

SCALING THE PRIOR: SIZE-CONSISTENT GEOMETRIC DIFFUSION FOR 3D MOLECULAR GENERATION

Anonymous authors

Paper under double-blind review

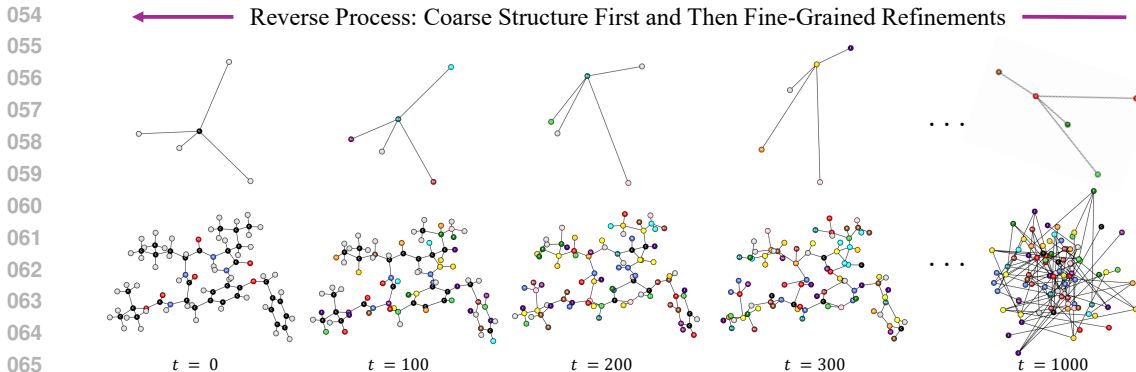
ABSTRACT

Diffusion models usually operate in fixed-dimensional metric spaces. In contrast, geometric molecular data naturally vary in dimensionality as molecules have different sizes (numbers of atoms). As a simple adaptation, existing diffusion models for geometric molecular generation employ network architectures that can handle variable-sized inputs, such as graph neural networks and transformers. **However, these approaches overlook the fact that the molecular size also determines the spatial scale of the atomic coordinates, which in turn induces inconsistent behaviors in the generative trajectories across different molecular sizes.** The generative process of geometric diffusion for 3D molecular generation can be viewed as first establishing a coarse structural target, followed by progressively refining the precise atomic positions. In particular, larger molecules tend to establish coarse structures earlier than smaller molecules due to their larger spatial scales relative to that of the noise. As a result, the reverse process becomes inconsistent across molecular sizes, with the denoising trajectories relying heavily on molecular sizes rather than on a unified generative pattern. In this work, we are the first to identify and analyze this size-induced inconsistency through a decomposition of the denoising dynamics, which reveals how spatial scale affects the progression of molecular formation, in both 3D structures and atom types. Building on this insight, we propose Scaling the Prior (StP), a simple yet effective approach that normalizes the learning and generative process across molecular sizes by rescaling the prior distribution based on molecular sizes. This adjustment harmonizes the denoising trajectories, enabling the model to learn a unified generative pattern and produce consistently high-quality molecules.

1 INTRODUCTION

In recent years, diffusion models (Ho et al., 2020; Song et al., 2020; 2021) have achieved great advances in both theory and practical applications (Nichol & Dhariwal, 2021; Rombach et al., 2022; Nie et al., 2025; Zhang et al., 2023). These models consist of a forward process and a reverse process. The forward process is a fixed Markov chain that incrementally corrupts data samples by adding noise until the original signal is lost. The reverse process is modeled by a denoising neural network, which is trained to iteratively undo the corruption and reconstruct the original data; once trained, new samples are generated by running the reverse process, beginning from pure Gaussian noise. Motivated by their success in vision tasks such as image generation, researchers have extended diffusion models to various domains, including 3D molecular generation (Hoogetboom et al., 2022; Xu et al., 2023). In particular, 3D molecular generation includes generating both the continuous 3D atomic configurations as well as discrete categorical features such as atom types. Diffusion models usually operate in fixed-dimensional metric spaces, such as $\mathbb{R}^{H \times W \times C}$ for images, where H , W , and C are the height, width, and number of channels, respectively. However, 3D molecular data varies in dimensionality because molecules have different numbers of atoms. A straightforward strategy taken by existing work is to adopt network architectures that are inherently suited to variable-sized¹ inputs, such as graph neural networks (GNNs) and transformers (Hoogetboom et al., 2022; Ding & Hofmann, 2025). However, 3D molecular data presents a unique challenge that has yet to be fully identified and addressed.

¹In this work, size always denotes the number of atoms; the spatial extent of a molecule is called *scale*.



067 Figure 1: The diffusion/forward process of two 3D molecules of varying sizes from the GEOM-
068 Drugs dataset (Axelrod & Gomez-Bombarelli, 2022), using the common noise schedule from prior
069 work (Hoogetboom et al., 2022). Colors represent different atom types. **Existing works apply the**
070 **same Gaussian prior across molecules of all sizes. However, it leads to size-induced inconsis-**
071 **tencies.** The geometry of the 5-atom molecule on the top is already unrecognizable after 200 steps,
072 while the 91-atom molecule on the bottom still preserves its overall geometry. **This size-induced**
073 **inconsistency makes larger molecules stabilize earlier in the reverse process and show better**
074 **performance**, as shown in Sec. 3.1 and Sec. 3.2. Atom types are represented as one-hot vectors with
075 the same scale across molecular sizes. However, inconsistencies in spatial scales also propagate to
076 atom types as shown in Sec. 3.2. These two molecules differ greatly in their spatial scales; they are
077 drawn to the same scale for visualization purposes.

078 **The generative process in 3D molecular diffusion can be viewed as first establishing a coarse**
079 **structural target, followed by progressively refining atomic positions.** This challenge arises from
080 a distinctive property of molecular data; that is, molecules vary in size (number of atoms), which
081 naturally induces different spatial scales. In contrast, pixel values in image data are always 0 to 255,
082 independent of the image dimension. For smaller molecules, the spatial extent is small, so the same
083 level of noise corresponds to larger perturbations on a relative scale compared to larger molecules.
084 This discrepancy leads to fundamentally different denoising behaviors across molecular sizes. For
085 larger molecules, the generative process establishes a coarse structural target and begins fine-grained
086 positional adjustments early, while for small molecules, at the same reverse time step, it remains fo-
087 cused on forming the coarse structures. An illustration is shown in Fig. 1. **Therefore, the generative**
088 **process behaves inconsistently across molecular sizes**, leading to suboptimal performance.

089 In this work, we analyze denoising dynamics through a decomposition of shape and atom types,
090 which separates the generative process into a shape-preserving radial part of the 3D configuration
091 and a categorical component governing atom types. This decomposition provides a quantitative lens
092 to reveal the scale-dependent inconsistency, showing how **molecules of different sizes establish**
093 **their coarse structural targets at different rates during the reverse process.** As we aim to
094 mitigate this issue, a natural question arises: *Can we simply normalize the data so that the 3D*
095 *atomic configurations of molecules of different sizes have the same scale?* Unfortunately, directly
096 normalizing the molecular data would distort important chemical information such as bond lengths,
097 which encode the type and strength of chemical interactions (Kindermans & Müller, 2018; Qu et al.,
098 2025). We provide more discussion in Sec. 3.3 and ablation studies in Sec. 4.4. Based on these
099 insights, we propose a simple yet effective solution: Scaling the Prior (StP). Instead of rescaling the
100 molecular data itself, which would irreparably distort chemically meaningful structural information,
101 StP rescales the prior distribution based on molecular sizes. **By aligning the noise scale with**
102 **the spatial extent of a molecule, this adjustment harmonizes the denoising trajectories across**
103 **different molecular sizes.** As a result, the model learns a unified generative pattern, establishing
104 structural target and refining local details in a consistent manner regardless of size.

104 **In summary, our contribution can be summarized as follows:** ① We are the first to identify and
105 introduce an important, yet often overlooked, size-dependent inconsistency in diffusion-based 3D
106 molecular generation. ② We propose principled metrics to quantify the structural and atom type
107 changes in the diffusion process, which can provide future insights. ③ We provide a simple yet
effective solution, StP, which rescales the noise based on the molecular size, harmonizing denoising

trajectories across molecule sizes while preserving chemically meaningful structural information. ④ Experimental results demonstrate that StP significantly improves validity, stability, and overall quality of generated molecules, without introducing additional architectural complexity.

2 PRELIMINARIES AND RELATED WORK

2.1 DIFFUSION MODELS

Diffusion models learn a data distribution by inverting a forward diffusion process. Given a data point \mathbf{x} , the forward diffusion process adds increasing levels of Gaussian noise to it:

$$q(\mathbf{z}_t | \mathbf{x}) = \mathcal{N}(\mathbf{z}_t | \alpha_t \mathbf{x}, \sigma_t^2 \mathbf{I}), \quad t = 0, 1, 2, \dots, T, \quad (1)$$

where \mathbf{z}_t denotes a noisy version of \mathbf{x} at time t , $\alpha_t > 0$ controls how much of the original signal is retained, and $\sigma_t > 0$ controls how much noise is injected. A special case is the variance preserving process in which $\alpha_t^2 + \sigma_t^2 = 1$ (Ho et al., 2020; Song et al., 2021). In general, $\alpha_t \approx 1$ at $t = 0$ and then monotonically decreases to 0 at $t = T$, corresponding to the progressive corruption of the data into pure noise. The diffusion process is Markovian and can equivalently be expressed by its transition distributions:

$$q(\mathbf{z}_t | \mathbf{z}_s) = \mathcal{N}(\mathbf{z}_t | \alpha_{t|s} \mathbf{z}_s, \sigma_{t|s}^2 \mathbf{I}), \quad (2)$$

where $T \geq t > s \geq 0$, $\alpha_{t|s} = \alpha_t / \alpha_s$, and $\sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2$. Given equation 1 and equation 2, we can derive the posterior of the transitions conditioned on \mathbf{x} :

$$q(\mathbf{z}_s | \mathbf{z}_t, \mathbf{x}) = \mathcal{N}(\mathbf{z}_s | \boldsymbol{\mu}_{s|t}(\mathbf{z}_t, \mathbf{x}), \sigma_{s|t}^2 \mathbf{I}), \quad (3)$$

where $\boldsymbol{\mu}_{s|t}(\mathbf{z}_t, \mathbf{x})$ and $\sigma_{s|t}^2$ can be derived analytically:

$$\boldsymbol{\mu}_{s|t}(\mathbf{z}_t, \mathbf{x}) = \frac{\alpha_{t|s} \sigma_s^2}{\sigma_t^2} \mathbf{z}_t + \frac{\alpha_s \sigma_{t|s}^2}{\sigma_t^2} \mathbf{x}, \quad \sigma_{s|t}^2 = \frac{\sigma_{t|s}^2 \sigma_s^2}{\sigma_t^2}. \quad (4)$$

The posterior distribution defines the reverse process (the generative process). The posterior distribution is generally unknown in the generative process as it is conditioned on unknown \mathbf{x} . In practice, a neural network parameterization ϕ is trained to approximate the true posterior:

$$p_\phi(\mathbf{z}_s | \mathbf{z}_t) = q(\mathbf{z}_s | \mathbf{z}_t, \mathbf{x}_\phi(\mathbf{z}_t, t)). \quad (5)$$

A variational lower bound on the log-likelihood of \mathbf{x} can be derived for the diffusion model as:

$$\log p(\mathbf{x}) \geq \mathcal{L}_0 + \mathcal{L}_{\text{prior}} + \sum_{t=1}^T \mathcal{L}_t, \quad (6)$$

where $\mathcal{L}_0 = \log p(\mathbf{x} | \mathbf{z}_0)$ models the reconstruction error, $\mathcal{L}_{\text{prior}} = -\text{KL}(q(\mathbf{z}_T | \mathbf{x}) | p(\mathbf{z}_T))$ models the divergence between the prior standard normal distribution and the final latent distribution $q(\mathbf{z}_T | \mathbf{x})$, and $\mathcal{L}_t = -\text{KL}(q(\mathbf{z}_s | \mathbf{x}, \mathbf{z}_t) | p(\mathbf{z}_s | \mathbf{z}_t))$ quantifies the divergence between the learned posterior distribution and the true posterior distribution. Following prior work (Ho et al., 2020), instead of directly predicting \mathbf{x} , optimization is easier when predicting the noise instead. In particular, $\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})$, then the neural network ϕ outputs $\hat{\boldsymbol{\epsilon}} = \phi(\mathbf{z}_t, t)$, so that $\hat{\mathbf{x}} = (1/\alpha_t) \mathbf{z}_t - (\sigma_t/\alpha_t) \hat{\boldsymbol{\epsilon}}$. Let $\text{SNR}(t)$ be the *signal-to-noise ratio* (Kingma et al., 2023), defined as $\text{SNR}(t) = \alpha_t^2 / \sigma_t^2$, the term \mathcal{L}_t can be further expanded as:

$$\mathcal{L}_t = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{1}{2} \left(\frac{\text{SNR}(t-1)}{\text{SNR}(t)} - 1 \right) \|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}\|^2 \right]. \quad (7)$$

2.2 3D MOLECULAR DIFFUSION

We are interested in generating 3D molecules from scratch. Specifically, a molecule with N atoms can be represented as (\mathbf{x}, \mathbf{h}) , where $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \in \mathbb{R}^{N \times 3}$ corresponds to the atomic coordinates in the 3D space and $\mathbf{h} = (\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N) \in \mathbb{R}^{N \times d}$ corresponds to d -dimensional features of the atoms, e.g., one-hot encoding of atom types. Pioneering works in 3D molecular generation with diffusion models, such as EDM (Hooeboom et al., 2022) and GeoLDM (Xu et al., 2023),

emphasize the importance of $E(3)$ -invariance of the molecular distribution: Euclidean transformations (translations and rotations) do not change the underlying molecule, and thus the probability density should remain unchanged. Translation invariance is enforced by operating in a zero-mean linear subspace, in which the coordinates of both molecules and noise are centered at the origin: $\mathbf{x} := (\mathbf{x}_1 - \bar{\mathbf{x}}, \mathbf{x}_2 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}}) \in \mathbb{R}^{(N-1) \times 3}$, where $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$. Rotation invariance is further ensured by employing a denoising network that is equivariant to rotations (Hoogeboom et al., 2022; Xu et al., 2022). Nevertheless, recent studies suggest that rotation equivariance is not always necessary (Ding & Hofmann, 2025; Joshi et al., 2025), and that employing equivariant denoising networks may introduce substantial computational overhead. Besides diffusion models, flow-based approaches have also attracted significant attention (Dunn & Koes, 2024; Irwin et al., 2025; Song et al., 2023; 2024; Hong et al., 2025). These methods typically employ a Gaussian prior over atomic positions, and it has been shown that, under a Gaussian prior, they are mathematically equivalent to diffusion models (Albergo et al., 2023; Ma et al., 2024; Gao et al., 2025).

The theory of diffusion models and their stochastic differential equation (SDE) formulation are defined over fixed-dimensional metric spaces. This creates a fundamental mismatch for molecules, whose sizes vary and thus cannot be naturally embedded into a single fixed ambient space. In practice, a workaround is to employ denoising networks that can work with variable-sized inputs, such as GNNs and transformers (Hoogeboom et al., 2022; Xu et al., 2023; Joshi et al., 2025; Ding & Hofmann, 2025). One way to view this workaround is through the lens of amortization: there is a diffusion process for molecules of each size; however, a unified denoising network is trained to amortize across these processes, learning to handle variable molecular sizes within a single model. In this view, the SDE framework remains conceptually defined for a fixed-dimensional ambient space, but the denoising network provides a shared parametrization that generalizes across different molecule sizes. While this amortization enables the treatment of variable-size data at the architectural level, it overlooks another critical factor: Molecular size also affects the spatial scale of the coordinates, which in turn leads to inconsistencies in the diffusion process.

3 SCALING THE PRIOR: SIZE-CONSISTENT GEOMETRIC DIFFUSION

3.1 OBSERVING SIZE-INDUCED INCONSISTENCIES

To begin, we present an intriguing and counterintuitive phenomenon in 3D molecular generation with diffusion models. **As shown in Fig. 2, there is a clear trend that as molecular size increases, sampling quality increases, despite the fact that large molecules have higher structural complexity and lower data availability.** Similar inconsistencies can be observed in various diffusion-based 3D molecular generation methods and datasets; we present the experimental details and additional results in Appendix A.

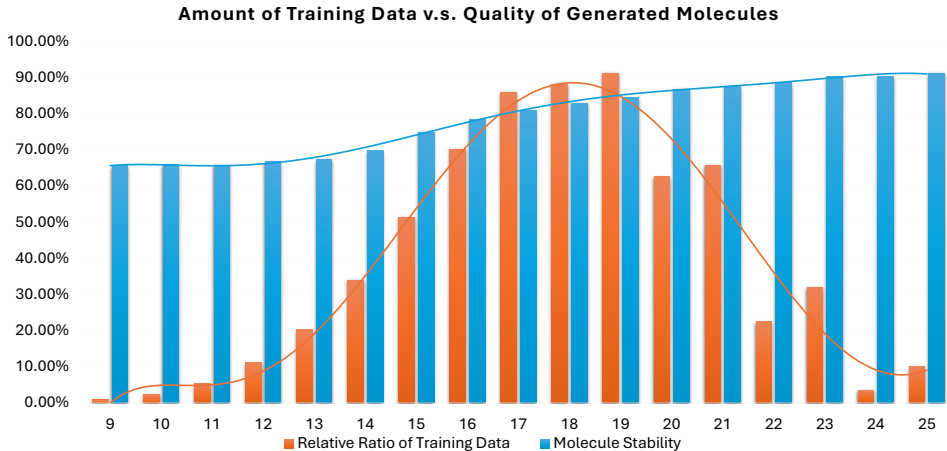


Figure 2: Sampling quality (y -axis) versus molecular size (x -axis) with EDM (Hoogeboom et al., 2022) on the QM9 dataset (Ramakrishnan et al., 2014). Clearly, larger molecules achieve higher sampling quality despite the scarcity of corresponding training samples. For visualization, the relative ratio of training data for the molecular size with the most training samples is normalized to the highest stability value, and all other sizes are scaled relative to this.

This phenomenon calls for careful analysis to uncover its underlying causes and implications in diffusion-based 3D molecular generation. Therefore, we delve deeper to uncover its underlying cause. In Sec. 3.2, we identify the root cause from the unique characteristic of 3D molecular data that different sizes of molecules have different scales, which leads to inconsistencies in the generative process and generated sample qualities. Building on this analysis, in Sec. 3.3, we propose a solution tailored to mitigate the observed inconsistencies and improve generative performance across molecular sizes.

Discussion on Latent Diffusion Models. Latent diffusion models encode the coordinates and atom features (\mathbf{x}, \mathbf{h}) into a lower-dimensional latent space. Existing latent diffusion models, e.g., GeoLDM and RADM, still separate coordinates and atom features, keeping coordinates in the 3D Euclidean space in the latent representation as $\mathbf{x}' \in \mathbb{R}^{(N-1) \times 3}$, $\mathbf{h}' \in \mathbb{R}^{N \times d'}$, where typically $d' < d$. Most importantly, coordinates rarely change in the latent space (up to $SO(3)$ transformations for RADM), and size-induced inconsistencies remain in these latent diffusion models. The analysis and techniques presented in this work are still applicable. A more detailed discussion and concrete examples are provided in Appendix B.

3.2 UNVEILING SIZE-INDUCED INCONSISTENCIES

Diffusion models can be viewed as an approximate autoregression in the frequency domain (Risänen et al., 2023). In the forward process, noise gradually destroys high-frequency details while low-frequency overall structures persist longer. In generation, the model first constructs the coarse structure and then progressively adds higher-frequency details towards the generation target.

Remark 3.1. The generation target refers to all possible clean data samples that can be inferred from noisy samples. For example, given a mildly noisy image, we may infer its cleaned form.

Importantly, the generation target is not always a single clean data point. For example, at $t = T$, pure noise provides no information about any specific clean sample; the generation target is ambiguous and can contain any data sample in the distribution. In contrast, at values of t closer to 0, the slightly noised input corresponds only to a limited subset of similar data samples (e.g., molecules with similar 3D configurations). **As t goes to 0, the noise level decreases, and the generation target becomes more precise and eventually converges to a concrete clean data point (e.g., a specific molecule).**

To this end, the generation target of a noisy sample (z_t^x, z_t^h) , which contains both noisy coordinates and atom types, can be defined as

$$\text{Target}(z_t^x, z_t^h) = \left\{ (\mathbf{x}, \mathbf{h}) \mid p_\phi(\mathbf{x}, \mathbf{h} \mid z_t^x, z_t^h) \geq \tau \cdot \max_{\mathbf{x}', \mathbf{h}'} p_\phi(\mathbf{x}', \mathbf{h}' \mid z_t^x, z_t^h) \right\}, \quad (8)$$

where $0 < \tau \leq 1$ is a threshold parameter. Intuitively, in this definition, the generation target corresponds to the clean data samples at $t = 0$ whose likelihood remains within a fraction τ of the most probable sample under the posterior (the noisy sample).

For molecular diffusion, the generative process first finds a coarse atomic structure and then refines the precise atomic positions toward a concrete molecule as the generation target becomes more precise. However, we notice that there is an inconsistency across molecular sizes: For larger molecules, the generation target becomes precise faster, and the generative process initiates fine-grained positional adjustments toward concrete molecules earlier. In contrast, for smaller molecules at the same reverse time step, it still continues to focus on forming the coarse structures. Although directly comparing the generation target in equation 8 for different time t is intractable; we decompose it into two generation target alignment components, the *shape alignment coefficient* γ_t and the *atom type alignment ratio* β_t , to quantify this phenomenon. Concretely, given a noisy sample $z_t = (z_t^x, z_t^h)$, we estimate a representative of $\text{Target}(z_t^x, z_t^h)$ as $(\hat{\mathbf{x}}_t, \hat{\mathbf{h}}_t) = (1/\alpha_t) z_t - (\sigma_t/\alpha_t) \hat{\epsilon}^2$, which is the predicted clean data at time t . For a generative trajectory $\{z_t\}_{t=0}^T$, we take $(\hat{\mathbf{x}}_0, \hat{\mathbf{h}}_0)$ as its final generation target. By comparing the predicted clean data at time t : $(\hat{\mathbf{x}}_t, \hat{\mathbf{h}}_t)$ with the predicted clean data at time 0: $(\hat{\mathbf{x}}_0, \hat{\mathbf{h}}_0)$, we can assess the rate at which the generation target converges to a concrete

²Following prior work, the 3D coordinates and atom features are concatenated, which are then processed by a single denoising network to produce the combined noise prediction.

molecule with respect to time t . The convergence in generation target is quantified by:

$$\gamma_t = \frac{\langle \hat{\mathbf{x}}_t, \hat{\mathbf{x}}_0 \rangle}{\|\hat{\mathbf{x}}_0\|^2}, \quad \beta_t = \frac{1}{N} \sum_{n=1}^N \mathbf{1} [\hat{\mathbf{h}}_{t,n} = \hat{\mathbf{h}}_{0,n}], \quad (9)$$

where $\mathbf{1}[\cdot]$ is the indicator function checking if the n -th atom types are the same, γ_t is the radial coefficient indicating the alignment of the shape between $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{x}}_0$ (how much of $\hat{\mathbf{x}}_t$ lies in the direction of $\hat{\mathbf{x}}_0$), and β_t is the ratio of atoms whose types remain unchanged. We adopt the radial projection to quantify the convergence of generation target in atomic coordinates because the Euclidean norm difference $\|\hat{\mathbf{x}}_t - \hat{\mathbf{x}}_0\|$ is itself sensitive to the data scale, while the radial projection emphasizes directional alignment and relative magnitude consistency. **Overall, γ_t and β_t reflect how quickly the denoising trajectory “lines up” with the final 3D geometric structure and atom types, respectively.**

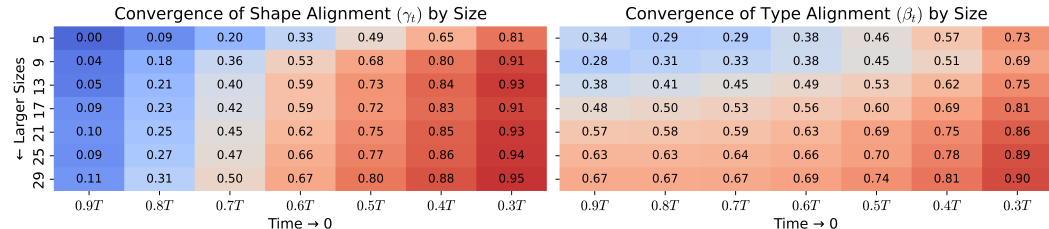


Figure 3: Inconsistent convergence rates across molecule sizes on QM9 with EDM. Values closer to 1 indicate better alignment with the final generation target. Larger molecules tend to stabilize more quickly than smaller ones, resulting in discrepancies in generation behavior. It is important to note that atom types are represented as one-hot vectors, which remain on the same scale across different molecule sizes; however, they are still affected by the stabilization of the 3D geometry.

By construction, γ_t and β_t provide complementary views of the convergence of the generation target. As the reverse diffusion proceeds toward $t \rightarrow 0$, both γ_t and β_t should approach 1. We now track these two quantities, γ_t and β_t , across molecules of different sizes to study how generation target emerges during generation. As shown in Fig. 3, there is a clear difference in the rate of convergence between different molecular sizes; **as the molecular size grows, both γ_t and β_t saturate earlier in the reverse trajectory.** Details and additional results on GEOM-Drugs and with other backbone diffusion models are shown in Appendix A.2. This indicates that larger molecules quickly stabilize to a coarse generation target, while smaller molecules continue to undergo structural adjustments for longer. Such inconsistencies have significant implications for molecular generation quality.

3.3 SCALING THE PRIOR: HARMONIZING GENERATIVE DYNAMICS

To mitigate the size inconsistency in the generative process, a straightforward approach is to ensure that molecules of different sizes are on the same scale. Atom types are categorical features and, in principle, share the same scale across molecular sizes. However, we still observe systematic differences because atom-type predictions are *not independent* of coordinates. Both modalities are processed simultaneously by the same denoising network. As a result, the varying scales of coordinates across molecular sizes are the main issue. To this end, we propose to Scale the Prior. Instead of using a unified standard Gaussian as the prior for coordinates, we introduce a size-dependent variance for the prior. Specifically, the marginal distribution in the forward process becomes:

$$q_\gamma(z_t^x | \mathbf{x}) = \mathcal{N}(z_t^x | \alpha_t \mathbf{x}, \gamma_N^2 \sigma_t^2 \mathbf{I}), \quad (10)$$

$$q(z_t^h | \mathbf{h}) = \mathcal{N}(z_t^h | \alpha_t \mathbf{h}, \sigma_t^2 \mathbf{I}),$$

where q_γ denotes the marginal distribution from StP to be distinguished from that of the standard diffusion in equation 1, γ_N is a size-varying scaling factor, and σ_t is from the original noise schedule as in equation 1. For clarity, we write the forward process separately for 3D positions and atom features; however, since the Gaussian is isotropic, this is equivalent to writing them jointly.

Proposition 3.2. *The marginal distributions at any time t of Scaling the Prior are equivalent to those of normalizing the data \mathbf{x} by a factor of $\frac{1}{\gamma_N}$ before the forward process, and subsequently unnormalizing after sampling.*

Proposition 3.2 is straightforward to prove; nevertheless, we include the proof in Appendix C. This result establishes a direct connection between StP and normalization. In point cloud diffusion in computer vision, the data is often already approximately on the same scale (Luo & Hu, 2021), or per-sample normalization is applied (Tyszkiewicz et al., 2023). Here, per-sample normalization refers to schemes that are not consistent across the dataset; for example, normalizing molecules differently depending on their sizes. **A key distinction between StP and standard normalization is that StP preserves structural information, such as bond lengths, consistently across the dataset, whereas per-sample normalization does not. Structural information, such as bond lengths, is fundamental to molecular modeling (Kindermans & Müller, 2018; Qu et al., 2025).** Naive per-sample normalization can alter chemically important distances, leading to generative models that struggle to capture consistent and chemically meaningful geometry. **Concretely, in StP, the forward process in equation 10 scales the entire coordinate space by α_t and injects isotropic Gaussian noise. As a result, the expected interatomic distances, including bond lengths, are $\mathbb{E}[|z_{t,i}^x - z_{t,j}^x|] = \alpha_t |\mathbf{x}_i - \mathbf{x}_j|$ (i, j denote the i -th and j -th atom, respectively), which shows that the expectation of bond lengths is purely determined by the global scaling factor α_t and remains consistent across time steps. On contrast, in per-sample normalization, the expected interatomic distances $\mathbb{E}[|z_{t,i}^x - z_{t,j}^x|] = \frac{\alpha_t}{\gamma_N} |\mathbf{x}_i - \mathbf{x}_j|$, which introduces an arbitrary size-dependent factor γ_N , altering the expected bond geometry. In addition, we provide a learning perspective of StP in Appendix A.3, which offers a high-level illustration of why StP is beneficial for learning the underlying molecular distribution.**

Normalizing the Scales. Since it is now clear that StP is equivalent to normalization but preserves chemically meaningful structural information, we can define γ_N so that it effectively normalizes molecules of different sizes to the same scale. For a 3D configuration $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \in \mathbb{R}^{N \times 3}$, assuming zero-mean as described in Sec. 2.2, we define the *scale* of \mathbf{x} as:

$$s(\mathbf{x}) = \frac{1}{|\text{Hull}(\mathbf{x})|} \sum_{v \in \text{Hull}(\mathbf{x})} \|v\|_2, \quad (11)$$

where $\text{Hull}(\mathbf{x})$ denotes the set of vertices (atoms) of the convex hull of \mathbf{x} and $|\text{Hull}(\mathbf{x})|$ is the number of convex hull vertices. For reference, we provide the definition of convex hull in Appendix A.4. To obtain a normalization factor that is consistent across molecules of different sizes, for molecules of size N , we define the *size-specific scale* as the average scale over all the molecules of size N . We now can define γ_N as the in proportion to the size-specific scale. Mathematically:

$$\gamma_N = \frac{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_N}[s(\mathbf{x})]}{Z}, \quad (12)$$

where \mathcal{D}_N is the subset of molecules containing exactly N atoms in the training set, and Z is a normalization factor that is the same in all molecular sizes and will not affect the relative scale across different sizes. An empirical choice can be $Z = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[s(\mathbf{x})]$ with \mathcal{D} being the entire dataset.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets. We evaluate our approach on the QM9 (Ramakrishnan et al., 2014) and GEOM-Drugs (Axelrod & Gomez-Bombarelli, 2022) benchmarks, which are standard and widely used benchmarks in 3D molecular generation. The QM9 dataset contains 130K molecules with up to 29 atoms (including hydrogens). The GEOM-Drugs dataset comprises 430K molecules with up to 181 atoms and an average size of 44.4 atoms. For both datasets, we follow the exact data splits and setup used in Hoogeboom et al. (2022).

Baselines. We use EDM, RADM (DiT-B), and GeoLDM, which are state-of-the-art 3D molecular diffusion models, as backbone models for our proposed StP and compare with the original without StP. To ensure a fair comparison, we strictly follow the implementation details provided in the original works, except that we might use a different batch size due to computational constraints. We also provide results for several **non-diffusion-based** 3D molecular generation models for reference, including G-SchNet (Gebauer et al., 2019), ENF (Satorras et al., 2021), EDM-bridge (Wu et al., 2022), EquiFM (Song et al., 2023), and GeoBFN (Song et al., 2024). All diffusion and flow baselines are reported with 1,000 steps. GeoBFN also has a 2,000-step variant; we report the 1,000-step version for consistency. The baseline results are directly obtained from their original works.

Evaluation Metrics. Following prior work, the bond types are based on the pairwise atomic distance and the atom types. We report the following evaluation metrics: (1) *Atom Stability*: the percentage of atoms in all generated molecules with the correct valence; (2) *Molecule Stability*: the percentage of generated molecules in which all atoms are stable; and (3) *Validity*×*Uniqueness* of the generated molecules as measured by RDKit. For the GEOM-Drugs dataset, following prior work, we only report atom stability and validity, since molecule stability is close to 0 and uniqueness is close to 1 for all methods. **Following pioneering works in 3D molecular generation (Hoogetboom et al., 2022; Xu et al., 2023), we do not use explicit bond information, and we do not apply any post-hoc refinement using computational chemistry software (e.g., Open Babel) to improve the generated molecules. There are works that explicitly use bond information or apply post-hoc refinements (Irwin et al., 2025; Dunn & Koes, 2025); these approaches naturally achieve better reported performance. Therefore, our work is not directly comparable to these works.** Following prior work, for all experiments in this section, 10,000 molecules are sampled, **the experiments are repeated 3 times with the mean values reported**, and the standard division is not reported if it is negligible.

4.2 STP: IMPROVED GENERATION QUALITIES

Table 1: Results of different 3D molecular generation pipelines. Three SOTA 3D molecular diffusion models, EDM, GeoLDM, and RADM, are adapted into their StP (our method) variants. The standard deviations are reported for QM9 after \pm ; they are negligible after rounding for GEOM-Drugs. Whenever an StP variant surpasses its original model, the result is highlighted in orange. The overall best performing model is highlighted in red.

Dataset	QM9				GEOM-Drugs	
	Atom Stab (%)	Molecule Stab (%)	Valid (%)	Valid× Unique (%)	Atom Stab (%)	Valid (%)
Dataset	99.00	95.20	97.70	97.70	86.50	99.90
G-SchNet	95.70	68.10	85.50	80.30	-	-
ENF	85.00	84.90	40.20	39.40	-	-
EDM-bridge	98.80	84.60	92.00	90.70	82.40	92.80
EquiFM	98.90	88.30	94.70	93.50	84.10	98.90
GeoBFN	99.08	90.87	95.31	92.96	85.60	92.08
EDM	98.70	82.00	91.90	90.70	81.30	92.60
EDM-StP	98.83±0.03	88.07±0.22	94.41±0.08	92.63±0.14	84.11	95.59
RADM	98.50	87.30	94.10	91.70	85.00	99.30
RADM-StP	98.59±0.01	87.62±0.10	94.19±0.17	91.51±0.15	85.27	99.49
GeoLDM	98.90	89.40	93.80	92.70	84.40	99.30
GeoLDM-StP	99.08±0.05	90.70±0.22	95.41±0.16	93.49±0.16	86.78	99.37

Baseline results are taken from original works, some only have one decimal places available.

We first generate molecules unconditionally with trained models. The results are presented in Table 1. For EDM, StP enhances the earliest EDM to surpass GeoLDM in generating valid molecules and outperforms the most recent RADM on QM9 across all metrics. For RADM, StP improves the quality of generated molecules with higher atom stability, molecular stability, and validity; the uniqueness is lower, potentially because the transformer backbone of RADM fits the training data strongly. For GeoLDM, StP improves performance across all metrics; in fact, **GeoLDM-StP outperforms the complicatedly designed SOTA flow-based baseline, GeoBFN, and achieves the overall best results³, establishing a new SOTA performance for 3D molecular diffusion**. It is worth noting that GeoLDM achieves higher atom stability than the dataset distribution on both QM9 and GEOM-Drugs, similar to GeoBFN on QM9. This may be because most atoms in the dataset are stable, and the model implicitly denoises toward the mode of the distribution (i.e., stable atoms). This is desirable in stable molecular generation. From a generative modeling perspective, a good model should closely match the data distribution. In this regard, GeoLDM-StP is still the closest to the dataset distribution compared to other methods (the difference for GeoLDM-StP is the lowest; 0.08% on QM9 and 0.28% on GEOM-Drugs). **In addition, we provide the same per-size visualization as in Fig. 2 but with StP in Fig. 11 in Appendix A.3. Clearly, StP improves the quality of generated molecules across all sizes and significantly reduces the size-induced inconsistency. In conclusion, StP consistently enhances the quality of generated molecules by improving both stability and validity across backbone diffusion models and benchmark datasets, demonstrat-**

³No single model achieves the best results across all columns; however, GeoLDM-StP demonstrates the best overall performance.

Table 2: Results on conditional generation tasks. The reported metric is MAE (lower is better). The overall best result is highlighted in **red**. Clearly, StP enhances the ability to generate molecules with desired properties and achieves the best performance for all 6 tasks.

Property Unit	α Bohr ³	$\Delta\varepsilon$ meV	$\varepsilon_{\text{HOMO}}$ meV	$\varepsilon_{\text{LUMO}}$ meV	μ D	C_v $\frac{\text{cal}}{\text{mol}}\text{K}$
QM9 (Lower Bound)	0.10	64	39	36	0.043	0.040
EDM	2.76	655	356	584	1.111	1.101
EquiFM	2.41	591	337	530	1.106	1.033
GeoLDM	2.37	587	340	522	1.108	1.025
GeoBFN	2.34	577	328	516	0.998	0.949
RADM	1.98	458	290	383	0.814	0.869
RADM-StP	1.96	425	282	360	0.792	0.861

ing its effectiveness and generalizability. These results and findings highlight the importance, effectiveness, and generalizability of StP.

Efficient Sampling. The sampling process is unnecessarily slow in the original implementations of EDM, GeoLDM, and RADM, especially for GEOM-Drugs. We provide an efficient implementation that accelerates the sampling process by at most $12.21\times$ (from 12.94 seconds per sample to 1.06 seconds). The details and numerical results are provided in Appendix D.1. We do not claim any novel technical contribution; it is provided solely as a resource for the community.

4.3 CONDITIONAL GENERATION

Following prior work Hoogeboom et al. (2022); Xu et al. (2023), we assess the effectiveness of StP in conditional generation tasks in QM9, specifically generating molecules with target properties. We report the Mean Absolute Error (MAE) using the pretrained predictor; a smaller deviation from QM9 (Lower Bound) indicates better performance. More details are described in D.2. We need to train a separate diffusion model for each conditional generation task. Due to computational constraints, we only evaluate the best-performing diffusion model on conditional tasks, RADM, with StP. The results are presented in Table 2. **Clearly, StP not only generates more valid and stable molecules, but also enhances the ability to generate molecules with desired property values; StP improves the performance of RADM on all six conditional generation tasks and establishes a new SOTA in conditional generation performance.**

4.4 ABLATION STUDY: COMPARISON WITH DIRECT NORMALIZATION

As discussed in Sec. 3.3, directly normalizing molecules using different normalization constants across molecular sizes distorts important structural information. To demonstrate its impact, we conduct an ablation study comparing the proposed StP (-StP) with direct data normalization on the QM9 dataset (-N). In particular, the atomic coordinates are normalized by γ_N as discussed in Sec. 3.3. The results are presented in Table 3. For EDM, the results may at first seem counterintuitive: direct normalization leads to lower atom stability but higher molecule stability and overall validity compared with the original model. This occurs because ① there are more unstable atoms in larger molecules due to distortions in structural information (e.g., bond length), and no matter how many unstable atoms there are, it is just one unstable molecule, and ② the benefits of harmonizing generative dynamics outweigh the harms of structural distortion for smaller molecules. For GeoLDM, similarly to EDM, we observe the same atom stability but higher molecule stability and validity. For RADM, we observe even worse performance with direct normalization; this is potentially due to the transformer architecture, causing the harms of structural distortion outweigh the benefits of reducing size-inconsistencies. **Overall, for all three diffusion models, StP consistently demonstrates the best performance. In conclusion, these re-**

Table 3: Results on comparison with direct normalization. The best result among the original, direct normalization, and StP is highlighted in **orange**. Clearly, StP achieves the best performance.

	Atom Stab (%)	Mol. Stab (%)	Valid (%)
EDM	98.7	82.0	91.9
EDM-N	98.6	85.8	93.1
EDM-StP	98.8	88.1	94.4
RADM	98.5	87.3	94.1
RADM-N	98.4	86.2	93.6
RADM-StP	98.6	87.6	94.2
GeoLDM	98.9	89.4	93.8
GeoLDM-N	98.9	90.1	94.3
GeoLDM-StP	99.1	90.7	95.4

486 **sults clearly illustrate the importance of StP, which not only harmonizes generative dynamics**
487 **across different molecular sizes but also preserves essential chemical structural information.**
488

489 5 CONCLUSION

490
491 In this work, we identified a fundamental yet previously overlooked issue in diffusion-based 3D
492 molecular generation: the inconsistency of denoising dynamics across molecular sizes. By decom-
493 posing the generative process into shape and atom type components, we revealed how molecules
494 of different sizes establish structural targets at different rates, leading to invalid or unstable outputs.
495 To address this, we proposed StP, Scaling the Prior, which rescales the noise distribution relative to
496 molecular sizes. Unlike direct data normalization, StP preserves chemically meaningful information
497 such as bond lengths while harmonizing denoising trajectories across molecular sizes. Our analy-
498 sis and experiments demonstrate that StP significantly improves the validity, stability, and overall
499 quality of generated molecules, without requiring additional architectural complexity.

500 **Limitations and Future Work.** While our work addresses size-induced inconsistency in diffusion-
501 based 3D molecular generation, our work can be extended to flow-based 3D molecular generation.
502 Flow-based 3D molecular generation usually uses the standard Gaussian distribution as the prior
503 distribution for coordinates as well (Song et al., 2023; 2024; Dunn & Koes, 2024). Practically, it can
504 work the same way that we use different variance levels for the prior Gaussian. Furthermore, ap-
505 plying to flow models, StP theoretically could be justified from the perspective of optimal transport,
506 where StP can reduce the cost of transporting noise distributions across molecules of different sizes.
507 **In addition, the normalization factors used in Sec. 3.3 are not derived from physical or chemical**
508 **principles. Exploring normalization factors with physical justification or chemical meaning could**
509 **be an interesting direction for future work.**
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

ETHICS STATEMENT

This research is centered on the development and evaluation of diffusion models for 3D molecular generation. The study does not involve human participants, personal data, or any sensitive information that could raise concerns regarding privacy, security, or fairness. Moreover, no conflicts of interest, issues of legal compliance, or potentially harmful applications have been identified in connection with this work.

REPRODUCIBILITY

The datasets used in this work are open-source. Upon acceptance of the paper, we will release the source code together with detailed instructions for dataset preparation, pre-trained models, and configuration files necessary to reproduce the main experiments. Along with the code, we will provide the values of scales, the values of γ_N (see Sec. 3.3), and evaluation results for different molecular sizes. Comprehensive guidelines, including command-line examples for training and evaluation, will also be made available. All theoretical claims in this work are rigorously supported by proofs presented in the appendix.

REFERENCES

- Michael S. Albergo, Nicholas M. Boffi, and Eric Vanden-Eijnden. Stochastic Interpolants: A unifying framework for flows and diffusions, 2023. URL <https://arxiv.org/abs/2303.08797>.
- Simon Axelrod and Rafael Gomez-Bombarelli. GEOM, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data*, 9(1):185, 2022.
- Mark De Berg, Otfried Cheong, Marc Van Kreveld, and Mark Overmars. *Computational geometry: algorithms and applications*. Springer, 2008.
- Yuhui Ding and Thomas Hofmann. Scalable non-equivariant 3d molecule generation via rotational alignment. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=15KpQ5MmaD>.
- Ian Dunn and David R. Koes. Flowmol3: Flow matching for 3d de novo small-molecule generation, 2025. URL <https://arxiv.org/abs/2508.12629>.
- Ian Dunn and David Ryan Koes. Mixed continuous and categorical flow matching for 3d de novo molecule generation, 2024. URL <https://arxiv.org/abs/2404.19739>.
- Ruiqi Gao, Emiel Hoogeboom, Jonathan Heek, Valentin De Bortoli, Kevin Patrick Murphy, and Tim Salimans. Diffusion models and gaussian flow matching: Two sides of the same coin. In *The Fourth Blogpost Track at ICLR 2025*, 2025. URL <https://openreview.net/forum?id=C8Yyg9wy0s>.
- Niklas Gebauer, Michael Gastegger, and Kristof Schütt. Symmetry-adapted generation of 3d point sets for the targeted discovery of molecules. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems* 32, pp. 7566–7578. Curran Associates, Inc., 2019.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020.
- Haokai Hong, Wanyu Lin, and Kay Chen Tan. Accelerating 3d molecule generation via jointly geometric optimal transport. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion for molecule generation in 3d. In *International conference on machine learning*, pp. 8867–8887. PMLR, 2022.

- 594 Ross Irwin, Alessandro Tibo, Jon Paul Janet, and Simon Olsson. SemlaFlow—Efficient 3d molecular
595 generation with latent attention and equivariant flow matching. In *The 28th International Con-*
596 *ference on Artificial Intelligence and Statistics*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=bee2G6pEh0)
597 [forum?id=bee2G6pEh0](https://openreview.net/forum?id=bee2G6pEh0).
- 598
599 Chaitanya K. Joshi, Xiang Fu, Yi-Lun Liao, Vahe Gharakhanyan, Benjamin Kurt Miller, Anuroop
600 Sriram, and Zachary W. Ulissi. All-atom Diffusion Transformers: Unified generative modelling
601 of molecules and materials, 2025. URL <https://arxiv.org/abs/2503.03965>.
- 602
603 P-J Kindermans and K-R Müller. Schnet—a deep learning architecture for molecules and materials.
604 *The Journal of chemical physics*, 148(24), 2018.
- 605
606 Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models,
607 2023. URL <https://arxiv.org/abs/2107.00630>.
- 608
609 Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Pro-*
610 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.
611 2837–2845, June 2021.
- 612
613 Nanye Ma, Mark Goldstein, Michael S. Albergo, Nicholas M. Boffi, Eric Vanden-Eijnden, and
614 Saining Xie. SiT: Exploring flow and diffusion-based generative models with scalable inter-
615 polant transformers. In *Computer Vision – ECCV 2024: 18th European Conference, Milan,*
616 *Italy, September 29–October 4, 2024, Proceedings, Part LXXVII*, pp. 23–40, Berlin, Heidelberg,
617 2024. Springer-Verlag. ISBN 978-3-031-72979-9. doi: 10.1007/978-3-031-72980-5_2. URL
618 https://doi.org/10.1007/978-3-031-72980-5_2.
- 619
620 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models,
621 2021. URL <https://openreview.net/forum?id=-NEXDKk8gZ>.
- 622
623 Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin,
624 Ji-Rong Wen, and Chongxuan Li. Large language diffusion models, 2025. URL [https://](https://arxiv.org/abs/2502.09992)
625 arxiv.org/abs/2502.09992.
- 626
627 Jingxiang Qu, Wenhan Gao, Jiaying Zhang, Xufeng Liu, Hua Wei, Haibin Ling, and Yi Liu. RISE:
628 Radius of influence based subgraph extraction for 3d molecular graph explanation. In *Proceedings*
629 *of the 42nd International Conference on Machine Learning*, 2025.
- 630
631 Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum
632 chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- 633
634 Severi Rissanen, Markus Heinonen, and Arno Solin. Generative modelling with inverse heat dis-
635 sipation. In *The Eleventh International Conference on Learning Representations*, 2023. URL
636 <https://openreview.net/forum?id=4PJUBT9f20l>.
- 637
638 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. High-
639 Resolution Image Synthesis with Latent Diffusion Models . In *2022 IEEE/CVF Conference on*
640 *Computer Vision and Pattern Recognition (CVPR)*, pp. 10674–10685, Los Alamitos, CA, USA,
641 June 2022. IEEE Computer Society. doi: 10.1109/CVPR52688.2022.01042. URL [https://](https://doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.01042)
642 doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.01042.
- 643
644 Victor Garcia Satorras, Emiel Hoogeboom, Fabian Bernd Fuchs, Ingmar Posner, and Max Welling.
645 E(n) equivariant normalizing flows. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman
646 Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL [https://](https://openreview.net/forum?id=N5hQI_RowVA)
647 openreview.net/forum?id=N5hQI_RowVA.
- 648
649 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
650 *preprint arXiv:2010.02502*, 2020.
- 651
652 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
653 Poole. Score-based generative modeling through stochastic differential equations. In *Intern-*
654 *ational Conference on Learning Representations*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=PXTIG12RRHS)
655 [forum?id=PXTIG12RRHS](https://openreview.net/forum?id=PXTIG12RRHS).

- 648 Yuxuan Song, Jingjing Gong, Minkai Xu, Ziyao Cao, Yanyan Lan, Stefano Ermon, Hao Zhou,
649 and Wei-Ying Ma. Equivariant flow matching with hybrid probability transport for 3d molecule
650 generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL
651 <https://openreview.net/forum?id=hHUZ5V9XFu>.
- 652 Yuxuan Song, Jingjing Gong, Hao Zhou, Mingyue Zheng, Jingjing Liu, and Wei-Ying Ma. Unified
653 generative modeling of 3d molecules with bayesian flow networks. In *The Twelfth International
654 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?
655 id=NSVtmnzeRB](https://openreview.net/forum?id=NSVtmnzeRB).
- 656 Michał J Tyszkiewicz, Pascal Fua, and Eduard Trulls. GECCO: Geometrically-conditioned point
657 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
658 pp. 2128–2138, 2023.
- 659 Lemeng Wu, Chengyue Gong, Xingchao Liu, Mao Ye, and qiang liu. Diffusion-based molecule
660 generation with informative prior bridges. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,
661 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL
662 <https://openreview.net/forum?id=TJUNtiZiTKE>.
- 663 Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. GeoDiff: A
664 geometric diffusion model for molecular conformation generation. In *International Confer-
665 ence on Learning Representations*, 2022. URL [https://openreview.net/forum?id=
666 PzcvxEMzvQC](https://openreview.net/forum?id=PzcvxEMzvQC).
- 667 Minkai Xu, Alexander S Powers, Ron O Dror, Stefano Ermon, and Jure Leskovec. Geometric latent
668 diffusion models for 3d molecule generation. In *International Conference on Machine Learning*,
669 pp. 38592–38610. PMLR, 2023.
- 670 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
671 diffusion models. In *2023 IEEE/CVF Conference on International Conference on Computer
672 Vision (ICCV)*, pp. 3813–3824, 2023.
- 673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

LLM USAGE

LLMs have been used to assist with polishing the writing (e.g., grammar, typos, readability) and to help with writing code when given clear and detailed instructions. All LLM-generated texts and codes are verified before use by the authors (humans). All conceptual innovations, methodological designs, and research contributions are developed solely by the authors.

A SIZE-INDUCED INCONSISTENCIES

A.1 GENERATION QUALITY ACROSS DIFFERENT SIZES

Experimental Details. We assess the difference in generation quality by sampling 3,000 molecules of each size with 3 SOTA diffusion models for 3D molecular generation, namely EDM, RADDM, and GeoLDM. We follow the exact same setups and directly use the pretrained weights provided by the respective works.

Additional Results. We present size-induced inconsistencies in the **molecule stability of generated molecules** on QM9 with RADDM and GeoLDM in Fig. 4 and Fig. 5, respectively, **size-induced inconsistencies in the validity of generated molecules on QM9** for all 3 diffusion models in Fig. 6, Fig. 7, and Fig. 8, respectively, and **size-induced inconsistencies in atom stability** on GEOM-Drugs for all 3 diffusion models in Fig. 9. **It is clear that size-induced inconsistencies are present in various diffusion pipelines for 3D molecular generation.** Note that we do not present validity results for the GEOM-Drugs dataset because the validity is already very close to 1 across all molecule sizes. This aligns with the observation in Hooigeboom et al. (2022) that validity remains near 1 when evaluated over all possible sizes.

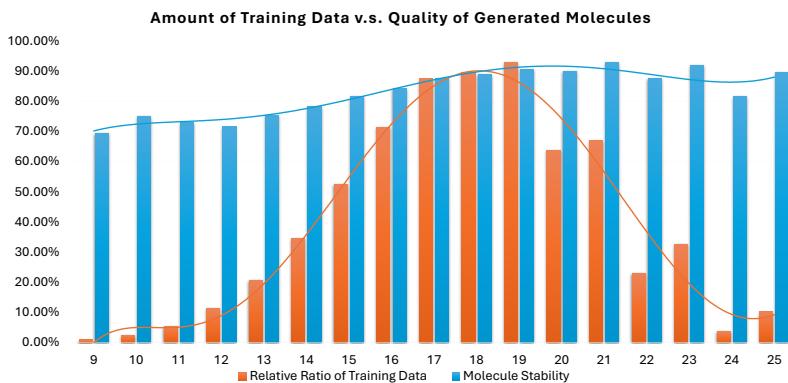


Figure 4: **Molecule stability of generated molecules** (y -axis) versus molecular size (x -axis) with RADDM on QM9.

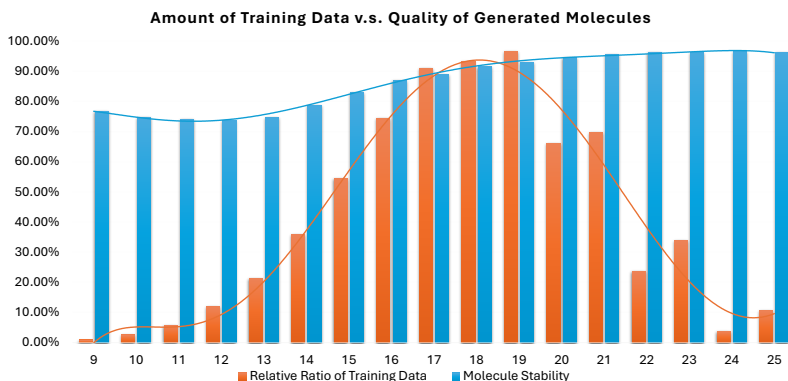
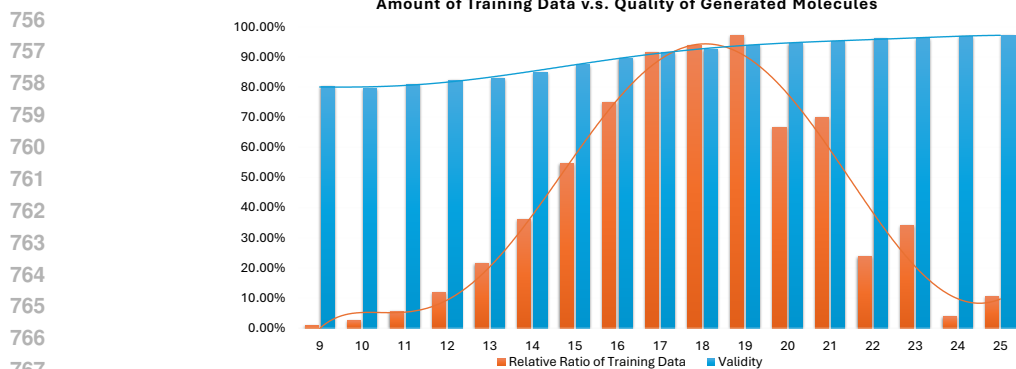


Figure 5: **Molecule stability of generated molecules** (y -axis) versus molecular size (x -axis) with GeoLDM on QM9.

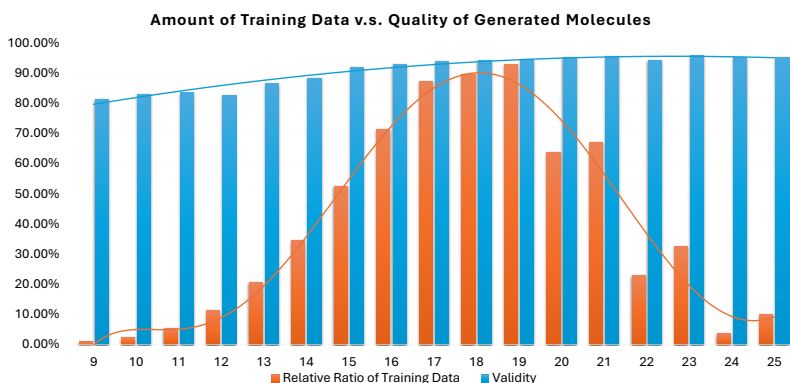
A.2 GENERATION TARGET ALIGNMENT ACROSS DIFFERENT SIZES

Experimental Details. We sample 1,000 generation trajectories for each size with 3 SOTA diffusion models for 3D molecular generation, namely EDM, RADDM, and GeoLDM. For latent diffusion



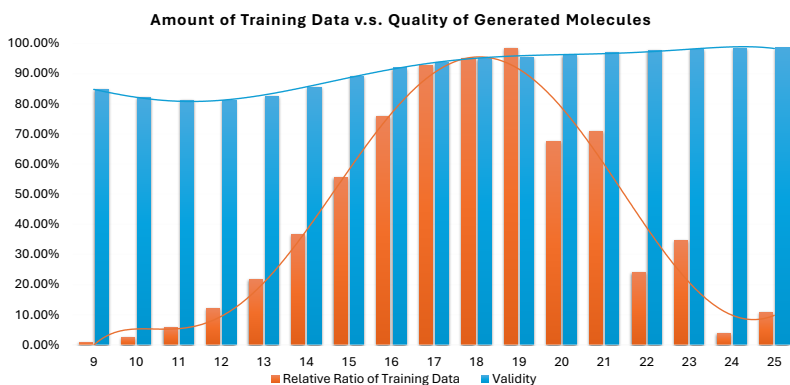
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781

Figure 6: Validity of generated molecules (y -axis) versus molecular size (x -axis) with EDM on QM9.



782
783
784
785
786
787
788
789
790
791
792
793
794
795

Figure 7: Validity of generated molecules (y -axis) versus molecular size (x -axis) with RADM on QM9.



796
797
798
799
800
801
802
803
804
805
806

Figure 8: Validity of generated molecules (y -axis) versus molecular size (x -axis) with GeoLDM on QM9.

807
808
809

models (RADM and GeoLDM), the coordinates are represented directly in the 3D Euclidean latent space. For atom types, we obtain the concrete atom types by passing their latent representations through the decoder, which are then used to compute the alignment ratio. Note that tracking the full trajectory is memory intensive, so only 100 trajectories are sampled each time, and we average over 10 runs. It is observed that the alignment convergence patterns remain consistent across different runs. We follow the exact same setups and directly use the pretrained weights provided by the respective works.

Additional Results. We present size-induced inconsistencies in the generation target alignment for EDM in Fig. 12, for RADM in Fig. 13, and for GeoLDM in Fig. 14. Note that these figures are placed at the end of the paper due to space constraints. **It is clear that there are size-induced inconsistencies for all three SOTA diffusion models.** The inconsistencies in atom type alignment

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

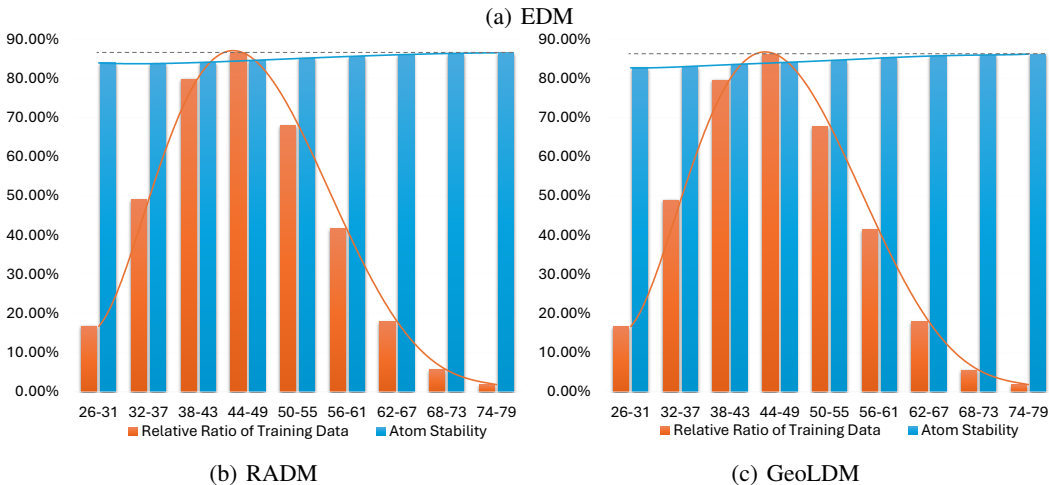
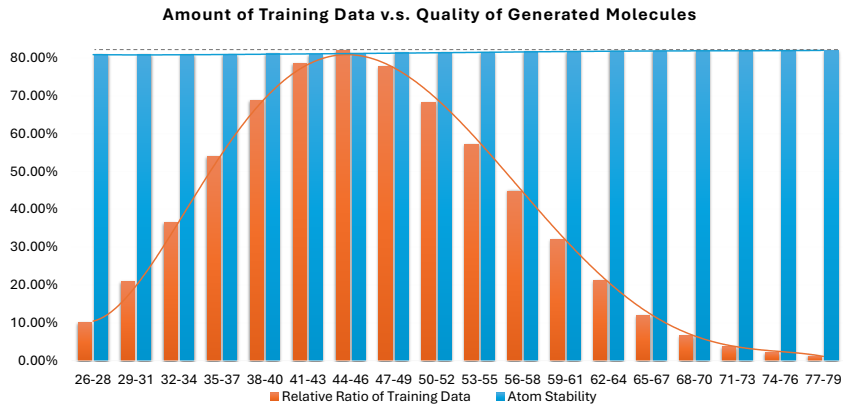


Figure 9: Atom stability of generated molecules (y -axis) versus molecular size (x -axis) with all 3 SOTA diffusion models on GEOM. For visualization, molecular sizes are grouped into inclusive bins. Note that the reported metric for sampling quality is now atom stability. While the numerical differences appear smaller compared to molecular stability, they are nonetheless pronounced and important.

observed on the GEOM-Drugs dataset for GeoLDM are less pronounced than that in EDM and RADM, because atom types in latent models are determined much slower, after the shapes have already stabilized; this is potentially because the autoencoder is untrained and the latent space is highly compact, as indicated by the authors in their GitHub repository with the provided pretrained weights. However, we can still observe the inconsistencies, especially by comparing the values of small sizes (e.g. 10 and 12) with that of large sizes (e.g. 45 and larger).

A.3 STP: A LEARNING RESPECTIVE

For notational simplicity, we consider only the 3D coordinates so that $\mathbf{x} \in \mathbb{R}^{3 \times N}$. In Sec. 2.1, the denoiser network learns to map the noisy sample to the noise added; mathematically, $\phi : (z_t, t) \mapsto \epsilon$, where $z_t^x = \alpha_t \mathbf{x} + \sigma_t \epsilon$. However, this denoiser must implicitly learn a size-dependent scaling to account for the variation in molecular spatial scales. For example, at $t = 1$, starting from pure Gaussian noise, the denoiser needs to learn to contract coordinates for smaller molecules and expand them for larger ones, as illustrated in Fig. 10. Depending on the normalization of the data, this may also manifest as expanding more or contracting less for larger molecules. This inconsistency in the learning dynamics leads to a size-induced inconsistency in the generation quality as demonstrated in Fig. 2 in the main paper. StP mitigates this issue by scaling the prior Gaussian distribution so that the generative process learns consistent expansion or contraction across molecular sizes.

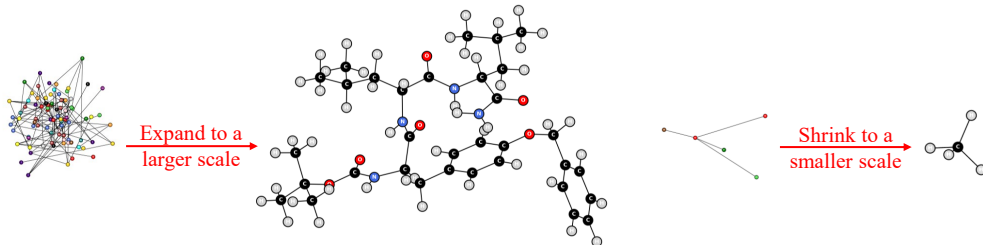


Figure 10: Illustration of different spatial scales between a molecule of size 91 and a molecule of size 5, compared to their respective Gaussian noise of the same size. The generative process and the denoiser network are trained to map pure Gaussian noise to realistic and valid molecular structures. However, they must learn a size-dependent behavior, contracting coordinates for smaller molecules and expanding them for larger ones. This requirement introduces inconsistency in the learning process, leading to size-induced inconsistencies in generation quality. StP mitigates this issue by scaling the spatial extent of the Gaussian prior for different sizes so that the generative model learns the same dynamics consistently across molecular sizes.

In Fig. 11, we provide the same per-size visualization as in Fig. 2 of the main paper, but with StP, to illustrate its effectiveness in reducing size-induced inconsistency. As shown in Fig. 11, the generation quality with StP remains roughly the same for molecular sizes greater than 17, whereas it continues to increase without StP. In addition, StP significantly improves the generation quality for smaller molecules. Larger molecules still exhibit better performance under the same data scarcity, potentially because the backbone neural networks generalize more effectively for larger molecules due to the increased number of interatomic interactions.

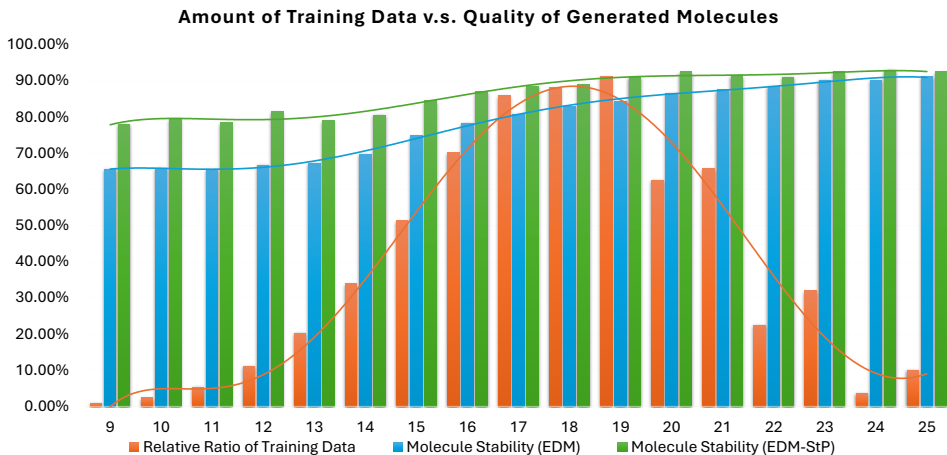


Figure 11: Molecule stability of generated molecules (y -axis) versus molecular size (x -axis) with EDM and EDM-StP on the QM9 dataset (Ramakrishnan et al., 2014). Clearly, StP mitigates the size-induced inconsistency in generation quality. The generation quality with StP remains roughly the same for molecular sizes greater than 17, whereas it continues to increase without StP. In addition, StP significantly improves the generation quality for smaller molecules.

A.4 CONVEX HULL

Given a set of 3D coordinates $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \in \mathbb{R}^{N \times 3}$, the *convex hull* is the smallest convex polytope that contains all points \mathbf{x}_i . Formally, the convex hull is defined as

$$\text{Hull}(\{\mathbf{x}_1, \dots, \mathbf{x}_N\}) = \left\{ \sum_{i=1}^N \lambda_i \mathbf{x}_i \mid \lambda_i \geq 0, \sum_{i=1}^N \lambda_i = 1 \right\}. \quad (13)$$

A key property is that the convex hull can always be described by a subset of the original points, called vertices or extreme points. These are the points that cannot be written as convex combinations of the others. **Geometrically, they lie on the “outer surface” of the configuration. All other points in the set $\{x_i\}$ are contained strictly inside the convex hull** and can be expressed as convex combinations of the vertices. The readers are referred to De Berg et al. (2008) for further details on convex hulls and the algorithms to compute convex hull vertices. In this work, to find convex hull vertices, we use the algorithm implemented in the *scipy* Python library with *scipy.spatial.ConvexHull()*. Specifically, it performs the QHull (divide and conquer) algorithm.

B DISCUSSION ON LATENT DIFFUSION MODELS

Latent diffusion models encode molecular coordinates and atom features (\mathbf{x}, \mathbf{h}) into a lower-dimensional latent space. Existing approaches, such as GeoLDM and RADM, still treat coordinates and atom features separately, keeping coordinates in the 3D Euclidean space within the latent representation as $\mathbf{x}' \in \mathbb{R}^{(N-1) \times 3}$, $\mathbf{h}' \in \mathbb{R}^{N \times d'}$, where typically $d' < d$.

GeoLDM aims to respect the rotational symmetries inherent in molecular data (translation symmetry can always easily be ensured by zero-mean). Specifically, GeoLDM employs an equivariant encoder $\mathcal{E}(\cdot)$ and decoder $\mathcal{D}(\cdot)$ such that

$$\mathcal{E}(R \cdot \mathbf{x}) = R \cdot \mathcal{E}(\mathbf{x}) \quad \text{and} \quad \mathcal{D}(R \cdot \mathbf{x}) = R \cdot \mathcal{D}(\mathbf{x}), \quad (14)$$

where R is a representation of the $SO(3)$ group. For simplicity, we ignore translations since they can be easily handled by centering coordinates to zero mean, and we omit atomic features since they are scalar values that remain unchanged under Euclidean transformations. In essence, equivariance guarantees that when the molecular point cloud is rotated, the latent representation rotates in the same way. To achieve this, the latent space must carry a valid representation of $SO(3)$. The smallest nontrivial irreducible representation (irrep) of $SO(3)$ is the 3D vector representation. **Consequently, maintaining rotational equivariance requires the latent representation of coordinates to remain in 3D coordinate form.**

On the other hand, RADM aims to standardize molecular configurations into a fixed orientation in an unsupervised manner. After standardization, the 3D coordinates can be encoded into lower dimensional spaces. However, RADM also still remains in 3D coordinate form. Despite consideration of the rotation symmetry of the molecular data, the final generation goal is to generate molecules with their full 3D atomic configurations. Molecular configurations are very sensitive: small errors in 3D coordinates can lead to totally invalid or unstable molecules. If 3D coordinates were compressed into a lower-dimensional latent space (e.g., one or two dimensions), this would inevitably result in loss of geometric information. Since such information cannot be perfectly recovered when mapping back to 3D space, the generated structures would risk being distorted or chemically unrealistic. **Therefore, to minimize the reconstruction error in atomic coordinates, it is desirable to have the coordinates still represented in 3D in the latent space.**

Due to similar reasons, it is also undesirable for the latent space to rely on overly complicated transformations of the coordinates. If the encoding step introduces highly nonlinear or abstract latent representations of molecular geometry, it may obscure or distort the essential spatial relationships between atoms. Since the generative model must ultimately recover precise 3D configurations, such complexity can make the reconstruction inaccurate. While small inaccuracies may be tolerable in other application domains, such as image generation, they are intolerable for molecular generation. **Therefore, it is most effective to maintain a latent representation that stays close to the original 3D coordinate configuration. In fact, in both GeoLDM and RADM, it can be seen that, while the atomic features are encoded into lower-dimensions, the coordinates rarely change in the latent space.** Quantitatively, we can measure the relative l_2 difference between the original coordinates and their latent representations: $\frac{\|\mathbf{x} - \mathcal{E}(\mathbf{x})\|_2}{\|\mathbf{x}\|_2}$. We report the average relative l_2 difference over the entire training dataset in Table 4. Clearly, the encoders are almost the identity function for coordinates; the atomic coordinates rarely change in latent space.

Table 4: Relative L_2 difference between the original atomic coordinates and the latent coordinates. For RADM, we also apply the rotator to the original coordinates to ensure rotational alignment, so that the error reflects structural discrepancies rather than arbitrary orientation differences.

Auto-Encoder	QM9	GEOM-Drugs
RADM	2.87%	0.03%
GeoLDM	1.95%	1.23%

C PROOF TO PROPOSITION

Proposition 3.2. *The marginal distributions at any time t of Scaling the Prior are equivalent to those of normalizing the data \mathbf{x} by a factor of $\frac{1}{\gamma_N}$ before the forward process, and subsequently unnormalizing after sampling.*

Proof. Mathematically, it is equivalent to proving the following:

$$q_\gamma(\mathbf{z}_t^x | \mathbf{x}) = q\left(\frac{1}{\gamma_N}\mathbf{z}_t^x \mid \frac{1}{\gamma_N}\mathbf{x}\right). \quad (15)$$

Note that:

$$\begin{aligned} q_\gamma(\mathbf{z}_t^x | \mathbf{x}) &= \mathcal{N}\left(\mathbf{z}_t^x; \sqrt{\alpha_t}\mathbf{x}, \gamma_N^2\sigma_t^2\mathbf{I}\right) \\ &\propto \exp\left(-\frac{\|\mathbf{z}_t^x - \sqrt{\alpha_t}\mathbf{x}\|^2}{2\gamma_N^2\sigma_t^2}\right) \\ &= \exp\left(-\frac{1/\gamma_N^2\|\mathbf{z}_t^x - \sqrt{\alpha_t}\mathbf{x}\|^2}{2\sigma_t^2}\right) \\ &= \exp\left(-\frac{\left\|\frac{1}{\gamma_N}\mathbf{z}_t^x - \frac{1}{\gamma_N}\sqrt{\alpha_t}\mathbf{x}\right\|^2}{2\sigma_t^2}\right) \\ &\propto \mathcal{N}\left(\frac{1}{\gamma_N}\mathbf{z}_t^x; \frac{\sqrt{\alpha_t}}{\gamma_N}\mathbf{x}, \sigma_t^2\mathbf{I}\right) \\ &= q\left(\frac{1}{\gamma_N}\mathbf{z}_t^x \mid \frac{1}{\gamma_N}\mathbf{x}\right). \end{aligned} \quad (16)$$

□

D EXPERIMENTAL DETAILS AND ADDITIONAL RESULTS

D.1 EFFICIENT SAMPLING

In the implementation of sampling in EDM, as well as in RADM and GeoLDM which inherit from EDM, the sampling process becomes unnecessarily slow when the maximum molecular size is much larger than that of most molecules in the dataset (e.g. the average molecular size in GEOM-Drugs is 44.4 while the maximum is 181). This is unnecessarily slow because during sampling, the method constructs graphs of the dataset’s maximum size N_{\max} , padding with dummy nodes for parallelization as described in the pseudo code in Algorithm 1. **However, within a given batch, the actual maximum size $N_{\text{batch_max}}$ is often much smaller than N_{\max} . Noticing this bottleneck, we can accelerate the sampling process by constructing graphs only up to the batch maximum size $N_{\text{batch_max}}$. More importantly, when sampling a large number of molecules, we can first sample molecular sizes and then sort them so that each batch contains molecules with similar sizes, thereby avoiding unnecessary computation.** If array storage and future parallel process are needed, the generated results can be padded with dummy nodes to be the maximum size from the overall sampled molecular sizes or the dataset max N_{\max} for simplicity. The pseudo code for accelerated sampling is presented in Algorithm 2. Notice that this sampling is mathematically equivalent to the original sampling in the implementation in EDM; it is just an efficient implementation.

Table 5: Generation time per sample for original vs. efficient methods.

Model	Original (s)	Efficient (s)	Speedup (\times)
EDM	12.94	1.06	12.21
RADM	2.11	0.59	3.58
GeoLDM	13.03	1.07	12.18

We follow the standard sampling practice used in all three SOTA diffusion models. Specifically, we generate 10,000 molecules with a batch size of 100 on GEOM-Drugs. Noticeably, this accelerates the generation by $12.21\times$ for EDM, $3.58\times$, for RADM, almost $12.18\times$ for GeoLDM. We provide the statistics in Table 5. The timing results are recorded on a single NVIDIA RTX A6000 48GB Graphics Card. **This acceleration is crucial, as it enables size-dependent studies in this work that would otherwise be computationally prohibitive.**

Training. Note that during training, the graphs are constructed based on the batch maximum molecular size; thus, it is already efficient. Although it is possible to not randomly shuffle the training set and sort it, this may introduce systematic bias in the training process. Moreover, in most cases, the batch maximum is still relatively small and all molecular sizes within the batch are close, so the additional computational overhead is tolerable.

Algorithm 1 Original EDM Sampling Implementation

Require: M (total sample size), M_{batch} (batch size), P_{size} (distribution of molecule sizes), N_{max} (maximum molecular size), ψ (trained diffusion model with sampling method)

- 1: Initialize `all_samples` $\leftarrow []$
- 2: `num_batches` $\leftarrow M/M_{\text{batch}}$ \blacktriangleright Assume M_{batch} divides M
- 3: **for** $i = 1$ to `num_batches` **do**
- 4: $\mathbf{N}_{\text{batch}} \in \mathbb{Z}^{M_{\text{batch}}} \leftarrow \text{Sample}(P_{\text{size}}, M_{\text{batch}})$ \blacktriangleright Sample M_{batch} times from P_{size}
- 5: Construct M_{batch} initial graphs \mathcal{G} , all of size N_{max} with node masks based on $\mathbf{N}_{\text{batch}}$
- 6: `batch_samples` $\leftarrow \psi.\text{sample}(\mathcal{G})$
- 7: `all_samples` $\leftarrow \text{all_samples} \cup \text{batch_samples}$ \blacktriangleright Append batch samples
- 8: **end for**
- 9: **return** `all_samples`

Algorithm 2 Efficient Sampling Implementation

Require: M (total sample size), M_{batch} (batch size), P_{size} (distribution of molecule sizes), N_{max} (maximum molecular size), ψ (trained diffusion model with sampling method)

- 1: Initialize `all_samples` $\leftarrow []$
- 2: $\mathbf{N} \in \mathbb{Z}^M \leftarrow \text{Sort}(\text{Sample}(P_{\text{size}}, M))$ \blacktriangleright Sample M times from P_{size} and sort
- 3: `num_batches` $\leftarrow M/M_{\text{batch}}$ \blacktriangleright Assume M_{batch} divides M
- 4: **for** $i = 1$ to `num_batches` **do**
- 5: $\mathbf{N}_{\text{batch}} \leftarrow \mathbf{N}[(i-1) \cdot M_{\text{batch}} + 1 : i \cdot M_{\text{batch}}]$ \blacktriangleright Get a batch of molecular sizes by indexing
- 6: $N_{\text{batch_max}} = \mathbf{N}_{\text{batch}} \leftarrow \mathbf{N}[i \cdot M_{\text{batch}}]$ \blacktriangleright Maximum molecular size in the current batch
- 7: Construct M_{batch} initial graphs \mathcal{G} , all of size $N_{\text{batch_max}}$ with node masks based on $\mathbf{N}_{\text{batch}}$
- 8: `batch_samples` $\leftarrow \psi.\text{sample}(\mathcal{G})$
- 9: `batch_samples` $\leftarrow \text{Pad}(\text{batch_samples}, \mathbf{N}[M])$ \blacktriangleright Pad generated samples to max size
- 10: `all_samples` $\leftarrow \text{all_samples} \cup \text{batch_samples}$ \blacktriangleright Append batch samples
- 11: **end for**
- 12: **return** `all_samples`

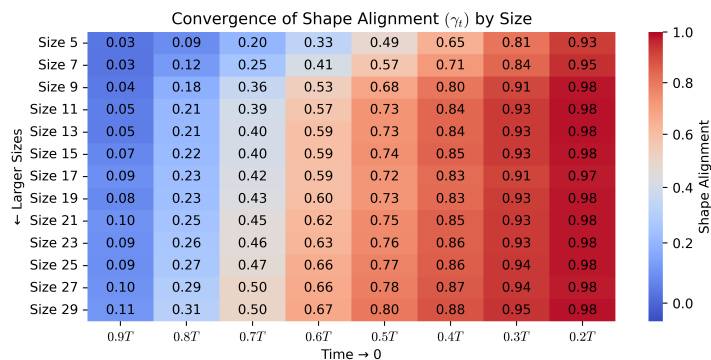
D.2 CONDITIONAL GENERATION

Following prior work, EDM (Hoogeboom et al., 2022), the objective of conditional generation is to generate molecules with target properties. Mathematically, this conditional generation in EDM is to generate $\mathbf{x}, \mathbf{h} \sim p(\mathbf{x}, \mathbf{h} \mid c)$ given some desired property c . The optimization lower bound for the conditional case can be defined as $\log p(\mathbf{x}, \mathbf{h} \mid c) \geq \mathcal{L}_{c,0} + \mathcal{L}_{c,\text{base}} + \sum_{t=1}^T \mathcal{L}_{c,t}$. The main difference is that the denoiser $\hat{\epsilon}_t = \phi(\mathbf{z}_t, [t, c])$ takes an additional input, a property c , which is concatenated to the node features. Given a trained conditional model, sampling is done by first sampling the

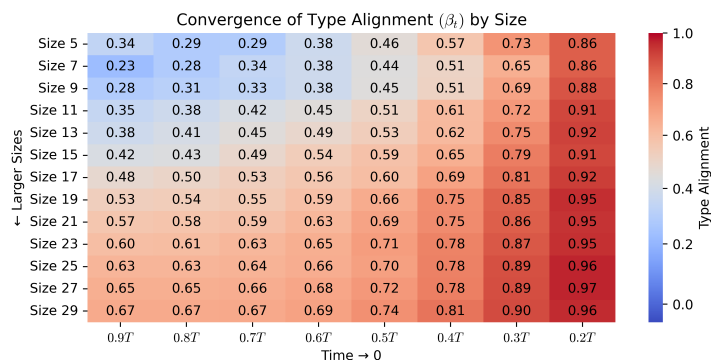
1080 molecular size (number of atoms M) and a property value c from a distribution $c, M \sim p(c, M)$
1081 that is inferred from the training partition. Then, given c and M , we can generate molecules using
1082 conditional models.

1083 **Implementation Details.** We follow the exact experimental setup as in EDM; RADM and Ge-
1084 oLDM all follow the exact same setup as well. Specifically, the QM9 training set is evenly split:
1085 an EGNN is trained on the first half as a property predictor, while the generative model is trained
1086 on the second half. After training, the predictor is used to evaluate the molecules generated by the
1087 model. We use the exact same EGNN model and hyperparameters as in EDM to train the property
1088 predictor.
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

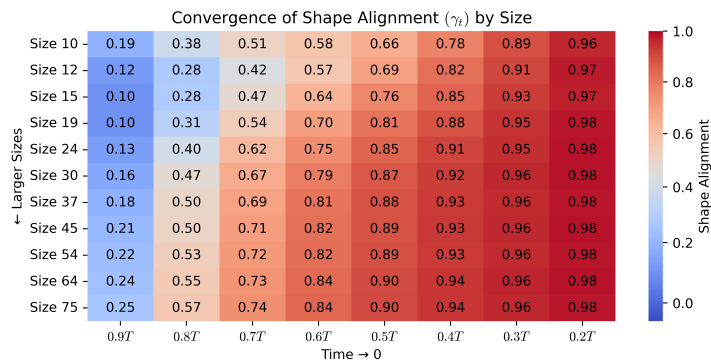
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187



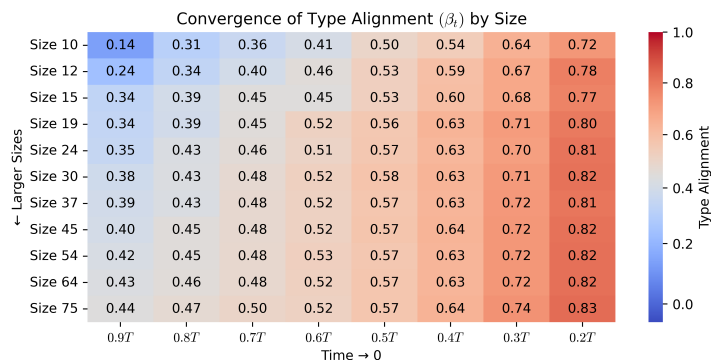
(a) Shape alignment on QM9 with EDM.



(b) Type alignment on QM9 with EDM.



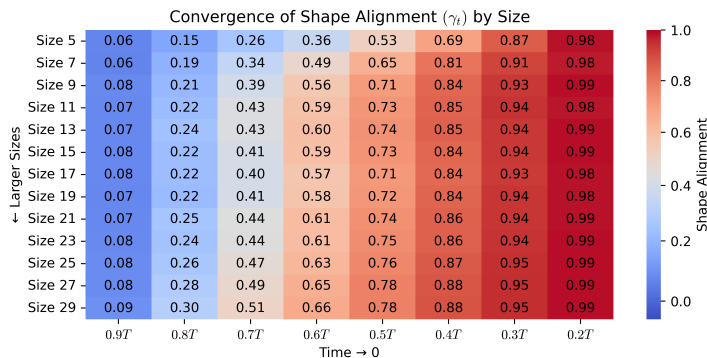
(c) Shape alignment on GEOM-Drugs with EDM.



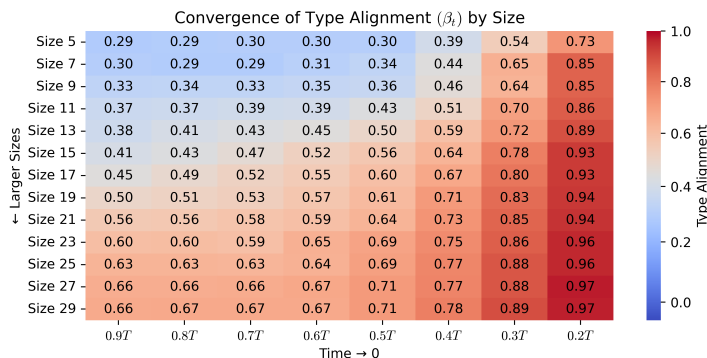
(d) Type alignment on GEOM-Drugs with EDM.

Figure 12: Convergence of generation target alignment with EDM.

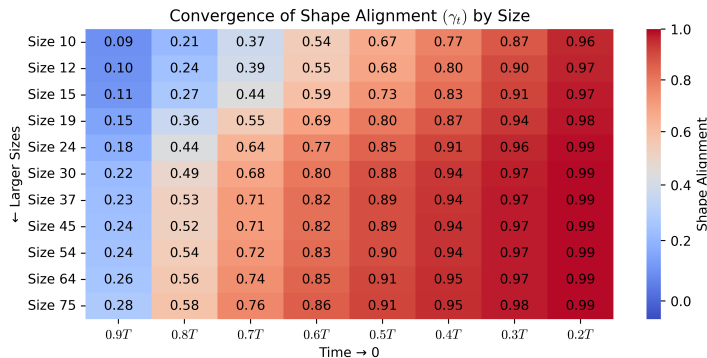
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241



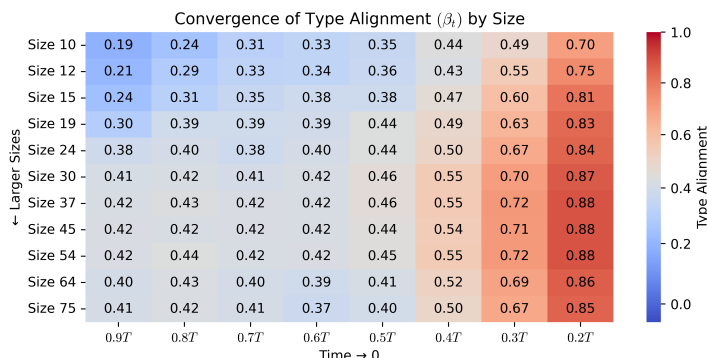
(a) Shape alignment on QM9 with RADM.



(b) Type alignment on QM9 with RADM.



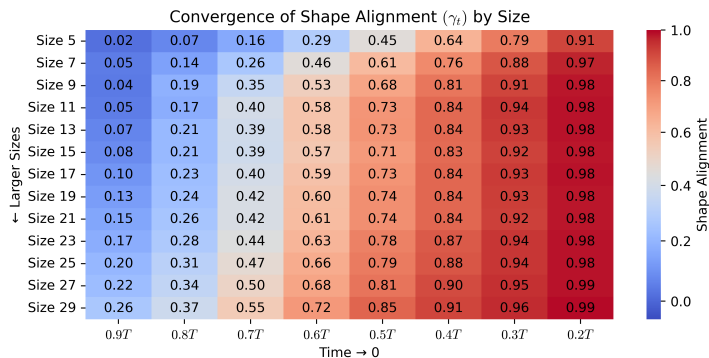
(c) Shape alignment on GEOM-Drugs with RADM.



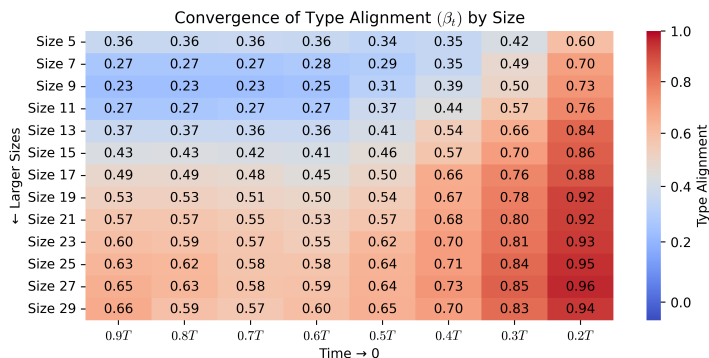
(d) Type alignment on GEOM-Drugs with RADM.

Figure 13: Convergence of generation target alignment with RADM.

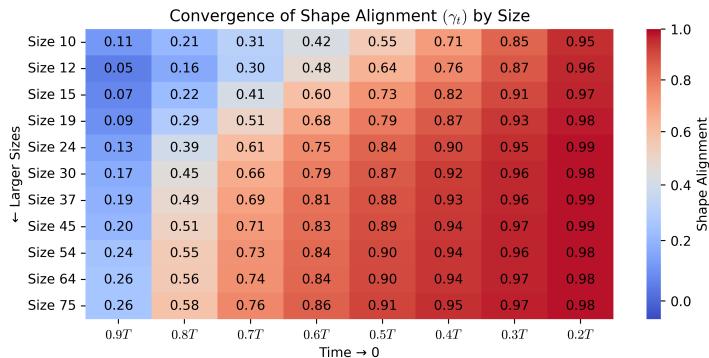
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295



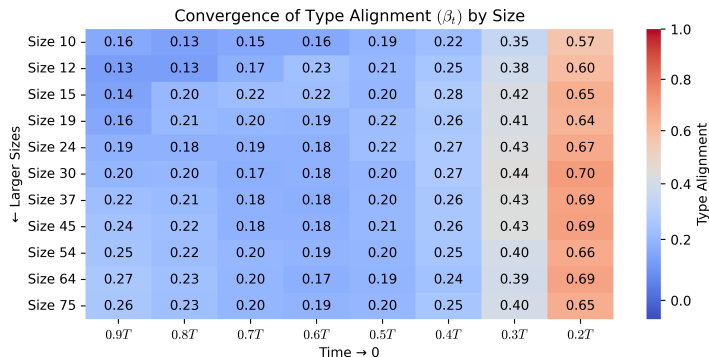
(a) Shape alignment on QM9 with GeoLDM.



(b) Type alignment on QM9 with GeoLDM.



(c) Shape alignment on GEOM-Drugs with GeoLDM.



(d) Type alignment on GEOM-Drugs with GeoLDM.

Figure 14: Convergence of generation target alignment with GeoLDM.