# Enhancing Generalizability in Molecular Conformation Generation with METRIZATION-Informed Geometric Diffusion Pretraining

**Xiaozhuang Song[1], Yuzhao Tu[1], Hangting Ye[2], Wei Fan[3],**
**Qingquan Zhang[4], Xiaoxue Wang[5], Tianshu Yu[1],**

[1]The Chinese University of Hong Kong, Shenzhen,
[2]Jilin University,
[3]University of Oxford,
[4]Southern University of Science and Technology,
[5]ChemLex Technology Co., Ltd.
xiaozhuangsong1@link.cuhk.edu.cn, yutianshu@cuhk.edu.cn

## Abstract

Diffusion-based generative models have recently excelled in generating molecular conformations but struggled with the generalization issue – models trained on one dataset may produce meaningless conformations on out-of-distribution molecules. On the other hand, distance geometry serves as a generalizable tool for the traditional computational chemistry methods of molecular conformation, which is predicated on the assumption that it is possible to adequately define the set of all potential conformations of any non-rigid molecular system using purely geometric constraints. In this work, we for the first time explicitly incorporate distance geometry constraints into pretraining phase of diffusion-based molecular generation models to improve the generalizability. Inspired by the classical distance geometry solution designed for solving the molecular distance geometry problem, we propose MiGDiff, a Metrization-Informed Geometric Diffusion framework. MiGDiff injects distance geometry constraints by pretraining the deep geometric diffusion backbone within the Metrization sampling approach, yielding a "Metrization-driven pretraining + Data-driven finetuning" paradigm. Experimental results demonstrate that MiGDiff outperforms state-of-the-art methods and possesses strong generalization capabilities, particularly on generating previously unseen molecules, revealing the vast untapped potential of combining traditional computational methods with deep generative models for 3D molecular generation.

**Code** — https://github.com/ShawnKS/MIGDiff

## Introduction

The generation of 3D molecular conformations is a cornerstone in various scientific domains (Wang et al. 2023a), including drug discovery (Anderson 2003; Ingraham et al. 2023) and material science (Rosei et al. 2003; Anstine and Isayev 2023). Traditionally, two primary methodologies have been employed to tackle this challenge: data-driven generative methods (Guo et al. 2024) and computational methods grounded in quantum chemistry principles (Havel and Wüthrich 1985; Kuszewski, Nilges, and Brünger 1992; An and Tao 2003; Sit, Wu, and Yuan 2009; Sit and Wu 2011).

Each approach has its own set of advantages and limitations. Pure data-driven generation methods have recently gained prominence due to their remarkable ability to capture the overall structural distribution of molecules (Gómez-Bombarelli et al. 2018; Fu et al. 2021; Hoffman et al. 2022), wherein diffusion-based generative models stand out with exemplary performance (Xu et al. 2022; Guo et al. 2024; Zhou et al. 2024; Wang et al. 2023b). The high-level idea of diffusion-based molecular generation methods undergo a transition from their stable equilibrium conformations to a state of increased disorder through a series of controlled diffusion steps, and then learn a model to reverse the diffusion process (Song and Ermon 2019; Ho, Jain, and Abbeel 2020; Xu et al. 2022). These methods leverage vast datasets to learn and predict molecular structures (Ryan, Lengyel, and Shatruk 2018; Ross et al. 2022; Siebenmorgen et al. 2024). However, this reliance on pre-existing data poses a significant limitation: generated molecules tend to bear a strong resemblance to training samples, resulting in a lack of generalization when applied to novel, unseen molecular structures (Keith et al. 2021; Dou et al. 2023; van Tilborg et al. 2024). Conversely, traditional computational methods, particularly those based on quantum chemistry and related mathematical tools (Havel and Wüthrich 1985; An and Tao 2003; McArdle et al. 2020), are celebrated for their robust generalization capabilities. Distance geometry is the mathematical basis for a geometric theory of computational methods (Havel 1998; Riniker and Landrum 2015). Specifically, the problem involves determining a 3-dimensional spatial structure that meets the given upper and lower bounds on distances between points. This problem has been proven to be *strong NP-hard* (Saxe 1979). To solve this problem, the METRIZATION algorithm (Kuszewski, Nilges, and Brünger 1992; Spellmeyer et al. 1997) was proposed as an efficient distance relationships sampling method given the bounds. By combining the METRIZATION algorithm with Schoenberg's theorem (Schoenberg 1935), it enables the embedding of bounds into 3-dimensional spatial coordinates. These methods can adeptly handle a wide array of molecular structures, thanks to their foundation in fundamental physical principles. Nevertheless, they are often criticized for their low solution efficiency and an inherent difficulty in capturing the holistic structure of complex molecules (Heinen et al. 2020).
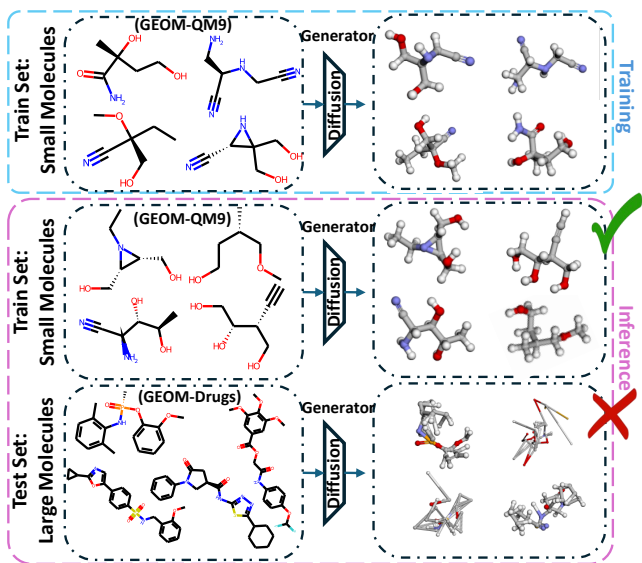
Figure 1: The generalization issue of current molecular conformation diffusion models.

In response to the aforementioned challenges, we propose a novel training paradigm for diffusion-based molecular conformation generation named METRIZATION-Informed Geometric Diffusion (MIGDIFF). This method aims to synthesize the strengths of both data-driven and physical rule-driven paradigms, via injecting physical constraints into deep geometric diffusion models by pretraining the geometric backbone within the METRIZATION sampling approach. Experimentally, MIGDIFF demonstrates strong generalization performance. The results indicate that diffusion models trained solely on real-world datasets perform inadequately when generating conformations from unseen distributions, and are highly constrained by the quality and quantity of the available data. In contrast, geometric diffusion models enhanced with METRIZATION pretraining can still generate excellent and stable conformation structures even without supervision from the target data distribution. By integrating physical rules encoded into geometric constraints, this framework significantly enhances the generalization ability of the generation model, enabling it to produce diverse and accurate molecular conformations beyond the scope of the training data.

## Related Work

### Diffusion-based Molecular Generation

Diffusion models have gained significant traction as powerful tools for drug discovery applications (Guo et al. 2024; Xu et al. 2022; Hua et al. 2024; Zhou et al. 2024). Diffusion-based methods for generating molecular structures typically start by transitioning from stable equilibrium conformations to a state of heightened disorder via a sequence of regulated diffusion steps (Yang et al. 2023; Guo et al. 2024; Cao et al. 2024). Subsequently, a model is trained to reverse this diffusion process (Xu et al. 2022). Recent studies have explored the integration of equivariant graph models within diffusion models for molecule generation (Xu et al. 2022; Hoogeboom

et al. 2022). For example, (Xu et al. 2022) introduces a diffusion model with an equivariant GNN, enabling the model to work jointly on atom features and coordinates (Gilmer et al. 2017; Schütt et al. 2017). Overall, incorporating equivariant graph models into diffusion models for molecule generation has demonstrated promising results in recent research (Satorras, Hoogeboom, and Welling 2021; Han et al. 2022).

### Molecular Distance Geometry Optimization

Distance Geometry is the study of geometry with the basic entity being distance (Schoenberg 1935; Alfakih, Khandani, and Wolkowicz 1999; Liberti et al. 2014; Liberti and Lavor 2016) and is the mathematical basis for a geometric theory of molecular conformation (Havel 1998; Kuntz 1992; Liberti et al. 2011; Trinajstic 2018). The awarding of the Nobel Prize to Wüthrich for the application of Nuclear Magnetic Resonance (NMR) techniques in protein analysis has significantly highlighted the role of Distance Geometry in the field of structural bioinformatics (Havel and Wüthrich 1985). Various algorithms have emerged through the design of combining different structures and different optimization methods (Saxe 1979; An and Tao 2003; Sit and Wu 2011; Abuter et al. 2019; Meng et al. 2020; McArdle et al. 2020). Among them, the METRIZATION algorithm, combined with the postprocess design of coordinate refinement (Havel 1998), has been widely applied in modern computational libraries (Landrum et al. 2006; Riniker and Landrum 2015; Liberti and Lavor 2018; Sobez and Reiher 2020).

## Preliminaries

### Problem Formulation and Notations

**Notations.** In this paper, each molecule with $n$ atoms is represented as an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{v_i\}_{i=1}^n$ is the set of vertices representing atoms and $\mathcal{E} = \{e_{i,j} \mid (i,j) \subseteq |\mathcal{V}| \times |\mathcal{V}|\}$ is the set of edges representing inter-atomic bonds. Each node $v_i \in \mathcal{V}$ describes the atomic attributes, e.g., the element type. Each edge $e_{i,j} \in \mathcal{E}$ describes the corresponding connection between $v_i$ and $v_j$, and is labeled with its chemical type. In addition, we also assign the unconnected edges with a virtual type. The position of each atom $\mathcal{V}$ is represented by a coordinate vector $x \in \mathbb{R}^3$ in the 3-dimensional space, and the full set of positions (i.e., the conformation) can be represented as a matrix $C = [x_1, x_2, \ldots, x_n] \in \mathbb{R}^{n \times 3}$. Distance bounds for all pairs of atoms in the molecule are determined and represented in distance bounds $u_{ij}$ and $l_{ij}$, respectively, specifying the lower bound and the upper bound of the distance between atoms $v_i$ and $v_j$. These bounds are derived from empirical information, including ideal bond lengths, bond angles, and torsion angles.

**Problem Definition.** The task of molecular conformation generation is a conditional generative problem, where we are interested in generating stable conformations for a provided molecular graph $\mathcal{G}$. Given multiple graphs $\mathcal{G}$, and for each $\mathcal{G}$ given its conformations $C$ as independent and identically distributed (i.i.d) samples from an underlying Boltzmann distribution, our goal is learning a generative model $p_\theta(C \mid \mathcal{G})$

that allows for easy sampling and approximates the Boltzmann function.

## Molecular Distance Geometry Basics

**Schoenberg's Theorem.** An $n \times n$ symmetric matrix $\mathbf{D} = (d_{ij})$ is a valid Euclidean distance matrix (EDM) of $n$ points $X = \{x_1, \ldots, x_n\} \subset \mathbb{R}^r$ if and only if the Gram matrix $\mathbf{G} = -\frac{1}{2}\mathbf{J}\mathbf{D}^2\mathbf{J}$ is positive semi-definite (PSD) of rank $r$, where $\mathbf{J} = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^T$. Given the eigendecomposition $\mathbf{G} = \mathbf{P}\Lambda\mathbf{P}^\top$, where $\Lambda$ is a diagonal matrix of eigenvalues and $\mathbf{P}$ is the matrix of corresponding orthonormal eigenvectors, a realization of the points $X$ in $\mathbb{R}^r$ can be obtained as $X = \mathbf{P}\sqrt{\Lambda}$. This theorem is widely used in distance geometry computation problems (Havel 1998; Sit, Wu, and Yuan 2009; Sit and Wu 2011).

For $x_1, \ldots, x_l$ and $x_{l+1}$ as the coordinate vectors of atoms $1, \ldots, l+1$, if the distances $d_{i,j}$ are given, then $\|x_i - x_j\| = d_{i,j}$ for all $i, j = 1, \ldots, l+1$, and $\|x_i\|^2 - 2x_i^\top x_j + \|x_j\|^2 = d_{i,j}^2$, $i, j = 1, \ldots, l+1$. Setting the origin at $x_{l+1} = (0,0,0)^\top$, we get $\|x_i\| = d_{i,l+1}$ and

$$d_{i,l+1}^2 - 2x_i^\top x_j + d_{j,l+1}^2 = d_{i,j}^2, \quad i, j = 1, \ldots, l. \quad (1)$$

$$\mathbf{G} = \left\{ \frac{d_{i,l+1}^2 - d_{i,j}^2 + d_{j,l+1}^2}{2} : i, j = 1, \ldots, l \right\}. \quad (2)$$

Let $\{d_{i,j} : i, j = 1, \ldots, l+1\}$ be a set of distances in $\mathbb{R}^k$, for some $k < l$, $\mathbf{G}$ has maximum rank $k$. For the singular value decomposition $\mathbf{G} = U\Sigma U^\top$, according to Schoenberg's Theorem, if $\mathbf{G}$ is a matrix of rank $\leq k$, then $X = V\Sigma^{1/2}$ solves $XX^\top = \mathbf{G}$, where $V = U(:, 1:k)$ and $\Sigma = \Sigma(1:k, 1:k)$. Note that distances may have errors since it is challenging to randomly generate a distance matrix that perfectly aligns with the embedding dimension. As such, the distance matrix $\tilde{\mathbf{D}}$ obtained through METRIZATION is conventionally regarded as a realization of an approximation (Havel 1998). Though $\tilde{\mathbf{D}}$ might not be a valid Euclidean distance matrix and thus $\tilde{\mathbf{G}} = -\frac{1}{2}J\tilde{\mathbf{D}}^2 J$ may have negative eigenvalues, we can compute its positive eigenvalues $\lambda_1, \ldots, \lambda_H$ and corresponding eigenvectors to recover an approximate realization $x \in \mathbb{R}^H$ of $\tilde{\mathbf{D}}$.

Then, a buildup procedure can be implemented as follows. First, construct a Gram matrix $\mathbf{G}$ from the distances $d_{i,j}$ among $l+1$ atoms. Next, compute the singular value decomposition of $\mathbf{G}$, such that $\mathbf{G} = U\Sigma U^\top$. From this, obtain $X = \mathbf{P}\sqrt{\Lambda}$, where $\mathbf{P} = U(:, 1:3)$ and $\Lambda = \Sigma(1:3, 1:3)$. Set the coordinates of atom $l+1$ at $(0,0,0)^\top$. This procedure can be applied to any set of $l+1$ atoms in $k$-dimensional Euclidean space, making it a standard algorithm for solving distance geometry problems when the distances for all pairs of atoms are given. For a rigorous introduction and proof involving Schoenberg's Theorem for 3-dimensional coordinate embeddings, one might refer to (Sit, Wu, and Yuan 2009; Sit and Wu 2011; Liberti and Lavor 2016).

## Methodology

The MIGDIFF framework enhances molecular geometry prediction through two main components: METRIZATION-based molecular geometry sampling and a Deep Geometric Diffusion Network. The METRIZATION component generates accurate molecular geometries, ensuring compliance with physical and chemical constraints. These geometries provide a robust foundation for the Deep Geometric Diffusion Network, which undergoes a pretrain-finetune process. In the pretraining phase, the network learns from the METRIZATION-generated data. It is then fine-tuned using real-world datasets with positional records, enhancing the network's ability to accurately predict molecular structures. Following the introduction of these components, we will detail their organization and training within the MIGDIFF pipeline.

### Metrization-based Molecular Geometry Sampling

**Metrization.** The problem that the METRIZATION algorithm (Kuszewski, Nilges, and Brünger 1992) attempts to solve can be summarized as:

*Given the independent upper and lower bounds for each atom pair in a molecule, how to generate atomic distance relationships that comply with global geometric constraints?*

This problem has been proven to be *strongly NP-hard* (Saxe 1979). The intuitive design like parallel sampling each edge and then obtaining atomic distance relationships that meet the bounds could be highly inefficient and ineffective, since a determined edge will alter the edges of adjacent triangular relationships, which then propagate to affect the global structural bounds. Therefore, METRIZATION employs the following design to achieve simple and efficient sampling. In the METRIZATION process, one of the distances $d_{ij}$ is randomly set to a value between its lower and upper limits, $d_{ij} \in [l_{ij}, u_{ij}]$. The bounds are then updated to this value, making $l_{ij} = u_{ij} = d_{ij}$. Using these modified bounds as input, the triangle inequality limits are re-computed. This procedure is repeated for each distance $d_{ij}$, iterating until the lower and upper triangle inequality limits converge to the desired distance matrix that satisfies both the triangle inequalities and the original bounds (Havel 1998). Then, using aforementioned Schoenberg's solution (Schoenberg 1935), we calculate the Gram matrix of the sampled distance matrix from METRIZATION method as $\tilde{\mathbf{G}} = -\frac{1}{2}J\tilde{\mathbf{D}}^2 J$. Since $\tilde{\mathbf{G}}$ is a PSD matrix of rank $k$, we can obtain the coordinates in the corresponding $k$-dimensional space as $X = \tilde{\mathbf{P}}\sqrt{\tilde{\Lambda}}$. This results in $X$ being the coordinates in a Euclidean space that best preserve the distances given by $\tilde{\mathbf{D}}$. Due to the approximate nature of $\tilde{\mathbf{G}}$, $\tilde{\Lambda}$ may have more than 3 eigenvalues, including negative ones. The solution is to keep the largest 3 eigenvalues and set the negative ones to zero.

**Conjugate Gradient Optimization for Coordinate Refinement.** The coordinates obtained from metrization and embedding often exhibit significant constraint violations, particularly in covalent bond lengths, where even minor deviations can have severe energetic consequences. These violations are corrected by minimizing an error function that measures the total violation of both distance and chirality constraints.

Following (Havel 1998), we use the following error function $E_d(C)$ to refine the coordinates:

$$E_d(C) = \sum_{\{i,j\}} \max \left( 0, \frac{(\|x_i^t - x_j^t\|^2 - u_{ij}^2)}{\epsilon_u + u_{ij}^2} \right)^2$$
$$+ \sum_{\{i,j\}} \max \left( 0, \frac{(l_{ij}^2 - \|x_i^t - x_j^t\|^2)}{\epsilon_l + \|x_i^t - x_j^t\|^2} \right)^2,$$

where $x_i^t$ and $x_j^t$ are the coordinates of atoms $i$ and $j$, respectively, and $\epsilon_l$, $\epsilon_u$ are parameters to ensure numerical stability.

## Geometric Diffusion Network

In this section, we introduce the backbone geometric diffusion network used for training. Our backbone network design is primarily based on GeoDiff (Xu et al. 2022).

**Forward Process.** The deep geometric diffusion network for molecular generation undergoes a transition from their stable equilibrium conformations $C^0$ to a state of increased disorder through a series of controlled diffusion steps. A common practice for this progression is simulated by a forward diffusion process modeled as a Markov chain (Xu et al. 2022), where variance parameters $\beta_t$ dictate the scale of noise injection at each step:

$$q(C^{1:T}|C^0) = \prod_{t=1}^{T} q(C^t|C^{t-1}),$$
$$q(C^t|C^{t-1}) = \mathcal{N}(C^t; \sqrt{1 - \beta_t}C^{t-1}, \beta_t I).$$

This ensures a gradual increase in randomness, culminating in a highly chaotic state akin to a white noise distribution after $T$ iterations.

**Reverse Process.** Given molecular graphs $\mathcal{G}$, a reverse generative process is implemented to reconstruct conformations $C^0$ from white noise $C^T$. This procedure, beginning with noisy particles $C^T \sim p(C^T)$, functions as the inverse of the previously described diffusion dynamics. Formally, the transition probabilities between states in this reverse trajectory are defined through a learnable conditional Markov chain as:

$$p_\theta(C^{0:T-1}|\mathcal{G}, C^T) = \prod_{t=1}^{T} p_\theta(C^{t-1}|\mathcal{G}, C^t),$$
$$p_\theta(C^{t-1}|\mathcal{G}, C^t) = \mathcal{N}(C^{t-1}; \mu_\theta(\mathcal{G}, C^t, t), \sigma_t^2 I).$$

where the means are estimated through $\mu_\theta$, and $\sigma_t$ represents a variance that can be freely defined. Thus, starting from a standard Gaussian distribution $p(C^T)$, the 3D structure generation for molecular graph $\mathcal{G}$ proceeds through two phases: initially drawing chaotic particles $C^T$ from $p(C^T)$, then progressively refining them via reverse Markov kernels $p_\theta(C^{t-1}|\mathcal{G}, C^t)$ in an iterative manner.

## Roto-translational Equivariant Network

The essence of geometric diffusion in molecular generation is to establish a density $p_\theta(C^0)$ that is invariant to rotations and translations, allowing the construction of a roto-translationally equivariant $\epsilon_\theta$ for predicting the noise to reconstruct the conformations. We adopt the design of the recent
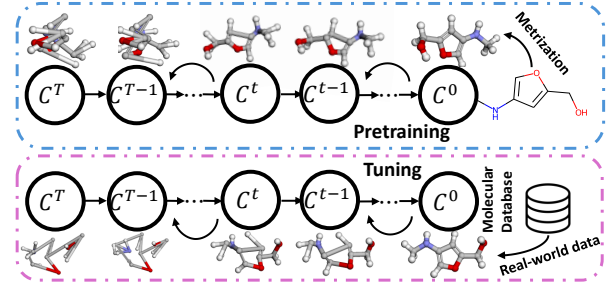


Figure 2: An overview of the model training for MIGDIFF

equivariant network model (Xu et al. 2022) to create an equivariant convolutional network. At each layer $l$, the network processes node embeddings $\mathbf{h}^l$ and spatial coordinates $x^t$, then generates updated embeddings $\mathbf{h}^{l+1}$ and coordinates $x^{t+1}$, which can be expressed as:

$$\tau_{ij}^{l+1} = \Phi_m \left( \mathbf{h}_i^l, \mathbf{h}_j^l, \|x_i^l - x_j^l\|^2, \tau_{ij}^l; \theta_m \right), \quad (3)$$
$$\mathbf{h}_i^{l+1} = \Theta_h \left( \mathbf{h}_i^l, \tau_i^{l+1}; \theta_h \right), \quad (4)$$
$$x_i^{t-1} = \sum_{j \in \mathcal{N}(i)} \frac{1}{d_{ij}^t}(x_i^t - x_j^t)\Phi_x(\tau_{ij}^{l+1}; \theta_x), \quad (5)$$

where $\Phi_m$, $\Phi_h$ and $\Phi_x$ are feed-forward networks and $d_{ij}$ denotes interatomic distances. $\tau_{ij}$ represents the hidden representation of edge between atom $i$ and atom $j$. $\mathbf{h}_i^{l+1}$ is the atom embedding of atom $i$ at layer $l + 1$ and $x_i^{t-1}$ is the coordinates at step $t-1$ for atom $i$. The embedding $\tau$ inherits its invariance properties through the invariant feature $\mathbf{h}$ and the equivariant coordinate $x$. Since $\tau$ is constructed exclusively from these invariant elements, it maintains invariance throughout the transformation. The transformation of $x$ incorporates spatial differences weighted by invariant attributes, thereby preserving translation invariance and rotation equivariance. Thus, the architecture of $\epsilon_\theta$ with $L$ successive layers guarantees equivariance with respect to $C^t$ through compositional structure (Satorras, Hoogeboom, and Welling 2021; Xu et al. 2022).

Followed the design of (Xu et al. 2022), we empirically employ two types of graph convolutional neural networks for modeling different structural information: a message passing neural network (MPNN) (Gilmer et al. 2017) that aggregates features from bond-linked atoms, and SchNet (Schütt et al. 2017), which aggregates features based on edge features and intra-atomic distances. These two types of graph modeling can be written as:

$$\mathbf{h}_i'^{l+1} = \sum_{j \notin \mathcal{N}(i)} \mathbf{W}\tau_{ij}^l d_{ij}^t v_j, \quad (6)$$
$$\mathbf{h}_i''^{l+1} = \sum_{j \in \mathcal{N}(i)} \tau_{ij}^l v_j, \quad (7)$$

and the update for coordinates $x^{l+1}$ can be written as:

$$x_i^{t-1} = \sum_{j \in \mathcal{N}(i)} \frac{1}{d_{ij}^t}(x_i^t - x_j^t)\Phi_x(\tau_{ij}'^{\,l+1}; \theta_x) +$$
$$\sum_{j \notin \mathcal{N}(i)} \frac{1}{d_{ij}^t}(x_i^t - x_j^t)\Phi_x(\tau_{ij}''^{\,l+1}; \theta_x),$$

where $\mathbf{h}'$ and $\mathbf{h}''$ represent the node embeddings from MPNN and SchNet, respectively. $\tau_{ij}'^{\,l+1}$ and $\tau_{ij}''^{\,l+1}$ are the hidden representations of edge $ij$ from MPNN and SchNet obtained from Eq. (3) using the node and edge embeddings.

## Model Training

**Training Objective**   As the log-likelihood $\mathbb{E}[\log p_\theta(C^0|\mathcal{G})]$ is intractable, a common practice is to maximize its evidence lower bound (ELBO) as:

$$\mathbb{E}[\log p_\theta(C^0|\mathcal{G})] = \mathbb{E}\left[\log \frac{p_\theta(C^{0:T}|\mathcal{G})}{q(C^{1:T}|C^0)}\right]$$
$$\geq -\mathbb{E}_q\left[\sum_{t=1}^{T} D_{\mathrm{KL}}(q(C^{t-1}|C^t, C^0)\|p_\theta(C^{t-1}|C^t, \mathcal{G}))\right],$$

where $q(C^t|C^0)$ for any timestep $t$ can be expressed as $q(C^t|C^0) = \mathcal{N}(C^t; \sqrt{\bar{\alpha}_t}C^0, (1-\bar{\alpha}_t)I)$, $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$. This leads to an analytically tractable $q(C^{t-1}|C^t, C^0)$. Building upon the insights from (Ho, Jain, and Abbeel 2020; Xu et al. 2022), the ELBO of the diffusion model can be further simplified by expressing the KL divergence between Gaussians as a weighted $L_2$ distance between $\epsilon_\theta$ and $\epsilon$. Thus, we can independently sample chaotic conformations at various timesteps from $q(C^{t-1}|C^t, C^0)$ to optimize the objective more efficiently and use $\epsilon_\theta$ to model the noise vector $\epsilon$ (Xu et al. 2022).

**Metrization Distance Geometry Pretraining**   The pretraining phase leverages METRIZATION-based distance geometry to initialize the diffusion network's understanding of molecular distance geometry structures. During this phase, *synthetic data* is generated using predefined bounds and constraints to form initial distance matrices. The goal is to train the network to respect geometric constraints and accurately predict distance relationships. By using METRIZATION, we ensure that the network learns to generate physically plausible molecular conformations even before being exposed to real-world molecular data. This step provides a robust foundation for subsequent training phases, enhancing the network's ability to handle complex molecular geometries.

**Real-world Molecular Structures Tuning**   Following the pretraining warm-up phase, the diffusion network undergoes fine-tuning using data from real-world molecular datasets. This stage adapts the METRIZATION-Informed diffusion network to actual molecular geometries. The tuning aims to refine the network's ability to accurately conform to real molecular geometries, enhancing its predictions. Simultaneously, the data-driven GNN model is also adjusted to improve its understanding and processing of molecular structures based on learned distance geometry principles. This comprehensive fine-tuning ensures that both theoretical knowledge and empirical data are effectively integrated, optimizing the network's capability to generate accurate and physically plausible molecular configurations.

**Conformation Sampling**   We can generate the stable conformation $C^0$ for a given graph $\mathcal{G}$ through the following steps with the trained reverse dynamics $\epsilon_\theta(\mathcal{G}, C^t, t)$ and transition mean $\mu_\theta(\mathcal{G}, C^t, t)$: The process begins with randomly sampling noisy states $C^T$ from distribution $p(C^T)$. Then, following a Markovian chain, we can iteratively draw samples $C^{t-1} \sim p_\theta(C^{t-1}|\mathcal{G}, C^t)$ as $t$ decreases from $T$ to 1. This sequential process gradually refines the particle positions from chaotic states to equilibrium states.

# Experiments

Following the setting used in previous studies (Shi et al. 2021; Xu et al. 2022), we empirically evaluate the performance of the proposed framework MIGDIFF on tasks: **molecular conformation generation** and **property prediction**. Our tests are conducted on both small molecules and drug-like compounds to ensure robustness and applicability across different chemical domains. Specifically, we further evaluate the molecular conformation generation capability in molecular conformation generation tasks, that is, to test whether the model learned from a specific data distribution can be generalized to another, which is an important issue in practical applications (Rosei et al. 2003; Gibson et al. 2009).

## Experiment Setup

**Datasets and Experimental Setup**   In line with previous studies (Xu et al. 2022), we utilize the following two datasets for performance evaluation: GEOM-QM9 for small molecules and GEOM-DRUGS for medium-sized organic compounds (Ramakrishnan et al. 2014; Axelrod and Gomez-Bombarelli 2022). The dataset split follows the setting used in (Shi et al. 2021), each dataset contains 200,000 conformations with 40,000 molecules, 5 conformations each. The valid split has 5,000 molecules with 5 conformations each. The test split has 200 molecules with 22,408 conformations for GEOM-QM9 and 14,324 conformations for GEOM-DRUGS.

For the MIGDIFF setup, both MPNNs and SchNets are implemented with 4 layers and a hidden embedding dimension of 128. For pretraining, we leverage METRIZATION-based molecular geometry sampling to sample 4 conformations for each molecules in the pretraining set. More experimental details can be found at supplementary files.

**Baselines**   We compare MIGDIFF with methods from both traditional computational methods and recent or established state-of-the-art deep learning baselines. For the deep learning approaches, we test the following models with highest reported performance: GEODIFF (Xu et al. 2022), GRAPHDG (Simm and Hernandez-Lobato 2020), CGCF (Xu et al. 2021b), CONFVAE (Xu et al. 2021a), GENMOL (Ganea et al. 2021) and CONFGF (Shi et al. 2021). we include 2 variant models in the evaluation: MIGDIFF-P, which only pretrains the diffusion backbone of MIGDIFF without any finetuning. GEODIFF-CR, which uses coordinate refinement as the post-process operation for GEODIFF. The vanilla
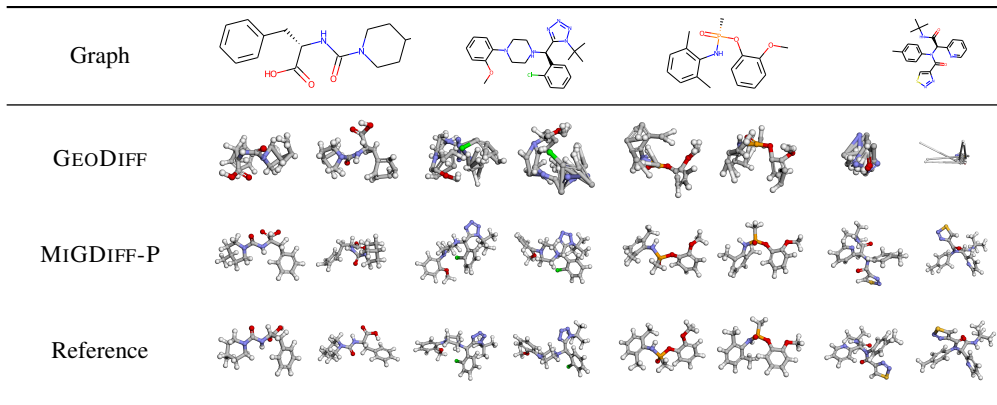
Table 1: Visualization results on the GEOM-DRUGS test set for the model trained without supervision from the GEOM-DRUGS training set.

| Models | COV-R (%) ↑ | | MAT-R (Å) ↓ | | COV-P (%) ↑ | | MAT-P (Å) ↓ | |
|---|---|---|---|---|---|---|---|---|
| | Mean | Median | Mean | Median | Mean | Median | Mean | Median |
| RDKIT | 60.91 | 65.70 | 1.2026 | 1.1252 | 72.22 | 88.72 | 1.0976 | 0.9539 |
| GEODIFF | 18.93 | 0.00 | 2.6392 | 2.6889 | 12.46 | 0.00 | 3.4616 | 3.4423 |
| GEODIFF-CR | 23.99 | 0.00 | 2.3199 | 2.4951 | 21.94 | 0.00 | 2.6264 | 2.5799 |
| GEODIFF + FF | 85.48 | 97.98 | 0.9404 | 0.9275 | 54.89 | 64.57 | 1.4651 | 1.3453 |
| METRIZATION | 45.71 | 40.72 | 1.3887 | 1.3292 | 44.93 | 34.90 | 1.4418 | 1.3391 |
| MIGDIFF-P | 73.51 | 80.85 | 1.0732 | 1.0501 | 58.34 | 62.85 | 1.2755 | 1.2358 |
| MIGDIFF-P + FF | **89.85** | **100.00** | **0.8056** | **0.7932** | **81.31** | **94.24** | **1.0374** | **0.9650** |

Table 2: Results on the GEOM-DRUGS test set for the model trained without supervision from the GEOM-DRUGS training set.

METRIZATION (Kuszewski, Nilges, and Brünger 1992) is also included as a method in the comparison. For MIGDIFF and MIGDIFF-P, we perform METRIZATION sampling for molecules present in the training datsets for pretraining.

## Conformation Generation

**Evaluation Metrics**   We evaluate the generation quality using the diversity and quality of the generated conformations with 4 RMSD-based metrics as outlined in (Xu et al. 2022; Ganea et al. 2021) The generated quality is evaluated by the coverage (COV) and matching (MAT) of the generated conformations and reference conformations (Kabsch 1976). Let $S_g$ and $S_r$ represent the sets of generated and reference conformations. The recall metrics of the coverage and matching scores can be defined as:

$$\text{COV-R}(S_g, S_r) = \frac{\left| \left\{ C \in S_r \mid \Delta(C, \hat{C}) \leq \delta, \hat{C} \in S_g \right\} \right|}{|S_r|},$$

(8)

$$\text{MAT-R}(S_g, S_r) = \frac{1}{|S_r|} \sum_{C \in S_r} \min_{\hat{C} \in S_g} \Delta(C, \hat{C}),$$

(9)

where $\delta$ represents a predefined threshold, and $\Delta$ denotes the RMSD (Root Mean Square Deviation). Following (Xu et al. 2021a; Ganea et al. 2021), $\delta$ is set to 0.5Å for GEOM-QM9 and 1.25Å for GEOM-DRUGS datasets. $S_g$ is set to twice the

size of $S_r$ for each molecule. The precision scores COV-P and MAT-P are calculated by swapping the positions of reference and generated sets in the above formulas. In general, the precision metrics focus on quality, while the Recall metrics emphasize diversity. COV scores measure the percentage of structures in one set covered by another, with covering meaning the RMSD between two conformations is within $\delta$. MAT scores measure the average RMSD of conformations in one set with its closest neighbor in another. Higher COV or lower MAT scores indicate more realistic conformations.

**Result Discussion**   A major focus of our experiments is to evaluate generalization ability for the deep geometric diffusion models on conformation generation. We test it by training GEODIFF within GEOM-QM9 datasets while using GEOM-DRUGS for generation evaluation, and compare it with MIGDIFF, the deep geometric diffusion model fully pretrained on METRIZATION generated conformations.

The results are shown in Tab. 1 and Tab. 2. We found that there are several interesting observations. First, we can observe that MIGDIFF-P performs significantly better compared with GEODIFF, while GEODIFF trained with small molecules struggled severely with generalization problems. The image in Tab.1 further illustrates this point: GEODIFF, trained only on GEOM-QM9, struggles to generate effective conformations on GEOM-DRUGS, whereas MIGDIFF produces relatively superior conformations. The results demonstrate the importance of incorporating the METRIZATION

| Models | COV-R (%) ↑ | | MAT-R (Å) ↓ | | COV-P (%) ↑ | | MAT-P (Å) ↓ | |
|---|---|---|---|---|---|---|---|---|
| | Mean | Median | Mean | Median | Mean | Median | Mean | Median |
| CVGAE | 0.00 | 0.00 | 3.0702 | 2.9937 | - | - | - | - |
| GRAPHDG | 8.27 | 0.00 | 1.9722 | 1.9845 | 2.08 | 0.00 | 2.4340 | 2.4100 |
| CGCF | 53.96 | 57.06 | 1.2487 | 1.2247 | 21.68 | 13.72 | 1.8571 | 1.8066 |
| CONFVAE | 55.20 | 59.43 | 1.2380 | 1.1417 | 22.96 | 14.05 | 1.8287 | 1.8159 |
| GEOMOL | 67.16 | 71.71 | 1.0875 | 1.0586 | - | - | - | - |
| CONFGF | 62.15 | 70.93 | 1.1629 | 1.1596 | 23.42 | 15.52 | 1.7219 | 1.6863 |
| GEODIFF | 89.13 | 97.88 | 0.8629 | 0.8529 | 61.47 | 64.55 | 1.1712 | 1.1232 |
| MIGDIFF | **89.76** | **97.97** | **0.8604** | **0.8478** | **63.61** | **66.36** | **1.1203** | **1.1062** |

Table 3: Results on the GEOM-DRUGS test set for the model trained with supervision from the GEOM-DRUGS training set.

knowledge into data-driven generative models, which can significantly boost the generalization performance. Secondly, while the generation results obtained through METRIZATION show only average performance, the MIGDIFF-P model, which is trained solely on this data, significantly outperforms models that are trained on METRIZATION data alone. This superior performance indicates that leveraging METRIZATION data can lead to the development of more robust models. These findings highlight the potential advantages of rule-based data training paradigms, suggesting that they can enhance the overall effectiveness and reliability of machine learning models.

We can also notice that GEODIFF-CR only slightly improved GEODIFF's performance in this setting, which emphasizes the importance of METRIZATION. We also tested the performance of RDKIT in this setting. The results show that MIGDIFF-P performs better on COV-R and MAT-R. However, similar to the observations in (Xu et al. 2022), MIGDIFF-P either does not show comparable performance with RDKIT on COV-P and MAT-P. This indicates that ML models tend to explore more possible representatives, while RDKIT focuses on the most common conformations, emphasizing quality over diversity. Since RDKIT's optimization uses the force field method (Halgren 1996), we also tested the performance of molecules post-processed with force field optimization with MIGDIFF-P and GEODIFF.

Besides, the evaluation of conventional molecular conformations in Tab. 3 shows MIGDIFF consistently outperforming other baselines, demonstrating the efficacy of METRIZATION in enhancing in-distribution performance as well. More results and visualizations can be found at supplementary files.

**Property Prediction**

**Evaluation metrics** This task evaluates molecular *ensemble properties* (Axelrod and Gomez-Bombarelli 2022) using generated conformations to assess sample quality. We follow the property prediction task setup in (Xu et al. 2022; Shi et al. 2021), which uses 30 molecules from GEOM-QM9 dataset and generates 50 samples for each. We use the PSI4 toolkit to calculate the energy ($E$) and HOMO-LUMO gap ($\epsilon$) of each conformer, with comparisons made to the ground truth for average energy ($\bar{E}$) and average gap ($\Delta\bar{\epsilon}$). The evaluation is divided into two scenarios: in-distribution (ID) where the model is trained on the GEOM-QM9 dataset, and out-

| | $\overline{E}$ | | $\Delta\bar{\epsilon}$ | |
|---|---|---|---|---|
| Method | ID | OOD | ID | OOD |
| RDKIT | 0.9233 | | 0.3698 | |
| MIGDIFF-P | 0.1751 | | 0.3506 | |
| CONFGF | 0.1765 | 1.3560 | 0.4688 | 2.1339 |
| GEODIFF | 0.1551 | 0.3210 | 0.3091 | 0.4395 |
| MIGDIFF | **0.1522** | **0.1707** | **0.2879** | **0.3296** |

Table 4: Results of Property Prediction task.

of-distribution (OOD) where it is trained using the GEOM-DRUGS dataset. We use the mean absolute errors (MAE) between the calculated properties and the ground truth.

**Result Analysis** The results are shown in Tab. 4, we can see that both GEODIFF and MIGDIFF performs well on ID property prediction. It also shows that the prediction performance with MIGDIFF and MIGDIFF-P generated conformations has a more robust performance when faced with OOD setting. Since the properties are highly related and sensitive to the geometric structure, the results show that the incorporation of METRIZATION in MIGDIFF helps to produce more accurate conformation and more generalizable property predictions.

## Conclusion

In conclusion, we introduced a novel diffusion-based molecular generation framework, METRIZATION-informed Diffusion Generation (MIGDIFF), that effectively combines data-driven and distance geometry rule-driven paradigms by pre-training the deep geometric diffusion backbone within the METRIZATION sampling approach. This integration allows MIGDIFF to embed distance geometry within the diffusion framework, significantly enhancing its generalization capabilities. Experimental results show that MIGDIFF demonstrates robust generalization performance, effectively generating stable and accurate molecular conformations across diverse datasets. This work underscores the value of incorporating physical-informed constraints into diffusion models, paving the way for further advancements in building deep learning models for computational biochemistry.

## Acknowledgments

## References

Abuter, R.; Amorim, A.; Bauböck, M.; Berger, J.; Bonnet, H.; Brandner, W.; Clénet, Y.; Du Foresto, V. C.; De Zeeuw, P.; Dexter, J.; et al. 2019. A geometric distance measurement to the Galactic center black hole with 0.3% uncertainty. *Astronomy & Astrophysics*, 625: L10.

Alfakih, A. Y.; Khandani, A.; and Wolkowicz, H. 1999. Solving Euclidean distance matrix completion problems via semidefinite programming. *Computational optimization and applications*, 12: 13–30.

An, L. T. H.; and Tao, P. D. 2003. Large-scale molecular optimization from distance matrices by a dc optimization approach. *SIAM Journal on Optimization*, 14(1): 77–114.

Anderson, A. C. 2003. The process of structure-based drug design. *Chemistry & biology*, 10(9): 787–797.

Anstine, D. M.; and Isayev, O. 2023. Generative models as an emerging paradigm in the chemical sciences. *Journal of the American Chemical Society*, 145(16): 8736–8750.

Axelrod, S.; and Gomez-Bombarelli, R. 2022. GEOM, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data*, 9(1): 185.

Cao, H.; Tan, C.; Gao, Z.; Xu, Y.; Chen, G.; Heng, P.-A.; and Li, S. Z. 2024. A survey on generative diffusion models. *IEEE Transactions on Knowledge and Data Engineering*.

Dou, B.; Zhu, Z.; Merkurjev, E.; Ke, L.; Chen, L.; Jiang, J.; Zhu, Y.; Liu, J.; Zhang, B.; and Wei, G.-W. 2023. Machine learning methods for small data challenges in molecular science. *Chemical Reviews*, 123(13): 8736–8780.

Fu, T.; Xiao, C.; Li, X.; Glass, L. M.; and Sun, J. 2021. Mimosa: Multi-constraint molecule sampling for molecule optimization. *AAAI*, 35(1): 125–133.

Ganea, O.; Pattanaik, L.; Coley, C.; Barzilay, R.; Jensen, K.; Green, W.; and Jaakkola, T. 2021. Geomol: Torsional geometric generation of molecular 3d conformer ensembles. *NeurIPS*.

Gibson, D. G.; Young, L.; Chuang, R.-Y.; Venter, J. C.; Hutchison, C. A.; and Smith, H. O. 2009. Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nature methods*, 6(5): 343–345.

Gilmer, J.; Schoenholz, S. S.; Riley, P. F.; Vinyals, O.; and Dahl, G. E. 2017. Neural message passing for quantum chemistry. In *ICML*.

Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; and Aspuru-Guzik, A. 2018. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2): 268–276.

Guo, Z.; Liu, J.; Wang, Y.; Chen, M.; Wang, D.; Xu, D.; and Cheng, J. 2024. Diffusion models in bioinformatics and computational biology. *Nature reviews bioengineering*, 2(2): 136–154.

Halgren, T. A. 1996. Merck molecular force field. V. Extension of MMFF94 using experimental data, additional computational data, and empirical rules. *Journal of Computational Chemistry*, 17(5-6): 616–641.

Han, J.; Rong, Y.; Xu, T.; and Huang, W. 2022. Geometrically equivariant graph neural networks: A survey. *arXiv preprint arXiv:2202.07230*.

Havel, T. F. 1998. Distance geometry: Theory, algorithms, and chemical applications. *Encyclopedia of Computational Chemistry*, 120: 723–742.

Havel, T. F.; and Wüthrich, K. 1985. An evaluation of the combined use of nuclear magnetic resonance and distance geometry for the determination of protein conformations in solution. *Journal of molecular biology*, 182(2): 281–294.

Heinen, S.; Schwilk, M.; von Rudorff, G. F.; and von Lilienfeld, O. A. 2020. Machine learning the computational cost of quantum chemistry. *Machine Learning: Science and Technology*, 1(2): 025002.

Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *NeurIPS*.

Hoffman, S. C.; Chenthamarakshan, V.; Wadhawan, K.; Chen, P.-Y.; and Das, P. 2022. Optimizing molecules using efficient queries from property evaluations. *Nature Machine Intelligence*, 4(1): 21–31.

Hoogeboom, E.; Satorras, V. G.; Vignac, C.; and Welling, M. 2022. Equivariant diffusion for molecule generation in 3d. In *ICML*, 8867–8887. PMLR.

Hua, C.; Luan, S.; Xu, M.; Ying, Z.; Fu, J.; Ermon, S.; and Precup, D. 2024. Mudiff: Unified diffusion for complete molecule generation. In *Learning on Graphs Conference*, 33–1.

Ingraham, J. B.; Baranov, M.; Costello, Z.; Barber, K. W.; Wang, W.; Ismail, A.; Frappier, V.; Lord, D. M.; Ng-Thow-Hing, C.; Van Vlack, E. R.; et al. 2023. Illuminating protein space with a programmable generative model. *Nature*, 623(7989): 1070–1078.

Kabsch, W. 1976. A solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography*, 32(5): 922–923.

Keith, J. A.; Vassilev-Galindo, V.; Cheng, B.; Chmiela, S.; Gastegger, M.; Muller, K.-R.; and Tkatchenko, A. 2021. Combining machine learning and computational chemistry for predictive insights into chemical systems. *Chemical reviews*, 121(16): 9816–9872.

Kuntz, I. D. 1992. Structure-based strategies for drug design and discovery. *Science*, 257(5073): 1078–1082.

Kuszewski, J.; Nilges, M.; and Brünger, A. T. 1992. Sampling and efficiency of metric matrix distance geometry: a novel partial metrization algorithm. *Journal of biomolecular NMR*, 2: 33–56.

Landrum, G.; et al. 2006. RDKit: Open-source cheminformatics.

Liberti, L.; and Lavor, C. 2016. Six mathematical gems from the history of distance geometry. *International Transactions in Operational Research*, 23(5): 897–920.

Liberti, L.; and Lavor, C. 2018. Open research areas in distance geometry. *Open problems in optimization and data analysis*, 183–223.

Liberti, L.; Lavor, C.; Maculan, N.; and Mucherino, A. 2014. Euclidean distance geometry and applications. *SIAM review*, 56(1): 3–69.

Liberti, L.; Lavor, C.; Mucherino, A.; and Maculan, N. 2011. Molecular distance geometry methods: from continuous to discrete. *International Transactions in Operational Research*, 18(1): 33–51.

McArdle, S.; Endo, S.; Aspuru-Guzik, A.; Benjamin, S. C.; and Yuan, X. 2020. Quantum computational chemistry. *Reviews of Modern Physics*, 92(1): 015003.

Meng, G.; Lam, N. Y.; Lucas, E. L.; Saint-Denis, T. G.; Verma, P.; Chekshin, N.; and Yu, J.-Q. 2020. Achieving site-selectivity for C–H activation processes based on distance and geometry: a carpenter's approach. *Journal of the American Chemical Society*, 142(24): 10571–10591.

Ramakrishnan, R.; Dral, P. O.; Rupp, M.; and Von Lilienfeld, O. A. 2014. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1): 1–7.

Riniker, S.; and Landrum, G. A. 2015. Better informed distance geometry: using what we know to improve conformation generation. *Journal of chemical information and modeling*, 55(12): 2562–2574.

Rosei, F.; Schunack, M.; Naitoh, Y.; Jiang, P.; Gourdon, A.; Laegsgaard, E.; Stensgaard, I.; Joachim, C.; and Besenbacher, F. 2003. Properties of large organic molecules on metal surfaces. *Progress in Surface Science*, 71(5-8): 95–146.

Ross, J.; Belgodere, B.; Chenthamarakshan, V.; Padhi, I.; Mroueh, Y.; and Das, P. 2022. Large-scale chemical language representations capture molecular structure and properties. *Nature Machine Intelligence*, 4(12): 1256–1264.

Ryan, K.; Lengyel, J.; and Shatruk, M. 2018. Crystal structure prediction via deep learning. *Journal of the American Chemical Society*, 140(32): 10158–10168.

Satorras, V. G.; Hoogeboom, E.; and Welling, M. 2021. E (n) equivariant graph neural networks. In *ICML*, 9323–9332. PMLR.

Saxe, J. B. 1979. Embeddability of weighted graphs in k-space is strongly NP-hard. In *17th Allerton Conf. Commun. Control Comput.*, 480–489.

Schoenberg, I. J. 1935. Remarks to Maurice Frechet's article"sur la definition axiomatique d'une classe d'espace distances vectoriellement applicable sur l'espace de hilbert. *Annals of Mathematics*, 724–732.

Schütt, K.; Kindermans, P.-J.; Sauceda Felix, H. E.; Chmiela, S.; Tkatchenko, A.; and Müller, K.-R. 2017. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *NeurIPS*.

Shi, C.; Luo, S.; Xu, M.; and Tang, J. 2021. Learning gradient fields for molecular conformation generation. In *ICML*.

Siebenmorgen, T.; Menezes, F.; Benassou, S.; Merdivan, E.; Didi, K.; Mourão, A. S. D.; Kitel, R.; Liò, P.; Kesselheim, S.; Piraud, M.; et al. 2024. MISATO: machine learning dataset of protein–ligand complexes for structure-based drug discovery. *Nature Computational Science*, 1–12.

Simm, G.; and Hernandez-Lobato, J. M. 2020. A Generative Model for Molecular Distance Geometry. In *ICML*, 8949–8958.

Sit, A.; and Wu, Z. 2011. Solving a generalized distance geometry problem for protein structure determination. *Bulletin of Mathematical Biology*, 73: 2809–2836.

Sit, A.; Wu, Z.; and Yuan, Y. 2009. A geometric buildup algorithm for the solution of the distance geometry problem using least-squares approximation. *Bulletin of mathematical biology*, 71(8): 1914–1933.

Sobez, J.-G.; and Reiher, M. 2020. Molassembler: Molecular graph construction, modification, and conformer generation for inorganic and organic molecules. *Journal of Chemical Information and Modeling*, 60(8): 3884–3900.

Song, Y.; and Ermon, S. 2019. Generative modeling by estimating gradients of the data distribution. *NeurIPS*.

Spellmeyer, D. C.; Wong, A. K.; Bower, M. J.; and Blaney, J. M. 1997. Conformational analysis using distance geometry methods. *Journal of Molecular Graphics and Modelling*, 15(1): 18–36.

Trinajstic, N. 2018. *Chemical graph theory*. CRC press.

van Tilborg, D.; Brinkmann, H.; Criscuolo, E.; Rossen, L.; Özçelik, R.; and Grisoni, F. 2024. Deep learning for low-data drug discovery: hurdles and opportunities. *Current Opinion in Structural Biology*, 86: 102818.

Wang, H.; Fu, T.; Du, Y.; Gao, W.; Huang, K.; Liu, Z.; Chandak, P.; Liu, S.; Van Katwyk, P.; Deac, A.; et al. 2023a. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972): 47–60.

Wang, Y.; Elhag, A. A.; Jaitly, N.; Susskind, J. M.; and Bautista, M. A. 2023b. Generating molecular conformer fields. *arXiv preprint arXiv:2311.17932*.

Xu, M.; Luo, S.; Bengio, Y.; Peng, J.; and Tang, J. 2021a. Learning neural generative dynamics for molecular conformation generation. In *ICLR*.

Xu, M.; Wang, W.; Luo, S.; Shi, C.; Bengio, Y.; Gomez-Bombarelli, R.; and Tang, J. 2021b. An end-to-end framework for molecular conformation generation via bilevel programming. In *ICML*.

Xu, M.; Yu, L.; Song, Y.; Shi, C.; Ermon, S.; and Tang, J. 2022. GeoDiff: A Geometric Diffusion Model for Molecular Conformation Generation. In *ICLR*.

Yang, L.; Zhang, Z.; Song, Y.; Hong, S.; Xu, R.; Zhao, Y.; Zhang, W.; Cui, B.; and Yang, M.-H. 2023. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4): 1–39.

Zhou, Z.; Liu, R.; Zheng, J.; Wang, X.; and Yu, T. 2024. On Diffusion Process in SE (3)-invariant Space. *arXiv preprint arXiv:2403.01430*.