
Deep Interactions for Multimodal Molecular Property Prediction

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Multi-modal learning by means of leveraging both 2D graph and 3D point cloud
2 information has become a prevalent method to improve model performance in
3 molecular property prediction. However, many recent techniques focus on specific
4 pre-training tasks such as contrastive learning, feature blending, and atom/subgraph
5 masking in order to learn multi-modality even though design of model architecture
6 is also impactful for both pre-training and downstream task performance. Rely-
7 ing on pre-training tasks to align 2D and 3D modalities lacks direct interaction
8 which may be more effective in multimodal learning. In this work, we propose
9 MOLINTERACT, which takes a simple yet effective architecture-focused approach
10 to multimodal molecule learning which addresses these challenges. MOLINTER-
11 ACT leverages an interaction layer for fusing 2D and 3D information and fostering
12 cross-modal alignment, showing strong results using even the simplest pre-training
13 methods such as predicting features of the 3D point cloud and 2D graph. MOLIN-
14 TERACT exceeds several current state-of-the-art multimodal pre-training techniques
15 and architectures on various downstream 2D and 3D molecule property prediction
16 benchmark tasks.

17 1 Introduction

18 AI-assisted drug discovery has driven recent research interest in utilizing neural networks for molecule
19 learning. The machine learning community has become especially interested in developing high-
20 quality representations for molecules, which are crucial for predicting molecular properties for a
21 variety of downstream cheminformatic tasks. Self-supervised learning (SSL) on molecular data
22 has emerged as a prevalent research direction to achieve this, leveraging the 2D graph structures
23 of molecules [22, 56, 66]. In parallel, many SSL strategies for 3D point cloud representations of
24 molecules have also been developed [36, 15]. More recent works demonstrate the effectiveness of mul-
25 timodal SSL techniques which combine 2D and 3D modalities to create better unified representations
26 [54, 34, 36, 77, 60].

27 Many of these recent and successful multimodal SSL methods make use of SSL techniques from
28 other fields of machine learning. For example, many works leverage attribute and atom/subgraph
29 masking & prediction [22, 66, 75, 24, 34, 35] similar to how masked language-modeling is used to
30 pre-train large transformer-based networks such as BERT [8]. Other works [34, 54, 62] prefer to
31 leverage contrastive learning [5] in order to align the 2D and 3D views of molecules in a unified
32 embedding space similar to how CLIP [45] aligns caption and image embeddings and SimCLR [3]
33 aligns views of images.

34 These recent approaches typically consider improving molecular SSL via specific pre-training
35 strategies and tasks, but not the underlying architecture. For instance, a common approach [34, 35, 54]
36 is to take two separate models for encoding 3D and 2D structures and then design a pre-training

37 task to align their output embeddings. Alternatively, other works take a single, modality-agnostic
38 model, usually a transformer [58], and task it with predicting multimodal properties such as bond
39 angles [60, 77] or shortest-path distances [73]. Both of these approaches rely on a chosen pre-training
40 task to align 2D and 3D views of molecules.

41 However, it is not clear to what extent such approaches are able to fully learn cross-modal interactions.
42 For example, contrastive approaches using separate encoders seek to maximize the mutual information
43 between coarse-grained molecule embeddings, and so they may fail to capture key fine-grained
44 relationships. On the other hand, predictive approaches using a single backbone typically accept
45 only atom identities as input, leaving the pre-training task as the only source of multi-modality,
46 potentially missing features which can be extracted by modality-specific encoders. An additional
47 issue is that many pre-training tasks are complex and may require extensive tuning. For example,
48 GraphMVP [34] and MoleculeSDE [35] use a variational autoencoder [27] and diffusion model [19],
49 respectively, to reconstruct the original 2D and 3D structures, which risk mode collapse and are
50 sensitive to hyperparameters such as noise schedules. Other techniques like MoleBlend [73] require
51 both coordinate denoising [15] and prediction of blended multimodal features, where the ratio of
52 noised nodes and 2D & 3D features needs to be tuned.

53 In order to achieve fine-grained multi-modal information with simpler pre-training tasks, we turn our
54 focus to the role of architectural design for more effective SSL. This work introduces MOLINTERACT,
55 a deep learning architecture designed to fuse 2D and 3D modalities of molecules to better foster
56 multimodal performance. MOLINTERACT uses a series of interaction layers to learn how to combine
57 2D and 3D embeddings. Specifically, MOLINTERACT consists of two message-passing endpoints for
58 2D and 3D data to produce corresponding 2D and 3D embeddings which are then fused and split apart
59 repeatedly in order to exchange unimodal information during pre-training. We pair MOLINTERACT
60 with a set of simple pre-training tasks from the existing literature, such as bond and dihedral angle
61 prediction, which are both sensible in a molecular context and require virtually no tuning. We
62 show that, even with such straightforward pre-training tasks and architecture, MOLINTERACT is
63 able to yield strong multimodal performance, emphasizing the impact of directly fusing 2D and 3D
64 atom embeddings in a model-based approach to improving SSL. We conduct various experiments to
65 demonstrate state-of-the-art performance across various downstream 2D and 3D benchmark tasks.

66 2 Background and Related Work

67 **SSL for molecules.** Self-supervised learning (SSL) [30, 37] has been adopted in a wide range of
68 domains to obtain high quality representations for downstream tasks. A slew of recent works have
69 emerged attempting to apply the same pre-train-then-finetune paradigm to molecule learning. In
70 particular, research has been aimed at molecular SSL with the primary downstream task of molecule
71 property prediction [4, 75] in mind due to the potential of saving tremendous amounts of time
72 screening new drugs and compounds. However, the success of molecule property prediction requires
73 a comprehensive extraction of molecular features from various modalities, which becomes especially
74 important when only one modality is available for a given real-world molecule. For example, in
75 certain cases, only a compound’s 2D structure may be known, but there may be little to no data on its
76 equilibrium conformers. In light of this, it has become important to solve the challenge of learning
77 informative molecular representations using all kinds of modalities, particularly 2D graphs and 3D
78 point clouds.

79 **Existing work on multimodal SSL.** In order to fuse multimodal representations, works such as
80 GraphMVP [37] and 3DInfomax [54] aim to maximize the mutual information between the 2D and
81 3D views of molecules, treating 2D graphs with their corresponding 3D conformations as positive
82 samples and all other pairs as negative samples. Alternatively, another line of work proposes to
83 incorporate both modalities via prediction of the original data. MoleculeSDE [35] generates 3D
84 SE(3)-equivariant conformations from the 2D graph, and vice versa, and Zhu et al. [78] use a single
85 model to reconstruct the input 2D graph from the 3D point cloud and vice versa. Similarly, many
86 works [34, 35] task unimodal models with predicting masked sets of 2D and 3D atoms. Other works
87 such as ChemRL-GEM [10] and 3D-PGT [61] propose to predict internal coordinates such as bond
88 length, bond angle, and dihedral angles in order to distill 2D and 3D information. In contrast to these
89 methods, MOLINTERACT seeks to supplement a set of predictive pre-training tasks by leveraging a
90 fusion layer to force interactions between 2D and 3D embeddings to facilitate multimodal learning.

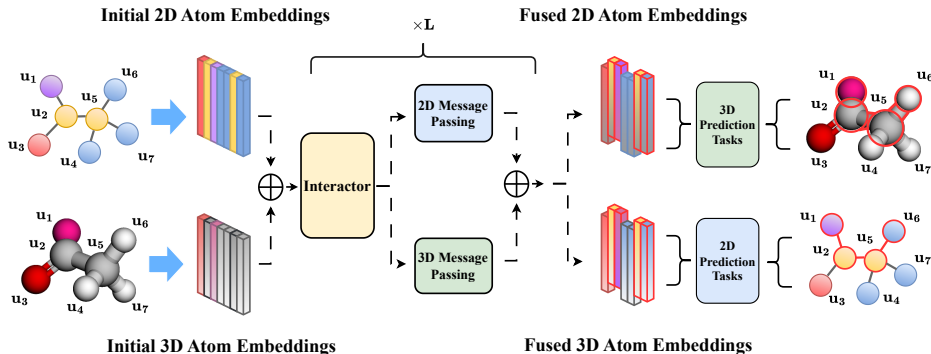


Figure 1: The proposed method’s pre-training pipeline. From left-to-right, the input molecule’s 2D and 3D graphs are used to derive initial 2D and 3D atom embeddings via message-passing layers. These embeddings are then mixed by an interaction layer before being fed back into the unimodal message-passing branches of the architecture. This process of message-passing followed by interaction repeats L times before the final embeddings from each tower are used for pre-training tasks of the opposite modality, e.g. predicting 3D quantities using the 2D encoder embeddings. Not shown are residual connections between each interaction block to preserve lower-order information.

91 **Modality interaction.** In order to fully leverage the synergy between different modalities, recent
 92 works from other fields propose to learn more fine-grained modality alignment through deep inter-
 93 active architectures. GreaseLM [74] and Dragon [71] propose to align language models and graph
 94 neural networks on knowledge graphs through an interaction token, aiming to integrate text and graph
 95 modalities to better identify relevant relations and entities in a given passage. Other works [6, 31] de-
 96 sign similar deep interaction layers in various domains such as social networks and recommendation.
 97 MOLINTERACT takes inspiration from these methods, proposing to interact 2D and 3D modalities on
 98 a fine-grained level in order to better facilitate pre-training and create high-performing multimodal
 99 representations.

100 3 Method - Deep Interactions

101 3.1 Notation and Preliminaries

102 We consider molecules in terms of their 2D graph and 3D point cloud modalities. For simplicity, we
 103 will use the term “graph” to refer to both the typical 2D node-edge formalism as well as a molecule’s
 104 3D point cloud. We denote the 2D graph of a molecule with n atoms by $G_{2D} = (V, E, \mathbf{X}, \mathbf{B})$
 105 where V is a set of its atoms (nodes), E is a set of its bonds (edges), $\mathbf{X} \in \mathbb{R}^{n \times d_V}$ is a 2D feature
 106 matrix corresponding to the atoms of the molecule with features specific to the 2D graph, such
 107 as membership in a ring [21], and $\mathbf{B} \in \mathbb{R}^{|E| \times d_E}$ is an edge feature matrix corresponding to edge
 108 information such as bond type. For simplicity, we let $d = d_V = d_E$. We also define the 3D graph of
 109 a molecule by $G_{3D} = (\mathbf{R}, \mathbf{X})$ where $\mathbf{R} \in \mathbb{R}^{n \times 3}$ is the molecule’s position matrix where rows are
 110 (x, y, z) coordinates in 3D space. Unless otherwise specified, we use $z_{2D}^{(\ell, i)}$ and $z_{3D}^{(\ell, i)}$
 111 to refer to the i th 2D and 3D atom embedding resulting from the ℓ th layer of a neural network. For simplicity, we
 112 assume that all embeddings are of dimension d , and we use $\mathbf{W}_{(\cdot)}$ to refer to a linear layer with the
 113 bias omitted. All classification-based loss functions use CE to stand for Cross-Entropy.

114 MOLINTERACT is comprised of two central components: (1) an architecture which fosters deep
 115 multimodal interactions, and (2) a pre-training scheme which leverages this architecture to enforce
 116 multimodal understanding similar to previous works. We introduce each component one-by-one in
 117 the following sections.

118 3.2 Model Architecture

119 During pre-training and multi-modal fine-tuning, MOLINTERACT receives a molecule’s 2D and 3D
 120 views G_{2D} and G_{3D} . These views are then passed through a two-tower architecture which alternates

121 between phases of message-passing and interaction. Each tower is a 2D and 3D modality-specific
122 stack of encoders periodically conjoined by interactor layers as visualized in Figure 1.

123 **2D and 3D atom encoders.** In order to compute 2D atom embeddings, we follow previous work
124 on multimodal pre-training [34, 54, 36, 77, 60] and use message-passing graph neural network [13]
125 (MPNN) layers as 2D encoders. Given a molecular graph G_{2D} and one of its nodes i , its 2D node
126 embedding h_i at the $(\ell + 1)$ th layer of an MPNN is given by

$$h_i^{(\ell+1)} = \text{Update} \left(h_i^{(\ell)}, \text{Agg}_{j \in \mathcal{N}(i)} \psi \left(h_i^{(\ell)}, h_j^{(\ell)}, e_{ij} \right) \right) \quad (1)$$

127 where Update is a function which updates the node embedding, Agg is a permutation-invariant
128 function on the neighbors of i , and ψ is a function which computes "messages", or interactions,
129 between the node i and its neighbor j with the edge between them as context. In our case, we use
130 layers from the GINE [22] architecture, which is a variant of GIN [68] that incorporates edge features
131 during message-passing. We choose GINE due to its simplicity and 1-WL-expressivity [64], although
132 we note that MOLINTERACT places no restrictions on its 2D backbone, and one can easily replace the
133 MPNN with a more powerful 2D model such as a graph transformer [9, 48, 40] as other multimodal
134 works [60, 73] do to improve performance.

135 To compute 3D atom embeddings, we opt to use the continuous convolutional layers from SchNet [52].
136 These layers conduct message-passing according to relative distances between atoms, incorporating
137 both geometric and atom information into the resulting embeddings. Similar to the 2D encoder,
138 MOLINTERACT is agnostic to the choice of 3D encoder, and so one may choose to opt for more
139 expressive 3D backbones [12, 28, 38, 59, 53]. Despite the limited expressivity of SchNet and GIN,
140 we find that MOLINTERACT is able to outperform several state-of-the-art methods, which we show in
141 Section 4.

142 3.3 Multimodal Interaction Layer

143 At the ℓ th layer of unimodal message-passing, we take the 2D and 3D atom embeddings $z_{2D}^{(\ell,i)}$ and
144 $z_{3D}^{(\ell,i)}$ and pass them to an interaction layer $\phi^{(\ell)}$. Then, the updated atom embeddings $z_{2D}^{(\ell+1,i)}$ and
145 $z_{3D}^{(\ell+1,i)}$ are decoded from the output of $\phi^{(\ell)}$ and fed back into their respective unimodal message-
146 passing towers. There are a variety of options for choosing $\phi^{(\ell)}$, such as an attention-based aggregation
147 between the embeddings, or the aggregation of representative nodes [74], such as a virtual node for
148 the 2D graph [13, 44, 25] or a center-of-mass node for the 3D graph. However, for simplicity, we
149 use a 2-layer MLP of dimension $2d$ with Swish activation [18] for $\phi^{(\ell)}$, and feed it the concatenated
150 unimodal embeddings. We run ablations testing different functions for each $\phi^{(\ell)}$ in Table 2. For
151 decoding, we simply split the output of the MLP in half along the channel dimension to retrieve the
152 updated 2D and 3D embeddings. With this, we do not risk information loss via pooling or choosing
153 a representative token for the whole molecule, attaining more granular, node-level interactions.
154 Formally, the multimodal embeddings at layer $\ell + 1$ are given by

$$w^{(\ell+1,i)} = \phi^{(\ell)}(z_{2D}^{(\ell,i)}, z_{3D}^{(\ell,i)}) = \text{MLP}^{(\ell)} \left(z_{2D}^{(\ell,i)} \parallel z_{3D}^{(\ell,i)} \right) \quad (2)$$

155 where \parallel denotes concatenation in column-major order. Then, our updated atom embeddings can be
156 written in the following implementation-friendly way:

$$z_{2D}^{(\ell+1,i)}, z_{3D}^{(\ell+1,i)} = w_{:,;d}^{(\ell+1,i)}, w_{:,d}^{(\ell+1,i)} \quad (3)$$

157 where the subscripts of $w^{(\ell+1,i)}$ denote Python-like indices. this way, the 2D embeddings are a
158 fusion of both 2D and 3D features, and similarly for the 3D embeddings with each subsequent
159 iteration of message-passing and interaction, encoding higher-order multimodal features. Unlike
160 molecule-level approaches like 3D Infomax [54], GraphMVP [34], and MoleculeSDE [35], and
161 unlike modality-agnostic backbone-based methods like MoleBlend [73], MOLINTERACT benefits
162 from both fine-grained, atom-level interactions as well as modality-specific encoders to create more
163 powerful multimodal representations. See Appendix F to see a UMAP visualization of the molecule
164 embeddings from MOLINTERACT compared with MoleculeSDE.

165 The pre-training tasks for MOLINTERACT are designed to be heterogeneous, meaning 3D quantities
 166 are predicted using $z_{2D}^{(\ell,i)}$, and 2D quantities are predicted using $z_{3D}^{(\ell,i)}$. By using embeddings of one
 167 modality to predict features of the opposite modality, we maximize a lower bound of the mutual
 168 information between the modalities. As Liu et al. [34] note, if X_{3D}, X_{2D} denote random variables
 169 from 2D and 3D spaces, then their mutual information (MI) $I(X_{3D}, X_{2D})$ is bounded below by
 170 $-\frac{1}{2} (H(X_{3D}|X_{2D}) + H(X_{2D}|X_{3D}))$ where H denotes entropy. Visibly, this bound is maximized
 171 when $p(x_{3D}|x_{2D})$ and $p(x_{2D}|x_{3D})$ are also maximized. This motivates the pre-training pipeline in
 172 this work where predicting 2D quantities from 3D embeddings and 3D quantities from 2D embeddings
 173 maximizes the MI between 2D and 3D information in MOLINTERACT.

174 Intuitively, during pre-training, the interaction layer ϕ serves as an exchange pathway between the
 175 unimodal towers. As the only point of contact between the 2D and 3D message-passing layers, ϕ
 176 must effectively route the most important cross-modal information relevant to the pre-training task.
 177 During fine-tuning, each ϕ serves as an aggregator of multimodal features, learning to fuse 2D and
 178 3D information effectively for the given downstream task.

179 3.4 Pre-training Tasks

180 MOLINTERACT’s architecture is complemented by a set of simple pre-training tasks in order to
 181 facilitate multimodal learning using its fused atom embeddings. Specifically, for our 3D tasks, we
 182 choose to predict interatomic distances, bond angles, and dihedral angles. We then predict edge type,
 183 shortest-path distance, and centrality ranking as our 2D tasks. We choose these specific pre-training
 184 tasks due to (1) their hierarchical relationships with each other in their respective modalities and
 185 (2) their ease of computation. The intuition behind (1) is that, in order to learn a comprehensive
 186 multimodal representation, both lower-order and higher-order geometric/topological features must
 187 be learned and infused in the post-interaction atom embeddings. We value (2) for efficiency’s sake,
 188 noting that computing these quantities is fairly straightforward and require little to no tuning to be
 189 effective. Such prediction tasks are among the simplest in the molecular SSL literature compared
 190 to diffusion models and substructure masking, which highlights the effectiveness of the interaction
 191 layers during pre-training.

192 **3D Tasks.** For our 3D tasks, similar to [60] and [10], we choose to predict interatomic distances,
 193 bond angles, and dihedral angles: three of the some of the basic internal coordinates in 3D molecular
 194 graphs. Not only are such 3D graphs uniquely identified by these primitives [38, 59], but they are also
 195 crucially linked to energy-based properties of molecules, making them important geometric priors to
 196 encode in 2D embeddings. Furthermore, these angles are prime examples of hierarchical quantities
 197 since interatomic distances can be used to compute bond angles, and bond angles can be used to
 198 compute dihedral angles, ascending from 2-tuple atom information to 4-tuple atom information.

199 Given an L -layer instance of MOLINTERACT, We introduce our interatomic distance loss and bond
 200 angle losses as follows:

$$\mathcal{L}_{\text{Inter}} = \frac{1}{|I|} \sum_{(i,j) \in I} \left(\mathbf{W}_{\text{Inter}} \left(z_{2D}^{(L,i)} + z_{2D}^{(L,j)} \right) - \alpha_{ij} \right)^2 \quad (4)$$

201 where I is a set of sampled atom index pairs, and α_{ij} is the ground-truth interatomic distance between
 202 atoms i and j . Note that we add the 2D embeddings in this loss formulation due to the requirement
 203 that distances between atoms are symmetric ($\alpha_{ij} = \alpha_{ji}$), and therefore their encodings should be
 204 permutation-invariant. Our loss function for predicting bond angles is similarly defined:

$$\mathcal{L}_{\text{B}} = \frac{1}{|B|} \sum_{(i,j,k) \in B} \left(\mathbf{W}_{\text{B}} \left(z_{2D}^{(L,i)} \| z_{2D}^{(L,j)} \| z_{2D}^{(L,k)} \right) - \beta_{ijk} \right)^2 \quad (5)$$

205 where B is a set of computed bond angles of adjacent atoms, and i, j, k denote indices of the respective
 206 anchor, source, and destination atoms. Note that we concatenate the indexed atom embeddings rather
 207 than sum them since bond angles differ depending on which atoms are chosen as anchors.

208 Finally, while our dihedral angle prediction loss could be analogously defined, we found that our
 209 model had difficulty predicting dihedral angles directly, with MSE barely reducing from 0.475 in
 210 the first epoch to 0.45 in the 50th when pre-training on PCQM4Mv2, indicating that no useful angle
 211 information was being learned. To improve learning, we replace the angle regression task with an

212 angle classification task, where the task becomes to categorize quadruplets of atom embeddings
 213 according to what bin the corresponding dihedral angle belongs to. This is a relatively easier task
 214 than direct prediction, and we subsequently saw an increase in performance. Formally, our loss is the
 215 cross-entropy term

$$\mathcal{L}_D = -\frac{1}{|D|} \sum_{(i,j,k,l) \in D} \text{CE} \left(\mathcal{D}_b, \mathbf{W}_D \left(z_{2D}^{(L,*)} \right) \right) \quad (6)$$

216 where D is a set of dihedral angles, $z_{2D}^{(L,*)}$ is short for $z_{2D}^{(L,i)} || z_{2D}^{(L,j)} || z_{2D}^{(L,k)} || z_{2D}^{(L,\ell)}$ where i, j, k, l are
 217 indices of the atoms which form a dihedral angle from D , and each \mathcal{D}_b indexes the bin to which each
 218 angle belongs. In our case, we use $|D| = 20$. We note that Guha et al. [17] take a similar approach to
 219 turning a regression problem into a classification problem, although they do this to aid in conformer
 220 prediction rather than angle prediction.

221 **2D Tasks.** The ways in which 2D topological quantities are relevant for molecular property pre-
 222 diction are more subtle than in the 3D case. Given that G_{3D} is missing key information regarding
 223 atom-atom relationships such as bond types, we first task MOLINTERACT with classifying edges
 224 from G_{2D} according to the cross-entropy loss term

$$\mathcal{L}_{\text{Edge}} = -\frac{1}{|E|} \sum_{(i,j) \in E} \text{CE} \left(\mathbf{B}_{ij}, \mathbf{W}_{\text{Edge}} \left(z_{3D}^{(L,i)} + z_{3D}^{(L,j)} \right) \right) \quad (7)$$

225 where \mathbf{B}_{ij} indexes the label of the corresponding edge (i, j) . With this loss, we aim to install
 226 precise atom relational information in the 3D embeddings. Next, a logically higher-order task is
 227 to determine the shortest-path distances (SPDs) between atoms, similar to Transformer-M [40] and
 228 MOLEBLEND [73], which encodes a global characterization of the molecule’s topology. Further,
 229 bond type information may serve as a useful preliminary task given that edge classification implicitly
 230 informs the modeling of existing edges, meaning that SPD prediction becomes a task of counting
 231 which of the said edges are incident. Formally, our SPD loss is formulated as

$$\mathcal{L}_{\text{SPD}} = -\frac{1}{|C|} \sum_{(i,j) \in C} \text{CE} \left(\mathbf{D}_{ij}, \mathbf{W}_{\text{SPD}} \left(z_{3D}^{(L,i)} + z_{3D}^{(L,j)} \right) \right) \quad (8)$$

232 where $C \subseteq V \times V$ is a set of sampled node pairs, and \mathbf{D}_{ij} corresponds to the SPD between atoms i
 233 and j .

234 Our final 2D pre-training task is centrality ranking, which aims to use SPD information from the
 235 preceding pre-training task to capture global structure. Centrality [1, 2] is a concept from network
 236 science which quantifies node importance. In the molecular case, centrality might be used as an
 237 indicator of structural importance such as acting as a bridge between a ring or functional group [43].
 238 Furthermore, centrality may act as a proxy for structure/subgraph membership since atoms which
 239 participate in chemically relevant substructures are likely to have similar centrality measures. This
 240 may serve as informative signal for the 3D tower of MOLINTERACT which is ignorant of the 2D
 241 graph topology. In this way, learning to rank nodes by centrality may be thought of as a proxy task to
 242 more complex structural pre-training tasks such as subgraph masking, replacing subgraph sampling
 243 steps with the simple cross-entropy loss term

$$\mathcal{L}_{\text{Cent}} = -\frac{1}{|C|} \sum_{(i,j) \in C} \text{CE} \left(\mathcal{C}_{i,j}, \mathbf{W}_{\text{Cent}} \left(z_{3D}^{(L,i)} + z_{3D}^{(L,j)} \right) \right) \quad (9)$$

244 where $\mathcal{C}_{i,j} = 1$ if node i has a higher centrality than node j and 0 otherwise. Among the various
 245 definitions of centrality, we experiment with betweenness centrality and eigenvector centrality [63, 11].
 246 Betweenness centrality of a node v is defined as $\sum_{s,t \in V} \frac{\sigma(s,t|v)}{\sigma(s,t)}$, where $\sigma(s,t|v)$ stands for the
 247 number of shortest paths between s and t which pass through v , which would appear to reuse
 248 information from \mathcal{L}_{SPD} . However, in practice, we observe superior performance using the eigenvector
 249 centrality of each node u , defined as the u th entry of the eigenvector corresponding to the largest
 250 eigenvalue of the 2D graph’s adjacency matrix. Intuitively, information learned during SPD prediction
 251 may still be used since a node’s eigenvector centrality is proportional to the number of infinite random
 252 walks passing through that node. We provide some visualizations of betweenness and eigenvector
 253 centrality on molecular graphs in Appendix B.

Method	α	$\Delta\mathcal{E}$	$\mathcal{E}_{\text{HOMO}}$	$\mathcal{E}_{\text{LUMO}}$	μ	C_v	G	H	R^2	U	U_0	ZPVE
Stock SchNet (Schütt et al. [52]), 8 layers	0.076	51.28	33.17	26.53	0.032	0.031	17.86	15.77	0.146	17.88	18.24	1.605
Distance Prediction (Liu et al. [36])	0.065	45.87	27.61	23.34	0.031	0.033	14.83	15.81	0.248	15.07	15.01	1.837
3D InfoGraph (Liu et al. [36])	0.062	45.96	29.29	24.60	0.028	0.030	13.93	13.97	0.133	13.55	13.47	1.644
3D InfoMax (Stärk et al. [54])	0.057	42.09	25.90	21.60	0.028	0.030	13.73	13.62	0.141	13.81	13.30	1.670
GraphMVP (Liu et al. [34])	0.056	41.99	25.75	21.58	0.027	0.029	13.43	13.31	0.136	13.03	13.07	1.609
MoleculeSDE (Liu et al. [35])	0.054	41.77	25.74	21.41	0.026	0.028	13.07	12.05	0.151	12.54	12.04	1.587
MOLEBLEND (Yu et al. [73])	0.060	34.75	<u>21.47</u>	<u>19.23</u>	0.037	0.031	12.44	11.97	0.417	12.02	11.82	1.580
MOLINTERACT (no pre-training)	<u>0.048</u>	37.66	21.87	19.45	<u>0.022</u>	<u>0.026</u>	9.54	8.84	<u>0.119</u>	8.77	8.421	1.396
MOLINTERACT ($\mathcal{L}_{\text{Simple}}$)	0.047	35.92	21.54	<u>18.34</u>	0.021	0.025	<u>9.13</u>	<u>8.26</u>	0.097	<u>8.16</u>	<u>8.17</u>	<u>1.365</u>
MOLINTERACT (\mathcal{L}_{All})	0.047	<u>35.58</u>	20.60	17.88	0.021	0.025	8.56	8.24	0.136	7.92	7.72	1.327

Table 1: Performance on QM9 measured in MAE. Lower is better.

Pre-training method	α	$\Delta\mathcal{E}$	$\mathcal{E}_{\text{HOMO}}$	$\mathcal{E}_{\text{LUMO}}$	μ	C_v	G	H	R^2	U	U_0	ZPVE
Only 3D tasks	0.048	36.58	21.12	18.54	0.024	0.026	9.45	8.80	0.152	8.62	8.70	1.448
Betweenness	0.050	37.78	22.18	18.90	0.025	0.027	10.08	10.30	0.171	9.71	10.50	1.508
Mean Interactor	0.061	46.82	30.89	24.18	0.035	0.030	11.93	12.06	0.126	11.67	11.99	1.543
Self-Attention Interactor	0.074	52.73	32.27	27.74	0.042	0.034	14.80	14.09	0.116	14.41	13.98	1.749
$\mathcal{L}_{\text{Simple}}$, 3D structures only	0.057	42.00	24.78	20.77	0.022	0.028	13.47	12.87	0.163	12.75	12.36	1.480

Table 2: Ablations of MOLINTERACT on QM9.

254 Finally, the total loss formulation during pre-training is $\mathcal{L}_{\text{All}} = \mathcal{L}_{\text{Inter}} + \mathcal{L}_{\text{B}} + \mathcal{L}_{\text{D}} + \mathcal{L}_{\text{Edge}} + \mathcal{L}_{\text{SPD}} + \mathcal{L}_{\text{Cent}}$.
255 In practice, we see that each loss term exhibits varying influence since terms like $\mathcal{L}_{\text{edge}}$ are naturally
256 easier to minimize than more complex terms like \mathcal{L}_{D} . Therefore, in our experiments, we compare
257 both MOLINTERACT using \mathcal{L}_{All} and $\mathcal{L}_{\text{Simple}} = \mathcal{L}_{\text{Inter}} + \mathcal{L}_{\text{SPD}}$, and find comparable performance. We
258 find that $\mathcal{L}_{\text{Inter}}$ and \mathcal{L}_{SPD} work best together, achieving the best overall performance among all the
259 tasks considered. Intuitively, interatomic distances with shortest-path distances give the minimum
260 description of the topology of the 3D and 2D graphs. For a more detailed analysis of the behavior of
261 these losses, see Appendix E. In summary, each of these 3D and 2D pre-training tasks play a role
262 in forming a unified molecular representation in terms of geometric and topological quantities in
263 increasing levels of complexity.

264 4 Experiments

265 4.1 Datasets and Experimental Setup

266 We pre-train an 8-layer version of MOLINTERACT with 9M parameters for 50 epochs on
267 PCQM4Mv2 [21], which contains over 3.3M molecules with their DFT-computed 3D structures
268 from the PubChemQC [42] project. For evaluation, we evaluate MOLINTERACT on 12 tasks from
269 QM9 [46] and 8 datasets from MoleculeNet [65] in order to compare our method with works in the
270 multimodal molecular SSL literature. We compare MOLINTERACT on QM9 and MoleculeNet with
271 the same baselines as reported by the comprehensive study by Yu et al. [73]. All best metrics are
272 **bolded**, and second-best metrics are underlined. Results for QM9 are measured in mean absolute
273 error (MAE), and results for MoleculeNet are measured in ROC AUC. All metrics reported are from
274 each dataset’s test split using the weights which perform best on a validation set. We also include
275 results on QM8 [47, 51] in Appendix D due to space limitations.

276 In datasets where both 2D and 3D information are available such as QM9 and QM8, we provide both
277 2D and 3D structures to MOLINTERACT, aggregate the resulting embeddings with mean pooling,
278 and then input their concatenation to a 2-layer MLP head. Otherwise, when only one modality is
279 available, as in MoleculeNet, only the corresponding unimodal branch of our method is activated
280 while the frozen atom embeddings of the other modality are used as input to the interaction layers in
281 place of the embeddings produced by the disabled complementary branch.

Method	BBBP	Tox21	ToxCast	SIDER	ClinTox	MUV	HIV	BACE	Average
AttrMask (Hu et al. [22])	65.0±2.3	74.8±0.2	62.9±0.1	61.2±0.1	87.7±1.1	73.4±2.0	76.8±0.5	79.7±0.3	72.68
ContextPred (Hu et al. [22])	65.7±0.6	74.2±0.0	62.5±0.3	62.2±0.5	77.2±0.8	75.3±1.5	77.1±0.8	76.0±2.0	71.28
GraphCL (You et al. [72])	69.7±0.6	73.9±0.6	62.4±0.5	60.5±0.8	76.0±2.6	69.8±2.6	78.5±1.2	75.4±1.4	70.78
InfoGraph (Sun et al. [56])	67.5±0.1	73.2±0.4	63.7±0.5	59.9±0.3	76.5±1.0	74.1±0.7	75.1±0.9	77.8±0.8	70.98
GROVER (Rong et al. [50])	<u>70.0±0.10</u>	74.3±0.1	65.4±0.4	<u>64.8±0.6</u>	81.2±3.0	67.3±1.8	62.5±0.9	82.6±0.7	71.01
MolCLR (Wang et al. [62])	66.6±1.8	73.0±0.1	62.9±0.3	57.5±1.7	86.1±0.9	72.5±2.3	76.2±1.5	71.5±3.1	70.79
GraphLoG (Xu et al. [69])	72.5±0.8	75.7±0.5	63.5±0.7	61.2±1.1	76.7±3.3	76.0±1.1	77.8±0.8	<u>83.5±1.2</u>	73.40
MGSSL (Zhang et al. [75])	69.7±0.9	76.5±0.3	64.1±0.7	61.8±0.8	80.7±2.1	<u>78.7±1.5</u>	78.8±1.2	79.1±0.9	73.70
GraphMAE (Hou et al. [20])	72.0±0.6	75.5±0.6	64.1±0.3	60.3±1.1	82.3±1.2	76.3±2.4	77.2±1.0	83.1±0.9	73.85
Mole-BERT (Xia et al. [66])	71.9±1.6	76.8±0.5	64.3±0.2	62.8±1.1	78.9±3.0	78.6±1.8	78.2±0.8	80.8±1.4	74.04
3D InfoMax (Stärk et al. [54])	69.1±1.0	74.5±0.7	64.4±0.8	60.6±0.7	79.9±3.4	74.4±2.4	76.1±1.3	79.7±1.5	72.34
GraphMVP (Liu et al. [34])	68.5±0.2	74.5±0.4	62.7±0.1	62.3±1.6	79.0±2.5	75.0±1.4	74.8±1.4	76.8±1.1	71.69
MoleculeSDE (Liu et al. [35])	71.8±0.7	76.8±0.3	65.0±0.2	60.8±0.3	87.0±0.5	80.9±0.3	78.8±0.9	79.5±2.1	<u>75.07</u>
MOLEBLEND (Yu et al. [73])	73.0±0.8	77.8±0.8	66.1±0.0	64.9±0.3	<u>87.6±0.7</u>	77.2±2.3	<u>79.0±0.8</u>	83.7±1.4	76.16
MOLINTERACT ($\mathcal{L}_{\text{Simple}}$)	67.2±3.9	76.4±0.2	64.5±0.2	62.5±0.4	86.1±0.4	78.6±0.4	78.6±0.8	82.4±1.7	74.52
MOLINTERACT (\mathcal{L}_{All})	68.5±1.3	<u>77.3±0.5</u>	<u>65.4±0.2</u>	62.9±0.4	88.4±1.0	77.1±3.1	79.5±0.4	79.1±0.03	74.77

Table 3: Performance on MoleculeNet measured in ROC AUC. Higher is better.

282 4.2 3D Datasets - QM9

283 In QM9, we follow Thölke and Fabritiis [57] and finetune on 110K random molecules and use the
284 remaining 10K and 10.8K molecules as validation and test sets, respectively. QM9 is a dataset of
285 small molecules designed to test models’ abilities to predicting various quantum and thermodynamic
286 properties, which crucially depend on 3D information. Performance is measured in terms of MAE
287 in order to determine how closely models can match DFT-level approximations of various quantum
288 properties of small molecules. Per Table 1, both versions of MOLINTERACT exhibit a substantial
289 lead in performance compared to baseline 3D pre-training methods, even without pre-training,
290 demonstrating the effectiveness of the deep multimodal interaction layers. Including the 3D and 2D
291 pre-training tasks, we see that MOLINTERACT’s performance improves across the board, exceeding
292 the most recent state-of-the-art by as much as 34% (U), further validating the use of deep interactions
293 for improving the quality of learned features even with simple predictive SSL tasks.

294 4.3 2D Datasets - MoleculeNet

295 MoleculeNet is a set of 2D-only datasets with property prediction tasks ranging from toxicity
296 prediction to drug reactivity. ROC-AUC is used in order to evaluate each model’s ability to correctly
297 determine these properties. Following [73], we report the mean and standard deviation across three
298 random seeds and use the Bemis-Murcko scaffolds recommended in DeepChem [49]. Per Table 3,
299 MOLINTERACT’s performance on MoleculeNet is competitive with its multimodal and unimodal
300 GNN-based peers, among which it ranks only behind MoleculeSDE [35] in average ROC AUC and
301 even exceeds MoleculeSDE in 5/8 datasets. A possible explanation is that MoleculeSDE is tasked
302 directly with reconstructing the original equilibrium state 3D conformer during pre-training, granting
303 it highly detailed 3D knowledge which is especially useful in determining properties which may
304 benefit from a precise understanding of the 3D geometry of a molecule like blood-brain barrier
305 permeability [55]. It is also possible that the kind of 3D information learned by MOLINTERACT
306 is not immediately useful for MoleculeNet tasks unlike quantum metrics, as shown in Table 4.
307 Regarding other baselines, MOLINTERACT primarily falls behind non-GNN-based methods like
308 GROVER [50] and MOLEBLEND [73] which use more powerful transformer backbones. We also see
309 that MOLINTERACT with $\mathcal{L}_{\text{Simple}}$ performs worse than with \mathcal{L}_{All} , which is expected given its smaller
310 ensemble of loss functions.

311 4.4 Ablation Studies

312 **3D Transfer Performance in QM9.** Table 4 shows performance on QM9 where only 2D graphs are
313 provided to MOLINTERACT in order to test the degree of information transfer. Such an application
314 may be valuable in cases where 3D structures are not consistently available, such as high-throughput

Method	α	$\Delta\mathcal{E}$	$\mathcal{E}_{\text{HOMO}}$	$\mathcal{E}_{\text{LUMO}}$	μ	C_v	R^2	ZPVE
PNA Corso et al. [7]	0.3972	123.08	82.10	85.72	0.4133	0.1670	22.14	15.08
GraphCL (You et al. [72])	0.3295	120.08	79.57	80.81	0.3937	0.1422	21.84	12.39
AttrMask (You et al. [72])	0.3570	116.21	80.58	84.93	0.4626	0.1587	29.23	25.91
GPT-GNN (Hu et al. [23])	0.3732	131.99	93.11	99.84	0.3975	0.1795	29.21	11.17
GraphMVP (Liu et al. [34])	0.3227	101.84	68.62	70.23	0.3489	0.1287	17.03	7.96
3D InfoMax (Stärk et al. [54])	0.3268	101.71	68.96	<u>69.51</u>	0.3507	0.1306	17.39	7.96
3D-PGT (Wang et al. [60])	<u>0.3121</u>	<u>101.53</u>	<u>68.24</u>	69.73	<u>0.3409</u>	<u>0.1217</u>	<u>16.89</u>	<u>7.92</u>
MOLINTERACT (2D)	0.2145	86.49	59.77	57.55	0.3055	0.0830	12.03	5.11

Table 4: Performance on QM9 using 2D-only models to study the degree of 3D information transfer.

preliminary drug screening [54, 16]. Under this restriction, MOLINTERACT outperforms several 2D baselines benchmarked by [54], demonstrating substantially stronger 3D performance given only 2D graphs, suggesting a high degree of 3D-to-2D information transfer despite such a simple suite of pre-training tasks. This emphasizes the importance of the architecture of MOLINTERACT, showing the strength of the interaction layers even when one modality is missing.

Impact of Pre-training Tasks and Interactor Types. In Table 2, we investigate the impact of the pre-training tasks in the architecture of MOLINTERACT on performance in the QM9 dataset. “Only 3D tasks” refers to the method pre-trained only on interatomic distance, bond angle, and dihedral angle prediction. “Betweenness” refers to the \mathcal{L}_{All} setting swapping centrality ranking loss with betweenness ranking loss. “Mean” and “Self-Attention Interactor” refer to the $\mathcal{L}_{\text{Simple}}$ setting except with the averaging operation and a separate single-head self-attention modules for ϕ layer, respectively. Finally, “3D structures only” refers to the setting where only 3D graphs are supplied to MOLINTERACT.

In these ablations, we observe a noticeable decline in performance across the board compared to the final version of MOLINTERACT. First, the “Only 3D tasks” ablation confirms that the 2D pre-training tasks indeed play a role in enhancing multimodal performance on downstream tasks even when they are not as directly related to properties such as $\Delta\mathcal{E}$, which are primarily reliant on 3D features. Next, the worse performance of betweenness centrality compared to eigenvector centrality suggests that the latter is more chemically meaningful. This is expected since eigenvector centrality is directly related to Laplacian eigenvector positional encodings [9, 29], which have been shown to enhance performance on molecular graphs by breaking the symmetries of WL-indistinguishable nodes [29]. The “Mean” and “Self-Attention” ablations show the superiority of a simple MLP-based interactor as a balance between a parameter-free and complex interactor. In our training runs, the self-attention-based interactor exhibited extensive over-fitting. Finally, given only 3D structures, we see that MOLINTERACT is competitive with contemporary methods, exceeds the current state-of-the-art in μ , C_v , and ZPVE, and vastly outperforms a stock SchNet on all metrics by as much as 18% ($\Delta\mathcal{E}$), suggesting that multimodal information is successfully utilized during interaction.

5 Conclusion

In this work, we introduce MOLINTERACT an architectural approach to improving multimodal self-supervised learning that leverages deep interactions to fuse 2D and 3D representations of molecules. With this deep interaction mechanism, our method is able to access fine-grained cross-modal information without sacrificing rich embeddings from modality-specific backbones, allowing for more effective interplay between 2D and 3D information when paired with even a simple set of predictive pre-training tasks, achieving new state-of-the-art performance on benchmark datasets as a result and contributing to the growing field of multimodal property prediction for small molecules.

References

- 350
- 351 [1] Phillip Bonacich. Power and centrality: A family of measures. *American Journal of Sociology*,
352 92(5):1170–1182, 1987. ISSN 00029602, 15375390.
- 353 [2] Stephen P. Borgatti. Centrality and network flow. *Social Networks*, 27(1):55–71, 2005. ISSN
354 0378-8733. doi: <https://doi.org/10.1016/j.socnet.2004.11.008>.
- 355 [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework
356 for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- 357 [4] Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta: large-scale self-
358 supervised pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885*,
359 2020.
- 360 [5] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with
361 application to face verification. In *2005 IEEE Computer Society Conference on Computer
362 Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 539–546 vol. 1, 2005. doi:
363 10.1109/CVPR.2005.202.
- 364 [6] Nurendra Choudhary, Edward W Huang, Karthik Subbian, and Chandan K Reddy. An inter-
365 pretable ensemble of graph and language models for improving search relevance in e-commerce.
366 In *Companion Proceedings of the ACM on Web Conference 2024*, pages 206–215, 2024.
- 367 [7] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal
368 neighbourhood aggregation for graph nets. In *Advances in Neural Information Processing
369 Systems*, 2020.
- 370 [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
371 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
372 Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter
373 of the Association for Computational Linguistics: Human Language Technologies, Volume 1
374 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association
375 for Computational Linguistics. doi: 10.18653/v1/N19-1423.
- 376 [9] Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs.
377 *AAAI Workshop on Deep Learning on Graphs: Methods and Applications*, 2021.
- 378 [10] Xiaomin Fang, Lihang Liu, Jieqiong Lei, Donglong He, Shanzhuo Zhang, Jingbo Zhou,
379 Fan Wang, Hua Wu, and Haifeng Wang. Geometry-enhanced molecular representation
380 learning for property prediction. *Nature Machine Intelligence*, pages 1–8, 2022. doi:
381 10.1038/s42256-021-00438-4.
- 382 [11] Linton C. Freeman. A set of measures of centrality based on betweenness. *Sociometry*, 40(1):
383 35–41, 1977. ISSN 00380431.
- 384 [12] Johannes Gasteiger, Janek Groß, and Stephan Günnemann. Directional message passing for
385 molecular graphs. In *International Conference on Learning Representations*, 2020.
- 386 [13] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl.
387 Neural message passing for quantum chemistry. In Doina Precup and Yee Whye Teh, edi-
388 tors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of
389 *Proceedings of Machine Learning Research*, pages 1263–1272. PMLR, 06–11 Aug 2017.
- 390 [14] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward
391 neural networks. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth
392 International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of
393 Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May
394 2010. PMLR.
- 395 [15] Jonathan Godwin, Michael Schaarschmidt, Alexander L Gaunt, Alvaro Sanchez-Gonzalez, Yulia
396 Rubanova, Petar Veličković, James Kirkpatrick, and Peter Battaglia. Simple GNN regularisation
397 for 3d molecular property prediction and beyond. In *International Conference on Learning
398 Representations*, 2022.
- 399 [16] Christoph Gorgulla, AkshatKumar Nigam, Matt Koop, Süleyman Selim Çınaroğlu, Christopher
400 Secker, Mohammad Haddadnia, Abhishek Kumar, Yehor Malets, Alexander Hasson, Minkai
401 Li, Ming Tang, Roni Levin-Konigsberg, Dmitry Radchenko, Aditya Kumar, Minko Gehev,
402 Pierre-Yves Aquilanti, Henry Gabb, Amr Alhossary, Gerhard Wagner, Alán Aspuru-Guzik,

- 403 Yurii S. Moroz, Konstantin Fackeldey, and Haribabu Arthanari. Virtualflow 2.0 - the next
404 generation drug discovery platform enabling adaptive screens of 69 billion molecules. *bioRxiv*,
405 2023. doi: 10.1101/2023.04.25.537981.
- 406 [17] Etash Kumar Guha, Shlok Natarajan, Thomas Möllenhoff, Mohammad Emtiyaz Khan, and Eu-
407 gene Ndiaye. Conformal prediction via regression-as-classification. In *International Conference*
408 *on Learning Representations*, 2024.
- 409 [18] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with
410 gaussian error linear units. *CoRR*, abs/1606.08415, 2016.
- 411 [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In
412 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural*
413 *Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- 414 [20] Zhenyu Hou, Xiao Liu, Yukuo Cen, Yuxiao Dong, Hongxia Yang, Chunjie Wang, and Jie Tang.
415 Graphmae: Self-supervised masked graph autoencoders. In *Proceedings of the 28th ACM*
416 *SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 594–604, 2022.
- 417 [21] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
418 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
419 *CoRR*, abs/2005.00687, 2020.
- 420 [22] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure
421 Leskovec. Strategies for pre-training graph neural networks. In *International Conference on*
422 *Learning Representations*, 2020.
- 423 [23] Ziniu Hu, Yuxiao Dong, Kuansan Wang, Kai-Wei Chang, and Yizhou Sun. Gpt-gnn: Generative
424 pre-training of graph neural networks. In *Proceedings of the 26th ACM SIGKDD Conference*
425 *on Knowledge Discovery and Data Mining*, 2020.
- 426 [24] Eric Inae, Gang Liu, and Meng Jiang. Motif-aware attribute masking for molecular graph
427 pre-training, 2023.
- 428 [25] Xiaofei He Junying Li, Deng Cai. Learning graph-level representation for drug discoveryk.
429 *arXiv preprint arXiv:1709.03741*, 2017.
- 430 [26] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
431 *Conference on Learning Representations*, San Diego, CA, USA, 2015.
- 432 [27] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *International*
433 *Conference on Learning Representations*, 2014.
- 434 [28] Johannes Klicpera, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional
435 graph neural networks for molecules. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman
436 Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- 437 [29] Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, and Prudencio Tossou.
438 Rethinking graph transformers with spectral attention. In M. Ranzato, A. Beygelzimer,
439 Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information*
440 *Processing Systems*, volume 34, pages 21618–21629. Curran Associates, Inc., 2021.
- 441 [30] Rayan Krishnan, Pranav Rajpurkar, and Eric J Topol. Self-supervised learning in medicine and
442 healthcare. *Nature Biomedical Engineering*, 6(12):1346–1352, 2022.
- 443 [31] Zhenyu Lei, Herun Wan, Wenqian Zhang, Shangbin Feng, Zilong Chen, Jundong Li, Qinghua
444 Zheng, and Minnan Luo. BIC: Twitter bot detection with text-graph interaction and semantic
445 consistency. In *Proceedings of the 61st Annual Meeting of the Association for Computational*
446 *Linguistics*, 2023.
- 447 [32] John W. Liebeschutz, Jana Hennemann, Tjelvar S. G. Olsson, and Colin R. Groom. The good,
448 the bad and the twisted: a survey of ligand geometry in protein crystal structures. *Journal of*
449 *Computer-Aided Molecular Design*, 26:169 – 183, 2012.
- 450 [33] Shengchao Liu, Mehmet F Demirel, and Yingyu Liang. N-gram graph: Simple unsupervised
451 representation for graphs, with applications to molecules. In H. Wallach, H. Larochelle,
452 A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information*
453 *Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 454 [34] Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang.
455 Pre-training molecular graph representation with 3d geometry. In *International Conference on*
456 *Learning Representations*, 2022.

- 457 [35] Shengchao Liu, Weitao Du, Zhi-Ming Ma, Hongyu Guo, and Jian Tang. A group symmetric
458 stochastic differential equation model for molecule multi-modal pretraining. In *International*
459 *Conference on Machine Learning*, pages 21497–21526. PMLR, 2023.
- 460 [36] Shengchao Liu, Hongyu Guo, and Jian Tang. Molecular geometry pretraining with SE(3)-
461 invariant denoising distance matching. In *International Conference on Learning Representations*,
462 2023.
- 463 [37] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang.
464 Self-supervised learning: Generative or contrastive. *IEEE transactions on knowledge and data*
465 *engineering*, 35(1):857–876, 2021.
- 466 [38] Yi Liu, Limei Wang, Meng Liu, Yuchao Lin, Xuan Zhang, Bora Oztekin, and Shuiwang Ji.
467 Spherical message passing for 3d molecular graphs. In *International Conference on Learning*
468 *Representations*, 2022.
- 469 [39] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In
470 *International Conference on Learning Representations*, 2017.
- 471 [40] Shengjie Luo, Tianlang Chen, Yixian Xu, Shuxin Zheng, Tie-Yan Liu, Liwei Wang, and Di He.
472 One transformer can understand both 2d & 3d molecular data. *arXiv preprint arXiv:2210.01765*,
473 2022.
- 474 [41] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation
475 and projection for dimension reduction, 2018. cite arxiv:1802.03426Comment: Reference
476 implementation available at <http://github.com/lmcinnes/umap>.
- 477 [42] Maho Nakata and Tomomi Shimazaki. Pubchemqc project: A large-scale first-principles
478 electronic structure database for data-driven chemistry. *Journal of Chemical Information and*
479 *Modeling*, 57(6):1300–1308, 2017. doi: 10.1021/acs.jcim.7b00083. PMID: 28481528.
- 480 [43] Nirmala Parisutham. How do centrality measures help to predict similarity patterns in molecular
481 chemical structural graphs? *Artificial Intelligence Chemistry*, 1(2):100007, 2023. ISSN
482 2949-7477. doi: <https://doi.org/10.1016/j.aichem.2023.100007>.
- 483 [44] Trang Pham, Truyen Tran, Khanh Hoa Dam, and Svetha Venkatesh. Graph classification via
484 deep learning with virtual nodes. *CoRR*, abs/1708.04357, 2017.
- 485 [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
486 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
487 Sutskever. Learning transferable visual models from natural language supervision, 2021.
- 488 [46] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. von Lilienfeld. Quantum chemistry structures
489 and properties of 134 kilo molecules. *Sci Data*, 1:140022, 2014.
- 490 [47] Raghunathan Ramakrishnan, Mia Hartmann, Enrico Tapavicza, and O. Anatole von Lilienfeld.
491 Electronic spectra from TDDFT and machine learning in chemical space. *The Journal of*
492 *Chemical Physics*, 143(8):084111, 08 2015. ISSN 0021-9606. doi: 10.1063/1.4928757.
- 493 [48] Ladislav Rampásek, Mikhail Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and
494 Dominique Beaini. Recipe for a General, Powerful, Scalable Graph Transformer. *Advances in*
495 *Neural Information Processing Systems*, 35, 2022.
- 496 [49] Bharath Ramsundar, Peter Eastman, Patrick Walters, Vijay Pande, Karl Leswing, and Zhenqin
497 Wu. *Deep Learning for the Life Sciences*. O’Reilly Media, 2019.
- 498 [50] Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying WEI, Wenbing Huang, and Junzhou
499 Huang. Self-supervised graph transformer on large-scale molecular data. In H. Larochelle,
500 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information*
501 *Processing Systems*, volume 33, pages 12559–12571. Curran Associates, Inc., 2020.
- 502 [51] Lars Ruddigkeit, Ruud van Deursen, Lorenz C. Blum, and Jean-Louis Reymond. Enumeration
503 of 166 billion organic small molecules in the chemical universe database gdb-17. *Journal of*
504 *Chemical Information and Modeling*, 52(11):2864–2875, 2012. doi: 10.1021/ci300415d. PMID:
505 23088335.
- 506 [52] Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Saucedo Felix, Stefan Chmiela, Alexandre
507 Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolutional neural network
508 for modeling quantum interactions. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach,
509 R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing*
510 *Systems*, volume 30. Curran Associates, Inc., 2017.

- 511 [53] Kristof Schütt, Oliver Unke, and Michael Gastegger. Equivariant message passing for the
512 prediction of tensorial properties and molecular spectra. In Marina Meila and Tong Zhang,
513 editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of
514 *Proceedings of Machine Learning Research*, pages 9377–9388. PMLR, 18–24 Jul 2021.
- 515 [54] Hannes Stärk, Dominique Beaini, Gabriele Corso, Prudencio Tossou, Christian Dallago, Stephan
516 Günnemann, and Pietro Lió. 3D infomax improves GNNs for molecular property prediction. In
517 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato,
518 editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of
519 *Proceedings of Machine Learning Research*, pages 20479–20502. PMLR, 17–23 Jul 2022.
- 520 [55] Claudia Suenderhauf, Felix Hammann, and Jörg Huwlyer. Computational prediction of blood-
521 brain barrier permeability using decision tree induction. *Molecules*, 17(9):10429–10445, 2012.
522 ISSN 1420-3049. doi: 10.3390/molecules170910429.
- 523 [56] Fan-Yun Sun, Jordan Hoffman, Vikas Verma, and Jian Tang. Infograph: Unsupervised and
524 semi-supervised graph-level representation learning via mutual information maximization. In
525 *International Conference on Learning Representations*, 2019.
- 526 [57] Philipp Thölke and Gianni De Fabritiis. Equivariant transformers for neural network based
527 molecular potentials. In *International Conference on Learning Representations*, 2022.
- 528 [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
529 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg,
530 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural
531 Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 532 [59] Limei Wang, Yi Liu, Yuchao Lin, Haoran Liu, and Shuiwang Ji. ComENet: Towards complete
533 and efficient message passing for 3d molecular graphs. In Alice H. Oh, Alekh Agarwal, Danielle
534 Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*,
535 2022.
- 536 [60] Xu Wang, Huan Zhao, Wei-wei Tu, and Quanming Yao. Automated 3d pre-training for
537 molecular property prediction. In *Proceedings of the 29th ACM SIGKDD Conference on
538 Knowledge Discovery and Data Mining*, KDD ’23, page 2419–2430, New York, NY, USA, 2023.
539 Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599252.
- 540 [61] Xu Wang, Huan Zhao, Wei-wei Tu, and Quanming Yao. Automated 3d pre-training for
541 molecular property prediction. In *Proceedings of the 29th ACM SIGKDD Conference on
542 Knowledge Discovery and Data Mining*, pages 2419–2430, 2023.
- 543 [62] Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive
544 learning of representations via graph neural networks. *Nature Machine Intelligence*, pages 1–9,
545 2022. doi: 10.1038/s42256-022-00447-x.
- 546 [63] Stanley Wasserman and Katherine Faust. *Social Network Analysis: Methods and Applications*.
547 Structural Analysis in the Social Sciences. Cambridge University Press, 1994.
- 548 [64] B.Yu. Weisfeiler and A.A. Leman. The reduction of a graph to canonical form and the algebra
549 which appears therein. *Nauchno-Tekhnicheskaya Informatsiya*, 2(9):12–16, 1968. Translation
550 from Russian to English by Grigory Ryabov.
- 551 [65] Zhenqin Wu, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S.
552 Pappu, Karl Leswing, and Vijay S. Pande. Moleculenet: a benchmark for molecular machine
553 learning. *Chemical Science*, 9:513 – 530, 2017.
- 554 [66] Jun Xia, Chengshuai Zhao, Bozhen Hu, Zhangyang Gao, Cheng Tan, Yue Liu, Siyuan Li, and
555 Stan Z. Li. Mole-BERT: Rethinking pre-training graph neural networks for molecules. In
556 *International Conference on Learning Representations*, 2023.
- 557 [67] Zhaoping Xiong, Dingyan Wang, Xiaohong Liu, Feisheng Zhong, Xiaozhe Wan, Xutong Li,
558 Zhaojun Li, Xiaomin Luo, Kaixian Chen, Hualiang Jiang, and Mingyue Zheng. Pushing the
559 boundaries of molecular representation for drug discovery with the graph attention mechanism.
560 *Journal of Medicinal Chemistry*, 63(16):8749–8760, 2020. doi: 10.1021/acs.jmedchem.9b00959.
561 PMID: 31408336.
- 562 [68] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
563 networks? In *International Conference on Learning Representations*, 2019.

- 564 [69] Minghao Xu, Hang Wang, Bingbing Ni, Hongyu Guo, and Jian Tang. Self-supervised graph-
565 level representation learning with local and global structure. In Marina Meila and Tong Zhang,
566 editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of
567 *Proceedings of Machine Learning Research*, pages 11548–11558. PMLR, 18–24 Jul 2021.
- 568 [70] Kevin Yang, Kyle Swanson, Wengong Jin, Connor Coley, Philipp Eiden, Hua Gao, Angel
569 Guzman-Perez, Timothy Hopper, Brian Kelley, Miriam Mathea, Andrew Palmer, Volker Settels,
570 Tommi Jaakkola, Klavs Jensen, and Regina Barzilay. Analyzing learned molecular representa-
571 tions for property prediction. *Journal of Chemical Information and Modeling*, 59(8):3370–3388,
572 2019. doi: 10.1021/acs.jcim.9b00237. PMID: 31361484.
- 573 [71] Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D. Manning,
574 Percy Liang, and Jure Leskovec. Deep bidirectional language-knowledge graph pretraining. In
575 *Neural Information Processing Systems*, 2022.
- 576 [72] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen.
577 Graph contrastive learning with augmentations. In H. Larochelle, M. Ranzato, R. Hadsell, M. F.
578 Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
579 pages 5812–5823. Curran Associates, Inc., 2020.
- 580 [73] Qiyang Yu, Yudi Zhang, Yuyan Ni, Shikun Feng, Yanyan Lan, Hao Zhou, and Jingjing Liu.
581 Multimodal molecular pretraining via modality blending. In *International Conference on*
582 *Learning Representations*, 2024.
- 583 [74] Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D
584 Manning, and Jure Leskovec. Greaselm: Graph reasoning enhanced language models. In
585 *International Conference on Learning Representations*, 2021.
- 586 [75] Zaixi Zhang, Qi Liu, Hao Wang, Chengqiang Lu, and Chee-Kong Lee. Motif-based graph
587 self-supervised learning for molecular property prediction. *Advances in Neural Information*
588 *Processing Systems*, 34:15870–15882, 2021.
- 589 [76] Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Linfeng
590 Zhang, and Guolin Ke. Uni-mol: A universal 3d molecular representation learning framework.
591 In *International Conference on Learning Representations*, 2023.
- 592 [77] Jinhua Zhu, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Wengang Zhou, Houqiang Li, and
593 Tie-Yan Liu. Unified 2d and 3d pre-training of molecular representations. In *Proceedings of*
594 *the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22,
595 page 2626–2636, New York, NY, USA, 2022. Association for Computing Machinery. ISBN
596 9781450393850. doi: 10.1145/3534678.3539368.
- 597 [78] Jinhua Zhu, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Wengang Zhou, Houqiang Li, and
598 Tie-Yan Liu. Unified 2d and 3d pre-training of molecular representations. In *Proceedings of the*
599 *28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 2626–2636,
600 2022.

601 **A Broader Impact**

602 This work proposes a more effective method for developing multimodal representations of molecules
603 for molecular property prediction. As such, it follows a line of work that has the potential to accelerate
604 the drug and compound discovery process, making the development of new therapeutics easier and
605 more cost-efficient. At the same time, there is potential for this work to be misused in order to aid in
606 the development of compounds which negatively impact humanity in the form of harmful drugs, for
607 example. We support the extensive usage of expert-guided control and regulation in order to steer the
608 use of this technology and similar AI-assisted drug discovery techniques for social good.

609 **B Limitations**

610 There are generally two major limitations of MOLINTERACT. First, like other multimodal approaches
611 like MOLEBLEND [73], MOLINTERACT only takes into account the geometry of single 3D molecular
612 conformations, such as those in QM9 [46] and QM8 [51] which are in equilibrium state. In this way,
613 our method may not learn a comprehensive 3D representation of molecules and the wide spectrum
614 of possible conformers which make up a valid 3D geometry for a given compound. However,
615 other works [34, 35] tackle this problem by also tasking their architecture with generating 3D
616 conformations directly, which is an SSL task which may be adopted for our architecture as well.
617 Second, MOLINTERACT is limited in that even though the multimodal representations it learns are
618 effective, it still finds optimal performance when both modalities are available to the model, suggesting
619 a slight dependence on both modalities being provided to its interaction layers for downstream task
620 performance. However, this may be remedied by experimenting with tasking the architecture with
621 reproducing 2D/3D topology/geometry similar to MoleculeSDE [35] in order to make use of the
622 deactivated unimodal branch during finetuning. Further, in the 2D-only case, MOLINTERACT already
623 demonstrates an advance in the state-of-the-art when only 2D information is provide, and when 3D
624 structures are involved, 2D structure is generally easily recoverable.

625 Figure 2 plots eigenvector centrality versus betweenness centrality on two different molecules from
626 PCQM4Mv2. While the two centrality measures are similar, eigenvector centrality is able to highlight
627 nodes which are not only towards the middle of the molecule but also parts of certain substructures,
628 such as the central ring in the top molecule or the top ring in the bottom molecule. In other words, it
629 appears eigenvector centrality is a better measure of “communities” in the molecular graph, assigning
630 nodes in substructures more similar centralities. Being able to discern which atoms have higher
631 centrality may be a useful proxy for learning higher-order structure in molecular graphs.

632 **C Hyperparameters and implementation details**

633 Hyperparameters for pre-training on PCQM4Mv2 and finetuning on QM9, QM8, and MoleculeNet
634 are in Table 5.

635 **D Pre-training computational cost**

636 Table 6 shows wall times to pre-train various baseline pre-training methods as included in the appendix
637 of MoleculeSDE [35]. Unfortunately, MOLEBLEND [73] does not yet have code available, and so
638 we could not include it in this benchmark. The wall-times for all methods besides MOLINTERACT
639 are reported from a machine using a single Nvidia V100 GPU. Due to access issues, we could
640 not attain a V100 GPU, and so our reported time is from a SLURM cluster node equipped with
641 an Nvidia A100 SXM4 80GB GPU. We recognize the generational gap in hardware, and so we
642 hypothesize that MOLINTERACT will almost certainly train slower on a V100. Reducing the number
643 of pre-training tasks will likely reduce pre-training wall time. Regrading memory requirements,
644 pre-training MOLINTERACT under our settings required at least 30GB VRAM.

645 **E Performance on QM8**

646 We also evaluate MOLINTERACT on 12 tasks from QM8 [47, 51]. QM8 is a smaller dataset than QM9
647 (20K vs 134K) with the task of predicting the electronic spectra of small organic molecules. Both 2D

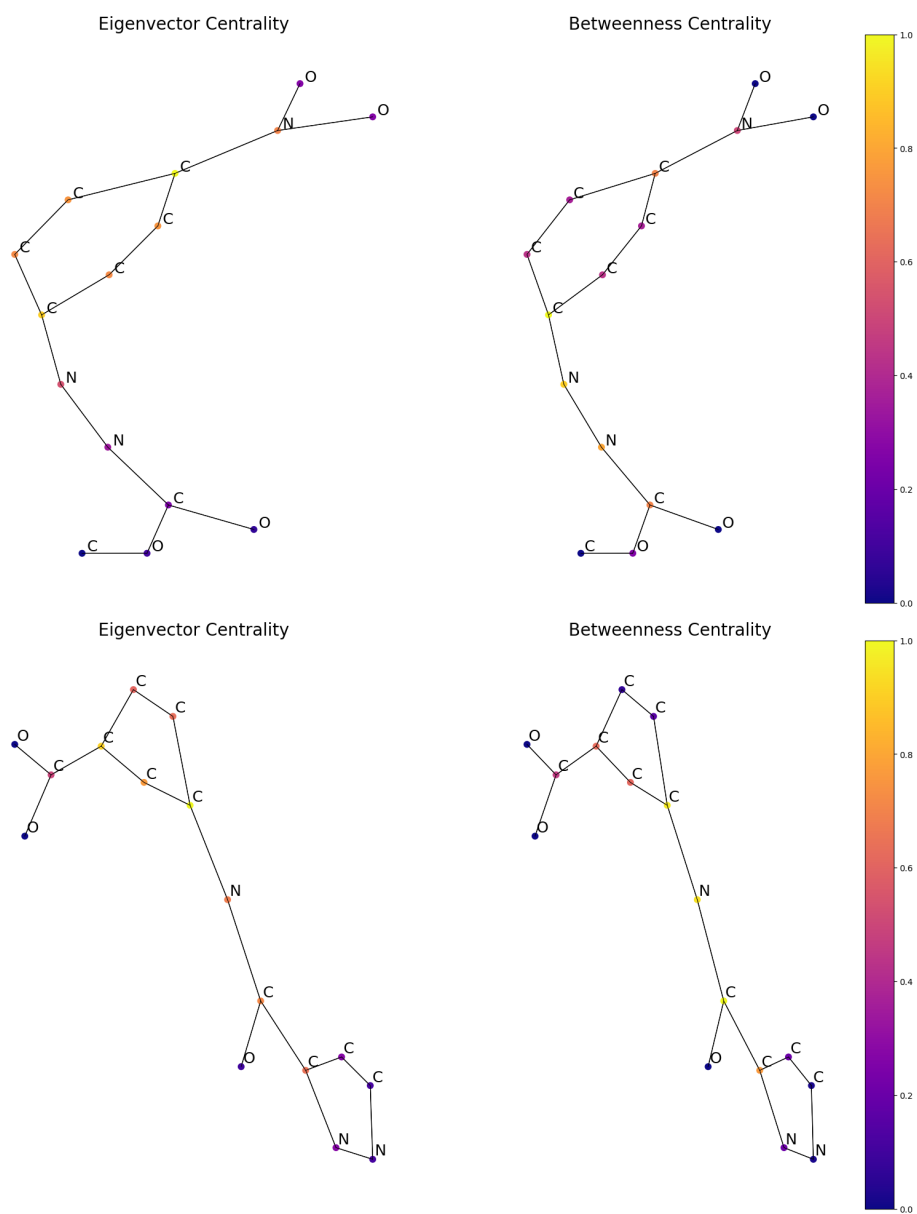


Figure 2: Comparing eigenvector and betweenness centrality on a molecule from PCQM4Mv2.

Hyperparameter	PCQM4Mv2	QM9	QM8	MoleculeNet
Optimizer	Adam [26]	Adam [26]	Adam [26]	Adam [26]
Initialization	Glorot uniform [14]	-	-	-
Learning rate scheduler	Cosine annealing [39]	Cosine annealing [39]	Cosine annealing [39]	Cosine annealing [39]
Adam betas	(0.9, 0.999)	(0.9, 0.999)	(0.9, 0.999)	(0.9, 0.999)
Batch size	1024	128	128	{32, 64, 128, 256}
Max learning rate	1e-4	1e-4	1e-4	{1e-3, 3e-4, 5e-4, 1e-5}
Min learning rate	0	0	0	0
Epochs	50	1000	40	{40, 60, 80, 100}
Weight decay	0.0	0.0	0.0	0.0
All embedding dimensions	300	300	300	300
Number of layers	8	8	8	8
Interactor activation	Swish	Swish	Swish	Swish
Interactor Batch norm	None	None	None	None
Interactor Layer norm	None	None	None	None
Number of SchNet filters	128	128	128	128
Number of SchNet Gaussians	51	51	51	51
GIN learnable ϵ	True	True	True	True
GIN Jumping knowledge	Last	Last	Last	{Last, Mean, Sum}
Dropout	0.0	0.0	0.0	{0.0, 0.1, 0.15}

Table 5: Hyperparameters for pre-training (PCQM4Mv2) and finetuning (QM9, QM8, MoleculeNet)

Pre-training algorithm	Min/epoch	GPU
AttrMask	5.5	Nvidia V100 32GB
ContextPred	14	Nvidia V100 32GB
InfoGraph	6	Nvidia V100 32GB
MolCLR	10	Nvidia V100 32GB
Distance Prediction	6.7	Nvidia V100 32GB
3D InfoGraph	7.5	Nvidia V100 32GB
3D InfoMax	8.6	Nvidia V100 32GB
GraphMVP	11	Nvidia V100 32GB
MoleculeSDE	30	Nvidia V100 32GB
MOLINTERACT (\mathcal{L}_{All})	17.8	Nvidia A100 80GB

Table 6: Wall time to pre-train MOLINTERACT compared to other pre-training algorithms.

Pre-training method	α	$\Delta\mathcal{E}$	\mathcal{E}_{HOMO}	\mathcal{E}_{LUMO}	μ	C_v	G	H	R^2	U	U_0	ZPVE
$\mathcal{L}_B + \mathcal{L}_{Cent}$	0.048	36.98	21.59	18.97	<u>0.022</u>	<u>0.026</u>	9.16	8.31	0.109	8.58	8.31	1.399
$\mathcal{L}_D + \mathcal{L}_{Cent}$	0.046	36.54	<u>20.85</u>	<u>18.33</u>	0.023	<u>0.026</u>	9.42	8.64	0.160	8.35	8.59	1.400
$\mathcal{L}_{Inter} + \mathcal{L}_{Cent}$	0.048	36.36	21.50	18.52	<u>0.022</u>	<u>0.026</u>	9.40	8.79	0.130	8.41	8.19	1.418
$\mathcal{L}_B + \mathcal{L}_{Edge}$	<u>0.047</u>	36.34	21.28	18.52	0.023	0.027	9.55	8.87	0.134	9.15	8.81	1.441
$\mathcal{L}_D + \mathcal{L}_{Edge}$	0.048	36.65	21.02	18.16	0.024	<u>0.026</u>	9.75	9.16	0.166	8.64	8.70	1.414
$\mathcal{L}_{Inter} + \mathcal{L}_{Edge}$	0.048	36.96	21.48	18.50	<u>0.022</u>	<u>0.026</u>	9.26	8.72	<u>0.102</u>	8.46	8.33	1.401
$\mathcal{L}_B + \mathcal{L}_{SPD}$	<u>0.047</u>	<u>36.32</u>	21.40	17.96	0.021	0.025	<u>9.13</u>	8.37	0.111	<u>8.36</u>	8.12	1.387
$\mathcal{L}_D + \mathcal{L}_{SPD}$	<u>0.047</u>	36.48	20.83	18.00	<u>0.022</u>	0.025	9.00	8.55	0.161	8.54	8.44	1.307
$\mathcal{L}_{Inter} + \mathcal{L}_{SPD}$ (\mathcal{L}_{Simple})	<u>0.047</u>	35.92	21.54	18.34	0.021	0.025	<u>9.13</u>	<u>8.26</u>	0.097	8.16	<u>8.17</u>	<u>1.365</u>

Table 7: Ablations of MOLINTERACT on QM9 with different combinations of 2D and 3D loss terms.

Method	Average MAE
D-MPNN (Yang et al. [70])	0.0190 \pm 0.0001
Attentive FP (Xiong et al. [67])	0.0179 \pm 0.001
N-Gram _{RF} (Liu et al. [33])	0.0236 \pm 0.0006
N-Gram _{XGB} (Liu et al. [33])	0.0215 \pm 0.0005
Pretrained GNN (Hu et al. [22])	0.0200 \pm 0.0001
GROVER _{base} (Rong et al. [50])	0.0218 \pm 0.0004
GROVER _{large} (Rong et al. [50])	0.0224 \pm 0.0003
MolCLR (Wang et al. [62])	0.0178 \pm 0.0003
ChemRL-GEM (Fang et al. [10])	0.0171 \pm 0.0001
UniMol (Zhou et al. [76])	0.0156\pm0.0001
MOLINTERACT (base)	0.0161 \pm 0.0005
MOLINTERACT ($\mathcal{L}_{\text{Simple}}$)	0.0158 \pm 0.0002
MOLINTERACT (\mathcal{L}_{All})	<u>0.0157\pm0.0002</u>

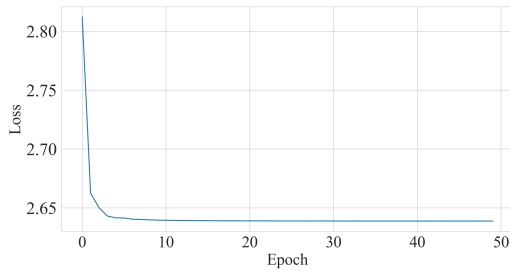
Table 8: Multi-task performance on QM8 measured in average MAE across 12 tasks. Lower is better.

648 and 3D structures are provided. Following Zhou et al. [76], we use an 80%/10%/10% scaffold split,
649 and train for only 40 epochs. We compare with baselines reported by Zhou et al. [76] and report the
650 average MAE of 12 tasks in a multi-task setting across three random seeds. Table 8 demonstrates the
651 effectiveness of MOLINTERACT, which not only outperforms pre-trained methods that leverage angle
652 information such as Fang et al. [10], but also competes with Uni-Mol, a large 3D model pre-trained
653 on over 200M molecular conformations, a dataset which is around 60 times larger and more diverse
654 than PCQM4Mv2. This shows that MOLINTERACT is able to use significantly less pre-training data,
655 which may be attributable to its utilization of both 2D and 3D information from modality-specific
656 encoders. Even when only using $\mathcal{L}_{\text{Simple}}$, MOLINTERACT achieves comparable results.

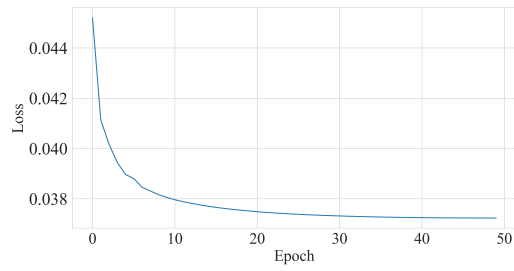
657 F Pre-training loss function behavior

658 In this section, we show loss curves for each loss function term in \mathcal{L}_{All} during pre-training on
659 PCQM4Mv2 for MOLINTERACT. In Figures 3a, 5a, and 5b, we see that lower-order quantities such
660 as interatomic distances and edge types, low loss and high accuracy are easily achieved by epoch
661 10 and begin to plateau thereafter. More complex quantities, such as bond angles and SPDs, exhibit
662 similar elbow-shaped curves but saturate more slowly as shown in Figures 3b, 6a, and 6b. Finally,
663 dihedral angle and eigenvector centrality classification are the hardest quantities to predict during
664 pre-training, with both losses and accuracies improving much more slowly per Figures 4a, 4b, 7a,
665 and 7b. This is expected given that the dihedral angle distribution in each molecule are complex
666 in comparison [32], and learning to rank nodes by eigenvector centrality distills global structural
667 patterns.

668 We also show comprehensive ablations for each combination of individual 2D and 3D pre-training
669 tasks in Table 7. We see that $\mathcal{L}_{\text{Simple}}$ performs the best overall out of each combination of tasks, with
670 $\mathcal{L}_{\text{D}} + \mathcal{L}_{\text{SPD}}$ following closely. Notably, the highest metrics usually occur for losses which include
671 \mathcal{L}_{SPD} , lending to the idea that shortest-path distances may contain the most useful 2D graph feature
672 information. This is somewhat surprising since SPDs do not include edge types, missing important
673 features such as whether an edge is a single or double bond, for example. A plausible explanation is
674 that edge information is already incorporated into the node embeddings during 2D message-passing
675 due to GINE’s edge feature-aware convolution. Meanwhile, interatomic distance and dihedral angle
676 prediction take turns as the most effective 3D tasks with bond angle regression lagging behind. While
677 all three quantities are related to the overall equilibrium state of a molecule, a possible explanation
678 for their performance difference is that interatomic distances give a more complete description of the
679 overall 3D structure of a molecule, and dihedral angles may offer more fine-grained information than
680 bond angles with more complex distributions.

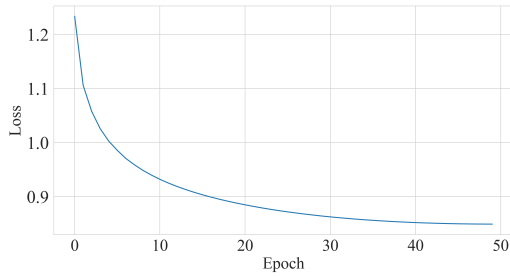


(a) Interatomic distance regression loss.

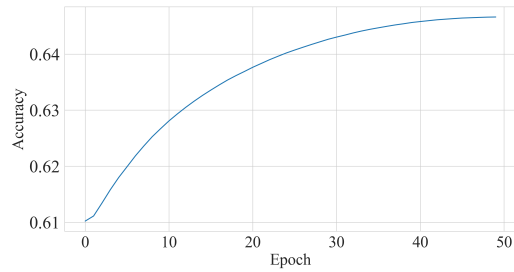


(b) Bond angle regression loss.

Figure 3: Interatomic distance and bond angle regression loss.

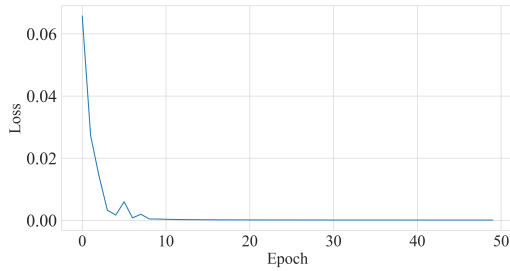


(a) Classification loss.

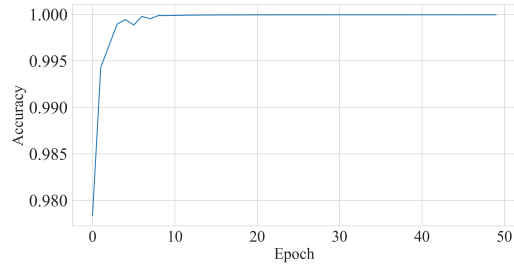


(b) Classification accuracy.

Figure 4: Dihedral angle classification loss and accuracy.

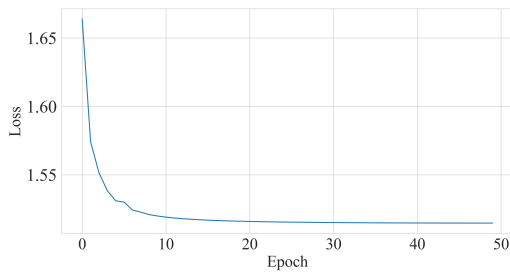


(a) Classification loss.

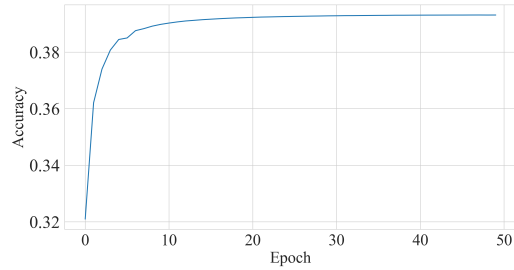


(b) Classification accuracy.

Figure 5: Edge type classification loss and accuracy.



(a) Classification loss.



(b) Classification accuracy.

Figure 6: SPD classification loss and accuracy.

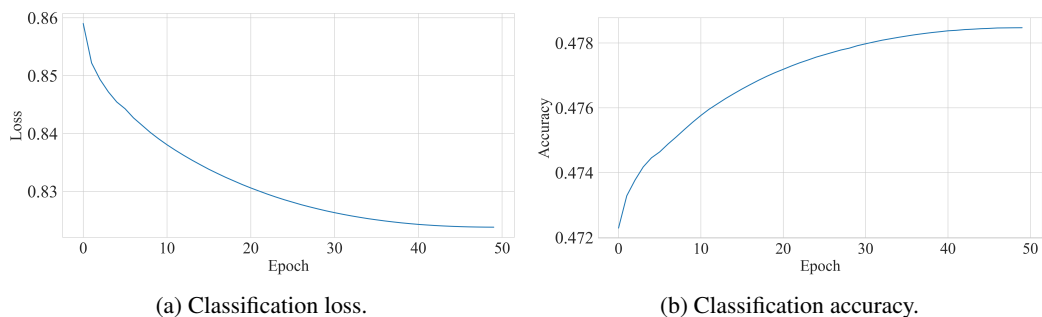


Figure 7: Centrality ranking loss and accuracy.

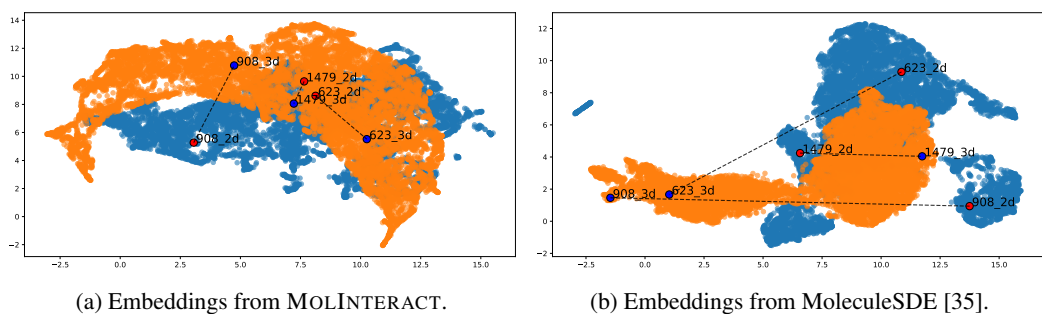


Figure 8: UMAP projection of QM9 molecule embeddings.

681 **G UMAP visualization**

682 In Figure 8a and Figure 8b, we select 3 random test molecules from QM9 and plot them on their
 683 respective UMAP [41] projections. We see that MOLINTERACT exhibits more faithful multimodal
 684 molecule representations with 2D and 3D embeddings being more closely co-located than in the
 685 embedding space for MoleculeSDE. The 2D and 3D latent spaces of MOLINTERACT are therefore
 686 more well-aligned, contributing to its effectiveness in downstream tasks.

687 **NeurIPS Paper Checklist**

688 **1. Claims**

689 Question: Do the main claims made in the abstract and introduction accurately reflect the
690 paper's contributions and scope?

691 Answer: [Yes]

692 Justification: Our work is focused on providing a new approach to multimodal molecular
693 self-supervised learning, which is reflected in the main content of the paper.

694 Guidelines:

- 695 • The answer NA means that the abstract and introduction do not include the claims
696 made in the paper.
- 697 • The abstract and/or introduction should clearly state the claims made, including the
698 contributions made in the paper and important assumptions and limitations. A No or
699 NA answer to this question will not be perceived well by the reviewers.
- 700 • The claims made should match theoretical and experimental results, and reflect how
701 much the results can be expected to generalize to other settings.
- 702 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
703 are not attained by the paper.

704 **2. Limitations**

705 Question: Does the paper discuss the limitations of the work performed by the authors?

706 Answer: [Yes]

707 Justification: Please see Appendix B for a discussion of limitations of our method.

708 Guidelines:

- 709 • The answer NA means that the paper has no limitation while the answer No means that
710 the paper has limitations, but those are not discussed in the paper.
- 711 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 712 • The paper should point out any strong assumptions and how robust the results are to
713 violations of these assumptions (e.g., independence assumptions, noiseless settings,
714 model well-specification, asymptotic approximations only holding locally). The authors
715 should reflect on how these assumptions might be violated in practice and what the
716 implications would be.
- 717 • The authors should reflect on the scope of the claims made, e.g., if the approach was
718 only tested on a few datasets or with a few runs. In general, empirical results often
719 depend on implicit assumptions, which should be articulated.
- 720 • The authors should reflect on the factors that influence the performance of the approach.
721 For example, a facial recognition algorithm may perform poorly when image resolution
722 is low or images are taken in low lighting. Or a speech-to-text system might not be
723 used reliably to provide closed captions for online lectures because it fails to handle
724 technical jargon.
- 725 • The authors should discuss the computational efficiency of the proposed algorithms
726 and how they scale with dataset size.
- 727 • If applicable, the authors should discuss possible limitations of their approach to
728 address problems of privacy and fairness.
- 729 • While the authors might fear that complete honesty about limitations might be used by
730 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
731 limitations that aren't acknowledged in the paper. The authors should use their best
732 judgment and recognize that individual actions in favor of transparency play an impor-
733 tant role in developing norms that preserve the integrity of the community. Reviewers
734 will be specifically instructed to not penalize honesty concerning limitations.

735 **3. Theory Assumptions and Proofs**

736 Question: For each theoretical result, does the paper provide the full set of assumptions and
737 a complete (and correct) proof?

738 Answer: [NA]

739 Justification: This work does not include theoretical results.

740 Guidelines:

- 741 • The answer NA means that the paper does not include theoretical results.
- 742 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 743 referenced.
- 744 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 745 • The proofs can either appear in the main paper or the supplemental material, but if
- 746 they appear in the supplemental material, the authors are encouraged to provide a short
- 747 proof sketch to provide intuition.
- 748 • Inversely, any informal proof provided in the core of the paper should be complemented
- 749 by formal proofs provided in appendix or supplemental material.
- 750 • Theorems and Lemmas that the proof relies upon should be properly referenced.

751 4. Experimental Result Reproducibility

752 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

753 perimental results of the paper to the extent that it affects the main claims and/or conclusions

754 of the paper (regardless of whether the code and data are provided or not)?

755 Answer: [Yes]

756 Justification: We report all the necessary experimental details of our method in Section 3.4

757 and Appendix C.

758 Guidelines:

- 759 • The answer NA means that the paper does not include experiments.
- 760 • If the paper includes experiments, a No answer to this question will not be perceived
- 761 well by the reviewers: Making the paper reproducible is important, regardless of
- 762 whether the code and data are provided or not.
- 763 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 764 to make their results reproducible or verifiable.
- 765 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 766 For example, if the contribution is a novel architecture, describing the architecture fully
- 767 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 768 be necessary to either make it possible for others to replicate the model with the same
- 769 dataset, or provide access to the model. In general, releasing code and data is often
- 770 one good way to accomplish this, but reproducibility can also be provided via detailed
- 771 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 772 of a large language model), releasing of a model checkpoint, or other means that are
- 773 appropriate to the research performed.
- 774 • While NeurIPS does not require releasing code, the conference does require all submis-
- 775 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 776 nature of the contribution. For example
- 777 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 778 to reproduce that algorithm.
- 779 (b) If the contribution is primarily a new model architecture, the paper should describe
- 780 the architecture clearly and fully.
- 781 (c) If the contribution is a new model (e.g., a large language model), then there should
- 782 either be a way to access this model for reproducing the results or a way to reproduce
- 783 the model (e.g., with an open-source dataset or instructions for how to construct
- 784 the dataset).
- 785 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 786 authors are welcome to describe the particular way they provide for reproducibility.
- 787 In the case of closed-source models, it may be that access to the model is limited in
- 788 some way (e.g., to registered users), but it should be possible for other researchers
- 789 to have some path to reproducing or verifying the results.

790 5. Open access to data and code

791 Question: Does the paper provide open access to the data and code, with sufficient instruc-

792 tions to faithfully reproduce the main experimental results, as described in supplemental

793 material?

794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845

Answer: [Yes]

Justification: We include all code and data and instructions to reproduce our results in our supplementary material as a zipped archive.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We list all training and test details needed to understand the results in Section 3.4 and Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We follow the existing literature on multimodal molecular SSL and report the mean and standard deviation for ROC AUC and MAE performance on MoleculeNet and QM8 from three random seeds. For QM9, we also follow the codebases from existing literature and report MAE from a single random seed (42).

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- 846 • The method for calculating the error bars should be explained (closed form formula,
847 call to a library function, bootstrap, etc.)
- 848 • The assumptions made should be given (e.g., Normally distributed errors).
- 849 • It should be clear whether the error bar is the standard deviation or the standard error
850 of the mean.
- 851 • It is OK to report 1-sigma error bars, but one should state it. The authors should
852 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
853 of Normality of errors is not verified.
- 854 • For asymmetric distributions, the authors should be careful not to show in tables or
855 figures symmetric error bars that would yield results that are out of range (e.g. negative
856 error rates).
- 857 • If error bars are reported in tables or plots, The authors should explain in the text how
858 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

860 Question: For each experiment, does the paper provide sufficient information on the com-
861 puter resources (type of compute workers, memory, time of execution) needed to reproduce
862 the experiments?

863 Answer: [Yes]

864 Justification: We include information on the amount of compute needed to run our experi-
865 ments in Appendix D.

866 Guidelines:

- 867 • The answer NA means that the paper does not include experiments.
- 868 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
869 or cloud provider, including relevant memory and storage.
- 870 • The paper should provide the amount of compute required for each of the individual
871 experimental runs as well as estimate the total compute.
- 872 • The paper should disclose whether the full research project required more compute
873 than the experiments reported in the paper (e.g., preliminary or failed experiments that
874 didn't make it into the paper).

9. Code Of Ethics

876 Question: Does the research conducted in the paper conform, in every respect, with the
877 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

878 Answer: [Yes]

879 Justification: This work did not involve human subjects, use compromising datasets, or
880 engage in behavior that breaches the Code of Ethics.

881 Guidelines:

- 882 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 883 • If the authors answer No, they should explain the special circumstances that require a
884 deviation from the Code of Ethics.
- 885 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
886 eration due to laws or regulations in their jurisdiction).

10. Broader Impacts

888 Question: Does the paper discuss both potential positive societal impacts and negative
889 societal impacts of the work performed?

890 Answer: [Yes]

891 Justification: We discuss the broader impacts of our work in Appenxdix A.

892 Guidelines:

- 893 • The answer NA means that there is no societal impact of the work performed.
- 894 • If the authors answer NA or No, they should explain why their work has no societal
895 impact or why the paper does not address societal impact.

- 896
- 897
- 898
- 899
- 900
- 901
- 902
- 903
- 904
- 905
- 906
- 907
- 908
- 909
- 910
- 911
- 912
- 913
- 914
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

915 11. Safeguards

916 Question: Does the paper describe safeguards that have been put in place for responsible
917 release of data or models that have a high risk for misuse (e.g., pretrained language models,
918 image generators, or scraped datasets)?

919 Answer: [NA]

920 Justification: The paper does not pose any risks of the kind described.

921 Guidelines:

- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

932 12. Licenses for existing assets

933 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
934 the paper, properly credited and are the license and terms of use explicitly mentioned and
935 properly respected?

936 Answer: [Yes]

937 Justification: We credit the original authors of QM9, QM8, MoleculeNet, and all competing
938 baselines.

939 Guidelines:

- 940
- 941
- 942
- 943
- 944
- 945
- 946
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 947
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 948
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 949
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 950
- 951
- 952
- 953
- 954

955 13. **New Assets**

956 Question: Are new assets introduced in the paper well documented and is the documentation
957 provided alongside the assets?

958 Answer: [\[Yes\]](#)

959 Justification: We include documented code and details for our model.

960 Guidelines:

- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968

969 14. **Crowdsourcing and Research with Human Subjects**

970 Question: For crowdsourcing experiments and research with human subjects, does the paper
971 include the full text of instructions given to participants and screenshots, if applicable, as
972 well as details about compensation (if any)?

973 Answer: [\[NA\]](#)

974 Justification: This paper did not involve crowdsourcing nor research with human subjects.

975 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 976
- 977
- 978
- 979
- 980
- 981
- 982
- 983

984 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 985 Subjects**

986 Question: Does the paper describe potential risks incurred by study participants, whether
987 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
988 approvals (or an equivalent approval/review based on the requirements of your country or
989 institution) were obtained?

990 Answer: [\[NA\]](#)

991 Justification: This paper did not involve crowdsourcing nor research with human subjects.

992 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- 993
- 994
- 995
- 996
- 997

998
999
1000
1001
1002

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.