
Learning Over Molecular Conformer Ensembles: Datasets and Benchmarks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Molecular Representation Learning (MRL) has proven impactful in numerous
2 biochemical applications such as drug discovery and enzyme design. While Graph
3 Neural Networks (GNNs) are effective at learning molecular representations from
4 a 2D molecular graph or a single 3D structure, existing works often overlook the
5 flexible nature of molecules, which continuously interconvert across conforma-
6 tions via chemical bond rotations and minor vibrational perturbations. To better
7 account for molecular flexibility, some recent works formulate MRL as an en-
8 semble learning problem, focusing on explicitly learning from a set of conformer
9 structures. However, most of these studies have limited datasets, tasks, and models.
10 In this work, we introduce the first Molecular AR Conformer Ensemble Learning
11 (MARCEL) benchmark to thoroughly evaluate the potential of learning on con-
12 former ensembles and suggest promising research directions. MARCEL includes
13 four datasets covering diverse molecule- and reaction-level properties of chemically
14 diverse molecules including organocatalysts and transition-metal catalysts, extend-
15 ing beyond the scope of common GNN benchmarks that are confined to drug-like
16 molecules. In addition, we conduct a comprehensive empirical study, which bench-
17 marks representative 1D, 2D, and 3D molecular representation learning models,
18 along with two strategies that explicitly incorporate conformer ensembles into 3D
19 MRL models. Our findings reveal that direct learning from an accessible conformer
20 space can improve performance on a variety of tasks and models.

21 1 Introduction

22 Recent years have seen the emergence of Molecular Representation Learning (MRL) as a promising
23 approach for modeling molecules with machine learning. In the typical formulation, MRL maps
24 discrete molecular objects to continuous features in a data-driven manner, encoding complex chemical
25 structures into representation vectors that can subsequently be utilized in different downstream tasks.
26 In particular, MRL now underpins a variety of biochemical applications spanning molecular property
27 prediction to the design of novel drug candidates [1–3].

28 Traditional approaches often encode chemical compounds with fingerprints, such as extended-
29 connectivity fingerprints [4, 5], which indicate the existence of certain substructures as binary bits in a
30 fixed-length sequence. Such line-based representations are concise and efficient, but have limited ex-
31 pressive power and have difficulty in capturing 3D structural information such as bonding geometries
32 and global shapes, which can be important for analyzing molecular properties and chemical reactiv-
33 ity [6, 7]. Recently, Graph Neural Networks (GNNs) have become an increasingly popular method of
34 learning molecular representations by treating molecules as graph-structured objects. Existing GNN
35 models for MRL can be broadly classified into two categories: 2D topological models [8–11] and
36 3D geometric models [12–17]. 2D GNNs typically model the molecular connectivity as a flat 2D
37 graph with atoms as nodes and bonds as edges, learning representations of chemical environments
38 by iteratively passing messages between neighboring atoms. Although powerful in the absence of

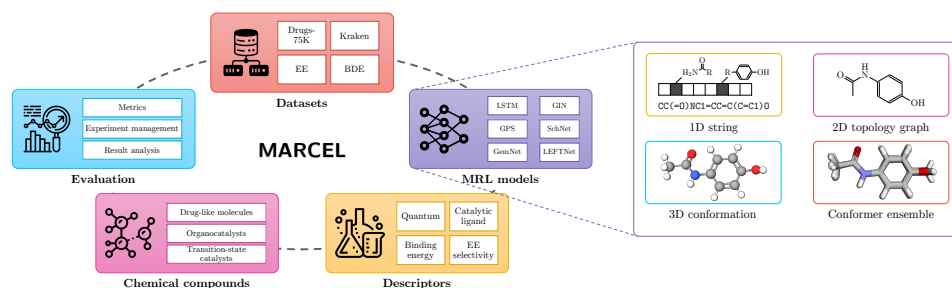


Figure 1: We present a MARCEL benchmark that comprehensively evaluates the potential of learning on conformer ensembles across a diverse set of molecules, datasets, and models.

39 structural information, 2D GNNs may fail to capture key conformational effects or stereochemical
 40 properties like chirality [18, 19], which is critical for modeling molecular interactions in areas such as
 41 drug design or chemical catalysis. Conversely, 3D GNNs are designed to model molecular conformers
 42 (conformations), which describe the structure of molecules in 3D space. Thus, these models have
 43 found widespread adoption for modeling electronic properties, predicting conformer energies and
 44 forces, and scoring interactions between ligands and proteins, amongst other applications.

45 In almost all applications, benchmarks, and demonstrations, 3D GNN models focus on encoding
 46 *individual* conformer structures. It is critical to recognize that in reality molecules are not rigid,
 47 static objects; rather, thermodynamically-permissible rotations of chemical bonds, small vibrational
 48 motions, and dynamic intermolecular interactions cause molecules to continuously convert between
 49 different conformations [20]. As a consequence, many experimentally observable chemical properties
 50 depend on the full distribution of thermodynamically-accessible conformers. For example, a molecule
 51 needs to be arranged into a particular pose to bind to a target protein, and this binding conformation
 52 changes depending on the dynamic interaction between the molecule and the target [21]. Also, it is
 53 often challenging to determine *a priori* the conformers that predominantly contribute to molecular
 54 properties without doing prohibitively expensive simulations. Therefore, a natural question arises:
 55 can we leverage the *collective* power of many different conformer structures lying on the local minima
 56 of the potential energy surface, also known as the *conformer ensemble*, to improve MRL models?

57 As shown by the empirical evidence from various studies, learning from an explicit conformer
 58 ensemble can prove to be advantageous for many tasks, including property and energy prediction [22–
 59 24], key conformer pose identification [25], and RNA sequence design [26]. However, these studies
 60 have been mostly confined to small-scale datasets, a limited set of tasks, and a restricted set of model
 61 architectures. As a result, it remains unclear (1) to what extent 2D GNNs can implicitly model
 62 molecular flexibility and (2) whether the *explicit* encoding of conformer ensembles can improve the
 63 performance of 3D models that traditionally encode only one single conformer.

64 In this paper, we present the first Molecular AR Conformer Ensemble Learning (MARCEL) benchmark.
 65 As shown in Figure 1, MARCEL covers a diverse range of chemical space, which focuses on four
 66 chemically-relevant tasks for both molecules and reactions, with an emphasis on Boltzmann-averaged
 67 properties of conformer ensembles computed at the Density-Functional Theory (DFT) level. Our
 68 datasets encompass a variety of compounds with high-quality conformers, including organocatalysts
 69 and transition-metal catalysts, extending beyond the scope of conventional GNN benchmarks which
 70 are often restricted to drug-like molecules. Moreover, we implement a comprehensive benchmark
 71 suite that enables extensive empirical studies across representative 1D, 2D, and 3D MRL models. We
 72 further explore the advantages of leveraging conformer ensembles through two straightforward strate-
 73 gies: (1) augmenting training samples by randomly selecting one conformer from the ensemble for
 74 each molecule and (2) applying an explicit multi-instance ensemble learning layer, which aggregates
 75 individual conformer embeddings.

76 Our experimental results confirm the potential effectiveness of incorporating conformer ensembles in
 77 MRL, highlighting the improvements over conventional single-conformation 3D networks. However,
 78 it is important to understand the heterogeneity of outcomes based on different dataset characteristics,
 79 task objectives, and model choices. Our investigation yields three key findings: (1) Leveraging
 80 molecular conformers by incorporating explicit set encoders, as a part of conformer ensemble learning
 81 strategies, can improve single-conformer 3D MRL models performance. (2) Data augmentation
 82 through conformer sampling may offer potential benefits, evidenced by improved results in the BDE
 83 dataset, suggesting a method to increase model robustness against imprecise structures. (3) Model
 84 selection for MRL depends on dataset sizes and tasks, with traditional 1D fingerprints and 2D models
 85 preferred for smaller datasets and 3D models for larger or reaction-focused tasks.

86 2 Problem Formulation

87 We represent a 2D molecular graph as a tuple $G = (\mathcal{V}, \mathcal{E}, \mathbf{X}, \mathbf{W})$, where $\mathcal{V} = \{v_i\}_{i=1}^{|\mathcal{V}|}$ is the node
88 set with each node corresponding to an atom, and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the edge set representing chemical
89 bonds as edges between nodes. Further, $\mathbf{X} \in \mathbb{R}^{d_v \times |\mathcal{V}|}$ contains vector attributes for each node, and
90 $\mathbf{W} \in \mathbb{R}^{d_w \times |\mathcal{E}|}$ contains attributes for each edge. When modeling chemical reactions, we represent a
91 molecule-molecule complex as a pair of graphs (G_1, G_2) . In this case, the conformation describes the
92 combined structure of the interacting molecules. For a given molecule or molecular complex, we
93 assume that its geometry can be effectively characterized by a representative set of discrete, sampled
94 conformers from the thermodynamically-accessible conformer distribution. Formally, this set can
95 be denoted as $\mathcal{C} = \{C_i\}_{i=1}^{|\mathcal{C}|}$, where $C_i \in \mathbb{R}^{|\mathcal{V}| \times 3}$ represents one conformer structure in 3D space.
96 In reality, the conformer distribution is continuous; \mathcal{C} in our study contains representative samples
97 of the infinite set. Each conformer in the sampled ensemble is associated with a statistical weight
98 given by $p_i = \frac{\exp\left(-\frac{e_i}{k_B T}\right)}{\sum_j \exp\left(-\frac{e_j}{k_B T}\right)}$, which corresponds to its probability under experimental conditions.
99 Here, e_i is the energy of the conformer C_i , k_B is the Boltzmann constant, and T is the temperature.
100 Notably, p_i is not prior information to the models analyzed in this benchmark. Rather, we use a
101 discrete approximation of p_i to compute the ground-truth labels for our regression tasks.

102 3 Datasets and Tasks

103 **MARCEL** contains four small-to-large-scale datasets involving nine regression tasks with consider-
104 ably diverse chemistry. **Drugs-75K** and **Kraken** focus on molecular properties, while **EE** and **BDE**
105 focus on reaction-centric properties. **MARCEL** includes molecules with high structural flexibility,
106 evidenced by an average number of rotatable bonds exceeding 5. Table 1 summarizes the datasets.

107 **Drugs-75K** is a subset of the **GEOM-Drugs** [27] dataset, which includes 75,099 molecules with
108 at least 5 rotatable bonds. For each molecule, we focus on three important quantum chemical
109 descriptors: ionization potential, electron affinity, and electronegativity [28]. The tasks are to predict
110 the Boltzmann-averaged value of each property across the conformer ensemble $\langle y \rangle_{k_B} = \sum_{C_i \in \mathcal{C}} p_i y_i$,
111 where y_i is a conformer-specific property. We are given each C_i , and the goal is to predict $\langle y \rangle_{k_B}$
112 from the molecular graph G , a single conformer $C_i \in \mathcal{C}$, or the set \mathcal{C} .

113 **Kraken** [29] is a dataset of 1,552 monodentate organophosphorus (III) ligands along with their
114 DFT-computed conformer ensembles. In this study, we consider four 3D ligand descriptors exhibiting
115 significant variance among conformers: **Sterimol B₅**, **Sterimol L**, **buried Sterimol B₅**, and **buried**
116 **Sterimol L**. These descriptors quantify the steric features of each ligand in units of Å and are often
117 employed for Quantitative Structure-Activity Relationship (QSAR) modeling in catalysis design.

118 As in the **Drugs-75K** tasks, the goal is to predict the Boltzmann-averaged value of each property
119 across the conformer ensemble from the molecular graph G , a single conformer $C_i \in \mathcal{C}$, or the set \mathcal{C} .

120 **EE** [30] is a dataset of 872 catalyst-substrate pairs involving 253 Rh-bound atropisomeric catalysts
121 derived from chiral bisphosphine, with 10 enamides as substrates. The dataset includes conformations
122 of catalyst-substrate transition state complexes in two separate pro-S and pro-R configurations. The
123 task is to predict the Enantiomeric Excess (EE) of the chemical reaction involving the substrate.
124 Unlike properties in **Drugs-75K** and **Kraken**, EE depends on the conformer ensembles of *each* pro-R
125 and pro-S complex. The goal is to predict EE from the graphs of the catalyst and substrate $(G_{\text{cat}}, G_{\text{sub}})$,
126 a conformer $C_i^{(R)} \in \mathcal{C}^{(R)}$ and $C_i^{(S)} \in \mathcal{C}^{(S)}$ for each complex, or the ensembles $\mathcal{C}^{(R)}$ and $\mathcal{C}^{(S)}$.

127 **BDE** [31] is a dataset containing 5,915 organometallic catalysts ML_1L_2 consisting of a metal center
128 coordinated to two flexible organic ligands. The data includes conformations of each unbound catalyst,
129 as well as conformations of the catalyst when bound to ethylene and bromide after oxidative addition
130 with vinyl bromide. Each catalyst has an electronic binding energy to be predicted. Although the
131 binding energies are computed via DFT, the conformers provided for modeling are initially generated
132 with **Open Babel** [32] followed by further geometry optimization, which ensures that the 3D structures
133 are likely to be the global minimum energy conformers at the force field level [31]. This dataset
134 realistically represents the setting in which precise conformer ensembles are unknown at inference.
135 The task is to predict the binding energy from the graphs of the unbound and bound catalyst, sampled
136 conformers $C_i^{(\text{unbound})} \in \mathcal{C}^{(\text{unbound})}$ and $C_i^{(\text{bound})} \in \mathcal{C}^{(\text{bound})}$, or the ensembles $\mathcal{C}^{(\text{unbound})}$ and $\mathcal{C}^{(\text{bound})}$.

Table 1: Statistics of the four datasets. The numbers of heavy atoms and rotatable bonds (“rot. bonds”) are averaged per conformer.

Dataset	# Molecules	# Conformers	# Heavy atoms	# Rot. bonds	# Targets	Atomic species
Drugs-75K	75,099	558,002	30.56	7.53	3	H, C, N, O, F, Si, P, S, Cl
Kraken	1,552	21,287	23.70	9.05	4	H, B, C, N, O, F, Si, P, S, Cl, Fe, Se, Br, Sn, I
Dataset	# Reactions	# Conformers	# Heavy atoms	# Rot. bonds	# Targets	Atomic species
EE	872	Pro-R: 14,807 Pro-S: 13,999	59.32	18.57	1	H, C, N, O, F, P, Cl, Br, Rh
BDE	5,915	Ligand: 73,834 Complex: 40,264	29.62 32.38	6.99 6.99	1	H, C, N, O, F, P, Cl, Ni, Cu, Br, Pd, Ag, Pt, Au

137 **Dataset Preparation.** We implement several preprocessing steps to ensure the quality and validity of
138 our datasets and facilitate their integration into machine learning models.

- 139 • **Conformer deduplication.** To eliminate redundant conformers in each ensemble \mathcal{C} , we first
140 align every pair of conformers using RDKit [33], accounting for symmetric atom permutations.
141 Subsequently, we employ Butina clustering [34] based on the Root Mean Square Deviation (RMSD)
142 values derived from conformer alignment. Within each cluster, we select the conformer with the
143 lowest energy. Note that Boltzmann-averaged regression labels are computed *before* deduplication.
- 144 • **Selection of molecules.** We focus on modeling flexible molecules, for which conformer ensemble
145 learning may be relevant to capture their properties. Hence, we only retain molecules with more
146 than 5 rotatable bonds. We also remove molecules with missing 3D geometries or 2D graphs.

147 4 Benchmarking Molecular Representation Learning Models

148 The representation of molecular data is crucial for applying machine learning models to problems in
149 chemistry and biology. These representations typically include 1D strings, 2D topological graphs,
150 and 3D geometric graphs. For a comprehensive benchmark for MRL models, our MARCEL includes
151 a diverse representative selection of models for each of the aforementioned molecular representations.
152 In this section, we provide an overview of these models and describe how they are tailored to our
153 tasks. We also introduce two strategies of explicitly encoding conformer ensembles using 3D models.

154 4.1 1D Models

155 Our 1D baselines include Random Forest [35] models operating on molecular fingerprints [33, 36,
156 37]. Fingerprints convert a molecular graph into a bit array indicating the presence of chemical
157 substructures and are widely used for cheminformatics and QSAR modeling in the low-data regime.
158 Additionally, we include Long Short-Term Memory (LSTM) [38] and Transformer [39] models,
159 popular sequence-based neural network architectures, operating on SMILES strings. For the BDE and
160 EE datasets, we concatenate the SMILES of each molecule or complex with a “.” symbol and use a
161 single sequence encoder. Further details on model implementations can be found in Appendix B.1.

162 4.2 2D Graph Neural Networks

163 We employ four widely-used GNN models as 2D baseline methods, including Graph Isomorphism
164 Network (GIN) [40], GIN with Virtual Node (GIN-VN) [41], ChemProp [42], and GraphGPS [43].
165 Following OGB protocols [41], we employ a diverse set of atomic features such as aromaticity and
166 hybridization for nodes, as well as bond features like ring information for edges (Appendix B.2). For
167 the EE and BDE datasets, we employ a two-tower architecture with two separate 2D GNN models:
168 for EE, since both pro-S and pro-R complexes share the same 2D graph, we leverage two separate
169 GNNs to encode the catalyst and substrate; for BDE, we also encode the unbound and bound catalysts
170 separately. We then concatenate these together to obtain the system-level embeddings.

171 4.3 3D Graph Neural Networks

172 We include six representative 3D GNNs that encompass diverse modeling perspectives. For invariant
173 networks, our experiments involve SchNet [12], DimeNet++ [13], and GemNet [14]. For equivariant
174 networks, we include PaiNN [15], ClofNet [16], and LEFTNet [17].

175 We use atom types as the sole atom features for the 3D models. For both training and inference on
176 Drug-75K, Kraken, and EE datasets, all the single-conformer 3D models encode the lowest-energy
177 conformer of each conformer ensemble, which has the largest Boltzmann weight and hence provides
178 the strongest model. Since imprecise conformers are encoded for the BDE task, we use a fixed,
179 randomly sampled conformer for each unbound- and bound-catalyst during training and inference.

180 The 3D models also employ a two-tower architecture for the EE and BDE datasets. Two separate
181 3D GNNs are used to encode representations for each pro-S and pro-R complex in EE, or for each
182 catalyst and bound complex in BDE, which are then concatenated to form the final representations.

183 We note that although using the lowest-energy conformer will yield the strongest performance, this
184 setting can be unrealistic: it is often not possible to identify the lowest energy conformer without
185 searching the entire conformer space. The lowest energy conformer can also depend on the force
186 field used for geometry optimization, which may neglect experimental conditions such as solvents.

187 4.4 Incorporating Conformer Ensembles into Molecular Representations

188 3D geometric models primarily focus on learning representations from individual 3D structures.
189 Although some models preserve global symmetries such as SE(3)-equivariance, these models do
190 not learn representations that capture conformational flexibility which is caused by internal degrees
191 of freedom such as bond rotations. Here, we describe two straightforward strategies that model
192 conformational flexibility by explicitly leveraging conformer ensembles.

193 4.4.1 Strategy 1: Training-Time Data Augmentation via Conformer Sampling

194 A direct approach to modeling conformer flexibility is to simply enrich the training data by randomly
195 sampling a conformer from the ensemble during each training epoch. Formally, for a given molecule
196 G and its conformer ensemble \mathcal{C} , we randomly select a conformer with uniform probability $p = 1/|\mathcal{C}|$
197 while using the same training label for each conformer. Note that during inference, the lowest-energy
198 conformer is used to evaluate the model performance. This strategy aligns with learning represen-
199 tations invariant to conformational changes, thus implicitly encoding the flexibility of molecular
200 structures, and has been shown to be useful for learning chirality-sensitive 3D representations [19].
201 When conformer ensembles are available, the strategy is computationally efficient as it maintains
202 the same complexity as the base 3D model. Unlike the other ensemble methods, this strategy can be
203 used if conformer ensembles are only available at training time. In Appendix C, we evaluate two
204 alternative scenarios where conformer ensembles are also available during evaluation.

205 4.4.2 Strategy 2: Ensemble Learning with Explicit Set Encoders

206 The second strategy utilizes a set encoder to simultaneously encode the entire conformer ensemble \mathcal{C}
207 at both training and inference time. Inspired by the multi-instance learning framework [44–46], this
208 strategy first employs 3D GNNs to generate individual conformer embeddings and then aggregates
209 these embeddings using a set encoder, as illustrated in Figure 2.

210 Formally, for each conformer $C_i \in \mathcal{C}$, we obtain its corresponding embedding $z_i = f(G, C_i) \in$
211 \mathbb{R}^d , where f is a single-conformer 3D model and d is the embedding dimension. Note that the
212 embedding z is a (3D) graph-level representation resulting from a pooling function over the node-
213 level embeddings after message passing. To further aggregate these embeddings $\mathcal{Z} = \{z_i\}_{i=1}^{|\mathcal{C}|}$ into a
214 single molecular representation, we consider the following three set encoders:

- 215 • **Mean pooling** simply computes the mean of all the conformer embeddings.
- 216 • **DeepSets** [47] utilizes a permutation-invariant function to process a set of inputs. It first applies a
217 MultiLayer Perceptron (MLP) h to each conformer embedding and then aggregates the transformed
218 embeddings using sum pooling followed by another MLP g :

$$s^{\text{DS}} = g \left(\sum_{i=1}^{|\mathcal{C}|} h(z_i) \right). \quad (1)$$

219 This method retains more discernible information from individual embeddings compared to mean
220 pooling at a cost of two non-linear functions.

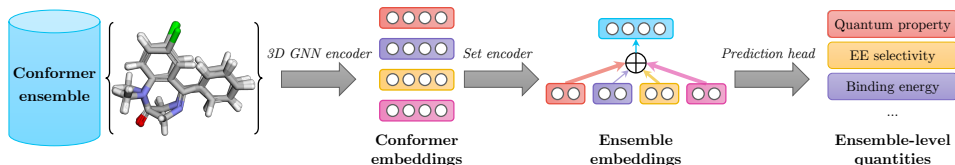


Figure 2: Conformer ensemble learning with explicit set encoders (Strategy 2). Individual conformer embeddings are first obtained via 3D GNN encoders. Then, a set encoder is employed to aggregate conformer embeddings. Finally, a linear projection head is used to generate the prediction.

- 221 • **Self-attention** [48] further computes a weighted sum of the embeddings, where the weights are
 222 obtained by applying a softmax function to the dot product of the embeddings:

$$s^{\text{ATT}} = \sum_{i=1}^{|\mathcal{C}|} c_i, \quad \text{where } c_i = g \left(\sum_{j=1}^{|\mathcal{C}|} \alpha_{ij} h(z_j) \right), \quad \alpha_{ij} = \frac{\exp((\mathbf{W}h(z_i))^{\top} (\mathbf{W}h(z_j)))}{\sum_{k=1}^{|\mathcal{C}|} \exp((\mathbf{W}h(z_i))^{\top} (\mathbf{W}h(z_k)))}. \quad (2)$$

223 Here, $\mathbf{W} \in \mathbb{R}^{d \times d}$ is a learnable weight matrix. This approach can capture conformer interactions.

224 By employing these set encoders, we can learn a model that is more sensitive to the full range of
 225 conformer variations present in the ensemble. After obtaining the ensemble embeddings, we further
 226 apply a linear projection head to generate the final prediction.

227 5 Experiments

228 5.1 Experimental Configurations

229 Each dataset is partitioned randomly into three subsets: 70% for training, 10% for validation, and
 230 20% for test. Each model is trained over 2,000 epochs using the Adam optimizer [49] with early
 231 stopping triggered if there is no improvement on the training loss over 200 epochs. For all nine
 232 regression targets, experiments are repeated three times, and the results reported correspond to the
 233 model that performs best on the validation set in terms of Mean Absolute Error (MAE).

234 The Boltzmann-averaged targets are computed over all available conformers. For ensemble learning
 235 models, we cap the number of encoded conformers per molecule to a maximum of 20, which
 236 empirically improves training stability and leads to better performance. To ensure a fair comparison,
 237 the hidden dimension size is uniformly set to 128 for all models. Other settings mostly follow the
 238 original configurations as described in the respective papers. We specify all hyperparameters and
 239 describe experimental environments in Appendix B.3.

240 5.2 Results and Analysis

241 We summarize the performance of the 1D, 2D, and 3D MRL models in Table 2. Figure 3 reports the
 242 *performance changes* in Mean Absolute Error (MAE) for each 3D model when applying the ensemble
 243 learning strategies. The raw performance data with standard deviation and the parameter size of each
 244 model can be found in Appendix D. In summary, although performance varies across the datasets,
 245 tasks, and models, the ensemble learning strategies improve upon 3D models that only encode
 246 one conformer in 48 out of 54 experiments across 9 tasks and 6 base models, demonstrating the
 247 effectiveness of conformer ensemble learning. Our analysis leads to the following key observations.

248 **Observation 1. The conformer ensemble learning strategy with explicit set encoders frequently
 249 yields improved performance.**

250 Figure 3 indicates that encoding conformer ensembles can substantially reduce test error, achieving
 251 improvements in 108 experiments across all 9 tasks, 6 base models, and 3 set encoders, most notably
 252 on the tasks in the smaller-sized Kraken dataset. This, however, does not always extend to larger
 253 datasets like Drugs-75K. We conjecture that for Drugs-75K, the computational burden of encoding all
 254 conformers in each ensemble alters the learning dynamics of the underlying model, making training
 255 more challenging. A similar finding was reported by Axelrod and Gómez-Bombarelli [23].

256 Among the three set encoders, DeepSets consistently demonstrates significant improvements in 42 out of
 257 54 experiments across 9 tasks and 6 base 3D models. We conjecture that this superior performance

Table 2: Performance of 1D, 2D, and 3D baseline MRL models and the best results from ensemble learning strategies on 3D GNNs. The metric used is the Mean Absolute Error (MAE, \downarrow). The **bold** indicates the best-performing model, while underlined denotes the second-best.

Category	Model	Drugs-75K			Kraken				EE	BDE
		IP	EA	χ	B ₅	L	BurB ₅	BurL		
1D	Random forest	0.4987	0.4747	0.2732	0.4760	0.4303	0.2758	0.1521	61.2963	3.0335
	LSTM	0.4788	0.4648	0.2505	0.4879	0.5142	0.2813	0.1924	64.0088	2.8279
	Transformer	0.6617	0.5850	0.4073	0.9611	0.8389	0.4929	0.2781	62.0816	10.0771
2D	GIN	0.4354	0.4169	0.2260	0.3128	0.4003	0.1719	0.1200	62.3065	2.6368
	GIN+VN	0.4361	0.4169	0.2267	0.3567	0.4344	0.2422	0.1741	62.3815	2.7417
	ChemProp	0.4595	0.4417	0.2441	0.4850	0.5452	0.3002	0.1948	61.0336	2.6616
	GraphGPS	0.4351	0.4085	0.2212	0.3450	0.4363	0.2066	0.1500	61.6251	2.4827
3D	SchNet	0.4394	0.4207	0.2243	0.3293	0.5458	0.2295	0.1861	17.7421	2.5488
	DimeNet++	0.4441	0.4233	0.2436	0.3510	0.4174	0.2097	0.1526	14.6414	1.4503
	GemNet	0.4069	<u>0.3922</u>	0.1970	0.2789	0.3754	0.1782	0.1635	18.0338	1.6530
	PaiNN	0.4505	0.4495	0.2324	0.3443	0.4471	0.2395	0.1673	20.2359	2.1261
	ClofNet	0.4393	0.4251	0.2378	0.4873	0.6417	0.2884	0.2529	33.9473	2.6057
	LEFTNet	0.4174	0.3964	0.2083	0.3072	0.4493	0.2176	0.1486	19.7974	1.5328

is due to its ability of effectively modeling set objects at a relatively minor computational overhead of two non-linear transformations. On the other hand, the simple mean pooling approach loses discriminative power across the conformers in the ensemble, resulting in inferior performance. It is also noteworthy that the attention models exhibit mixed results compared to the vanilla 3D models, despite theoretically being the most powerful set encoders. This inconsistency might be attributable to the computational intricacy of the self-attention layer, which models the pairwise relationship among conformers in each ensemble and hence could require more sophisticated training strategies. Future research should consider developing better neural architectures that are specifically designed to more efficiently encode structural information from conformer ensembles.

Observation 2. Sampling conformers at training time can improve performance, especially on imprecise conformer structures.

We observe that data augmentation improves performance on 34 experiments, especially on the challenging BDE dataset, where the other ensemble learning strategies often do not help. Note that the conformers in the BDE dataset originate from Open Babel, as opposed to the golden-standard DFT-level conformers present in other datasets. This suggests that training with randomly sampled conformers might offer robustness to noise in the imprecise structures. On other tasks, randomly sampling the conformers at each epoch may help the model learn an invariance to conformational changes, but does not always increase performance for all 3D models. This might be because the sampling probability is uniform across the entire conformer set, which does not respect the underlying Boltzmann weight of each conformer. In future work, it may be worthwhile to investigate whether more physics-informed sampling strategies could lead to more consistent performance gains.

Observation 3. No model consistently outperforms the rest, with substantial task dependencies.

The results in Table 2 suggest that no single model outperforms the others across all tasks. Of the 1D models, LSTM outperforms Random Forest and Transformer models on Drugs-75K and BDE, demonstrating the effectiveness of SMILES-based representations of molecules on large-scale datasets. For small datasets such as Kraken and EE, Random Forests outperform sequence models at a lower computational cost, indicating that traditional models are superior in the low-data regime.

Amongst 2D models, GIN delivers the best performance on four tasks compared to all other models; GraphGPS also demonstrates strong performance on several tasks (B₅, L, and BurL). Surprisingly, the 2D models are also competitive with some 3D models on the large-scale Drugs-75K tasks. This is possibly due to the fact that the electronic properties in Drugs-75K are not as sensitive to conformational changes, thus explicitly modeling the structures in 3D may not be necessary. However, all 2D models perform worse as compared to the 3D models in the reaction datasets EE and BDE, indicating the important role of spatial interactions in determining reaction-related properties.

For 3D models, GemNet and LEFTNet excel in IP, EA, and χ . The complexity of these two equivariant models may especially benefit from the large dataset size of Drugs-75K. For Kraken and the two reaction datasets, DimeNet++ — an invariant model — achieves promising performance, suggesting that highly-complex 3D models may be less useful for chemical applications with small-to-medium sized datasets. On EE, we observe that 3D models remarkably outperform 1D and 2D models, likely

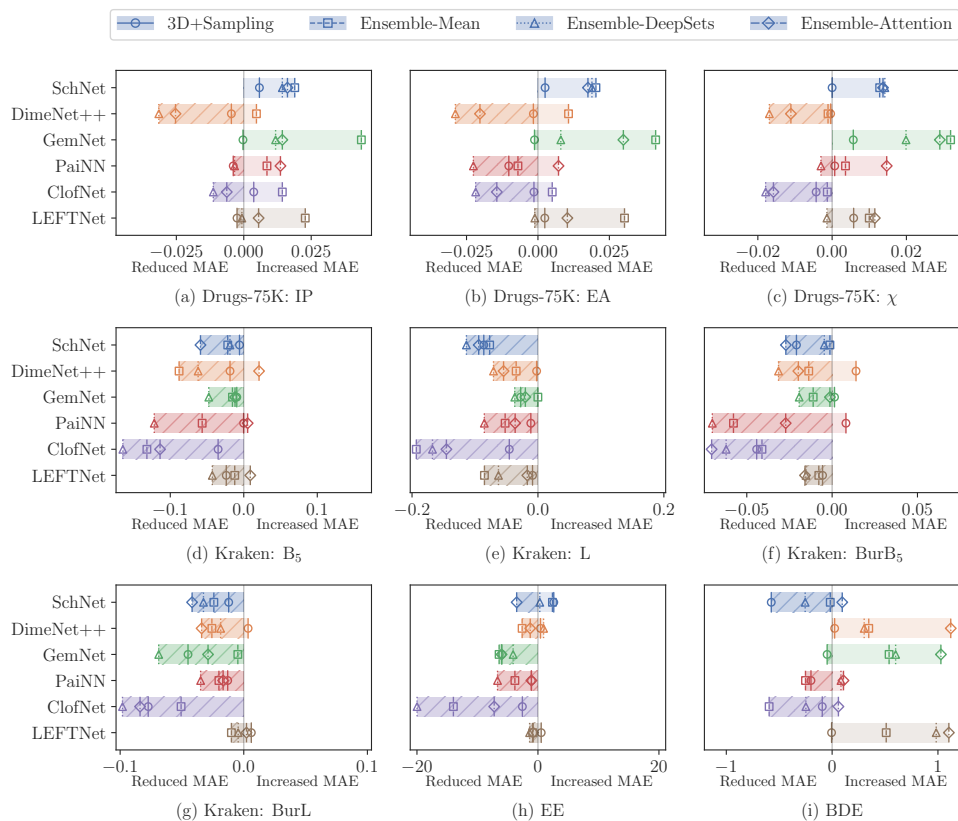


Figure 3: *Performance changes* of four conformer ensemble learning strategies on the basis of six 3D graph models. Here, negative values (marked in *hatch patterns*) denote *reduced* Mean Absolute Error (MAE), signifying a performance improvement due to the incorporation of conformer ensembles.

297 because enantioselectivity depends on subtle spatial interactions. When predicting binding energies,
 298 using 3D models also leads to modest improvements.

299 Overall, model performance varies substantially across tasks, even within the same dataset, emphasizing
 300 the diversity of the tasks in MARCEL. Generally, 1D and 2D models perform well on small-scale
 301 molecular datasets, while 3D models excel on large datasets and reaction-centric tasks. MARCEL
 302 also highlights the benefits of explicitly encoding multiple conformers to improve MRL.

303 6 Discussions and Conclusions

304 In this work, we present the first Molecular Conformer Ensemble Learning benchmark (MARCEL)
 305 to evaluate the potential of learning from a set of conformer structures. Through two conformer
 306 ensemble learning strategies, we discover performance improvements across various tasks. However,
 307 there are some limitations that require further consideration. First, our studied ensemble learning
 308 strategies do not universally improve performance across all tasks and datasets. This highlights the
 309 need for more tailored approaches that integrate with domain expertise to better model specific tasks
 310 and datasets of practical interest. Second, the computational cost of encoding all conformers within
 311 the ensembles, especially for larger datasets, suggests the need to further study the trade-offs between
 312 model complexity and efficiency. Finally, our datasets only contain regression tasks and do not cover
 313 all of the relevant chemical space, which might limit the generalization of our experimental findings.

314 Despite these challenges, we envision that our work will prompt further research in the geometric
 315 deep learning community on how to make use of conformer ensembles for molecular property
 316 prediction. For instance, future research could explore new model architectures that can efficiently
 317 encode ensemble-level information or more sophisticated conformer sampling strategies. We also
 318 hope that our work will stimulate collaborative research across the machine learning and chemistry
 319 fields, with the ultimate goal of pushing the boundaries of predictive molecular modeling and aligning
 320 algorithmic advancements with the practical needs of the chemistry community.

321 References

- 322 [1] Jun Xia, Yanqiao Zhu, Yuanqi Du, and Stan Z. Li. A Systematic Survey of Chemical Pre-trained Models.
323 2023. [1](#)
- 324 [2] Oliver Wieder, Stefan Kohlbacher, Méline Kuenemann, Arthur Garon, Pierre Ducrot, Thomas Seidel, and
325 Thierry Langer. A Compact Review of Molecular Property Prediction with Graph Neural Networks. *Drug*
326 *Discov. Today Technol.*, 37:1–12, 2020. [1](#)
- 327 [3] W. Patrick Walters and Regina Barzilay. Applications of Deep Learning in Molecule Generation and
328 Molecular Property Prediction. *Acc. Chem. Res.*, 54(2):263–270, 2021. [1](#)
- 329 [4] H. L. Morgan. The Generation of a Unique Machine Description for Chemical Structures — A Technique
330 Developed at Chemical Abstracts Service. *J. Chem. Doc.*, 5(2):107–113, 1965. [1](#)
- 331 [5] Robert C. Glem, Andreas Bender, Catrin H. Arnbj, Lars Carlsson, Scott Boyer, and James Smith. Circular
332 Fingerprints: Flexible Molecular Descriptors with Applications from Physical Chemistry to ADME.
333 *IDrugs*, 9(3):199–204, 2006. [1](#)
- 334 [6] G. Skoraczynski, P. Dittwald, B. Miasojedow, S. Szymkuć, E. P. Gajewska, B. A. Grzybowski, and
335 A. Gambin. Predicting the Outcomes of Organic Reactions via Machine Learning: Are Current Descriptors
336 Sufficient? *Sci. Rep.*, 7(1):1–9, 2017. [1](#)
- 337 [7] Zhen Liu, Yurii S. Moroz, and Olexandr Isayev. The Challenge of Balancing Model Sensitivity and
338 Robustness in Predicting Yields: A Benchmarking Study of Amide Coupling Reactions. *chemrxiv.org*,
339 2023. [1](#)
- 340 [8] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Networks.
341 In *ICLR*, 2017. [1](#)
- 342 [9] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural Message
343 Passing for Quantum Chemistry. In *ICML*, pages 1263–1272, 2017. [1](#)
- 344 [10] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka.
345 Representation Learning on Graphs with Jumping Knowledge Networks. In *ICML*, pages 5453–5462,
346 2018. [1](#)
- 347 [11] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio.
348 Graph Attention Networks. In *ICLR*, 2018. [1](#)
- 349 [12] Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Sauceda Felix, Stefan Chmiela, Alexandre Tkatchenko,
350 and Klaus-Robert Müller. SchNet: A Continuous-Filter Convolutional Neural Network for Modeling
351 Quantum Interactions. In *NIPS*, pages 991–1001, 2017. [1](#), [4](#)
- 352 [13] Johannes Klicpera, Janek Groß, and Stephan Günnemann. Directional Message Passing for Molecular
353 Graphs. In *ICLR*, 2020. [1](#), [4](#)
- 354 [14] Johannes Gasteiger, Florian Becker, and Stephan Günnemann. GemNet: Universal Directional Graph
355 Neural Networks for Molecules. In *NeurIPS*, pages 6790–6802, 2021. [1](#), [4](#)
- 356 [15] Kristof Schütt, Oliver T. Unke, and Michael Gastegger. Equivariant Message Passing for the Prediction of
357 Tensorial Properties and Molecular Spectra. In *ICML*, pages 9377–9388, 2021. [1](#), [4](#)
- 358 [16] Weitao Du, He Zhang, Yuanqi Du, Qi Meng, Wei Chen, Nanning Zheng, Bin Shao, and Tie-Yan Liu. SE(3)
359 Equivariant Graph Neural Networks with Complete Local Frames. In *ICML*, pages 5583–5608, 2022. [1](#), [4](#)
- 360 [17] Weitao Du, Yuanqi Du, Limei Wang, Dieqiao Feng, Guifeng Wang, Shuiwang Ji, Carla Gomes, and
361 Zhi-Ming Ma. A New Perspective on Building Efficient and Expressive 3D Equivariant Graph Neural
362 Networks. *arXiv.org*, 2023. [1](#), [4](#)
- 363 [18] Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan,
364 and Martin Grohe. Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. In *AAAI*,
365 pages 4602–4609, 2019. [2](#)
- 366 [19] Keir Adams, Lagnajit Pattanaik, and Connor W. Coley. Learning 3D Representations of Molecular Chirality
367 with Invariance to Bond Rotations. In *ICLR*, 2022. [2](#), [5](#)
- 368 [20] Bharath Ramsundar, Peter Eastman, Patrick Walters, and Vijay Pande. *Deep Learning for the Life Sciences:
369 Applying Deep Learning to Genomics, Microscopy, Drug Discovery, and More*. O’Reilly Media, 2019. [2](#)

- 370 [21] Emanuele Perola and Paul S. Charifson. Conformational Analysis of Drug-Like Molecules Bound to
371 Proteins: An Extensive Study of Ligand Reorganization upon Binding. *J. Med. Chem.*, 47(10):2499–2510,
372 2004. 2
- 373 [22] Andrew F. Zahrt, Jeremy J. Henle, Brennan T. Rose, Yang Wang, William T. Darrow, and Scott E. Denmark.
374 Prediction of Higher-Selectivity Catalysts by Computer-Driven Workflow and Machine Learning. *Science*,
375 363(6424):eaau5631, 2019. 2
- 376 [23] Simon Axelrod and Rafael Gómez-Bombarelli. Molecular Machine Learning with Conformer Ensembles.
377 *arXiv.org*, 2020. 2, 6
- 378 [24] Jan Weinreich, Nicholas J. Browning, and O. Anatole von Lilienfeld. Machine Learning of Free Energies
379 in Chemical Compound Space Using Ensemble Representations: Reaching Experimental Uncertainty for
380 Solvation. *J. Chem. Phys.*, 154(13):134113, 2021. 2
- 381 [25] Kangway V. Chuang and Michael J. Keiser. Attention-Based Learning on Molecular Ensembles. *arXiv.org*,
382 2020. 2
- 383 [26] Chaitanya K. Joshi, Arian R. Jamasb, Ramon Viñas, Charles Harris, Simon Mathis, and Pietro Liò.
384 Multi-State RNA Design with Geometric Multi-Graph Neural Networks. *arXiv.org*, 2023. 2
- 385 [27] Simon Axelrod and Rafael Gómez-Bombarelli. GEOM, Energy-Annotated Molecular Conformations for
386 Property Prediction and Molecular Generation. *Sci. Data*, 9(1):185, 2022. 3, 12
- 387 [28] Andrzej M. Żurański, Jason Y. Wang, Benjamin J. Shields, and Abigail G. Doyle. Auto-QChem: An
388 Automated Workflow for the Generation and Storage of DFT Calculations for Organic Molecules. *React.*
389 *Chem. Eng.*, 7(6):1276–1284, 2022. 3, 12
- 390 [29] Tobias Gensch, Gabriel dos Passos Gomes, Pascal Friederich, Ellyn Peters, Théophile Gaudin, Robert
391 Pollice, Kjell Jorner, AkshatKumar Nigam, Michael Lindner-D’Addario, Matthew S. Sigman, and Alán
392 Aspuru-Guzik. A Comprehensive Discovery Platform for Organophosphorus Ligands for Catalysis. *J. Am.*
393 *Chem. Soc.*, 144:1205–1217, 2022. 3, 13
- 394 [30] Anthony R. Rosales, Jessica Wahlers, Elaine Limé, Rebecca E. Meadows, Kevin W. Leslie, Rhona Savin,
395 Fiona Bell, Eric Hansen, Paul Helquist, Rachel H. Munday, Olaf Wiest, and Per-Ola Norrby. Rapid Virtual
396 Screening of Enantioselective Catalysts Using CatVS. *Nat. Catal.*, 2(1):41–45, 2019. 3, 13, 14
- 397 [31] Benjamin Meyer, Boodsarin Sawatlon, Stefan Heinen, O. Anatole von Lilienfeld, and Clémence Cormin-
398 boeuf. Machine Learning Meets Volcano Plots: Computational Discovery of Cross-Coupling Catalysts.
399 *Chem. Sci.*, 9:7069–7077, 2018. 3, 14
- 400 [32] Noel M. O’Boyle, Michael Banck, Craig A. James, Chris Morley, Tim Vandermeersch, and Geoffrey R.
401 Hutchison. Open Babel: An Open Chemical Toolbox. *J. Cheminformatics*, 3(1):1–14, 2011. 3, 14
- 402 [33] Greg Landrum, Paolo Tosco, Brian Kelley, Ric, sriniker, gedec, Riccardo Vianello, NadineSchneider,
403 Eisuke Kawashima, Andrew Dalke, Dan N, David Cosgrove, Brian Cole, Matt Swain, Samo Turk,
404 AlexanderSavelyev, Gareth Jones, Alain Vaucher, Maciej Wójcikowski, Ichiru Take, Daniel Probst, Kazuya
405 Ujihara, Vincent F. Scalfani, guillaume godin, Axel Pahl, Francois Berenger, JLVarjo, strets123, JP, and
406 DoliathGavid. rdkit/rdkit: 2022_03_2 (q1 2022) release, 2022. 4, 14
- 407 [34] Darko Butina. Unsupervised Data Base Clustering Based on Daylight’s Fingerprint and Tanimoto Similarity:
408 A Fast and Automated Way To Cluster Small and Large Data Sets. *J. Chem. Inf. Comput. Sci.*, 39(4):
409 747–750, 1999. 4
- 410 [35] Leo Breiman. Random Forests. *Mach. Learn.*, 45(1):5–32, 2001. 4
- 411 [36] David Rogers and Mathew Hahn. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model.*, 50:742–754,
412 2010. 4, 14
- 413 [37] Joseph L. Durant, Burton A. Leland, Douglas R. Henry, and James G. Nourse. Reoptimization of MDL
414 Keys for Use in Drug Discovery. *J. Chem. Inf. Comput. Sci.*, 42(6):1273–1280, 2002. 4, 14
- 415 [38] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Comp.*, 9(8):1735–1780,
416 1997. 4
- 417 [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Uszkoreit
418 Kaiser, and Illia Polosukhin. Attention is All You Need. In *NIPS*, pages 5998–6008, 2017. 4, 14

- 419 [40] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How Powerful Are Graph Neural Networks?
420 In *ICLR*, 2019. 4
- 421 [41] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and
422 Jure Leskovec. Open Graph Benchmark: Datasets for Machine Learning on Graphs. In *NeurIPS*, pages
423 22118–22133, 2020. 4, 15
- 424 [42] Kevin Yang, Kyle Swanson, Wengong Jin, Connor W. Coley, Philipp Eiden, Hua Gao, Angel Guzman-Perez,
425 Timothy Hopper, Brian Kelley, Miriam Mathea, Andrew Palmer, Volker Settels, Tommi S. Jaakkola, Klavs F.
426 Jensen, and Regina Barzilay. Analyzing Learned Molecular Representations for Property Prediction. *J. Chem. Inf. Model.*, 59(8):3370–3388, 2019. 4
- 428 [43] Ladislav Rampásek, Mikhail Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique
429 Beaini. Recipe for a General, Powerful, Scalable Graph Transformer. In *NeurIPS*, pages 14501–14515,
430 2022. 4
- 431 [44] Thomas G. Dietterich, Richard H. Lathrop, and Tomás Lozano-Pérez. Solving the Multiple Instance
432 Problem with Axis-Parallel Rectangles. *Artif. Intell.*, 89(1-2):31–71, 1997. 5
- 433 [45] Oded Maron and Tomás Lozano-Pérez. A Framework for Multiple-Instance Learning. In *NIPS*, pages
434 570–576, 1997. 5
- 435 [46] Maximilian Ilse, Jakub M. Tomczak, and Max Welling. Attention-based Deep Multiple Instance Learning.
436 In *ICML*, pages 2132–2141, 2018. 5
- 437 [47] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabás Póczos, Ruslan R. Salakhutdinov, and
438 Alexander J. Smola. Deep Sets. In *NIPS*, pages 3391–3401, 2017. 5
- 439 [48] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural Machine Translation by Jointly Learning
440 to Align and Translate. In *ICLR*, 2015. 6
- 441 [49] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015. 6, 15
- 442 [50] Zhen Liu, Tetiana Zubatiuk, Adrian Roitberg, and Olexandr Isayev. Auto3D: Automatic Generation of the
443 Low-Energy 3D Structures with ANI Neural Network Potentials. *J. Chem. Inf. Model.*, 62(22):5373–5382,
444 2022. 12
- 445 [51] Roman Zubatyuk, Justin S. Smith, Benjamin T. Nebgen, Sergei Tretiak, and Olexandr Isayev. Teaching a
446 Neural Network to Attach and Detach Electrons From Molecules. *Nat. Commun.*, 12(1):1–11, 2021. 12
- 447 [52] Qiyuan Zhao, Sai Mahit Vaddadi, Michael Woulfe, Lawal A. Ogunfowora, Sanjay S. Garimella, Olexandr
448 Isayev, and Brett M. Savoie. Comprehensive Exploration of Graphically Defined Reaction Spaces. *Sci. Data*,
449 10(1):1–10, 2023. 12
- 450 [53] Carlo Adamo and Vincenzo Barone. Toward Reliable Density Functional Methods Without Adjustable
451 Parameters: The PBE0 Model. *J. Chem. Phys.*, 110(13):6158–6170, 1999. 12
- 452 [54] A. Verloop, W. Hoogenstraaten, and J. Tipker. Development and Application of New Steric Substituent
453 Parameters in Drug Design. In E. J. Ariëns, editor, *Drug Design*, volume 11 of *Medicinal Chemistry: A*
454 *Series of Monographs*, pages 165–207. Academic Press, Amsterdam, 1976. 13
- 455 [55] G. P. Moss. Basic Terminology of Stereochemistry (IUPAC Recommendations 1996). *Pure Appl. Chem.*,
456 68(12):2193–2222, 1996. 14
- 457 [56] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel,
458 Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake VanderPlas, Alexandre Passos,
459 David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay. Scikit-learn: Machine
460 Learning in Python. *J. Mach. Learn. Res.*, 12:2825–2830, 2011. 14
- 461 [57] Philip Gage. A New Algorithm for Data Compression. *C Users J.*, 12(2):23–38, 1994. 14
- 462 [58] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
463 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang,
464 Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie
465 Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In
466 *NeurIPS*, pages 8024–8035, 2019. 15
- 467 [59] Matthias Fey and Jan Eric Lenssen. Fast Graph Representation Learning with PyTorch Geometric. In
468 *RLGM@ICLR*, 2019. 15

Supplementary Material for MARCEL

469

470	A Dataset Description	12
471	A.1 Drugs-75K	12
472	A.2 Kraken	13
473	A.3 EE	13
474	A.4 BDE	14
475	B Implementation Details	14
476	B.1 Implementation of 1D Models	14
477	B.2 Featurizations of Molecules for 2D Models	15
478	B.3 Hyperparameter Specifications and Experimental Environments	15
479	C Additional Experiments on Evaluation Schemes of the Conformer Sampling Strategy	15
480	D Raw Data	16

481 A Dataset Description

482 MARCEL include four datasets that cover a diverse range of chemical space, which focuses on four
483 chemically-relevant tasks for both molecules and reactions, with an emphasis on Boltzmann-averaged
484 properties of conformer ensembles computed at the Density-Functional Theory (DFT) level. Detailed
485 information regarding dataset access, data formatting, and loading procedures can be found at our
486 GitHub repository <https://anonymous.4open.science/r/MARCEL-4813>. Any subsequent
487 updates will also be posted on this repository.

488 A.1 Drugs-75K

489 Drugs-75K is a subset of the GEOM-Drugs [27] dataset, which includes 75,099 drug-like molecules
490 with at least 5 rotatable bonds. The original GEOM-Drugs dataset was constructed using semi-
491 empirical DFT methods, which is less accurate than full DFT. To curate the Drugs-75K subset,
492 Auto3D [50] is used to generate and optimize the conformer ensembles for each molecule and
493 AIMNet-NSE [51] is used to calculate three important DFT-based reactivity descriptors: ionization
494 potential, electron affinity, and electronegativity [28].

495 Auto3D [50] efficiently generates high-quality conformers, with a mean RMSD at around 0.2 Å when
496 compared with DFT conformers. It has been used in other large conformer dataset generation [52].
497 Regarding the neural network surrogate AIMNET-NSE [51], it mimics the PBE0/ma-def2-SVP
498 method of DFT, which is widely used in the chemistry community. Investigating their accuracy is out
499 of the scope of this paper, but are readily accessible from multiple sources [51, 53].

500 **Objectives.** The tasks are to predict the Boltzmann-averaged value of each property across the
501 conformer ensemble $\langle y \rangle_{k_B} = \sum_{C_i \in \mathcal{C}} p_i y_i$, where y_i is a conformer-specific property. We are given
502 each C_i , and the goal is to predict $\langle y \rangle_{k_B}$ from the molecular graph G , a single conformer $C_i \in \mathcal{C}$, or
503 the set \mathcal{C} .

504 **Dataset preparation.** In preparing the 75K version of GEOM-Drugs, we begin with the original
505 SMILES strings of the molecules. We first exclude molecules that have less than 5 rotatable bonds.
506 To enable the utilization of AIMNet-NSE for descriptor computation, we retain only those molecules
507 containing atoms of H, C, N, O, F, Si, P, S, and Cl. Further, we generate DFT-level conformers
508 and compute their energies with Auto3D. Based on these conformers, we compute three chemical
509 bond energy descriptors using AIMNet-NSE. We exclude conformers that Auto3D fails to converge
510 and charged molecules that are unable to be processed by AIMNet-NSE, which results in 75,099
511 molecules. Subsequently, we compute molecular-level Boltzmann-averaged descriptors based on

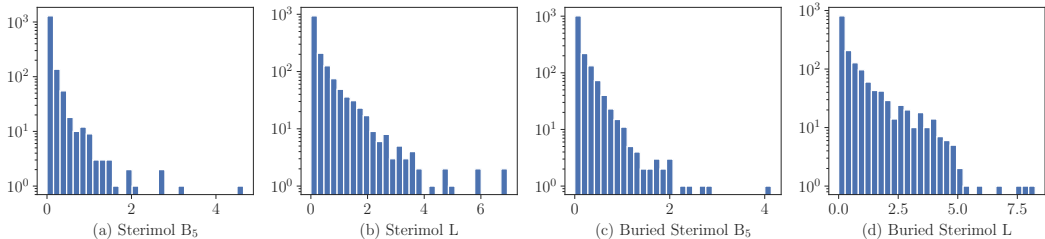


Figure S1: Histogram of the ratio of the variance of each conformer property to the variance of each Boltzmann-averaged property in the Kraken dataset.

512 conformer-level descriptors. Finally, we undertake a deduplication process as outlined in Section 3
 513 with a RMSD threshold of 2.0, which yields a total of 558,002 distinct conformers.

514 **Data availability and license.** The original GEOM-Drugs dataset is publicly available at <https://github.com/learningmatter-mit/geom> but no license is specified. Our Drugs-75K can be
 515 accessed at [https://anonymous.4open.science/r/MARCEL-4813/datasets/Drugs/README](https://anonymous.4open.science/r/MARCEL-4813/datasets/Drugs/README.md)
 516 [E.md](https://anonymous.4open.science/r/MARCEL-4813/datasets/Drugs/README.md). As for the conformer ensembles and descriptors that we generated, they are licensed under the
 517 Apache License.
 518

519 A.2 Kraken

520 Kraken [29] is a dataset of 1,552 monodentate organophosphorus (III) ligands along with their
 521 DFT-computed conformer ensembles. In this study, we consider four 3D catalytic ligand descriptors
 522 exhibiting significant variance among conformers: Sterimol B₅, Sterimol L, buried Sterimol B₅, and
 523 buried Sterimol L. These descriptors quantify the steric size of a substituent in Å, and are commonly
 524 employed for Quantitative Structure-Activity Relationship (QSAR) modeling. The buried Sterimol
 525 variants describe the steric effects within the first coordination sphere of a metal [54].

526 **Objectives.** As in the Drugs-75K tasks, the goal is to predict the Boltzmann-averaged value of each
 527 property across the conformer ensemble from the molecular graph G , a single conformer $C_i \in \mathcal{C}$, or
 528 the set \mathcal{C} .

529 **Dataset preparation.** In this study, we utilize the original 3D geometry structures of molecules and
 530 their corresponding Boltzmann-averaged properties provided in the Kraken dataset. Among the 78
 531 physical-organic properties listed in the original dataset, we select four properties that demonstrate
 532 high variance across conformer ensembles, as illustrated in Figure S1.

533 **Data availability and license.** The Kraken dataset is publicly accessible at [https://kraken.cs.](https://kraken.cs.toronto.edu)
 534 [toronto.edu](https://kraken.cs.toronto.edu). Its copyright is retained by the original authors. Under the permission of the original
 535 authors, the Kraken dataset with the conformer ensembles and the four conformer-level descriptors
 536 used in this study can be accessed at [https://anonymous.4open.science/r/MARCEL-4813/da](https://anonymous.4open.science/r/MARCEL-4813/datasets/Kraken/README.md)
 537 [tasetts/Kraken/README.md](https://anonymous.4open.science/r/MARCEL-4813/datasets/Kraken/README.md).

538 A.3 EE

539 EE [30] is a dataset of 872 catalyst-substrate pairs involving 253 Rhodium (Rh)-bound atropisomeric
 540 catalysts derived from chiral bisphosphine, with 10 enamides as substrates. The dataset includes
 541 conformations of catalyst-substrate transition state complexes in two separate pro-S and pro-R
 542 configurations. The task is to predict the Enantiomeric Excess (EE) of the chemical reaction involving
 543 the substrate, defined as the absolute ratio between the concentration of each enantiomer in the
 544 product distribution.

545 **Objectives.** EE depends on the conformer ensembles of *each* pro-R and pro-S complex. The goal is
 546 to predict EE from the graphs of the catalyst and substrate (G_{cat} , G_{sub}), a conformer $C_i^{(R)} \in \mathcal{C}^{(R)}$ and
 547 $C_i^{(S)} \in \mathcal{C}^{(S)}$ for each complex, or the ensembles $\mathcal{C}^{(R)}$ and $\mathcal{C}^{(S)}$.

548 **Dataset preparation.** The conformer ensembles are generated with Q2MM, which automatically
 549 generates Transition State Force Fields (TSFFs) in order to simulate the conformer ensembles of each

550 prochiral transition state complex. Then, the EE values are computed from the conformer ensembles
551 by Boltzmann-averaging the activation energies for the competing transition states [30, 55]. Finally,
552 we conduct the same conformer deduplication process as described in Section 3 with a RMSD
553 threshold of 1.0.

554 **Data availability and license.** As of now, the EE dataset is proprietary, given that the publication
555 addressing the conformer ensembles is still under preparation. Therefore, access to the EE dataset is
556 restricted to review purposes only. We anticipate making the EE dataset publicly accessible following
557 the acceptance of the corresponding paper.

558 A.4 BDE

559 BDE [31] is a dataset containing 5,915 organometallic catalysts ML_1L_2 consisting of a metal center
560 ($M = Pd, Pt, Au, Ag, Cu, Ni$) coordinated to two flexible organic ligands (L_1 and L_2), each selected
561 from a 91-membered ligand library. The data includes conformations of each unbound catalyst, as
562 well as conformations of the catalyst when bound to ethylene and bromide after oxidative addition
563 with vinyl bromide. Each catalyst has an electronic binding energy, computed as the difference
564 in the minimum energies of the bound-catalyst complex and unbound catalyst, following the DFT-
565 optimization of their respective conformer ensembles.

566 Although the binding energies are computed via DFT, the conformers provided for modeling are
567 initially generated with Open Babel [32], followed by further geometric optimization steps, which
568 ensures that the generated 3D structures are likely to be the global minimum energy conformers at
569 the force field level [31, Supplementary Information]. We also note that obtaining DFT-optimized
570 conformers for BDE is not feasible given the time-consuming nature of the process — a single
571 geometric search using DFT can take 2 to 3 days. Therefore, this realistically represents the setting in
572 which precise conformer ensembles are unknown at inference.

573 **Objectives.** The task is to predict the binding energy from the graphs of the unbound and bound cata-
574 lyst, sampled conformers $C_i^{(\text{unbound})} \in \mathcal{C}^{(\text{unbound})}$ and $C_i^{(\text{bound})} \in \mathcal{C}^{(\text{bound})}$, or the ensembles $\mathcal{C}^{(\text{unbound})}$
575 and $\mathcal{C}^{(\text{bound})}$.

576 **Dataset preparation.** We employ Open Babel [32] to produce conformers for each unbound catalyst
577 and each bound complex. In order to avoid redundancy, we follow a deduplication process as outlined
578 in Section 3. For the unbound catalysts, a RMSD threshold value of 0.5 is applied, whereas for the
579 bound complexes, a threshold of 1.0 is used.

580 **Data availability and license.** The binding energy descriptors can be accessed at <https://archive.materialscloud.org/record/2018.0014/v1> under the Creative Commons Attribution 4.0
581 International license. The conformers are publicly available at <https://anonymous.4open.science/r/MARCEL-4813/datasets/BDE/README.md> under the Apache license.
582
583

584 B Implementation Details

585 B.1 Implementation of 1D Models

586 For the random forest model that operates on fingerprints, we employ three molecular finger-
587 print schemes: the Molecular ACCess System (MACCS) [37], Extended-Connectivity Fingerprints
588 (ECFP) [36], and RDKit topological fingerprints [33]. Then, we concatenate their outputs into a
589 single vector, which might lead to some feature redundancy, given the possible overlaps in these three
590 fingerprint representations of the molecular structure. To tackle this issue, we remove any features
591 that exhibit a high correlation exceeding 90% with the other features. For implementation, we employ
592 Scikit-Learn [56] and compute fingerprints with RDKit [33].

593 For both LSTM and Transformer models that operate on SMILES strings, we use a Byte-Pair
594 Encoding (BPE)-based tokenizer [57] that is pretrained on PubChem10M, which strikes a balance
595 among character- and word-level representations and allows to handle large vocabularies in molecular
596 corpora. For the Transformer model, we further follow the positional embedding scheme [39] to
597 capture the positional relationship among tokens in the SMILES string.

Table S1: A summary of node and edge features used in 2D GNN models.

	Feature	Explanation
Node	AtomicNum	Atomic number, representing the type of atom.
	ChiralTag	Indicator of chirality, a property of asymmetry.
	TotalDegree	Sum of implicit and explicit bonds of an atom.
	FormalCharge	Charge of an atom assuming equal sharing of bonding electrons.
	TotalNumHs	Total number of hydrogen atoms bonded to the atom.
	NumRadicalElectrons	Count of unpaired electrons in an atom.
	Hybridization	Type of atomic orbital hybridization in the atom.
	IsAromatic	Boolean indicating if the atom is part of an aromatic ring.
	IsInRing	Boolean indicating if the atom is part of any ring structure.
Edge	BondType	Type of the bond (e.g., single, double, triple, aromatic).
	Stereo	Stereochemistry of the bond (e.g., "none", "any", "Z", or "E" for double bonds).
	IsConjugated	Boolean indicating if the bond is part of a conjugated system.

598 B.2 Featurizations of Molecules for 2D Models

599 Following OGB [41], we employ a rich set of features for atoms (nodes) and bonds (edges) for 2D
600 GNN models. A complete list of node and features can be found in Table S1.

601 B.3 Hyperparameter Specifications and Experimental Environments

602 Each model is trained over 2,000 epochs using the Adam optimizer [49] with early stopping triggered
603 if there is no improvement in the training loss over 200 epochs. To ensure a fair comparison, the
604 hidden dimension size is uniformly set to 128 for all models. Other hyperparameters mostly follow
605 the original configurations as described in the respective papers. The complete hyperparameter set of
606 each model can be found in <https://anonymous.4open.science/r/MARCEL-4813/benchmarks/params>.
607

608 We utilize PyTorch [58] and PyTorch-Geometric [59] to implement all deep learning models. Most
609 of the experiments are conducted on servers equipped with NVIDIA A100 GPUs, each with 40GB of
610 memory. For memory-intensive models such as GemNet and LEFTNet, we use servers with NVIDIA
611 H100 GPUs, each with 80GB memory. The cumulative computation time across all experiments
612 amounts to approximately 6,000 single GPU hours.

613 C Additional Experiments on Evaluation Schemes of the Conformer 614 Sampling Strategy

615 In this section, we conduct one additional experiment on the conformer ensemble learning strategies.
616 We assess all 3D models on five tasks: Ionization Potential (IP) from the Drugs-75K dataset, B₅ and
617 BurB₅ from the Kraken dataset, and tasks from the EE and BDE datasets.

618 In our previous setup, we evaluate the conformer sampling strategy using the lowest-energy conformer
619 of each molecule at evaluation time, to provide a direct comparison to the single-conformer 3D
620 models that are trained and tested with the lowest energy conformation. In these experiments, we
621 continue to sample a random conformer uniformly from the conformer ensemble during training
622 time, but consider two additional evaluation schemes: (1) evaluating model performance when
623 encoding a randomly sampled conformer, and (2) evaluating model -performance when averaging the
624 per-conformer predictions across the entire conformer ensemble.

625 The results of these experiments are summarized in Table S2. In the table, we refer to the original
626 evaluation scheme as "fixed", and the additional schemes as "random" and "all", respectively. We
627 find that across all three schemes, using the lowest-energy conformer for evaluation consistently
628 yields the best performance. This is expected, as the lowest-energy conformer contributes the most
629 to ensemble-level descriptors. The random conformer evaluation scheme generally yields the worst
630 performance, which is likely due to the introduction of noise from less relevant conformers at test

Table S2: Performance comparison of three conformer sampling variants with different evaluation strategies. All models are trained with a randomly sampled conformer from the ensemble. The last column summarizes the average rank across all datasets for each base model.

Model	Evaluation Strategy	Drugs-75K	Kraken		EE	BDE	Average Rank
		IP	B ₅	BurB ₅			
SchNet	Fixed	0.4452	0.3235	0.2086	20.3595	1.9737	1
	Random	0.4498	0.3682	0.2454	22.0380	2.4416	3
	All	0.4428	0.3856	0.2407	18.0296	2.0106	2
DimeNet++	Fixed	0.4395	0.3323	0.2237	15.0596	1.4741	= 2
	Random	0.4555	0.3549	0.2222	13.5681	1.4688	= 2
	All	0.4479	0.3282	0.2001	12.3562	1.6270	1
GemNet	Fixed	0.4066	0.2694	0.1796	12.0541	1.6059	1
	Random	0.4250	0.4034	0.2534	16.1709	1.7894	3
	All	0.4320	0.4523	0.2481	14.3952	1.6660	2
PaiNN	Fixed	0.4466	0.3441	0.2476	19.1521	1.9262	1
	Random	0.4770	0.3756	0.2478	21.3553	1.9411	3
	All	0.4478	0.3458	0.2342	19.1955	1.8696	2
ClofNet	Fixed	0.4430	0.4524	0.2442	31.3733	2.5126	1
	Random	0.4530	0.4689	0.2736	31.3675	2.6310	= 2
	All	0.4363	0.4749	0.2855	34.3203	2.0271	= 2
LEFTNet	Fixed	0.4149	0.2834	0.2120	20.3358	1.5276	1
	Random	0.4518	0.3177	0.2344	20.3740	1.5842	3
	All	0.4274	0.3152	0.2170	18.8945	1.8663	2

631 time. Interestingly, we observe occasional performance improvement when averaging the predictions
632 across all conformers in the ensemble, indicating that explicitly using ensemble-level information
633 during evaluation can be beneficial.

634 D Raw Data

635 The raw performance data with standard deviation of Table 2 and Figure 3 is summarized in Table S3.

Table S3: Raw performance data (mean ± standard deviation) of representative 1D, 2D, 3D, and conformer ensemble MRL models in terms of absolute test error.

Category	Model	Drugs-75K			Kraken				EE	BDE	
		IP	EA	χ	B ₅	L	BurB ₅	BurL			
1D	Random forest	0.4987 _{±0.0037}	0.4747 _{±0.0022}	0.2732 _{±0.0031}	0.4760 _{±0.0041}	0.4303 _{±0.0090}	0.2758 _{±0.0180}	0.1521 _{±0.0149}	61.2963 _{±2.8640}	3.0335 _{±0.2693}	
	LSTM	0.4788 _{±0.0024}	0.4648 _{±0.0002}	0.2505 _{±0.0050}	0.4879 _{±0.0280}	0.5142 _{±0.0411}	0.2813 _{±0.0041}	0.1924 _{±0.0028}	64.0088 _{±2.3708}	2.8279 _{±0.0728}	
	Transformer	0.6617 _{±0.0023}	0.5850 _{±0.0031}	0.4073 _{±0.0006}	0.9611 _{±0.0813}	0.8389 _{±0.0431}	0.4929 _{±0.0369}	0.2781 _{±0.0207}	62.0816 _{±2.1789}	10.0771 _{±0.6457}	
2D	GIN	0.4354 _{±0.0029}	0.4169 _{±0.0032}	0.2260 _{±0.0017}	0.3128 _{±0.0264}	0.4003 _{±0.0341}	0.1719 _{±0.0031}	0.1200 _{±0.0040}	62.3065 _{±2.9010}	2.6368 _{±0.2276}	
	GIN-VN	0.4361 _{±0.0059}	0.4169 _{±0.0083}	0.2267 _{±0.0002}	0.3567 _{±0.0031}	0.4344 _{±0.0416}	0.2422 _{±0.0033}	0.1741 _{±0.0109}	62.3815 _{±2.1882}	2.7417 _{±0.2446}	
	ChemProp	0.4595 _{±0.0028}	0.4417 _{±0.0045}	0.2441 _{±0.0012}	0.4850 _{±0.0968}	0.5452 _{±0.0454}	0.3002 _{±0.0086}	0.1948 _{±0.0138}	61.0336 _{±2.9715}	2.6616 _{±0.1429}	
	GraphGPS	0.4351 _{±0.0049}	0.4085 _{±0.0035}	0.2212 _{±0.0054}	0.3450 _{±0.0324}	0.4363 _{±0.0133}	0.2066 _{±0.0115}	0.1500 _{±0.0138}	61.6251 _{±1.3743}	2.4827 _{±0.1992}	
3D	SchNet	0.4394 _{±0.0062}	0.4207 _{±0.0021}	0.2243 _{±0.0089}	0.3293 _{±0.0068}	0.5458 _{±0.0341}	0.2295 _{±0.0111}	0.1861 _{±0.0095}	17.7421 _{±1.0899}	2.5488 _{±0.0050}	
	DimeNet++	0.4441 _{±0.0087}	0.4233 _{±0.0072}	0.2436 _{±0.0075}	0.3510 _{±0.0107}	0.4174 _{±0.0397}	0.2097 _{±0.0160}	0.1526 _{±0.0072}	14.6414 _{±2.2791}	1.4503 _{±0.0370}	
	GemNet	0.4069 _{±0.0007}	0.3922 _{±0.0024}	0.1970 _{±0.0039}	0.2789 _{±0.0125}	0.3754 _{±0.0086}	0.1782 _{±0.0099}	0.1635 _{±0.0063}	18.0338 _{±2.4777}	1.6530 _{±0.3081}	
	PaiNN	0.4505 _{±0.0041}	0.4495 _{±0.0054}	0.2324 _{±0.0040}	0.3443 _{±0.0388}	0.4471 _{±0.0324}	0.2395 _{±0.0176}	0.1673 _{±0.0088}	20.2359 _{±1.2128}	2.1261 _{±0.0920}	
	ClofNet	0.4393 _{±0.0084}	0.4251 _{±0.0066}	0.2378 _{±0.0020}	0.4873 _{±0.0093}	0.6417 _{±0.0362}	0.2884 _{±0.0166}	0.2529 _{±0.0052}	33.9473 _{±1.4633}	2.6057 _{±0.0236}	
LEFTNet	0.4174 _{±0.0007}	0.3964 _{±0.0099}	0.2083 _{±0.0054}	0.3072 _{±0.0012}	0.4493 _{±0.0261}	0.2176 _{±0.0010}	0.1486 _{±0.0095}	19.7974 _{±1.4097}	1.5328 _{±0.0567}		
3D +Sampling	SchNet	0.4452 _{±0.0080}	0.4232 _{±0.0042}	0.2243 _{±0.0022}	0.3235 _{±0.0147}	0.4598 _{±0.0041}	0.2086 _{±0.0111}	0.1739 _{±0.0142}	20.3595 _{±1.5260}	1.9737 _{±0.0125}	
	DimeNet++	0.4395 _{±0.0032}	0.4217 _{±0.0040}	0.2432 _{±0.0048}	0.3323 _{±0.0320}	0.4153 _{±0.0208}	0.2237 _{±0.0122}	0.1561 _{±0.0241}	15.0596 _{±0.2867}	1.4741 _{±0.0349}	
	GemNet	0.4066 _{±0.0015}	0.3910 _{±0.0004}	0.2027 _{±0.0013}	0.2694 _{±0.0221}	0.3488 _{±0.0252}	0.1796 _{±0.0098}	0.1184 _{±0.0033}	12.0541 _{±0.7735}	1.6059 _{±0.1094}	
	PaiNN	0.4466 _{±0.0087}	0.4393 _{±0.0045}	0.2331 _{±0.0037}	0.3441 _{±0.0161}	0.4358 _{±0.0343}	0.2476 _{±0.0070}	0.1543 _{±0.0022}	19.1521 _{±0.2386}	1.9262 _{±0.0188}	
	ClofNet	0.4430 _{±0.0074}	0.4237 _{±0.0005}	0.2335 _{±0.0090}	0.4524 _{±0.0935}	0.5962 _{±0.0074}	0.2442 _{±0.0109}	0.1756 _{±0.0112}	31.3733 _{±1.9892}	2.5126 _{±0.2366}	
LEFTNet	0.4149 _{±0.0019}	0.3988 _{±0.0048}	0.2141 _{±0.0084}	0.2834 _{±0.0068}	0.4407 _{±0.0531}	0.2120 _{±0.0097}	0.1547 _{±0.0101}	20.3358 _{±0.6614}	1.5276 _{±0.0088}		
Ensemble	Mean	0.4583 _{±0.0019}	0.4410 _{±0.0018}	0.2371 _{±0.0098}	0.3075 _{±0.0151}	0.4691 _{±0.0234}	0.2282 _{±0.0206}	0.1619 _{±0.0062}	20.1392 _{±1.5748}	2.5312 _{±0.0246}	
	SchNet	DeepSet	0.4537 _{±0.0065}	0.4396 _{±0.0010}	0.2385 _{±0.0066}	0.3105 _{±0.0381}	0.4322 _{±0.0464}	0.2249 _{±0.0234}	0.1535 _{±0.0076}	18.0495 _{±1.2846}	2.2941 _{±0.2229}
	Attention	0.4556 _{±0.0075}	0.4382 _{±0.0125}	0.2380 _{±0.0007}	0.2704 _{±0.0187}	0.4517 _{±0.0132}	0.2024 _{±0.0183}	0.1443 _{±0.0043}	14.2238 _{±0.5451}	2.6445 _{±0.0031}	
	Mean	0.4488 _{±0.0086}	0.4340 _{±0.0079}	0.2425 _{±0.0060}	0.2630 _{±0.0122}	0.3828 _{±0.0331}	0.1960 _{±0.0059}	0.1268 _{±0.0060}	12.0259 _{±0.8933}	1.7964 _{±0.1260}	
	DimeNet++	DeepSet	0.4126 _{±0.0076}	0.3944 _{±0.0034}	0.2267 _{±0.0047}	0.2889 _{±0.0069}	0.3468 _{±0.0090}	0.1783 _{±0.0110}	0.1339 _{±0.0087}	15.5754 _{±2.6294}	1.7533 _{±0.0163}
	Attention	0.4188 _{±0.0024}	0.4030 _{±0.0075}	0.2325 _{±0.0028}	0.3718 _{±0.0300}	0.3628 _{±0.0259}	0.1899 _{±0.0081}	0.1185 _{±0.0105}	13.3643 _{±1.4309}	2.5714 _{±0.2149}	
	Mean	0.4505 _{±0.0052}	0.4334 _{±0.0023}	0.2289 _{±0.0032}	0.2635 _{±0.0053}	0.3753 _{±0.0036}	0.1671 _{±0.0154}	0.1587 _{±0.0029}	11.6142 _{±1.7271}	2.1914 _{±0.0605}	
	GemNet	DeepSet	0.4187 _{±0.0022}	0.4002 _{±0.0012}	0.2169 _{±0.0036}	0.2313 _{±0.0026}	0.3386 _{±0.0269}	0.1589 _{±0.0068}	0.0947 _{±0.0012}	13.9273 _{±1.8656}	2.2532 _{±0.2106}
	Attention	0.4212 _{±0.0017}	0.4221 _{±0.0097}	0.2260 _{±0.0056}	0.2670 _{±0.0026}	0.3554 _{±0.0147}	0.1769 _{±0.0153}	0.1346 _{±0.0075}	12.0249 _{±1.8418}	2.6810 _{±0.0223}	
	Mean	0.4591 _{±0.0024}	0.4425 _{±0.0064}	0.2360 _{±0.0032}	0.2877 _{±0.0252}	0.3950 _{±0.0233}	0.1817 _{±0.0091}	0.1472 _{±0.0039}	16.4239 _{±0.0743}	1.8744 _{±0.1657}	
	PaiNN	DeepSet	0.4471 _{±0.0071}	0.4269 _{±0.0033}	0.2294 _{±0.0065}	0.2225 _{±0.0218}	0.3619 _{±0.0192}	0.1693 _{±0.0111}	0.1324 _{±0.0091}	13.5570 _{±0.5505}	2.2097 _{±0.0586}
	Attention	0.4641 _{±0.0016}	0.4567 _{±0.0094}	0.2471 _{±0.0049}	0.3496 _{±0.0140}	0.4109 _{±0.0167}	0.2123 _{±0.0005}	0.1506 _{±0.0029}	19.1556 _{±2.2765}	2.2335 _{±0.1255}	
	Mean	0.4536 _{±0.0030}	0.4301 _{±0.0007}	0.2365 _{±0.0075}	0.3555 _{±0.0193}	0.4485 _{±0.0053}	0.2473 _{±0.0076}	0.2022 _{±0.0212}	19.9710 _{±0.7745}	2.0106 _{±0.0856}	
	ClofNet	DeepSet	0.4280 _{±0.0056}	0.4033 _{±0.0024}	0.2199 _{±0.0073}	0.3228 _{±0.0020}	0.4742 _{±0.0161}	0.2263 _{±0.0249}	0.1548 _{±0.0039}	13.9647 _{±1.2753}	2.3576 _{±0.0496}
	Attention	0.4330 _{±0.0071}	0.4107 _{±0.0048}	0.2220 _{±0.0084}	0.3734 _{±0.0267}	0.4963 _{±0.0286}	0.2178 _{±0.0186}	0.1690 _{±0.0281}	26.7133 _{±1.7225}	2.6652 _{±0.1438}	
	Mean	0.4402 _{±0.0062}	0.4267 _{±0.0026}	0.2183 _{±0.0007}	0.2949 _{±0.0001}	0.3643 _{±0.0352}	0.2098 _{±0.0146}	0.1386 _{±0.0007}	18.9245 _{±2.0136}	2.0440 _{±0.0076}	
	LEFTNet	DeepSet	0.4167 _{±0.0043}	0.3953 _{±0.0009}	0.2069 _{±0.0022}	0.2644 _{±0.0130}	0.3866 _{±0.0270}	0.2023 _{±0.0026}	0.1441 _{±0.0042}	18.4189 _{±1.8922}	2.5165 _{±0.3077}
Attention	0.4229 _{±0.0059}	0.4067 _{±0.0047}	0.2198 _{±0.0011}	0.3161 _{±0.0116}	0.4324 _{±0.0292}	0.2017 _{±0.0023}	0.1508 _{±0.0075}	18.9988 _{±1.6904}	2.6361 _{±0.1560}		