

ALPHAMOL: RIGID BODY REPRESENTATION OF MOLECULAR STRUCTURES FOR PREDICTION OF SMALL MOLECULE GROUND-STATE

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent success of AlphaFold2 (Jumper *et al.*, 2021) in predicting structures of proteins from multiple sequence alignments (MSA) raises the question: can we generalize this approach to other important types of molecules? The positive answer to this question opens a door to overcoming the lack of structural data needed to train the model for predicting structures of RNA, proteins with non-standard amino-acids and proteins with post-translational modifications. In this work we presented a new model for predicting molecular structures, that generalizes AlphaFold2 approach to predicting structures of proteins. Two key contributions this work provides is a new representation of molecules as a collection of neighborhoods that behave as rigid bodies and a way to encode the bonds between rigid bodies into a prediction algorithm. We test this approach on the task of predicting ground-state structures of small molecules. Code available at AlphaMol repository.

1 INTRODUCTION

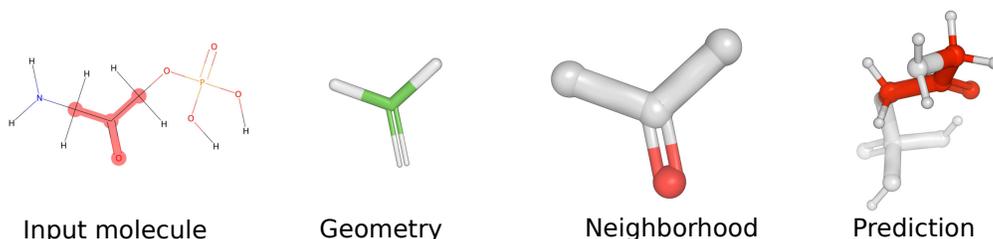
The overarching goal of computational structural biology is to develop precise predictive models capable of elucidating all chemical and structural transformations occurring within a cell. Similarly, the ultimate goal of machine learning is the development of a generally intelligent model capable of operating with comprehensive and various modalities of physical reality. The cutting edge research is currently centered on the creation of multi-modal large language models, leveraging diverse data sources such as audio, images, or even active robotic arm manipulators Driess *et al.* (2023); Huang *et al.* (2023). These seminal studies underscore the remarkable capacity of models trained on cross-modal datasets to transfer knowledge to novel tasks. Over the past decade, the application of machine learning has driven significant advancements in computational structural biology. Notably, the recent breakthrough in protein structure prediction by AlphaFold2 Jumper *et al.* (2021) is based on the transformer architectures which have emerged as dominant architectures in deep learning. Subsequent research focused on refining and extending AlphaFold2’s capabilities, for example, predicting of protein-protein complexes Evans *et al.* (2021), utilizing single sequences as inputs instead of multiple sequence alignments (MSA) Lin *et al.* (2022), and implementing diffusion-based learning Ingraham *et al.* (2022); Watson *et al.* (2022). Moreover, recent successes in predicting the structures of non-coding RNA Pearce *et al.* (2022) and elucidating protein-ligand interactions based on AlphaFold2 descriptors Hekkelman *et al.* (2023) highlight the possibility of constructing a multi-modal model that works across various molecular domains. Given physico-chemical diversity of macromolecules (proteins, nucleic acids) as well as small molecules, the researches have also attempted to unify atomic description of aminoacid residues, drug-like organic molecules (ligands), nucleic acids and non-standard residues Krishna *et al.* (2023). However, in practice, protein and nucleic acid representations are typically retained as sequences of tokens, while ligands and non-standard amino acids are represented as atom graphs. This problem forced the authors to perform atomization of residue tokens by representing aminoacids or nucleotides as ligands. On the other hand, one of the remarkable features of AlphaFold2 is the coarse-grain protein representation, where each amino acid residue is described as a rigid body, while the inner degrees of freedom are predicted separately. This coarse-graining method allowed extracting the orientation of aminoacid residues from MSA, that plays a big role during the co-evolution of the neighboring protein residues, and used in the structural module of AlphaFold2 to predict structures more accurately compared to the

054 other methods. In this work we provide foundation for consistent model architecture that allow one to
055 eliminate token-based representation of polymers, including proteins and nucleic acids. We proposed
056 a novel representation of an arbitrary type molecule as a set of rigid bodies with additional constraints
057 and demonstrated the utility of this representation by considering one of the most complex chemical
058 type, namely, small molecular structures. We further improved the Evoformer block architecture
059 from AlphaFold2, such that it incorporates explicit positions of rigid bodies eliminating the need
060 for a separate structural module. We have tested our method on the task of predicting ground-
061 state molecular structures using the curated dataset Molecule3d Xu *et al.* (2021c) extracted from
062 PubChemQC Nakata and Shimazaki (2017) of $4 \cdot 10^6$ molecules. Our model achieves the average root
063 mean square deviation (RMSD) of 0.83 over $7 \cdot 10^5$ molecules in the testing subset of the Molecule3d
064 dataset, outperforming RDKit Landrum (2020) and other published methods. More specifically, the
065 mean average error of the predicted distance matrix between the heavy atoms of the molecules is
066 0.30 versus 0.53 for the RDKit ETKDG algorithm Riniker and Landrum (2015) and 0.66 for the
067 DeeperGCN-DAGNN + DistanceXu *et al.* (2021c).
068
069
070

071 1.1 PREVIOUS WORK

072

073
074 Existing methods to predict small molecule structures typically rely on graph-based methods, which
075 treat molecules as a graph of nodes (atoms) connected by edges (bonds). One of the first such works
076 is CVGAE Mansimov *et al.* (2019), however, the quality of structure predictions was impractical and
077 required further optimization using molecular dynamics force fields. Another work, GraphDGSimm
078 and Hernández-Lobato (2019), extends the molecular graph to the second and third atomic neighbors
079 and transforms the graph into a distance matrix. Such graph extension was needed to fix the dihedral
080 angles in the molecule. The prediction target of this method is the extended distance matrix of the
081 conformer, and the 3D structure is then obtained by a non-differentiable Euclidian distance geometry
082 method (EDG, Riniker and Landrum (2015)). This method relies on a standard message-passing
083 neural network in conjunction with a conditional variational autoencoder and works entirely in the
084 internal coordinates of a molecule. Thus, it generates predictions that are invariant with respect to
085 rotation and translation. The following work CGCFXu *et al.* (2021b) improved training by using
086 energy-based learning and neural ordinary differential equation approach, while using the same
087 molecular representation. The next improvement, ConfVAE Xu *et al.* (2021a), was achieved by
088 allowing gradient propagation through the EDG step using bilevel programming approach. This was
089 the first method that allowed end-to-end training coupled with translational and rotational invariance.
090 The ConfGF method Shi *et al.* (2021) utilized a different approach: they estimate the gradient of the
091 logarithm of the probability distribution of the interatomic distances using an energy-based method
092 and then derive an atomic gradient field using the chain rule. This work proposed the first equivariant
093 end-to-end differentiable structure prediction algorithm for small molecules. However, because of the
094 difficulties in training energy-based models, one forced to additionally sample structures perturbed
095 using Gaussian noise. The other approach for end-to-end differentiable molecular structure prediction,
096 GeomolGanea *et al.* (2021), directly predicts local neighborhoods of atoms and then assembles the
097 whole structure by predicting torsion angles along the shared bonds between neighborhoods. Of note,
098 this was the first method, that accounted for enantiomers of a molecule and still invariant with respect
099 to the rotations and translations. One of the most recent methods, DMCG Zhu *et al.* (2022), that
100 directly predicts the atomic coordinates, was the first to consider molecular graph isomorphisms in
101 the loss function. Another branch of methods rely on diffusion-based learning to directly generate
102 coordinates of atoms in molecules started by GeoDiffXu *et al.* (2022). SDEGen Zhang *et al.* (2023)
103 applies diffusion to the distance matrices and Jing *et al.* (2022) in the torsional angle space. Finally,
104 EC-Conf Fan *et al.* (2023) improves the sampling efficiency of the diffusion method for this task.
105 Despite the considerable efforts that went into this field of machine learning, it is important to note,
106 that the utility of the methods and evaluation metrics are still questionable. For example, Gengmo
107 Zhou Zhou *et al.* (2023) and coauthors showed that the standard EDKTG algorithm implemented in
RDKit with minor modifications outperforms all the previously mentioned methods on the GEOM-
QM9 and GEOM-DRUGS datasets. Appendix Table 2 provides a list of the previously developed
methods for small molecular structure prediction.



117 Figure 1: The outline of our approach to the problem of molecular structure prediction. For each
118 atomic neighborhood in a molecule, we construct a rigid body, that corresponds to the fixed bond
119 geometry and inter-atomic bond distances. The model then predicts rotations and translations of
120 neighborhoods.

121 122 2 METHODS

123
124 This section is organized as follows. First, we introduce one of the key distinctions of our approach,
125 which is representation of a molecule as a set of rigid neighborhoods. Then we describe featurization
126 of this molecular representation. This is followed by addressing the symmetry and chirality problem
127 in prediction of spatial transforms of rigid bodies. Next, we touch on the equivariance of the spatial
128 transforms predictions. Finally, we describe the model architecture, introducing a modified Evoformer
129 block, and the loss functions.

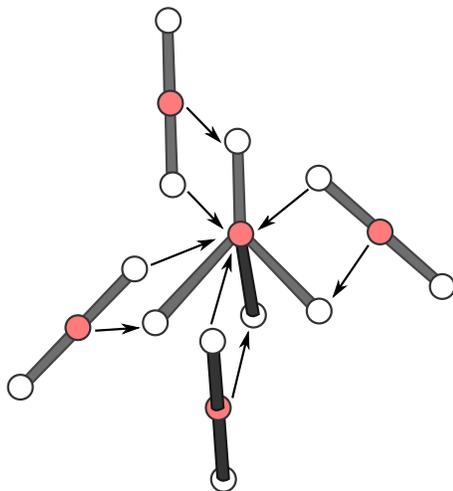
130 131 2.1 MOLECULE REPRESENTATION

132
133 The idea is to decompose the input molecule as a set of atomic neighborhoods, and for each atomic
134 neighborhood to use a predefined bond geometry and atomic distances, such that one can explicitly
135 compute the atomic coordinates of a neighborhood. Given this representation, the model is trained
136 to predict the rotation and translation of the whole neighborhood by matching the correct structure.
137 Figure 1 shows an example of a rigid neighborhood. To construct the rigid bodies constituting the
138 atomic neighborhoods we develop an algorithm that combine interatomic distances and geometries
139 (see Appendix section A.1.), resulting in reconstruction of molecular neighborhoods with acceptable
140 precision (Figure 1). We observed that local bond geometries can be rigorously approximated by the
141 rigid bodies, except for some isolated cases (see Appendix Table 6). We classified bond geometries
142 based on the hybridization of valent electron orbitals in the central atom and a set of bonds between
143 the central atom and its neighbors. This entails that each geometry has a symmetry, which we compute
144 by enumerating permutations of bonds and realigning them to the initial geometry. We observed, that
145 only 22 different geometries are enough to cover most of the molecules in the Molecule3D dataset.
146 Table 6 shows examples of the extracted arrangements and the RMSD distributions for the molecules
147 from the dataset. Although, for some molecular structures neighborhoods deviate from the standard
148 22 geometries and average atomic distances (see Appendix Figure 10, row 4), these examples are
149 rare and the algorithm is able to correct such small deviations (See Appendix section A.2.4).

150 151 2.2 FEATURES

152 Similarly to the AlphaFold2 model, the features we passed into the model are divided into single and
153 pairwise ones. Single features encode each rigid body separately, while pairwise features encode
154 spatial relationships between the rigid bodies. The algorithm to calculate the features is listed in
155 Appendix Table 8. Each neighborhood is described by the ordered set of vectors corresponding to
156 the concatenated central atom features, bond features, neighbor atom features, and neighbor atom
157 coordinates in the neighborhood. We embed each neighborhood using a 4-layer fully connected
158 neural network. Note, that the pairwise features consist of shared bond features between the two
159 neighborhoods and alignment matrices and translations between them. The necessity of including
160 specific transformation between the pairs of molecular neighborhoods is illustrated in Figure 2. If the
161 geometric information about the specific bond shared between the two neighborhoods is not passed
to the model, it leaves room for ambiguity and degrades model predictions. To obtain the alignment
matrix R_{ij} from a neighborhood i to the neighborhood j , we compute the matrices H_{ij} that align the

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177



178
179
180
181
182
183
184
185
186
187
188
189
190
191

Figure 2: Example of some of the possible pairings between two neighborhoods. Red circles are the neighborhood central atoms.

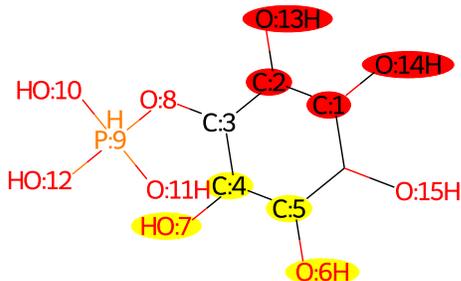


Figure 3: Example of an isomorphism of a molecule. Red and yellow highlight atoms swapping which the molecule graph does not change.

shared bond parallel to the x-axis. The alignment matrix for neighborhood i along the shared bond of neighborhood j is then:

$$R_{ij} = H_{ji}^T H_{flip} H_{ij} \quad (1)$$

, where H_{flip} is the matrix that flips x-axis. Similarly, we compute alignment translations t_{ij} between the neighborhood pairs. These features depend only on the bond geometries and distances described in the section 2.1. Appendix section A.2 provides detailed description of the data processing pipeline and Appendix algorithms 9 and 10 shows the computation of the R_{ij} matrices and t_{ij} vectors.

192
193

2.3 SYMMETRY AND CHIRALITY

194
195
196
197
198
199
200

To predict rotations and translations of rigid bodies corresponding to a molecule, one has to take into account permutation symmetries of the molecule graph. Figure 3 shows a molecule, where swapping atoms 1,2,13,14 with the corresponding atoms 5,4,7,6 does not change the molecular graph. Therefore, for example the neighborhood comprises C:2, O:13, C:1, and C:3 may correspond to two different rotations and translations. During the training process of our model, we generate all possible isomorphisms and then compute different combinations of target rotations and translations of neighborhoods.

201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

The chirality of molecules strongly affects their properties and their interactions with biological molecules, therefore it has to be addressed carefully. In our case the same bond geometries with two different chiralities do not align. Therefore one needs to specify chirality of each neighborhood to get the correct bond geometry for prediction. During the training, we extract this data from the target structures using standard RDKit utilities. First, we assign ranks to atoms using CIP rules Cahn *et al.* (1966). Then for each neighborhood central atom, we arrange its neighbors in the order defined by the atomic rank (see Appendix algorithm 6). We assign a chiral label as the sign of the volume built on the vectors, formed by the bonds from the central atom of a neighborhood. During the extraction of the bond geometries, we verify that all geometries and allowed permutations have the positive chiral label (see Appendix algorithm 2). During the construction of the neighborhoods (see in Figure 1), we reflect the bonds with respect to the central atom, if the chirality of the neighborhood does not correspond to the label assigned by the RDKit (see Appendix algorithms 9 and 11). Note, that for some molecules, an isomorphism can swap atoms of a chiral center. If such permutation changes the chirality, we can choose it arbitrarily at the cost of the isomorphism. In this case, we fix the chirality and choose all the isomorphisms for which the neighborhood aligns with the ground truth structure (see Appendix algorithm 11). Also note, that for some molecules the number of isomorphisms is exceedingly large to store the coordinates of atoms for each one of them. Therefore, we factorize

isomorphisms into local and global ones, where local isomorphisms act only on one neighborhood, leaving every other atom in place (see Appendix section A.2.2).

2.4 EQUIVARIANCE

The rotations and translations predicted by the model should be equivariant, that is they should rotate or translate when the input is rotated or translated accordingly, formally written as:

$$F(g \circ x) = r \circ F(x)$$

where F is the prediction model, g and r are the elements of the SE(3) group: the group of all rotations and translations of the 3-dimensional space. Here, we used a simplified definition of the equivariance, where $g = r$ and used rotation matrices and translation vectors as inputs and outputs.

Previous approaches to the construction of equivariant neural networks can be broadly separated in two classes: generally or specific equivariant networks. The first class of approaches can represent any F that can be learned by a neural network. The methods belonging to this class either represented tensors as a decomposition into spherical harmonics Thomas *et al.* (2018) or require lifting R^3 space into high-dimensional Lie group space Hutchinson *et al.* (2021). Both of these approaches require recomputing decomposition into spherical harmonics which can be numerically unstable and computationally costly. The second class of approaches to equivariance is exemplified by the SchNet Schütt *et al.* (2017) and AlphaFold2 structure module. These methods include equivariant operations in the initial R^3 space into the model, which limits the function F such a model can learn. Recently, a new approach emerged in the class of generally equivariant neural networks Du *et al.* (2022a); Wang and Zhang (2022), where the key idea is to construct a set of reference frames and project input tensors onto them. The predictions are then given as a set of decomposition coefficients that are transformed back into tensors. The seminal work by Du *et al.* Du *et al.* (2022a) dealt with systems composed of sets of points, thus, the frame construction may lead to numerical instability in case of points being close in space or belonging to the same plane. Here, we work with a system comprised of rigid bodies with natural frames of reference associated with them, therefore avoiding the aforementioned problem.

Let $X = (R_1, T_1, \dots, R_N, T_N) \in R^{12N}$ be a many-body system embedded into R^3 space, where N is the number of rigid bodies. For rigid body i , we use $R_i(t) \in R^9$ and $T_i(t) \in R^3$ to denote its rotation and translation at iteration t , respectively. The rows of the matrix $B_i^k = R_i = (a_i^1, a_i^2, a_i^3) \in R^3$ correspond to the basis vectors of each local frame associated with the neighborhood i (these local frames make a complete basis for vectors). Additionally, we construct the basis for rotation matrixes using the approach outlined by Du *et al.* Du *et al.* (2022a), namely $B_{ij}^{kl} = a_i^k \otimes a_j^l \in R^9$, where $k, l \in 1, 2, 3$ and $i, j \in 1 \dots N$. To preserve translation equivariance we first centralize the translations $T_i^{(c)} = T_i - \frac{1}{N} \sum_i T_i$. Then we take each vector input for the rigid body i (like $T_i^{(c)}$) and compute its decomposition coefficients in the basis B_i^k . Similarly, for the input rotation matrixes defined on the pairs of rigid bodies (see R_{ij} in section 2.2) we use a basis composed of matrixes B_{ij}^{kl} . Additionally, we have input vectors defined on the pairs of neighborhoods (t_{ij} see section 2.2), that we have to scalarize. To avoid constructing a new basis we instead transform these vectors into tensors: $t_{ij}^{(t)} = t_{ij} \otimes t_{ij}$ and then decompose them into basis B_{ij}^{kl} . The resulting scalars then can be used in a neural network without breaking the SE(3)-equivariance.

To obtain updated rotations and translations of the neighborhoods we predict translation and rotation update vectors and matrixes, correspondingly. For the vector prediction, we interpret predictions as the decomposition coefficients in the basis B_i^k . Thus, translations updates are computed as $dT_i = \sum_k F_{ik} B_i^k$, where F_{ik} is the k -th network output for i -th neighborhood. In principle, one can use the same vectorization procedure for computing rotation matrixes updates using the basis B_{ij}^{kl} . However, in our case, the space of all possible rotation matrixes forms a manifold in the space of all possible decompositions in the basis B_{ij}^{kl} . On the other hand, we need to predict only a valid rotation matrix, that is orthogonal matrix. For this, we use the 6D-rotation representation first proposed in the work Zhou *et al.* (2019), where the idea is that any rotation matrix can be represented using just two vectors (b_1, b_2) and to construct the complete orthogonal matrix from these vectors one follows the

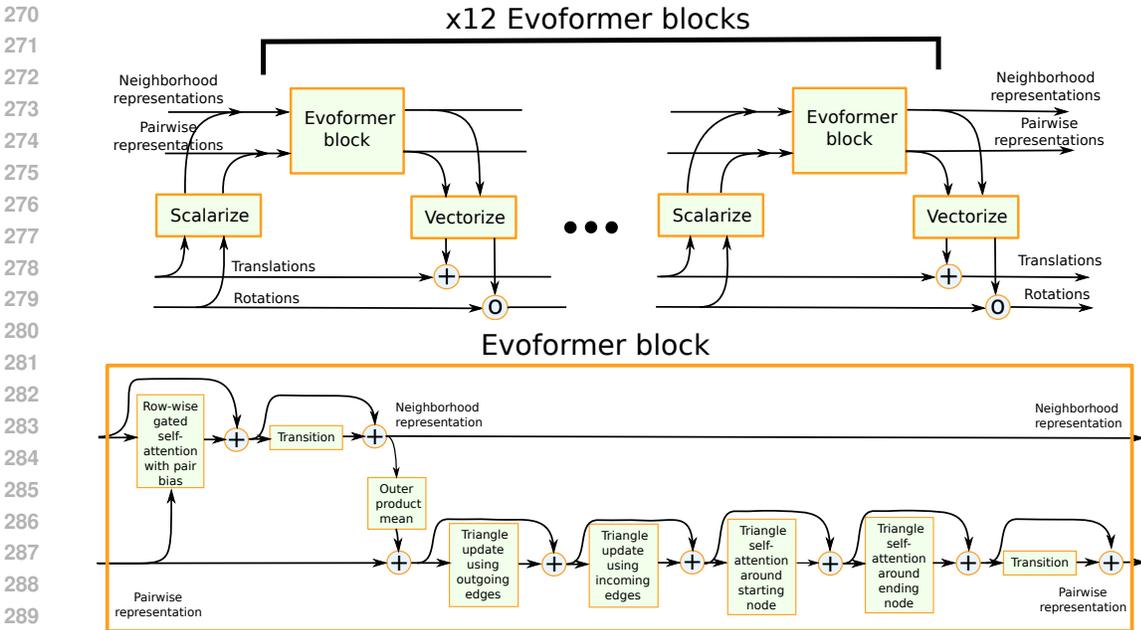


Figure 4: Schematic representation of the AlphaMol model along with the modified Evoformer block. Scalarization and vectorization blocks are described in section 2.4. "+" symbol denotes component-wise sum of feature tensors and "o" symbol denotes matrix multiplication.

Graham-Schmidt process:

$$GS \left(\begin{bmatrix} | & | \\ b_1 & b_2 \\ | & | \end{bmatrix} \right) = \begin{bmatrix} | & | & | \\ c_1 & c_2 & c_3 \\ | & | & | \end{bmatrix} = \begin{cases} c_1 = N(b_1) \\ c_2 = N(b_2 - (c_1 \cdot b_2)c_1) \\ c_3 = c_1 \times c_2 \end{cases} \quad (2)$$

, where N is the vector normalization operation. Using this rotation representation we need to predict two vectors instead of a matrix, so we can use the same basis B_i^k , that we employed for predicting translation updates. The detailed description of vectorization and scalarization algorithms developed in this work presented in Appendix section A.3, particularly, Appendix algorithms 15 - 18.

2.5 MODEL

The model consists of sequential Evoformer blocks that predict rotations and translations updates to the initial rotations and translations of the molecular neighborhoods (see Figure 4). Each updated transform is computed as follows:

$$T_i = T_i^{pred} + T_i^{prev} \quad (3)$$

$$R_i = R_i^{pred} \circ R_i^{prev} \quad (4)$$

, where subscripts *pred* and *prev* corresponds to predicted transforms from each Evoformer module and previous total transform, respectively; and i denotes the neighborhood index. We used identity matrices and zero vectors as the initial ones. Note, that the molecular representation comprises pairwise features to encode relations between the rigid neighborhoods and the single features to encode characteristics of the rigid neighborhoods itself. Therefore, the operations within the Evoformer block, treating the prediction of molecular structures as a graph inference problem in R^3 , must reflect chemically feasible molecular geometries in the R^3 space. We would like to note, that operations on the pairwise representations are one of the key innovations introduced in AlphaFold2: they are based on the intuition, that the pairwise representations contain information that must satisfy the triangle inequality on distances. Therefore the corresponding update functions operate on triangles of edges involving three different nodes; the missing edge of the triangle is included using logit bias to axial attention coupled with the multiplicative update, which uses two edges to update the missing third

edge. The axial attention is also modeled using the bias to the row-wise attention, that is coming from previous pairwise representations. This completes the information flow from pairwise representations to individual molecular neighborhoods.

To improve the training stability of the model we followed the recent observation by Shuangfei et al Zhai *et al.* (2023) that training instability is usually accompanied by the low entropy of the attention layer. We implemented σ Reparam algorithm to regularize attention weights in all the layers that use attention in the Evoformer block (Appendix algorithm 19). Additionally, following recent efforts by the Ziyao et al Li *et al.* (2022), we replaced all ReLU activation units with Gaussian Error Linear Units (GELU) and added postprocessing layer to the output of the OuterProductMean module.

2.6 LOSSES

We used several loss terms to improve predictions and training stability of the model. To compare predicted structures with the ground truth ones we use frame-aligned point error loss (FAPE) first proposed by the AlphaFold2 team. Given a set of predicted coordinates x_i , $i \in 1, \dots, N$ and reference frames F_j , $j \in 1, \dots, M$, along with a set of true coordinates x_i^{gt} and reference frames F_j^{gt} , the loss is computed by transforming coordinates of each point x_i into the frame F_j and comparing them to the corresponding ground truth coordinates x_i^{gt} in the ground truth frames F_j^{gt} :

$$FAPE(F, x, F^{gt}, x^{gt}) = \frac{1}{NM} \sum_{ij} \|F_i \circ x_j - F_i^{gt} \circ x_j^{gt}\|$$

, where N is the number of atoms in a molecule and M is the number of rigid bodies comprising it (see Appendix algorithm 21). The detailed description of the whole algorithm is given in Appendix section A.4.1 and Appendix algorithm 22. The second loss term penalizes the clashes of atoms in the predicted structure. To construct it, we used Van-der-Waals radii for each atom Mantina *et al.* (2009) and assign minimum distance between any two atoms in a molecule as $r_{ij}^{min} = r_i^{VW} + r_j^{VW}$. However, we want to exclude atoms that belong to the same neighborhoods from the loss. Therefore, we set $r_{ij}^{min} = 0$ if atoms i and j are the second-order neighbors in the molecular graph (Appendix algorithm 23). After predicting atomic coordinates for a molecule we calculate the clash penalty as:

$$L_{clash}(x) = \frac{1}{\sum_{ij} r_{ij}^{min} > 0} \sum_{ij} ReLU(r_{ij}^{min} - \|x_i - x_j\|)$$

It is also worth to note, that some atomic positions are predicted more than once, because of the overlap of rigid neighborhoods along the covalent bonds of a molecule (See Fig. 2). Therefore, we averaged each atomic position across all the rigid neighborhoods containing given atom (Appendix algorithm 20). Finally, during the prediction of rotation matrices (Eq.2) one can face co-linearity problem for the vectors b_1 and b_2 , which are used to parameterize rotation. Hence, we added the co-linearity loss, that minimizes scalar product between the normalized vectors b_1 and b_2 (Appendix algorithm 24). One of the key features of AlphaFold2 is its ability to predict the quality of the resulting structures, which we aim to retain in our work. For this, we predict the per-neighborhood accuracy of the structure (pLDDT) using a small neural network from the final neighborhood representations (Appendix algorithm 25). All the losses, except for pLDDT are predicted for each iteration of the Evoformer block and then averaged. To take into account isomorphisms of the molecular graph, we compute the minimum FAPE and pLDDT losses over all the isomorphisms.

3 RESULTS

The dataset used for training and test is based on PubChemQC Nakata and Shimazaki (2017) consisting of approximately 4 million molecules, where each molecule is represented with Simplified Molecular Input Line Entry Specification (SMILES) description, IUPAC International Chemical Identifier (InChI), and the ground-state and excited-state 3D geometries of these molecules. Here we used the pre-processed version of this dataset, Molecule3D, provided by Xu *et al.* (2021c). The dataset is split into 60% / 20% / 20% subsets for training/validation/test, and two different splitting strategies is used: random and scaffold-based splitting, where a scaffold refers to a molecule’s core component consisting of connected rings without branches. Note, that scaffold-based split leads to a distribution shift between the training and test subsets. We used the scaffold-based split, which forces a model to

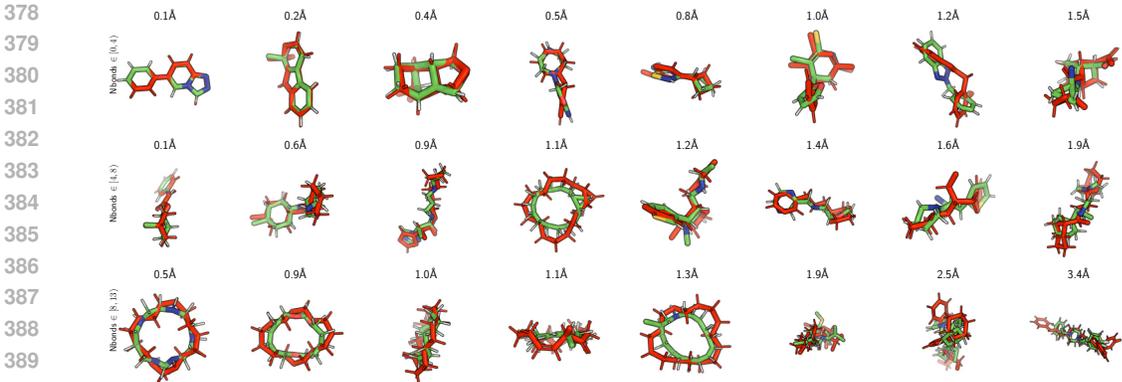


Figure 5: Examples of predicted structures from the test set. Red: true structures; Blue: predicted structures. We split the test set into sets of molecules depending on the number of the rotatable bonds and the RMSD of the prediction. Then for each number of rotatable bonds, we split RMSD region into 8 equal intervals and picked one example for each.

Table 1: The evaluation of the trained models on the test set of Molecule3D benchmark.

Model	MAE	RMSD	FAPE	HOMO-LUMO
AlphaMol/8	0.365	0.886	1.565	0.435
AlphaMol/12	0.353	0.893	1.509	0.337
AlphaMol/24	0.304	0.830	1.434	0.382
RDKit ETKDG	0.532	-	-	0.1524
DeeperGCN-DAGNN + Distance	0.660	-	-	0.2000
DeeperGCN-DAGNN + Coordinates	0.763	-	-	0.2371

capture such distribution shifts in chemical space and measures the out-of-distribution generalization ability of the model. Additionally, we filtered the dataset to exclude molecules containing only one neighborhood, molecules, whose graph has disconnected components, molecules that have pentavalent atoms, and some other cases (see Appendix section A.2). In total, we excluded $\approx 4,000$ structures across the training, validation, and test sets.

We have trained models with different number of Evoformer blocks: small (8), medium(12), and big(24). The model was trained using Adam optimizer with the learning rate $1.5 \cdot 10^{-4}$, without a learning rate schedule. The training of each model was carried out on one node with 4xV100 for $1 \cdot 10^6$ iterations, which approximately correspond to 5, 7 and 15 days of node-time for small, medium and big models. To compare our results with the previous algorithms we compute the mean absolute error (MAE) performance metric:

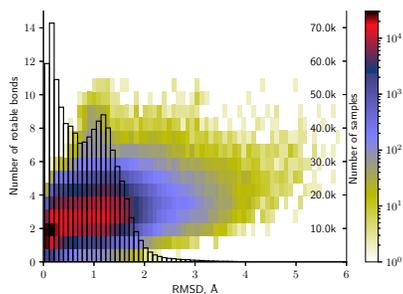
$$MAE = \frac{1}{N^2} \sum_{i,j=1..N, i \neq j} |d_{ij}^{pred} - d_{ij}^{data}| \quad (5)$$

where d_{ij}^{pred} is the distance between atom i and atom j in the prediction and d_{ij}^{data} is the actual distance between the same atoms.

We compared our models with the DeeperGCN-DAGNN model Liu *et al.* (2021) model, which predicts either distances between atoms or 3D atomic coordinates Xu *et al.* (2021c). Additionally, we used the ETKDG algorithm Riniker and Landrum (2015) implemented in RDKit Landrum (2020) as a baseline. The results are shown in Table 1. One can see, that even the smallest model outperforms previously published methods in predicting 3D structures of ground states of molecules.

The existing methods typically report degradation of the prediction quality with respect to the number of rotatable bonds in a molecule. We did not observed such a drawback for our method, and Figure 6 shows the distribution of RMSD of the predicted structures in the test set versus the number of rotatable bonds. As one can see, the RMSD distribution has the second peak at the RMSD value of \sim

432
433
434
435
436
437
438
439
440
441
442



443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Figure 6: Left axis: distribution of RMSD of predicted structures depending on the number of rotatable bonds. Color scheme uses logarithmic scale. Right axis: summary distribution of RMSD scores.

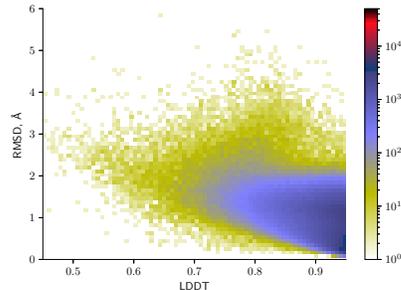


Figure 7: Distribution of RMSD of predicted structures depending on average predicted LDDT of the molecule. Color scheme uses logarithmic scale.

1.3 Å corresponding to molecules with up to 11 rotatable bonds, while the worst case predictions have about 7 rotatable bonds. Figure 5 shows examples of the predicted molecular structures sampled from three groups with different number of rotatable bonds and 8 groups corresponding to the different RMSD values.

It is important to note, that one of the advantages of our model is ability to estimate the confidence of the predictions, which could be useful in the downstream tasks. Figure 7 shows the correlation between the average predicted LDDT and the RMSD of the predicted structure. Expectedly, we observed a negative correlation between the predicted LDDT and RMSD.

4 DISCUSSION

In this study we developed a method to predict ground state small molecule structures based on novel representation of molecular structures as a set of rigid neighborhoods. We introduced how to compute loss functions over all isomorphisms of a molecular structure, as well as how to fix the chirality of the predicted structures. These features are especially relevant for biological molecules; for example, chirality may imply very different bioactivity properties between the small molecule enantiomers and structural properties of protein chains. To the best of our knowledge, the proposed approach of factorizing isomorphisms of the molecules is the only one available that can deal with this long tail molecules. For example, lipids have long tails of carbohydrates for which other approaches to account for isomorphisms fail, because their number grows exponentially with the length of the molecule. To successfully train the model we introduced some changes to the Evoformer block, namely, σ Reparam method Zhai *et al.* (2023), scalarization approach to equivariance Du *et al.* (2022a) and Graham-Schmidt process Zhou *et al.* (2019) to represent rotations. These changes allowed us to exclude the learning rate scheduling and running exponential averaging and significantly stabilized the training.

REFERENCES

- Baek, M. *et al.* (2021). Accurate prediction of protein structures and interactions using a three-track neural network. *Science*, **373**(6557), 871–876.
- Cahn, R. S. *et al.* (1966). Specification of molecular chirality. *Angewandte Chemie International Edition in English*, **5**(4), 385–415.
- Chan, L. *et al.* (2020). Bokei: Bayesian optimization using knowledge of correlated torsions and expected improvement for conformer generation. *Physical Chemistry Chemical Physics*, **22**(9), 5211–5219.
- Driess, D. *et al.* (2023). Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*.

- 486 Du, W. *et al.* (2022a). Se (3) equivariant graph neural networks with complete local frames. In
487 *International Conference on Machine Learning*, pages 5583–5608. PMLR.
- 488
- 489 Du, Y. *et al.* (2022b). Molgensurvey: A systematic survey in machine learning models for molecule
490 design. *arXiv preprint arXiv:2203.14500*.
- 491 Evans, R. *et al.* (2021). Protein complex prediction with alphafold-multimer. *bioRxiv*, pages 2021–10.
492
- 493 Fan, Z. *et al.* (2023). Ec-conf: A ultra-fast diffusion model for molecular conformation generation
494 with equivariant consistency. *arXiv preprint arXiv:2308.00237*.
- 495 Ganea, O. *et al.* (2021). Geomol: Torsional geometric generation of molecular 3d conformer
496 ensembles. *Advances in Neural Information Processing Systems*, **34**, 13757–13769.
497
- 498 Hekkelman, M. L. *et al.* (2023). Alphafill: enriching alphafold models with ligands and cofactors.
499 *Nature Methods*, **20**(2), 205–213.
- 500 Huang, S. *et al.* (2023). Language is not all you need: Aligning perception with language models.
501 *arXiv preprint arXiv:2302.14045*.
502
- 503 Hutchinson, M. J. *et al.* (2021). Lietransformer: Equivariant self-attention for lie groups. In
504 *International Conference on Machine Learning*, pages 4533–4543. PMLR.
- 505 Ingraham, J. *et al.* (2022). Illuminating protein space with a programmable generative model. *bioRxiv*,
506 pages 2022–12.
507
- 508 Jing, B. *et al.* (2022). Torsional diffusion for molecular conformer generation. *Advances in Neural*
509 *Information Processing Systems*, **35**, 24240–24253.
- 510 Jumper, J. *et al.* (2021). Highly accurate protein structure prediction with alphafold. *Nature*,
511 **596**(7873), 583–589.
512
- 513 Krishna, R. *et al.* (2023). Generalized biomolecular modeling and design with rosettafold all-atom.
514 *bioRxiv*, pages 2023–10.
- 515 Landrum, G. (2020). Rdkit: Open-source cheminformatics.
516
- 517 Li, Z. *et al.* (2022). Uni-fold: an open-source platform for developing protein folding models beyond
518 alphafold. *bioRxiv*, pages 2022–08.
- 519 Lin, Z. *et al.* (2022). Evolutionary-scale prediction of atomic level protein structure with a language
520 model. *bioRxiv*, pages 2022–07.
521
- 522 Liu, M. *et al.* (2021). Fast quantum property prediction via deeper 2d and 3d graph networks. *arXiv*
523 *preprint arXiv:2106.08551*.
- 524 Luo, S. *et al.* (2021). Predicting molecular conformation via dynamic graph score matching. *Advances*
525 *in Neural Information Processing Systems*, **34**, 19784–19795.
526
- 527 Mansimov, E. *et al.* (2019). Molecular geometry prediction using a deep generative graph neural
528 network. *Scientific reports*, **9**(1), 20381.
- 529 Mantina, M. *et al.* (2009). Consistent van der waals radii for the whole main group. *The Journal of*
530 *Physical Chemistry A*, **113**(19), 5806–5812.
531
- 532 Nakata, M. and Shimazaki, T. (2017). Pubchemqc project: a large-scale first-principles electronic
533 structure database for data-driven chemistry. *Journal of chemical information and modeling*, **57**(6),
534 1300–1308.
- 535 Pearce, R. *et al.* (2022). De novo rna tertiary structure prediction at atomic resolution using geometric
536 potentials from deep learning. *bioRxiv*, pages 2022–05.
537
- 538 Riniker, S. and Landrum, G. A. (2015). Better informed distance geometry: using what we know
539 to improve conformation generation. *Journal of chemical information and modeling*, **55**(12),
2562–2574.

- 540 Schütt, K. *et al.* (2017). Schnet: A continuous-filter convolutional neural network for modeling
541 quantum interactions. *Advances in neural information processing systems*, **30**.
542
- 543 Shi, C. *et al.* (2021). Learning gradient fields for molecular conformation generation. In *International*
544 *Conference on Machine Learning*, pages 9558–9568. PMLR.
- 545 Simm, G. N. and Hernández-Lobato, J. M. (2019). A generative model for molecular distance
546 geometry. *arXiv preprint arXiv:1909.11459*.
547
- 548 Thomas, N. *et al.* (2018). Tensor field networks: Rotation-and translation-equivariant neural networks
549 for 3d point clouds. *arXiv preprint arXiv:1802.08219*.
- 550 Wang, X. and Zhang, M. (2022). Graph neural network with local frame for molecular potential
551 energy surface. *arXiv preprint arXiv:2208.00716*.
552
- 553 Watson, J. L. *et al.* (2022). Broadly applicable and accurate protein design by integrating structure
554 prediction networks and diffusion generative models. *bioRxiv*, pages 2022–12.
555
- 556 Wu, J. *et al.* (2021). Se (3)-equivariant energy-based models for end-to-end protein folding. *bioRxiv*,
557 pages 2021–06.
- 558 Wu, R. *et al.* (2022). High-resolution de novo structure prediction from primary sequence. *BioRxiv*,
559 pages 2022–07.
560
- 561 Xu, M. *et al.* (2021a). An end-to-end framework for molecular conformation generation via bilevel
562 programming. In *International Conference on Machine Learning*, pages 11537–11547. PMLR.
- 563 Xu, M. *et al.* (2021b). Learning neural generative dynamics for molecular conformation generation.
564 *arXiv preprint arXiv:2102.10240*.
565
- 566 Xu, M. *et al.* (2022). Geodiff: A geometric diffusion model for molecular conformation generation.
567 *arXiv preprint arXiv:2203.02923*.
- 568 Xu, Z. *et al.* (2021c). Molecule3d: A benchmark for predicting 3d geometries from molecular graphs.
569 *arXiv preprint arXiv:2110.01717*.
570
- 571 Zhai, S. *et al.* (2023). Stabilizing transformer training by preventing attention entropy collapse. In
572 *International Conference on Machine Learning*, pages 40770–40803. PMLR.
- 573
- 574 Zhang, H. *et al.* (????). Equivariant vector field network for many-body system modeling.
- 575
- 576 Zhang, H. *et al.* (2023). Sdegen: learning to evolve molecular conformations from thermodynamic
577 noise for conformation generation. *Chemical Science*, **14**(6), 1557–1568.
- 578
- 579 Zhou, G. *et al.* (2023). Do deep learning methods really perform better in molecular conformation
580 generation? *arXiv preprint arXiv:2302.07061*.
- 581
- 582 Zhou, Y. *et al.* (2019). On the continuity of rotation representations in neural networks. In *Proceedings*
583 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5745–5753.
- 584
- 585 Zhu, J. *et al.* (2022). Direct molecular conformation generation. *arXiv preprint arXiv:2202.01356*.

586 A APPENDIX

587 A.1 BOND GEOMETRIES

589 In this work we first remove terminal hydrogens of the molecule using the algorithm 1. We do this in
590 because we want to preserve hydrogens for carbon atoms in SP3 hybridization, while keeping the
591 same number of neighborhoods as with the hydrogen-free molecular representation. In particular we
592 remove all hydrogens bonded to heavy atoms that are connected only to one other heavy atom.
593

We extract bond geometries and bond distances from the dataset using the algorithm 2.

Algorithm 1 Removing terminal hydrogens

```

594
595
596 Require: Mol
597   done ← False
598   while not done do
599     done ← True
600     for atom_i ∈ Mol.GetAtoms() do
601       if atom_i.Symbol ≠ H then
602         Neighb ← atom_i.Neighbors()
603          $N_{heavy} = \sum_{atom_j \in Neighb} \begin{cases} 1, & atom_j.Symbol \neq H \\ 0, & otherwise \end{cases}$  ▷ Number of heavy atoms
604         connected to atom_i
605          $N_{hyd} = len(Neighb) - N_{heavy}$  ▷ Number of hydrogens connected to atom_i
606         if  $N_{heavy} == 1$  &&  $N_{hyd} > 0$  then ▷ If atom_i hydrogens are terminal
607           for atom_j ∈ Neighb do
608             if atom_j.Symbol == H then
609               Mol.RemoveAtom(atom_j)
610               done ← False
611             end if
612           end for
613         end if
614       end for
615     end while
616     return Mol

```

Algorithm 2 An overview of bond geometries extraction

```

620 Require: Dataset
621   mol ∈ Dataset
622   for mol ∈ Dataset do
623     Hoods ← Neighborhoods(mol)
624     for hood ∈ Hoods do
625       HoodKey ← Cat(hybridization, num_single_bonds, num_double_bonds, ...)
626       UnitVecs ← []
627       for (i, j) ∈ HoodBonds do ▷ Index i is always the root atom of the neighborhood
628         BondKey ← Cat(mol[i].Symbol, mol[j].Symbol) ▷ mol[i].Symbol is the atomic
        symbol of i-th atom
629          $BondDist \leftarrow \sqrt{|mol[i].Coords - mol[j].Coords|^2}$ 
630         Store[BondKey] ← BondDist
631         UnitVecs.append((mol[j].Coords - mol[i].Coords)/BondDist)
632       end for
633       if  $len(UnitVecs) \geq 3$  && volume_sign(UnitVecs) < 0 then
634          $UnitVecs \leftarrow \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{pmatrix} \cdot UnitVecs$  ▷ Making sure that all geometries have
635         positive volume sign
636       end if
637       perms ← []
638       for perm ∈ Permutations(range(len(UnitVecs))) do ▷ Permutations of bond indices
639         if RMSD(UnitVecs, UnitVecs[perm]) < 0.1 then ▷ Only counting permutations
640         for one chirality
641           perms.append(perm)
642         end if
643       end for
644       Store[HoodKey] ← UnitVecs, perms
645     end for
646   end for

```

Table 2: Previous work on molecular structure prediction

Method	Year	Inner representation	Equivariance	Training framework
GraphDG Simm and Hernández-Lobato (2019)	2020	Distances	Scalars	VAE
ConfVAE Xu <i>et al.</i> (2021a)	2021	Distances	Scalars	VAE
ConfGF Shi <i>et al.</i> (2021)	2021	Distances	Scalars	Energy-based
DGSM Luo <i>et al.</i> (2021)	2021	Distances	Scalars	Energy-based
CGCF Xu <i>et al.</i> (2021b)	2021	Distances	Scalars	Energy-based, NeuralODE
SDEGen Zhang <i>et al.</i> (2023)	2023	Distances	Scalars	Diffusion
GeoMol Ganea <i>et al.</i> (2021)	2021	Angles and dist.	Scalars	OT
BOKEI Chan <i>et al.</i> (2020)	2020	Torsions	Scalars	BO
Torsional Diffusion Jing <i>et al.</i> (2022)	2022	Torsions	Scalars	Diffusion
CVGAE Mansimov <i>et al.</i> (2019)	2019	Coordinates	Fixed frame	VAE
DMCG Zhu <i>et al.</i> (2022)	2022	Coordinates	Fixed frame	VAE
EVFN Zhang <i>et al.</i> (????)	–	Coordinates	Scalarization	Energy-based
GeoDiff Xu <i>et al.</i> (2022)	2022	Coordinates	Scalars	Diffusion
EC-Conf Fan <i>et al.</i> (2023)	2023	Coordinates	Irr. repr.	Diffusion

Table 3: Previous methods that predict structures of molecules based on molecular graph description. Updated version of the survey Du *et al.* (2022b). We also classify score-matching algorithms as energy-based.

Table 4: Previous work on protein structure prediction

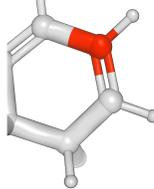
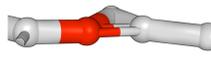
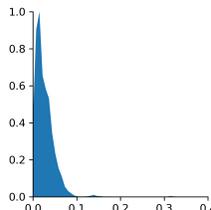
Method	Inner representation	Equivariance	Training framework
AlphaFold2 Jumper <i>et al.</i> (2021)	Rigid bodies	Invariant Point Attention	OT
RoseTTAFold Baek <i>et al.</i> (2021)	Atomic coordinates	SE(3) Transformer	OT
OmegaFold Wu <i>et al.</i> (2022)	Rigid bodies	Invariant Point Attention	OT
ESMFold Lin <i>et al.</i> (2022)	Rigid bodies	Invariant Point Attention	OT
SE(3)-Fold Wu <i>et al.</i> (2021)	Coordinates	Scalars	Energy-based
RFDiffusion Watson <i>et al.</i> (2022)	Atomic coordinates	SE(3) Transformer	Diffusion
Chroma Ingraham <i>et al.</i> (2022)	Rigid bodies	Relative transforms	Correlated Diffusion

Table 5: Previous methods that predict structures of proteins.

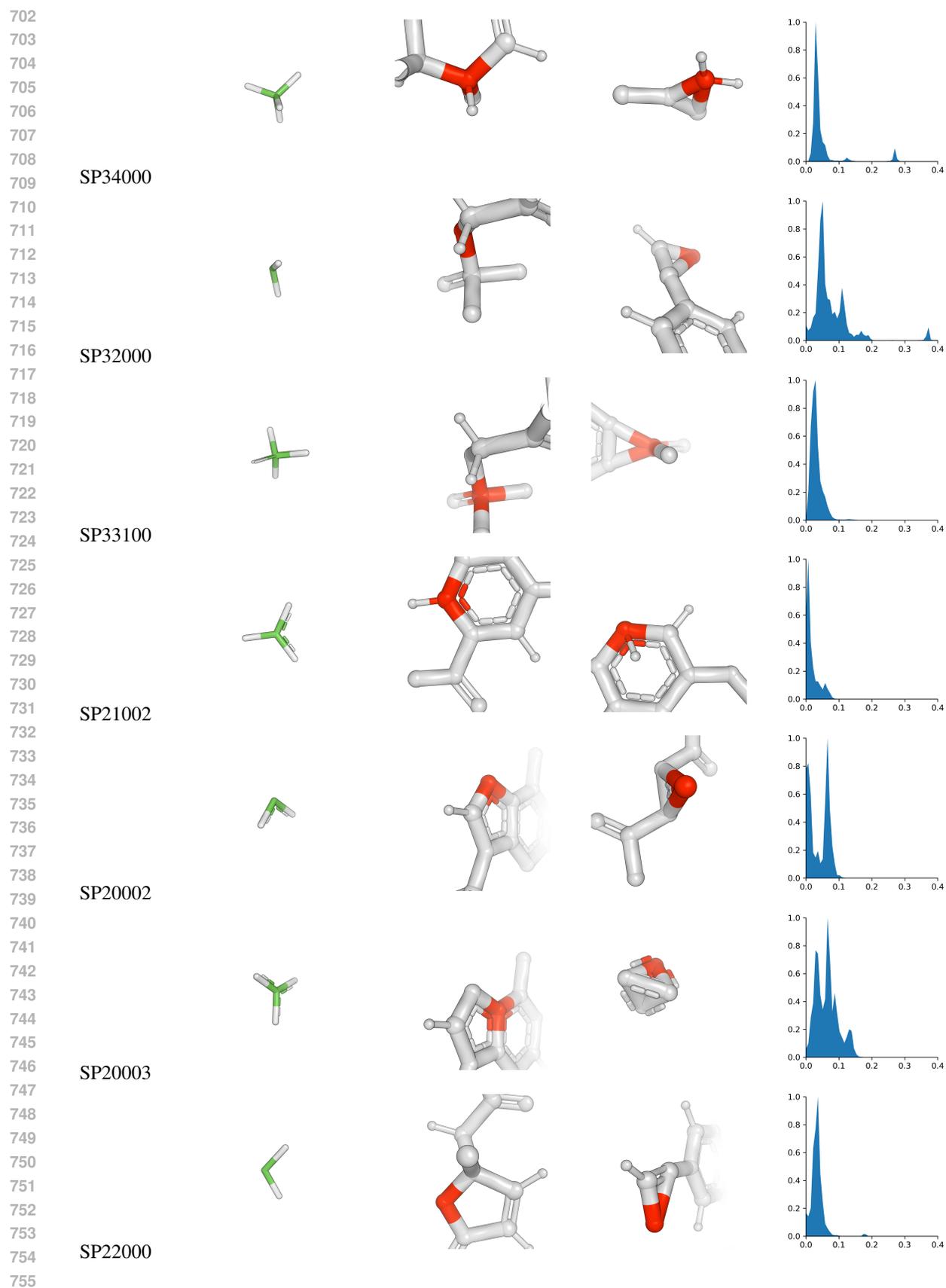
We use algorithm 3 to extract molecular neighborhoods and expression 6 to obtain the volume sign. As we can see we treat a bond geometry as an object that is equal to other geometry up to a permutation of bonds. I.e. only the hybridization of the root atom and the number of single/double/triple bonds counts. We store the first bond length and bond geometry of a particular class without averaging over all the structures in the dataset. Table 6 shows resulting geometries and the distribution of RMSD values when aligning them to the structures from the dataset. It also contains examples of aligned geometries as well as outliers in terms of alignment RMSD.

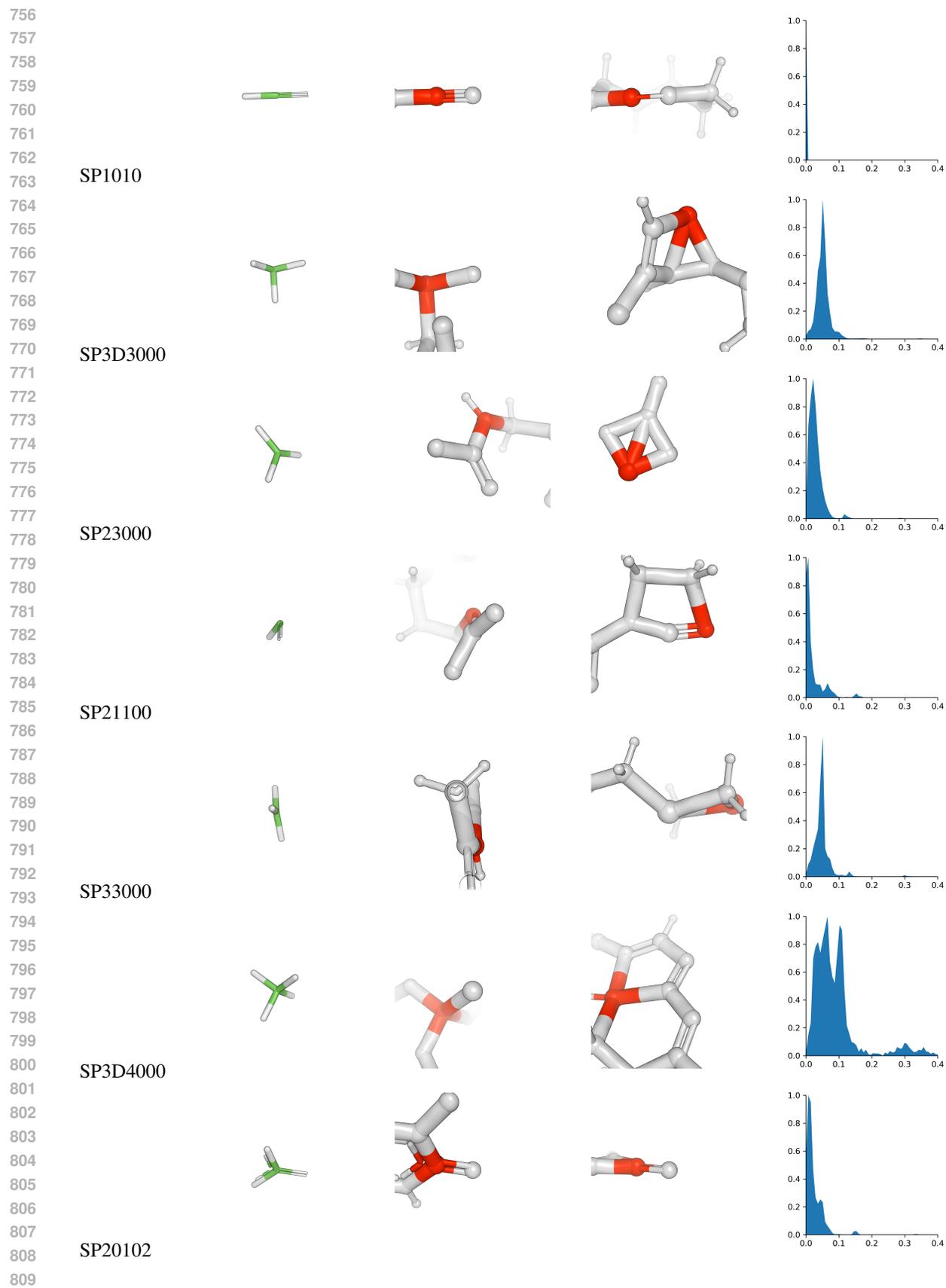
$$\begin{cases} \vec{v}_i = \vec{r}_i, & \text{3-neighbor} \\ \vec{v}_i = \vec{r}_4 - \vec{r}_i, & \text{4-neighbor} \end{cases}, i \in [1, 2, 3]$$

$$\text{volume_sign} = \text{sign}(v_1, v_2 \times v_3)$$

Signature	Geometry	Initial molecule	Outlier	RMSD distribution
				

SP22100





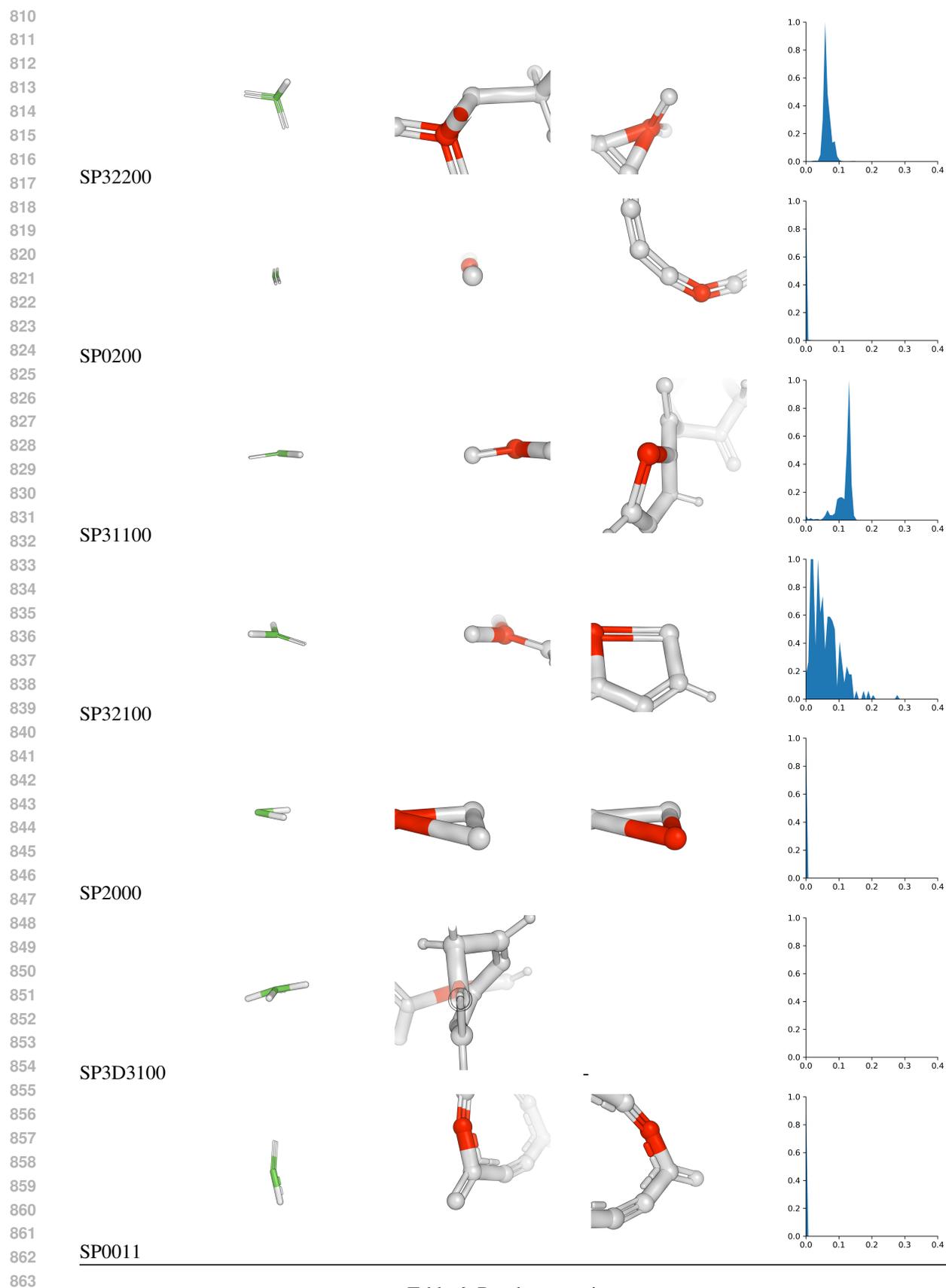


Table 6: Bond geometries.

Algorithm 3 Algorithm for extracting molecular neighborhoods

```

864 Require: mol, rank
865  $A_{ij \in [0, N)}, \leftarrow \text{Adjacency}(mol)$   $\triangleright$  Getting adjacency matrix of the molecule
866  $B_{ij \in [0, N)}, \leftarrow \text{BondIndex}(mol)$   $\triangleright$  Bond indices for adjacency matrix
867  $HoodIdx \leftarrow 0$ 
868 for  $k < N$  do
869   if  $\sum_i A_{ki} > 1$  then  $\triangleright$  Degree of the node k
870      $HoodRootAtom[HoodIdx] \leftarrow k$ 
871      $neighbors \leftarrow []$ 
872     for  $l < N$  do
873       if  $A_{kl} == 1$  then
874          $neighbors.append(l, rank[l])$ 
875       end if
876     end for
877      $neighbors \leftarrow \text{sort}(neighbors, \text{lambda } x : x[1])$   $\triangleright$  Sorting neighbors according to the
878     CIP ranks
879      $AtomIdx \leftarrow 0$ 
880     for  $l \in neighbors$  do
881        $HoodAtomIndices[HoodIdx, AtomIdx] \leftarrow l$ 
882        $HoodBondIndices[HoodIdx, AtomIdx] \leftarrow B_{kl}$ 
883        $AtomIdx \leftarrow AtomIdx + 1$ 
884     end for
885      $HoodIdx \leftarrow HoodIdx + 1$ 
886   end if
887 end for return  $HoodRootAtom, HoodAtomIndices, HoodBondIndices$ 

```

A.2 DATA PROCESSING

The data processing pipeline is outlined in algorithms 4 and 5. The pipeline 4 processes molecular graphs without any knowledge of ground truth atomic coordinates. We extract atomic and bond feature sets described in the manuscript first. Then we get the indices of atoms comprising molecular neighborhoods and rearrange per-atom features into per-neighborhood and pairwise features. Afterwards we compute initial coordinates for each neighborhood atoms according to the bonds geometries, that we described in section A.1. Afterwards we use these coordinates to get pairwise transforms between neighborhoods. And finally, we enumerate isomorphisms in the molecule.

The processing of atomic coordinates is outlined in algorithm 5. We generate coordinates for each isomorphisms of the molecular graph and extract transforms between initial neighborhood coordinates and generated ground truth coordinates giving us a set of ground truth transforms.

Algorithm 4 Pipeline of molecular graph pre-processing

```

904 Require: mol  $\triangleright$  Molecular graph
905  $F_{i \in [0, N)}, \leftarrow \text{AtomicFeatures}(mol)$   $\triangleright$  Getting atomic features
906  $B_{j \in [0, M)}, \leftarrow \text{BondFeatures}(mol)$   $\triangleright$  Getting bond features
907  $NeighbIdx_{ij} \leftarrow \text{GetNeighborhoods}(mol)$   $\triangleright$  Getting indices of atoms in neighborhoods,
908  $i \in [0, M), j \in [0, m_i]$ 
909  $HoodFeat_{i \in [0, N)} \leftarrow \text{GetHoodFeatures}(F_i, B_{ij}, NeighbIdx_{ij})$   $\triangleright$  Getting neighborhood
910 features
911  $PairFeat_{ij \in [0, N)} \leftarrow \text{GetPairFeatures}(mol)$   $\triangleright$  Getting pairwise neighborhood features
912  $x_{ij}^{init} \leftarrow \text{GetInitCoords}(mol)$   $\triangleright$  Getting initial coordinates of atoms in neighborhoods
913  $T_{ij \in [0, N)} \leftarrow \text{GetPairTransforms}(mol)$   $\triangleright$  Getting pairwise transforms of neighborhoods
914  $global\_iso, local\_iso \leftarrow \text{GetIsomorphisms}(mol)$   $\triangleright$  Getting isomorphisms of the molecule

```

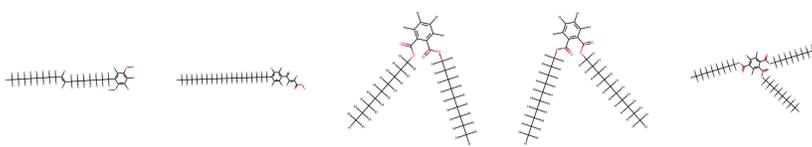
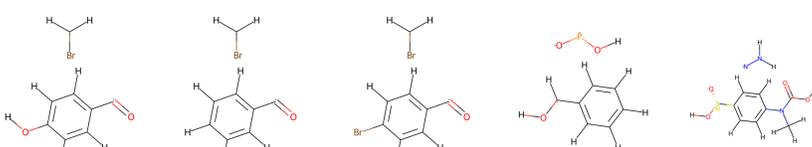
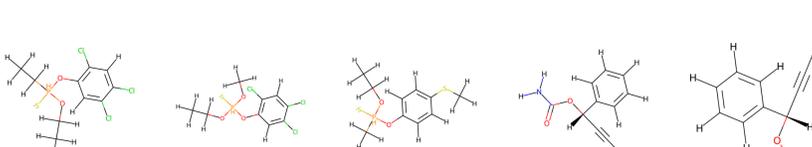
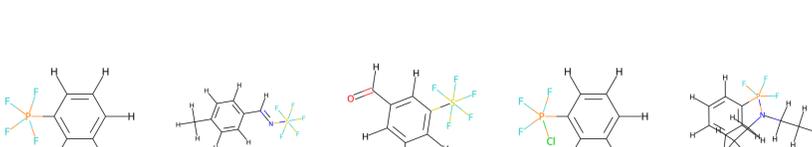
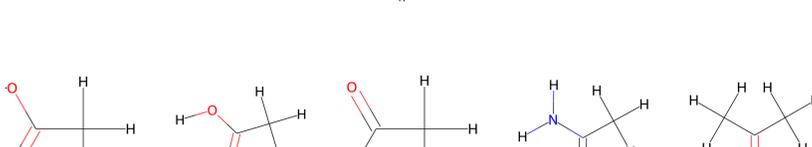
This pipeline has several special cases shown in the Table 7 along with the number of such cases in the dataset. The first case is the excessive number of isomorphisms: when we try to predict a molecule with hydrogens, these hydrogens can switch places between themselves without changing

Algorithm 5 Pipeline of atomic coordinates pre-processing

Require: $atom_coords$, $global_iso$, $local_iso$, x^{init} , $NeighborIdx$
for $k \in num_global_iso$ **do**
 for $l \in num_local_iso$ **do**
 $iso_neighbor_indices_{kl} \leftarrow global_iso[k] \cdot local_iso[l] \cdot NeighborIdx$
 $gt_neighbor_positions_{kl} \leftarrow atom_coords[iso_neighbor_indices_{kl}]$
 end for
end for
 $gt_neighb_transforms \leftarrow \mathbf{FitNeighbTransforms}(x^{init}, gt_neighbor_positions)$

the graph. Table 7 show that the number of isomorphism tend to increase exponentially when we have long hydrocarbon tails in some molecules. We describe our solution to this problem in section A.2.2.

Some molecules have several connected components in their graphs. In this case we select the largest one and continue processing. Molecules containing penta-valent silicon compounds (... Si_5 ...) and penta-valent phosphorous (... CPF_4 ...) are excluded from the dataset. We also exclude molecules containing only one neighborhood, single atom or atoms rare in biologically relevant compounds (Be, Kr, Ar etc).

Exception key	Examples	Number
Isomorphisms		4733
Disconnected		1547
Alignment		21676
Pentavalent		800
Single hood		3841

972						
973						
974		Br	Ca	Cl	HH	HH
975						K
976						
977	Single atom					43
978						
979						
980		Kr	O	Be	Ne	Ar
981						He
982						
983	Atom type					32

Table 7: Exclusion set.

A.2.1 NEIGHBORHOODS

As outlined in the data processing pipeline Alg 4, we first obtain atomic indices for each neighborhood. To make them consistent with chirality of the molecule we sort the indices for each neighborhood according to the CIP ranks of atoms in the molecule (Alg. 6). Then we compute the features of

Algorithm 6 GetNeighborhoods

994 **Require:** A_{ij} , $i, j \in [0, N)$ ▷ Adjacency matrix of the molecule
 995 **Require:** CIP_i ▷ CIP ranks of atoms in the molecule
 996 $RootIdx = \{i : i \text{ is not leaf node}\}$ ▷ Every non-leaf atom makes a neighborhood
 997 **for** $i \in RootIdx$ **do**
 998 $NeighbIdx_i \leftarrow \{j : A_{ij} = 1\}$
 999 **Sort**($NeighbIdx_i$, $key = CIP_i$) ▷ Sorting indexes using CIP ranks of the atoms
 1000 **end for**
 1001 **return** $NeighbIdx_i$

the neighborhoods following Alg 7. The atomic and bond features (F and B) are described in the manuscript. Pairwise features for neighborhoods are computed using algorithm 8, that is similar to the

Algorithm 7 GetHoodFeatures

1007 **Require:** $edge_j$, $j \in [0, M)$ ▷ Indexes of atoms connected by the M bonds
 1008 **Require:** F_i , $i \in [0, N)$ ▷ Atomic features
 1009 **Require:** B_j , $j \in [0, M)$ ▷ Bond features
 1010 **Require:** $RootIdx$, $NeighbIdx$ ▷ Root and neighbor indices for each neighborhood
 1011 **for** $i \in RootIdx$ **do**
 1012 **for** $j \in NeighbIdx_i$ **do**
 1013 $bond_index \leftarrow \{k : edge_k = (i, j)\}$
 1014 $HoodFeat_{ij} \leftarrow \text{cat}(F_i, F_j, B_{bond_index})$
 1015 **end for**
 1016 **end for**
 1017 **return** $HoodFeat_{ij}$

1018 previous algorithm.

1020 Additionally, we have to provide geometric input features to our model. We compute the initial
 1021 geometry of a neighborhood by combining bond geometry and bond lengths, extracted from the
 1022 data and described in the section A.1. Algorithm 9 outlines the procedure we use to assign initial
 1023 coordinates to the atoms of neighborhoods. We first assign signatures to the neighborhoods in the
 1024 same way we did it in section A.1, then we load bond order, allowed permutations and bond lengths
 1025 from the data geometric data. We then find permutation that matches the current bond order to the
 one in the data and assign the coordinates according to this permutation.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

Algorithm 8 GetPairFeatures

Require: $edge_j, j \in [0, M)$ ▷ Indexes of atoms connected by the M bonds
Require: $F_i, i \in [0, N)$ ▷ Atomic features
Require: $B_j, j \in [0, M)$ ▷ Bond features
Require: $RootIdx$ ▷ Root indices for each neighborhood
for $i \in RootIdx$ **do**
 for $j \in RootIdx$ **do**
 $bond_index \leftarrow \{k : edge_k = (i, j)\}$
 if $bond_index \neq \emptyset$ **then**
 $PairFeat_{ij} \leftarrow \text{cat}(F_i, F_j, B_{bond_index})$
 end if
 end for
end for
return $PairFeat_{ij}$

Algorithm 9 GetInitCoords

Require: $edge_j, j \in [0, M)$ ▷ Indexes of atoms connected by the M bonds
Require: $BondType_j, j \in [0, M)$ ▷ Bond types
Require: $Permutations_i, i \in [0, N)$ ▷ Allowed permutations for each neighborhood
Require: $BondOrder_i, i \in [0, N)$ ▷ Bond order for each neighborhood
Require: $Coords_{ij}, i \in [0, N), j \in [0, K_i)$ ▷ Coordinates of unit vectors of geometry of each neighborhood
Require: $BondLength_{ij}, i \in [0, N), j \in [0, K_i)$ ▷ Bond lengths of each bond in neighborhoods
Require: $RootIdx, NeighbIdx$ ▷ Root and neighbor indices for each neighborhood
for $i \in RootIdx$ **do**
 $bond_order \leftarrow \{\}$
 for $j \in NeighbIdx_i$ **do** ▷ Extracting bond order of the current neighborhood
 $bond_index \leftarrow \{k : edge_k = (i, j)\}$
 $bond_order \leftarrow bond_order \cup \{BondType_{bond_index}\}$
 end for
 $perm \leftarrow \{p : p(bond_order) = BondOrder_i, p \in Permutations_i\}$ ▷ Matching the bond orders
 $init_coords_i \leftarrow \{(0, 0, 0)\}$ ▷ Placing the root atom in the center
 for $k \in perm(NeighbIdx_i)$ **do** ▷ Computing the initial coordinates
 $init_coords_{ik} \leftarrow Coords_{ik} * BondLength_{ik}$
 end for
end for
return $init_coords_{ik}$

Table 8: Each one-hot feature is an encoding of a given property + 1 bit that indicates that the property is abnormal or incorrectly assigned.

Feature	Encryption	Size
Atom type	one-hot	35
Is aromatic	1/0	1
Atom degree	one-hot	8
Atom hybridization	one-hot	6
Atom implicit valence	one-hot	8
Atom formal charge	one-hot	4
Size of the ring atom belongs to	one-hot	6
Number of rings atom belongs to	one-hot	5
Atom chirality to	1/0	1
Bond type	one-hot	4
Bond in ring	1/0	1
Bond is conjugated	1/0	1
Bond is aromatic	1/0	1
Bond chirality	1/0	1

Finally we obtain pair transforms using algorithm 10. For each pair of neighborhoods i and j that are connected by a bond we take the coordinates of atom in the first neighborhood vec_src that constitutes the bond ij . Similarly we take the coordinates of atom from the second neighborhood j vec_dst that belongs to the bond ji between neighborhoods. Then we compute the transform that aligns second neighborhood along the bond vec_src .

Algorithm 10 GetPairTansforms

Require: $edge_j, j \in [0, M)$ ▷ Indexes of atoms connected by the M bonds
Require: $init_coords_{ik}, i \in [0, N), k \in [0, K_i)$ ▷ Initial coordinates of neighborhood atoms
Require: $RootIdx, NeighbIdx$ ▷ Root and neighbor indices for each neighborhood

```

for  $i \in RootIdx$  do
  for  $j \in NeighbIdx$  do
    if  $\{k : edge_k = (i, j)\} \neq \emptyset$  then
       $vec\_src \leftarrow \{init\_coords_{ik} : NeighbIdx_{ik} = j\}$  ▷ Initial coordinates of atom from
      neighborhood  $i$  to  $j$ 
       $vec\_dst \leftarrow \{init\_coords_{jk} : NeighbIdx_{jk} = i\}$  ▷ Initial coordinates of atom from
      neighborhood  $j$  to  $i$ 
       $H_{src} \leftarrow AlignX(vec\_src)$  ▷ Alignment matrix of vector to the X-axis
       $H_{dst} \leftarrow AlignX(vec\_dst)$ 
       $T_{ij}^{rot} \leftarrow H_{src}^T \cdot H_{flip} \cdot H_{dst}$  ▷ Rotational part of the pair transform
       $T_{ij}^{trans} \leftarrow vec\_src$  ▷ Translational part of the pair transform
    end if
  end for
end for
return  $T_{ij}$ 

```

The final set of features of each neighborhood is given in the Table 8.

A.2.2 ISOMORPHISMS FACTORIZATION

Lets denote a set of all isomorphisms of a molecule graph as I . Any one element in the set I is a permutation of atomic indices of a molecule. We need to enumerate this set in order to compute the final loss of the model prediction:

$$L = \min_{iso \in I} FAPE(T_{pred}, X_{pred}, T_{gt}(iso), X_{gt}(iso)) \quad (6)$$

where T_{pred}, T_{gt} are the predicted and ground truth transforms of the neighborhoods and X_{pred}, X_{gt} are the predicted and ground truth atomic coordinates of the molecule. The Eq.6 represents the

form we use in the current work, alternatively we can take the minimum over $T_{pred}(iso^{-1})$ and $X_{pred}(iso^{-1})$. However this alternative formulation makes algorithm more computationally heavy. The $X_{gt}(iso)$ can be written as $X_{gt}(iso) = \{x_{gt}^{iso(i)}, i \in [0, N]\}$. $T_{gt}(iso)$ are the transforms from initial coordinates of the neighborhood atoms to the ground truth coordinates permuted using an element of the set I .

Naively we can enumerate all isomorphisms of a molecule graph by first coloring the graph vertexes using the atom type and graph edges using the bond type. However, this procedure will yield the number of isomorphisms that exceed 10^5 for all molecules in the first row of Table 7. The reason for this is that each hydrogen bonded to a carbon using single bond can be swapped for any other hydrogen of the same atom, so that the number of such swap combinations grows exponentially with the number of carbons in the molecule.

In this work we deal with this problem for the case of terminal hydrogen atoms only. To circumvent combinatorial explosion of the isomorphism set, we factorize the set into local and global isomorphisms. Local isomorphisms $I_i^{(l)}$ are computed for each neighborhood i . Each local isomorphism should only permute hydrogens, leaving heavy atoms in their respective places. Global isomorphisms $I^{(g)}$ are the isomorphisms of the molecule without hydrogens, that are extended to the added hydrogens with identity permutation.

Proposition: $\forall i \in I \exists l \in I^{(l)}, m \in I^{(g)} : i = l \cdot m$ We leave out the proof of this proposition, however it should be trivial. If the proposition holds, we compute the loss of the model prediction in the following way:

$$iso_i^{(l)} = \operatorname{argmin}_{iso \in I_i^{(l)}} FAPE(T_{pred}, X_{pred}, T_{gt}(iso), X_{gt}(iso)), i \in [0, M] \quad (7)$$

$$L = \min_{iso \in I^{(g)}} FAPE \left(T_{pred}, X_{pred}, T_{gt} \left(iso \times \prod_{i \in [0, M]} iso_i^{(l)} \right), X_{gt} \left(iso \times \prod_{i \in [0, M]} iso_i \right) \right) \quad (8)$$

Despite this method, there are still some ligands that have highly symmetric molecular graphs. Although the number of global isomorphisms they have is on the order of 10^4 , we still exclude 91 molecules from the dataset that have more than 512 global isomorphisms for convenience. Figure 8 shows some examples of such molecules.

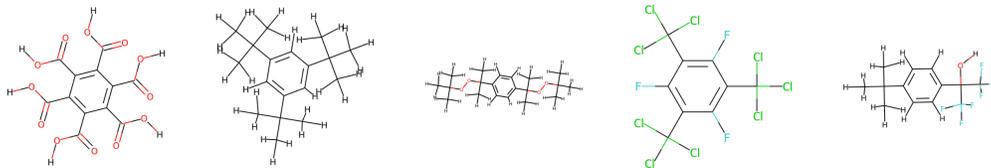


Figure 8: Examples of molecules with the graph that has more than 512 global isomorphisms.

A.2.3 GROUND TRUTH

While obtaining neighborhood transforms from neighbor positions we run into two special cases: the chirality of a neighborhood is not set and the corresponding bond geometry fits the coordinates badly ($RMSD - \min RMSD > 0.3$). The algorithm 11 outlines the way we detect and treat these special cases. If the *chirality_mask* is true, then we assign chirality to this atom and reprocess the data. In practice, we observed that quatro-valent phosphorus in certain ligands is not labeled as chiral. If the *fit_mask* contains true values, we leave transforms that we obtained in this algorithm and unmask all its local isomorphisms.

A.2.4 DATASET STATISTICS

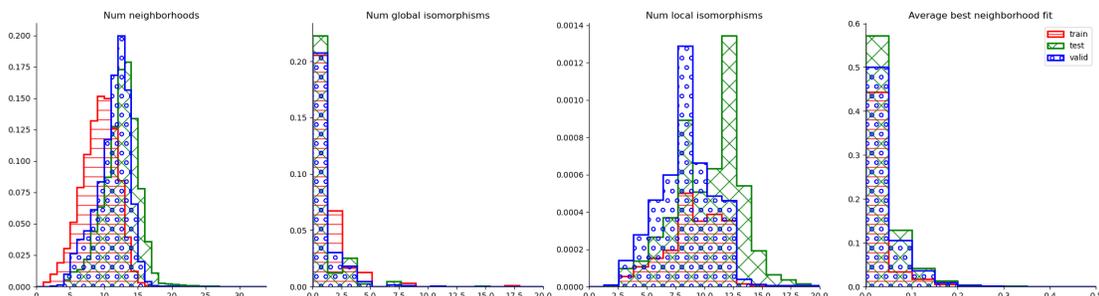
Overall out of 3,982,254 molecules in the dataset we exclude 5105. Out of which are 801 compounds containing penta-valent atoms, 4137 single-neighborhood compounds, 43 examples that have atoms

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

Algorithm 11 FitNeighbTransforms

Require: $init_coords_{ik}$, $i \in [0, N)$, $k \in [0, K_i)$ \triangleright Initial coordinates of atoms in neighborhoods
Require: $gt_positions_{glik}$, $g \in global_iso$, $l \in local_iso$, $i \in [0, N)$, $k \in [0, K_i)$ \triangleright Ground truth coordinates
 $rmsd_{gli} \leftarrow RMSD(init_coords, gt_positions)$ \triangleright Aligning all initial coordinates to all ground truth ones
 $min_rmsd_i \leftarrow \min_{g^i}(rmsd_{gli})$ \triangleright Min RMSD for each neighborhood
 $iso_mask_{gli} \leftarrow (rmsd_{gli} - min_rmsd_i) < 0.3$ \triangleright Flagging true all variants that are close to the best fit
 $loc_iso_mask_{gi} \leftarrow \exists_l iso_mask_{gli}$
if $\forall_g \neg(\forall_i loc_iso_mask_{gi})$ **then**
 \triangleright If all global isomorphisms have at least one hood that does not fit any local isomorphisms
 $W_{gi} \leftarrow \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2loc_iso_mask_{gi} - 1 \end{pmatrix}$
 \triangleright We flip chirality of neighborhoods that did not fit any local isomorphism
 $rmsd_{gli}^{chir} \leftarrow RMSD(W \cdot init_coords, gt_positions)$
 $loc_iso_mask_{gi}^{chir} \leftarrow \exists_l rmsd_{gli}^{chir} < 0.82$
 \triangleright If fit changes when we flip the chirality then it is assigned incorrectly
 $chirality_mask \leftarrow (\exists_g loc_iso_mask_{gi}) \oplus (\exists_g loc_iso_mask_{gi}^{chir})$
if $\exists_i chirality_mask_i$ **then**
 Flip chiral center i and redo data processing
end if
 \triangleright If the fit does not change, then we have non-standard neighborhood geometry
 $fit_mask \leftarrow (\neg \exists_g loc_iso_mask_{gi}) \wedge (\neg \exists_g loc_iso_mask_{gi}^{chir})$
if $\exists_i fit_mask_i$ **then**
 $iso_mask[fit_mask] \leftarrow True$
end if
end if
 \triangleright If there is no local isomorphism for at least one neighborhood, we mask such global isomorphism
 $glob_iso_mask_g \leftarrow \forall_i \exists_l iso_mask_{gli}$
 $iso_mask[\neg glob_iso_mask_g] \leftarrow False$
 $loc_iso_mask_{gi} \leftarrow \exists_l iso_mask_{gli}$
if $\forall_g \neg(\forall_i loc_iso_mask_{gi})$ **then**
 \triangleright If after all the changes we still have problematic neighborhoods we retain those close to the optimal one
 $loc_iso_min_idx_{gi} = \underset{l}{\operatorname{argmin}}(rmsd_{gli})$
 $opt_rmsd_g \leftarrow \sum_i \min_l(rmsd_{gli})$ \triangleright Min RMSD for each isomorphism
 $glob_iso_mask_g \leftarrow opt_rmsd_g - \min_g(opt_rmsd_g) < 0.3$
 $iso_mask_{gli} \leftarrow False$
 $iso_mask_{glob_iso_min_idx_{gi}} \leftarrow True$
 $iso_mask_{g-\neg glob_iso_mask_g} \leftarrow False$
end if
 \triangleright Finally we obtain rotations and translations from alignment matrix U
return $rot \leftarrow rmsd_{gli}.U^T$
return $trans \leftarrow gt_positions_{gl0k}$
return iso_mask

1242 without neighbors, 32 compounds with rare atom types and 1 compound that is disconnected and has
 1243 non-trivial chirality. Additionally we remove 91 examples because of excessive number of global
 1244 isomorphisms.
 1245



1257 Figure 9: Distribution of number of neighborhoods, global and local isomorphisms and average
 1258 neighborhood best fit rmsd in training, validation and test sets.
 1259

1260 Figure 9 shows the distribution of the number of neighborhoods in train, test and validation subsets,
 1261 as well as distribution of number of global and local isomorphisms and average neighborhood rmsd
 1262 fit. We see that test set is slightly more challenging than the training and validation subsets. We
 1263 also can see in Fig. 9 that the overall best fit of the dataset using our molecule representation is good
 1264 enough to assume the validity of such a representation.
 1265

1266 Additionally, Figure 10 shows molecules that are outliers in the statistics shown on the Fig 9. The
 1267 first row shows largest molecules in the dataset with the labels corresponding to the numbers of
 1268 neighborhoods in each molecule. The second row corresponds to the molecules with the biggest
 1269 numbers of global isomorphisms. The third row shows molecules with the biggest number of
 1270 neighborhoods with non-trivial local isomorphisms (that are not filtered out based on the structure).
 1271 Finally, the fourth row shows the molecules that have the worst fit using rigid-body bond geometries
 1272 along with the average rmsd over all neighborhoods.

1273 A.3 MODEL

1274 Algorithm 12 gives an overview of the model. First, we generate initial transforms of the neighbor-
 1275 hoods and embed input features. Then we iterate over Evoformer blocks, and iteratively refine initial
 1276 translations and rotations. Iteration parameters (Linear layers and Evoformer parameters) are unique.
 1277 During the iteration we first scalarize the geometric features, pass them as a transformer input to
 1278 the Evoformer and then vectorize the transformed Evoformer output, which we use to update the
 1279 geometry of a molecule.
 1280

1281 In this work we follow the approach by the AlphaFold2 team and use the same initialization of the rigid
 1282 body transforms. Algorithm 13 shows that we place all the neighborhoods in the frame origin and
 1283 assign them the same rotation.

1284 For feature embedding we take the output features of the data processing algorithm *HoodFeat* and
 1285 *PairFeat*. First we append initial atomic coordinates of the neighbors in neighborhoods and linearly
 1286 transform them, obtaining neighbor embeddings *neighb_embed* as a result. Next we concatenate
 1287 neighbor embedding for each neighborhood and transform these features obtaining embeddings
 1288 for neighborhoods *hood_embed*. The pairwise neighborhood features are passed through a linear
 1289 transform followed by ReLU and another linear transform to obtain pairwise embeddings *pair_embed*.
 1290 Algorithm 14 shows the pseudocode of the whole process.

1291 To make Evoformer application on the geometric features equivariant we first centralize the coordinates
 1292 of neighborhoods (Algorithm 15). This gives us translation invariance. The rotation equivariance is
 1293 archived by transforming geometric features into scalars in certain basis and after applying Evoformer,
 1294 transforming the output back into geometric features using the same basis. To do this we first construct
 1295 a set of basis for a molecule, one for each neighborhood. We use rotation of the neighborhood as a
 vector basis (Algorithm 16, *node_basis*). To scalarize rotation matrixes we also construct the basis

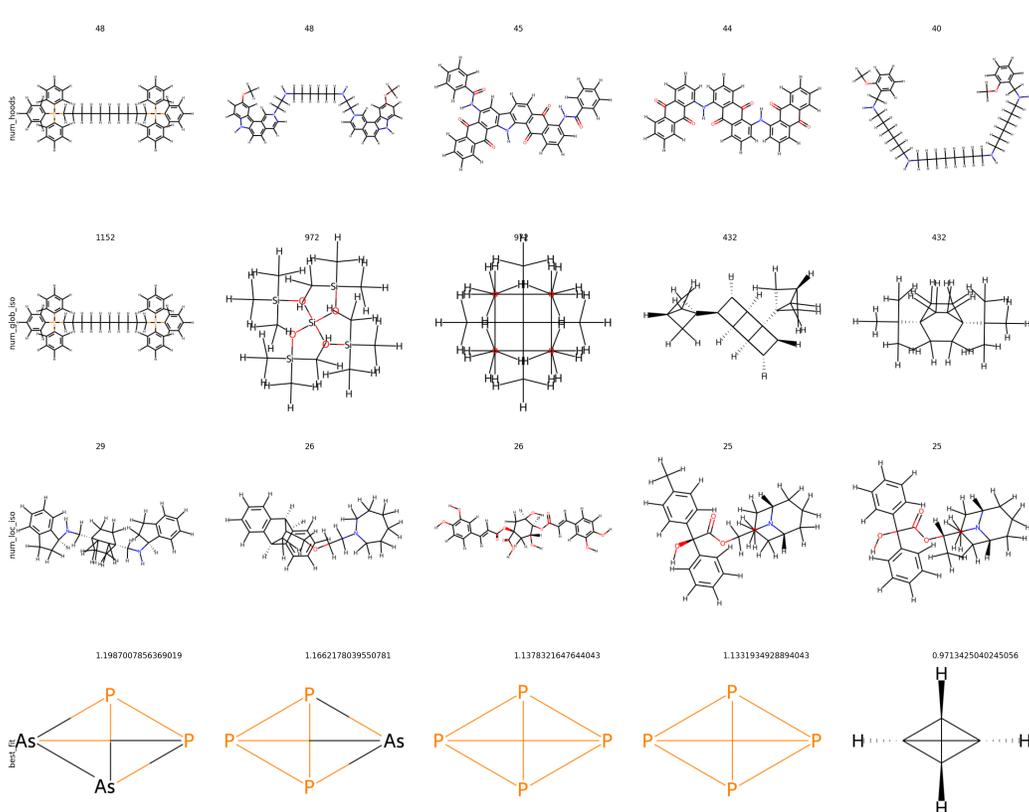


Figure 10: Outliers from the datasets. Labels on the left denote the category, labels above denote the parameter, that this category has. The values of parameters correspond to the data on Fig 9.

in the space of (1,1)-tensors. We do it by taking outer product of vector basis of each node with all the other nodes. This way we obtain N_{hoods}^2 basis for pairwise tensor features (Algorithm 16, *edge_basis*).

Algorithm 15 Centralize

Require: $trans_{bi} \in R^3$, $i \in [0, N_{hoods})$
 $center_b \in R^3 \leftarrow \frac{1}{N_{hoods}} \sum_i trans_{bi}$
return $trans_{bi} - center_b$

Algorithm 16 GetBasis

Require: $rot_{bi} \in R^9$, $i \in [0, N_{hoods})$
 $node_basis_{bik} \in R^3 \leftarrow rot_{bi}^T$, $k \in [0, 3)$
 $edge_basis_{bijkl} \in R^9 \leftarrow node_basis_{bik} \otimes node_basis_{bjl}$, $k, l \in [0, 3)$
return $node_basis_{bik}$, $edge_basis_{bijkl}$

To obtain scalar features from translations, rotations of each neighborhood and the pairwise translations and rotations features we first project translations onto vectors of the *node_basis* (Algorithm 17) obtaining *node_scalars*. For the pairwise translations we first compute relative translations *rel_trans* and then transform them into translation tensor by taking outer product of their components (Algorithm 17, *edge_trans*). Then we project the relative translations tensor onto the pairwise tensor basis. Similarly we obtain scalar pairwise rotations. Finally we concatenate pairwise scalars into *edge_scalars*.

Algorithm 17 Scalarize

Require: $trans_{bi} \in R^3, rot_{bi} \in R^9 \quad i \in [0, N_{hoods})$
Require: $pair_trans_{bij} \in R^3, pair_rot_{bij} \in R^9 \quad i, j \in [0, N_{hoods})$
Require: $node_basis_{bim} \in R^3, edge_basis_{bijkl} \in R^9 \quad i, j \in [0, N_{hoods}), \quad m \in [0, 3), \quad k, l \in [0, 9)$
 $node_scalars_{bi} \in R^3 \leftarrow (trans_{bi}, node_basis_{bim})$
 $rel_trans_{bij} \in R^3 \leftarrow (trans_{bi} - trans_{bj}) + pair_trans_{bij}$
 $edge_trans_{bij} \in R^9 \leftarrow rel_trans_{bij} \otimes rel_trans_{bij}$
 $scalar_edge_trans_{bij} \in R^9 \leftarrow (rel_trans_{bij}, edge_basis_{bijkl})$
 $scalar_edge_rot_{bij} \in R^9 \leftarrow (pair_rot_{bij}, edge_basis_{bijkl})$
 $edge_scalars_{bij} \in R^{18} \leftarrow edge_trans_{bij} \oplus scalar_edge_rot_{bij}$
return $node_scalars_{bi} \in R^3, edge_scalars_{bij} \in R^{18}$

Algorithm 18 describes our vectorization process of the output features $transform_{bim}$. Here we obtain three vectors for each neighborhood by treating $transform_{bim}$ as the decomposition coefficients into the $node_basis_{bi}$. One of these three vectors is then interpreted as a translation update of a neighborhood $translation_{bi}$. The other two vectors are used to construct rotation using Gram-Schmidt process.

Algorithm 18 Vectorize

Require: $transform_{bim} \in R^3, \quad i \in [0, N_{hoods}), \quad m \in [0, 3)$
Require: $node_basis_{bim} \in R^3, edge_basis_{bijkl} \in R^9 \quad i, j \in [0, N_{hoods}), \quad m \in [0, 3), \quad k, l \in [0, 9)$
 $vectors_{bik} \in R^3 \leftarrow \sum_m (node_basis_{bik}, transform_{bim})$
 $translation_{bi} \in R^3 \leftarrow vectors_{bi0}$
 $a_{bi}^{(1)} \in R^3 \leftarrow vectors_{bi1}$
 $a_{bi}^{(2)} \in R^3 \leftarrow vectors_{bi2}$

▷ Gram-schmidt process

$b_{bi}^{(1)} \leftarrow \frac{a_{bi}^{(1)}}{|a_{bi}^{(1)}|}$
 $b_{bi}^{(2)} \leftarrow a_{bi}^{(2)} - (b_{bi}^{(1)}, a_{bi}^{(2)})b_{bi}^{(1)}$
 $b_{bi}^{(2)} \leftarrow \frac{b_{bi}^{(2)}}{|b_{bi}^{(2)}|}$
 $b_{bi}^{(3)} \leftarrow b_{bi}^{(1)} \times b_{bi}^{(2)}$
 $rotation_{bi} \in R^9 \leftarrow b_{bi}^{(1)} \oplus b_{bi}^{(2)} \oplus b_{bi}^{(3)}$
return $translation_{bi} \in R^3, rotation_{bi} \in R^9$

A.3.1 EVOFORMER

In this work we use the same Evoformer block, as the one in AlphaFold2 with few modifications. First, we do not need column-wise attention, because the input single features have only one row. The second important change is that throughout the Evoformer block we use GELU activation function instead of ReLU. The major change we made is the addition of spectral normalization in the attention layer. Algorithm 19 shows the changes to the gated self-attention with pair bias in bold. Similar changes are done to the triangle attention modules.

A.4 LOSSES

First we compute the iterative atomic structures of a molecule based on rotations and translations (all_rot, all_trans) output of the model. Algorithm 20 shows the outline of molecule reconstruction using the scatter operation. Effectively we predict atom positions belonging to the bond between neighborhoods twice and then average over these predictions. We omit the technical details of tensor manipulation. Similar procedure is performed for the single representations of the neighborhoods. In

Algorithm 19 Gated self-attention with pair bias

1458 **Require:** $m_{bi} \in R^{N_{feat}}$, $i \in [0, N_{hoods})$
1459 **Require:** $z_{bij} \in R^{N_{feat}}$, $i \in [0, N_{hoods})$
1460 \triangleright Iteratively compute maximum eigenvalue of the matrix $K^T Q$
1461 $\mathbf{W} \leftarrow \mathbf{K}^T \mathbf{Q}$
1462 $\mathbf{u} \leftarrow \mathbf{W} \cdot \mathbf{u}$ $\triangleright u$ is the parameter of this module, saved for the next step
1463 $\mathbf{u} \leftarrow \frac{\mathbf{u}}{|\mathbf{u}|}$
1464 $\mathbf{v} \leftarrow \mathbf{W} \cdot \mathbf{v}$ $\triangleright v$ is the parameter of this module, saved for the next step
1465 $\mathbf{v} \leftarrow \frac{\mathbf{v}}{|\mathbf{v}|}$
1466 $\sigma^h \leftarrow \sum_{dc} \mathbf{u}_d^h \mathbf{W}_{dc}^h \mathbf{v}_c^h$ \triangleright Standard gated self-attention with pair bias
1467 $m_{bi} \leftarrow LayerNorm(m_{bi})$
1468 $q_{bi}^h, k_{bi}^h, v_{bi}^h \leftarrow LinearNoBias^{QKV}(m_{bi})$ $\triangleright Q, K, V$ are matrixes of the linear transform
1469 $b_{bij}^h \leftarrow LinearNoBias(LayerNorm(z_{bij}))$
1470 $g_{bi}^h \leftarrow Sigmoid(Linear(m_{bi}))$
1471 $a_{bij}^h \leftarrow Softmax_j(\frac{1}{\sigma^h \sqrt{c}} q_{bi}^h k_{bj}^h + b_{bij}^h)$ \triangleright Additional factor σ^h
1472 $o_{bi}^h \leftarrow g_{bi}^h \cdot \sum_j a_{bij}^h v_{bi}^h$
1473 $\tilde{m}_{bi} \leftarrow Linear(concat_h(o_{bi}^h))$
1474 **return** \tilde{m}

1475 the end we have $atom_positions_{lbk}$ tensor, where l indexes evoformer blocks outputs, b corresponds
1476 to the molecule index in a batch and k enumerates the atoms in a molecule. Additionally we obtain
1477 $single_act_{bk}$ representation for each atom in the batch of molecules.

Algorithm 20 Structure reconstruction

1484 **Require:** $all_rot_{lbi} \in R^9, all_trans_{lbi} \in R^3$, $l \in [0, num_evoformer_blocks)$, $b \in$
1485 $[0, batch_size)$, $i \in [0, N_{hoods})$
1486 **Require:** $init_coords_{bik} \in R^3$, $i \in [0, N_{hoods})$, $k \in [0, N_{neighb})$
1487 **Require:** $atom_mask_{bik}$, $i \in [0, N_{hoods})$, $k \in [0, N_{neighb})$ \triangleright 1/0 depending on whether atom k
1488 is present in neighborhood i
1489 **Require:** $atom_indices_{bik}$, $i \in [0, N_{hoods})$, $k \in [0, N_{neighb})$ \triangleright global index of atom k in
1490 neighborhood i
1491 $neighbor_positions_{lbik} \leftarrow all_rot_{lbi} \cdot init_coords_{bik} + all_trans_{lbi}$
1492 $num_atoms_m \leftarrow \sum_{bik} atom_mask_{bik} \delta(atom_indices_{bik} - m)$ \triangleright Scatter operation: we sum
1493 over $atom_mask$ to the cell with indices of $atom_index$
1494 $atom_positions_{lm} \leftarrow \frac{1}{num_atoms_m} \sum_{bik} neighbor_positions_{lbik} \delta(atom_indices_{bik} - m)$
1495 $atom_positions_{lbk} \in R^3 \leftarrow atom_positions_{lm}$, $k \in [0, num_atoms_b)$ \triangleright Rearranging tensor
1496 to the batch of molecules
1497 $single_act_m \leftarrow \frac{1}{num_atoms_m} \sum_{bik} single_act_{bik} \delta(atom_indices_{bik} - l)$
1498 $single_act_{bk} \in R^{N_{feat}} \leftarrow single_act_m$, $k \in [0, num_atoms_b)$
1499 **return** $atom_positions_{lbk}, single_act_{bk}, neighbor_positions_{lbik}$

A.4.1 STRUCTURAL LOSSES

1504 The first loss that we compute is Frame-Aligned Point Error (FAPE, Algorithm 21) of the reconstructed
1505 structure with respect to the ground truth structure. However, because we have our ground truth data
1506 in the factorized form we first have to find the best matching global and local isomorphism. Algorithm
1507 22 outlines our implementation. The key feature of this algorithm is that it is an approximation
1508 of the Eq.7. We do not reconstruct structures for each isomorphism, instead we compute FAPE
1509 for neighborhoods for each isomorphism and then approximate FAPE of the whole structure by
1510 the sum FAPE over all neighborhoods. In practice this approximation gives us the same minimum
1511 as the whole-structure FAPE. Afterwards we reconstruct ground truth structure for the selected
isomorphisms and compute the correct FAPE score for the whole reconstructed structure.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

Algorithm 21 Frame-Aligned Point Error(FAPE)

Require: $pred_T_i \in R^{12}, pred_pos_j \in R^3, gt_T_i \in R^{12}, gt_pos_j \in R^3$

$x_{ij} \leftarrow pred_T_i^{-1} \circ pred_pos_j$

$gt_x_{ij} \leftarrow gt_T_i^{-1} \circ gt_pos_j$

$d_{ij} = \sqrt{\|x_{ij} - gt_x_{ij}\|^2}$

return $\frac{1}{10N_i N_j} \sum_{ij} (\min(10, d_{ij}))$

Algorithm 22 Structure loss

Require: $gt_neighbor_positions_{bglin} \in R^3, gt_rot_{bgli} \in R^9, gt_trans_{glt} \in R^3, b \in [0, batch_size), g \in [0, N_{glob.iso}), l \in [0, N_{loc.iso}), i \in [0, N_{hoods})$

Require: $all_rot_{mbi} \in R^9, all_trans_{mbi} \in R^3, atom_positions_{mbk} \in R^3, m \in [0, num_evoformer_blocks)$

Require: $neighbor_positions_{mbin}$

Require: $atom_mask_{bin}, i \in [0, N_{hoods}), n \in [0, N_{neighbor}) \triangleright 1/0$ depending on whether atom n is present in neighborhood i

Require: $atom_indices_{bin}, i \in [0, N_{hoods}), n \in [0, N_{neighbor}) \triangleright$ global index of atom n in neighborhood i

\triangleright Computing FAPE loss for each neighborhood and isomorphism between the prediction and the ground truth

$pred_rigids_{mbi} \in R^{12} \leftarrow all_rot_{mbi} \oplus all_trans_{mbi}$

$gt_rigids_{bgli} \in R^{12} \leftarrow gt_rot_{bgli} \oplus gt_trans_{bgt}$

$neighb_fape_{mbgli} \leftarrow FAPE(pred_rigids_{mbi}, gt_rigids_{bgt}, neighbor_positions_{mbin}, gt_neighbor_positions_{bgtin})$

\triangleright Getting indices of local and global isomorphisms

$local_iso_idx_{mbgi}, min_local_iso_{mbgi} \leftarrow \underset{l}{\operatorname{argmin}}(neighb_fape_{mbgli}), \underset{l}{\min}(neighb_fape_{mbgli})$

$global_iso_idx_{mb} \leftarrow \underset{g}{\operatorname{argmin}}(\sum_i min_local_iso_{mbgi})$

$local_iso_idx_{mbi} \leftarrow local_iso_idx_{m,b,global_iso_idx_{mb},i}$

\triangleright Selecting rigid transforms and neighbor positions of the ground truth based on local and global isomorphism

$gt_rigids_{mbi} \leftarrow gt_rigids_{b,global_iso_idx_{mb},local_iso_idx_{mbi},i}$

$gt_neighbor_positions_{mbin} \leftarrow gt_neighbor_positions_{b,global_iso_idx_{mb},local_iso_idx_{mbi},i,n}$

\triangleright Reconstructing atomic positions for each molecule from neighbor positions, same operation as in Structure reconstruction algorithm

$num_atoms_f \leftarrow \sum_{bin} atom_mask_{bin} \delta(atom_indices_{bin} - f)$

$gt_atom_positions_{mf} \leftarrow \frac{1}{num_atoms_f} \sum_{bin} gt_neighbor_positions_{mbin} \delta(atom_indices_{bin} - f)$

$gt_atom_positions_{mbk} \in R^3 \leftarrow gt_atom_positions_{mf}, k \in [0, num_atoms_b) \triangleright$ Rearranging tensor to the batch of molecules

\triangleright Computing FAPE loss for the reconstructed ground truth for each evofomer block and molecule in the batch

$struct_fape_{mb} \leftarrow FAPE(pred_rigids_{mbi}, gt_rigids_{mbi}, atom_positions_{mbk}, gt_atom_positions_{mbk})$

return $\frac{\sum_{mb} struct_fape_{mb}}{num_evoformer_blocks \cdot batch_size}, gt_atom_positions_{mbk}$

1566 Additionally we penalize the clashes between atoms in the structure. Algorithm 23 shows our
 1567 procedure for computing clash loss. Importantly, we exclude first and second neighbors from the
 1568 loss, because it is guaranteed that some of these neighbors belong to the same neighborhoods and
 1569 our algorithm treats them as rigid bodies. To compute second-order neighbors we use well known
 1570 formula $A^{second} = A^{first}(A^{first})^T > 0$, where A^{first} is the adjacency matrix of the molecular
 1571 graph.

1572 Algorithm 23 Clash loss

1574 **Require:** $atom_positions_{bk} \in R^3$, $b \in [0, batch_size)$, $k \in [0, N_{atoms})$
 1575 **Require:** adj_{bkl} , $l \in [0, N_{atoms})$ \triangleright Adjacency matrix
 1576 **Require:** r_{bk} , $k \in [0, N_{atoms})$ \triangleright Atomic radius
 1577 $second_adj_{bkl} \leftarrow (\sum_m adj_{bkm} adj_{blm}) > 0$ \triangleright Adjacency for second neighbors
 1578 $min_dist_{bkl} \leftarrow (r_{bk} + r_{bl})(1 - second_adj_{bkl})$
 1579 $d_{bkl} \leftarrow \sqrt{\|atom_positions_{bk} - atom_positions_{bl}\|^2}$
 1580 $L = \frac{1}{\sum_{bkl}(min_dist_{bkl} > 0)} \sum_{bkl} ReLU(min_dist_{bkl} - d_{bkl})$
 1581 **return** L

1582
 1583 Finally, collinearity of the predictions is also penalized as described in the Algorithm 24. Specifically
 1584 we save collinearity values during the vectorization stage described by the Algorithm 18. After
 1585 computing all the Evoformer iterations, we average over the batch and iterations and obtain the loss.
 1586

1587 Algorithm 24 Collinearity loss

1588 **Require:** $b_{mbi}^{(1)} \in R^3$, $a_{mbi}^{(2)} \in R^3$ $m \in [0, num_evoformer_blocks)$, $b \in [0, batch_size)$, $k \in$
 1589 $[0, N_{hoods})$
 1590 $coll_{mbi} \leftarrow (b_{mbi}^{(1)}, \frac{a_{mbi}^{(2)}}{\|a_{mbi}^{(2)}\|_1})$
 1591 $L \leftarrow \frac{1}{num_evoformer_blocks \cdot batch_size \cdot N_{hoods}} \sum_{mbi} coll_{mbi}$
 1592 **return** L

1593
 1594
 1595 Similar to AlphaFold2 we predict the model confidence over its own structure prediction. In our
 1596 case we predict per-neighborhood IDDT scores. As the Algorithm 25 shows, our implementation has
 1597 almost no changes from the one used in AlphaFold2.
 1598

1599 Algorithm 25 pLDDT loss

1600 **Require:** $atom_positions_{bk}, gt_atom_positions_{bk}, single_act_{bi}$, $b \in [0, batch_size)$, $i \in$
 1601 $[0, N_{hoods})$, $k \in [0, N_{atoms})$
 1602 **Require:** $neighbor_atom_indices_{bin}$, $n \in [0, N_{neighbor})$
 1603 **Require:** $v^{(bins)} \in R^{N_{bins}}$ \triangleright Vector of bin cutoff values, f.e [1, 3, 5, ... 99]
 1604 \triangleright Computing ground truth LDDT, based on the predicted and gt atomic structures, then averaging
 1605 atomic IDDT over neighborhoods
 1606 $d_{bkl} \leftarrow \sqrt{\|atom_positions_{bk} - atom_positions_{bl}\|^2}$
 1607 $gt_d_{bkl} \leftarrow \sqrt{\|gt_atom_positions_{bk} - gt_atom_positions_{bl}\|^2}$
 1608 $L1_{bkl} \leftarrow |d_{bkl} - gt_d_{bkl}|$
 1609 $score_{bkl} \leftarrow \frac{1}{4} ((L1_{bkl} < 0.5) + (L1_{bkl} < 1.0) + (L1_{bkl} < 2.0) + (L1_{bkl} < 4.0))$
 1610 $gt_LDDT_{bk} \leftarrow \frac{1}{\sum_l (gt_d_{bkl} < 15)} \sum_l score_{bkl} (gt_d_{bkl} < 15)$
 1611 $gt_LDDT_{bi} \leftarrow \frac{1}{N_{neighbor}} \sum_n gt_LDDT_{b, neighbor_atom_indices_{bin}}$
 1612 \triangleright Computing predicted LDDT and using cross-entropy loss to compare it to ground-truth LDDT
 1613 $y_{bi} \in R^{N_{bins}} \leftarrow relu(Linear(relu(Linear(LayerNorm(single_act_{bi}))))))$
 1614 $p_{bi} \leftarrow SoftMax(Linear(y_{bi}))$
 1615 $gt_p_{bi} \leftarrow OneHot(gt_LDDT_{bi}, v^{(bins)})$
 1616 $pLDDT_{bi} \leftarrow p_{bi}^T v^{(bins)}$
 1617 $L \leftarrow \frac{1}{batch_size \cdot N_{hoods}} \sum_{ib} (gt_p_{bi}^T \log p_{bi})$
 1618 **return** $pLDDT_{bi}, L$

A.4.2 PROPERTY PREDICTION

We added property prediction head to check whether the atom-wise single activations can be used in the end-to-end fashion for predictions. We use standard SchNet architecture with the assumption of a fully-connected atomic graph. Algorithms 26, 27 and 28 summarize the architecture of the neural network used for HOMO-LUMO gap prediction as well as the loss computation.

Algorithm 26 Property prediction head

Require: $atom_positions_{bk} \in R^3$, $single_act_{bk}$, $b \in [0, batch_size)$, $k \in [0, N_{atoms})$
Require: gt_b ▷ Ground-truth homo-lumo gap

$h_{bk} \leftarrow LayerNorm(Linear(single_act_{bk}))$
 $d_{bkl} \leftarrow \sqrt{\|atom_positions_{bk} - atom_positions_{bl}\|^2}$
 $rbf_{bkl} \leftarrow RBF(d_{bkl})$
for $i \in [0, N_{interact})$ **do**
 $h_{bk} \leftarrow h_{bk} + Interaction_i(h_{bk}, rbf_{bkl})$
end for
 $pred_b \leftarrow \sum_k Linear(ShiftedSoftPlus(Linear(h_{bk})))$
 $L \leftarrow \frac{1}{batch_size} \sum_b |pred_b - gt_b|$
return L

Algorithm 27 Radial basis function

Require: d_{bkl} , $b \in [0, batch_size)$, $k, l \in [0, N_{atoms})$
Require: $x_0 = 0.0$, $x_1 = 5.0$, $N_{gaussians} = 50$

$r_m \leftarrow x_0 + \frac{m}{N_{gaussians}}(x_1 - x_0)$, $m \in [0, N_{gaussians})$
 $dec_{bklm} \leftarrow e^{-\frac{(d_{bkl} - r_m)^2}{2(r_1 - r_0)^2}}$
return $dec_{bkl} \in R^{N_{gaussians}}$

Algorithm 28 Interaction block

Require: $h_{bk} \in R^{N_{feat}}$, $rbf_{bkl} \in R^{N_{gaussians}}$, $b \in [0, batch_size)$, $k, l \in [0, N_{atoms})$

$h \leftarrow Linear(h)$ ▷ CF convolution module

$W_{bkl} \in R^{N_{feat}} \leftarrow ShiftedSoftPlus(Linear(ShiftedSoftPlus(Linear(rbf_{bkl}))))$
 $h_{bk} \leftarrow \sum_l (h_{bk} \cdot W_{bkl})$ ▷ Output of interaction block

$h_{bk} \leftarrow Linear(Softplus(Linear(h_{bk})))$
return $h_{bk} \in R^{N_{feat}}$
