

---

# 3DLinker: An E(3) Equivariant Variational Autoencoder for Molecular Linker Design

---

Yinan Huang<sup>1</sup> Xingang Peng<sup>2</sup> Jianzhu Ma<sup>3,1</sup> Muhan Zhang<sup>3,1</sup>

## Abstract

Deep learning has achieved tremendous success in designing novel chemical compounds with desirable pharmaceutical properties. In this work, we focus on a new type of drug design problem—generating a small “linker” to physically attach two independent molecules with their distinct functions. The main computational challenges include: 1) the generation of linkers is *conditional* on the two given molecules, in contrast to generating complete molecules from scratch in previous works; 2) linkers heavily depend on the anchor atoms of the two molecules to be connected, which are not known beforehand; 3) 3D structures and orientations of the molecules need to be considered to avoid atom clashes, for which equivariance to E(3) group are necessary. To address these problems, we propose a conditional generative model, named 3DLinker, which is able to predict anchor atoms and jointly generate linker graphs and their 3D structures based on an E(3) equivariant graph variational autoencoder. So far as we know, no previous models could achieve this task. We compare our model with multiple conditional generative models modified from other molecular design tasks and find that our model has a significantly higher rate in recovering molecular graphs, and more importantly, accurately predicting the 3D coordinates of all the atoms.

## 1. Introduction

The biological functions of most small molecule drugs are to inhibit the activity of the target protein by binding its active sites. In drug discovery, designing new molecule

drugs with desired pharmacophoric properties remains challenging due to the discreteness and enormity of the search space (Polishchuk et al., 2013). To address this problem, many machine learning methods have been developed to embed molecules in a compact hidden space, making promising progress in multiple downstream computational tasks such as molecular de-novo design, molecular optimization, and chemical property prediction.

Molecules are generally represented by graphs with atoms and bonds represented as nodes and edges, respectively. Graph generative models (Liu et al., 2018; Shi et al., 2019; Jin et al., 2018; 2020) are commonly applied to model the marginal probability for FDA-approved drug molecules and it is expected that the newly sampled molecules from the model have similar or better pharmacophoric properties.

However, in complex diseases such as cancer, mutations of amino acids could significantly impact the binding affinity between drugs and target proteins. The drug might fall off the drug target when a particular amino acid mutates with a certain probability due to the weak binding affinity, making the patient drug-resistant. To solve this problem, more recently, an alternative drug mechanism named Proteolysis targeting chimera (PROTAC) is developed to inhibit the protein functions by prompting complete degradation of the target protein. PROTAC is a unique molecule composed of two *fragment* molecules and a *linker* molecule: one fragment binds the target protein, the other binds another molecule that can degrade the target protein, and the linker attaches the two fragments together. Because PROTAC needs only to bind their targets with high selectivity (rather than inhibit the target protein’s activity), many efforts are devoted to retooling previously ineffective inhibitor molecules as PROTAC for developing the next-generation drugs. Even though PROTAC owns promising potential, it has not been broadly pushed into clinical trial stages. One of the key challenges is the design of linker, which has a critical influence on the ultimate degradation of the target protein. To date, linker design still relies on the expertise of structural biologists and thus is very time-intensive. Therefore, there are increasing efforts to develop deep learning methods to address linker design problems (Imrie et al., 2020; Yang et al., 2020).

---

<sup>1</sup>Beijing Institute for General Artificial Intelligence <sup>2</sup>Tsinghua University <sup>3</sup>Institute for Artificial Intelligence, Peking University. Correspondence to: Muhan Zhang <muhan@pku.edu.cn>, Jianzhu Ma <majianzhu@pku.edu.cn>.

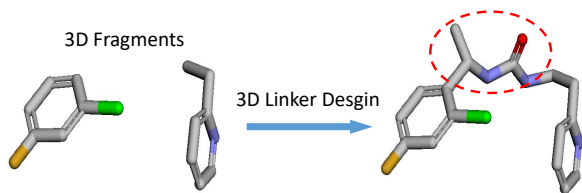


Figure 1: 3D linker design problem: given two fragments’ graph with 3D coordinates (left), the goal is to generate a linker graph with 3D coordinates to link these two fragments (right). The 3D coordinates of the generated linker must align with the two fragments, otherwise they cannot link.

A critical challenge for computational linker design stems from its strong 3D spatial constraints compared to classical graph generation tasks. It is known that a successful fragment linker should not disturb the spatial configurations of the two fragments (Ichihara et al., 2011; Klön, 2015). In addition, the anchors between fragments and linker also have correlations with their spatial poses. The linker design problem should be extended to include the 3D information, and a 3D-aware generative model is needed to generate realistic linkers, rather than invoking a graph generative model. In this paper, we propose a conditional generative model, named 3DLinker, that jointly models the 2D molecular graphs and 3D structures of linker for solving the 3D linker design problem.

Given the graph and spatial coordinates of two fragments, 3DLinker can jointly generate graphs and spatial coordinates of the linker. Notably, it does not rely on pre-determined anchors and can accurately predict anchors based on the two observed fragments to be connected. More importantly, 3DLinker can predict 3D coordinates directly and at the same time keeps equivariant to rotations, translations and reflections, which makes it insensitive to choices of the coordinate system. Finally, since the generative model is based on the variational autoencoder (VAE) framework (Kingma & Welling, 2013), it can be used as an unsupervised representation learning method whose latent representations are fed into downstream tasks such as drug-likeness prediction. To the best of our knowledge, 3DLinker is the first trial that simultaneously predicts equivariant graph and 3D coordinates for the linker design problem.

## 2. Background

In this section we introduce the definition of the 3D linker design problem and basic concepts of E(3) equivariance.

### 2.1. 3D Linker Design

Linker design is to generate a small molecule that can link two given molecular fragments at specific anchors (bind-

ing atoms). Instead of modeling it as a 2D graph generation problem, it is crucial to take 3D information into account, since the designed fragment should satisfy spatial constraints such as not disturbing the relative poses or causing any atom clashes between the two fragments. Therefore, it requires the linker design algorithm to be able to generate both the chemical graph and 3D coordinates given the two fragments and their spatial coordinates (Figure 1).

Mathematically, a molecule can be viewed as a graph with atoms as nodes and chemical bonds as edges. A 3D molecule can be represented by a graph  $G = (V; E; X)$  with 3D coordinates  $r = (x; y; z)$ , where  $V$  is the set of nodes,  $|V|$  is the number of nodes,  $E \subseteq V \times V$  are edges,  $X$  are node types and  $R$  is a matrix whose  $i$ -th row is  $r_i^T$  ( $\supset$  stands for transpose). In the 3D linker design, two fragments are defined by two unlinked subgraphs  $G_F = (G_{F,1}; G_{F,2})$  with geometry  $R_F$ , and a linker is denoted by  $G_L$  with geometry  $R_L$ . Let  $G; R$  be the graph and geometry of the ground truth linked molecule containing both fragments and linker. A 3D linker design model is a conditional generative model that completes the ground truth molecule graph as well as its geometry given the two fragments:

$$p(G; R | G_F; R_F): \quad (1)$$

### 2.2. O(3), E(3) Groups and Equivariance

Group is a set of operations equipped with multiplications, associativity, the identity element and the inverse element. Group of all 3D rotations and reflections is called 3D Orthogonal Group or O(3), and group of all 3D rotations, translations and reflections is called 3D Euclidean Group or E(3). Let  $X$  be the input space,  $Y$  be output space and  $GL(X)$  be all invertible linear transformations from  $X$  to  $X$  (similar for  $GL(Y)$ ). A function  $f: X \rightarrow Y$  is called *equivariant* to the group  $G$ , if for all group element  $g \in G$  and all  $x \in X$ , there exists group representations  $\rho^X: G \rightarrow GL(X)$  and  $\rho^Y: G \rightarrow GL(Y)$  such that

$$\rho^Y(g) \circ f(x) = f(\rho^X(g)x): \quad (2)$$

If  $\rho^Y(g)$  is the identity function for all  $g \in G$ , then we say  $f$  is *invariant* to group  $G$ . In 3D linker design, it is known that a molecule graph  $G$  should not depend on a specific coordinate system and  $R$  should change equivariantly to transformations of the coordinate system. Therefore it raises a constraint that for any  $g \in E(3)$ , the generative model  $p(G; R | G_F; R_F)$  should satisfy

$$p(G; \rho(g)R | \rho(g)G_F; \rho(g)R_F) = p(G; R | G_F; R_F): \quad (3)$$

where  $\rho(g)$  can be any rotation, translations or reflections matrix in the 3D space.

### 3. Related Work

**Graph-based Molecular Generation.** Variational Graph Auto-Encoders are the most popular models for molecular generations. Early approaches mainly focus on embedding the 2D chemical graphs into low dimensional space and sample new molecules by perturbing the hidden values. The representative works include GraphVAE (Simonovsky & Komodakis, 2018), CGVAE (Liu et al., 2018), JT-VAE (Jin et al., 2018), GraphNVP (Madhawa et al., 2019) and so on. Although none of these methods designs molecular linkers, graph-based VAE models serve as the basic building block of our entire architecture. All these models could be plugged into our framework by transforming the generative model into a conditional one. From the perspective of training techniques, auto-regression is also widely adopted to train graph-based deep learning models, such as GraphRNN (You et al., 2018), DeepGMG (Li et al., 2018), GraphAF (Shi et al., 2019). Most of these models generate nodes and edges in a sequential manner.

**Point-cloud-based Molecular Models.** An important component of our work is to model and design the 3D structures of molecular fragments, in which maintaining the equivariant properties is the crucial computational challenge. One typical solution is to model the molecules as 3D point clouds using equivariant neural networks (Sic et al., 2017; Klicpera et al., 2020; Liu et al., 2021; Satorras et al., 2021b; Thomas et al., 2018; Fuchs et al., 2020; Deng et al., 2021; Jing et al., 2020). To train such models, auto-regression is a more common solution, such as G-SchNet (Gebauer et al., 2019), G-SphereNet (Anonymous, 2022), and (Luo et al., 2021). But low-based model ENF (Satorras et al., 2021a) and reinforcement learning (Simm et al., 2020a;b) could also be applied. The main limitation of point-cloud-based models is that they cannot directly generate discrete graph structures, which makes it difficult to model chemical constraints like valency (maximal number of hydrogen atoms one can combine with).

**Molecular Linker Design.** DeLinker (Imrie et al., 2020) is the first attempt to apply deep learning methods to the linker design problem. It constructs a conditional graph generative model that generates linker given two fragments. It adapts CGAVE (Liu et al., 2018), generating edges step by step starting with fragments and two known anchor nodes as the binding sites. The spatial distance and angle between two fragments are provided to the model as side information to guide the generation. DEVELOP (Imrie et al., 2021) improves DeLinker by encoding the spatial information of fragments using CNN. SyntaLinker (Yang et al., 2020) is a text-based transformer that directly transforms the SMILES (a text representation of molecular graphs) of input fragments into ground-truth ones. However, none of them have a detailed atom-level description of molecule geometry, which

is insufficient to express the fragments' geometry. In addition, anchor nodes are either not considered or assumed to be known in advance, the latter of which is rare in a real-world application. Most importantly, they are only able to generate graph representations of linker without 3D coordinates. Gen3D (Roney et al., 2021), though generating both graphs and coordinates equivariantly, is not designed for linker design.

### 4. Methodology

In this section, we present our 3DLinker, a conditional VAE-based generative model that generates both invariant graphs and equivariant absolute coordinates of linkers given two 3D fragments.

**Notations.** Let  $G_F; G_L; G$  be graphs of fragments, linker and full molecule (ground truth) respectively, and similarly for coordinates  $R_F; R_L; R$  as in section 2.1. As we complete the full molecule graph step by step, we use  $G_t$  and  $R_t$  to denote the current (existing) graph and coordinates at timestamp  $t$  where  $G_0 = G_F; R_0 = R_F$ . The encoder embeds each node  $v$  with both invariant features  $h_i \in \mathbb{R}^{n_h}$  (for embedding the graph) and equivariant features  $v_i \in \mathbb{R}^{n_v \times 3}$  (for embedding the coordinates), which are further used for sampling invariant latent variables  $z_i^h \in \mathbb{R}^{m_h}$  and equivariant latent variables  $z_i^v \in \mathbb{R}^{m_v \times 3}$ . Symbol  $h; v; z^h; z^v$  without subscripts refer to that variable for all nodes in a general sense. For column vectors  $a \in \mathbb{R}^c$  and  $b \in \mathbb{R}^c$ , we use  $a \odot b \in \mathbb{R}^c$  to denote point-wise multiplication and  $\text{diag}(a) \in \mathbb{R}^{c \times c}$  to denote a matrix whose diagonal is  $a$  and zero otherwise.

**Equivariant Features For Coordinates Predictions.** The equivariant nature makes it difficult to predict absolute coordinates directly. Many existing works (Gebauer et al., 2019; Anonymous, 2022; Xu et al., 2021) tackle this problem by encoding coordinate information as invariant node features and predicting invariant quantities such as distances and edge angles. However, these indirect methods are either computationally intensive (need transformation to local coordinate system (Anonymous, 2022)) or introduce extra error from the second nonconvex optimization (for translating distance matrices into absolute coordinates (Xu et al., 2021)). Instead, we propose to generate absolute coordinates directly while preserving equivariance. To generate equivariant coordinates, only leveraging invariant features is not enough: we cannot produce an equivariant quantity arbitrarily by combining invariant quantities. Therefore, in addition to invariant node features, we need to introduce extra equivariant node features that can be directly used for composing equivariant coordinates. We use notations  $h$  for invariant features and  $v$  for equivariant features.

**Vector Neurons.** Classical fully connected neural networks

or MLPs cannot preserve equivariance and thus is not suitable for transforming equivariant features. In this regard, vector neuron networks or VN-MLP (Deng et al., 2021) propose a ReLU-like nonlinear function for equivariant features. Concretely, given an equivariant input  $v \in \mathbb{R}^{n_v \times 3}$ , Vector-ReLU learns two weight matrices  $W \in \mathbb{R}^{n_v \times n_v}$  and  $U \in \mathbb{R}^{n_v \times n_v}$  to map  $v$  to output  $v^0 \in \mathbb{R}^{n_v \times 3}$  via

$$q = W v \in \mathbb{R}^{n_v \times 3}; \quad k = U v \in \mathbb{R}^{n_v \times 3}; \quad (4a)$$

$$v^0 = q \cdot \text{diag}(\mathbb{1}_{h_q; k_i < 0}) \cdot h_q; \quad \frac{k}{\|k\|} \cdot i \cdot \frac{k}{\|k\|}; \quad (4b)$$

where  $h_q; k_i \in \mathbb{R}^{n_v}$  is the inner product in the last axis,  $\mathbb{1}_{h_q; k_i < 0}$  is the indicator function, and  $\|k\| \in \mathbb{R}^{n_v}$  is the norm of  $k$  over the last axis. It is easy to verify that this is equivariant, since both  $q$  and  $k$  are linear combinations of equivariant input  $v$  while coefficients  $h_q; k_i$  is invariant. Intuitively, Vector-ReLU projects  $q$  to the orthogonal plane of a learnable direction  $k$  if  $q$  lies in the other side of the plane, which is analogous to the cutoff in classic ReLU. This nonlinearity enhances the expressive power while preserving the equivariance. We use VN-MLP to denote a neural network stacked by multiple Vector-ReLU units.

**Mixed-Features Message Passing** Now we are ready to introduce our Mixed-Features Message Passing (MF-MP) scheme. MF-MP performs message passing for invariant features  $h$  and equivariant features  $v$  simultaneously, and in each step the two types of features are properly mixed so that 1) their respective invariance and equivariance properties are preserved, and 2) one type of feature helps the update of the other type and vice versa.

In the first step, invariant features  $h \in \mathbb{R}^{n_h}$  and equivariant features  $v \in \mathbb{R}^{n_v \times 3}$  are transformed and mixed to construct new expressive intermediate features  $h^0, v^0$  by

$$h_j^0 = \text{VN-MLP}_1(v_j) \in \mathbb{R}^{n_h}; \quad (5a)$$

$$h_j^{00} = \text{VN-MLP}_2(v_j) \in \mathbb{R}^{n_v}; \quad (5b)$$

$$v_j^0 = \text{diag}(\mathbb{1}_{h_j^0}) \cdot \text{VN-MLP}_3(v_j) \in \mathbb{R}^{n_v \times 3}; \quad (5c)$$

Next, point convolution (Thomas et al., 2018; Stret et al., 2017; 2021) is applied to linearly transform the mixed features  $h^0, h^{00}, v^0$  into messages:

$$m_{ij}^h = \text{Ker}_1(kr_{ij}; k) \cdot h_j^0; \quad (6a)$$

$$m_{ij}^v = \text{diag}(\text{Ker}_2(kr_{ij}; k)) \cdot v_j^0 + \text{Ker}_3(kr_{ij}; k) \cdot h_j^{00} \cdot r_{ij}^{\geq}; \quad (6b)$$

where  $r_{ij} = r_i - r_j$  is the relative displacement,  $\text{Ker}$  are learnable kernels such as RBFs that transform a scalar distance into a multi-dimensional output vector using different shape parameters, making the messages geometry-aware. Intuitively, it reflects discrete levels of physical interactions

(short-range, long-range) at different distances. More details on  $\text{Ker}$  are given in Appendix B.

Finally, Gated Recurrent Units (GRU) (Li et al., 2015) and VN-MLP are applied as powerful nonlinear transformations to update the node features with the messages:

$$h_i = \text{GRU}(h_i; \sum_{j \in \mathcal{N}(i)} m_{ij}^h); \quad (7a)$$

$$v_i = \text{VN-MLP}_4(v_i; \sum_{j \in \mathcal{N}(i)} m_{ij}^v); \quad (7b)$$

Here  $\mathcal{N}(i)$  stands for neighbors of  $i$ . Our Mixed-Features Message Passing (MF-MP) above effectively mixes invariant and equivariant features in each step to help the update of each other with powerful nonlinear functions. Proof of MF-MP's equivariance w.r.t. E(3) is included in Appendix A. We also discuss how it relates to and differs from Tensor Field Networks (Thomas et al., 2018) in appendix B.

Now we describe details about the encoder and decoder of 3DLinker using MF-MP as building blocks. 3DLinker is a conditional latent generative model  $(G; R; jG_F; R_F)$  including an encoder  $q(z^h; z^v; jG_F; G; R_F; R)$ , a decoder  $p(G; R; jG_F; R_F; z^h; z^v)$  and a prior  $p(z^h; z^v; jG_F; R_F)$ .

#### 4.1. Encoder

The encoder  $q(z^h; z^v; jG_F; G; R_F; R)$  computes node-level latent distributions utilizing MF-MP. Initially, invariant features  $h$  are embeddings of node types and we let equivariant features  $v = 0$ . After applying several times of MF-MP, we obtain the final node features  $h$  and  $v$ . The latent variables are sampled by  $z_i^h \sim \mathcal{N}(h_i; (\sigma_i^h)^2 I)$ ,  $z_i^v \sim \mathcal{N}(v_i; (\sigma_i^v)^2 I)$  for linker nodes only, where the means and variances are computed from the final node features:

$$\sigma_i^h \in \mathbb{R}^{V_L}; \quad h_i = \text{MLP}_4(h_i); \quad (\sigma_i^h)^2 = \text{MLP}_5(h_i); \quad (8a)$$

$$\sigma_i^v = \text{VN-MLP}_5(v_i); \quad (\sigma_i^v)^2 = \text{MLP}_6(h_i); \quad (8b)$$

Note that for equivariant latent variables  $z^v$  the covariance  $(\sigma_i^v)^2 I$  assign the same variance to  $x, y, z$  directions, which is the simplest way to preserve equivariance. Since fragments  $(G_F; R_F)$  are given during generation, there is no need to sample their latent variables. Instead, we run the same (weight-sharing) MF-MP network again on fragments only, which gives another set of final node features  $h^F, v^F$  for fragment nodes. Latent variables of fragment nodes are deterministically obtained by

$$\sigma_i^F \in \mathbb{R}^{V_F}; \quad z_i^h = \text{MLP}_7(h_i^F); \quad z_i^v = \text{VN-MLP}_6(v_i^F); \quad (9a)$$

Figure 2: Illustration of overall encoding and decoding process. For encoding, ground truth is sent into a MF-MP encoder to get node-level representations. Those representations of nodes in fragments are discarded and replaced by representations that are computed separately on the fragments graph only. For decoding, two anchor nodes are predicted as the binding sites for the linker. Node Types of the linker are simultaneously predicted before linking. With two anchor nodes and node types of the linker, edges and coordinates are sequentially predicted, as demonstrated in Figure 3.

#### 4.2. Decoder

The decoder  $p(G; R | G_F; R_F; z^h; z^v)$  constructs  $(G; R)$  from fragments  $(G_F; R_F)$  in a sequential manner. In the decoding process we incorporate the valency rules of molecules by masking out impossible edges and anchor nodes. The decoding process consists of the following steps:

- (1) Anchor Node Prediction: predict anchor nodes  $(a_1; a_2)$  for the two fragments. These two anchor nodes are served as the binding points for linker to connect.
- (2) Node Type Prediction: predict node type  $X$  for all linker nodes.
- (3) Edge and Coordinate Prediction: Put the two anchor nodes in a queue. Then do the following until the queue is empty:
  - (i) Pop a node  $e$  from the queue and denote it as the current focus node
  - (ii) Predict an edge between the focus node and another node. The connected nodes are added to the queue. If node  $e$  is a linker node and is connected to the existing graph for the first time, predict its coordinates.
  - (iii) Repeat (ii) and (iii) until an artificial stop node is connected. Update the coordinates of all nodes in the current linker. The focus node is then marked as closed, which cannot be added to the queue or connected anymore. Then go back to (i).

Mathematically we factorize the joint probability into:

$$p(G; R | G_F; R_F; z^h; z^v) = p(E; X; R | E_F; X_F; R_F; z^h; z^v) \cdot p(a_1; a_2 | z^h; z^v) \cdot p(X | z^h) \cdot p(\{z\} | z^h; z^v) \quad (10)$$

Anchor
Node Types
Edges and Coordinates

where  $E_T = E$  and  $R_T = R$ . Details of each component are explained in the following.

**Anchor Node Prediction.** To jointly predict two anchor nodes, we further factorize the joint probability into, the probability of anchor  $a_1$  on the first fragment  $G_{F,1}$ , and  $p_{a_2}$ , the probability of anchor  $a_2$  on the second fragment  $G_{F,2}$  conditioning on  $a_1$ :

$$p(a_1; a_2 | z^h; z^v) = p(a_1 | f(z_i^h; z_i^v)_{i \in V_{F,1}}) \cdot p(a_2 | z_{a_1}^h; z_{a_1}^v; f(z_i^h; z_i^v)_{i \in V_{F,2}}) \quad (11)$$

Concretely, each node on the first fragment will get a score  $c_i = g(z_i^h; kA_1 z_i^v k)$ , where  $A_1 \in \mathbb{R}^{n_v \times n_v}$  is a learnable linear transformation. The scores are then passed to a softmax to compute the anchor probability for nodes of the first fragment:  $p_{a_1} = \exp(c_{a_1}) / \sum_{i \in V_{F,1}} \exp(c_i)$ . Then the latent variables of this predicted anchor node as well as nodes of the second fragment are used to compute another group of scores  $s_i = g(z_i^h; kA_1 z_i^v k; z_{a_1}^h; A_1 z_{a_1}^v)$ ,

and the probability of the second anchor  $p_{a_2} = \exp(c_{a_2}^0) / \sum_{i \in \mathcal{V}_{F,2}} \exp(c_i^0)$ .

**Node Type Prediction.** Node types of the linker are directly predicted using their latent variables. We leverage the self-attention mechanism (Vaswani et al., 2017) to obtain new node features, which are then passed to an MLP to get the logits of node types. After node types are sampled, the types' embeddings are concatenated to the corresponding latent variables  $\mathbf{z}^h$  for the latter procedures.

outputs equivariant coordinates for node  $i$ :

$$p_{i,j} = \sigma_{11}(\mathbf{z}_i^h; \mathbf{z}_j^h; \mathbf{h}A_3 \mathbf{z}_i^v; A_4 \mathbf{z}_j^v); \quad (13a)$$

$$q_{i,j} = \sigma_{12}(\mathbf{z}_i^h; \mathbf{z}_j^h; \mathbf{h}A_5 \mathbf{z}_i^v; A_6 \mathbf{z}_j^v); \quad (13b)$$

$$\mathbf{r}_t = \mathbf{r} + \sum_{j \in \mathcal{V}_t} p_{i,j} (\mathbf{r}_j - \mathbf{r}) + \text{VN-MLP}_7 \left( \sum_{j \in \mathcal{V}_t} q_{i,j} \text{VN-MLP}_8(\mathbf{z}_i^v; \mathbf{z}_j^v) \right) \quad (13c)$$

for any generic coordinate  $\mathbf{r}$  scaled reference point. Here VN-MLP taking two inputs means concatenation along the first axis ( $\mathbf{r}_v$ ). The idea is to compute pair-wise interactions (13a, 13b) and predict a deviation from reference point  $\mathbf{r}$  (13c). If a linker node  $i$  is first connected to the graph, we use the mass center of the current graph  $\mathbf{r}_t = \sum_{j \in \mathcal{V}_t} \mathbf{r}_j / |\mathcal{V}_t|$  as the reference point and predict its absolute coordinates  $\mathbf{r}_i = \text{pred}_i(\mathbf{z}_i^h; \mathbf{z}_i^v; \mathbf{R}_t; \mathbf{r}_t)$ ; once the stop node is chosen, all linker nodes in the current graph will update their coordinates using their current coordinates as reference points,  $\mathbf{r}_i = \text{updt}_i(\mathbf{z}_i^h; \mathbf{z}_i^v; \mathbf{R}_t; \mathbf{r}_i)$ . Note that  $\text{pred}$  and  $\text{updt}$  have distinct network weights.

Figure 3: Illustration of sequential predictions of edges and coordinates. We first pick up on a node to focus. Then we

sample an edge between the focus node and other nodes (including an artificial stop node). If a linker node is first connected to the existing graph, its coordinates will be predicted. Each time before prediction MF-MP is applied to capture information from the existing graph. We keep adding edges until the stop node is selected, and then coordinates of all link nodes in the existing graph will be simultaneously updated. We then refocus on a new node and repeat. The procedure continues until all nodes in the linker have been focused.

**Edge and Coordinate Prediction.** The edge and coordinate generation path  $(E_0; R_0) \rightarrow (E_1; R_1) \rightarrow \dots \rightarrow (E; R)$  is defined by a sequence of node focusing, edge prediction and coordinate prediction procedures. The node focusing and edge connecting order is pre-determined by a Breath-First Search to enable teacher-forcing training. When a focus node is picked, we add new edges to it until connecting to the stop node. Concretely, at each step we first apply MF-MP (with different weights from that in the encoder) to obtain updated node features  $\mathbf{z}^h; \mathbf{z}^v$  from the initial latent variables  $\mathbf{z}^h; \mathbf{z}^v$ . Then we compute the probability for edge  $E_{i,f}$  by

$$s_{i,f} = \sigma_{10}(\mathbf{z}_i^h; \mathbf{z}_f^h; \mathbf{k}A_2 \mathbf{z}_i^v; A_2 \mathbf{z}_f^v; \mathbf{z}_i^h; \mathbf{z}_f^h);$$

$$p(E_{i,f}) = \frac{\exp(s_{i,f})}{\sum_{j \in \mathcal{V}} \exp(s_{j,f})}; \quad (12)$$

To predict coordinates, we define  $\mathbf{r}_i(\mathbf{z}_i^h; \mathbf{z}_i^v; \mathbf{R}_t; \mathbf{r})$  that

### 4.3. Training Using ELBO

Variational autoencoder (Kingma & Welling, 2013) is trained by maximizing the Evidence Lower Bound (ELBO):

$$\mathcal{L}(\theta; \phi) = \mathbb{E}_{z \sim q} \log p(G; R | G_F; R_F; z^h; z^v) - D_{KL}(q \parallel p); \quad (14)$$

where the prior  $p(z^h; z^v | G_F; R_F)$  simply takes standard Gaussian for all linker nodes. The reconstruction error term  $\mathbb{E}_{z \sim q} \log p(G; R | G_F; R_F; z^h; z^v)$  is approximated by one Monte-Carlo sampling, and we apply teacher forcing (Kolen & Kremer, 2001) for anchor, node type and edge prediction following the pre-determined order. The loss of anchor node prediction, node type prediction and edge prediction are standard cross entropy while for loss of coordinate prediction we use log-MSE used by (Yu, 2020).

### 4.4. Generation

During generation, a maximum number of linker nodes is set and we sample this maximum number of latent variables  $\mathbf{z}^h; \mathbf{z}^v$  for linker nodes (though some nodes might never be included). Then the generation follows the same procedure as the decoder except that there is no teacher forcing.

## 5. Experiments

### 5.1. Experiment Setup

**Dataset.** To evaluate our model, we choose a subset of ZINC (Sterling & Irwin, 2015). For each molecule, we

Table 1: Performance metrics for generated molecules on ZINC dataset. Uncertainty is estimated by deviation.

Metrics (%)	Valid	Recover	Pass 2D	Iters	$S_{\text{Tanim}}$	RMSD (Å)	Unique	Novel
3DLinker (given anchor)	99.20 <sub>0:04</sub>	95.24 <sub>1:17</sub>	90.39 <sub>0:11</sub>	55.06 <sub>0:09</sub>	0.081 <sub>0:005</sub>	0.081 <sub>0:004</sub>	29.20 <sub>0:11</sub>	32.16 <sub>0:12</sub>
3DLinker	98.68 <sub>0:03</sub>	93.75 <sub>0:36</sub>	90.32 <sub>0:18</sub>	54.83 <sub>0:08</sub>	0.081 <sub>0:004</sub>	0.081 <sub>0:004</sub>	29.34 <sub>0:17</sub>	32.53 <sub>0:12</sub>
DeLinker+ConfVAE	98.35 <sub>0:07</sub>	80.35 <sub>2:58</sub>	89.92 <sub>0:53</sub>	52.86 <sub>0:04</sub>	1.345 <sub>0:028</sub>	1.345 <sub>0:028</sub>	44.53 <sub>0:42</sub>	39.53 <sub>0:58</sub>
GraphAF+ConfVAE	34.25 <sub>0:17</sub>	21.23 <sub>1:82</sub>	82.00 <sub>2:21</sub>	38.00 <sub>0:14</sub>	1.263 <sub>0:124</sub>	1.263 <sub>0:124</sub>	84.06 <sub>0:31</sub>	78.33 <sub>0:12</sub>
GraphVAE+ConfVAE	63.07 <sub>0:59</sub>	1.40 <sub>1:37</sub>	86.17 <sub>0:26</sub>	51.31 <sub>0:07</sub>	1.523 <sub>0:478</sub>	1.523 <sub>0:478</sub>	43.59 <sub>0:35</sub>	91.59 <sub>1:47</sub>
SyntaLinker+ConfVAE	80.14 <sub>0:31</sub>	85.47 <sub>2:46</sub>	97.41 <sub>0:11</sub>	54.80 <sub>0:05</sub>	1.402 <sub>0:016</sub>	1.402 <sub>0:016</sub>	41.51 <sub>0:75</sub>	13.17 <sub>0:27</sub>
Gen3D	75.92 <sub>0:10</sub>	37.99 <sub>0:67</sub>	81.41 <sub>0:12</sub>	48.67 <sub>0:06</sub>	1.415 <sub>0:272</sub>	1.415 <sub>0:272</sub>	49.86 <sub>0:49</sub>	64.17 <sub>0:21</sub>

perform 20 times of MMFF force field optimization using RDKit (Landrum) and choose the one with the lowest energy as the ground truth. Following the same procedure from (Hussain & Rea, 2010), the (fragments, linker) pairs are produced by enumerating all double cuts of acyclic single bonds that are not within any functional groups. In total, we obtain 365,749 (fragments, linker, coordinates) triplets and randomly split them into training (365,039), validation (351) and test (358).

Evaluation. We evaluate the generated molecules for multiple 2D (graph) and 3D (coordinates) metrics, including the standard ones such as validity, uniqueness and novelty (Brown et al., 2019). In addition, we also evaluate the percentage of generated molecules passing 2D properties, including synthetic accessibility (Ertl & Schuffenhauer, 2009), ring aromaticity, and pan-assay interference compounds (PAINS) (Baell & Holloway, 2010). After filtering by validity and 2D property filters, the recovery rate is calculated to report the percentage of generated molecules that perfectly recover the ground truth molecule graphs. We also compare Tanimoto similarity using the Morgan fingerprint provided by RDKit, which estimates the similarity between ground-truth and generated molecular graphs. To evaluate the quality of the 3D structures, the predicted 3D structures are compared to the ground truth using root-mean-square deviation (RMSD). Note that RMSD is only computed for generated molecules that perfectly recover the ground truth molecular graphs (including their isomorphic variants), since only recovered molecules have atom-to-atom alignment to ground truth. Also note that we do not apply MMFF to the generated structures since hydrogen is not explicitly included in the dataset for computational efficiency. Following DeLinker, we compute another 3D metric, named shape-and-color similarity score ( $SC_{\text{RDKit}}$ ). Appendix C contains more details about the evaluation standards.

Baselines. Though there are works of molecular graph generative models and molecular geometry prediction given molecular graph, they rarely focus on fragment linking or jointly modeling both graph and geometry. Existing molecule generative models either do not work on conditional (linker) generation, or cannot predict 3D coordinates. Therefore, we implement multiple baselines by adapting

Table 2:  $SC_{\text{RDKit}}$  score distribution (%) and averaged score on ZINC.

Metrics	$SC_{\text{RDKit}}$ Fragments			
	> 0.7	> 0.8	> 0.9	Average
3DLinker (given anchor)	43.10	16.09	2.60	0.684
3DLinker	42.55	15.85	2.49	0.683
DeLinker+ConfVAE	39.96	13.39	1.93	0.675
GraphAF+ConfVAE	19.33	3.36	0.32	0.624
GraphVAE+ConfVAE	13.17	2.15	0.00	0.601

multiple generative models to conditional generative models to generate 2D linker graphs given two molecular fragments. The generated graphs are then taken by a molecular geometry prediction model, ConfVAE (Xu et al., 2021), to predict the 3D coordinates of each atom. Our baselines include DeLinker+ConfVAE, GraphAF+ConfVAE, GraphVAE+ConfVAE, SyntaLinker+ConfVAE and Gen3D. DeLinker and SyntaLinker are existing baselines for conditional linker generation. GraphAF is an autoregressive model, and GraphVAE is a VAE-based model with graph-level encodings. Gen3D is an autoregressive model using EGNN that can jointly generate graphs and coordinates. The latter three are adapted to conditional graph generation. For ConfVAE, we modify its decoder to predict linker coordinates conditionally. Please see appendix C for implementation details.

## 5.2. Results

We trained 3DLinker for 20 epochs using Adam optimizer with a learning rate 0.006, batch size 48 and KL trade-off = 0.6. Training details for other baselines are included in Appendix C. Each model generates 250 samples per two fragments, which leads to in total  $250 \times 358 = 89500$  samples. We conduct such generations three times independently for uncertainty estimation of metrics in the main table 1. Note that since DeLinker takes anchor nodes as known ground truth, we add another comparison where anchor nodes are known to 3DLinker when generation, denoted as 3DLinker (given anchor). Results in Table 1 and Table 2 show that 3DLinker could generate valid and similar linker graph structures with a higher recovery rate, and at the same time achieve accurate predictions of the 3D coordinates of each atom. An interesting observation is that although focusing on 3D structures, 3DLinker achieves superior recovering

Figure 4: An example of fragment linking. The top-5 similar  $S_{RDKit}$  Fragments proposed by 3DLinker (1st row) and DeLinker+ConfVAE (second row) are shown. Generations from 3DLinker are more realistic and similar to ground truth in terms of  $S_{RDKit}$  and 3D geometry.

Table 3: Root Mean Squared Error of QED prediction on ZINC dataset.

Metrics	RMSE <sub>QED</sub>
3DLinker	0.0833
DeLinker+ConfVAE	0.1077
GraphVAE+ConfVAE	0.1179

accuracy of the 2D molecular graphs. Our interpretation is that incorporating 3D constraints benefits the reconstruction of 2D graphs, though it might influence the diversity (low novelty and uniqueness). This also explains why the novelty and uniqueness are relatively low because the 3D constraints significantly reduce the valid chemical compound structures. In addition, GraphAF and GraphVAE are not able to obey valency rules during training, which explains their low valid and recovery rates. In terms of 3D structure predictions, the low RMSD of 3DLinker demonstrates the effectiveness of

both equivariant features and coordinate update strategies. We also notice that since SyntaLinker has no constraints on chemical valency or whether output contains the original fragments, it turns out to have low validity. Nevertheless, it has a higher rate of passing ring aromaticity filter (one of 2D filters), which might be because operating on SMILES does not easily break aromaticity, while operating on nodes and edges sometimes does. GEN3D turns out to have a poor RMSD, and we find it is partially because GEN3D predicts distance matrices between atoms: a small perturbation to the distance matrix may lead to significant changes of conformation in long-chain structures, and the error scales up quickly with the increasing of number of nodes. This may also explain the poor performance of ConfVAE, which also makes predictions based on distance matrices.

### 5.3. Ablation Study

In the ablation study, we focus on two key components of 3DLinker: (1) equivariant features vs. invariant features alone; (2) update of all coordinates until the generation finishes, instead of fixing them in the 3D space one by one. Results display a significant performance decrease after removing either equivariant features or coordinate updates, especially for the RMSD and recovery rate. See Appendix D for details.

### 5.4. Molecular Property Prediction

We show a downstream task of molecule property prediction. We use the learned latent variables from VAE models to predict the Quantitative Estimate of Drug-Likeness (QED) by a gated sum:

$$QED = \sum_{i \in G} \sigma \left( \frac{U_{13}(z_i^h; kV z_i^y k) + U_{14}(z_i^h; kU z_i^y k)}{2} \right); \quad (15)$$

where  $U_{13}$  and  $U_{14}$  are two separate neural networks,  $U$  are two linear transform matrices and  $\sigma(\cdot)$  is a sigmoid function. QED is widely adopted to quantify the potential for a small molecule to be a drug while the function of molecular linkers is to connect two existing drugs. Here we adopt QED to check whether 3DLinker produces biological and biochemical meaningful molecules. Training and test sets are in the same setup as in Section 5.1, with QED scores computed by RDKit. For each pre-trained model, we train a QED predictor for 20 epochs with a learning rate of 0.006 and batch size of 48.

As shown in Table 3, 3DLinker achieves the lowest Root Mean Square Error (RMSE) among all the models, suggesting a good expressing power of its learned latent representations.

### 5.5. Visualization

In the end, we present multiple examples of linker design by visualizations (see more examples of visualization in Ap-



pendix E.). In Figure 4, we show the top-5 molecules with highest  $SC_{RDKit}$  generated by 3DLinker (rst row in the middle) and DeLinker+ConfVAE (second row in the middle). It is obvious that molecules from 3DLinker are generally more similar to the ground truth 2D chemical graph, and have a better spatial alignment with the ground truth 3D structure. In practice, when the ground truth is unknown, the “best” linker could be estimated by the generation likelihood.

## 6. Conclusions

We developed 3DLinker, a conditional variational autoencoder that is able to jointly model graph and 3D representations, predict anchor nodes and samples linkers. Experiments show that 3DLinker is able to generate linkers with both a high recovery rate and precise geometry.

There are still limitations to be considered in the future. First, models should be able to sample number of linker nodes instead of setting a maximal number of nodes in advance. Second, though it is known that the spatial configuration of fragments should not be disturbed, in practice slight differences exist between different linkers, which is an avenue for future works to take into account.

## Software and Data

We implement 3DLinker based on the released code of DeLinker (<https://github.com/mrie/DeLinker>). Our code and data are available at <https://github.com/GraphPKU/3DLinker>.

## References

- Anonymous. An autoregressive flow model for 3d molecular geometry generation from scratch. Submitted to The Tenth International Conference on Learning Representations 2022. URL <https://openreview.net/forum?id=C03Ajc-NS5W>. under review.
- Baell, J. B. and Holloway, G. A. New substructure filters for removal of pan assay interference compounds (pains) from screening libraries and for their exclusion in bioassays. *Journal of medicinal chemistry*, 53(7):2719–2740, 2010.
- Brown, N., Fiscato, M., Segler, M. H., and Vaucher, A. C. Guacamol: benchmarking models for de novo molecular design. *Journal of chemical information and modeling*, 59(3):1096–1108, 2019.
- Deng, C., Litany, O., Duan, Y., Poulencard, A., Tagliasacchi, A., and Guibas, L. Vector neurons: A general framework for so(3)-equivariant networks. *arXiv preprint arXiv:2104.12229*, 2021.

Ertl, P. and Schuffenhauer, A. Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *Journal of cheminformatics*, 1(1):1–11, 2009.

Fuchs, F. B., Worrall, D. E., Fischer, V., and Welling, M. Se(3)-transformers: 3d roto-translation equivariant attention networks. *arXiv preprint arXiv:2006.10503*, 2020.

Gebauer, N. W., Gastegger, M., and SchöK, T. Symmetry-adapted generation of 3d point sets for the targeted discovery of molecules. *arXiv preprint arXiv:1906.00957*, 2019.

Griffiths, D. J. and Schroeter, D. Introduction to quantum mechanics. Cambridge University Press, 2018.

Hussain, J. and Rea, C. Computationally efficient algorithm to identify matched molecular pairs (mmps) in large data sets. *Journal of chemical information and modeling*, 50(3):339–348, 2010.

Ichihara, O., Barker, J., Law, R. J., and Whittaker, M. Compound design by fragment-linking. *Molecular Informatics*, 30(4):298–306, 2011.

Imrie, F., Bradley, A. R., van der Schaar, M., and Deane, C. M. Deep generative models for 3d linker design. *Journal of chemical information and modeling*, 60(4):1983–1995, 2020.

Imrie, F., Hadfield, T. E., Bradley, A. R., and Deane, C. M. Deep generative design with 3d pharmacophoric constraints. *bioRxiv*, 2021.

Jin, W., Barzilay, R., and Jaakkola, T. Junction tree variational autoencoder for molecular graph generation. In International conference on machine learning, pp. 2323–2332. PMLR, 2018.

Jin, W., Barzilay, R., and Jaakkola, T. Hierarchical generation of molecular graphs using structural motifs. In International Conference on Machine Learning, pp. 4839–4848. PMLR, 2020.

Jing, B., Eismann, S., Suriana, P., Townshend, R. J., and Dror, R. Learning from protein structure with geometric vector perceptrons. *arXiv preprint arXiv:2009.01411*, 2020.

Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Klicpera, J., Groß, J., and Günnemann, S. Directional message passing for molecular graphs. *arXiv preprint arXiv:2003.03123*, 2020.

Klön, A. E. Fragment-Based Methods in Drug Discovery. Springer, 2015.

- Kolen, J. F. and Kremer, S. *A field guide to dynamical recurrent networks* John Wiley & Sons, 2001.
- Landrum, G. Rdkit: Open-source cheminformatics <http://www.rdkit.org/>. Accessed: 2022/1/13.
- Landrum, G. A., Penzotti, J. E., and Putta, S. Feature-map vectors: a new class of informative descriptors for computational drug discovery. *Journal of computer-aided molecular design* 20(12):751–762, 2006.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)* pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Li, Y., Tarlow, D., Brockschmidt, M., and Zemel, R. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493* 2015.
- Li, Y., Vinyals, O., Dyer, C., Pascanu, R., and Battaglia, P. Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324* 2018.
- Liu, Q., Allamanis, M., Brockschmidt, M., and Gaunt, A. Constrained graph variational autoencoders for molecular design. *Advances in Neural Information Processing Systems* 31:7795–7804, 2018.
- Liu, Y., Wang, L., Liu, M., Zhang, X., Oztekin, B., and Ji, S. Spherical message passing for 3d graph networks. *arXiv preprint arXiv:2102.05013* 2021.
- Luo, S., Guan, J., Ma, J., and Peng, J. A 3d generative model for structure-based drug design. *Advances in Neural Information Processing Systems* 34, 2021.
- Madhawa, K., Ishiguro, K., Nakago, K., and Abe, M. Graph-nvp: An invertible flow model for generating molecular graphs. *arXiv preprint arXiv:1905.11600* 2019.
- Polishchuk, P. G., Madzhidov, T. I., and Varnek, A. Estimation of the size of drug-like chemical space based on gdb-17 data. *Journal of computer-aided molecular design* 27(8):675–679, 2013.
- Putta, S., Landrum, G. A., and Penzotti, J. E. Conformation mining: an algorithm for finding biologically relevant conformations. *Journal of medicinal chemistry* 48(9):3313–3318, 2005.
- Roney, J. P., Maragakis, P., Skopp, P., and Shaw, D. E. Generating realistic 3d molecules with an equivariant conditional likelihood model. 2021.
- Satorras, V. G., Hoogeboom, E., Fuchs, F. B., Posner, I., and Welling, M. E(n) equivariant normalizing flows for molecule generation in 3d. *arXiv preprint arXiv:2105.09016* 2021a.
- Satorras, V. G., Hoogeboom, E., and Welling, M. E(n) equivariant graph neural networks. *arXiv preprint arXiv:2102.09844* 2021b.
- Schütt, K., Kindermans, P.-J., Felix, H. E. S., Chmiela, S., Tkatchenko, A., and Müller, K.-R. SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. In *NIPS* 2017.
- Schütt, K. T., Unke, O. T., and Gastegger, M. Equivariant message passing for the prediction of tensorial properties and molecular spectra. *arXiv preprint arXiv:2102.03150* 2021.
- Shi, C., Xu, M., Zhu, Z., Zhang, W., Zhang, M., and Tang, J. Graphaf: a flow-based autoregressive model for molecular graph generation. *International Conference on Learning Representations* 2019.
- Simm, G., Pinsler, R., and Hernández-Lobato, J. M. Reinforcement learning for molecular design guided by quantum mechanics. *International Conference on Machine Learning* pp. 8959–8969. PMLR, 2020a.
- Simm, G. N., Pinsler, R., Gysi, G., and Hernández-Lobato, J. M. Symmetry-aware actor-critic for 3d molecular design. *arXiv preprint arXiv:2011.12747* 2020b.
- Simonovsky, M. and Komodakis, N. Graphvae: Towards generation of small graphs using variational autoencoders. In *International conference on artificial neural networks* pp. 412–422. Springer, 2018.
- Sterling, T. and Irwin, J. J. Zinc 15–ligand discovery for everyone. *Journal of chemical information and modeling* 55(11):2324–2337, 2015.
- Thomas, N., Smidt, T., Kearnes, S., Yang, L., Li, L., Kohlhoff, K., and Riley, P. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:1802.08219* 2018.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, ., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems* pp. 5998–6008, 2017.
- Xu, M., Wang, W., Luo, S., Shi, C., Bengio, Y., Gomez-Bombarelli, R., and Tang, J. An end-to-end framework for molecular conformation generation via bilevel programming. *arXiv preprint arXiv:2105.07246* 2021.
- Yang, Y., Zheng, S., Su, S., Zhao, C., Xu, J., and Chen, H. Syntalinker: automatic fragment linking with deep conditional transformer neural network. *Chemical science* 11(31):8312–8322, 2020.

You, J., Ying, R., Ren, X., Hamilton, W., and Leskovec, J. Graphrnn: Generating realistic graphs with deep autoregressive models. International conference on machine learning pp. 5708–5717. PMLR, 2018.

Yu, R. A tutorial on vaes: From bayes' rule to lossless compression. arXiv preprint arXiv:2006.10273, 2020.

## A. Proof of Invariance/Equivariance

In this section we prove that 3DLinker satisfies E(3) equivariance, namely for all  $G \in E(3)$ , we have  $p(G; (g)RjG_F; (g)R_F) = p(G; RjG_F; R_F)$ . For simplicity, we denote E(3) invariance and equivariance as E(3)-inv and E(3)-eqv, and O(3) invariance and equivariance as O(3)-inv and O(3)-eqv, and translational invariance as T-inv.

We first prove  $h$  is E(3)-inv while  $v$  is O(3)-eqv/T-inv after MF-MP (5, 6, 7)

Lemma A.1. In MF-MP,  $h$  is E(3)-inv and  $v$  is O(3)-eqv, if  $h$  is E(3)-inv and  $v$  is O(3)-eqv.

Proof. First let us show that in (5),  $h^0, h^{00}$  are E(3)-inv and  $v^0$  is O(3)-eqv/T-inv. Note that VN-MLP( $v$ ) is E(3)-eqv and T-inv due to the properties of vector neurons, and thus the kernel  $\text{Ker}(\text{VN-MLP}(v))$  is O(3)-inv and T-inv (because O(3) preserves inner product), namely E(3)-inv. Therefore in (5a, 5b) and  $h^0, h^{00}$  are both E(3)-inv. Equation (5c) is essentially an O(3)-eqv VN-MLP( $h_v(v_j)$ ) scaled by E(3)-inv features  $\text{diag}_{h_v}(h_j)g$ , and thus its output  $v_j^0$  is also O(3)-eqv/T-inv.

Then in message function (6a, 6b), message  $g$  is E(3)-inv since kernel  $\text{Ker}$  is only a function of E(3)-inv distance  $r_{ij}$ . For (6b), The first term of RHS is a O(3)-eqv/T-inv feature scaled by E(3)-inv kernel, which makes it O(3)-eqv/T-inv. Similarly for the second term is a O(3)-eqv/T-inv relative displacement scaled by both E(3)-inv kernels and features. Therefore  $h$  is E(3)-inv and  $v$  is O(3)-eqv/T-inv.

Finally, in (7a)  $h$  is E(3)-inv since all its inputs are E(3)-inv. In (7b)  $v$  is O(3)-eqv/T-inv due to vector neurons.  $\square$

Theorem A.2. The generative model  $p(G; RjG_F; R_F)$  (8, 9, 10) satisfies equivariance condition (3).

Proof. By lemma A.1, the encoded latent variables  $z$  and  $v$  from (8, 9) is E(3)-inv and O(3)-eqv/T-inv respectively. In the decoding process, anchor nodes, node types and edges are all predicted from  $z^v$ . Thus the probability of graph  $G$  is E(3)-inv. The coordinates are predicted through (13). Note that  $q_{ij}$  in (13a, 13b) are E(3)-inv and in (13c) is E(3)-eqv (mass center), which implies terms  $\sum_{j \in V_t} p_{ij}(r_j - r)$  and  $\text{VN-MLP}_7(\sum_{j \in V_t} q_{ij} \text{VN-MLP}_8(z_j^v; z_j^v))$  are all O(3)-eqv/T-inv. Finally we can conclude that  $r$  is E(3)-eqv because an O(3)-eqv/T-inv quantity plus an E(3)-eqv quantity results in an E(3)-eqv quantity). Therefore  $R$  is E(3)-eqv.  $\square$

## B. Point Convolution

In the message function (6), kernel  $\text{Ker}(kr_{ij} | k)$  assigns weights relying on distance  $r_{ij} | k$ . Concretely in our implementation, our kernel first applies Gaussian functions with different means and perform a learnable affine transform:

$$\text{Gaussian}(kr_{ij} | k) = \begin{matrix} 0 \\ \vdots \\ 10 \end{matrix} \begin{matrix} \exp(-k(r_{ij} - \mu_1)^2) \\ \exp(-k(r_{ij} - \mu_2)^2) \\ \vdots \\ \exp(-k(r_{ij} - \mu_{10})^2) \end{matrix} \begin{matrix} g_1 \\ g_2 \\ \vdots \\ g_{10} \end{matrix} \begin{matrix} C \\ A \end{matrix}; \quad \text{Ker}(kr_{ij} | k) = \text{Linear}(\text{Gaussian}(kr_{ij} | k)) + \text{Bias} \quad (16)$$

where  $\mu_1 < \mu_2 < \dots < \mu_{10}$  are hyper-parameters and  $k$  is a learnable parameter. The intuition behind is to capture the intensity of interactions between nodes (atoms) with different distances. The largest mean  $\mu_1$  is basically the maximal correlation length: if two nodes are separated by distance beyond  $\mu_1$ , their interaction is nearly neglectable.

Mathematically, message functions (6) can be seen as tensor products using spherical harmonics, as described in Tensor Field Networks (Thomas et al., 2018). In group representation theory, invariant features and equivariant features are called type-0 and type-1 tensors respectively, and the theory (Griffiths & Schroeter, 2018) tells us the correct way to construct new type-0 features or type-1 features using type-1 spherical harmonics  $Y^l(r)$  and  $h; v$ . In our case since we only have type-0 and type-1 features,  $Y^0(r) / r^0$  and  $Y^1(r) / r^1$  are all we need, which explains the design of (6).

Note that there are several differences between our method and tensor products in Tensor Field Networks (TFN). First TFN is based on SO(3) equivariant, while we seek for O(3) equivariant. So terms like cross product are discarded since they violate mirror symmetry. Besides, we mix different types of features before convolution, in contrast to convolution on raw features. Finally, TFN uses simple activation functions like scaling with norm, while we leverage Vector Neuron for novel non-linearity.

## C. Experiment Details

Some evaluation metrics. Validity is defined by percentage of generated molecules that both obey chemical constraints (valency) and successfully links two fragments into connected graphs. Invalid molecules are discarded for the following evaluations. Uniqueness means percentage of non-duplicate generated molecules. Novelty refers to percentage of generated molecules whose linkers are not present in training set.

SC<sub>RDKit</sub> uses two RDKit built-in functions as described in (Putta et al., 2005) and (Landrum et al., 2006) to compute color similarity scores between two 3D molecules based on the overlap of their pharmacophoric features. And the shape similarity score is a simple volumetric comparison between the two 3D molecules. Both scores are between 0 (no match) and 1 (perfect match), which are averaged to produce a final score between 0 and 1. Scores above 0.7 indicate a good match, while scores above 0.9 suggest an almost perfect match. Following DeLiSC<sub>RDKit</sub> is measured only on fragments (we re-generate their coordinates together with the linker), which embodies the capability to generate linkers without disturbing fragments.

3DLinker. We train 3DLinker for 20 epochs with kl trade-off beta 0.6. Note that the anchor node prediction is asymmetric to the permutation  $a_1$  and  $a_2$ , and thus we apply a permutation to enhance our model.

GraphAF. Originally GraphAF is an autoregressive flow model  $p(G) = \prod_t p(G_{t+1} | G_t)$ . To model a conditional probability  $p(G | G_F)$ , all we need is to mask out loss of  $G_F$  and only compute loss starting from  $G_F$  to  $G$ . We trained GraphAF for 170 epochs, and other hyper-parameters are consistent with its source code (<https://github.com/DeepGraphLearning/GraphAF>).

GraphVAE. GraphVAE represents a graph as node type  $F$ , adjacency  $A$  and edge type  $E$  and maps them into a graph-level representation  $z$ . Then  $z$  is decoded into  $F$ ,  $A$ ,  $E$  with an additional graph matching to compute loss. To modify it into a conditional generative models, a naive approach is to build a generative model for  $G$  directly, and then predict the anchor nodes that connects to fragments. Concretely, let  $a_1, a_2$  be anchor nodes of fragments, and  $b_1, b_2$  be the corresponding anchor nodes of linker, the decoder model  $p(G | G_F; z)$  is:

$$p(G | G_F; z) = p(G_L | z_F; z) p(a_1; a_2; b_1; b_2 | z_F; z) p(E_{a_1; b_1}; E_{a_2; b_2} | z_F; z); \quad (17)$$

where  $z_F$  is a graph-level encoding of fragments and  $a_1, b_1; a_2, b_2$  are two edges that connects fragments and linker. The realization is based on code provided by <https://github.com/snap-stanford/GraphRNN>. In experiments, all hyper-parameters are unchanged except KL trade-off beta 0.6.

ConfVAE. ConfVAE is a VAE for geometry generation given graph  $(R | G)$ . It encodes graph  $G$  and distance matrix  $D = f_{d_{ij}} = k r_i - r_j k$ ;  $j \in G$  into latent variables  $z$ , and uses a flow model to decode distance matrix  $D = f(G; z)$ . Concretely, its flow is

$$D = f(G; z) = D(0) + \int_0^z g(G; D(\cdot); z) dz; \quad (18)$$

where  $g$  is a Message Passing Neural Networks (MPNN)  $D(0) \sim N(0; I)$  is a base distribution. To transform  $(R | G)$  to a conditional model  $p(R | G; R_0)$ , we modify the flow to

$$\begin{matrix} D_0 \\ D_L \end{matrix} = f(G; z; D_0) = \begin{matrix} D_0 \\ D_L(0) \end{matrix} + \int_0^z \begin{matrix} 0 & 0 \\ 0 & 1 \end{matrix} g(G; D_0; D_L(\cdot); z) dz \quad (19)$$

where  $D_0 = f_{d_{ij}} = k r_i - r_j k$ ;  $j \in G_0$  are distances we already knew while  $D_L = f_{d_{ij}} = k r_i - r_j k$ ;  $j \in G \setminus G_0$  are distances between nodes with at least one is in linker. After distance matrix predicted, we transform it into absolute coordinates by optimizing the following:

$$\min_{\{r_i\}_{i \in G \setminus G_0}} \sum_{i, j \in G \setminus G_0} (D_{ij} - k r_i - r_j k)^2; \quad (20)$$

Note that we only need to optimize  $\{r_i\}_{i \in G \setminus G_0}$  since  $R_0 = f_{r_i\}_{i \in G_0}$  are given. Also there is no need for coordinates alignment since the coordinate system is well-defined by coordinates of fragments. Code is provided by <https://github.com/MinkaiXu/ConfVAE-ICML21>. Note that although ConfVAE can be trained in an end-to-end manner (both equation 19 and 20), it only makes a few improvements compared to training the flow alone (see the experiments of its original paper (Xu et al., 2021)). Thus we choose only to train the flow model alone in our experiments. For each graph, we sample one geometry that is optimized by (20) with ten times random initialization and 300 steps of gradient descent.

**Gen3D.** Gen3D is an autoregressive model using EGNN as node embedder. Following the same logic of modifying GraphAF, we count loss starting from the two fragments. There are two additional modifications we find useful to make: (1) we predict coordinates only relying on distance distribution; (2) we model the distances as Gaussian distribution (regression) instead of discrete distribution on grids (classification) proposed in the original paper;

## D. Ablation Study

We conduct ablation study on two aspects: removing equivariant features and coordinates update strategy. We train these three models for 20 epochs and evaluate them by the same methods in previous experiments. Results are shown in Table 4 and 5. We can see a dramatic drop of performances after moving either equivariant features or coordinates update.

Table 4: Ablation Study. (eqv-) stands for removing equivariant features while (update-) means removing coordinates update strategy.

Metrics	Valid (%)	Recovered (%)	Pass 2D filters (%)	RMSD	Unique (%)	novel (%)
3DLinker	98.67	93.58	90.37	0.079	29.42	32.48
3DLinker (eqv-)	99.42	86.59	92.68	1.352	34.58	27.02
3DLinker (update-)	98.85	39.94	62.81	0.399	55.93	72.25

Table 5: Ablation study. (eqv-) stands for removing equivariant features while (update-) means removing coordinates update strategy.

Metrics	SC <sub>RDKit</sub> Fragments			Average
	> 0:7 (%)	> 0:8 (%)	> 0:9 (%)	
3DLinker	42.55	15.85	2.49	0.683
3DLinker (eqv-)	38.51	13.15	1.76	0.672
3DLinker (update-)	37.34	10.87	1.13	0.670

Especially, both these two modules contribute greatly to the prediction of coordinates, leading to a decrease of RMSD by 1.3 and 0.3 respectively. Also it is interesting to see that coordinates update has a significant impact on graph quality. A possible reason is that updating the coordinates results in a flexible intermediate coordinates, which may increase the expressive capacity of features. In some sense it is similar to EGNN (Satorras et al., 2021b), who also update intermediate coordinates in the forward pass.

## E. More Visualizations

For some of the ground-truth fragments, we randomly select 8 generated samples for visualization. As shown in the following figures, each row contains the ground-truth (the most left one) and other 8 generated samples. Note that there are examples showing missing fragments, which indicates the failure of fragment linking (invalid generation).

