MRL aims at  
28 embedding a molecule into a vector in latent space as a foundation model, on top of which the learned  
29 representations could be used for a variety of downstream tasks, such as target identification [60],  
30 retrosynthetic analysis [56], search of antibiotics [45], virtual screening [36] for drug discovery, etc.

31 The challenge, however, is that existing MRL methods are mostly based on an underlying hypothesis  
32 that training and testing molecules are independently sampled from an identical environment, yet  
33 real-world environments are often dynamic and uncertain, which requires the model to effectively  
34 handle distribution shifts. In fact, the available experimental molecule data are rather limited  
35 while the candidate molecules to be tested are often diverse, coming from unknown environments.  
36 Taking the virtual screening [36] as an example (which is a common protocol in drug discovery and  
37 usually for target identification), the prediction model is typically trained on some known target

38 proteins. However, some unpredictable events like COVID-19 may occur, bringing new targets from  
39 unknown distributions. Similar scenarios where training and testing data are sampled from different  
40 distributions are common in real world, posing an urgency for strengthening current MRL methods  
41 regarding out-of-distribution (OOD) generalization [37, 5, 39].



52 Figure 1: An example: the  
53 shared substructure hydroxy  
54 (-OH) invariantly contrib-  
55 utes to the water solubility  
56 of the two molecules which  
57 contain different scaffolds,  
58 i.e. sampled from different  
59 environments by definition.

60 the former is 3C-ring and the latter is 6C-ring. The common practice specifies environments as some  
61 prominent information of the molecules e.g. scaffold pattern [29, 20]. Thus, the data-generating  
62 environments and the induced distributions which these two molecules are sampled from can be  
63 considered different [20]. Though sampled from different distributions, they are both readily soluble  
64 in water due to the invariant substructure hydroxy [21] shared across different environments. Hence,  
65 a promising paradigm would be to learn the causal data-generating invariance from the substructures  
66 across environments, regarding a certain property, for the OOD generalization purpose.

67 Another important observation for consideration is that existing specifications for environments are  
68 often handcrafted or rule-based and not structured, which could provide insufficient information  
69 for capturing the fundamental relations across domains from the casual data-generating perspective.  
70 Besides, some studies [24, 14] show that directly utilizing such environment labels as input when  
71 adapting existing OOD generalization methods to MRL tasks can be problematic. Furthermore,  
72 manual specifications of environments may be unavailable in reality. Hence, we aim to develop a  
73 label-free model that does not rely on the above ad-hoc environment labels. As shown later, our model  
74 can infer the environment labels in an unsupervised manner, namely for environment clustering.

75 To achieve robust molecule representation for OOD generalization and overcome potentially unreli-  
76 able environment labels, we devise a new MRL framework without explicitly using the environment  
77 label information. We first formulate OOD generalization for molecular property prediction by  
78 introducing a latent variable for environments that affect the data generation. Then we analyze  
79 the essential cause behind the failure of existing MRL models and propose a new learning scheme  
80 based on the invariance principle [52]. The training procedures contain two steps: 1) optimize an  
81 environment inference model from training data; 2) optimize a molecule encoder and a predictor.  
82 Our general framework can integrate existing GNN backbones and achieve improvements on four  
83 OGB molecular property prediction tasks [20], as shown in our experimental results. As for a newly  
84 released benchmark for drug-oriented OOD learning [24], even without access to environment labels,  
85 our method can still outperform state-of-the-art models that rely on environment labels for training in  
86 five out of six datasets. **The contributions of this paper are:**

- 87 • We formulate the OOD problem for molecule representation learning, by particularly incorporating  
88 an important observation that the substructure of molecule can convey invariant casual information  
89 across environments, regarding certain property prediction tasks. To our best knowledge, this is  
90 the first work that explicitly models such substructures for OOD molecule representation.

<sup>1</sup>As a 2-D structural molecular framework [4], the scaffold reduces the chemical structure of a molecule to its core components, which can be obtained by removing side chains and only reserving the rings and parts connecting rings [58]. The scaffold can be an indicator to define a specific environment [29, 20].

- 91 • Under the environment-invariance principle with specific substructure invariance priors, we propose  
92 a new learning objective to learn robust representations. In particular, our model does not require  
93 environment labels which in fact can be noisy and unreliable, but instead achieve environment  
94 inference in an unsupervised manner. To our best knowledge, this is the first work for environment  
95 clustering for molecule representation, which has not been fulfilled in existing OOD literature  
96 beyond molecules (thanks to the specific substructure-property invariance priors).
- 97 • We conduct extensive experiments on ten public datasets. Results demonstrate that our model  
98 yields consistent and significant improvements over various existing MRL methods as backbones  
99 and also achieves competitive or even superior prediction compared to state-of-the-art models  
100 tailored to OOD learning with environment labels used as extra inputs in both training and testing.  
101 Particularly, our method achieves up to 5.9% higher ROC-AUC on public OGB molecular property  
102 prediction benchmarks than the counterpart model trained with traditional objective. Besides,  
103 for drug-oriented benchmarks DrugOOD, when environment labels are not used, our model still  
104 outperforms several SOTA approaches tailored for general OOD learning (using environment  
105 labels as extra training information) by up to 3.9% w.r.t. ROC-AUC.

## 106 2 Backgrounds and Related Works

107 **Out-of-Distribution Generalization.** Deep neural networks are prone to suffering significant per-  
108 formance degradation under distribution shifts, motivating a surge of works on OOD generalization.  
109 Recent studies [42, 2, 7, 52] assume that there is a potential environment variable  $e$  accounting for  
110 the distribution shift between the training and testing data. In general cases the goal is to predict the  
111 target label  $y$  given the associated input  $x$ . Then, the OOD problem could be formally formulated as:

$$\min_f \max_{e \in \mathcal{E}} \mathbb{E}_{(x,y) \sim p(x,y|e=e)} [l(f(x), y)|e], \quad (1)$$

112 where  $\mathcal{E}$  denotes the support of environments,  $f(\cdot)$  is the prediction model and  $l(\cdot, \cdot)$  represents a  
113 loss function. Notice that  $\mathbb{E}_{(x,y) \sim p(x,y|e=e)} [l(f(x), y)|e]$  is called the **risk function** under a given  
114 environment  $e$  and denoted as  $\mathcal{R}_e(\mathbf{x}^e, \mathbf{y}^e)$  [30].

115 **Invariant Learning.** There is an emerging line of research [42, 2, 8, 10] regarding invariant pre-  
116 dictor learning, for solving the OOD generalization problem. These methods propose to find an  
117 invariant predictor that could uncover invariant relationships between inputs and targets across all  
118 environments [30]. The invariant predictor should learn an invariant representation satisfying such a  
119 **invariance principle**: 1) **sufficiency**: shows sufficient predictive power for the target, 2) **invariance**:  
120 contributes to equal (optimal) performance for the downstream tasks across all environments. A recent  
121 work [52] leverages invariance principle to handle distribution shifts on graphs, but it is designed for  
122 node-level prediction. In contrast, we focus on molecular graph-level classification tasks.

123 **Molecule Representation Learning.** Existing molecule representation learning methods can be  
124 classified into two categories. The first is SMILES-based methods where SMILES refers to Sim-  
125 plified Molecular Input Line Entry System [1]. They use language models to process the textual  
126 representation (SMILES) of a molecule, for example, Transformer [48] or BERT [13]. SMILES is a  
127 linear encoding for molecules and highly depends on the traverse order of molecule graphs. There-  
128 fore its expressiveness is limited for problems like medication recommendation which we believe  
129 calls for fine-grained molecular structure extraction. Beyond the above linear encoding protocol,  
130 structure-based methods are also developed, which can be further classified into fingerprint-based  
131 and graph neural networks (GNN)-based methods. The molecular fingerprint techniques date back to  
132 the Morgan fingerprints [38]. However, those fingerprint-based methods are often handcrafted and  
133 not trained in an end-to-end fashion [23]. Since molecules can be viewed as structured graphs, graph  
134 neural networks have been widely used to learn molecule representation [25, 16, 21].

135 Existing general OOD methods [2, 46, 15, 61, 43] are not tailored to such non-Euclidean structured  
136 data, i.e. molecules. In several recent works on molecule property classification tasks [62, 50], the  
137 importance of molecular substructures has been emphasized and such inductive bias is incorporated  
138 into the design of those models. However, they are still based on the i.i.d. assumption and do not  
139 leverage those invariant substructure across different environments to achieve robust representations.  
140 In this paper, we propose a general framework orthogonal to these MRL studies to bridge OOD and  
141 MRL, which can adopt any existing MRL methods as the backbone to improve their robustness.

## 142 3 Methodology

### 143 3.1 Problem Formulation

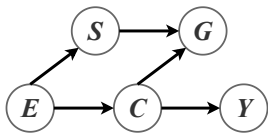
144 We propose a OOD generalization framework tailored for molecule representation learning, entitled  
145 MoleOOD. All the random variables and the corresponding realizations are denoted as bold and thin  
146 letters, respectively. We first formulate the OOD generalization problem for MRL.

147 **OOD Generalization Problem on Molecule Representation Learning.** A molecular graph can be  
148 represented as  $G = (V, E)$ , where  $V$  is the graph’s node set corresponding to atoms constituting the  
149 molecule and  $E$  denotes the graph’s edge sets corresponding to chemical bonds. The training and test-  
150 ing molecule graph datasets are denoted as  $\mathcal{G}^{train} = \{(G_i, y_i)\}_{i=1}^{N^{train}}$  and  $\mathcal{G}^{test} = \{(G_i, y_i)\}_{i=1}^{N^{test}}$ .  
151 Notice that the test dataset is drawn outside the distribution of the training dataset. The goal of  
152 molecule representation learning task is to predict the target label  $y$  given the associated input  
153 molecule  $G$ . Based on Eq. 1, we can formulate the OOD problem on MRL tasks as:

$$\min_f \max_{e \in \mathcal{E}} \mathbb{E}_{(G_i, y_i) \sim p(\mathbf{G}, \mathbf{y} | e=e)} [l(f(G_i), y_i) | e]. \quad (2)$$

154 The difficulty of this problem is that the training data only cover very limited environments in  $\mathcal{E}$  while  
155 the model is expected to perform well on all the environments. Before delving into our solution, we  
156 first provide a causal perspective to understand the data-generating process and shed light on the  
157 limitations of existing MRL models.

### 158 3.2 Understanding the Data-generating Process



164 Figure 2: SCM

159 We elaborate our approach in the context of molecule property classification  
160 tasks in this paper. Recalling the two molecules *Cyclopropanol* ( $C_3H_6O$ ) and *1,4-Cyclohexanediol* ( $C_6H_{12}O_2$ ) used for illustration in  
161 Sec. 1, they are sampled from different environments. Because both of  
162 them contain the hydroxy ( $-OH$ ), which we can call invariant or causal  
163 substructure in this case, these two molecules are readily soluble in water.  
164 We formalize such a data-generating process of molecule property predic-  
165 tion with a general Structural Causal Model (SCM) [17, 40] in Fig. 2. The abstract data variables  
166 are denoted by the nodes and the directed arrows represent the causalities. This SCM illustrates  
167 the causalities among variables:  $E$  as the environment,  $S$  as the spurious substructures,  $C$  as the  
168 invariant/causal substructures w.r.t  $Y$ ,  $G$  as the instance molecule graph,  $Y$  as the ground-truth label.

- 170 •  $S \leftarrow E \rightarrow C$ : the environmental variable impacts the underlying data generating distribution.  
171 Furthermore, substructures could be divided into causal and spurious ones across all environments.
- 172 •  $S \rightarrow G \leftarrow C$ : an instance molecule graph is made up of the causal and spurious substructures.
- 173 •  $C \rightarrow Y$ :  $Y$ , the ground-truth label, is only determined by  $C$ . This causation is the focus of our work.

174 Taking *Cyclopropanol* ( $C_3H_6O$ ) as an example, we can specify  $E$  as the *3C-ring* scaffold,  $C$  as the  
175 substructure hydroxy ( $-OH$ ),  $S$  as the substructures aside from hydroxy,  $G$  as *Cyclopropanol*,  $Y$  as  
176 good water solubility. The good water solubility  $Y$  is only attributed to the invariant substructure  
177 hydroxy, i.e.  $C$ , rather than other spurious substructures  $S$ .

178 Existing MRL methods do not differentiate invariant and spurious substructures. Hence, the spurious  
179 correlations between irrelevant substructures  $S$  and the target label  $Y$  will be encoded to learned  
180 molecular representations. When tested on unseen environments, the downstream classifier will be  
181 easily misled by these spurious correlations [53]. With the knowledge that (bio)chemical properties  
182 of a molecule are usually associated with a few privileged substructures [28, 41, 63, 26], we aim  
183 to suppress such spurious correlations and leverage environment-invariant substructures that more  
184 stably relate with the labels across environments to learn invariant molecular representations. Notice  
185 that the learned invariant molecular representations should satisfy the invariance principle mentioned  
186 in Sec. 2. We next introduce our method formally and then give the instantiation of our model.

### 187 3.3 Model Formulation

188 The framework contains two parts, the fronted molecule encoder  $\Phi$  for learning an “invariant  
189 representation” of the input molecule graph and the back-end predictor  $\omega$  for final prediction. Solving  
190 the formulation in Eq. 2 directly is intractable in practice since we cannot know all the environments,

191 i.e. obtain a complete support set  $\mathcal{E}$ . We resort to minimizing the expectation of risks from different  
 192 environments known in the training data,

$$\min_{\omega, \Phi} \mathbb{E}_{\mathbf{e}}[\mathcal{R}_{\mathbf{e}}(\mathbf{G}^{\mathbf{e}}, \mathbf{y}^{\mathbf{e}})], \text{ s.t. } \mathbf{y} \perp\!\!\!\perp \mathbf{e} \mid \Phi(\mathbf{G}), \quad (3)$$

193 where  $f = \omega \circ \Phi$  and  $\perp\!\!\!\perp$  denotes probabilistic independence. All learnable parameters of the molecule  
 194 encoder  $\Phi$  and the predictor  $\omega$  are included in  $\theta$ . Different from Eq. 2, we add an extra invariance  
 195 constraint  $\mathbf{y} \perp\!\!\!\perp \mathbf{e} \mid \Phi(\mathbf{G})$ , which is used to suppress spurious correlations [8]. Since assessing  
 196 causality is challenging, we could rethink the problem on the basis of information theory. Recall  
 197 that we hope to let the molecule encoder leverage environment-invariant substructures and learn a  
 198 molecular representation  $\Phi(G)$  given a molecule  $G$ . Our goal is to maximize the predictive power  
 199 of  $\Phi(\mathbf{G})$  on  $\mathbf{y}$ , which can be measured by mutual information between  $\Phi(\mathbf{G})$  and  $\mathbf{y}$ . Meanwhile,  
 200 probabilistic independence between  $\mathbf{y}$  and  $\mathbf{e}$  given  $\Phi(\mathbf{G})$  can be achieved via minimizing their mutual  
 201 information. For convenience, we denote  $\Phi(\mathbf{G})$  as  $\mathbf{z}$  and Eq. 3 can be approximately solved by:

$$\max_{\omega, \Phi} \mathbb{I}(\mathbf{z}; \mathbf{y}), \text{ s.t. } \min_{\omega, \Phi} \mathbb{I}(\mathbf{y}; \mathbf{e} \mid \mathbf{z}). \quad (4)$$

202 Treating the outputs of  $\omega$  and  $\Phi$  as distribution  $q_{\theta}(\mathbf{z} \mid \mathbf{G})$  and  $q_{\theta}(\mathbf{y} \mid \mathbf{z})$  respectively, Eq. 4 can be  
 203 specified as:

$$\max_{q_{\theta}(\mathbf{y} \mid \mathbf{z}), q_{\theta}(\mathbf{z} \mid \mathbf{G})} \mathbb{I}(\mathbf{z}; \mathbf{y}), \text{ s.t. } \min_{q_{\theta}(\mathbf{y} \mid \mathbf{z}), q_{\theta}(\mathbf{z} \mid \mathbf{G})} \mathbb{I}(\mathbf{y}; \mathbf{e} \mid \mathbf{z}). \quad (5)$$

204 Now, we have arrived at a clearer but still intractable optimization objective. Before specifying the  
 205 practical instantiation of Eq. 5, let’s discuss on the environment variable  $\mathbf{e}$  first.

206 In practice, due to the non-trivial efforts to label the molecular environments, manual specifications  
 207 of the environments may be unavailable in many cases. We may directly label molecules to different  
 208 environments in terms of their scaffolds when the environment label is unavailable. But this is  
 209 unreasonable in practice, because the final total environment number will be too large. Taking the  
 210 dataset HIV for molecule property prediction tasks released by Open Graph Benchmark [20] for  
 211 an example, OGB uses scaffold to split the molecules into different environments. Assuming that  
 212 we regard each scaffold as an environment directly, 41,127 molecules in HIV are partitioned into  
 213 19,076 environments (see details in Appendix E). This environment count is much larger than other  
 214 OOD datasets from other domains, e.g. Camelyon17<sup>2</sup> [3], CivilComments<sup>3</sup> [6], etc. Even though  
 215 some datasets may provide manual specifications of environments, the environment counts are also  
 216 too large, which is unfriendly to existing OOD models [24, 14]. Therefore, we propose to design an  
 217 environment-inference model  $\psi$  to partition the molecule into different environments with a relatively  
 218 smaller environment count. We denote the environment count as a hyper-parameter  $k$ .

219 Given prior  $p(\mathbf{e} \mid \mathbf{G})$ , we need to maximize the log likelihood of  $p_{\tau}(\mathbf{y} \mid \mathbf{G})$  and then obtain the posterior  
 220  $p_{\tau}(\mathbf{e} \mid \mathbf{G}, \mathbf{y})$ , which are parameterized by  $\tau$ . Since there is no analytical solutions to the true posterior,  
 221 here we use variational inference (VI) to approximate it. Specifically, we introduce a variational  
 222 distribution  $q_{\kappa}(\mathbf{e} \mid \mathbf{G}, \mathbf{y})$  parameterized by  $\kappa$  to approximate  $p_{\tau}(\mathbf{e} \mid \mathbf{G}, \mathbf{y})$ .

223 **Proposition 1.** *The Evidence Lower Bound (ELBO) of the observed molecule graph and correspond-*  
 224 *ing label tuple  $(G, y)$ :  $\mathcal{L}(\tau, \kappa; (G, y)) = \mathbb{E}_{q_{\kappa}}[\log p_{\tau}(y \mid G, e)] - D_{KL}(q_{\kappa}(e \mid G, y) \parallel p_{\tau}(e \mid G))$ .*

225 Our goal is to minimize the Kullback-Leibler (KL) divergence between  $q_{\kappa}(\mathbf{e} \mid \mathbf{G}, \mathbf{y})$  and  $p_{\tau}(\mathbf{e} \mid \mathbf{G}, \mathbf{y})$ ,  
 226 i.e.  $D_{KL}(q_{\kappa}(\mathbf{e} \mid \mathbf{G}, \mathbf{y}) \parallel p_{\tau}(\mathbf{e} \mid \mathbf{G}, \mathbf{y}))$ , which is equivalent to maximizing the ELBO in Proposition 1.  
 227 Then, the objective used to train this environment-inference model is transformed to:

$$\mathcal{L}_{elbo}(\tau, \kappa; \mathcal{G}) = \frac{1}{|\mathcal{G}|} \sum_{(G, y) \in \mathcal{G}} [\mathbb{E}_{q_{\kappa}}[\log p_{\tau}(y \mid G, e)] - D_{KL}(q_{\kappa}(e \mid G, y) \parallel p(e \mid G))]. \quad (6)$$

228

229 Let’s look back to the objective given in Eq. 5 and give an equivalent tractable objective in practical  
 230 instantiation, which involves the environment-inference model defined above:

$$\mathcal{L}_{inv}(\theta; \mathcal{G}, \tau) = \underbrace{\frac{1}{|\mathcal{G}|} \sum_{(G, y) \in \mathcal{G}} |\log q_{\theta}(y \mid G) - \mathbb{E}_{p(\mathbf{e} \mid \mathbf{G})}[\log p_{\tau}(y \mid G, e)]|}_{\textcircled{1}} + \beta \mathbb{E}_{\mathbf{e}} \left[ \underbrace{\frac{1}{|\mathcal{G}^{\mathbf{e}}|} \sum_{(G, y) \in \mathcal{G}^{\mathbf{e}}} [-\log q_{\theta}(y \mid G)]}_{\textcircled{2}} \right], \quad (7)$$

231 where  $\mathcal{G}^{\mathbf{e}}$  consists of the pairs of molecular graph  $G$  and corresponding label  $y$  under environment  $\mathbf{e}$ .

<sup>2</sup>Camelyon17 is for tumor prediction, partitioning 455,954 issue slides into 5 environments.

<sup>3</sup>CivilComments is for toxicity prediction, partitioning 448,000 online comments into 16 environments.

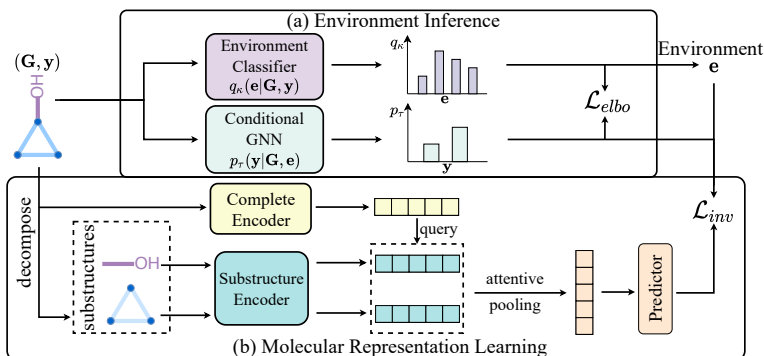


Figure 3: Overview of our model. The whole training procedure is divided into two stages: 1) Optimize the environment-inference model. Given an input molecule  $(\mathbf{G}, \mathbf{y})$ , we first infer the latent environment variable  $\mathbf{e}$ . This stage is trained under the guidance of  $\mathcal{L}_{elbo}$ . 2) Optimize the molecule encoder and the final predictor guided by  $\mathcal{L}_{inv}$ .

232 **Theorem 1.** With  $q_{\theta}(\mathbf{y}|\mathbf{z})$  treated as a variational distribution, minimizing term ① in Eq. 7  
 233 contributes to  $\min_{q_{\theta}(\mathbf{y}|\mathbf{z}), q_{\theta}(\mathbf{z}|\mathbf{G})} \mathbb{I}(\mathbf{y}; \mathbf{e}|\mathbf{z})$ , letting  $\mathbf{z}$  show equal performance for the downstream  
 234 tasks across all environments, i.e.  $p(\mathbf{y}|\mathbf{z}, \mathbf{e}) = p(\mathbf{y}|\mathbf{z})$ .

235 **Theorem 2.** Regarding  $q_{\theta}(\mathbf{y}|\mathbf{z})$  as a variational distribution, minimizing term ② in Eq. 7 equals  
 236 to  $\max_{q_{\theta}(\mathbf{y}|\mathbf{z}), q_{\theta}(\mathbf{z}|\mathbf{G})} \mathbb{I}(\mathbf{z}; \mathbf{y})$ , letting  $\mathbf{z}$  show sufficient predictive power for downstream tasks.

237 Serving as theoretical justifications, Th. 1 and Th. 2 reveal that optimizing the objective in Eq. 7 forces  
 238 the learned representation  $\mathbf{z}$  to satisfy the invariance principle mentioned in Sec. 2, thus ensuring a  
 239 valid solution for OOD problem defined in Eq. 2. Due to the limited space, the detailed proofs can be  
 240 found in Appendix B.

### 241 3.4 Model Instantiations and Training

242 **Environment-inference Module.** For the approximate posterior model  $q_{\kappa}(\mathbf{e}|\mathbf{G}, \mathbf{y})$ , in principle we  
 243 should design a module, entitled **Environment Classifier**, that takes  $(\mathbf{G}, \mathbf{y})$  as the input and outputs  
 244 the probabilistic distribution of  $\mathbf{e}$ . We use a Graph Isomorphism Network (GIN) [54], to learn a graph  
 245 representation given  $\mathbf{G}$ . Then, the concatenation of this graph representation and label vector is fed  
 246 to a feed-forward network to obtain a probabilistic distribution with regard to  $\mathbf{e}$ . We could set the  
 247 prior  $p_{\tau}(\mathbf{e}|\mathbf{G})$  to *Uniform distribution*, encouraging the learned environment-partition to be uniform.  
 248 As for  $p_{\tau}(\mathbf{y}|\mathbf{G}, \mathbf{e})$ , we also choose a GNN model followed by a softmax activation function to model  
 249 it. We call this module **Conditional GNN** because it conditions on  $\mathbf{e}$ . It takes  $(\mathbf{G}, \mathbf{e})$  as the input  
 250 and outputs the probabilistic distribution of  $\mathbf{y}$ .

251 **The Molecule Encoder & The Final Predictor.** Recall that we aim to learn an invariant substructure-  
 252 aware molecular representation. Given a molecule  $\mathbf{G}$ , we can choose any molecule representation  
 253 learning method to learn a representation  $\mathbf{r}_{\mathbf{G}}$  for the complete molecular graph. This part is en-  
 254 titled **Complete Encoder**. Meanwhile, we decompose the input molecule into a set of chemical  
 255 substructures using a molecule segmentation method, e.g. *breaking retrosynthetically interesting*  
 256 *chemical substructures* (BRICS) [12], which is available as an API in RDKit [31]. For each sub-  
 257 structure, we consider using a simple GNN to learn a corresponding representation. We call this  
 258 GNN **Substructure Encoder**. Then, considering  $\mathbf{r}_{\mathbf{G}}$  as a query with regard to substructures, we  
 259 operate attentive pooling on these substructure representations to obtain a new substructure-aware  
 260 molecular representation. We then use this substructure-aware representation for downstream task.  
 261 Guided by our proposed learning objective, we can encode some invariant relationships between  
 262 certain substructures and target properties into this representation. The Complete Encoder, the  
 263 Substructure encoder and the attentive pooling operation constitute our **Molecule Encoder**  $\Phi$ . As for  
 264 the **Predictor**  $\omega$ , we implement it with a multi-layer perceptron, followed by a softmax function. The  
 265 overview of our model is demonstrated in Fig. 3.

266 **Training.** We adopt a simple yet efficient two-stage training strategy to search for optimal parameters:

267 1) **optimizing the environment-inference model:**  $\kappa^*, \tau^* \leftarrow \arg \max_{\kappa, \tau} \mathcal{L}_{elbo}(\tau, \kappa; \mathcal{G}^{train})$ .

268 2) **optimizing the molecule encoder and the predictor:**  $\theta^* \leftarrow \arg \min_{\theta} \mathcal{L}_{inv}(\theta; \mathcal{G}^{train}, \tau)$ .

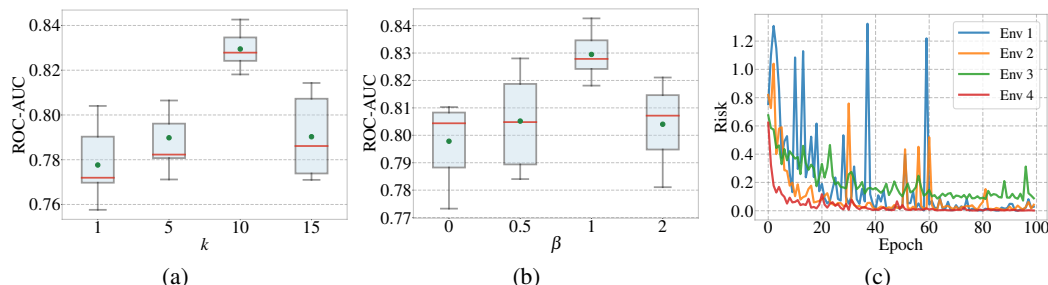


Figure 4: (a) Varying the specified environment number  $k$ . (b) Varying the trading-off parameter  $\beta$  in Eq. 7. (c) Risk curves of environments in the training process. All results are from ‘GraphSAGE + ours.’ on BACE dataset.

## 269 4 Experiments

270 Experiments are performed on 10 benchmark datasets. Each experiment is repeated 5 times with  
 271 mean and standard deviation reported, running on a machine with i9-10920X CPU, RTX 3090 GPU  
 272 and 128G RAM. **The source code for reproducing the results will be made publicly available.**

### 273 4.1 Datasets and Setups

274 **Datasets and protocols.** The four datasets **BACE**, **BBBP**, **SIDER** and **HIV**, are from by Open  
 275 Graph Benchmark (OGB) [20]. We use the default train/val/test split with ratio 8:1:1. Each split  
 276 contains a set of scaffolds (almost) different to each other. Hence we believe that to a certain  
 277 degree, it provides an OOD test-bed as different scaffold often suggest different data-generation  
 278 environments. The other six datasets are generated by the dataset curator provided by DrugOOD [24].  
 279 DrugOOD provides more diverse splitting indicators than OGB, including assay, scaffold and size.  
 280 To comprehensively evaluate the performance of our method under different environment definitions,  
 281 we adopt these three different splitting schemes on categories IC50 and EC50 provided in DrugOOD.  
 282 Then we obtain six datasets, **EC50-\*** and **IC50-\***, where the suffix  $*$  specifies the splitting scheme i.e.  
 283 **IC50/EC50-assay/scaffold/size**. Notice that only the six datasets from DrugOOD provide manual  
 284 specified environment labels. Refer to Appendix E for more details of datasets.

285 **Metric.** As the concerned property prediction tasks all relate to classification, we report the ROC-  
 286 AUC score which is also in line with previous MRL works [59, 55, 49].

287 **Baselines.** Ideally, any MRL method can be adapted into our method as backbone to improve their  
 288 generalization ability against distribution shifts. We adapt three backbones: **GCN** [27], **GIN** [54]  
 289 and **GraphSAGE** [18] into our method. We compare the adapted version with the original method.  
 290 We also compare against another augmented version “+ virtual node” [16, 22, 35]. Furthermore,  
 291 we compare our method with six OOD generalization methods on MRL tasks: **ERM** [47], **IRM** [2],  
 292 **DeepCoral** [46], **DANN** [15], **MixUp** [61] and **GroupDro** [43]. Due to the fact that most of these  
 293 methods require the manual specification of environments in dataset, we report this comparison on  
 294 datasets from DrugOOD only. Each of the method is configured using the same parameters reported  
 295 in the original paper or selected by grid search. For the sake of fairness, the embedding size of all  
 296 methods are set to be equal in comparison. We specify the training details in the Appendix C.

### 297 4.2 Performance Comparison

298 **Improvements to existing MRL methods.** As demonstrated in Table 1, baselines obtain consistent  
 299 improvements after adapted to our methods across all the four datasets released by OGB in terms  
 300 of ROC-AUC. Our method also beats the augmented version, “+ virtual node”, of baselines on  
 301 all datasets, i.e. adding a virtual node. The results indicate that, orthogonal to prior studies on  
 302 MRL, our method is a general framework which can incorporate existing MRL methods and improve  
 303 their generalization ability for OOD data. We attribute the superior performance of our method in  
 304 molecular properties predictions under OOD setting to that, our proposed learning objective enforces  
 305 the model to learn environment-invariant representations against distribution shifts.

306 **Superiority to other OOD generalization methods.** Table 2 summarizes the results in comparison  
 307 with six state-of-the-art methods tailored for OOD learning, where we obtain the following observa-  
 308 tions. Except on IC50-size, our method outperforms all baselines across all datasets due to its ability

Table 1: Performance comparison with baselines on 4 out-of-distribution molecular property prediction datasets from Open Graph Benchmark (OGB) [20] in terms of ROC-AUC (%), namely, BACE, BBBP, SIDER and HIV. The best and the runner-up results are highlighted in **bolded** and underlined respectively. We emphasize the comparison against ‘\* + **virtual node**’, a variant of the original method augmented by an additional node connecting to all nodes in the raw graphs [16, 22, 35].

Methods	BACE	BBBP	SIDER	HIV
GCN [27]	80.01 ± 3.49	67.92 ± 1.07	58.90 ± 1.30	76.35 ± 2.01
GCN + <b>virtual node</b>	<u>77.51 ± 3.07</u>	<u>68.19 ± 1.86</u>	<u>60.71 ± 1.34</u>	75.76 ± 2.21
GCN + <b>ours.</b>	<b>84.33 ± 1.07</b>	<b>70.62 ± 0.99</b>	<b>63.38 ± 0.67</b>	<b>77.73 ± 0.76</b>
GIN [54]	77.83 ± 3.15	66.93 ± 2.31	59.05 ± 1.47	76.58 ± 1.02
GIN + <b>virtual node</b>	<u>79.64 ± 2.02</u>	<u>66.77 ± 0.95</u>	<u>59.12 ± 0.95</u>	77.11 ± 0.96
GIN + <b>ours.</b>	<b>81.09 ± 2.03</b>	<b>69.84 ± 1.84</b>	<b>61.63 ± 1.08</b>	<b>78.31 ± 0.24</b>
GraphSAGE [18]	77.41 ± 1.19	70.58 ± 0.58	58.00 ± 0.95	76.98 ± 1.13
GraphSAGE + <b>virtual node</b>	<u>78.34 ± 2.08</u>	<u>69.29 ± 0.99</u>	<u>59.48 ± 1.37</u>	77.28 ± 1.53
GraphSAGE + <b>ours.</b>	<b>82.95 ± 0.85</b>	<b>71.02 ± 0.75</b>	<b>61.09 ± 0.28</b>	<b>79.39 ± 0.51</b>

Table 2: Evaluation with other OOD generalization methods on 6 out-of-distribution datasets from DrugOOD [24] in terms of ROC-AUC (%). The best and the runner-up in each columns are highlighted in **bolded** and underlined respectively. Note the baselines except ERM and MixUp all require environment labels. All methods including ours use GIN [54] as backbones.

Dataset	IC50			EC50		
	Assay	Scaffold	Size	Assay	Scaffold	Size
ERM [47]	70.93 ± 2.10	67.31 ± 1.72	67.40 ± 0.56	69.35 ± 7.38	63.92 ± 2.09	60.94 ± 1.95
IRM [2]	<u>70.85 ± 2.41</u>	<u>66.06 ± 1.23</u>	58.46 ± 2.11	69.94 ± 1.03	63.74 ± 2.15	58.30 ± 1.51
DeepCoral [46]	69.82 ± 4.23	66.36 ± 2.57	59.21 ± 2.09	69.42 ± 3.35	63.66 ± 1.87	56.13 ± 1.77
DANN [15]	70.00 ± 1.03	63.61 ± 2.32	65.77 ± 0.47	66.97 ± 7.19	64.33 ± 1.82	61.11 ± 0.64
MixUp [61]	70.22 ± 3.66	66.43 ± 1.08	<b>67.77 ± 0.23</b>	70.62 ± 2.12	64.53 ± 1.66	62.67 ± 1.41
GroupDro [43]	69.98 ± 1.74	64.09 ± 2.05	58.46 ± 2.69	70.52 ± 3.38	64.13 ± 1.81	59.06 ± 1.50
<b>Ours.</b>	<b>71.38 ± 0.68</b>	<b>68.02 ± 0.55</b>	66.51 ± 0.55	<b>73.25 ± 1.24</b>	<b>66.69 ± 0.34</b>	<b>65.09 ± 0.90</b>

309 to enforce the molecule encoder to leverage environment-invariant substructures that more stably  
310 relate with the labels across environments. Our method ranks the third on IC50-size after MixUp and  
311 ERM. Different from the other methods, MixUp constructs more training examples and uses more  
312 data to train the model. That’s why MixUp obtains best performance among all methods on IC50-size  
313 in our analysis. As for ERM, [24, 14] have pointed out that simple ERM shows better performance  
314 compared to subsequent OOD methods when datasets have relatively large environment counts. Even  
315 though we have set the environment number  $k$  to a smaller value than the ground-truth number given  
316 by the dataset, we still need to prevent  $k$  from being too small (see discussion in Sec. 4.4), leading to  
317 our poorer performance than ERM on IC50-size.

### 318 4.3 Ablation Study

319 We analyze the contributions of different model components to the final performance in this section.  
320 Table 3 reports detailed ablation experimental results on EC50-assay, EC50-scaffold and EC50-size.

321 **Attention-based architecture.** We study the impact of the attention-based architecture introduced  
322 in Sec. 3.4 by assembling this architecture with ERM loss. We beat ERM and MixUp only with this  
323 architecture on three datasets. The results show that learning a representation for each substructure  
324 and then attentively aggregating these learned representations to obtain a final substructure-aware  
325 representation performs better than learning a representation for a complete molecular graph directly.  
326 This verifies our assumption that the substructure perspective is of importance to boosting performance  
327 of existing MRL methods. With the aid of such a substructure-grained learning architecture, the  
328 impact of our learning objective can be further strengthened.

329 **New learning objective.** To evaluate the impact of our proposed new learning objective, we equip  
330 GIN with this new objective. We can see compared to using the substructure-grained learning archi-  
331 tecture only, only using the proposed new learning objective can bring more significant improvement.  
332 Thus, we can attribute the main superiority of our full model to this new objective. Combined with  
333 the architecture discussed above, the new objective is able to better guide the molecule encoder to  
334 learn environment-invariant molecular representations against distribution shifts.

Table 3: Ablation study on EC50-\* by ROC-AUC (%). We show the results of MixUp that performs best among baselines on all EC50-\* datasets and the naive ERM, which minimizes the average empirical loss on training data, for comparison. Notice that ERM and MixUp don’t require manual specified environments labels. We also present the results of DANN, which requires manual specifications of environment and obtains competitive results with MixUp. All methods use GIN [54] as backbone.

Method	Assay	Scaffold	Size
ERM (GIN + ERM loss)	69.35 ± 7.38	63.92 ± 2.09	60.94 ± 1.95
MixUp	70.62 ± 2.12	64.53 ± 1.66	62.67 ± 1.41
DANN	66.97 ± 7.19	64.33 ± 1.82	61.11 ± 0.64
Our architecture + ERM loss	71.44 ± 2.02	65.99 ± 0.42	64.23 ± 0.71
GIN + new learning objective	72.07 ± 1.14	66.33 ± 1.38	64.43 ± 1.10
DANN using our inferred environment label	68.83 ± 2.44	64.95 ± 1.07	62.56 ± 1.54
Our model using given environment label	71.94 ± 2.77	66.29 ± 0.85	63.38 ± 1.20
<b>Our full model</b>	<b>73.25 ± 1.24</b>	<b>66.69 ± 0.34</b>	<b>65.09 ± 0.90</b>

335 **Environment inference.** Now we turn to investigate the performance with respect to our proposed  
336 environment-inference module. One motivation for this module is that in reality manual specifications  
337 of environments may be unavailable due to the high price for labeling environments by experts. But  
338 when environment labels are available, how will be performance be like if directly utilizing the given  
339 environment partition? An ablation study is targeted on this. Taking the EC50-assay dataset as an  
340 example, it has given the environment partition and it specifies 47 environments in total. We utilize  
341 the given environment partition directly and keep the remaining parts in line with our full model.  
342 The results show that utilizing the given environment label, our method still can beat ERM and  
343 MixUp. But compared to our full model where we set the environment number  $k$  to 20, it obtains  
344 inferior performance. Additionally, to further examine the effectiveness of our proposed environment  
345 inference method, we relabel the environment for each molecule for DANN according to our inferred  
346 environment partition. We can see that based on the new environment partition, DANN obtains better  
347 performance than using the initial given environment labels across three datasets. The reason why  
348 inferring environment instead can outperform directly using the given environment label is mainly  
349 due to the existing given partitions are often handcrafted-rule-based and not structured. In contrast,  
350 letting the model learn a environment partition by itself may be more effective to some degree.

#### 351 4.4 Hyper-parameter Sensitivity & Risk Dynamics

352 We investigate the sensitivity of our method to these two hyper-parameters: the specified number  
353 of environments  $k$ , the trading-off parameter  $\beta$  in Eq. 7. Fig. 4(a) shows the performance regarding  
354 different environment number  $k$ . It shows that the performance of our methods degrades when  $k$  is  
355 too small (e.g.  $k = 1, 5$ ) or too large (e.g.  $k = 15$ ). When  $k = 1$  i.e. we regard all training data  
356 as from only one environment, the performance is the poorest. This justifies that partitioning the  
357 training samples into different environments is necessary. Fig. 4(b) shows the results of our method  
358 by varying the trade-off parameter  $\beta$ . Our method obtains the worst performance when  $\beta = 0$ . This  
359 is mainly because Eq. 7 is reduced to the first term when  $\beta = 0$ . According to Theorem 2, without the  
360 second term of Eq. 7, the sufficiency condition of invariance principle cannot be satisfied, resulting  
361 in the performance degradation. Additionally, to shed insights of the ability of our method to lower  
362 the risks of different environments, we visualize the risk dynamic curve of some environments in  
363 Fig. 4(c). As is shown in Fig. 4(c), the difficulties of decreasing the risk on different environments  
364 are different. Though the risks of some environments vibrate violently at the beginning of training  
365 process (e.g. Env 1 and Env 2), with time elapsing, risks on all environments can decrease stably.

## 366 5 Conclusion

367 We have proposed a general framework which can incorporate any existing MRL method as backbone  
368 to improve their generalization ability against distribution shifts. Specifically, we devise a new  
369 learning scheme with its equivalent practical instantiation. We also develop an environment inference  
370 model to identify each molecule’s corresponding environment without need of manual specifications  
371 of environments. Extensive experimental results on ten datasets demonstrate that our model yields  
372 consistent and significant improvements over various existing MRL methods as backbones. Addi-  
373 tionally, our model achieves competitive or even superior performance compared to state-of-the-art  
374 models designed for OOD learning that require manual specified environment labels as extra inputs.

375 **References**

- 376 [1] E. Anderson, G. D. Veith, and D. Weininger. *SMILES, a line notation and computerized*  
377 *interpreter for chemical structures*. US Environmental Protection Agency, Environmental  
378 Research Laboratory, 1987.
- 379 [2] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. Invariant risk minimization. *arXiv*  
380 *preprint arXiv:1907.02893*, 2019.
- 381 [3] P. Bandi, O. Geessink, Q. Manson, M. Van Dijk, M. Balkenhol, M. Hermsen, B. E. Bejnordi,  
382 B. Lee, K. Paeng, A. Zhong, et al. From detection of individual metastases to classification of  
383 lymph node status at the patient level: the camelyon17 challenge. *IEEE transactions on medical*  
384 *imaging*, 38(2):550–560, 2018.
- 385 [4] G. W. Bemis and M. A. Murcko. The properties of known drugs. 1. molecular frameworks.  
386 *Journal of medicinal chemistry*, 39(15):2887–2893, 1996.
- 387 [5] G. Blanchard, G. Lee, and C. Scott. Generalizing from several related classification tasks to  
388 a new unlabeled sample. In *Advances in Neural Information Processing Systems (NeurIPS)*,  
389 pages 2178–2186, 2011.
- 390 [6] D. Borkan, L. Dixon, J. Sorensen, N. Thain, and L. Vasserman. Nuanced metrics for measuring  
391 unintended bias with real data for text classification. In *Companion proceedings of the 2019*  
392 *world wide web conference*, pages 491–500, 2019.
- 393 [7] P. Bühlmann. Invariance, causality and robustness. *Statistical Science*, 35(3):404–426, 2020.
- 394 [8] S. Chang, Y. Zhang, M. Yu, and T. Jaakkola. Invariant rationalization. In *International*  
395 *Conference on Machine Learning*, pages 1448–1458. PMLR, 2020.
- 396 [9] T. Cofala and O. Kramer. An evolutionary fragment-based approach to molecular fingerprint  
397 reconstruction. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages  
398 1156–1163, 2022.
- 399 [10] E. Creager, J.-H. Jacobsen, and R. Zemel. Environment inference for invariant learning. In  
400 *International Conference on Machine Learning*, pages 2189–2200. PMLR, 2021.
- 401 [11] L. David, A. Thakkar, R. Mercado, and O. Engkvist. Molecular representations in ai-driven  
402 drug discovery: a review and practical guide. *Journal of Cheminformatics*, 12(1):1–22, 2020.
- 403 [12] J. Degen, C. Wegscheid-Gerlach, A. Zaliani, and M. Rarey. On the art of compiling and  
404 using ‘drug-like’ chemical fragment spaces. *ChemMedChem*, 2008.
- 405 [13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional  
406 transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 407 [14] M. Ding, K. Kong, J. Chen, J. Kirchenbauer, M. Goldblum, D. Wipf, F. Huang, and T. Goldstein.  
408 A closer look at distribution shifts and out-of-distribution generalization on graphs, 2022.
- 409 [15] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and  
410 V. Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning*  
411 *research*, 17(1):2096–2030, 2016.
- 412 [16] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing  
413 for quantum chemistry. In *International conference on machine learning*, pages 1263–1272.  
414 PMLR, 2017.
- 415 [17] M. Glymour, J. Pearl, and N. P. Jewell. *Causal inference in statistics: A primer*. John Wiley &  
416 Sons, 2016.
- 417 [18] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs.  
418 *Advances in neural information processing systems*, 30, 2017.
- 419 [19] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *Journal*  
420 *of Machine Learning Research*, 2013.

- 421 [20] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec. Open  
422 graph benchmark: Datasets for machine learning on graphs. *Advances in neural information*  
423 *processing systems*, 33:22118–22133, 2020.
- 424 [21] S. Ishida, T. Miyazaki, Y. Sugaya, and S. Omachi. Graph neural networks with multiple feature  
425 extraction paths for chemical property estimation. *Molecules*, 26(11):3125, 2021.
- 426 [22] K. Ishiguro, S.-i. Maeda, and M. Koyama. Graph warp module: an auxiliary module for  
427 boosting the power of graph neural networks. *arXiv preprint arXiv:1902.01020*, 2019.
- 428 [23] S. Jaeger, S. Fulle, and S. Turk. Mol2vec: unsupervised machine learning approach with  
429 chemical intuition. *Journal of chemical information and modeling*, 2018.
- 430 [24] Y. Ji, L. Zhang, J. Wu, B. Wu, L.-K. Huang, T. Xu, Y. Rong, L. Li, J. Ren, D. Xue, et al. Drugood:  
431 Out-of-distribution (ood) dataset curator and benchmark for ai-aided drug discovery—a focus on  
432 affinity prediction problems with noise annotations. *arXiv preprint arXiv:2201.09637*, 2022.
- 433 [25] W. Jin, C. Coley, R. Barzilay, and T. Jaakkola. Predicting organic reaction outcomes with  
434 weisfeiler-lehman network. *Advances in neural information processing systems*, 30, 2017.
- 435 [26] Y.-T. Kao, S.-F. Wang, M.-H. Wu, S.-H. Her, Y.-H. Yang, C.-H. Lee, H.-F. Lee, A.-R. Lee, L.-C.  
436 Chang, and L.-H. Pao. A substructure-based screening approach to uncover n-nitrosamines in  
437 drug substances. *Journal of Food & Drug Analysis*, 30(1), 2022.
- 438 [27] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks.  
439 In *5th International Conference on Learning Representations, ICLR 2017*, 2017.
- 440 [28] J. Klekota and F. P. Roth. Chemical substructures that enrich for biological activity. *Bioinfor-*  
441 *matics*, 24(21):2518–2525, 2008.
- 442 [29] P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Ya-  
443 sunaga, R. L. Phillips, I. Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In  
444 *International Conference on Machine Learning*, pages 5637–5664. PMLR, 2021.
- 445 [30] D. Krueger, E. Caballero, J.-H. Jacobsen, A. Zhang, J. Binas, D. Zhang, R. Le Priol, and  
446 A. Courville. Out-of-distribution generalization via risk extrapolation (rex). In *International*  
447 *Conference on Machine Learning*, pages 5815–5826. PMLR, 2021.
- 448 [31] G. Landrum et al. Rdkit: Open-source cheminformatics. 2006.
- 449 [32] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh. Set transformer: A framework for  
450 attention-based permutation-invariant neural networks. In *International conference on machine*  
451 *learning*, pages 3744–3753. PMLR, 2019.
- 452 [33] A. Leman and B. Weisfeiler. A reduction of a graph to a canonical form and an algebra arising  
453 during this reduction. *Nauchno-Technicheskaya Informatsiya*, 2(9):12–16, 1968.
- 454 [34] X. Q. Lewell, D. B. Judd, S. P. Watson, and M. M. Hann. Recap retrosynthetic combinatorial  
455 analysis procedure: a powerful new technique for identifying privileged molecular fragments  
456 with useful applications in combinatorial chemistry. *Journal of chemical information and*  
457 *computer sciences*, 38(3):511–522, 1998.
- 458 [35] J. Li, D. Cai, and X. He. Learning graph-level representation for drug discovery. *arXiv preprint*  
459 *arXiv:1709.03741*, 2017.
- 460 [36] J. Lim, S. Ryu, K. Park, Y. J. Choe, J. Ham, and W. Y. Kim. Predicting drug–target interaction  
461 using a novel graph neural network with 3d structure-embedded graph representation. *Journal*  
462 *of chemical information and modeling*, 59(9):3981–3988, 2019.
- 463 [37] Y. Mansour, M. Mohri, and A. Rostamizadeh. Domain adaptation: Learning bounds and  
464 algorithms. In *Conference on Learning Theory (COLT)*, 2009.
- 465 [38] H. L. Morgan. The generation of a unique machine description for chemical structures—a  
466 technique developed at chemical abstracts service. *Journal of Chemical Documentation*, 1965.

- 467 [39] K. Muandet, D. Balduzzi, and B. Schölkopf. Domain generalization via invariant feature  
468 representation. In *International Conference on Machine Learning (ICML)*, pages 10–18, 2013.
- 469 [40] J. Pearl et al. Models, reasoning and inference. *Cambridge, UK: Cambridge University Press*,  
470 19:2, 2000.
- 471 [41] C. Phanus-umporn, W. Shoombuatong, V. Prachayasittikul, N. Anuwongcharoen, and C. Nan-  
472 tasanamat. Privileged substructures for anti-sickling activity via cheminformatic analysis. *RSC*  
473 *advances*, 2018.
- 474 [42] M. Rojas-Carulla, B. Schölkopf, R. Turner, and J. Peters. Invariant models for causal transfer  
475 learning. *The Journal of Machine Learning Research*, 19(1):1309–1342, 2018.
- 476 [43] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Distributionally robust neural networks  
477 for group shifts: On the importance of regularization for worst-case generalization. *arXiv*  
478 *preprint arXiv:1911.08731*, 2019.
- 479 [44] B. Sanchez-Lengeling, C. Outeiral, G. L. Guimaraes, and A. Aspuru-Guzik. Optimizing  
480 distributions over molecular space. an objective-reinforced generative adversarial network for  
481 inverse-design chemistry (organic). 2017.
- 482 [45] J. M. Stokes, K. Yang, K. Swanson, W. Jin, A. Cubillos-Ruiz, N. M. Donghia, C. R. MacNair,  
483 S. French, L. A. Carfrae, Z. Bloom-Ackermann, et al. A deep learning approach to antibiotic  
484 discovery. *Cell*, 180(4):688–702, 2020.
- 485 [46] B. Sun and K. Saenko. Deep coral: Correlation alignment for deep domain adaptation. In  
486 *European conference on computer vision*, pages 443–450. Springer, 2016.
- 487 [47] V. Vapnik. *The nature of statistical learning theory*. Springer science & business media, 1999.
- 488 [48] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and  
489 I. Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 490 [49] H. Wang, W. Li, X. Jin, K. Cho, H. Ji, J. Han, and M. Burke. Chemical-reaction-aware molecule  
491 representation learning. In *International Conference on Learning Representations*, 2022.
- 492 [50] S. Wang, Z. Li, S. Zhang, M. Jiang, X. Wang, and Z. Wei. Molecular property prediction based  
493 on a multichannel substructure graph. *IEEE Access*, 8:18601–18614, 2020.
- 494 [51] Y. Wang, R. Magar, C. Liang, and A. Barati Farimani. Improving molecular contrastive  
495 learning via faulty negative mitigation and decomposed fragment contrast. *Journal of Chemical*  
496 *Information and Modeling*, 2022.
- 497 [52] Q. Wu, H. Zhang, J. Yan, and D. Wipf. Handling distribution shifts on graphs: An invariance  
498 perspective. In *International Conference on Learning Representations*, 2022.
- 499 [53] Y. Wu, X. Wang, A. Zhang, X. He, and T.-S. Chua. Discovering invariant rationales for graph  
500 neural networks. In *International Conference on Learning Representations*, 2022.
- 501 [54] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? In *7th*  
502 *International Conference on Learning Representations, ICLR 2019*, 2019.
- 503 [55] M. Xu, H. Wang, B. Ni, H. Guo, and J. Tang. Self-supervised graph-level representation  
504 learning with local and global structure. In *International Conference on Machine Learning*,  
505 pages 11548–11558. PMLR, 2021.
- 506 [56] C. Yan, Q. Ding, P. Zhao, S. Zheng, J. Yang, Y. Yu, and J. Huang. Retroxpert: Decompose  
507 retrosynthesis prediction like a chemist. *Advances in Neural Information Processing Systems*,  
508 33:11248–11258, 2020.
- 509 [57] C. Yang, C. Xiao, F. Ma, L. Glass, and J. Sun. Safedrug: Dual molecular graph encoders for  
510 recommending effective and safe drug combinations. In *IJCAI*, pages 3735–3741, 2021.

- 511 [58] A. B. Yongye, J. Waddell, and J. L. Medina-Franco. Molecular scaffold analysis of natural  
512 products databases in the public domain. *Chemical biology & drug design*, 80(5):717–724,  
513 2012.
- 514 [59] Y. You, T. Chen, Y. Sui, T. Chen, Z. Wang, and Y. Shen. Graph contrastive learning with  
515 augmentations. *Advances in Neural Information Processing Systems*, 33:5812–5823, 2020.
- 516 [60] X. Zeng, S. Zhu, W. Lu, Z. Liu, J. Huang, Y. Zhou, J. Fang, Y. Huang, H. Guo, L. Li, et al. Target  
517 identification among known drugs by deep learning from heterogeneous networks. *Chemical  
518 Science*, 11(7):1775–1797, 2020.
- 519 [61] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk mini-  
520 mization. *arXiv preprint arXiv:1710.09412*, 2017.
- 521 [62] X. Zhao, B. Zong, Z. Guan, K. Zhang, and W. Zhao. Substructure assembling network for graph  
522 classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32,  
523 2018.
- 524 [63] J. Zhu, Y. Liu, C. Wen, and X. Wu. Dgdfs: Dependence guided discriminative feature selection  
525 for predicting adverse drug-drug interaction. *IEEE Transactions on Knowledge and Data  
526 Engineering*, 2020.

527 **Checklist**

- 528 1. For all authors...
- 529 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
530 contributions and scope? [Yes]
- 531 (b) Did you describe the limitations of your work? [Yes] See Appendix I.
- 532 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
533 Appendix J.
- 534 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
535 them? [Yes]
- 536 2. If you are including theoretical results...
- 537 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 538 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A  
539 and B.
- 540 3. If you ran experiments...
- 541 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
542 mental results (either in the supplemental material or as a URL)? [Yes] See Appendix C.
- 543 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
544 were chosen)? [Yes] See Section 4 and Appendix C.
- 545 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
546 ments multiple times)? [Yes] See Section 4.
- 547 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
548 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4.
- 549 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 550 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 4 and  
551 Appendix E.
- 552 (b) Did you mention the license of the assets? [Yes] See Section 4 and Appendix E.
- 553 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 554 (d) Did you discuss whether and how consent was obtained from people whose data you're  
555 using/curating? [Yes] See Appendix E.
- 556 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
557 information or offensive content? [Yes] See Appendix E.
- 558 5. If you used crowdsourcing or conducted research with human subjects...
- 559 (a) Did you include the full text of instructions given to participants and screenshots, if  
560 applicable? [N/A]
- 561 (b) Did you describe any potential participant risks, with links to Institutional Review  
562 Board (IRB) approvals, if applicable? [N/A]
- 563 (c) Did you include the estimated hourly wage paid to participants and the total amount  
564 spent on participant compensation? [N/A]

565 **A Proof for Proposition 1**

566 *Proof.* Our goal is to minimize the Kullback-Leibler (KL) divergence between  $q_\kappa(\mathbf{e}|\mathbf{G}, \mathbf{y})$  and  
 567  $p_\tau(\mathbf{e}|\mathbf{G}, \mathbf{y})$ . For the observed molecule graph and corresponding label tuple  $(G, y)$ ,

$$\begin{aligned}
 & D_{KL}(q_\kappa(\mathbf{e}|G, y) \parallel p_\tau(\mathbf{e}|G, y)) \\
 &= \int_{q_\kappa} q_\kappa(\mathbf{e}|G, y) \log \frac{q_\kappa(\mathbf{e}|G, y)}{p_\tau(\mathbf{e}|G, y)} d\mathbf{e} = \int_{q_\kappa} q_\kappa(\mathbf{e}|G, y) \log \frac{q_\kappa(\mathbf{e}|G, y)p_\tau(y|G)p_\tau(G)}{p_\tau(\mathbf{e}, G, y)} d\mathbf{e} \\
 &= \left( \int_{q_\kappa} q_\kappa(\mathbf{e}|G, y) \log q_\kappa(\mathbf{e}|G, y) d\mathbf{e} + \int_{q_\kappa} p_\kappa(\mathbf{e}|G, y) \log p_\tau(G) d\mathbf{e} \right. \\
 &\quad \left. - \int_{q_\kappa} \log p_\tau(\mathbf{e}, G, y) d\mathbf{e} \right) + \int_{q_\kappa} q_\kappa(\mathbf{e}|G, y) \log p_\tau(y|G) d\mathbf{e} \\
 &= \int_{q_\kappa} q_\kappa(\mathbf{e}|G, y) \log \frac{q_\kappa(\mathbf{e}|G, y)}{p_\tau(y|G, \mathbf{e})p_\tau(\mathbf{e}|G)} d\mathbf{e} + \log p_\tau(y|G) \\
 &= \mathbb{E}_{q_\kappa} [\log q_\kappa(\mathbf{e}|G, y) - \log p_\tau(y|G, \mathbf{e}) - \log p_\tau(\mathbf{e}|G)] + \log p_\tau(y|G) \\
 &= -\mathbb{E}_{q_\kappa} [\log p_\tau(y|G, \mathbf{e})] + \mathbb{E}_{q_\kappa} [\log q_\kappa(\mathbf{e}|G, y) - \log p_\tau(\mathbf{e}|G)] + \log p_\tau(y|G) \\
 &= -\underbrace{(\mathbb{E}_{q_\kappa} [\log p_\tau(y|G, \mathbf{e})] - D_{KL}(q_\kappa(\mathbf{e}|G, y) \parallel p_\tau(\mathbf{e}|G)))}_{\mathcal{L}(\tau, \kappa; G, y)} + \log p_\tau(y|G) \\
 &= -\mathcal{L}(\tau, \kappa; G, y) + \log p_\tau(y|G)
 \end{aligned} \tag{8}$$

568 Rearrange Eq. 8 and we can get,

$$\mathcal{L}(\tau, \kappa; G, y) = -D_{KL}(q_\kappa(\mathbf{e}|G, y) \parallel p_\tau(\mathbf{e}|G, y)) + \log p_\tau(y|G). \tag{9}$$

569 The defined  $\mathcal{L}(\tau, \kappa; G, y)$  is called *Evidence Lower BOund* (ELBO) [19]. According to Eq. 9,  
 570 maximizing this ELBO is equivalent to minimizing the KL divergence and maximizing  $\log p_\tau(y|G)$ .  
 571 For the observed molecule graph and corresponding label tuple  $(G, y)$ , we obtain the ELBO:

$$\mathcal{L}(\tau, \kappa; G, y) = \mathbb{E}_{q_\kappa} [\log p_\tau(y|G, \mathbf{e})] - D_{KL}(q_\kappa(\mathbf{e}|G, y) \parallel p(\mathbf{e}|G)). \tag{10}$$

572 We thus conclude the proof.  $\square$

573 **B Proofs for Theorems**

574 In this paper, we extend the invariance assumption [42, 2] to molecule representation learning:

575 **Assumption 1.** *Given a molecular graph  $G$ , there exists an encoder  $\Phi$  yielding a graph-level*  
 576 *representation  $r_G \in \mathbb{R}^d$ . Define  $\mathbf{r}$  as a random variable of  $r_G$  and it satisfies: 1) (Invariance*  
 577 *condition):  $p(\mathbf{y}|\mathbf{r}, \mathbf{e}) = p(\mathbf{y}|\mathbf{r})$ , and 2) (Sufficiency condition):  $\mathbf{y} = h(\mathbf{r}) + \mathbf{n}$ , where  $h$  is a non-*  
 578 *linear function,  $\mathbf{n}$  is a independent noise.*

579 With the terminology of information theory, we present a useful lemma [52] that interprets the  
 580 invariance and sufficiency conditions in Assumption 1:

581 **Lemma 1.** *In terms of information theory, the two conditions in Assumption 1 can be equivalently*  
 582 *expressed as 1) invariance:  $I(\mathbf{y}; \mathbf{e}|\mathbf{r}) = 0$  and 2) sufficiency:  $I(\mathbf{y}; \mathbf{r})$  is maximized.*

583 *Proof.* For the invariance, it can be obtained by the fact that

$$I(\mathbf{y}; \mathbf{e}|\mathbf{r}) = \mathbb{E}_{p(\mathbf{e}, \mathbf{r})} [D_{KL}(p(\mathbf{y}|\mathbf{e}, \mathbf{r}) \parallel p(\mathbf{y}|\mathbf{r}))] \tag{11}$$

584 For the sufficiency, we first prove that every triplet  $(\mathbf{G}, \mathbf{r}, \mathbf{y})$  satisfying that  $\mathbf{y} = h(\mathbf{r}) + \mathbf{n}$  would also  
 585 satisfy  $\mathbf{r} = \arg \max_{\mathbf{r}} I(\mathbf{y}; \mathbf{r})$ . We prove it by contradiction. Assume that  $\mathbf{r} \neq \arg \max_{\mathbf{r}} I(\mathbf{y}; \mathbf{r})$  and  
 586 there exists  $\mathbf{r}'$  with  $\mathbf{r}' = \arg \max_{\mathbf{r}} I(\mathbf{y}; \mathbf{r})$  with  $\mathbf{r} \neq \mathbf{r}'$ . Then there exists another random variable  
 587  $\tilde{\mathbf{r}}$  and a mapping function  $f_m$  such that  $\mathbf{r}' = f_m(\mathbf{r}, \tilde{\mathbf{r}})$ . Then we will have  $I(\mathbf{y}; \mathbf{r}') = I(\mathbf{y}; \mathbf{r}, \tilde{\mathbf{r}}) =$   
 588  $I(h(\mathbf{r}); \mathbf{r}, \tilde{\mathbf{r}}) = I(h(\mathbf{r}); \mathbf{r}) = I(\mathbf{y}; \mathbf{r})$ , which leads to contradiction.

589 Then we prove that every triplet  $(\mathbf{G}, \mathbf{r}, \mathbf{y})$  satisfying that  $\mathbf{r} = \arg \max_{\mathbf{r}} I(\mathbf{y}; \mathbf{r})$  would also satisfy  
 590  $\mathbf{y} = h(\mathbf{r}) + \mathbf{n}$  by contradiction. Suppose that  $\mathbf{y} \neq h(\mathbf{r}) + \mathbf{n}$  and there exists  $\mathbf{r}' \neq \mathbf{r}$  with  $\mathbf{y} = h(\mathbf{r}') + \mathbf{n}$ .  
 591 Then the inequality  $I(h(\mathbf{r}'); \mathbf{r}) < I(h(\mathbf{r}'); \mathbf{r}')$  holds. That means  $\mathbf{r}' = \arg \max_{\mathbf{r}} I(\mathbf{y}; \mathbf{r})$ , leading to  
 592 contradiction.  $\square$

593 **B.1 Proof for Theorem 1**

594 *Proof.* According to the dependency relationship  $\mathbf{z} \leftarrow \mathbf{G} \rightarrow \mathbf{y}$ , we have

$$\begin{aligned}
& I(\mathbf{y}; \mathbf{e}|\mathbf{z}) \\
&= D_{KL}(p(\mathbf{y}|\mathbf{z}, \mathbf{e}) \parallel p(\mathbf{y}|\mathbf{z})) \\
&= D_{KL}(p(\mathbf{y}|\mathbf{z}, \mathbf{e}) \parallel \mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y}|\mathbf{z}, \mathbf{e})]) \\
&= D_{KL}(q(\mathbf{y}|\mathbf{z}) \parallel \mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y}|\mathbf{G}, \mathbf{e})]) - D_{KL}(q(\mathbf{y}|\mathbf{z}) \parallel p(\mathbf{y}|\mathbf{z}, \mathbf{e})) \\
&\quad - D_{KL}(\mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y}|\mathbf{z}, \mathbf{e})] \parallel \mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y}|\mathbf{G}, \mathbf{e})]) \\
&\leq D_{KL}(q(\mathbf{y}|\mathbf{z}) \parallel \mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y}|\mathbf{G}, \mathbf{e})]).
\end{aligned} \tag{12}$$

595 Next, we have

$$\begin{aligned}
& D_{KL}(q(\mathbf{y}|\mathbf{z}) \parallel \mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y}|\mathbf{G}, \mathbf{e})]) \\
&= \mathbb{E}_{G \sim p(\mathbf{G})} \mathbb{E}_{y_G \sim p(\mathbf{y}|\mathbf{G}=G)} \mathbb{E}_{z_G \sim q(\mathbf{z}|\mathbf{G}=G)} \left[ \log \frac{q(\mathbf{y} = y_G | \mathbf{z} = z_G)}{\mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y} = y_G | \mathbf{G} = G, \mathbf{e} = e)]} \right] \\
&= \frac{1}{|\mathcal{G}|} \sum_{(G, y_G) \in \mathcal{G}} \mathbb{E}_{z_G \sim q(\mathbf{z}|\mathbf{G}=G)} \left[ \log \frac{q(\mathbf{y} = y_G | \mathbf{z} = z_G)}{\mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y} = y_G | \mathbf{G} = G, \mathbf{e} = e)]} \right].
\end{aligned} \tag{13}$$

596 Based on Jensen Inequality and Triangle Inequality, we can obtain that  
597  $D_{KL}(q(\mathbf{y}|\mathbf{z}) \parallel \mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [p(\mathbf{y}|\mathbf{G}, \mathbf{e})])$  is upper bounded by:

$$\frac{1}{|\mathcal{G}|} \sum_{(G, y) \in \mathcal{G}} |\log q_\theta(y|G) - \mathbb{E}_{p(\mathbf{e}|\mathbf{G})} [\log p_\tau(y|G, e)]|. \tag{14}$$

598 Thus we can prove that minimizing term ① in Eq. 7 is equivalent to  $\min_{q_\theta(\mathbf{y}|\mathbf{z}), q_\theta(\mathbf{z}|\mathbf{G})} I(\mathbf{y}; \mathbf{e}|\mathbf{z})$ .  
599  $\square$

600 **B.2 Proof for Theorem 2**

601 *Proof.* Given the dependency relationship  $\mathbf{z} \leftarrow \mathbf{G} \rightarrow \mathbf{y}$ , we hold  $\max_{q(\mathbf{z}|\mathbf{G})} I(\mathbf{y}; \mathbf{z})$  is equivalent to  
602  $\min_{q(\mathbf{z}|\mathbf{G})} I(\mathbf{y}; \mathbf{G}|\mathbf{z})$ . Also we have

$$\begin{aligned}
I(\mathbf{y}; \mathbf{G}|\mathbf{z}) &= D_{KL}(p(\mathbf{y}|\mathbf{G}, \mathbf{e}) \parallel p(\mathbf{y}|\mathbf{z}, \mathbf{e})) \\
&= D_{KL}(p(\mathbf{y}|\mathbf{G}, \mathbf{e}) \parallel q(\mathbf{y}|\mathbf{z})) - D_{KL}(p(\mathbf{y}|\mathbf{z}, \mathbf{e}) \parallel q(\mathbf{y}|\mathbf{z})) \\
&\leq D_{KL}(p(\mathbf{y}|\mathbf{G}, \mathbf{e}) \parallel q(\mathbf{y}|\mathbf{z})),
\end{aligned} \tag{15}$$

603 Based on this, we will have

$$I(\mathbf{y}; \mathbf{G}|\mathbf{z}) \leq \min_{q(\mathbf{y}|\mathbf{z})} D_{KL}(p(\mathbf{y}|\mathbf{G}, \mathbf{e}) \parallel q(\mathbf{y}|\mathbf{z})). \tag{16}$$

604 Then we can also derive the following inequality via Jensen Inequality:

$$\begin{aligned}
D_{KL}(p(\mathbf{y}|\mathbf{G}, \mathbf{e}) \parallel q(\mathbf{y}|\mathbf{z})) &= \mathbb{E}_{\mathbf{e}} \mathbb{E}_{G \sim p_e(\mathbf{G})} \left[ \mathbb{E}_{y_G \sim p_e(\mathbf{y}|\mathbf{G}=G)} \mathbb{E}_{z \sim q(\mathbf{z}|\mathbf{G}=G)} \left[ \log \frac{p_e(\mathbf{y} = y_G | \mathbf{G} = G)}{q(\mathbf{y} = y_G | \mathbf{z} = z_G)} \right] \right] \\
&\leq \mathbb{E}_{\mathbf{e}} \left[ \frac{1}{|\mathcal{G}^e|} \sum_{(G, y_G) \in \mathcal{G}^e} \log \frac{p_e(\mathbf{y} = y_G | \mathbf{G} = G)}{\mathbb{E}_{z \sim q(\mathbf{z}|\mathbf{G}=G)} q(\mathbf{y} = y_G | \mathbf{z} = z_G)} \right] \\
&= C + \mathbb{E}_{\mathbf{e}} \left[ -\frac{1}{|\mathcal{G}^e|} \sum_{(G, y_G) \in \mathcal{G}^e} \log q(\mathbf{y} = y_G | \mathbf{G} = G) \right],
\end{aligned} \tag{17}$$

605 where  $C$  is a constant. Then the problem  $\min_{q(\mathbf{y}|\mathbf{z})} D_{KL}(p(\mathbf{y}|\mathbf{G}, \mathbf{e}) \parallel q(\mathbf{y}|\mathbf{z}))$  can be solve by

$$\min \mathbb{E}_{\mathbf{e}} \left[ \frac{1}{|\mathcal{G}^e|} \sum_{(G, y_G) \in \mathcal{G}^e} [-\log q_\theta(\mathbf{y} = y_G | \mathbf{G} = G)] \right], \tag{18}$$

606 which means minimizing term ② in Eq. 7 contributes to  $\max_{q_\theta(\mathbf{y}|\mathbf{z}), q_\theta(\mathbf{z}|\mathbf{G})} I(\mathbf{z}; \mathbf{y})$ .  $\square$

## 607 C Implementation Details

### 608 C.1 Baselines

609 This section describes training configurations for all baselines, which are compared in this paper.

610 **Three backbones.** We adapt three backbones into our method, namely, **GCN** [27], **Graph-**  
611 **SAGE** [18] and **GIN** [54]. We also emphasize the comparison with their augmented versions,  
612 i.e. “+ **virtual node**” [16, 22, 35]. For GCN and GIN, we use the implementations provided by Open  
613 Graph Benchmark [20]<sup>4</sup>. We implement GraphSAGE and its corresponding augmented version by  
614 ourselves. For these baselines, grid search of learning rate over  $\{1e-2, 5e-3, 1e-3, 5e-4, 1e-4\}$   
615 and dropout rate over  $\{0.1, 0.3, 0.5\}$  is performed to select the best parameters. The embedding size  
616 of all methods including ours are all set to 256 for the sake of fairness.

- 617 • **GCN** [27] is a scalable approach on graph-structured data that is based on an efficient variant of  
618 convolutional neural networks.
- 619 • **GIN** [54] generalizes the Weisfeiler-Lehman (WL) graph isomorphism test [33] and hence achieves  
620 maximum discriminative power among GNNs.
- 621 • **GraphSAGE** [18] learns a function that generates embeddings by sampling and aggregating  
622 features from a node’s local neighborhood.
- 623 • **GCN/GIN/GraphSAGE + virtual node** [16, 22, 35] is a variant of the original method augmented  
624 by an additional node connecting to all nodes in the raw graph.

625 **Models tailored for OOD learning.** We compare our method against six state-of-the-art methods:  
626 **ERM** [47], **IRM** [2], **DeepCoral** [46], **DANN** [15], **MixUp** [61] and **GroupDro** [43]. We use the  
627 implementations of these six method provided by DrugOOD<sup>5</sup>. We search for the optimal hyper-  
628 parameters by ranging learning rate over  $\{1e-3, 5e-4, 1e-4, 5e-5, 1e-5\}$  and dropout rate  
629 over  $\{0.1, 0.3, 0.5\}$ . The embedding size of all models including ours are all set to 128 for fairness.

- 630 • **ERM** [47] minimizes the average empirical loss on training data.
- 631 • **IRM** [2] penalizes feature distributions for environments that have different optimum predictors.  
632 We set the penalty weight and the penalty anneal iteration to 10 and 500, respectively.
- 633 • **DeepCoral** [46] penalizes differences in the means and covariances of the feature distributions for  
634 each environment, which are exactly the distribution of last layer activations in a neural network.  
635 The penalty weight is set to 0.001.
- 636 • **DANN** [15] encourages feature representations to be consistent across domains. We set to the  
637 inverse factor to 0.2.
- 638 • **MixUp** [61] constructs additional virtual samples for training from two examples which are  
639 randomly sampled from the training data. We set the probability and interpolate strength to 0.1.
- 640 • **GroupDro** [43] minimizes the worst-case training loss over a set of pre-defined environments.  
641 The step size is set to 0.001.

### 642 C.2 Our Method

643 We implement our method in Pytorch. As for experiments on OGB datasets, we implement the  
644 Environment Classifier, the Conditional GNN and the Substructure Encoder which are mentioned  
645 in Sec. 3.4 all in Graph Isomorphism Network (GIN) [54]. We use grid search on validation set for  
646 hyper-parameter tuning by ranging learning rate from  $\{1e-2, 5e-3, 1e-3, 5e-4, 1e-4, 5e-5,$   
647  $1e-5\}$ , dropout rate from  $\{0.1, 0.2, 0.3, 0.4, 0.5\}$ , the trading-off parameter  $\beta$  from  $\{0.5, 1, 2, 4\}$ ,  
648 the environment count  $k$  from  $\{5, 10, 15, 20, 40, 80\}$ . As for the prior  $p(\mathbf{e}|\mathbf{G})$ , we set it to a *Uniform*  
649 distribution or a discrete *Gaussian* distribution. We use CrossEntropyLoss for all models and the  
650 Adam optimizer is used for gradient-based optimization.

<sup>4</sup><https://github.com/snap-stanford/ogb>

<sup>5</sup><https://github.com/tencent-ailab/DrugOOD>

## 651 D Training Procedure

652 The training procedure of our method is summarized in Algorithm 1.

---

**Algorithm 1:** The training procedure.

---

**Input:** Dataset  $\mathcal{G}^{train} = \{(G_i, y_i)\}_{i=1}^{N^{train}}$ ; Number of training epochs for environment inference module  $E_1$ ; Number of training epochs for the molecule encoder and the predictor  $E_2$ ; Batch size  $B$ .

**Output:** Trained parameters  $\theta$ .

```
1 Initialize parameters  $\theta, \tau$  and  $\kappa$ ;  
2 for  $i \leftarrow 1$  to  $E_1$  do  
3   Sample data batches  $\mathcal{B} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_k\}$  from  $\mathcal{G}^{train}$  with batch size  $B$ ;  
4   for  $j \leftarrow 1$  to  $k$  do  
5     Compute batch loss  $\mathcal{L}_{elbo}(\tau, \kappa; \mathcal{G}_j)$  according to Eq. 6;  
6     Backpropagate  $-\mathcal{L}_{elbo}$  and optimize parameters  $\tau, \kappa$ ;  
7 Freeze the parameters  $\kappa, \tau$ ;  
8 for  $i \leftarrow 1$  to  $E_2$  do  
9   Sample data batches  $\mathcal{B} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_k\}$  from  $\mathcal{G}^{train}$  with batch size  $B$ ;  
10  for  $j \leftarrow 1$  to  $k$  do  
11    Determine the environment of each sample  $(G, y)$  in  $\mathcal{G}_k$  by  $\arg \max_e q_\kappa(e|G, y)$ ;  
12    Compute batch loss  $\mathcal{L}_{inv}(\theta; \mathcal{G}_k, \tau)$  according to Eq. 7;  
13    Backpropagate  $\mathcal{L}_{inv}$  and optimize parameters  $\theta$ ;  
14 Output the parameters  $\theta$ ;
```

---

## 653 E More Details of Datasets

654 In this paper, we use ten publicly available benchmark datasets in total. Four of them, namely,  
655 BACE, BBBP, SIDER and HIV are released by Open Graph Benchmark (OGB) [20]. The rest six are  
656 released by DrugOOD [24], i.e. IC50-assay, IC50-scaffold, IC50-size, EC50-assay, EC50-scaffold  
657 and EC50-size. We provide detailed descriptions for them as below.

- 658 • **BBBP** is a dataset of Brain-Blood Barrier Penetration. Each molecule has a label indicating  
659 whether it can penetrate through brain cell membrane to enter central nervous system.
- 660 • **BACE** is a dataset of binding affinity against human beta-secretase 1. Each molecule has a label  
661 indicating whether it binds to human beta-secretase 1.
- 662 • **SIDER** is a dataset of marked drugs and adverse drug reactions (ADRs). Molecules are grouped  
663 into 27 system organ classes.
- 664 • **HIV** is a dataset of HIV antiviral activity. Each molecule has an active or inactive label.
- 665 • **IC50/EC50-scaffold/assay/size** are datasets generated by the automated dataset curator provided  
666 by DrugOOD from the large-scale bioassay deposition website ChEMBL. The suffix specifies  
667 the splitting scheme. These six datasets target on ligand-based affinity prediction (LBAP). Each  
668 molecule has an active or inactive label.

669 Notice that all these ten datasets do not contain personally identifiable information or offensive  
670 content. Table 4 shows the detailed statistics of datasets. For all datasets, we adopt the default  
671 training-validation-test split as shown in Table 4. We use all molecules in the training set to optimize  
672 the model parameters. Then, we select hyper-parameters using the validation set, and we report the  
673 results on test molecule set for the model that achieves the best results on the validation set.

Table 4: **Summary of datasets used in this paper.** #Train/#Valid/#Test denotes the number of samples in the training/validation/test set, respectively. #Total is the sum of #Train, #Valid and #Test. #Tasks is the output dimensionality required for prediction. Additionally, we also list which split scheme is adopted and whether the manual specification of environments is available for each dataset.

	Dataset	#Train	#Valid	#Test	#Total	#Tasks	Split Scheme	Specify Environments?
OGB	BACE	1, 210	151	152	1, 513	1	Scaffold	✗
	BBBP	1, 631	204	204	2, 039	1	Scaffold	✗
	SIDER	1, 141	143	143	1, 427	27	Scaffold	✗
	HIV	32, 901	4, 113	4, 113	41, 127	1	Scaffold	✗
DrugOOD	EC50-assay	4, 540	2, 572	2, 490	9, 602	1	Assay	✓
	EC50-scaffold	2, 570	2, 532	2, 533	7, 635	1	Scaffold	✓
	EC50-size	4, 684	2, 313	2, 398	9, 395	1	Size	✓
	IC50-assay	34, 179	19, 028	19, 028	72, 235	1	Assay	✓
	IC50-scaffold	21, 519	19, 041	19, 048	59, 608	1	Scaffold	✓
	IC50-size	36, 597	17, 660	16, 415	70, 672	1	Size	✓

674 Table 5: We count the  
675 number of scaffolds that  
676 contain 1, 2, 3, 4 and 5  
677 samples, respectively.

Size	Number
1	14, 295
2	2, 330
3	862
4	449
5	255

686

Next, let’s discuss on the details of HIV dataset, which is released by Open Graph Benchmark (OGB) [20]. OGB adopts scaffold splitting scheme to split the HIV into train/validation/test set. We count the number of scaffolds that only contain 1, 2, 3, 4 and 5 molecules, respectively, and summarize the statistics in Table 5. Notice that HIV has 19, 076 scaffolds in total. We can see there are 18, 191 scaffolds containing less or equal to 5 molecules, accounting for 95.45% of the total scaffold count. HIV has a great deal of environments that contains few samples, which poses great challenge to directly applying some existing OOD generalization methods to datasets like HIV [24]. Thus, for datasets released by OGB, partitioning the molecules into different environments according to their scaffolds may not be suitable in practice. Such a observation motivates us to propose the environment-inference model.

## 687 F Notations

688 We summarize the notations used in this paper in Table 6.

## 689 G Sensitivity to Molecule Segmentation Method

690 For all experiments in our original paper, we all adopt *breaking retrosynthetically interesting chemical substructures* (BRICS) to segment molecule into substructures, which is widely used in other works related to molecules [51, 57, 9]. To investigate the sensitivity of our method to different decomposing strategies, we adopt another molecule segmentation method called *retrosynthetic combinatorial analysis procedure* (RECAP) [34], which is also available as an API in RDKit package. RECAP and BRICS decompose molecules based on two different rules. We conduct experiments on EC50-assay/scaffold/size three datasets and the comparisons are summarized in Table 7. We can see that RECAP and BRICS show competitive performance on our model and both outperform the baselines by large margins.

## 699 H Future Direction

700 Sometimes, bio-chemical properties are affected by interactions between substructures. To encode  
701 such interactions between substructures into the final learned molecular representation, we utilize the  
702 permutation equivariant Set Attention Block (SAB) proposed in Set Transformer [32]. SAB takes  
703 a representation set of any size as input and outputs a representation set of equal size. SAB is able  
704 to encode pairwise and higher-order interactions between elements in input sets into outputs. We  
705 add such a SAB after the Substructure Encoder. For each molecule, we feed the representations of

Table 6: Notations.

Notation	Description
$e$	an environment instance
$\mathbf{e}$	a random variable of $e$
$\mathcal{E}$	the support of environments
$l(\cdot)$	the loss function
$\mathcal{R}(\cdot)$	the risk function
$G$	a molecular graph instance
$\mathbf{G}$	a random variable of $G$
$y$	a ground-truth label instance
$\mathbf{y}$	a random variable of $y$
$\mathcal{G}$	a dataset set, i.e. $\{(G, y)\}$
$\psi$	the environment-inference model
$\Phi$	the molecule encoder
$\omega$	the final predictor
$\mathbf{z}$	the denotation of $\Phi(\mathbf{G})$
$f$	$\omega \circ \Phi$
$\kappa$	the learnable parameters of the Environment Classifier
$\tau$	the learnable parameters of the Conditional GNN
$\theta$	the learnable parameters of $\Phi$ and $\omega$
$k$	hyper-parameter: the environment count
$\beta$	hyper-parameter: the trading-off parameter in Eq. 7

Table 7: Comparisons on 3 out-of-distribution datasets in terms of ROC-AUC (%). The best and the runner-up in each columns are highlighted in **bolded** and underlined respectively. Note the baselines except ERM and MixUp all require environment labels. All methods including ours use GIN [54] as backbones. Each experiment is repeated 5 times with mean and standard deviation reported.

Dataset	EC50		
	Assay	Scaffold	Size
ERM [47]	69.35 $\pm$ 7.38	63.92 $\pm$ 2.09	60.94 $\pm$ 1.95
IRM [2]	69.94 $\pm$ 1.03	63.74 $\pm$ 2.15	58.30 $\pm$ 1.51
DeepCoral [46]	69.42 $\pm$ 3.35	63.66 $\pm$ 1.87	56.13 $\pm$ 1.77
DANN [15]	66.97 $\pm$ 7.19	64.33 $\pm$ 1.82	61.11 $\pm$ 0.64
MixUp [61]	70.62 $\pm$ 2.12	64.53 $\pm$ 1.66	62.67 $\pm$ 1.41
GroupDro [43]	70.52 $\pm$ 3.38	64.13 $\pm$ 1.81	59.06 $\pm$ 1.50
<b>Ours + RECAP</b>	<u>72.72 <math>\pm</math> 3.94</u>	<u>66.34 <math>\pm</math> 0.52</u>	<b>65.48 <math>\pm</math> 1.10</b>
<b>Ours + BRICS</b>	<b>73.25 <math>\pm</math> 1.24</b>	<b>66.69 <math>\pm</math> 0.34</b>	<u>65.09 <math>\pm</math> 0.90</u>

706 its substructures to SAB to obtain new substructure representations. In this way, the final molecule  
707 representation could model interactions between substructures. We conduct experiments on EC50-  
708 assay/scaffold/size to examine the performance of adding such a SAB. As demonstrated in Table 8,  
709 we can see that adding such a SAB further improves our model on EC50-scaffold. This design is a  
710 naive attempt but brings us some valuable insights.

## 711 I Limitations

712 Some studies [24, 14] have empirically shown that existing models designed for OOD learning may  
713 fail to outperform the simple ERM [47] model when the environment count is large. Though we can  
714 relabel the environment for each molecule according to the new environment partition inferred by our  
715 devised environment-inference module, we still need to set the environment count  $k$  to a relatively  
716 larger value than that of other OOD datasets from other domain, e.g. Camelyon17 [3], which only  
717 contains five environments. Thus, using our inferred environment partition, existing models designed  
718 for OOD learning might still be inferior to ERM in some cases.

Table 8: Comparisons on 3 out-of-distribution datasets in terms of ROC-AUC (%). The best and the runner-up in each columns are highlighted in **bolded** and underlined respectively. Note the baselines except ERM and MixUp all require environment labels. All methods including ours use GIN [54] as backbones. Each experiment is repeated 5 times with mean and standard deviation reported.

Dataset Environment	EC50		
	Assay	Scaffold	Size
<b>ERM</b> [47]	69.35 ± 7.38	63.92 ± 2.09	60.94 ± 1.95
<b>IRM</b> [2]	69.94 ± 1.03	63.74 ± 2.15	58.30 ± 1.51
<b>DeepCoral</b> [46]	69.42 ± 3.35	63.66 ± 1.87	56.13 ± 1.77
<b>DANN</b> [15]	66.97 ± 7.19	64.33 ± 1.82	61.11 ± 0.64
<b>MixUp</b> [61]	70.62 ± 2.12	64.53 ± 1.66	62.67 ± 1.41
<b>GroupDro</b> [43]	70.52 ± 3.38	64.13 ± 1.81	59.06 ± 1.50
<b>Ours</b>	<b>73.25 ± 1.24</b>	<u>66.69 ± 0.34</u>	<b>65.09 ± 0.90</b>
<b>Ours + SAB</b>	<u>73.15 ± 2.69</u>	<b>67.26 ± 1.54</b>	<u>64.83 ± 1.07</u>

## 719 J Potential Negative Impacts

720 As far as we are concerned, we have not identified any negative social impact of this work.