
Geometry Informed Tokenization of Molecules for Language Model Generation

Xiner Li^{*1} Limei Wang^{*1} Youzhi Luo¹ Carl Edwards² Shurui Gui¹ Yuchao Lin¹ Heng Ji² Shuiwang Ji¹

Abstract

We consider molecule generation in 3D space using language models (LMs), which requires discrete tokenization of 3D molecular geometries. Although tokenization of molecular graphs exists, that for 3D geometries is largely unexplored. Here, we attempt to bridge this gap by proposing the Geo2Seq, which converts molecular geometries into $SE(3)$ -invariant 1D discrete sequences. Geo2Seq consists of canonical labeling and invariant spherical representation steps, which together maintain geometric and atomic fidelity in a format conducive to LMs. Our experiments show that, when coupled with Geo2Seq, various LMs excel in molecular geometry generation, especially in controlled generation tasks. Our code has been released as part of the AIRS library (<https://github.com/divelab/AIRS/>).

1. Introduction

The generation of novel molecules with desired properties is an important step in drug discovery. Specifically, the design of three-dimensional (3D) molecular geometries is particularly important because 3D information plays a critical role in determining many molecular properties. Different generative models have been used for 3D molecule generation. Early studies such as G-SchNet (Gebauer et al., 2019) use autoregressive generative models to generate 3D molecules by sequentially placing atoms in 3D space. It was observed that these models often yield results with low chemical validity. Recently, diffusion models (Hoogeboom et al., 2022; Xu et al., 2023a) achieve better performance in 3D molecule generation tasks. However, they typically need thousands of diffusion steps, resulting in long generation time.

Language models (LMs) (Vaswani et al., 2017; Devlin et al., 2018; Brown et al., 2020; Gu et al., 2021), with their stream-

lined data processing and powerful generation capabilities, have shown success across various domains, particularly in natural language processing (NLP). Recently, large language models (LLMs) (Zhao et al., 2023b) show extraordinary capabilities in learning complex patterns (Zhang et al., 2024) and generating meaningful outputs (Touvron et al., 2023; Achiam et al., 2023; Chowdhery et al., 2023). Despite their potential, the application of LLMs to the direct generation of 3D molecules is largely under-explored. This is primarily due to the fact that geometric graph structures of molecular data are fundamentally different from texts. However, 3D geometric information is crucial in molecular tasks, since different conformations of the same molecule topology have different properties, such as per-atom forces. This gap reveals a unique challenge of how to make use of the powerful pattern recognition and generative capabilities of LLMs to handle complicated molecular graph structures, especially geometries. On the other hand, solutions to this challenge with model-level modifications cannot effectively leverage the rapidly developing power of LMs. These solutions require specific module designs, which needs to be done separately for each LM architecture and can be infeasible for modern LMs released via APIs.

In this work, we bridge this gap by applying LMs to the task of 3D molecule generation. We employ a novel approach translating the intricate geometry of molecules into a format that can be effectively processed by LMs. This is achieved by our proposed tokenization method Geo2Seq, which converts 3D molecular structures into $SE(3)$ -invariant one-dimensional (1D) discrete sequences. The transformation is based on canonical labeling, which allows dimension reduction with no information loss outside graph isomorphism groups, and invariant spherical representations, which guarantees $SE(3)$ -invariance under the equivariant global frame. By doing so, we harness the advanced sequence-processing capabilities and efficiency of LMs while retaining essential geometric and atomic information. Note that since Geo2Seq operates solely on input data, our method is agnostic to the subsequent LMs used. and can seamlessly adapt to any state-of-the-art sequence model, maximizing LM capabilities while avoiding additional architecture design or redundant computations. When combined with powerful modern LLMs, Geo2Seq can achieve highly accurate modeling of 3D molecular structures. In addition, Geo2Seq can benefit

^{*}Equal contribution ¹Texas A&M University ²University of Illinois Urbana-Champaign. Correspondence to: Shuiwang Ji <sjj@tamu.edu>.

conditional generation by including real-world chemical properties in sequences because modern LLMs are capable of capturing long-context correlations to comprehend global structure and information in sequences. Our experimental results demonstrate these advantages. We show that using different LMs with Geo2Seq can reliably produce valid and diverse 3D molecules and outperform the strong diffusion-based baselines by a large margin in conditional generation. These results validate the feasibility of using LMs for 3D molecule generation and highlight the potential to aid in the discovery of new molecules, paving the way for applications such as drug development and material science.

2. Preliminaries and Related Work

2.1. 3D Molecule Generation

In this work, we study the problem of generating 3D molecules from scratch. Note that this problem is different from the 3D molecular conformation generation problem studied in the literature (Mansimov et al., 2019; Simm & Hernandez-Lobato, 2020; Gogineni et al., 2020; Xu et al., 2021a;b; Shi et al., 2021; Ganea et al., 2021; Xu et al., 2022; Jing et al., 2022), where 3D molecular conformations are generated from 2D molecular graphs. We represent a 3D molecule with n atoms in the form of a 3D point cloud (*i.e.*, a set of points with different positions in 3D Euclidean space) as $G = (z, \mathbf{R})$. Here, $z = [z_1, \dots, z_n] \in \mathbb{Z}^n$ is the atom type vector where z_i is the atomic number (nuclear charge number) of the i -th atom, and $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_n] \in \mathbb{R}^{3 \times n}$ is the atom coordinate matrix, where \mathbf{r}_i is the 3D coordinate of the i -th atom. Note that 3D atom coordinates \mathbf{R} are commonly called 3D molecular conformations or geometries in chemistry. We aim to solve the following two generation tasks in this work:

- **Random generation.** Given a 3D molecule dataset $\mathcal{G} = \{G_j\}_{j=1}^m$, we aim to learn an unconditional generative model $p_\theta(\cdot)$ on \mathcal{G} so that the model can generate valid and diverse 3D molecules.
- **Controllable generation.** Given a 3D molecule dataset $\mathcal{G} = \{(G_j, s_j)\}_{j=1}^m$ where s_j is a certain property value of G_j , we aim to learn a conditional generative model $p_\theta(\cdot|s)$ on \mathcal{G} so that for a given s , the model can generate 3D molecules whose quantum property values are s .

A major technical challenge of 3D molecule generation lies in maintaining invariant to $SE(3)$ transformations, including rotation and translation. In other words, ideal models should assign the same probability to $G = (z, \mathbf{R})$ and $G' = (z, \mathbf{R}')$ if $\mathbf{R}' = \mathbf{Q}\mathbf{R} + \mathbf{b}\mathbf{1}^T$, where $\mathbf{1}$ is an n -dimensional vector whose elements are all one, $\mathbf{b} \in \mathbb{R}^3$ is an arbitrary translation vector, and $\mathbf{Q} \in \mathbb{R}^{3 \times 3}$ is a rotation matrix satisfying $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}, |\mathbf{Q}| = 1$. To achieve

$SE(3)$ -invariance in 3D molecule generation, existing studies have proposed various strategies. Early studies propose to generate 3D atom positions by $SE(3)$ -invariant features, such as interatomic distances, angles and torsion angles. They construct 3D molecular structures through either atom-by-atom generation (Gebauer et al., 2019; Luo & Ji, 2022) or generating full distance matrices (Hoffmann & Noé, 2019) in one shot. Recently, more and more studies have applied generative models to generate 3D atom coordinate directly. These studies include E-NFs (Satorras et al., 2021a) and EDM (Hooeboom et al., 2022), which combine equivariant atom coordinate alignment process with equivariant EGNN (Satorras et al., 2021b) model for 3D molecule generation. Following EDM, many other studies have proposed to improve diffusion-based 3D molecule generation frameworks by stochastic differential equation (SDE) based diffusion models (Wu et al., 2022; Bao et al., 2023) or latent diffusion models (Xu et al., 2023a). Besides, some recent studies (Qiang et al., 2023) have explored generating 3D molecules through generating and connecting fragments first, then aligning atom coordinates with software like RDKit. We refer readers to Du et al. (2022); Zhang et al. (2023b) for a comprehensive review.

While generating 3D molecules in the form of 3D point clouds have been well studied, few studies have tried applying powerful language models to this problem. In this work, different from mainstream methods, we convert 3D point clouds to $SE(3)$ -invariant 1D discrete sequences, and show that generating sequences by LMs achieves promising performance in the 3D molecule generation task.

2.2. Chemical Language Model

LMs have catalyzed significant advancements across a spectrum of fields. Recently, LLMs have revolutionized the landscape of NLP and beyond (Touvron et al., 2023; Achiam et al., 2023; Chowdhery et al., 2023). Drawing inspiration from NLP methodologies, chemical language models (CLMs) have emerged as a competent way for representing molecules (Bran & Schwaller, 2023; Janakarajan et al., 2023; Bajorath, 2024; Zhang et al., 2024). Due to the superiority LMs show in generation tasks, most CLMs are designed as generative models. Variants of LMs have been adapted for molecular science, producing a variety of works.

CLMs learn the chemical vocabulary and syntax used to represent molecules, as well as the conditional probabilities of character occurrence at given positions of sequences depending on preceding characters. This vocabulary covers all characters from the adopted molecule representation. All inputs including chemical structures and properties should be converted into sequence form and tokenized for compatibility with language models. Commonly, SMILES (Weininger, 1988) is used for this sequential representation, although

other formats like SELFIES (Krenn et al., 2019), atom type strings, and custom strings with positional or property values are also viable options. To learn representations, CLMs are usually pre-trained on extensive molecular sequences through self-supervised learning. Subsequently, models are fine-tuned on more focused datasets with desired properties, such as activity against a target protein. Generative CLMs generally adopt an autoregressive training approach of next token prediction, *i.e.*, iteratively predicting each subsequent token in a sequence based on the preceding tokens. Traditional autoregressive models use the Transformer architecture with causal self-attention (Brown et al., 2020) due to its superior efficacy, while other sequence models like recurrent neural networks (RNNs) and state space models (SSMs) (Gu et al., 2021; Özçelik et al., 2024; 2023) also show considerable functionality.

Given a dataset of sequences, $\mathbf{U} = \{U_1, U_2, \dots, U_N\}$, where U_i is transformed from the representation, property conditions and/or descriptions of a molecule G_i with n_i nodes, let $U_i = \{u_1, u_2, \dots, u_{n_i}\}$ and all tokens u_i belong to vocabulary V . An autoregressive CLM has parameters θ encoding a distribution with conditional probabilities of each token given its predecessors, $p(U_i; \theta) = \prod_{j=1}^{n_i} p(u_j | u_0 : u_{j-1}; \theta)$. The optimization process involves maximizing the probabilities of the entire dataset $p(\mathbf{U}; \theta) = \prod_{i=1}^N p(U_i; \theta)$. Each conditional distribution $p(u_j | u_0 : u_{j-1}; \theta)$ is a categorical distribution over the vocabulary size $|V|$; thus the loss for each term aligns with the standard cross-entropy loss. To generate new sequences, the model samples each token sequentially from these conditional distributions. To introduce randomness and control into generation, the sampling process is typically modulated with Top-K (k) and temperature (τ) hyperparameters, enabling a balance between adherence and diversity.

Most existing CLM works consider chemical structures as well as other modalities such as natural language captions (Bagal et al., 2021; Li et al., 2023a;b; Edwards et al., 2022; Xie et al., 2023; Chen et al., 2023b; Tysinger et al., 2023; Xu et al., 2023b; Chen et al., 2023a; Pei et al., 2023; Liu et al., 2023b; Wang et al., 2023), while some focus on pure text of chemical literature (Luo et al., 2022a) or molecule strings (Haroon et al., 2023; Mao et al., 2023b; Blanchard et al., 2023; Mazuz et al., 2023; Fang et al., 2023; Kyro et al., 2023; Izdebski et al., 2023; Yoshikai et al., 2023; Wu et al., 2023; Mao et al., 2023a). Notably, all these works solely consider 2D molecules for representation learning and downstream tasks, overlooking 3D geometric structures which is crucial in many molecular predictive and generative tasks. For example, different conformations of the same 2D molecule have different potentials and per-atom forces. In order to use pivotal 3D information, another line of work incorporate geometric models such as GNNs in parallel with the CLM (Xia et al., 2023; Zhang et al., 2023a; Cao

et al., 2023; Liang et al., 2023; Liu et al., 2023a; Frey et al., 2023), which requires additional design and training techniques to mitigate alignment issues. Some works extend the architecture of CLM to include 3D-geometric-model-like modules in the attention block (Fuchs et al., 2020; Shi et al., 2022; Liao & Smidt, 2022; Thölke & De Fabritiis, 2021; Luo et al., 2022b; Masters et al., 2022; Ünlü et al., 2023; Zhao et al., 2023a), capturing 3D information as positional encodings with considerable computations and framework design. In contrast, Flam-Shepherd & Aspuru-Guzik (2023) make an initial attempt showing language models trained directly on contents of XYZ format chemical files can generate molecules with three coordinates, implying pure LMs’ potential to directly explore 3D chemical space. In this work, we propose an invariant 3D molecular sequencing algorithm, Geo2Seq, to empower CLMs with structural completeness and geometric invariance, showing LMs’ capabilities of understanding molecules precisely in 3D space. We extend beyond the conventional Transformer architecture of CLMs and additionally employ SSMs as LM backbones. Furthermore, Geo2Seq operates solely on the input data, which allows independence from model architecture and training techniques and provides reuse flexibility.

3. Tokenization of 3D Molecules

A fundamental difference between LMs and other models is that LMs use discrete inputs, *i.e.*, tokens. In this section, we introduce our tokenization method to map input 3D molecules with atomic coordinates to discrete token sequences appropriate for LM learning.

A main challenge in tokenization design is to develop bijective mappings between 3D molecules and token sequences, *i.e.*, obtaining the same token sequence for the same input 3D molecule, while obtaining different sequences for different inputs. In this section, we present our solutions to tackle this challenge. We first reorder the atoms in the input molecule to a canonical order (Section 3.1), such that any two isomorphic graphs result in the same canonical form, and any non-isomorphic graphs yield different canonical forms. We then convert 3D Cartesian coordinates to $SE(3)$ -invariant spherical representations, including distances and angles (Section 3.2). Combining them together, we obtain our geometry informed tokenization method Geo2Seq (Section 3.3). We provide rigorous proof of all theorems supporting the bijective mapping relation in Appendix B.

3.1. Serialization via Canonical Ordering

As the first step in 3D molecule tokenization, we need to transform a graph to a 1D sequential representation. We resort to canonical labeling as a solution for dimension reduction without information loss.

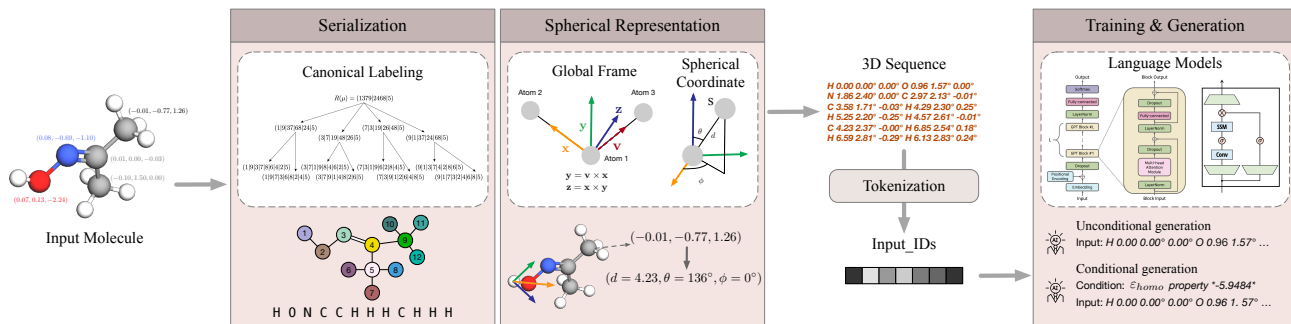


Figure 1: Overview of Geo2Seq. We use the canonical labeling order to arrange nodes in a row, fill in the place of each node with vector $[z_i, d_i, \theta_i, \phi_i]$, and concatenate all elements into a sequence. Each node vector contains atom type and spherical coordinates. Notably, the spherical coordinates are $SE(3)$ -invariant.

Canonical labeling (CL), in the context of graph theory, is a process to assign a unique form to each graph in a way that two graphs receive the same canonical form only if they are isomorphic (McKay et al., 1981). The canonical form is a re-indexed version of a graph, which is unique for the whole isomorphism class of a graph. The new indexes naturally establish the order of nodes in the graph. The order, which we refer to as canonical labels, is not necessarily unique if the graph has symmetries and thus has an automorphism group larger than 1. However, all canonical labels are strictly equivalent when used for serialization. The canonical label essentially re-assigns an index ℓ_i to each node originally indexed with i in graph G . Since canonical labeling can precisely distinguish non-isomorphic graphs, it fully contains the structure information of a graph G . Thus, by arranging nodes with attributes in the labeling order ℓ_1, ℓ_2, \dots , we obtain a sequential representation of attributed graphs with all structural information preserved.

The Nauty algorithm (McKay & Piperno, 2014), tailored for CL and computing graph automorphism groups, presents a rigorous formulation of CL. In this paper, we adopt the Nauty algorithm for CL calculation, while **all analyses and derivations apply to other rigorous algorithms**. The bijective mapping between CL-obtained sequential representation and graph can be proved based on graph isomorphism. First, due to the geometric need here, we extend to define the isomorphism problem for attributed graphs.

Definition 3.1. [Graph Isomorphism] Let $G_1 = (V_1, E_1, A_1)$ and $G_2 = (V_2, E_2, A_2)$ be two graphs, where V_i denotes the set of vertices, E_i denotes the set of edges, and A_i denotes the node attributes of G_i for $i = 1, 2$. Let $\text{attr}(v)$ denote the node attributes of vertex v . The graphs G_1 and G_2 are said to be isomorphic, denoted as $G_1 \cong G_2$, if there exists a bijection $b : V_1 \rightarrow V_2$ such that for every vertex $v \in V_1$, $\text{attr}(v) \in A_1 = \text{attr}(b(v)) \in A_2$, and for every pair of vertices $u, v \in V_1$,

$$(u, v) \in E_1 \Leftrightarrow (b(u), b(v)) \in E_2.$$

CL processes can also be extended to node/edge-attributed graphs, leading us to the guarantee below.

Lemma 3.2. [Canonical Labeling for Colored Graph Isomorphism] Let $G_1 = (V_1, E_1, A_1)$ and $G_2 = (V_2, E_2, A_2)$ be two finite, undirected graphs where V_i denotes the set of vertices, E_i denotes the set of edges, and A_i denotes the node attributes of the graph G_i for $i = 1, 2$. Let $L : \mathcal{G} \rightarrow \mathcal{L}$ be a function that maps a graph $G \in \mathcal{G}$, the set of all finite, undirected graphs, to its canonical label $L(G) \in \mathcal{L}$, the set of all possible canonical labels, as produced by the Nauty algorithm. Then the following equivalence holds:

$$L(G_1) = L(G_2) \Leftrightarrow G_1 \cong G_2$$

where $G_1 \cong G_2$ denotes that G_1 and G_2 are isomorphic.

Lemma 3.2 indicates that the CL process is both complete (sufficient to distinguish non-isomorphic graphs) and sound (not distinguishing actually isomorphic graphs). Note that if $L(G)$ corresponds to multiple automorphic labels, we can randomly select one since they are all equivalent and produce the same sequence later through Geo2Seq, as detailed in Appendix B. However, this is a very uncommon case for real-world 3D attributed graphs like molecules.

3.2. Invariant Spherical Representations

In this section, we describe how to incorporate 3D structure information into our sequences. One main challenge here is to ensure the $SE(3)$ -invariance property described in Section 2.1. Specifically, given a 3D molecule, if it is rotated or translated in the 3D space, its 3D representation should be unchanged. Another challenge is to ensure no information loss (Liu et al., 2022; Wang et al., 2022). Specifically, given the 3D representation, we can recover the given 3D structure. If two 3D structures cannot be matched via a $SE(3)$ transformation, the representations should be different. This property is important to the discriminative ability of models.

We address these challenges by **spherical representations**, i.e., using spherical coordinates to represent 3D structures.

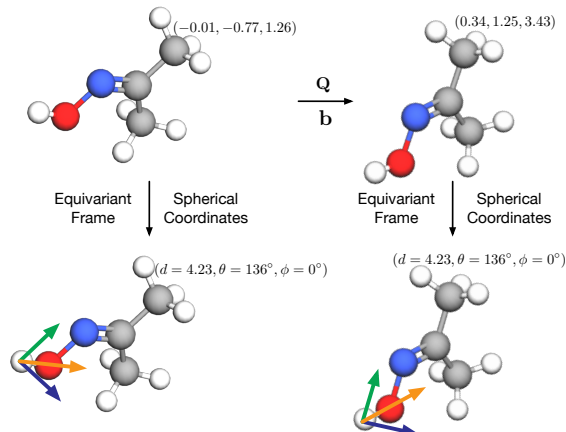


Figure 2: Illustrations of the equivariant frame and invariant spherical representations. If the molecule is rotated and translated by a rotation matrix \mathbf{Q} and a translation vector \mathbf{b} , the atom coordinates change accordingly. But our spherical representations remain invariant since the frame is equivariant to the $SE(3)$ -transformation.

Compared to Cartesian coordinates, spherical coordinate values are bounded in a smaller region, namely, a range of $[0, \pi]$ or $[0, 2\pi]$. This makes spherical coordinates advantageous in discretized representations and thus easier to be modeled by LMs. Given the same decimal place constraints, they require a smaller vocabulary size, and given the same vocabulary size, they lose less information. This is also supported by empirical results and analysis in Appendix C.

We propose to maintain $SE(3)$ -invariance while ensuring no information loss. Given a 3D molecule G with atom types \mathbf{z} and atom coordinates \mathbf{R} , we first build a **global coordinate frame** $\mathbf{F} = (\mathbf{x}, \mathbf{y}, \mathbf{z})$ based on the input. Specifically, as shown in Figure 1, the frame is built based on the first three non-collinear atoms in the canonical ordering $\mathbf{L}(G)$. Let ℓ_1, ℓ_2 , and ℓ_F be the indices of these three atoms. Then the global frame $\mathbf{F} = (\mathbf{x}, \mathbf{y}, \mathbf{z})$ is calculated as

$$\begin{aligned} \mathbf{x} &= \text{normalize}(\mathbf{r}_{\ell_2} - \mathbf{r}_{\ell_1}), \\ \mathbf{y} &= \text{normalize}((\mathbf{r}_{\ell_F} - \mathbf{r}_{\ell_1}) \times \mathbf{x}), \\ \mathbf{z} &= \mathbf{x} \times \mathbf{y}. \end{aligned} \quad (1)$$

Here $\text{normalize}(\cdot)$ is the function to normalize a vector to unit length. Note that the global frame is equivariant to the rotation and translation of the input molecule, as shown in Figure 2 and Appendix B.2. After obtaining the global frame, we use a function $f(\cdot)$ to convert the coordinates of each atom to **spherical coordinates** d, θ, ϕ under this frame. Specifically, for each node ℓ_i with coordinate \mathbf{r}_{ℓ_i} , the corresponding spherical coordinate is

$$\begin{aligned} d_{\ell_i} &= \|\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}\|_2, \\ \theta_{\ell_i} &= \arccos((\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{z} / d_{\ell_i}), \\ \phi_{\ell_i} &= \text{atan2}((\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{y}, (\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{x}). \end{aligned} \quad (2)$$

The spherical coordinates show the relative position of each

atom in the global frame \mathbf{F} . As shown in Figure 2, if the input coordinates are rotated by a matrix \mathbf{Q} and translated by a vector \mathbf{b} , the transformed spherical coordinates remain the same, so the spherical coordinates are $SE(3)$ -invariant.

Next, we demonstrate that there is no information loss in our method. We show that given our $SE(3)$ -invariant spherical representations, we can recover the given 3D structures. For each node ℓ_i , we convert the spherical coordinate $[d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}]$ to coordinate \mathbf{r}'_{ℓ_i} in 3D space as

$$[d_{\ell_i} \sin(\theta_{\ell_i}) \cos(\phi_{\ell_i}), d_{\ell_i} \sin(\theta_{\ell_i}) \sin(\phi_{\ell_i}), d_{\ell_i} \cos(\theta_{\ell_i})].$$

Note that our reconstructed coordinate \mathbf{r}'_{ℓ_i} may not be exactly the same as the original coordinate \mathbf{r}_{ℓ_i} . However, there exists a $SE(3)$ -transformation g , such that $g(\mathbf{r}'_{\ell_i}) = \mathbf{r}_{\ell_i}$ for all i . Note that the same transformation g is applied to all nodes. Formally, by applying the function $f(\cdot)$ to the 3D coordinate matrix \mathbf{R} , we can demonstrate the following properties of spherical representations.

Lemma 3.3. *Let $G = (\mathbf{z}, \mathbf{R})$ be a 3D graph with node type vector \mathbf{z} and node coordinate matrix \mathbf{R} . Let \mathbf{F} be the equivariant global frame of graph G built based on the first three non-collinear nodes in $\mathbf{L}(G)$. $f(\cdot)$ is our function that maps 3D coordinate matrix \mathbf{R} of G to spherical representations \mathbf{S} under the equivariant global frame \mathbf{F} . Then for any 3D transformation $g \in SE(3)$, we have $f(\mathbf{R}) = f(g(\mathbf{R}))$. Given spherical representations $\mathbf{S} = f(\mathbf{R})$, there exist a transformation $g \in SE(3)$, such that $f^{-1}(\mathbf{S}) = g(\mathbf{R})$.*

Lemma 3.3 indicates that our spherical representation is $SE(3)$ -invariant, and we can reconstruct (a transformation of) the original coordinates. Therefore, our method can convert 3D structures into $SE(3)$ -invariant representations with no information loss. Proofs are in Appendix B.

3.3. Geo2Seq: Geometry Informed Tokenization

In this section, we describe the process and properties of our 3D tokenization method, Geo2Seq. Equipped with canonical labeling that reduces graph structures to 1D sequences with no information loss regarding graph isomorphism, and $SE(3)$ -invariant spherical representations that ensure no 3D information loss, we develop Geo2Seq, a reversible transformation from 3D molecules to 1D sequences. Figure 1 shows an overview of Geo2Seq. Specifically, given a graph G with n nodes, Geo2Seq concatenates the node vector $[z_i, d_i, \theta_i, \phi_i]$ of every node in G to a 1D sequence by its canonical order, ℓ_1, \dots, ℓ_n . To formulate the properties of Geo2Seq, we extend the concept of graph isomorphism in Definition B.1 to 3D graphs.

Definition 3.4. [3D Graph Isomorphism] Let $G_1 = (\mathbf{z}_1, \mathbf{R}_1)$ and $G_2 = (\mathbf{z}_2, \mathbf{R}_2)$ be two 3D graphs, where \mathbf{z}_i is the node type vector and \mathbf{R}_i is the node coordinate matrix of the molecule G_i . Let V_i denote the set of vertices, A_i denote node attributes, and no edge exists. Two 3D graphs

G_1 and G_2 are **3D isomorphic**, denoted as $G_1 \cong_{3D} G_2$, if there exists a bijection $b : V_1 \rightarrow V_2$ such that $G_1 \cong G_2$ given $A_i = [z_i, \mathbf{R}_i]$, and there exists a 3D transformation $g \in SE(3)$ such that $\mathbf{r}_i^{G_1} = g(\mathbf{r}_{b(i)}^{G_2})$. If a small error ϵ is allowed such that $|\mathbf{r}_i^{G_1} - g(\mathbf{r}_{b(i)}^{G_2})| \leq \epsilon$, we call the two 3D graphs **ϵ -constrained 3D isomorphic**.

Considering Lemma 3.2, we specify $G = (V, E, A)$ with $A = [z, \mathbf{R}]$ and define the CL function for 3D molecules as \mathbf{L}_m , which extends the equivalence of Lemma 3.2 to \mathbf{L}_m with 3D isomorphism. We formulate Geo2Seq and our major theoretical derivations below.

Theorem 3.5. [Bijjective Mapping between 3D Graph and Sequence] Following Definition 3.4, let $G_1 = (z_1, \mathbf{R}_1)$ and $G_2 = (z_2, \mathbf{R}_2)$ be two 3D graphs. Let $\mathbf{L}_m(G)$ be the canonical label for 3D graph G and $f : \mathcal{R} \rightarrow \mathcal{S}$ be the function that maps 3D coordinates to its spherical representations. Given a graph G with n nodes and $\mathbf{X} = [x_1, \dots, x_n]^T \in \mathbb{R}^{n \times m}$, where $m \in \mathbb{Z}$, we define $\mathbf{L}_m(G) \otimes \mathbf{X} = \text{concat}(x_{\ell_1}, \dots, x_{\ell_n})$, where ℓ_i is the index of the node labeled i by $\mathbf{L}_m(G)$, and $\text{concat}(\cdot)$ concatenates elements as a sequence. We define

$$\text{Geo2Seq}(G) = \mathbf{L}_m(G) \otimes (z, f(\mathbf{R})) = \mathbf{L}_m(G) \otimes \mathbf{X},$$

where $x_i = [z_i, d_i, \theta_i, \phi_i]$. Then $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is a surjective function, and the following equivalence holds:

$$\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2) \Leftrightarrow G_1 \cong_{3D} G_2,$$

where $G_1 \cong_{3D} G_2$ denotes G_1 and G_2 are 3D isomorphic.

Theorem 3.5 establishes the following guarantees for Geo2Seq: (1) Given a 3D molecule, we can uniquely construct a 1D sequence using Geo2Seq. (2) If two molecules are 3D isomorphic, their sequence outputs from Geo2Seq are identical. (3) Given a sequence output of Geo2Seq, we can uniquely reconstruct a 3D molecule. (4) If two constructed sequences from Geo2Seq are identical, their corresponding molecules must be 3D isomorphic. This enable sequential tokenization of 3D molecules, preserving structural completeness and geometric invariance.

Due to the necessity of discreteness in serialization and tokenization for LMs, in reality, numerical values need to be discretized before concatenation. In practice, we round up numerical values to certain decimal places. Thus Theorem 3.5 can be extended with constraints, as below.

Corollary 3.6. [Constrained Bijjective Mapping between 3D Graph and Sequence] Following the notations and definitions of Theorem 3.5, let spherical coordinate values be rounded up to b decimal places. Then $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is a surjective function, and the following equivalence holds:

$$\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2) \Leftrightarrow G_1 \cong_{3D-|10^{-b}|/2} G_2,$$

where $G_1 \cong_{3D-|10^{-b}|/2} G_2$ denotes graphs G_1 and G_2 are $(|10^{-b}|/2)$ -constrained 3D isomorphic.

Corollary 3.6 extends Theorem 3.5’s guarantees for the practical use of Geo2Seq. If we allow a round-up error below $|10^{-b}|/2$ for coordinates when distinguishing 3D isomorphism, all properties still hold. This implies that the practical Geo2Seq implementation retains near-complete geometric information and invariance, with numerical precision of $\epsilon \leq |10^{-b}|/2$.

With discreteness incorporated, we can collect a finite vocabulary covering all accessible molecule samples to enable tokenization for LMs. Specifically, we use vocabularies of approximately 1K-16K tokens consisting of atom type tokens ‘C, N, O...’, and spherical coordinate tokens such as ‘−1.98’, ‘1.57°’ or ‘−0.032°’. Specifically, the vocabulary size is approximately 1.8K for the QM9 dataset, and 16K for the Geom-Drug dataset. Note that we consider chirality for atoms and use the special token suffixes ‘@’ and ‘@@’ to distinguish clockwise and counterclockwise chiral centers, for example, ‘C@’ and ‘C@@’. The numerical tokens range from the smallest to the largest distance and angle values with restricted precision of 2 or 3 decimal places. Experimental results show the benefits in using this level of tokenization, as detailed in Appendix C.

4. 3D Molecule Generation

Training and Sampling. Now that we have defined a canonical and robust sequence representation for 3D molecules, we turn to the method of modeling such sequences, U . Here, we attempt to train a model M with parameters θ to capture the distribution of such sequences, $p_\theta(U)$, in our dataset. As this is a well-studied problem within language modeling, we opt to use two language models, GPT (Radford et al., 2018) and Mamba (Gu & Dao, 2023), which have shown effective sequence modeling capabilities on a range of tasks. Both models are trained using a standard next-token prediction cross-entropy loss ℓ for all elements in the sequence:

$$\min_{\theta} \mathbb{E}_{u \in U} \left[\sum_{i=1}^{|u|-1} \ell(M_\theta(u_1, \dots, u_i), u_{i+1}) \right].$$

To sample from a trained model, we first select an initial atom token by sampling from the multinomial distribution of first-tokens in the training data (we note that in almost all cases this is ‘H’). We then perform a standard autoregressive sampling procedure by iteratively sampling from the conditional distribution $p_\theta(u_{i+1}|u_1, \dots, u_i)$ until the stop token or max length is reached. We sample from this distribution using top- k sampling (Fan et al., 2018) and a softmax temperature τ (Ackley et al., 1985; Fidler & Goldberg, 2017). Unless otherwise noted, $\tau = 0.7$ and $k = 80$.

Controllable Generation. For controllable generation, we follow Bagal et al. (2021) and use a conditioning token for the desired property. This token is created by projecting the

desired properties through a trainable linear layer to create a vector with the model’s initial token embedding space. This property token is then used as the initial element in the molecular sequence. Training and sampling are performed as before with this new sequence formulation. Sampling begins with the desired property’s token as input.

5. Experimental Studies

In this section, we evaluate the method of generating 3D molecules in the form of our proposed Geo2Seq representations by LLMs. We show that in the random generation task (see Section 2.1), the performance of Geo2Seq with GPT (Radford et al., 2018) or Mamba (Gu & Dao, 2023) models is better than or comparable with state-of-the-art 3D point cloud based methods, including EDM (Hooeboom et al., 2022) and GEOLDM (Xu et al., 2023a). In addition, in the controllable generation task (see Section 2.1), we show that Geo2Seq with Mamba models outperform previous 3D point cloud based methods by a large margin.

5.1. Random Generation

Data. We adopt two datasets, QM9 (Ramakrishnan et al., 2014) and GEOM-DRUGS (Axelrod & Gomez-Bombarelli, 2022), to evaluate performances in the random generation task. The QM9 dataset collects over 130k 3D molecules with 3D structures calculated by density functional theory (DFT). Each molecule in QM9 has less than 9 heavy atoms and its chemical elements all belong to H, C, N, O, F. Following Anderson et al. (2019), we split the dataset into train, validation and test sets with 100k, 18k and 12k samples, separately. The GEOM-DRUGS dataset consists of over 450k large molecules with 37 million DFT-calculated 3D structures. Molecules in GEOM-DRUGS has up to 181 atoms and 44.2 atoms on average. We follow Hooeboom et al. (2022) to select 30 3D structures with the lowest energies per molecule for model training.

Setup. On the QM9 dataset, we set the training batch size to 32, base learning rate to 0.0004, and train a 12-layer GPT model and a 26-layer Mamba model by AdamW (Loshchilov & Hutter, 2019) optimizers. On the GEOM-DRUGS dataset, we set the training batch size to 32, base learning rate to 0.0004, and train a 14-layer GPT model and a 28-layer Mamba model by AdamW optimizers. See Appendix D for more information about hyperparameters and other settings. When model training is completed, we randomly generate 10,000 molecules, and evaluate the performance on these molecules. Specifically, we first transform 3D molecular structures to 2D molecular graphs using the bond inference algorithm implemented in the official code of EDM. Then, we evaluate the performance by **atom stability**, which is the percentage of atoms with correct bond valencies, and **molecule stability**, which is the percentage of molecules whose all atoms have correct bond valencies.

In addition, we report the percentage of **valid** molecules that can be successfully converted to SMILES strings by RDKit, and the percentage of **valid and unique** molecules that can be converted to unique SMILES strings.

Baselines. We compare GPT and Mamba models with several strong baseline methods. Specifically, we compare with an autoregressive generation method G-SchNet (Gebauer et al., 2019) and an equivariant flow model based method E-NFs (Satorras et al., 2021a). We also compare with some recently proposed diffusion based methods, including EDM (Hooeboom et al., 2022), GDM (the non-equivariant variant of EDM) and GDM-AUG (GDM trained with random rotation as data augmentation). Besides, we compare with EDM-Bridge (Wu et al., 2022) and GEOLDM (Xu et al., 2023a), which are two latest 3D molecule generation methods improving EDM by SDE based diffusion models and latent diffusion models, respectively. To ensure that the comparison is fair, our methods and baseline methods use the same data split and evaluation metrics.

Results. We present the random generation results of different methods on QM9 and GEOM-DRUGS datasets in Table 1. Note that for GEOM-DRUGS dataset, all methods achieve nearly 0% molecule stability percentage and 100% uniqueness percentage. Thus, following previous studies, these two metrics are omitted. According to the results in Table 1, on QM9 dataset, generating 3D molecules in Geo2Seq representations with either GPT or Mamba models achieve better performance than all 3D point cloud based baseline methods in molecule stability and valid percentage, and achieves atom stability percentages close to the upper bound (99%). This demonstrates that our method can model 3D molecular structure distribution and capture the underlying chemical rules more accurately. It is worth noticing that our method does not achieve very high uniqueness percentage, showing that it is not easy for our method to generate a large number of diverse molecules. We believe this is due to that the conversion from real numbers to discrete tokens limits the search space of 3D molecular structures, especially on a small dataset like QM9, while it is easier to generate more diverse molecules for 3D point cloud based methods as they directly generate real numbers. This is reflected by the fact that our method achieves nearly 100% uniqueness percentage on the large GEOM-DRUGS dataset. On GEOM-DRUGS dataset, both GPT and Mamba models achieve reasonably high atom stability and valid percentage. The performance of our method is comparable with strong diffusion based baseline methods, showing that LLMs have the potential to model very complicated drug molecular structures well. We will explore further improving the performance on GEOM-DRUGS dataset with larger LLMs in the future.

See Appendix D.3 for additional experiments and met-

Table 1: Random generation performance on QM9 and GEOM-DRUGS datasets with 3 runs. Larger numbers indicate better performance. **bold** and underline highlight the best and second best performance, respectively. For GEOM-DRUGS dataset, molecule stability and unique percentage are close to 0% and 100% for all methods so they are not presented.

Method	QM9				GEOM-DRUGS	
	Atom Sta (%)	Mol Sta (%)	Valid (%)	Valid & Unique (%)	Atom Sta (%)	Valid (%)
Data	99.0	95.2	97.7	97.7	86.5	99.9
E-NFs	85.0	4.9	40.2	39.4	-	-
G-SchNet	95.7	68.1	85.5	80.3	-	-
GDM	97.0	63.2	-	-	75.0	90.8
GDM-AUG	97.6	71.6	90.4	89.5	77.7	91.8
EDM	98.7 \pm 0.1	82.0 \pm 0.4	91.9 \pm 0.5	90.7 \pm 0.6	81.3	92.6
EDM-Bridge	98.8 \pm 0.1	84.6 \pm 0.3	92.0 \pm 0.1	90.7 \pm 0.1	82.4	92.8
GEOLDM	98.9 \pm 0.1	89.4 \pm 0.5	93.8 \pm 0.4	92.7 \pm 0.5	84.4	99.3
Geo2Seq with GPT	98.3 \pm 0.1	90.3 \pm 0.1	94.8 \pm 0.2	80.6 \pm 0.4	82.6	87.4
Geo2Seq with Mamba	98.9 \pm 0.2	93.2 \pm 0.2	97.1 \pm 0.2	81.7 \pm 0.4	82.5	96.1

Table 2: Controllable generation performance on QM9 datasets. Smaller numbers indicate better performance.

Property (Units)	α (Bohr ³)	$\Delta\epsilon$ (meV)	ϵ_{HOMO} (meV)	ϵ_{LUMO} (meV)	μ (D)	C_v ($\frac{\text{cal}}{\text{mol K}}$)
Data	0.10	64	39	36	0.043	0.040
Random	9.01	1470	645	1457	1.616	6.857
N_{atoms}	3.86	866	426	813	1.053	1.971
EDM	2.76	655	356	584	1.111	1.101
GEOLDM	2.37	587	340	522	1.108	1.025
Geo2Seq with Mamba	0.46	98	<u>57</u>	<u>71</u>	<u>0.164</u>	0.275
Geo2Seq with GPT	<u>0.53</u>	<u>102</u>	48	53	0.097	<u>0.325</u>

rics (Huang et al., 2023; Vignac et al., 2023), Appendix C for ablation studies, Appendix D for complexity analysis, and Appendix F for visualization.

5.2. Controllable Generation

Data. In the controllable generation task, we train our models on molecules and their property labels in the QM9 (Ramakrishnan et al., 2014) dataset. Specifically, we try taking a certain quantum property value as the conditional input to LLMs, and train LLMs to generate molecules with the conditioned quantum property values. Following Hoogetboom et al. (2022), we split the training dataset of QM9 to two subsets where each subset has 50k samples, and train our conditional generation models and an EGNN (Satorras et al., 2021b) based quantum property prediction models on these two subsets, respectively. We conduct the controllable generation experiments on six quantum properties from QM9, including polarizability (α), HOMO energy (ϵ_{HOMO}), LUMO energy (ϵ_{LUMO}), HOMO-LUMO gap ($\Delta\epsilon$), dipole moment (μ) and heat capacity at 298.15K (C_v).

Setup. For the controllable generation experiment, we train 16-layer Mamba (Gu & Dao, 2023) models with the same hyperparameters as the random generation experiments in Section 5.1. To evaluate the performance, we sample 10000 quantum property values, generate molecules conditioned on these property values by trained models, and compute the mean absolute difference (**MAE**) between the given property values and the property values of the generated

molecules. Note that we use the trained EGNN based property prediction models to calculate the property values of the generated molecules.

Baselines. We compare our models with two equivariant diffusion models, EDM (Hoogetboom et al., 2022) and GEOLDM (Xu et al., 2023a). In addition, we use several baselines that are based on dataset molecules. One baseline (Data) is directly taking the molecules from the QM9 dataset and use their property values as conditions. The MAE metric simply reflects the prediction error of the trained property prediction model, which can be considered as a lower bound. The second baseline (Random) is taking the molecules from the dataset but uses the randomly shuffled property values as conditions, and its MAE can be considered as an upper bound. The third baseline (N_{atoms}) uses the molecules from the dataset but uses property values predicted from the number of atoms as conditions. Achieving better performance than this baseline shows that models can use conditional information beyond the number of atoms.

Results. Controllable generation results of different methods are summarized in Table 2. As shown in the table, among all six properties, our method outperforms the strong diffusion based baseline methods EDM and GEOLDM by a large margin. Our method moves a significant step in pushing the performance of controllable generation task towards the lower bound, *i.e.*, Data baseline. As we use the same training set as EDM and GEOLDM to train the conditional generation model, the good performance of our

method shows that LLMs have more powerful capacity in incorporating conditional information into the 3D molecular structure generation process. We believe that the powerful long-context correlation capturing structures from LLMs, *e.g.*, attention mechanism, play significant roles in achieving the good control of 3D molecule generation by the conditioned property values. The huge success of LLMs in controllable molecule generation will motivate broader applications of LLMs in goal-directed or constrained drug design. See Appendix F for visualization of molecules generated from some given polarizability values.

6. Conclusion and Discussion

Geo2Seq showcases the potential of pure LMs in revolutionizing molecular design and drug discovery when geometric information is properly transformed. The framework has certain limitations, particularly in the generalization abilities across the continuous domain of real numbers. Due to the discrete nature of vocabularies, LMs rely on large pre-training corpus, fine-grained tokenization or emergent abilities for better generalization, as a trade-off to high precision and versatility. Future works points towards several directions, such as expanding on conditional tasks and exploring advanced tokenization techniques.

Acknowledgments

This work was supported partially by National Science Foundation grant IIS-2006861 and National Institutes of Health grant U01AG070112 (to S.J.), and by the Molecule Maker Lab Institute (to H.J.): an AI research institute program supported by NSF under award No. 2019897. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

Impact Statement

This paper presents work whose goal is to advance the field of Deep Learning, particularly 3D molecule generation. While this research primarily contributes to technical advancements in generative modeling, it has potential implications in domains such as drug discovery. We acknowledge that generative models in science could be misused if not carefully applied. However, we emphasize the importance of responsible deployment and alignment with ethical guidelines in generative AI. Overall, our contributions align with the broader goal of machine learning methodologies, and we do not foresee any immediate ethical concerns beyond those generally associated with generative models.

References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Ackley, D. H., Hinton, G. E., and Sejnowski, T. J. A learning algorithm for boltzmann machines. *Cognitive science*, 9 (1):147–169, 1985.
- Anderson, B., Hy, T. S., and Kondor, R. Cormorant: Covariant molecular neural networks. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/03573b32b2746e6e8ca98b9123f2249b-Paper.pdf.
- Axelrod, S. and Gomez-Bombarelli, R. GEOM, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data*, 9(1):185, 2022.
- Bagal, V., Aggarwal, R., Vinod, P., and Priyakumar, U. D. MolGPT: molecular generation using a transformer-decoder model. *Journal of Chemical Information and Modeling*, 62(9):2064–2076, 2021.
- Bajorath, J. Chemical language models for molecular design. *Molecular Informatics*, 43(1):e202300288, 2024.
- Banck, M., Morley, C. A., Vandermeersch, T., and Hutchison, G. R. Open Babel: An open chemical toolbox, 2011. URL <https://doi.org/10.1186/1758-2946-3-33>.
- Bao, F., Zhao, M., Hao, Z., Li, P., Li, C., and Zhu, J. Equivariant energy-guided SDE for inverse molecular design. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=r0otLtOwYW>.
- Blanchard, A. E., Bhowmik, D., Fox, Z., Gounley, J., Glaser, J., Akpa, B. S., and Irle, S. Adaptive language model training for molecular design. *Journal of Cheminformatics*, 15(1):1–12, 2023.
- Bran, A. M. and Schwaller, P. Transformers and large language models for chemistry and drug discovery. *arXiv preprint arXiv:2310.06083*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.

- Buttenschoen, M., Morris, G. M., and Deane, C. M. PoseBusters: AI-based docking methods fail to generate physically valid poses or generalise to novel sequences, 2023.
- Cao, H., Liu, Z., Lu, X., Yao, Y., and Li, Y. InstructMol: Multi-modal integration for building a versatile and reliable molecular assistant in drug discovery. *arXiv preprint arXiv:2311.16208*, 2023.
- Chen, Y., Wang, Z., Zeng, X., Li, Y., Li, P., Ye, X., and Sakurai, T. Molecular language models: Rnns or transformer? *Briefings in Functional Genomics*, pp. elad012, 2023a.
- Chen, Y., Xi, N., Du, Y., Wang, H., Jianyu, C., Zhao, S., and Qin, B. From artificially real to real: Leveraging pseudo data from large language models for low-resource molecule discovery. *arXiv preprint arXiv:2309.05203*, 2023b.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. PaLM: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Cuthill, E. and McKee, J. Reducing the bandwidth of sparse symmetric matrices. In *Proceedings of the 1969 24th national conference*, pp. 157–172, 1969.
- Daigavane, A., Kim, S., Geiger, M., and Smidt, T. Symphony: Symmetry-equivariant point-centered spherical harmonics for molecule generation. *arXiv preprint arXiv:2311.16199*, 2023.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Dijkstra, E. W. A note on two problems in connexion with graphs. In *Edsger Wybe Dijkstra: His Life, Work, and Legacy*, pp. 287–290. 2022.
- Du, Y., Fu, T., Sun, J., and Liu, S. Molgensurvey: A systematic survey in machine learning models for molecule design. *arXiv preprint arXiv:2203.14500*, 2022.
- Edwards, C., Lai, T., Ros, K., Honke, G., Cho, K., and Ji, H. Translation between molecules and natural language. *arXiv preprint arXiv:2204.11817*, 2022.
- Fan, A., Lewis, M., and Dauphin, Y. Hierarchical neural story generation. *arXiv preprint arXiv:1805.04833*, 2018.
- Fang, Y., Zhang, N., Chen, Z., Fan, X., and Chen, H. Molecular language model as multi-task generator. *arXiv preprint arXiv:2301.11259*, 2023.
- Ficler, J. and Goldberg, Y. Controlling linguistic style aspects in neural language generation. *arXiv preprint arXiv:1707.02633*, 2017.
- Flam-Shepherd, D. and Aspuru-Guzik, A. Language models can generate molecules, materials, and protein binding sites directly in three dimensions as xyz, cif, and pdb files. *arXiv preprint arXiv:2305.05708*, 2023.
- Frey, N. C., Soklaski, R., Axelrod, S., Samsi, S., Gomez-Bombarelli, R., Coley, C. W., and Gadepally, V. Neural scaling of deep chemical models. *Nature Machine Intelligence*, 5(11):1297–1305, 2023.
- Fuchs, F., Worrall, D., Fischer, V., and Welling, M. Se(3)-transformers: 3d roto-translation equivariant attention networks. *Advances in neural information processing systems*, 33:1970–1981, 2020.
- Ganea, O.-E., Pattanaik, L., Coley, C. W., Barzilay, R., Jensen, K., Green, W., and Jaakkola, T. S. GeoMol: Torsional geometric generation of molecular 3D conformer ensembles. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, 2021. URL https://openreview.net/forum?id=af_hng9tuNj.
- Gebauer, N., Gastegger, M., and Schütt, K. Symmetry-adapted generation of 3D point sets for the targeted discovery of molecules. *Advances in neural information processing systems*, 32, 2019.
- Gogineni, T., Xu, Z., Punzalan, E., Jiang, R., Kammeraad, J., Tewari, A., and Zimmerman, P. TorsionNet: a reinforcement learning approach to sequential conformer search. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 20142–20153. Curran Associates, Inc., 2020.
- Gu, A. and Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- Gu, A., Goel, K., and Ré, C. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- Haroon, S., Hafsath, C., and Jereesh, A. Generative pre-trained transformer (GPT) based model with relative attention for de novo drug design. *Computational Biology and Chemistry*, 106:107911, 2023.
- Hilbert, D. and Hilbert, D. Über die stetige abbildung einer linie auf ein flächenstück. *Dritter Band: Analysis: Grundlagen der Mathematik- Physik Verschiedenes: Nebst Einer Lebensgeschichte*, pp. 1–2, 1935.

- Hoffmann, M. and Noé, F. Generating valid Euclidean distance matrices. *arXiv preprint arXiv:1910.03131*, 2019.
- Hoogeboom, E., Satorras, V. G., Vignac, C., and Welling, M. Equivariant diffusion for molecule generation in 3D. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 8867–8887. PMLR, 2022. URL <https://proceedings.mlr.press/v162/hoogeboom22a.html>.
- Huang, H., Sun, L., Du, B., and Lv, W. Learning joint 2d & 3d diffusion models for complete molecule generation. *arXiv preprint arXiv:2305.12347*, 2023.
- Izdebski, A., Weglarz-Tomczak, E., Szczurek, E., and Tomczak, J. M. De novo drug design with joint transformers. *arXiv preprint arXiv:2310.02066*, 2023.
- Janakarajan, N., Erdmann, T., Swaminathan, S., Laino, T., and Born, J. Language models in molecular discovery. *arXiv preprint arXiv:2309.16235*, 2023.
- Jing, B., Corso, G., Chang, J., Barzilay, R., and Jaakkola, T. S. Torsional diffusion for molecular conformer generation. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=w6fj2r62r_H.
- Kim, Y. and Kim, W. Y. Universal Structure Conversion Method for Organic Molecules: From Atomic Connectivity to Three-Dimensional Geometry. *Bulletin of the Korean Chemical Society*, 36(7):1769–1777, 2015. doi: <https://doi.org/10.1002/bkcs.10334>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/bkcs.10334>.
- Krenn, M., Häse, F., Nigam, A., Friederich, P., and Aspuru-Guzik, A. Selfies: a robust representation of semantically constrained graphs with an example application in chemistry. *arXiv preprint arXiv:1905.13741*, 1(3), 2019.
- Kyro, G. W., Morgunov, A., Brent, R. I., and Batista, V. S. Chemspaceal: An efficient active learning methodology applied to protein-specific molecular generation. *ArXiv*, 2023.
- 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.
- Lee, C. Y. An algorithm for path connections and its applications. *IRE transactions on electronic computers*, (3): 346–365, 1961.
- Li, J., Liu, Y., Fan, W., Wei, X.-Y., Liu, H., Tang, J., and Li, Q. Empowering molecule discovery for molecule-caption translation with large language models: A chatgpt perspective. *arXiv preprint arXiv:2306.06615*, 2023a.
- Li, Y., Gao, C., Song, X., Wang, X., Xu, Y., and Han, S. DrugGPT: A GPT-based strategy for designing potential ligands targeting specific proteins. *bioRxiv*, pp. 2023–06, 2023b.
- Liang, Y., Zhang, R., Zhang, L., and Xie, P. DrugChat: towards enabling ChatGPT-like capabilities on drug molecule graphs. *arXiv preprint arXiv:2309.03907*, 2023.
- Liao, Y.-L. and Smidt, T. Equiformer: Equivariant graph attention transformer for 3d atomistic graphs. *arXiv preprint arXiv:2206.11990*, 2022.
- Liu, S., Nie, W., Wang, C., Lu, J., Qiao, Z., Liu, L., Tang, J., Xiao, C., and Anandkumar, A. Multi-modal molecule structure–text model for text-based retrieval and editing. *Nature Machine Intelligence*, 5(12):1447–1457, 2023a.
- Liu, Y., Wang, L., Liu, M., Lin, Y., Zhang, X., Oztekin, B., and Ji, S. Spherical message passing for 3d molecular graphs. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=givsRXsOt9r>.
- Liu, Z., Zhang, W., Xia, Y., Wu, L., Xie, S., Qin, T., Zhang, M., and Liu, T.-Y. Molxpt: Wrapping molecules with text for generative pre-training. *arXiv preprint arXiv:2305.10688*, 2023b.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Luo, R., Sun, L., Xia, Y., Qin, T., Zhang, S., Poon, H., and Liu, T.-Y. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings in Bioinformatics*, 23(6):bbac409, 2022a.
- Luo, S., Chen, T., Xu, Y., Zheng, S., Liu, T.-Y., Wang, L., and He, D. One transformer can understand both 2d & 3d molecular data. *arXiv preprint arXiv:2210.01765*, 2022b.
- Luo, Y. and Ji, S. An autoregressive flow model for 3D molecular geometry generation from scratch. In *International Conference on Learning Representations*, 2022.
- Mansimov, E., Mahmood, O., Kang, S., and Cho, K. Molecular geometry prediction using a deep generative graph neural network. *Scientific reports*, 9(1):1–13, 2019.
- Mao, J., Wang, J., Zeb, A., Cho, K.-H., Jin, H., Kim, J., Lee, O., Wang, Y., and No, K. T. Deep molecular generative

- model based on variant transformer for antiviral drug design. *Available at SSRN 4345811*, 2023a.
- Mao, J., Wang, J., Zeb, A., Cho, K.-H., Jin, H., Kim, J., Lee, O., Wang, Y., and No, K. T. Transformer-based molecular generative model for antiviral drug design. *Journal of Chemical Information and Modeling*, 2023b.
- Masters, D., Dean, J., Klaser, K., Li, Z., Maddrell-Mander, S., Sanders, A., Helal, H., Beker, D., Rampásek, L., and Beaini, D. Gps++: An optimised hybrid mpnn/transformer for molecular property prediction. *arXiv preprint arXiv:2212.02229*, 2022.
- Mazuz, E., Shtar, G., Shapira, B., and Rokach, L. Molecule generation using transformers and policy gradient reinforcement learning. *Scientific Reports*, 13(1):8799, 2023.
- McKay, B. D. and Piperno, A. Practical graph isomorphism, ii. *Journal of symbolic computation*, 60:94–112, 2014.
- McKay, B. D. et al. Practical graph isomorphism. 1981.
- Özcelik, R., de Ruiter, S., and Grisoni, F. Structured state-space sequence models for de novo drug design. 2023.
- Özcelik, R., de Ruiter, S., Criscuolo, E., and Grisoni, F. Chemical language modeling with structured state spaces. 2024.
- Pei, Q., Zhang, W., Zhu, J., Wu, K., Gao, K., Wu, L., Xia, Y., and Yan, R. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. *arXiv preprint arXiv:2310.07276*, 2023.
- Qiang, B., Song, Y., Xu, M., Gong, J., Gao, B., Zhou, H., Ma, W.-Y., and Lan, Y. Coarse-to-fine: a hierarchical diffusion model for molecule generation in 3D. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 28277–28299. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/qiang23a.html>.
- Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. Improving language understanding by generative pre-training. 2018.
- Ramakrishnan, R., Dral, P. O., Rupp, M., and Von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- Satorras, V. G., Hoogeboom, E., Fuchs, F. B., Posner, I., and Welling, M. E(n) equivariant normalizing flows. In *Advances in Neural Information Processing Systems*, 2021a. URL https://openreview.net/forum?id=N5hQI_RowVA.
- Satorras, V. G., Hoogeboom, E., and Welling, M. E(n) equivariant graph neural networks. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9323–9332. PMLR, 18–24 Jul 2021b. URL <https://proceedings.mlr.press/v139/satorras21a.html>.
- Schaeffer, R., Miranda, B., and Koyejo, S. Are emergent abilities of large language models a mirage? *Advances in Neural Information Processing Systems*, 36, 2024.
- Shi, C., Luo, S., Xu, M., and Tang, J. Learning gradient fields for molecular conformation generation. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9558–9568. PMLR, 18–24 Jul 2021.
- Shi, Y., Zheng, S., Ke, G., Shen, Y., You, J., He, J., Luo, S., Liu, C., He, D., and Liu, T.-Y. Benchmarking graphormer on large-scale molecular modeling datasets. *arXiv preprint arXiv:2203.04810*, 2022.
- Simm, G. and Hernandez-Lobato, J. M. A generative model for molecular distance geometry. In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 8949–8958. PMLR, 13–18 Jul 2020.
- Thölke, P. and De Fabritiis, G. Equivariant transformers for neural network based molecular potentials. In *International Conference on Learning Representations*, 2021.
- Thomas, H. et al. Introduction to algorithms, 2009.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Tysinger, E. P., Rai, B. K., and Sinitskiy, A. V. Can we quickly learn to “translate” bioactive molecules with transformer models? *Journal of Chemical Information and Modeling*, 63(6):1734–1744, 2023.
- Uhrin, M. Through the eyes of a descriptor: Constructing complete, invertible descriptions of atomic environments. *Phys. Rev. B*, 104:144110, Oct 2021. doi: 10.1103/PhysRevB.104.144110. URL <https://link.aps.org/doi/10.1103/PhysRevB.104.144110>.
- Ünlü, A., Çevrim, E., Sarıgün, A., Çelikkilek, H., Güvenilir, H. A., Koyaş, A., Kahraman, D. C., Rifaioğlu, A., and Olğaç, A. Target specific de novo design of drug candidate molecules with graph transformer-based generative adversarial networks. *arXiv preprint arXiv:2302.07868*, 2023.

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Vignac, C., Osman, N., Toni, L., and Frossard, P. Midi: Mixed graph and 3d denoising diffusion for molecule generation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 560–576. Springer, 2023.
- Wang, L., Liu, Y., Lin, Y., Liu, H., and Ji, S. ComENet: Towards complete and efficient message passing for 3d molecular graphs. In *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=mCzMqEWSFJ>.
- Wang, Z., Wang, Z., Srinivasan, B., Ioannidis, V. N., Rangwala, H., and Anubhai, R. Biobridge: Bridging biomedical foundation models via knowledge graph. *arXiv preprint arXiv:2310.03320*, 2023.
- Weininger, D. SMILES, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.
- Wu, J.-N., Wang, T., Chen, Y., Tang, L.-J., Wu, H.-L., and Yu, R.-Q. Fragment-based t-smiles for de novo molecular generation. *arXiv preprint arXiv:2301.01829*, 2023.
- Wu, L., Gong, C., Liu, X., Ye, M., and qiang liu. Diffusion-based molecule generation with informative prior bridges. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=TJUNtiZITKE>.
- Xia, J., Zhu, Y., Du, Y., Liu, Y., and Li, S. A systematic survey of chemical pre-trained models. *IJCAI*, 2023.
- Xie, T., Wan, Y., Huang, W., Yin, Z., Liu, Y., Wang, S., Linghu, Q., Kit, C., Grazian, C., Zhang, W., et al. DARWIN series: Domain specific large language models for natural science. *arXiv preprint arXiv:2308.13565*, 2023.
- Xu, M., Luo, S., Bengio, Y., Peng, J., and Tang, J. Learning neural generative dynamics for molecular conformation generation. In *International Conference on Learning Representations*, 2021a. URL <https://openreview.net/forum?id=pAbm1qfheGk>.
- 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. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11537–11547. PMLR, 18–24 Jul 2021b.
- Xu, M., Yu, L., Song, Y., Shi, C., Ermon, S., and Tang, J. GeoDiff: A geometric diffusion model for molecular conformation generation. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=PzcvxEMzvQC>.
- Xu, M., Powers, A., Dror, R., Ermon, S., and Leskovec, J. Geometric latent diffusion models for 3D molecule generation. In *Proceedings of the 39th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2023a.
- Xu, Z., Luo, Y., Zhang, X., Xu, X., Xie, Y., Liu, M., Dickerson, K., Deng, C., Nakata, M., and Ji, S. Molecule3d: A benchmark for predicting 3d geometries from molecular graphs. *arXiv preprint arXiv:2110.01717*, 2021c.
- Xu, Z., Lei, X., Ma, M., and Pan, Y. Molecular generation and optimization of molecular properties using a transformer model. *Big Data Mining and Analytics*, 7(1): 142–155, 2023b.
- Yoshikai, Y., Mizuno, T., Nemoto, S., and Kusuhara, H. Difficulty in learning chirality for transformer fed with smiles. *arXiv preprint arXiv:2303.11593*, 2023.
- Zhang, Q., Ding, K., Lyv, T., Wang, X., Yin, Q., Zhang, Y., Yu, J., Wang, Y., Li, X., Xiang, Z., et al. Scientific large language models: A survey on biological & chemical domains. *arXiv preprint arXiv:2401.14656*, 2024.
- Zhang, W., Wang, X., Nie, W., Eaton, J., Rees, B., and Gu, Q. MoleculeGPT: Instruction following large language models for molecular property prediction. In *NeurIPS 2023 Workshop on New Frontiers of AI for Drug Discovery and Development*, 2023a.
- Zhang, X., Wang, L., Helwig, J., Luo, Y., Fu, C., Xie, Y., Liu, M., Lin, Y., Xu, Z., Yan, K., et al. Artificial intelligence for science in quantum, atomistic, and continuum systems. *arXiv preprint arXiv:2307.08423*, 2023b.
- Zhao, H., Liu, S., Ma, C., Xu, H., Fu, J., Deng, Z.-H., Kong, L., and Liu, Q. Gimlet: A unified graph-text model for instruction-based molecule zero-shot learning. *bioRxiv*, pp. 2023–05, 2023a.
- Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023b.

A. Broader Impacts and Limitations

Our work demonstrates the significant potential of pure language models (LMs) in revolutionizing molecular design and drug discovery by effectively transforming geometric information. The challenge of molecule design is particularly daunting when scientific experiments are cost-prohibitive or impractical. In many real-world scenarios, data collection is confined to specific chemical domains, yet the ability to generate molecules for broader tasks where experimental validation is difficult remains crucial. Traditional diffusion-based models fall short in terms of efficiency, scalability, and the ability to learn from extensive databases or transfer knowledge across different tasks. In contrast, LMs exhibit inherent advantages in these areas. We envision the development of efficient, large-scale models trained on vast chemical databases that can function across multiple datasets and molecular tasks. By introducing LMs into the 3D molecule generation field, we unlock substantial potential for broad scientific impact.

Our research adheres strictly to ethical guidelines, with no involvement of human subjects or potential privacy and fairness issues. This work aims to advance the field of Machine Learning and AI for drug discovery, with no immediate societal consequences requiring specific attention. We foresee no potential for malicious or unintended usage beyond known chemical applications. However, we recognize that all technological advancements carry inherent risks, and we advocate for ongoing evaluation of the broader implications of our methodology in various contexts.

We admit certain limitations, including that rounding up numerical values to certain decimal places bring information loss and discretized numbers impair generalization abilities across the continuous domain of real numbers. However, this is a trade-off between advantages brought by our model-agnostic framework. Due to the discrete nature of vocabularies, LMs depend on extensive pre-training corpora, fine-grained tokenization, or emergent abilities for better generalization, balancing high precision and versatility. Geo2Seq operates solely on the input data, which allows independence from model architecture and training techniques and provides reuse flexibility. This also means that we can effortlessly apply Geo2Seq on the latest generative language models, making seamless use of their capabilities. Future work points towards expanding on conditional tasks and exploring advanced tokenization techniques to enhance the model’s performance and applicability.

B. Proofs

B.1. Proof of Lemma 3.2

First, we define the isomorphism problem for attributed graphs as follows.

Definition B.1. [Graph Isomorphism] Let $G_1 = (V_1, E_1, A_1)$ and $G_2 = (V_2, E_2, A_2)$ be two graphs, where V_i denotes the set of vertices, E_i denotes the set of edges, and A_i denotes the node attributes of G_i for $i = 1, 2$. Let $attr(v)$ denote the node attributes of vertex v . The graphs G_1 and G_2 are said to be isomorphic, denoted as $G_1 \cong G_2$, if there exists a bijection $b : V_1 \rightarrow V_2$ such that for every vertex $v \in V_1$, $attr(v) \in A_1 = attr(b(v)) \in A_2$, and for every pair of vertices $u, v \in V_1$, $(u, v) \in E_1 \Leftrightarrow (b(u), b(v)) \in E_2$.

Next we prove Lemma 3.2.

Lemma (Colored Canonical Labeling for Graph Isomorphism). *Let $G_1 = (V_1, E_1, A_1)$ and $G_2 = (V_2, E_2, A_2)$ be two finite, undirected graphs where V_i denotes the set of vertices, E_i denotes the set of edges, and A_i denotes the node attributes of the graph G_i for $i = 1, 2$. Let $\mathbf{L} : \mathcal{G} \rightarrow \mathcal{L}$ be a function that maps a graph $G \in \mathcal{G}$, the set of all finite, undirected graphs, to its canonical labeling $\mathbf{L}(G) \in \mathcal{L}$, the set of all possible canonical labelings, as produced by the Nauty algorithm. Then the following equivalence holds:*

$$\mathbf{L}(G_1) = \mathbf{L}(G_2) \iff G_1 \cong G_2$$

where $G_1 \cong G_2$ denotes that the graphs G_1 and G_2 are isomorphic.

The Nauty algorithm, tailored for CL and computing graph automorphism groups, presents rigorous mathematical underpinnings to guarantee the CL properties. Here we leave out the proof of Nauty algorithm’s rigor for canonical labeling, which is detailed in the work of McKay & Piperno (2014). The key is the refinement process ensuring that the partitioning of the graph’s vertices is done in such a way that any two isomorphic graphs will end with the same partition structure.

B.2. Proof of Lemma 3.3

Lemma. *Let $G = (\mathbf{z}, \mathbf{R})$ be a 3D graph with node type vector \mathbf{z} and node coordinate matrix \mathbf{R} . Let \mathbf{F} be the equivariant global frame of graph G built based on the first three non-collinear nodes in $\mathbf{L}(G)$. $f(\cdot)$ is our function that maps 3D coor-*

dinate matrix \mathbf{R} to spherical representations \mathbf{S} under the equivariant global frame \mathbf{F} . Then for any $SE(3)$ transformation g , we have $f(\mathbf{R}) = f(g(\mathbf{R}))$. Given spherical representations $\mathbf{S} = f(\mathbf{R})$, there exist a $SE(3)$ transformation g , such that $f^{-1}(\mathbf{S}) = g(\mathbf{R})$.

Proof. Let ℓ_1, ℓ_2 , and ℓ_F be the indices of the first three non-collinear atoms in G . Then the global frame $\mathbf{F} = (\mathbf{x}, \mathbf{y}, \mathbf{z})$ is

$$\begin{aligned}\mathbf{x} &= \text{normalize}(\mathbf{r}_{\ell_2} - \mathbf{r}_{\ell_1}) \\ \mathbf{y} &= \text{normalize}((\mathbf{r}_{\ell_F} - \mathbf{r}_{\ell_1}) \times \mathbf{x}_1) \\ \mathbf{z} &= \mathbf{x} \times \mathbf{y}\end{aligned}$$

For a $SE(3)$ transformation g , let $\mathbf{R}' = g(\mathbf{R}) = \mathbf{Q}\mathbf{R} + \mathbf{b}$. Then the global frame $\mathbf{F}' = (\mathbf{x}', \mathbf{y}', \mathbf{z}')$ is

$$\begin{aligned}\mathbf{x}' &= \text{normalize}(\mathbf{r}_{\ell_2} - \mathbf{r}_{\ell_1}) = \text{normalize}(g(\mathbf{r}_{\ell_2}) - g(\mathbf{r}_{\ell_1})) = \mathbf{Q}\mathbf{x} \\ \mathbf{y}' &= \text{normalize}((\mathbf{r}_{\ell_F} - \mathbf{r}_{\ell_1}) \times \mathbf{x}_2) = \text{normalize}((g(\mathbf{r}_{\ell_F}) - g(\mathbf{r}_{\ell_1})) \times \mathbf{x}_2) = \mathbf{Q}\mathbf{y} \\ \mathbf{z}' &= \mathbf{x}' \times \mathbf{y}' = (\mathbf{Q}\mathbf{x}) \times (\mathbf{Q}\mathbf{y}) = \mathbf{Q}\mathbf{z}\end{aligned}$$

Thus $\mathbf{F}' = \mathbf{Q}\mathbf{F}$. Here $\text{normalize}(\cdot)$ is the function to normalize a vector to the corresponding unit vector. Then $\forall i$, the spherical representations $f(\mathbf{R})_{\ell_i}$ is

$$\begin{aligned}d_{\ell_i} &= \|\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}\|_2 \\ \theta_{\ell_i} &= \arccos((\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{z} / d_{\ell_i}) \\ \phi_{\ell_i} &= \text{atan2}((\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{y}, (\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{x})\end{aligned}$$

Similarly, the spherical representations $f(\mathbf{R}')_{\ell_i}$ is

$$\begin{aligned}d'_{\ell_i} &= \|\mathbf{r}'_{\ell_i} - \mathbf{r}'_{\ell_1}\|_2 = \|g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})\|_2 = d_{\ell_i} \\ \theta'_{\ell_i} &= \arccos((\mathbf{r}'_{\ell_i} - \mathbf{r}'_{\ell_1}) \cdot \mathbf{z}' / d'_{\ell_i}) = \arccos((g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})) \cdot \mathbf{z}' / d'_{\ell_i}) = \theta_{\ell_i} \\ \phi'_{\ell_i} &= \text{atan2}((\mathbf{r}'_{\ell_i} - \mathbf{r}'_{\ell_1}) \cdot \mathbf{y}', (\mathbf{r}'_{\ell_i} - \mathbf{r}'_{\ell_1}) \cdot \mathbf{x}') = \text{atan2}((g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})) \cdot \mathbf{y}', (g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})) \cdot \mathbf{x}') = \phi_{\ell_i}\end{aligned}$$

Therefore, we show that $f(\mathbf{R}) = f(g(\mathbf{R}))$. Next, we consider the function $f^{-1}(\cdot)$. For all i , the three terms in $f^{-1}(\mathbf{S})_{\ell_i}$ are

$$\begin{aligned}d_{\ell_i} \sin(\theta_{\ell_i}) \cos(\phi_{\ell_i}) \\ d_{\ell_i} \sin(\theta_{\ell_i}) \sin(\phi_{\ell_i}) \\ d_{\ell_i} \cos \theta_{\ell_i}\end{aligned} \tag{3}$$

Then we have $\mathbf{r}_{\ell_i} = f^{-1}(\mathbf{S})_{\ell_i}^T \mathbf{F} + \mathbf{r}_{\ell_0}$. Therefore, we show that there exist a $SE(3)$ transformation g , such that $g(f^{-1}(\mathbf{S})) = \mathbf{R}$. \square

B.3. Proof of Theorem 3.5

First we establish a lemma and provide its proof.

Lemma B.2. Let $G_1 = (\mathbf{z}_1, \mathbf{R}_1)$ and $G_2 = (\mathbf{z}_2, \mathbf{R}_2)$ be two 3D graphs, where \mathbf{z}_i is the node type vector and \mathbf{R}_i is the node coordinate matrix of the molecule G_i for $i = 1, 2$. Let $\mathbf{L}(G)$ be the canonical label of graph G . We have $G_1 \cong G_2$. Let ℓ_i and ℓ'_i denote the indexes of the node labeled i correspondingly in $\mathbf{L}(G_1)$ and $\mathbf{L}(G_2)$, respectively. Let \mathbf{F} be the equivariant global frame of graph G built based on the first three non-collinear atoms in $\mathbf{L}(G)$. Let $f : \mathcal{G} \rightarrow \mathcal{S}$ be a surjective function that maps a 3D graph $G \in \mathcal{G}$ to its spherical representations $\mathbf{S} = f(G) \in \mathcal{S}$ under the equivariant global frame \mathbf{F} . Then the following equivalence holds:

$$\forall i \in V_1, f(G_1)_{\ell_i} = f(G_2)_{\ell'_i} \iff G_1 \cong_{3D} G_2$$

where $G_1 \cong_{3D} G_2$ denotes that the graphs G_1 and G_2 are 3D isomorphic.

Proof. Let $\mathbf{L}(G)$ be the canonical labeling of graph G . Let ℓ_i and ℓ'_i denote the index of the node labeled i correspondingly in $\mathbf{L}(G_1)$ and $\mathbf{L}(G_2)$, respectively. We have

$$G_1 \cong_{3D} G_2 \iff \begin{cases} G_1 \cong G_2, \text{ and} \\ \text{there exists a 3D transformation } g \in SE(3) \text{ such that } \mathbf{r}_{\ell'_i}^{G_2} = g(\mathbf{r}_{\ell_i}^{G_1}). \end{cases}$$

Specifically, $g(\mathbf{r}_{\ell_i}) = \mathbf{Q}\mathbf{r}_{\ell_i} + \mathbf{b}$. Here \mathbf{Q} is a rotation matrix, and \mathbf{b} is a translation vector.

Let ℓ_1, ℓ_2 , and ℓ_F be the indices of the first three non-collinear atoms in G_1 . Then the equivariant global frame $F_1 = (\mathbf{x}_1, \mathbf{y}_1, \mathbf{z}_1)$ is

$$\begin{aligned}\mathbf{x}_1 &= \text{normalize}(\mathbf{r}_{\ell_2} - \mathbf{r}_{\ell_1}) \\ \mathbf{y}_1 &= \text{normalize}((\mathbf{r}_{\ell_F} - \mathbf{r}_{\ell_1}) \times \mathbf{x}_1) \\ \mathbf{z}_1 &= \mathbf{x}_1 \times \mathbf{y}_1\end{aligned}$$

Here $\text{normalize}(\cdot)$ is the function to normalize a vector to the corresponding unit vector. Then $\forall i$, the spherical representations $f(G_1)_{\ell_i}$ is

$$\begin{aligned}d_{\ell_i} &= \|\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}\|_2 \\ \theta_{\ell_i} &= \arccos((\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{z}_1 / d_{\ell_i}) \\ \phi_{\ell_i} &= \text{atan2}((\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{y}_1, (\mathbf{r}_{\ell_i} - \mathbf{r}_{\ell_1}) \cdot \mathbf{x}_1)\end{aligned}$$

Similarly, for G_2 , let ℓ'_1, ℓ'_2 , and ℓ'_F be the indices of the first three non-collinear atoms. Then the equivariant global frame $F_2 = (\mathbf{x}_2, \mathbf{y}_2, \mathbf{z}_2)$ is

$$\begin{aligned}\mathbf{x}_2 &= \text{normalize}(\mathbf{r}_{\ell'_2} - \mathbf{r}_{\ell'_1}) = \text{normalize}(g(\mathbf{r}_{\ell_2}) - g(\mathbf{r}_{\ell_1})) = Q\mathbf{x}_1 \\ \mathbf{y}_2 &= \text{normalize}((\mathbf{r}_{\ell'_F} - \mathbf{r}_{\ell'_1}) \times \mathbf{x}_2) = \text{normalize}((g(\mathbf{r}_{\ell_F}) - g(\mathbf{r}_{\ell_1})) \times \mathbf{x}_2) = Q\mathbf{y}_1 \\ \mathbf{z}_2 &= \mathbf{x}_2 \times \mathbf{y}_2 = (Q\mathbf{x}_1) \times (Q\mathbf{y}_1) = Q\mathbf{z}_1\end{aligned}$$

Then $\forall i$, the spherical representations $f(G_2)_{\ell'_i}$ is

$$\begin{aligned}d_{\ell'_i} &= \|\mathbf{r}_{\ell'_i} - \mathbf{r}_{\ell'_1}\|_2 = \|g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})\|_2 = d_{\ell_i} \\ \theta_{\ell'_i} &= \arccos((\mathbf{r}_{\ell'_i} - \mathbf{r}_{\ell'_1}) \cdot \mathbf{z}_2 / d_{\ell'_i}) = \arccos((g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})) \cdot \mathbf{z}_2 / d_{\ell_i}) = \theta_{\ell_i} \\ \phi_{\ell'_i} &= \text{atan2}((\mathbf{r}_{\ell'_i} - \mathbf{r}_{\ell'_1}) \cdot \mathbf{y}_2, (\mathbf{r}_{\ell'_i} - \mathbf{r}_{\ell'_1}) \cdot \mathbf{x}_2) = \text{atan2}((g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})) \cdot \mathbf{y}_2, (g(\mathbf{r}_{\ell_i}) - g(\mathbf{r}_{\ell_1})) \cdot \mathbf{x}_2) = \phi_{\ell_i}\end{aligned}$$

Therefore, we show that $G_1 \cong_{3D} G_2 \iff \forall i, \in V_1, f(G_1)_{\ell_i} = f(G_2)_{\ell'_i}$ holds. \square

Then we prove Theorem 3.5.

Theorem (Bijective mapping between 3D graph isomorphism and sequence). *Let $G_1 = (\mathbf{z}_1, \mathbf{R}_1)$ and $G_2 = (\mathbf{z}_2, \mathbf{R}_2)$ be two 3D graphs, where \mathbf{z}_j is the node type vector and \mathbf{R}_j is the node coordinate matrix of the molecule G_j for $j = 1, 2$. Let $\mathbf{L}_m(G)$ be the canonical label for 3D graph and $f : \mathcal{G} \rightarrow \mathcal{S}$ be the function that maps a 3D graph G to its spherical representations. Given graph G with n nodes and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times m}$, where $m \in \mathbb{Z}$, we define $\mathbf{L}_m(G) \otimes \mathbf{X} = \text{concat}(\mathbf{x}_{\ell_1}, \dots, \mathbf{x}_{\ell_n})$, where ℓ_i is the node index of the node labeled i in $\mathbf{L}_m(G)$, and $\text{concat}(\cdot)$ concatenates elements as a sequence. Define*

$$\text{Geo2Seq}(G_i) = \mathbf{L}_m(G) \otimes (\mathbf{z}, f(G)) = \mathbf{L}_m(G) \otimes \mathbf{X},$$

where $\mathbf{x}_i = [z_i, d_i, \theta_i, \phi_i]$. Then $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is a surjective function, and the following equivalence holds:

$$\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2) \iff G_1 \cong_{3D} G_2$$

where $G_1 \cong_{3D} G_2$ denotes that the graphs G_1 and G_2 are 3D isomorphic.

Proof. First, we prove that $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is a surjective function. Given the definition

$$\text{Geo2Seq}(G_i) = \mathbf{L}_m(G) \otimes (\mathbf{z}, f(G)) = \mathbf{L}_m(G) \otimes \mathbf{X},$$

where $\mathbf{x}_i = [z_i, d_i, \theta_i, \phi_i]$, we need to prove that all operations are deterministic. \otimes and \mathbf{z} are defined to be deterministic, and $f : \mathcal{G} \rightarrow \mathcal{S}$ is a function. $\mathbf{L}_m(G_j)$ outputs the automorphism group of G_j 's canonical label. By definition, the automorphism group contain different labels of the strictly identical graph. Let ℓ_i and ℓ'_i describe two different sets of labels of the same automorphism group with n nodes; since the graphs are identical,

$$[z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}] = [z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}] \text{ for } i = 1, \dots, n.$$

Thus $\text{concat}(\mathbf{x}_{\ell_1}, \dots, \mathbf{x}_{\ell_n}) = \text{concat}(\mathbf{x}_{\ell'_1}, \dots, \mathbf{x}_{\ell'_n})$, i.e., different labels of one automorphism group produce identical sequences with Geo2Seq . Therefore, $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is a well-defined function; given a 3D molecule, we can uniquely construct a 1D sequence from Geo2Seq .

Next we prove Geo2Seq 's surjectivity. Given any output sequence $q \in \mathcal{U}$ of Geo2Seq , the sequence is in the format

$$q = \text{concat}([z_1, d_1, \theta_1, \phi_1], \dots, [z_n, d_n, \theta_n, \phi_n]).$$

For the nodes in q , we denote with $S = [[d_1, \theta_1, \phi_1], \dots, [d_n, \theta_n, \phi_n]]$. Given the surjectivity of the spherical representation function $f : \mathcal{G} \rightarrow \mathcal{S}$ and the defined $f^{-1} : \mathcal{S} \rightarrow \mathcal{G}$, there must be a unique $G(\mathbf{z}, \mathbf{R}) \in \mathcal{G}$ where $S = f(G)$. Therefore, \forall

output sequence $q \in \mathcal{U}$ there exists $G(\mathbf{z}, \mathbf{R}) \in \mathcal{G}$ s.t. $q = \text{Geo2Seq}(G)$,

i.e., Geo2Seq is surjective; given a sequence output of Geo2Seq, we can uniquely reconstruct a 3D molecule.

Now we prove the equivalence $\text{Geo2Seq}(G_i) = \text{Geo2Seq}(G_j) \iff G_i \cong_{3D} G_j$, starting from right to left. Considering Lemma 3.2 for molecule $G = (\mathbf{z}, \mathbf{R})$, we specify $G = (V, E, A)$ with $A = [\mathbf{z}, \mathbf{R}]$ and define the CL function for 3D molecule graphs as \mathbf{L}_m , which extends the equivalence in Lemma 3.2 to \mathbf{L}_m on molecules with 3D isomorphism. If $G_1 \cong_{3D} G_2$, i.e., graphs G_1 and G_2 are 3D isomorphic, then from Lemma 3.2 we know the canonical forms $\mathbf{L}_m(G_1) = \mathbf{L}_m(G_2)$. Let graphs G_1 and G_2 have numbers of node n . Let ℓ_i and ℓ'_i be the denotations of a corresponding pair of canonical labelings from $\mathbf{L}_m(G_1)$ and $\mathbf{L}_m(G_2)$, respectively. Since graphs G_1 and G_2 are 3D isomorphic, from Def.3.4 we know $\forall i \in \mathcal{V}(G_1), z_{\ell_i} = z_{\ell'_i}$; and from Lemma B.2 we know $\forall i \in \mathcal{V}(G_1), f(G_1)_{\ell_i} = f(G_2)_{\ell'_i}$. Thus, we have

$$\begin{aligned} \text{Geo2Seq}(G_1) &= \mathbf{L}_m(G_1) \otimes (\mathbf{z}_1, f(G_1)) \\ &= \text{concat}_{z_j \in \mathbf{z}_1, d_j, \theta_j, \phi_j \in f(G_1), i=1, \dots, n}([z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}]) \\ &= \text{concat}_{z_j \in \mathbf{z}_2, d_j, \theta_j, \phi_j \in f(G_2), i=1, \dots, n}([z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}]) \\ &= \mathbf{L}_m(G_2) \otimes (\mathbf{z}_2, f(G_2)) = \text{Geo2Seq}(G_2). \end{aligned} \quad (4)$$

Note that if $\mathbf{L}_m(G_1)$ and $\mathbf{L}_m(G_2)$ contain automorphism groups larger than 1, we can include all possible labelings, which will all produce the same sequence later through Geo2Seq, as we have shown in detail above. However, this is a very rare case for real-world 3D graphs like molecules. Therefore, we have shown that if two molecules are 3D isomorphic considering atoms, bonds, and coordinates, their sequences resulting from Geo2Seq must be identical.

Finally, we prove the equivalence from left to right. We provide proof by contradiction. Given that $\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2)$, we assume that the graphs G_1 and G_2 are not 3D isomorphic. We denote with $G_1 = (\mathbf{z}_1, \mathbf{R}_1)$ and $G_2 = (\mathbf{z}_2, \mathbf{R}_2)$. If G_1 and G_2 are not even isomorphic for $A_i = z_i$, then from Def.B.1, there does not exist a node-to-node mapping from G_1 to G_2 , where each node is identically attributed and connected. And from Lemma 3.2, we know the canonical forms $\mathbf{L}_m(G_1) \neq \mathbf{L}_m(G_2)$. Thus for

$$\text{Geo2Seq}(G_1) = \text{concat}_{z_j \in \mathbf{z}_1, d_j, \theta_j, \phi_j \in f(G_1), i=1, \dots, n}([z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}]),$$

and

$$\text{Geo2Seq}(G_2) = \text{concat}_{z_j \in \mathbf{z}_2, d_j, \theta_j, \phi_j \in f(G_2), i=1, \dots, n}([z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}]),$$

there must be at least one pair of $z_{\ell_i}, z_{\ell'_i}$ where $z_{\ell_i} \neq z_{\ell'_i}$. Therefore, $\text{Geo2Seq}(G_1) \neq \text{Geo2Seq}(G_2)$, which is a contradiction to the initial condition that $\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2)$ and ends the proof.

If G_1 and G_2 are isomorphic for $A_i = z_i$, we continue with the following analyses. Let ℓ_i and ℓ'_i be the denotations of a corresponding pair of canonical labelings from $\mathbf{L}_m(G_1)$ and $\mathbf{L}_m(G_2)$, respectively. Let $f : \mathcal{G} \rightarrow \mathcal{S}$ be the surjective function mapping a 3D graph to its spherical representations. Since G_1 and G_2 are not 3D isomorphic, from Lemma B.2, we know there exists at least one

$$i \in V_1, \text{ s.t. } f(G_1)_{\ell_i} \neq f(G_2)_{\ell'_i};$$

otherwise, we would have

$$\forall i \in V_1, f(G_1)_{\ell_i} = f(G_2)_{\ell'_i} \Rightarrow G_1 \cong_{3D} G_2,$$

contradicting the above condition. Thus for

$$\text{Geo2Seq}(G_1) = \text{concat}_{z_j \in \mathbf{z}_1, d_j, \theta_j, \phi_j \in f(G_1), i=1, \dots, n}([z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}]),$$

and

$$\text{Geo2Seq}(G_2) = \text{concat}_{z_j \in \mathbf{z}_2, d_j, \theta_j, \phi_j \in f(G_2), i=1, \dots, n}([z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}]),$$

G_1 and G_2 are isomorphic, so

$$\forall i = 1, \dots, n, z_{\ell_i} = z_{\ell'_i};$$

at least one pair of spherical coordinates does not correspond, so there must be at least one pair of $(d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i})$ and $(d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i})$ where

$$(d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}) \neq (d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}).$$

Thus, $\text{Geo2Seq}(G_1) \neq \text{Geo2Seq}(G_2)$, which contradicts the initial condition that $\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2)$. Therefore, we have shown that if two constructed sequences from Geo2Seq are identical, their corresponding molecules must be 3D isomorphic considering atoms, bonds, and coordinates. This ends the proof. \square

B.4. Proof of Corollary 3.6

Corollary (Constrained bijective Mapping between 3D graph and sequence). *Let $G_1 = (\mathbf{z}_1, \mathbf{R}_1)$ and $G_2 = (\mathbf{z}_2, \mathbf{R}_2)$ be two 3D graphs, where \mathbf{z}_j is the node type vector and \mathbf{R}_j is the node coordinate matrix of the molecule G_j for*

$j = 1, 2$. Let $\mathbf{L}_m(G)$ be the canonical labeling for 3D graph and $f : \mathcal{G} \rightarrow \mathcal{S}$ be the function that maps a 3D graph G to its spherical representations. Given graph G with n nodes and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{n \times m}$, where $m \in \mathbb{Z}$, we define $\mathbf{L}_m(G) \otimes \mathbf{X} = \text{concat}(\mathbf{x}_{\ell_1}, \dots, \mathbf{x}_{\ell_n})$, where ℓ_i is the node index of the node labeled i in $\mathbf{L}_m(G)$, and $\text{concat}(\cdot)$ concatenates elements as a sequence. Define

$$\text{Geo2Seq}(G_i) = \mathbf{L}_m(G) \otimes (\mathbf{z}, f(G)) = \mathbf{L}_m(G) \otimes \mathbf{X},$$

where $\mathbf{x}_i = [z_i, d_i, \theta_i, \phi_i]$. Let the truncation of spherical coordinate values be after b decimal digits. Then $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is a surjective function, and the following equivalence holds:

$$\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2) \iff G_1 \cong_{3D - \lfloor \frac{10^{-b}}{2} \rfloor} G_2$$

where $G_1 \cong_{3D - \lfloor \frac{10^{-b}}{2} \rfloor} G_2$ denotes that the graphs G_1 and G_2 are $\frac{10^{-b}}{2}$ -constrained 3D isomorphic.

Proof. First, we prove that $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is a surjective function, which resembles the proof for Theorem 3.5. Given the definition

$$\text{Geo2Seq}(G_i) = \mathbf{L}_m(G) \otimes (\mathbf{z}, f(G)) = \mathbf{L}_m(G) \otimes \mathbf{X},$$

where $\mathbf{x}_i = [z_i, d_i, \theta_i, \phi_i]$, we need to prove that all operations are deterministic. \otimes and \mathbf{z} are defined to be deterministic, and $f : \mathcal{G} \rightarrow \mathcal{S}$ with truncation after certain decimal places is still a well-defined function. $\mathbf{L}_m(G_j)$ outputs the automorphism group of G_j 's canonical label. By definition, the automorphism group contain different labels of the strictly identical graph. Let ℓ_i and ℓ'_i describe two different sets of labels of the same automorphism group with n nodes; since the graphs are identical,

$$[z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}] = [z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}] \text{ for } i = 1, \dots, n.$$

Thus $\text{concat}(\mathbf{x}_{\ell_1}, \dots, \mathbf{x}_{\ell_n}) = \text{concat}(\mathbf{x}_{\ell'_1}, \dots, \mathbf{x}_{\ell'_n})$, i.e., different labels of one automorphism group produce identical sequences with Geo2Seq . Therefore, $\text{Geo2Seq} : \mathcal{G} \rightarrow \mathcal{U}$ is still a well-defined function; given a 3D molecule, we can uniquely construct a 1D sequence from Geo2Seq .

Next we prove Geo2Seq 's surjectivity. Given any output sequence $q \in \mathcal{U}$ of Geo2Seq , the sequence is in the format

$$q = \text{concat}([z_1, d_1, \theta_1, \phi_1], \dots, [z_n, d_n, \theta_n, \phi_n]).$$

For the nodes in q , we define $S_{\text{trun}} = [[d_1, \theta_1, \phi_1], \dots, [d_n, \theta_n, \phi_n]]$. Given the surjectivity of the spherical representation function $f : \mathcal{G} \rightarrow \mathcal{S}$ and the defined $f^{-1} : \mathcal{S} \rightarrow \mathcal{G}$, there must be a unique $G(\mathbf{z}, \mathbf{R}) \in \mathcal{G}$ where $S_{\text{trun}} = f(G)$. Therefore, \forall output sequence $q \in \mathcal{U}$ there exists $G(\mathbf{z}, \mathbf{R}) \in \mathcal{G}$ s.t. $q = \text{Geo2Seq}(G)$,

i.e., Geo2Seq is surjective; given a sequence output of Geo2Seq , we can uniquely reconstruct a 3D molecule.

Now we prove the equivalence $\text{Geo2Seq}(G_1) = \text{Geo2Seq}(G_2) \iff G_1 \cong_{3D} G_2$, starting from right to left. When a number is truncated after b decimal places, according to the rounding principle, the maximum error caused is $\epsilon \leq \frac{10^{-b}}{2}$. Considering Lemma 3.2 for molecule $G = (\mathbf{z}, \mathbf{R})$, we specify $G = (V, E, A)$ with $A = [\mathbf{z}, \mathbf{R}]$ and define the CL function for 3D molecule graphs as \mathbf{L}_m , which extends the equivalence in Lemma 3.2 to \mathbf{L}_m on molecules with 3D isomorphism. If $G_1 \cong_{3D - \lfloor \frac{10^{-b}}{2} \rfloor} G_2$, i.e., graphs G_1 and G_2 are $\frac{10^{-b}}{2}$ -constrained 3D isomorphic, then from Lemma 3.2 we know G_1 and G_2 are still isomorphic for $A_i = \mathbf{z}_i$, and the canonical forms $\mathbf{L}_m(G_1) = \mathbf{L}_m(G_2)$. Let graphs G_1 and G_2 have numbers of node n . Let ℓ_i and ℓ'_i be the denotations of a corresponding pair of canonical labelings from $\mathbf{L}_m(G_1)$ and $\mathbf{L}_m(G_2)$, respectively. Since graphs G_1 and G_2 are $\frac{10^{-b}}{2}$ -constrained 3D isomorphic, from Def.3.4 we know $\forall i \in \mathcal{V}(G_1), z_{\ell_i} = z_{\ell'_i}$; and from Lemma B.2 we know $\forall i \in \mathcal{V}(G_1), f(G_1)_{\ell_i} = f(G_2)_{\ell'_i}$ with $\frac{10^{-b}}{2}$ error range allowed for each numerical value. Thus, we still have

$$\begin{aligned} \text{Geo2Seq}(G_1) &= \mathbf{L}_m(G_1) \otimes (\mathbf{z}_1, f(G_1)) \\ &= \text{concat}_{z_j \in \mathbf{z}_1, d_j, \theta_j, \phi_j \in f(G_1), i=1, \dots, n} ([z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}]) \\ &= \text{concat}_{z_j \in \mathbf{z}_2, d_j, \theta_j, \phi_j \in f(G_2), i=1, \dots, n} ([z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}]) \\ &= \mathbf{L}_m(G_2) \otimes (\mathbf{z}_2, f(G_2)) = \text{Geo2Seq}(G_2). \end{aligned} \tag{5}$$

Note that if $\mathbf{L}_m(G_1)$ and $\mathbf{L}_m(G_2)$ contain automorphism groups larger than 1, we can include all possible labelings, which will all produce the same sequence later through Geo2Seq , as we have shown in detail above. However, this is a very rare case for real-world 3D graphs like molecules. Therefore, we have shown that if two molecules are 3D isomorphic considering atoms, bonds, and coordinates within the round-up error range $\frac{10^{-b}}{2}$, their sequences resulting from Geo2Seq must be identical.

Finally, we prove the equivalence from left to right. We provide proof by contradiction. Given that $\text{Geo2Seq}(G_1) =$

$Geo2Seq(G_2)$, we assume that the graphs G_1 and G_2 are not $\frac{|10^{-b}|}{2}$ -constrained 3D isomorphic. We denote with $G_1 = (z_1, R_1)$ and $G_2 = (z_2, R_2)$. If G_1 and G_2 are not even isomorphic for $A_i = z_i$, then from Def.B.1, there does not exist a node-to-node mapping from G_1 to G_2 , where each node is identically attributed and connected. And from Lemma 3.2, we know the canonical forms $L_m(G_1) \neq L_m(G_2)$. Thus for

$$Geo2Seq(G_1) = \text{concat}_{z_j \in z_1, d_j, \theta_j, \phi_j \in f(G_1), i=1, \dots, n}([z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}]),$$

and

$$Geo2Seq(G_2) = \text{concat}_{z_j \in z_2, d_j, \theta_j, \phi_j \in f(G_2), i=1, \dots, n}([z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}]),$$

there must be at least one pair of $z_{\ell_i}, z_{\ell'_i}$ where $z_{\ell_i} \neq z_{\ell'_i}$. Therefore, $Geo2Seq(G_1) \neq Geo2Seq(G_2)$, which is a contradiction to the initial condition that $Geo2Seq(G_1) = Geo2Seq(G_2)$ and ends the proof.

If G_1 and G_2 are isomorphic for $A_i = z_i$, we continue with the following analyses. Let ℓ_i and ℓ'_i be the denotations of a corresponding pair of canonical labelings from $L_m(G_1)$ and $L_m(G_2)$, respectively. Let $f : \mathcal{G} \rightarrow \mathcal{S}$ be the surjective function mapping a 3D graph to its spherical representations. Since G_1 and G_2 are not $\frac{|10^{-b}|}{2}$ -constrained 3D isomorphic, from Lemma B.2, we know there exists at least one

$$i \in V_1, s.t. f(G_1)_{\ell_i} \neq f(G_2)_{\ell'_i},$$

even with error range $\frac{|10^{-b}|}{2}$ allowed; otherwise, we would have

$$\forall i \in V_1, f(G_1)_{\ell_i} = f(G_2)_{\ell'_i} \Rightarrow G_1 \cong_{3D} G_2,$$

contradicting the above condition. Thus for

$$Geo2Seq(G_1) = \text{concat}_{z_j \in z_1, d_j, \theta_j, \phi_j \in f(G_1), i=1, \dots, n}([z_{\ell_i}, d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}]),$$

and

$$Geo2Seq(G_2) = \text{concat}_{z_j \in z_2, d_j, \theta_j, \phi_j \in f(G_2), i=1, \dots, n}([z_{\ell'_i}, d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}]),$$

G_1 and G_2 are isomorphic, so

$$\forall i = 1, \dots, n, z_{\ell_i} = z_{\ell'_i};$$

at least one pair of spherical coordinates does not correspond, so there must be at least one pair of $(d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i})$ and $(d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i})$ where

$$\min(|d_{\ell_i}, \theta_{\ell_i}, \phi_{\ell_i}| - |d_{\ell'_i}, \theta_{\ell'_i}, \phi_{\ell'_i}|) > \frac{|10^{-b}|}{2}.$$

Thus $Geo2Seq(G_1) \neq Geo2Seq(G_2)$, which contradicts the initial condition that $Geo2Seq(G_1) = Geo2Seq(G_2)$. Therefore, we have shown that if two constructed sequences from Geo2Seq are identical, their corresponding molecules must be 3D isomorphic considering atoms, bonds, and coordinates within the round-up error range $\frac{|10^{-b}|}{2}$. This ends the proof. \square

C. Ablation Studies

Table 3: Random generation performance with different atom generation orders.

Order	Atom Sta (%)	Mol Sta (%)	Valid (%)	Valid & Unique (%)
Canonical-locality	97.39	86.77	92.97	84.71
Canonical-nonlocality	96.45	81.36	90.89	83.37
Canonical-SMILES	97.35	85.86	92.97	84.05
DFS (Thomas et al., 2009)	95.95	81.54	90.45	82.48
BFS (Lee, 1961)	96.85	80.92	90.49	76.13
Dijkstra (Dijkstra, 2022)	95.29	77.25	88.97	73.52
Cuthill–McKee (Cuthill & McKee, 1969)	93.56	71.57	85.36	76.23
Hilbert-curve (Hilbert & Hilbert, 1935)	90.11	64.99	80.40	67.83
Random	64.87	20.14	43.16	38.44

To study the effects of atom order, 3D representations and tokenization of Geo2Seq on the generation performance of LLMs, we conduct a series of ablation experiments. Among all ablation experiments, we train 8-layer GPT models on QM9 dataset for 250 epochs with the same hyperparameters as Section 5.1 and use the random generation metrics in Section 5.1 to

Table 4: Random generation performance with different 3D representations.

3D representation	Atom Sta (%)	Mol Sta (%)	Valid (%)	Valid & Unique (%)
Original coordinates	91.1	58.1	75.6	55.1
Normalized coordinates	92.7	63.2	83.1	72.5
Invariant Cartesian coordinates	96.0	78.5	89.7	74.1
Inv-spherical coordinates	97.3	83.4	91.0	82.7
Inv-spherical coordinates-local distances	97.1	82.8	91.7	79.6

Table 5: Random generation performance with different tokenization.

Tokenization	Atom Sta (%)	Mol Sta (%)	Valid (%)	Valid & Unique (%)
Char-tokenization	90.5	43.7	71.5	71.0
BPE	85.3	55.3	74.4	57.6
Sub-tokenization	96.4	80.3	89.9	74.4
Comp-tokenization	97.0	82.2	91.0	75.5

compare the performance under different settings.

Ablation on atom order. First, we show that our proposed canonical order of atoms in Geo2Seq sequence representation is significant for LLMs to achieve good 3D molecular structure modeling. Specifically, we conduct an extended study of ordering algorithms, comparing our Geo2Seq with alternative canonicalization strategies as well as established traversing baselines. As we specified in Sec 3.1, theoretically, analyses and derivations apply to all rigorous CL algorithms. In the paper, we select Nauty Algorithm because its implementation has the best time efficiency among all existing CL algorithms. We implemented Nauty Algorithm for 3D molecules, where multiple strategies can be applied for the partitioning of graph vertices (a step in Nauty). We compare canonicalization strategies with/without locality considered. Canonicalization with locality considered can lead to better results, due to the importance of neighboring atom interactions in molecular evaluations. Given the similar nature, canonical SMILES produces a very similar ordering with "Nauty with locality", thus close in performances. The traversing baselines includes Breadth-First Search (BFS) (Lee, 1961), Depth-First Search (DFS) (Thomas et al., 2009), Dijkstra’s algorithm (Dijkstra, 2022), Cuthill–McKee algorithm (Cuthill & McKee, 1969), and Hilbert curve (Hilbert & Hilbert, 1935). We also compare with a Random sequence representation where atoms are randomly ordered. All the other settings of sequence representations remain the same. As Table 3 shows, canonicalization with locality considered can lead to better results, due to the importance of neighboring atom interactions in molecular evaluations. In addition, we can clearly observe that well-designed canonical ordering as in Geo2Seq significantly outperforms basic traverse strategies and the random order, which validates the significance of canonical order.

Advantage of Nauty Algorithm. Note that in the paper, we implement Nauty Algorithm for 3D molecules because: (1) its implementation has the best time efficiency among all existing CL algorithms; (2) it is naturally rigorous. The widely used canonical SMILES is based on the Morgan CL Algorithm, which is proven to be incomplete for isomorphism corner cases (such as two triangles versus one hexagon). While canonical SMILES solve corner cases by manual restrictions, Nauty Algorithm is elegantly rigorous. Still, we emphasize that all rigorous CL algorithms are usable for our method, while our contribution lies in achieving structural completeness and geometric invariance for LM learning of 3D molecules.

Ablation on 3D representation. Besides, we explore using different methods to represent 3D molecular structures. We compare the spherical coordinates in Geo2Seq with directly using the 3D Cartesian coordinates of atoms from QM9 xyz data files in sequences. We also study whether normalizing the xyz coordinates is effective by subtracting the xyz coordinates with the mass-center coordinates of each molecule. Additionally, we compare with using the $SE(3)$ -invariant Cartesian coordinates that are projected to the equivariant frame proposed in Section 3.2. We also explore adopting to manage distances in a more local scheme, which reduces the scale of the distances. We compare with "local distances", where our "distances to the global frame" are replaced with "relative distances to the previous atom" (except for the first atom) while the angles remain the same. Results in Table 4 demonstrate that LLMs achieve the best performance on spherical coordinates.

We believe this is due to that the numerical values of distances and angles of spherical coordinates lie in a smaller region than coordinates, which reduces outliers and makes it easier for LLMs to capture their correlation. Furthermore, both our spherical coordinates and that replaced with local distances achieve comparable results, while outperforming Cartesian coordinates. From these empirical results, we can analyze that the representation of azimuth and polar angles has brought sufficient advantage for LM learning over Cartesian coordinates, thus spherical representations with both distance schemes are showing promising performances. In addition, the similar performances could be attributed to that molecular systems often exhibit localized spatial structures (e.g., compact subunits or functional groups), which naturally constrain distances for most small molecules.

Advantage of invariant spherical representations. The above experiments show the superiority of invariant spherical coordinates over invariant Cartesian coordinates. While invariant Cartesian coordinates when our proposed equivariant frame is applied can also $SE(3)$ -invariance, spherical coordinates are advantageous in discretized representations. Compared to Cartesian coordinates, spherical coordinate values are bounded in a smaller region, namely, a range of $[0, \pi]$ or $[0, 2\pi]$. Given the same decimal place constraints, spherical coordinates require a smaller vocabulary size, and given the same vocabulary size, spherical coordinates present less information loss. This makes spherical coordinates advantageous in discretized representations and thus easier to be modeled by LMs. Lemma 3.3 and its proof aim to guarantee the validity that our proposed invariant spherical representations possess $SE(3)$ -invariance. We consider it as a part of our theoretical contribution towards the derivation of Theorem 3.5.

Ablation on tokenization. Finally, we explore other ways to tokenize real numbers in spherical coordinates. Instead of simply taking the complete real number as a token (Comp-tokenization), we try splitting it by the decimal point and treat every part as an individual token (Sub-tokenization). We also explore the common NLP tokenization method, including treating each character as a token (Char-tokenization) and Byte-Pair Encoding (BPE). We compare these tokenization methods in Table 5. Results show that our used Comp-tokenization leads to better performance. This shows that treating the complete real number as an individual token enables LLMs to capture 3D molecular structures more effectively.

Overall, through a series of ablation experiments, we show that canonical atom order, spherical coordinate representation and Comp-tokenization in Geo2Seq are all very useful in parsing 3D molecules to good sequence representations.

D. Experimental Details and Additional Results

D.1. Hyperparameters and Experimental Details

In the random generation experiment (Section 5.1), we apply two LMs, GPT (Radford et al., 2018) and Mamba (Gu & Dao, 2023), to our proposed Geo2Seq representations. For GPT models, we adopt the architecture of GPT-1, set the hidden dimension to 768, the number of attention head to 8, and the number of layers to 12 and 14 for QM9 and GEOM-DRUGS datasets, respectively. For Mamba models, we set the hidden dimension to 768 and the number of layers to 26 and 28 for QM9 and GEOM-DRUGS datasets, respectively. On QM9 dataset, we set the batch size to 32, base learning rate to 0.0004, the number of training epochs to 600 and 210 for GPT and Mamba models, respectively. On GEOM-DRUGS dataset, we set the batch size to 32, base learning rate to 0.0004, the number of training epochs to 20 and 25 for GPT and Mamba models, respectively. During model training, we use AdamW (Loshchilov & Hutter, 2019) optimizer and follow the commonly used linear warm up and cosine decay scheduler to adjust learning rates. Specifically, the learning rate first linearly increases from zero to the base learning rate 0.0004 when handling the first 10% of total training tokens, then gradually decreases to 0.00004 by the cosine decay scheduler. Besides, the tokenization of real numbers uses the precision of two and three decimal places for QM9 and GEOM-DRUGS datasets, respectively. In the controllable generation experiment (Section 5.2), we train 16-layer Mamba models for 200 epochs, and all the other hyperparameters and settings are the same as the random generation experiment. Based on data statistics, we set the context length to 512 for QM9 dataset and 744 for GEOM-DRUGS dataset throughout the experiments. All experiments on the QM9 dataset are conducted using a single NVIDIA A6000 GPU. Experiments on the GEOM-DRUGS dataset are deployed on 4 NVIDIA A100 GPUs.

D.2. Licenses

We strictly follow all licenses when using the public assets in this work. The QM9 dataset is under license CC-BY 4.0. The GEOM-DRUGS dataset is under license CC0 1.0. The code of EDM, GEOLDM, JODO, and MiDi is under MIT License.

D.3. Experiments on additional baselines and metrics

We extend our experiments with two more baselines, JODO (Huang et al., 2023) and MiDi (Vignac et al., 2023), which are diffusion models jointly generating 2D and 3D molecular information. We exclude them in experiments of the main paper, since the setting is not the same as ours. Our method follows works on 3D molecule generation without 2D information, such as bonds.

We extend the metrics of our evaluation for more comprehensive comparisons on random generation of QM9 dataset. We report the percentage of valid, unique and novel molecules, *i.e.*, that are not present in the training set. We also report the percentage of complete molecules in which all atoms are connected. Following JODO (Huang et al., 2023), we also include 2D metrics. Frechet ChemNet Distance (FCD) measures the distance between the test set and the generated set with the activation of the penultimate layer of ChemNet. Lower FCD values indicate more similarity between the two distributions. Similarity to the nearest neighbor (SNN) calculates an average Tanimoto similarity between the fingerprints of a generated molecule and its closest molecule in the test set. Fragment similarity (Frag) compares the distributions of BRICS fragments in the generated and test sets, and Scaffold similarity (Scaf) compares the frequencies of Bemis-Murcko scaffolds between them. Additionally, we include alignment metrics. For RDKit generated bonds, we compute the Maximum Mean Discrepancy (MMD) distances of the bond length (Bond), bond angle (Angle), and dihedral angle (Dihedral) distributions, and report their mean MMD distances. To ensure fair comparison, we evaluate the metrics of all methods on the generated 3D structures, and use RDKit to convert 3D structures to 2D graphs if needed. We use the same model and settings as the main paper for Geo2Seq, and follow the released codes for the baselines’ respective hyperparameter and settings. Table 6 reports the random generation results on QM9 dataset. According to the results, though our model is not designed to directly learn 2D information, the performance of our method is better than or comparable with baseline methods on all metrics including the 2D metrics, which demonstrates the effectiveness of our design.

Table 6: Additional random generation results on QM9 dataset.

Metric	EDM	GEOLDM	JODO	MiDi	Geo2Seq with Mamba
Atom Sta (%)	98.7	98.9	98.9	98.2	98.9
Mol Sta (%)	82.0	89.4	89.0	83.5	93.2
Valid (%)	91.9	93.8	94.9	95.2	97.1
Valid & Unique (%)	90.7	91.8	92.8	92.8	81.7
Valid & Unique & Novel (%)	83.0	83.1	85.2	85.5	71.2
Complete (%)	90.9	93.3	94.4	94.4	97.3
Bond Length MMD	0.18	0.12	0.27	1.09	0.08
Bond Angle MMD	0.04	0.04	0.05	0.05	0.04
Dihedral Angle MMD	0.003	0.003	0.0022	0.0033	0.0011
FCD	1.16	0.94	1.55	1.28	2.04
SNN	0.47	0.49	0.47	0.47	0.49
Frag	0.94	0.94	0.94	0.94	0.83
Scaf	0.29	0.33	0.25	0.26	0.38

Reporting the percentage of novel molecules is important in showing that language models can generate new molecules instead of merely memorizing the training dataset. Given our improvements on controllable generation is significant, we explore whether the generated molecules are different from the molecules in the training set. Thus we also extend the metric on controllable generation experiments. We use the same model and setting as the main paper. Table 7 presents the novelty results of controllable generation compared with EDM and JODO. Results show that our method achieves reasonably high novelty scores, which demonstrates that our method is not simply memorizing training data.

In addition, following Hoogetboom et al. (2022), we compare negative log-likelihood (NLL) performance on the random generation of QM9 dataset for Geo2Seq and baseline models that reports this metric. For this experiment, we use the same model and setting as the main paper. From Table 8, we can see the performance of our method is better than or comparable with all baseline methods, evidencing the validity of our model.

For more comprehensive comparisons, we also extend to include the metrics of Symphony (Daigavane et al., 2023) in our

Table 7: Additional controllable generation results for the percentage of valid, unique, and novel molecules on QM9 dataset.

Method	α	$\Delta\epsilon$	ϵ_{HOMO}	ϵ_{LUMO}	μ	C_v
EDM	87.0%	84.1%	79.8%	84.7%	73.0%	68.0%
JODO	86.5%	87.3%	86.7%	86.2%	86.8%	85.6%
Geo2Seq with Mamba	82.8%	82.8%	83.6%	83.0%	83.3%	83.6%

Table 8: Additional Negative Log Likelihood (NLL) comparisons of random generation on QM9 dataset.

Method	NLL
E-NF	-59.7
GDM	-94.7
EDM	-110.7
GEOLDM	-335.0
Geo2Seq with Mamba	<u>-242.0</u>

evaluation. As shown in Table 9,10,11, we compare the performances of baseline methods and Geo2Seq with Mamba on Symphony metrics. Multiple algorithms exist for bond order assignment: `xyz2mol` (Kim & Kim, 2015), OpenBabel (Banck et al., 2011) and a simple lookup table based on empirical pairwise distances in organic compounds (Hoogeboom et al., 2022). We perform the comparison between these algorithms for evaluating machine-learning generated 3D structures. In Table 9, we use each of these algorithms to infer the bonds and create a molecule from generated 3D molecular structure. A molecule is valid if the algorithm could successfully assign bond order with no net resulting charge. We also measure the uniqueness to see how many repetitions were present in the set of SMILES strings of valid generated molecules. Buttenschoen et al. (2023) showed that the predicted 3D structures from machine-learned protein-ligand docking models tend to be highly unphysical. Table 10 utilizes the PoseBusters framework to perform the following sanity checks to count how many of the predicted 3D structures are reasonable. The valid molecules from all models tend to be quite reasonable. Next, we evaluate models on how well they capture bonding patterns and the geometry of local environments found in the training set molecules as Table 11. We utilize the bispectrum (Uhrin, 2021) as a rotationally invariant descriptor of the geometry of local environments. Given a local environment with a central atom u , all of the neighbors of u are projected according to the inferred bonds onto the unit sphere S^2 . Then, the signal f is computed as a sum of Dirac delta distributions along the direction of each neighbor. The bispectrum $\mathcal{B}(f)$ of f is then defined as $\mathcal{B}(f) = \text{EXTRACTSCALARS}(f \otimes f \otimes f)$. Thus, f captures the distribution of atoms around u , and the bispectrum $\mathcal{B}(f)$ captures the geometry of this distribution. The bispectrum varies smoothly when f is varied and is guaranteed to be rotationally invariant. We follow Symphony and compute the bispectrum of local environments with at least 2 neighboring atoms, and exclude the pseudoscalars in the bispectra. For comparing discrete distributions, we use the symmetric Jensen-Shannon divergence (JSD) as Hoogeboom et al. (2022). Given the true distribution Q and the predicted distribution P , the Jensen-Shannon divergence between them is defined as: $D_{JS}(Q \parallel P) = \frac{1}{2}D_{KL}(Q \parallel M) + \frac{1}{2}D_{KL}(P \parallel M)$ where D_{KL} is the Kullback-Leibler divergence and $M = \frac{Q+P}{2}$ is the mean distribution. For continuous distributions, estimating the Jensen-Shannon divergence from samples is tricky without further assumptions on the distributions. We follow Symphony and use the MMD scores to compare samples from continuous distributions. Overall, the performance of our method is better than or comparable with baseline methods across the metrics, showing the effectiveness of our 3D molecule generation.

D.4. Generation Efficiency Analysis

We compare the generation efficiency of our method and the diffusion-based methods using a single NVIDIA A100 GPU and a batch size of 32. The results in Table 6 show that our method is much faster than diffusion-based methods, indicating the great efficiency of our method. Though we have take more memory compared to diffusion-based methods, our time efficiency is much better than diffusion-based methods. Throughput, or samples per second, is one of the most important metrics to measure generation efficiency. In particular, Geo2Seq with Mamba is more than 100 times faster than diffusion-based methods, indicating the high throughput of our method, a significant advantage in practical applications where speed is

Table 9: Additional validity and uniqueness percentages of molecules following Symphony.

Metric \uparrow	Symphony	EDM	G-SchNet	G-SphereNet	Geo2Seq
Validity via xyz2mol	83.50	86.74	74.97	26.92	95.42
Validity via OpenBabel	74.69	77.75	61.83	9.86	83.84
Validity via Lookup Table	68.11	90.77	80.13	16.36	97.55
Uniqueness via xyz2mol	97.98	99.16	96.73	21.69	98.88
Uniqueness via OpenBabel	99.61	99.95	98.71	7.51	99.91
Uniqueness via Lookup Table	97.68	98.64	93.20	23.29	98.95

Table 10: Percentage of valid molecules passing each PoseBusters test following Symphony.

Test \uparrow	Symphony	EDM	G-SchNet	G-SphereNet	Geo2Seq
All Atoms Connected	99.92	99.88	99.87	100.00	100.00
Reasonable Bond Angles	99.56	99.98	99.88	97.59	99.90
Reasonable Bond Lengths	98.72	100.00	99.93	72.99	100.00
Aromatic Ring Flatness	100.00	100.00	99.95	99.85	99.98
Double Bond Flatness	99.07	98.58	97.96	95.99	99.45
Reasonable Internal Energy	95.65	94.88	95.04	36.07	96.10
No Internal Steric Clash	98.16	99.79	99.57	98.07	99.33

Table 11: Additional comparison statistics of generated molecules to the training set for QM9 dataset following Symphony.

MMD of Bond Lengths \downarrow	Symphony	EDM	G-SchNet	G-SphereNet	Geo2Seq
C-H: 1.0	0.0739	0.0653	0.3817	0.1334	0.0488
C-C: 1.0	0.3254	0.0956	0.2530	1.0503	0.0705
C-O: 1.0	0.2571	0.0757	0.5315	0.6082	0.0712
C-N: 1.0	0.3086	0.1755	0.2999	0.4279	0.1056
N-H: 1.0	0.1032	0.1137	0.5968	0.1660	0.0965
C-O: 2.0	0.3033	0.0668	0.2628	2.0812	0.0667
O-N: 1.5	0.3707	0.1736	0.5828	0.4949	0.1570
O-H: 1.0	0.2872	0.1545	0.7899	0.1307	0.0990
C-C: 1.5	0.4142	0.1749	0.2051	0.8574	0.0832
C-N: 2.0	0.5938	0.3237	0.4194	2.1197	0.2676
MMD of Bispectra \downarrow	Symphony	EDM	G-SchNet	G-SphereNet	Geo2Seq
C: C2,H2	0.2165	0.1003	0.4333	0.6210	0.0955
C: C1,H3	0.2668	0.0025	0.0640	1.2004	0.0011
C: C3,H1	0.1111	0.2254	0.2045	1.1209	0.0867
C: C2,H1,O1	0.1500	0.2059	0.1732	0.8361	0.1058
C: C1,H2,O1	0.3300	0.1082	0.0954	1.6772	0.0802
O: C1,H1	0.0282	0.0056	0.0487	0.0030	0.0022
C: C2,H1,N1	0.1481	0.1521	0.1967	1.3461	0.1111
C: C2,H1	0.2525	0.0468	0.1788	0.2403	0.0851
C: C1,H2,N1	0.3631	0.2728	0.1610	0.9171	0.1285
N: C2,H1	0.0953	0.2339	0.2105	0.6141	0.1081
Jensen-Shannon Divergence \downarrow	Symphony	EDM	G-SchNet	G-SphereNet	Geo2Seq
Atom Type Counts	0.0003	0.0002	0.0011	0.0026	0.0002
Local Environment Counts	0.0039	0.0057	0.0150	0.1016	0.0035

crucial.

Table 12: Generation efficiency comparison between diffusion-based methods and our LM-based method.

Method	QM9			DRUG		
	Parameters	Memory	Sample/second	Parameters	Memory	Sample/second
EDM	5.3M	1.5GB	1.4	2.4M	7.4GB	0.1
GeoLDM	11.4M	1.5GB	1.4	5.5M	8.4GB	0.1
Geo2Seq with GPT	87.7M	2.4GB	8.3	105.4M	3.1GB	0.2
Geo2Seq with Mamba	91.8M	2.2GB	100.0	108.4M	2.6GB	16.7

D.5. Results with Pretraining

To show the advantage of pretraining, we compare the random generation performance on QM9 for models with and without pretraining on Molecule3D (Xu et al., 2021c) dataset, which includes around 4M molecules. Specifically, we conduct experiments on an 8-layer GPT model and a 20-layer Mamba model. The models are pretrained for 20 epochs and then finetuned for 200 epochs. The results in Table 13 demonstrate the advantage of pretraining. Future studies could explore pretraining on larger datasets.

Table 13: Random generation performance on QM9 for models with and without pretraining on Molecule3D dataset.

Method	Atom Sta (%)	Mol Sta (%)	Valid (%)	Valid & Unique (%)
Geo2Seq with GPT	97.0	82.2	91.0	75.5
Geo2Seq with GPT + pretraining	98.5	89.7	94.8	76.6
Geo2Seq with Mamba	97.4	86.8	93.0	78.8
Geo2Seq with Mamba + pretraining	98.3	89.4	94.9	83.5

E. Extended Studies

E.1. Scaling Laws

Scaling law refers to the relations between functional properties of interest, performance metrics in our case, and properties of the architecture or optimization process. In this section we explore the scaling laws of our models, specifically regarding parameter size, since they provide typical insights for LMs. Scaling laws in 3D molecule generation appears similar to that in NLP. We provide experiments on both GPT and Mamba in Table 14 and 15, respectively. As can be observed, LMs’ performances on molecules grow significantly with parameter size increase, similar to the emergence abilities widely-recognized in NLP tasks. As known from NLP studies (Schaeffer et al., 2024), model capabilities grow consistently with model size, while emergence abilities are largely caused by nonlinear metrics. This matches our observations, since the chemical metrics are hardly linear.

Table 14: Scaling laws on Geo2Seq with GPT model.

Parameter size - GPT	2556532	31309824	61650944	88012800	116342688
Atom sta(%)	76.2	89.6	96.5	98.3	98.5
Mol sta(%)	5.1	42.4	81.3	89.1	90.6
Valid(%)	45.5	73.1	90.9	94.3	95.1
Valid & Unique(%)	43.4	66.7	83.6	74.9	78.6

Note that we evaluate all models after 250 epochs for fairness concerns, while this fixed hyperparameter setting is not optimal for performances at all parameter sizes. Other settings are the same as the ablation studies.

Table 15: Scaling laws on Geo2Seq with Mamba model.

Parameter size - Mamba	2180352	31458048	61631232	93088512	121977600
Atom sta(%)	81.6	95.7	97.4	97.8	97.9
Mol sta(%)	13.6	79.2	86.8	88.3	89.0
Valid(%)	51.2	89.4	93	93.7	94.4
Valid & Unique(%)	49.6	78.7	78.8	82.6	83.5

E.2. Error Case Analysis

In the natural language domain, trained language models can produce error cases showing repetition or hallucinations. This is also a problem that often arises with LLMs. In this section, we provide the analysis of some error cases to introduce more insights into the field.

Similarly to NLP cases, our trained language models are showing repetition or hallucinations, especially when not trained to best convergence. This happens to both GPT and Mamba models. Below we show some error cases from a 16-layer Mamba model trained 150 epochs on the QM9 dataset. The error case below shows a typical repetition problem. The model generates repeated tokens for several periods, resulting in an invalid sample.

- H 0.00 0.00° 0.00° C 1.09 1.57° 0.00° N 2.02 2.15° 0.00° C 3.39 1.99° -0.02° H 3.98 2.10° 0.23° C 4.34 2.11° -0.35° H
4.43 2.38° -0.46° H 5.41 2.09° -0.29° H 4.29 1.96° -0.59° C 4.05 1.63° 0.04° H **4.09 4.09 4.09 4.09 4.09 4.09 4.09 4.09**
4.09 4.09 4.09 4.09 4.09 4.09 4.09 4.09 4.09 4.09 4.09 4.09 4.09 1.01° H 4.96 4.96 4.96 H 1.23° 1.23° 1.23° C 3.27
1.74° -0.03° H 4.02 1.83° -0.23° H 3.79 1.91° 0.17° H 3.79 1.53° -0.03°

For hallucination, our tokenization design actually prevents token-level hallucination by defining elements and whole-numerical-values as tokens, instead of using single characters. This prevents token-level hallucination, *i.e.*, non-existent elements or numbers such as ‘Hr’ or ‘-0.15’. However, there can still be sequence-level hallucinations, such as the error case below. The model generates distance values in the place the should be angle values (and vice versa).

- H 0.00 0.00° 0.00° N 1.01 1.57° 0.00° H 1.70 2.14° 0.00° C 2.06 1.13° -0.48° O 3.13 1.34° -0.42° N 2.49 0.62° -0.87°
C 2.94 0.10° -1.64° H 3.20 0.39° **3.91** H 2.81 0.33° -3.14° C 4.43 0.10° -1.77° H 5.15 0.22° 0.22° H 2.78 0.21° -0.95°
C **1.86** 1.86° **4.84** H **0.19**° 1.89° -2.06° H 8.28 1.94° -1.69° H 6.24 1.71° **5.70** C 3.97 0.67° -0.90° H 5.02 0.68° -0.77°
H 4.15 0.91° -1.11° C 2.93 0.50° **6.97** H 3.54 0.42° 0.42° H 2.74 0.88° 0.88°

These error cases will be rarer if the model well converges. When trained for 150 epochs, the model would generate $\sim 15\%$ of invalid samples, including the above discussed syntax problems. When trained for 250 epochs, the model would generate $< 2\%$ of invalid samples.

F. Visualization Results

F.1. Visualization of Generated Molecules

In this section, we provide visualizations of molecules generated from Geo2Seq with Mamba conditionally on the property of Polarizability α in Figure 3. The Polarizability of a molecule is the tendency to acquire an electric dipole moment when the molecule is subject to an external electric field. Large α values usually correspond to less isometrically molecular geometries. This is consistent with our generated examples.

In addition, we provide visualizations of molecules generated from Geo2Seq with Mamba trained on QM9 and DRUG in Figure 4 and Figure 5, respectively. These examples are randomly generated without any cherry pick. From the figures, we can see that the model can generate realistic molecular geometries for both small and large size molecules. However, similar to previous methods (Hoogetboom et al., 2022; Xu et al., 2023a), there are disconnected components, especially for larger molecules. A possible future direction is to apply fragment-based methods to reduce the sequence length, thus benefiting the training of language models.

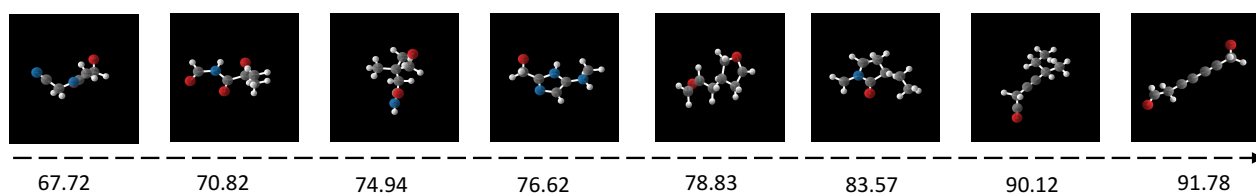


Figure 3: Visualization of generated molecules condition on the property of Polarizability α .

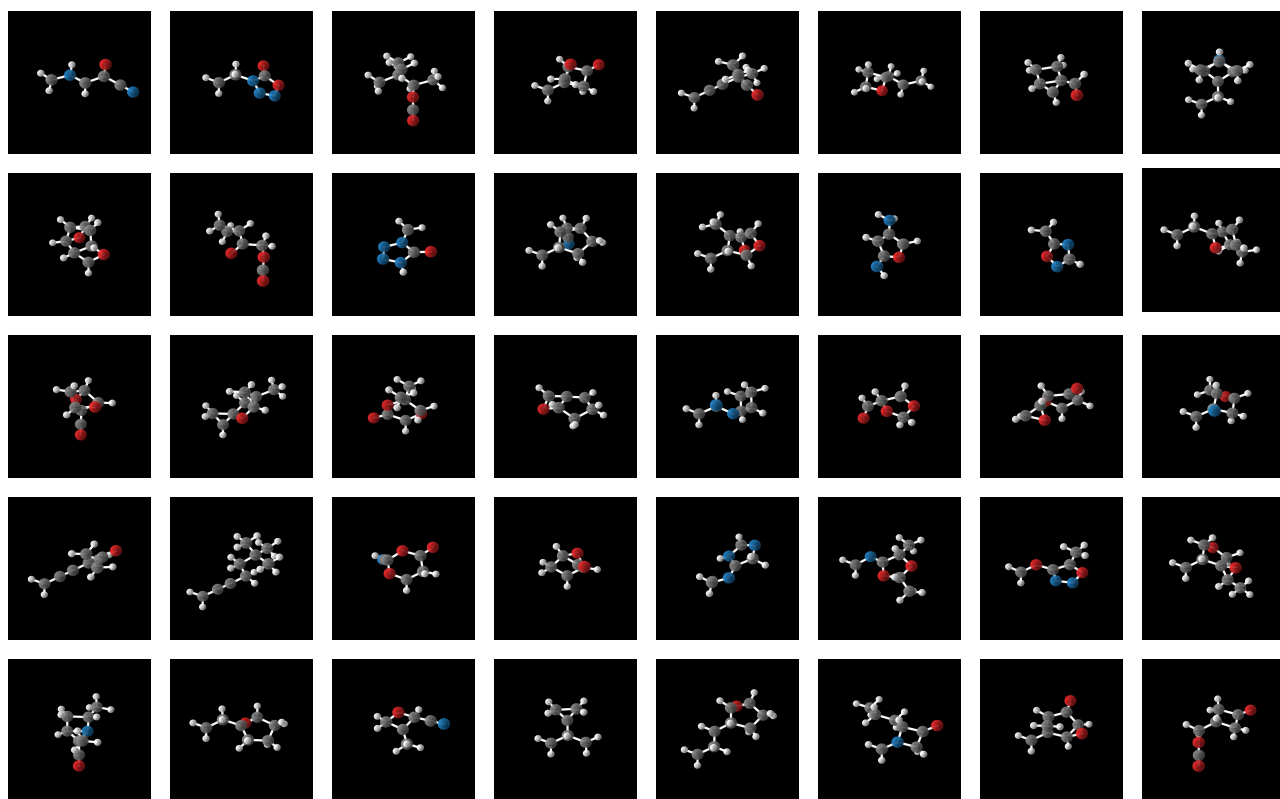


Figure 4: Visualization of molecules generated from Geo2Seq with Mamba trained on QM9.

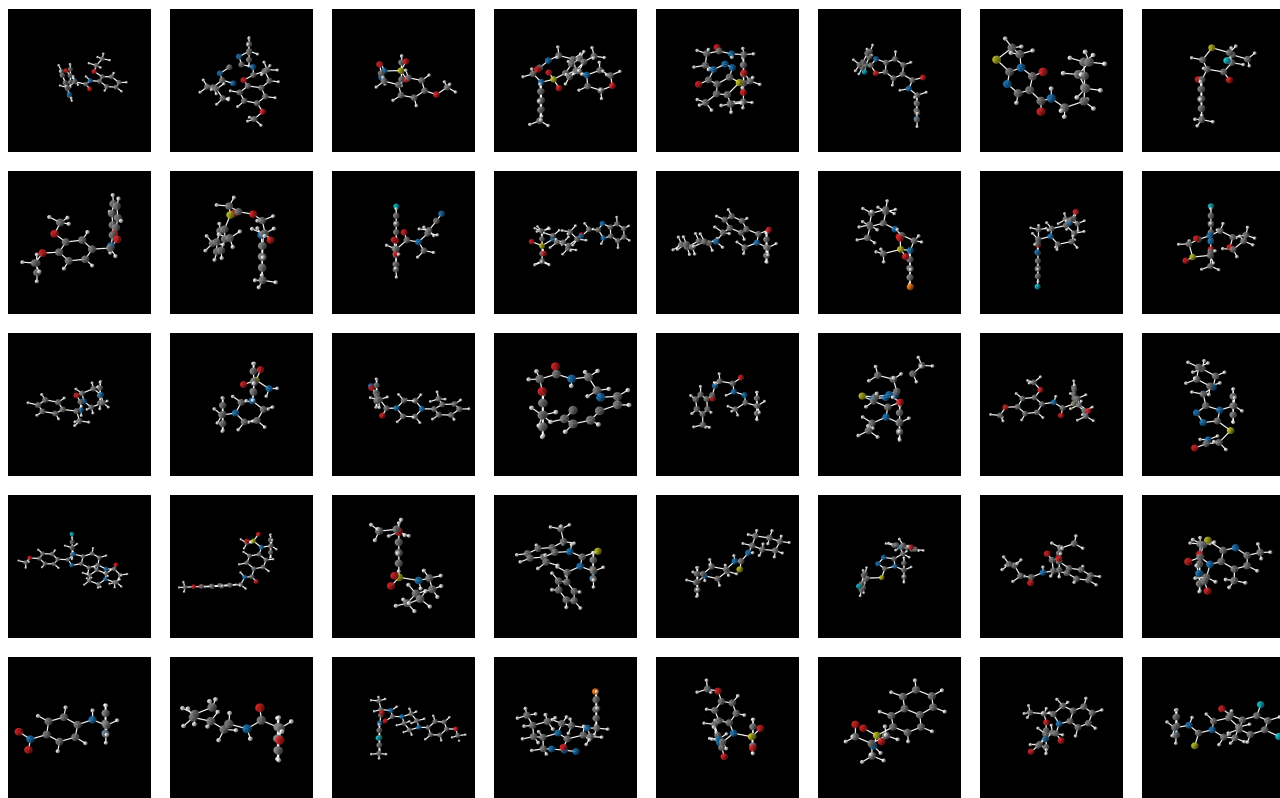


Figure 5: Visualization of molecules generated from Geo2Seq with Mamba trained on GEOM-DRUGS.

F.2. Visualization of Learned Token Embeddings

In this section, we provide UMAP visualizations of different (atom type, distance, and angle) token embeddings learned by Mamba models trained on QM9 and GEOM-DRUGS datasets. Patterns of the embeddings indicate that the model has successfully learned structure information from the sequence data, showcasing LMs’ capabilities to understanding molecules precisely in 3D space. For example, Figure 8 shows that similar angle tokens (*e.g.*, ‘1.41°’ and ‘1.42°’) are placed next to each other and the overall structure of all angles is a loop. Further, π -out-of-phase angles are placed near each other, such as ‘3.14°’, ‘-3.14°’, and ‘0°’. For atom type tokens, the model appears to capture the structure of the periodic table, although the rows and columns are not perfect in Figure 6. One reason is the limited atom types in the datasets (5 in QM9 and 16 in GEOM-DRUG), limiting the model’s capabilities to learn chemical patterns from the entire periodic table. We provide analyses of the visualization results in the caption of each figure as Figure 6 - Figure 10.

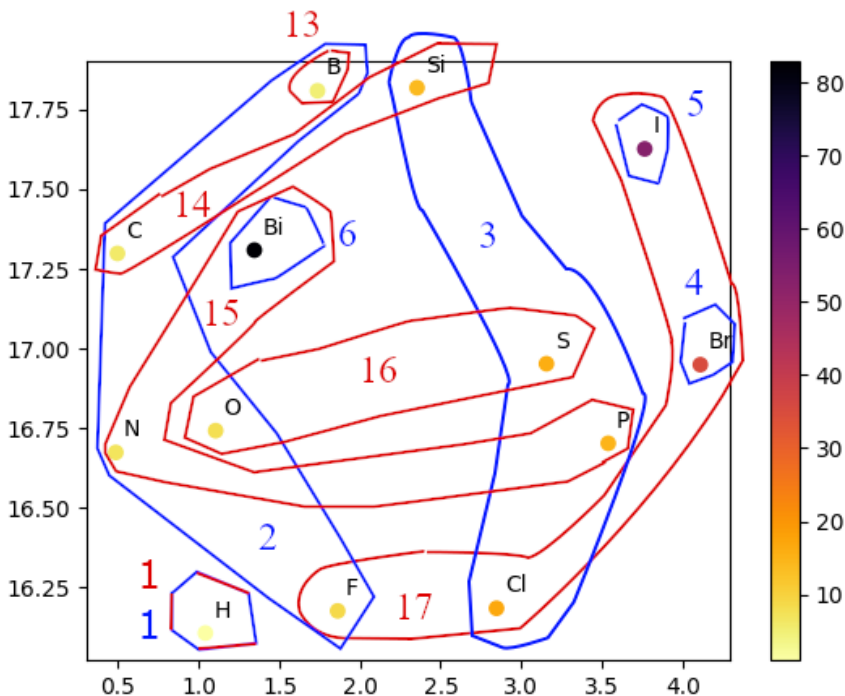


Figure 6: UMAP visualization of element token embeddings learned by a Mamba model trained on GEOM-DRUGS. Red groups indicate columns in the periodic table and blue groups indicate rows, which are both numbered. Points are colored by atomic weight. Overall, the model appears to capture the structure of the periodic table. The column generally increases from top to bottom, and the row generally increases from left to right.

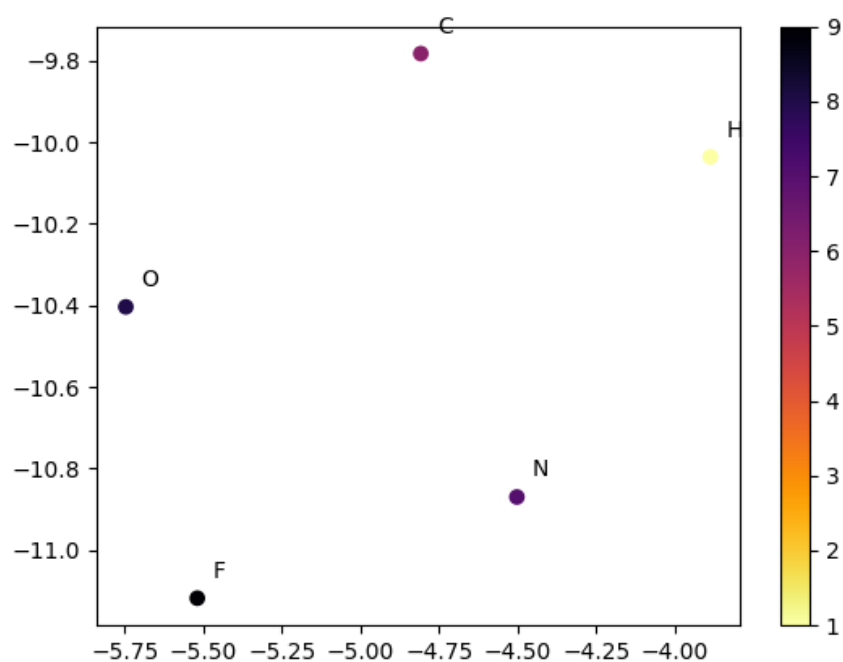


Figure 7: UMAP visualization of element token embeddings learned by a Mamba model trained on QM9. Points are colored by atomic weight. Overall, the model appears to distinguish well between different elements. All different elements are distributed distantly from each other in the embedding space.

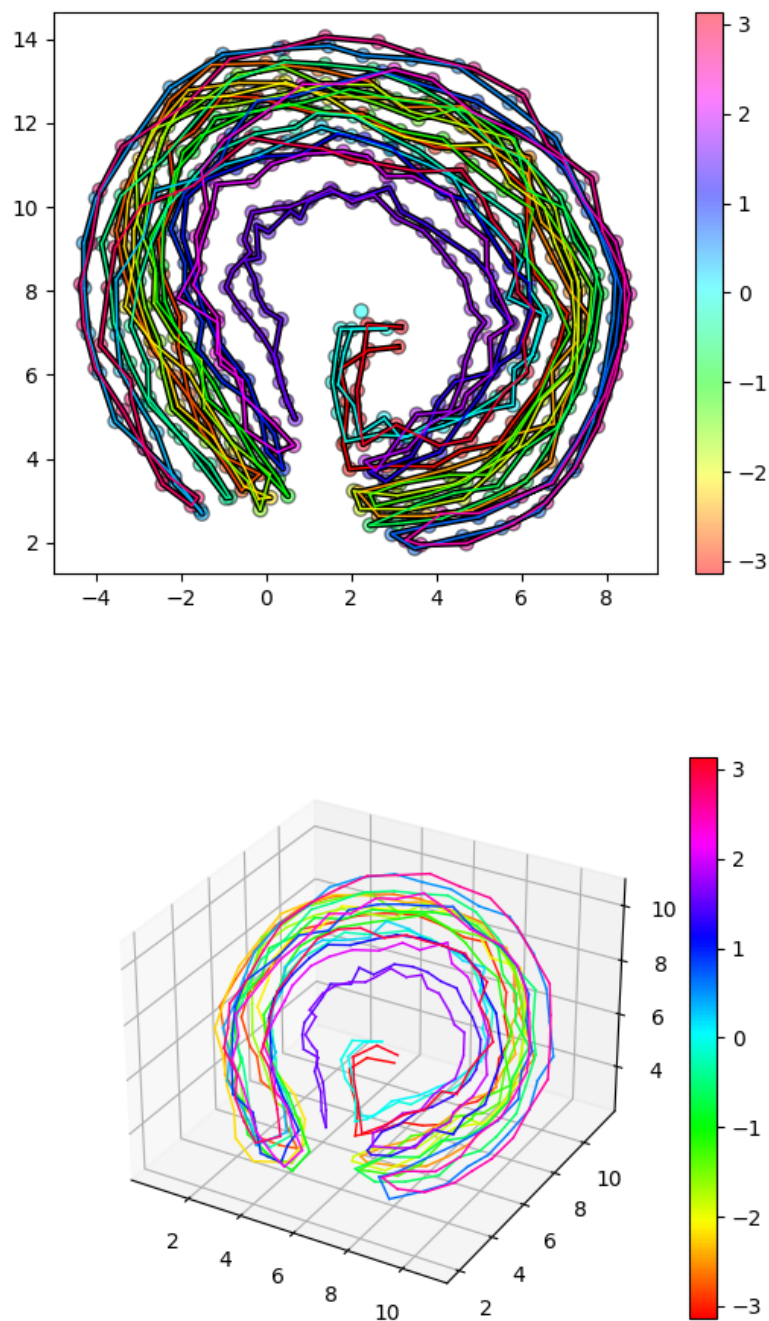


Figure 8: 2D and 3D UMAP visualization of angle token embeddings learned by a Mamba model trained on GEOM-DRUGS. It can be observed that similar tokens (e.g., ‘1.41°’ and ‘1.42°’) are placed next to each other and the overall structure is a loop. Further, π -out-of-phase angles are placed near each other, such as ‘3.14°’, ‘-3.14°’, and ‘0°’.

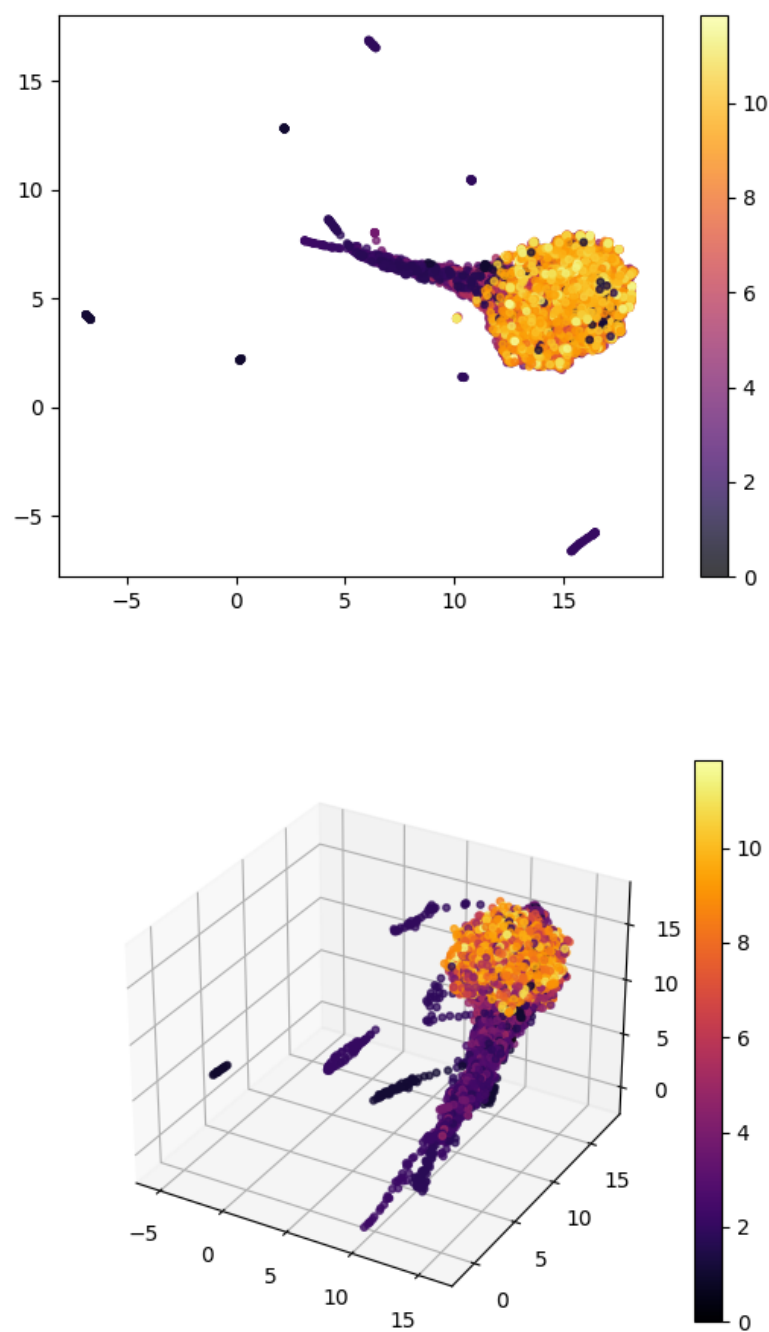


Figure 9: 2D and 3D UMAP visualization of distance token embeddings learned by a Mamba model trained on QM9. Representations of distances lower than 6 form relatively distinct patterns. This is likely because these values are much more frequently seen in the training data. Values over 20 cluster into a clump, suggesting that they are also recognized by the model.

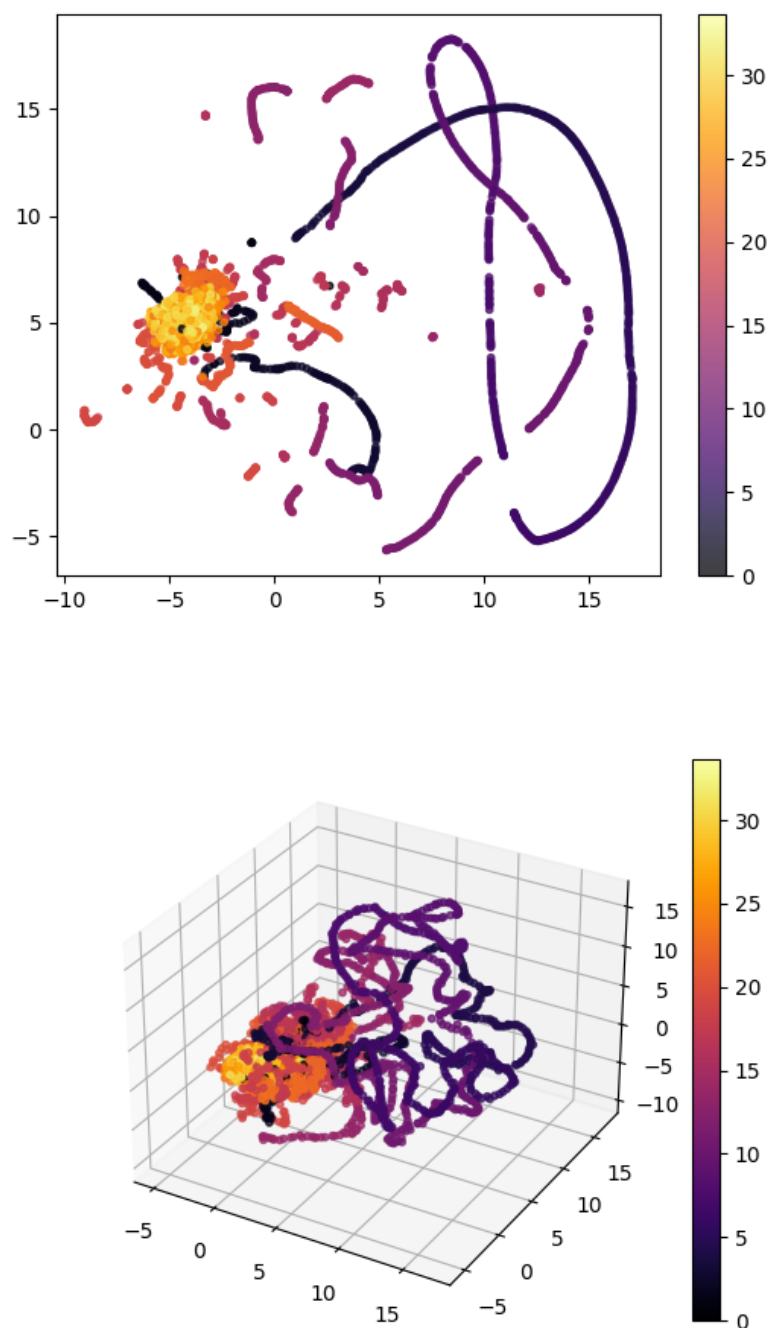


Figure 10: 2D and 3D UMAP visualization of distance token embeddings learned by a Mamba model trained on GEOM-DRUGS. It is notable that the best and most distinct representations seem to arise from between 5 and 20. This is likely because these values are much more frequently seen in the training data. Values over 20 form an indistinct clump. Interestingly, values > 20 are near values < 3 , which is initially unintuitive; however, they are likely placed in a similar location in the embedding space since both small and large distances are rarely seen in the data.