
Equivariant Neural Diffusion for Molecule Generation

François Cornet¹ Grigory Bartosh² Mikkel N. Schmidt¹ Christian A. Naesseth²

Abstract

We introduce Equivariant Neural Diffusion (END), a novel diffusion model for molecule generation in 3D that is equivariant to Euclidean transformations. Compared to current state-of-the-art equivariant diffusion models, the key innovation in END lies in its learnable forward process for enhanced generative modelling. Rather than pre-specified, the forward process is parameterized through a time- and data-dependent transformation that is equivariant to rigid transformations. Through a series of experiments on standard molecule generation benchmarks, we demonstrate that END improves on several strong baselines for both unconditional and conditional generation.

1. Introduction

The discovery of novel chemical compounds with relevant properties is critical to a number of scientific fields, such as drug discovery and materials design (Merchant et al., 2023). However, due to the large size and complex structure of the chemical space (Ruddigkeit et al., 2012), which combines continuous and discrete features, it is notably difficult to search. Additionally, *ab-initio* quantum mechanics methods for computing target properties are often computationally expensive, preventing brute-force enumeration. While some of these heavy computations can be amortized through learned surrogates, the need for innovative search methods remains, and generative models have recently emerged as a promising avenue (Anstine & Isayev, 2023). Generative models can learn complex data distributions, that, in turn, can be sampled from to obtain novel samples similar to the original data. Compared to other data modalities such as images or text, molecules have to adhere to strict chemical rules, and obey the symmetries of the 3D space.

¹Technical University of Denmark ²University of Amsterdam. Correspondence to: François Cornet <frjc@dtu.dk>, Grigory Bartosh <g.bartosh@uva.nl>, Mikkel N. Schmidt <mmsc@dtu.dk>, Christian A. Naesseth <c.a.naesseth@uva.nl>.

Accepted at the AI4Science Workshop of the 41st International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

Currently, the most promising directions for molecule generation in 3D are either auto-regressive models (Gebauer et al., 2019; 2022; Luo & Ji, 2022; Daigavane et al., 2024), building molecules one atom at a time, or Diffusion Models (DM) (Hoogeboom et al., 2022; Vignac et al., 2023; Le et al., 2024) that learn to revert a corruption mechanism that transforms the data distribution into noise. As both approaches directly operate in 3D space, they can leverage the numerous architectures designed for machine learned force fields (Unke et al., 2021), which were carefully developed to encode the symmetries inherent to the data (Schütt et al., 2017; 2021; Batzner et al., 2022; Batatia et al., 2022).

DM have not only been successful at molecule generation, but also on a variety of other data modalities (Yang et al., 2023). Nevertheless, most existing DM pre-specify the forward process, forcing the reverse process to comply with it. A recent line of work has sought to overcome that limitation and improve generation by replacing the fixed forward process with a learnable one (Bartosh et al., 2023; Nielsen et al., 2024; Bartosh et al., 2024).

Contributions In this paper, we present Equivariant Neural Diffusion (END), a novel diffusion model for molecule generation in 3D that (1) is equivariant to Euclidean transformations, and (2) features a learnable forward process. We demonstrate competitive performance in unconditional molecule generation on the QM9 and GEOM-Drugs benchmarks. For conditional generation driven by composition and substructure constraints, our approach exhibits a substantial performance gain compared to existing equivariant diffusion models, underscoring the utility of a learned forward model for effective conditional molecule generation.

2. Background

We begin by establishing the necessary background for generative modeling of geometric graphs. We first introduce the data representation and its inherent symmetries. We then discuss Diffusion Models (DM), and more specifically Equivariant Diffusion Model (EDM) (Hoogeboom et al., 2022). Finally, we present the Neural Flow Diffusion Models (NFDM) framework (Bartosh et al., 2024).

2.1. Equivariance

Molecules as geometric graphs in E(3) We consider geometric graphs embedded in 3-dimensional Euclidean space that represent molecules. Formally, each atomistic system can be described by a tuple $\mathbf{x} = (\mathbf{r}, \mathbf{h})$, where $\mathbf{r} = (\mathbf{r}_1, \dots, \mathbf{r}_M) \in \mathbb{R}^{M \times 3}$ form a collection of vectors in 3D representing the coordinates of the atoms, and $\mathbf{h} = (\mathbf{h}_1, \dots, \mathbf{h}_M) \in \mathbb{R}^{M \times D}$ are the associated scalar features (e.g. atomic types or charges).

When dealing with molecules, we are particularly interested in $E(3)$, the Euclidean group in 3 dimensions, generated by translations, rotations and reflections. Each group element in $E(3)$ can be represented as a combination of a translation vector $\mathbf{t} \in \mathbb{R}^3$ and an orthogonal matrix $\mathbf{R} \in O(3)$ encoding rotation or reflection. While scalar features \mathbf{h} remain invariant, coordinates \mathbf{r} transform under translation, rotation and reflection as $\mathbf{R}\mathbf{r} + \mathbf{t} = (\mathbf{R}\mathbf{r}_1 + \mathbf{t}, \dots, \mathbf{R}\mathbf{r}_M + \mathbf{t})$.

Equivariant functions A function $f: \mathcal{X} \rightarrow \mathcal{Y}$ is said to be equivariant to the action of a group G , or G -equivariant, if $g \cdot f(\mathbf{x}) = f(g \cdot \mathbf{x}), \forall g \in G$. It is said to be G -invariant, if $f(\mathbf{x}) = f(g \cdot \mathbf{x}), \forall g \in G$. In the case of a function $f: (\mathbb{R}^{M \times 3} \times \mathbb{R}^{M \times D}) \rightarrow (\mathbb{R}^{M \times 3} \times \mathbb{R}^{M \times D})$ operating on geometric graphs, the function is said to be $E(3)$ -equivariant if,

$$\mathbf{R}\mathbf{y}^{(r)} + \mathbf{t}, \mathbf{y}^{(h)} = f(\mathbf{R}\mathbf{r} + \mathbf{t}, \mathbf{h}),$$

$\forall \mathbf{R} \in O(3)$ and $\mathbf{t} \in \mathbb{R}^3$, where $\mathbf{y}^{(r)}$ and $\mathbf{y}^{(h)}$ denote the output related to \mathbf{r} and \mathbf{h} respectively. There exists a large variety of graph neural network architectures designed to be equivariant to the Euclidean group (Schütt et al., 2017; 2021; Batzner et al., 2022; Batatia et al., 2022).

Equivariant distributions A conditional distribution $p(\mathbf{y}|\mathbf{x})$ is equivariant to rotations and reflections when $p(\mathbf{y}|\mathbf{x}) = p(\mathbf{R}\mathbf{y}|\mathbf{R}\mathbf{x}), \forall \mathbf{R} \in O(3)$, while a distribution is said to be invariant when $p(\mathbf{x}) = p(\mathbf{R}\mathbf{x}), \forall \mathbf{R} \in O(3)$. Regarding translation, it is not possible to have a translation-invariant non-zero distribution, as it would require that $p(\mathbf{x}) = p(\mathbf{x} + \mathbf{t}), \forall \mathbf{t} \in \mathbb{R}^3, \mathbf{x} \in \mathbb{R}^{M \times 3}$, which would mean that $p(\mathbf{x})$ cannot integrate to 1 (Garcia Satorras et al., 2021). However, a translation invariant distribution can be constructed in the linear subspace R where the centre of gravity is fixed to $\mathbf{0}$ (i.e. zero CoG subspace): $R = \{\mathbf{r} \in \mathbb{R}^{M \times 3} : \frac{1}{M} \sum_{i=1}^M \mathbf{r}_i = \mathbf{0}\}$ (Xu et al., 2022). As R can be shown to be intrinsically equivalent to $\mathbb{R}^{(M-1) \times 3}$ (Bao et al., 2023), we will consider in what follows that \mathbf{r} is defined in $\mathbb{R}^{(M-1) \times 3}$ for ease of notation.

2.2. Equivariant Diffusion Models

Diffusion Models (DM) (Sohl-Dickstein et al., 2015; Ho et al., 2020) are generative models that learn distributions through a hierarchy of latent variables, corresponding to per-

turbed versions of the data at increasing noise scales. DM consist of a forward and a reverse (or generative) process. The Equivariant Diffusion Model (EDM) (Hooeboom et al., 2022) is a particular instance of a DM, where the learned marginal $p_\theta(\mathbf{x})$ is made invariant to the action of translations, rotations and reflections by construction. Intuitively, this means that the likelihood of a given molecule does not depend on its orientation.

Forward process The forward process perturbs samples from the data distribution, $\mathbf{x} \sim q(\mathbf{x})$, over time through noise injection, resulting in a trajectory of latent variables $(\mathbf{z}_t)_{t \in [0,1]}$, conditional on \mathbf{x} . The conditional distribution for $(\mathbf{z}_t)_{t \in [0,1]}$ given \mathbf{x} , can be described by an initial distribution $q(\mathbf{z}_0|\mathbf{x})$ and a Stochastic Differential Equation (SDE),

$$d\{\mathbf{z}_t^{(r)}, \mathbf{z}_t^{(h)}\} = f(t)[\mathbf{z}_t^{(r)}, \mathbf{z}_t^{(h)}] dt + g(t) d\{\mathbf{w}^{(r)}, \mathbf{w}^{(h)}\},$$

where the drift $f(t)$ and volatility $g(t)$ are scalar functions of time, and $\mathbf{w}^{(r)}$ and $\mathbf{w}^{(h)}$ are two independent standard Wiener processes defined in $\mathbb{R}^{(M-1) \times 3}$ and $\mathbb{R}^{M \times D}$ respectively. Specifically, EDM implements the Variance-Preserving SDE (VP-SDE) scheme (Song et al., 2020), with $f(t) = -\frac{1}{2}\beta(t)$ and $g(t) = \sqrt{\beta(t)}$ for a fixed schedule $\beta(t)$. Due to the linearity of the drift term, the conditional marginal distribution can be reconstructed as

$$q([\mathbf{z}_t^{(r)}, \mathbf{z}_t^{(h)}] | [\mathbf{r}, \mathbf{h}]) = q(\mathbf{z}_t^{(r)} | \mathbf{r}) q(\mathbf{z}_t^{(h)} | \mathbf{h}), \\ = \mathcal{N}(\mathbf{z}_t^{(r)}, |\alpha_t \mathbf{r}, \sigma_t^2 \mathbb{I}) \cdot \mathcal{N}(\mathbf{z}_t^{(h)}, |\alpha_t \mathbf{h}, \sigma_t^2 \mathbb{I}),$$

where $\alpha_t = \exp(-\frac{1}{2} \int_0^t \beta(s) ds)$ and $\sigma_t = 1 - \exp(-\frac{1}{2} \int_0^t \beta(s) ds)$. The conditional distribution evolves from a low-variance Gaussian centered around the data $q(\mathbf{z}_0|\mathbf{x}) \approx \mathcal{N}(\mathbf{z}_0|\mathbf{x}, \delta \mathbb{I})$ to an uninformative prior distribution (that contains no information about the data distribution), i.e. a unit Gaussian $q(\mathbf{z}_1|\mathbf{x}) \approx \mathcal{N}(\mathbf{z}_1|\mathbf{0}, \mathbb{I})$.

Reverse (generative) process Starting from the prior $[\mathbf{z}_1^{(r)}, \mathbf{z}_1^{(h)}] \sim \mathcal{N}(\mathbf{z}_1^{(r)}|\mathbf{0}, \mathbb{I}) \cdot \mathcal{N}(\mathbf{z}_1^{(h)}|\mathbf{0}, \mathbb{I})$, samples from $q(\mathbf{x})$ can be generated by reversing the forward process. This can be done by following the reverse-time SDE (Anderson, 1982),

$$d\{\mathbf{z}_t^{(r)}, \mathbf{z}_t^{(h)}\} = f^B(\mathbf{z}_t, t) dt + g(t) d\{\bar{\mathbf{w}}^{(r)}, \bar{\mathbf{w}}^{(h)}\} \\ = \left[f(t)\mathbf{z}_t - g^2(t)\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \right] dt + g(t) d\{\bar{\mathbf{w}}^{(r)}, \bar{\mathbf{w}}^{(h)}\},$$

where $\mathbf{z}_t = [\mathbf{z}_t^{(r)}, \mathbf{z}_t^{(h)}]$, $\bar{\mathbf{w}}^{(r)}$ and $\bar{\mathbf{w}}^{(h)}$ are independent standard Wiener processes defined in $\mathbb{R}^{(M-1) \times 3}$ and $\mathbb{R}^{M \times D}$, respectively, with time flowing backwards. DM approximate the reverse process by learning an approximation of the score function $\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)$ parameterized by a neural network $s_\theta(\mathbf{z}_t, t)$. With the learned score function $s_\theta(\mathbf{z}_t, t)$, a sample $\mathbf{z}_0 \sim p_\theta(\mathbf{z}_0) \approx q(\mathbf{z}_0) \approx q(\mathbf{x})$ can

be obtained by first sampling from the prior $[z_1^{(r)}, z_1^{(h)}] \sim \mathcal{N}(z_t^{(r)} | \mathbf{0}, \mathbb{I}) \cdot \mathcal{N}(z_t^{(h)} | \mathbf{0}, \mathbb{I})$, and then simulating the reverse SDE,

$$dz_t = \left[f(t)z_t - g^2(t)s_\theta(z_t, t) \right] dt + g(t) d\{\bar{w}^{(r)}, \bar{w}^{(h)}\},$$

where the true score function has been replaced by its approximation $s_\theta(z_t, t)$. In EDM, the approximate score is parameterized through an equivariant function: $s_\theta(z_t, t) = [s_\theta^{(r)}(z_t, t), s_\theta^{(h)}(z_t, t)]$ such that $s_\theta([\mathbf{R}z_t^{(r)}, z_t^{(h)}], t) = [\mathbf{R}s_\theta^{(r)}(z_t, t), s_\theta^{(h)}(z_t, t)]$, $\forall \mathbf{R} \in O(3)$. This is realised through the specific parameterization,

$$s_\theta(z_t, t) = \frac{\alpha_t \hat{x}_\theta(z_t, t) - z_t}{\sigma_t^2},$$

where the data point predictor \hat{x}_θ is implemented by an equivariant neural network.

Optimization The data point predictor \hat{x}_θ , or s_θ , is trained by optimizing the denoising score matching loss

$$\mathcal{L}_{\text{DSM}} = \mathbb{E}_{\substack{u(t) \\ q(\mathbf{x}, z_t)}} \left[\lambda(t) \left\| s_\theta(z_t, t) - \nabla_{z_t} \log q(z_t | \mathbf{x}) \right\|_2^2 \right],$$

where $\lambda(t)$ is a positive weighting function, and $u(t)$ is a uniform distribution over the interval $[0, 1]$.

2.3. Neural Flow Diffusion Models

Neural Flow Diffusion Models (NFDM) (Bartosh et al., 2024) are based on the observation that latent variables, z_t , in conventional DM are inferred through a pre-specified transformation. This potentially limits the flexibility of the latent space, and make the learning of the reverse (generative) process more challenging.

Forward process In contrast to conventional DMs, NFDM defines the forward process implicitly through a learnable transformation $F_\varphi(\varepsilon, t, \mathbf{x})$ of injected noise ε , time t , and data point \mathbf{x} . The latent variables z_t are obtained by transforming noise samples ε , conditional on data point \mathbf{x} and time step t : $z_t = F_\varphi(\varepsilon, t, \mathbf{x})$. If F_φ is differentiable with respect to ε and t , and invertible with respect to ε , then, for fixed \mathbf{x} and ε , samples from $q_\varphi(z_t | \mathbf{x})$ can be obtained by solving the following conditional Ordinary Differential Equation (ODE) until time t ,

$$dz_t = f_\varphi(z_t, t, \mathbf{x}) dt, \quad (1)$$

where $f_\varphi(z_t, t, \mathbf{x}) = \left. \frac{\partial F_\varphi(\varepsilon, t, \mathbf{x})}{\partial t} \right|_{\varepsilon = F_\varphi^{-1}(z_t, t, \mathbf{x})}$, and with $z_0 \sim q(z_0 | \mathbf{x})$. While F_φ and $q(\varepsilon)$ define the conditional marginal distribution $q_\varphi(z_t | \mathbf{x})$, we need a distribution over the trajectories $(z_t)_{t \in [0, 1]}$. NFDM obtains this through the introduction of a conditional SDE starting from z_0 and

running forward in time. Given access to the ODE in Equation (1) and the score function $\nabla_{z_t} \log q_\varphi(z_t | \mathbf{x})$, a conditional SDE with conditional marginal distribution $q_\varphi(z_t | \mathbf{x})$ is given by

$$dz_t = f_\varphi^F(z_t, t, \mathbf{x}) dt + g_\varphi(t) d\mathbf{w}, \quad (2)$$

where $f_\varphi^F(z_t, t, \mathbf{x}) = f_\varphi(z_t, t, \mathbf{x}) + \frac{g_\varphi^2(t)}{2} \nabla_{z_t} \log q_\varphi(z_t | \mathbf{x})$. The score function of $q_\varphi(z_t | \mathbf{x})$ is

$$\nabla_{z_t} \log q_\varphi(z_t | \mathbf{x}) = \nabla_{z_t} [\log q(\varepsilon) + \log |J_F^{-1}|], \quad (3)$$

with $\varepsilon = F^{-1}(z_t, t, \mathbf{x})$, and $J_F^{-1} = \frac{\partial F^{-1}(z_t, t, \mathbf{x})}{\partial z_t}$.

Reverse (generative) process A conditional reverse SDE that starts from $z_1 \sim q(z_1)$, runs backward in time, and reverses the conditional forward SDE from Equation (2) can be defined as

$$\begin{aligned} dz_t &= f_\varphi^B(z_t, t, \mathbf{x}) + g_\varphi(t) d\bar{\mathbf{w}} \\ &= \left[f_\varphi(z_t, t, \mathbf{x}) - \frac{g_\varphi^2(t)}{2} \nabla_{z_t} \log q_\varphi(z_t | \mathbf{x}) \right] dt + g_\varphi(t) d\bar{\mathbf{w}}. \end{aligned}$$

As we do not have access to \mathbf{x} when generating samples, we can rewrite Equation (2), with the prediction of \mathbf{x} ,

$$dz_t = \hat{f}_{\theta, \varphi}(z_t, t) dt + g_\varphi(t) d\bar{\mathbf{w}}, \quad (4)$$

where $\hat{f}_{\theta, \varphi}(z_t, t) = f_\varphi^B(z_t, t, \hat{x}_\theta(z_t, t))$, and \hat{x}_θ is a function that predicts the data point \mathbf{x} . Provided that the reconstruction distribution $q(z_0 | \mathbf{x})$ and the prior distribution $q(z_1)$ are defined, this fully specifies the reverse process.

Optimization The forward and reverse processes can be optimized jointly by matching the drift terms of the true and approximate conditional reverse SDEs,

$$\mathcal{L}_{\text{NFDM}} = \mathbb{E}_{\substack{u(t) \\ q_\varphi(\mathbf{x}, z_t)}} \left[\frac{1}{2g_\varphi^2(t)} \left\| f_\varphi^B(z_t, t, \mathbf{x}) - \hat{f}_{\theta, \varphi}(z_t, t) \right\|_2^2 \right]. \quad (5)$$

3. Equivariant Neural Diffusion

Equivariant Neural Diffusion (END) generalizes the Equivariant Diffusion Model (EDM) (Hoogeboom et al., 2022), by defining the forward process through a learnable transformation. Our approach is a synthesis of NFDM introduced in Section 2.3, and leverages ideas of EDM outlined in Section 2.2 to maintain the desired invariance of the learned marginal distribution $p_{\theta, \varphi}(z_0)$. By providing an equivariant learnable transformation F_φ and an equivariant data point predictor \hat{x}_θ , we show that it is possible to obtain a generative model with the desired properties. Finally, we propose a simple yet flexible parameterization meeting the requirements.

3.1. Formulation

The key innovation in END lies in its forward process, which is also leveraged in the reverse (generative) process. The forward process is defined through a learnable time- and data-dependent transformation $F_\varphi(\varepsilon, t, \mathbf{x})$, such that the latent \mathbf{z}_t transforms covariantly with the injected noise ε (i.e. a collection of random vectors) and the data point \mathbf{x} ,

$$F_\varphi(\mathbf{R}\varepsilon, t, \mathbf{R}\mathbf{x}) = \mathbf{R}F_\varphi(\varepsilon, t, \mathbf{x}) = \mathbf{R}\mathbf{z}_t.$$

We then define $\hat{\mathbf{x}}_\theta$ as another learnable equivariant function, implying that the predicted data point transforms covariantly with the latent \mathbf{z}_t , i.e. $\hat{\mathbf{x}}_\theta(\mathbf{R}\mathbf{z}_t, t) = \mathbf{R}\hat{\mathbf{x}}_\theta(\mathbf{z}_t, t)$. Finally, we choose the noise distribution $p(\varepsilon)$, and the prior distribution $p(\mathbf{z}_1)$ to be invariant to the considered symmetry group.

Invariance of the learned distribution With the following choices: (1) $p(\mathbf{z}_1)$ an invariant distribution, (2) F_φ an equivariant function that satisfies $F_\varphi(\mathbf{R}\varepsilon, t, \mathbf{R}\mathbf{x}) = \mathbf{R}F_\varphi(\varepsilon, t, \mathbf{x})$, and (3) $\hat{\mathbf{x}}_\theta$ an equivariant function, we have that the learned marginal $p_{\theta, \varphi}(\mathbf{z}_0)$ is invariant as desired. This can be shown by demonstrating that the reverse SDE is equivariant. We start by noting that the reverse SDE in END is given by

$$d\mathbf{z}_t = \hat{f}_{\theta, \varphi}(\mathbf{z}_t, t) dt + g_\varphi(t) d\bar{\mathbf{w}}.$$

As the Wiener process is isotropic, this boils down to showing that the drift term, $\hat{f}_{\theta, \varphi}(\mathbf{z}_t, t)$ is equivariant, i.e. $\hat{f}_{\theta, \varphi}(\mathbf{R}\mathbf{z}_t, t) = \mathbf{R}\hat{f}_{\theta, \varphi}(\mathbf{z}_t, t)$. As the drift is expressed as a sum of two terms, we inspect each of them separately. The first term is

$$f_\varphi(\mathbf{z}_t, t, \hat{\mathbf{x}}_\theta) = \left. \frac{\partial F_\varphi(\varepsilon, t, \hat{\mathbf{x}}_\theta)}{\partial t} \right|_{\varepsilon = F_\varphi^{-1}(\mathbf{z}_t, t, \hat{\mathbf{x}}_\theta)},$$

where $\hat{\mathbf{x}}_\theta = \hat{\mathbf{x}}_\theta(\mathbf{z}_t, t)$.

If F_φ is equivariant, then so is its time-derivative (see Appendix A.1). The same holds for its inverse with respect to ε (see Appendix A.2), such that we have $F_\varphi^{-1}(\mathbf{R}\mathbf{z}_t, t, \mathbf{R}\mathbf{x}) = \mathbf{R}F_\varphi^{-1}(\mathbf{z}_t, t, \mathbf{x}) = \mathbf{R}\varepsilon$. We additionally have that $\hat{\mathbf{x}}_\theta$ is equivariant by definition.

As the equivariance of F_φ implies the equivariance of q_φ , when looking at the second term of the drift, we can see that, for $\mathbf{y}_t = \mathbf{R}\mathbf{z}_t$, we have

$$\nabla_{\mathbf{y}_t} \log q_\varphi(\mathbf{y}_t | \hat{\mathbf{x}}_\theta(\mathbf{y}_t, t)) = \mathbf{R}\nabla_{\mathbf{z}_t} \log q_\varphi(\mathbf{z}_t | \hat{\mathbf{x}}_\theta(\mathbf{z}_t, t)).$$

In summary, in addition to an invariant prior, an equivariant F_φ and an equivariant $\hat{\mathbf{x}}_\theta$ ensure the equivariance of the reverse process, and hence the invariance of the learned distribution. In Appendix A.3, we additionally show that the objective function in Equation (5) is invariant, i.e. $\mathcal{L}_{\text{END}}(\mathbf{R}\mathbf{x}) = \mathcal{L}_{\text{END}}(\mathbf{x}), \forall \mathbf{R} \in O(3)$.

3.2. Parameterization

We now introduce a simple parameterization of F_φ that meets the requirements outlined above:

$$F_\varphi(\varepsilon, t, \mathbf{x}) = \mu_\varphi(\mathbf{x}, t) + U_\varphi(\mathbf{x}, t)\varepsilon, \quad (6)$$

where, due to the geometric nature of \mathbf{x} , $U_\varphi(\mathbf{x}, t) \in \mathbb{R}^{(M-1) \times 3 \times 3}$ is structured as a block-diagonal matrix where each block is a 3×3 matrix. This is the equivalent to a diagonal parameterization in the case of scalar features. Similarly to EDM, our parametrization of F_φ leads to a conditional marginal $q_\varphi(\mathbf{z}_t | \mathbf{x})$ that is a conditional Gaussian with (block-) diagonal covariance, with the notable difference that the mean and covariance are now learnable through F_φ ,

$$q_\varphi(\mathbf{z}_t | \mathbf{x}) = \mathcal{N}(\mathbf{z}_t | \mu_\varphi(\mathbf{x}, t), \Sigma_\varphi(\mathbf{x}, t)), \quad (7)$$

where $\Sigma_\varphi(\mathbf{x}, t) = U_\varphi(\mathbf{x}, t)U_\varphi^\top(\mathbf{x}, t)$ such that $\Sigma_\varphi(\mathbf{x}, t)$ is also block-diagonal.

As F_φ is linear in ε , both μ_φ and U_φ must be equivariant functions whose outputs transform covariantly with \mathbf{x} , to ensure the equivariance of $F_\varphi(\varepsilon, t, \mathbf{x})$,

$$\begin{aligned} F_\varphi(\mathbf{R}\varepsilon, t, \mathbf{R}\mathbf{x}) &= \mu_\varphi(\mathbf{R}\mathbf{x}, t) + U_\varphi(\mathbf{R}\mathbf{x}, t)\mathbf{R}\varepsilon, \\ &= \mathbf{R}\mu_\varphi(\mathbf{x}, t) + \mathbf{R}U_\varphi(\mathbf{x}, t)\mathbf{R}\varepsilon \\ &= \mathbf{R}F_\varphi(\varepsilon, t, \mathbf{x}), \end{aligned}$$

as desired. We can then readily check that q_φ is equivariant, as $\forall \mathbf{R} \in O(3)$, we have that

$$\begin{aligned} q_\varphi(\mathbf{z}_t | \mathbf{x}) &= \mathcal{N}(\mathbf{z}_t | \mu_\varphi(\mathbf{x}, t), \Sigma_\varphi(\mathbf{x}, t)) \\ &= \mathcal{N}(\mathbf{R}\mathbf{z}_t | \mathbf{R}\mu_\varphi(\mathbf{x}, t), \mathbf{R}\Sigma_\varphi(\mathbf{x}, t)\mathbf{R}^\top) \\ &= q_\varphi(\mathbf{R}\mathbf{z}_t | \mathbf{R}\mathbf{x}). \end{aligned}$$

Prior and Reconstruction While not strictly required, it can be advantageous to parameterize F_φ such that the prior and reconstruction losses do not need to be computed. To do so, we need to make sure that $F_\varphi(\varepsilon, \mathbf{x}, t)$ is such that (i) $q(\mathbf{z}_0 | \mathbf{x}) \approx \mathcal{N}(\mathbf{z}_0 | \mathbf{x}, \delta^2 \mathbb{I})$, and (ii) $q(\mathbf{z}_1 | \mathbf{x}) \approx \mathcal{N}(\mathbf{0}, \mathbb{I})$. We parameterize the mean function as

$$\mu_\varphi(\mathbf{x}, t) = (1-t)\mathbf{x} + t(1-t)\bar{\mu}_\varphi(\mathbf{x}, t), \quad (8)$$

which ensures that $\mu_\varphi(\mathbf{x}, 0) = \mathbf{x}$, and $\mu_\varphi(\mathbf{x}, 1) = \mathbf{0}$; whereas for $U_\varphi(\mathbf{x}, t)$, we use the following

$$\begin{aligned} U_\varphi(\mathbf{x}, t) &= (\delta^{1-t}\bar{\sigma}_\varphi(\mathbf{x}, t)^{t(1-t)})\mathbb{I} \\ &\quad + t(1-t)\bar{U}_\varphi(\mathbf{x}, t), \quad (9) \end{aligned}$$

which ensures that $\Sigma_\varphi(\mathbf{x}, 0) = \delta^2 \mathbb{I}$, and $\Sigma_\varphi(\mathbf{x}, 1) = \mathbb{I}$, while being unconstrained for $t \in]0, 1[$. We give additional details about F_φ in Appendix A.4.

Implementation In practice, F_φ is implemented as a neural network with an architecture similar to that of the data point predictor $\hat{x}_\theta(z_t, t)$, but with a specific readout layer that produces $[\bar{\mu}_\varphi(\mathbf{x}, t), \bar{\sigma}_\varphi, \bar{U}_\varphi(\mathbf{x}, t)]$. The mean output $\bar{\mu}_\varphi(\mathbf{x}, t)$ is similar to that of $\hat{x}_\theta(z_t, t)$. For $U_\varphi(\mathbf{x}, t)$, we need to ensure that $\Sigma_\varphi(\mathbf{x}, t) = U_\varphi(\mathbf{x}, t)U_\varphi^\top(\mathbf{x}, t)$ rotates properly: $\bar{\sigma}_\varphi(\mathbf{x}, t)$ is a positive invariant scalar, while $\bar{U}_\varphi(\mathbf{x}, t)$ is constructed as matrix whose columns are vectors that transform covariantly with \mathbf{x} .

For ease of notation, we introduced all notations in the linear subspace R , however in practice we work in the ambient space, i.e. $\mathbf{r} \in \mathbb{R}^{M \times 3}$. We detail in Appendix A.4.1, how working in ambient space is possible.

The training and sampling procedures are detailed in Algorithms 1 and 2 in the appendix.

3.3. Conditional Generation

While unconditional generation is a required stepping stone, many real-life applications require some form of controllability. As other generative models, DM can model conditional distributions $p(\mathbf{x}|c)$, where c is a given condition. While different methods exist for sampling from the conditional distribution, e.g. (Wu et al., 2024), the simplest approach consists in training a conditional model on pairs (\mathbf{x}, c) . In that setting, F_φ and \hat{x}_θ are simply provided with an extra input c representing the conditional information, such that they respectively become $F_\varphi(\varepsilon, t, \mathbf{x}, c)$, and $\hat{x}_\theta(z_t, t, c)$. It is important to note that, compared to conventional DM, the forward process of END is now condition-dependent.

4. Experiments

In this section, we demonstrate the benefits of END with a comprehensive set of experiments. In Section 4.1, we first display the advantages of END for unconditional generation on 2 standard benchmarks, namely QM9 (Ramakrishnan et al., 2014) and GEOM-DRUGS (Axelrod & Gomez-Bombarelli, 2022), including an ablation study on QM9. Then, in Section 4.2, we perform conditional generation in 2 distinct settings on QM9.

4.1. Unconditional Generation

Datasets The QM9 dataset (Ramakrishnan et al., 2014) contains 134 thousand small- and medium-sized organic molecules with up to 9 heavy atoms, and up to 29 when counting hydrogen atoms. GEOM-DRUGS (Axelrod & Gomez-Bombarelli, 2022) contains 430 thousand medium- and large-sized drug-like molecules with 44 atoms on average, and up to maximum 181 atoms. We follow the same setup as previous work (Hoogeboom et al., 2022).

Task and Evaluation For each method, we sample 10 000 molecules using the stochastic sampling procedure detailed in Algorithm 2, with the number of integration steps varying from 50 to 1000. We repeat each sampling for 3 seeds, and report averages along with standard deviations for each metric.

On QM9, we follow previous work (Hoogeboom et al., 2022; Xu et al., 2023), and first evaluate the chemical quality of the generated samples in terms of stability, validity, and uniqueness. We additionally evaluate the ability of the model to learn the atom and bond types distributions by measuring the total variation between the dataset’s and generated distributions. Finally, we evaluate the geometry of the generated molecules, by computing the MMD (Gretton et al., 2012) between the generated and true bond distances distributions for the 10 most common bonds. We provide additional details about the evaluation procedure in Appendix A.6.1.

On GEOM-DRUGS, in addition to atom stability and validity as commonly reported in previous work (Hoogeboom et al., 2022; Xu et al., 2023), we also compute connectivity and total variation for atom types. Connectivity accounts for the fact that validity can easily be increased by generating several disconnected fragments (where only the largest counts towards validity), while the total variation ensures that the model properly samples all atom types.

Baselines We compare END to several relevant baselines from the literature: the original EDM (Hoogeboom et al., 2022), EDM-BRIDGE (Wu et al., 2022) an improved version of EDM that adds a physics-inspired force guidance in the reverse process, and GEOLDM (Xu et al., 2023) an equivariant latent DM. For a fair comparison, we additionally implement our own EDM (Hoogeboom et al., 2022), and denote it EDM*. It features the exact same architecture as END, as well as the same amount of learnable parameters. In EDM*, \hat{x}_θ is made of 10 layers, while END comprises a 5-layer \hat{x}_θ and a 5-layer F_φ . We provide additional details in Appendix A.6.

Results on QM9 Our results on the QM9 dataset are summarized in Table 1, and a few illustrative samples are shown in Figure 1. In addition to comparing with baselines from the literature, we conduct an ablation study, where EDM* + γ is similar to EDM*, but with a learned SNR (Kingma et al., 2021) for each data modality (i.e. atomic types and coordinates), and END (μ_φ only) where only the mean of the conditional marginal is learned. Details about the compared models are provided in Table 5.

We observe that the addition of a learnable forward allows for improved generative modeling, as the two variants of END are shown to perform better than (or be on par with) all baselines across all metrics. Most notably, with as few

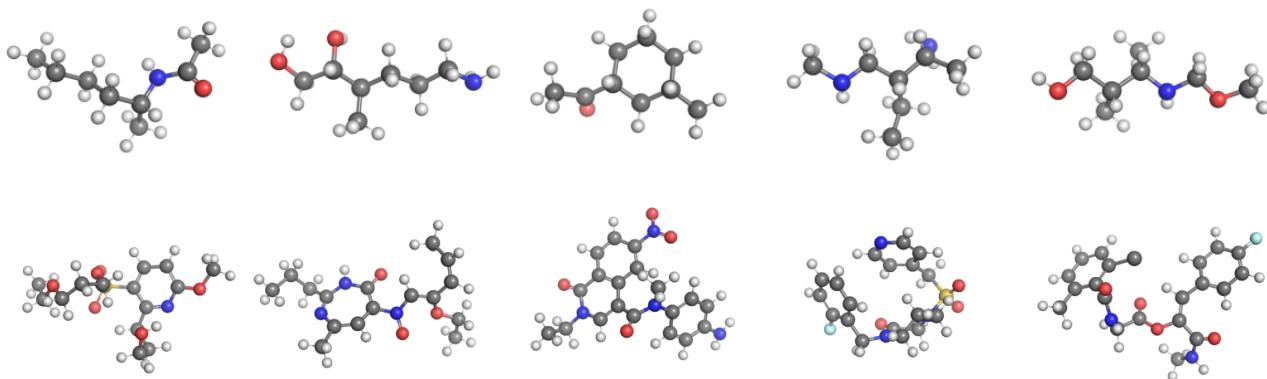


Figure 1: Representative samples generated by END on QM9 (top row), and GEOM-DRUGS (bottom row).

Table 1: Results on QM9. Metrics are obtained over 10000 samples, with mean/standard deviation across 3 sampling runs. For baselines from the literature, results are extracted from the respective papers. The two variants of END compare favorably to baselines across metrics, while offering competitive performance for reduced number of sampling steps.

Model	Steps	Stability (\uparrow)		Validity / Uniqueness (\uparrow)		Total Variation (\downarrow)		MMD (\downarrow)
		A [%]	M [%]	V [%]	V \times U [%]	A [10^{-2}]	B [10^{-3}]	[10^{-1}]
Data		99.0	95.2	97.7	97.7			
EDM (Hoogeboom et al., 2022)	1000	98.7	82.0	91.9	90.7			
EDM-BRIDGE (Wu et al., 2022)	1000	98.8	84.6	92.0	90.7			
GEOLDM (Xu et al., 2023)	1000	98.9 \pm .1	89.4 \pm .5	93.8 \pm .4	92.7 \pm .5			
EDM*	50	97.6 \pm .0	77.6 \pm .5	90.2 \pm .2	89.2 \pm .2	4.6 \pm .1	1.7 \pm .5	1.91 \pm .03
	100	98.1 \pm .0	81.9 \pm .4	92.1 \pm .2	90.9 \pm .2	3.5 \pm .1	1.4 \pm .3	1.67 \pm .02
	250	98.3 \pm .0	84.3 \pm .1	93.2 \pm .4	91.7 \pm .3	2.8 \pm .2	1.3 \pm .4	1.52 \pm .02
	500	98.4 \pm .0	85.2 \pm .5	93.5 \pm .2	92.2 \pm .3	2.6 \pm .2	1.3 \pm .4	1.50 \pm .04
	1000	98.4 \pm .0	85.3 \pm .3	93.5 \pm .1	91.9 \pm .1	2.5 \pm .1	1.4 \pm .4	1.51 \pm .02
EDM* + γ_φ	50	97.7 \pm .0	77.4 \pm .3	91.1 \pm .4	90.2 \pm .4	4.3 \pm .1	1.5 \pm .2	2.04 \pm .02
	100	98.2 \pm .0	82.6 \pm .2	92.9 \pm .2	91.6 \pm .2	3.2 \pm .1	1.2 \pm .2	1.66 \pm .01
	250	98.5 \pm .0	85.3 \pm .3	93.9 \pm .1	92.4 \pm .1	2.5 \pm .1	1.0 \pm .1	1.54 \pm .02
	500	98.5 \pm .1	86.1 \pm .4	94.1 \pm .2	92.5 \pm .2	2.2 \pm .1	1.0 \pm .3	1.50 \pm .02
	1000	98.5 \pm .0	86.1 \pm .3	94.1 \pm .2	92.4 \pm .2	2.1 \pm .1	1.1 \pm .1	1.45 \pm .03
END (μ_φ only)	50	98.5 \pm .0	83.9 \pm .2	95.2 \pm .2	93.8 \pm .3	1.4 \pm .1	1.9 \pm .4	2.80 \pm .06
	100	98.7 \pm .0	87.0 \pm .3	95.5 \pm .2	93.6 \pm .2	1.1 \pm .0	1.5 \pm .2	1.97 \pm .05
	250	98.9 \pm .0	89.0 \pm .2	95.8 \pm .2	93.8 \pm .2	1.0 \pm .0	0.5 \pm .1	1.48 \pm .03
	500	98.9 \pm .0	88.6 \pm .2	95.6 \pm .1	93.5 \pm .1	0.9 \pm .0	0.7 \pm .1	1.36 \pm .03
	1000	98.9 \pm .0	89.2 \pm .3	95.6 \pm .1	93.5 \pm .1	0.9 \pm .2	1.0 \pm .1	1.36 \pm .02
END	50	98.6 \pm .0	84.6 \pm .1	92.7 \pm .1	91.4 \pm .1	1.5 \pm .1	1.9 \pm .4	1.91 \pm .00
	100	98.8 \pm .0	87.4 \pm .2	94.1 \pm .0	92.3 \pm .2	1.3 \pm .0	1.8 \pm .3	1.63 \pm .02
	250	98.9 \pm .1	88.8 \pm .5	94.7 \pm .2	92.6 \pm .1	1.2 \pm .1	0.8 \pm .2	1.44 \pm .04
	500	98.9 \pm .0	88.8 \pm .4	94.8 \pm .2	92.8 \pm .2	1.2 \pm .1	0.8 \pm .5	1.41 \pm .01
	1000	98.9 \pm .0	89.1 \pm .1	94.8 \pm .1	92.6 \pm .2	1.2 \pm .1	0.8 \pm .5	1.37 \pm .04

Table 2: Results on GEOM-DRUGS. Metrics are obtained over 10000 samples, with mean/standard deviation across 3 sampling runs. For baselines from the literature, results are extracted from the respective papers. Most notably, END generates more connected samples.

Model	Steps	Stability (\uparrow)		Val. / Conn. (\uparrow)		TV (\downarrow)
		A [%]	V [%]	V \times C [%]	A [10^{-2}]	
Data		86.5	99.0			
EDM (Hoogeboom et al., 2022)	1000	81.3	92.6	—	—	—
EDM-BRIDGE (Wu et al., 2022)	1000	82.4	—	—	—	—
GEOLDM (Xu et al., 2023)	1000	84.4	99.3	—	—	—
EDM*	50	84.7 \pm .0	93.6 \pm .2	46.6 \pm .3	10.5 \pm .1	
	100	85.2 \pm .1	93.8 \pm .3	56.2 \pm .4	8.0 \pm .1	
	250	85.4 \pm .0	94.2 \pm .1	61.4 \pm .6	6.7 \pm .1	
	500	85.4 \pm .0	94.3 \pm .2	63.4 \pm .1	6.4 \pm .1	
	1000	85.3 \pm .1	94.4 \pm .1	64.2 \pm .6	6.2 \pm .0	
END	50	87.8 \pm .0	89.8 \pm .2	68.2 \pm .9	5.7 \pm .1	
	100	87.6 \pm .1	91.5 \pm .1	76.0 \pm .3	4.6 \pm .2	
	250	87.2 \pm .0	92.4 \pm .4	80.0 \pm .3	3.5 \pm .3	
	500	87.1 \pm .0	92.8 \pm .3	81.1 \pm .5	3.3 \pm .2	
	1000	87.0 \pm .0	92.9 \pm .3	82.2 \pm .2	3.0 \pm .3	

as 100 integration steps, END is able to generate samples that are qualitatively better than those generated by most baselines in 1000 steps.

Results on GEOM-DRUGS Our results on the GEOM-DRUGS dataset are presented in Table 2, and a few samples are displayed in Figure 1. We observe that END outperforms all baselines across all metrics but validity. This is due to the fact that the validity metric, as computed by (Hoogeboom et al., 2022; Xu et al., 2023), is obtained by only considering the largest fragment that can be extracted from a given sample. A method generating smaller disconnected fragments, can easily exploit and artificially increase this metric. When discarding disconnected samples, we observe that END does much better than the baseline, with an increase of around 20% on average.

4.2. Controllable Generation

Dataset and Setup We perform our experiments on the QM9 dataset, on 2 different tasks: composition-conditioned generation, and substructure-conditioned generation. Both tasks allow for direct validation with ground-truth properties without requiring expensive quantum mechanics calculations, or approximations with surrogate models. In each case, we train a conditional diffusion model as described in Section 3.3, i.e. where F_φ and \hat{x}_θ are provided with an extra input corresponding to the condition. Additional details are provided in Appendix A.6.4.

Task 1: composition-conditioned generation The model is tasked to generate a compound with a predefined composition, i.e. structural isomers of a given formula. The

Table 3: Results on composition-conditioned generation, where cEND offers nearly perfect composition controllability. Matching refers to the proportion of samples featuring the prompted composition.

Model	Steps	Matching [%] (\uparrow)
cEDM*	50	69.6 \pm .6
	100	73.0 \pm .6
	250	74.1 \pm 1.4
	500	76.2 \pm .6
	1000	75.5 \pm .5
cEND	50	89.2 \pm 0.8
	100	90.1 \pm 1.0
	250	91.2 \pm 0.8
	500	91.5 \pm 0.8
	1000	91.0 \pm 0.9

condition is specified as a vector $\mathbf{c} = (c_1, \dots, c_D) \in \mathbb{Z}^D$, where c_d denotes the number of atoms of type d that the sample should contain. To evaluate the model, we generate 10 samples per target formula, and compute the proportion of samples that match the provided composition. Our results are provided in Table 3, where we observe that conditional END significantly outperforms the baseline, and offers nearly perfectly controllable composition generation. Additionally, we can also see that reducing the number of sampling steps has a very limited impact on the controllability.

Task 2: substructure-conditioned generation We follow the same setup as (Bao et al., 2023) and train a con-

Table 4: Results on substructure-conditional sampling. For baselines from the literature, results are borrowed from the corresponding paper (Bao et al., 2023). cEND shows competitive performance, even surpassing EEGSDE that leverages an additional property predictor.

Model	Steps	Tanimoto Sim. (\uparrow)
CEDM (Bao et al., 2023)	1000	0.671 \pm .004
EEGSDE (Bao et al., 2023)	1000	0.750 \pm .003
cEDM*	50	0.601 \pm .000
	100	0.640 \pm .002
	250	0.663 \pm .002
	500	0.669 \pm .001
	1000	0.673 \pm .002
cEND	50	0.783 \pm .001
	100	0.807 \pm .001
	250	0.819 \pm .001
	500	0.825 \pm .001
	1000	0.828 \pm .001

ditional END, where the condition is a molecular fingerprint encoding structural information about the molecule. A fingerprint is a binary vector $c = (c_1, \dots, c_F) \in \{0, 1\}^F$, where c_f is set to 1 if substructure f is present in the molecule, or to 0 if not. Fingerprints are obtained using OPENBABEL (O’Boyle et al., 2011). To evaluate the ability of the compared models to leverage the provided structural information, we evaluate them by conditioning on unseen fingerprints (taken from test set) at sampling time. We then compute the similarity between the fingerprints computed on the generated samples and the fingerprints provided as inputs. We compare cEND to EEGSDE (Bao et al., 2023), an improved version of EDM (Hooeboom et al., 2022), that performs conditional generation by combining a conditional diffusion model and regressor guidance. Our results are presented in Table 4, along with a handful of samples in Figure 2. cEND offers better controllability than the compared baselines, as highlighted by the higher similarity.

5. Related Work

The main approaches to molecule generation in 3D are autoregressive models (Gebauer et al., 2019; Simm et al., 2020; Gebauer et al., 2022; Luo & Ji, 2022; Daigavane et al., 2024), flow-based models (Garcia Satorras et al., 2021), and diffusion models (Hooeboom et al., 2022; Igashov et al., 2024). A notable exception to the geometric graph representation of 3D molecules are voxels (Skalic et al., 2019; Ragoza et al., 2022; O Pinheiro et al., 2024), from which the 3D graph is extracted using some post-processing procedure. The closest work to ours is GEOLDM (Xu et al., 2023), a geometric latent diffusion model that performs

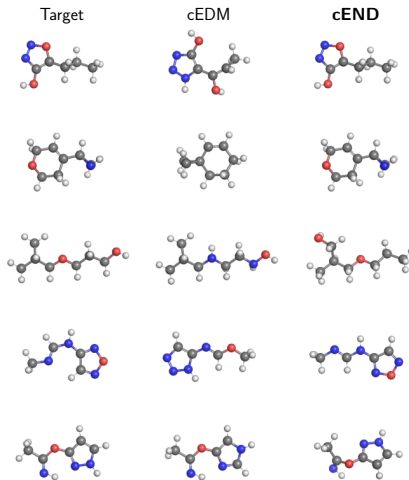


Figure 2: Excerpt of substructure-conditioned samples., where cEND can be seen to match the provided substructure better (in terms of compositions and local patterns).

diffusion in the latent space of an equivariant Variational Auto-Encoder (VAE). It corresponds to a particular instance of END with $F_\varphi(\varepsilon, t, \mathbf{x}) = \alpha_t E_\varphi(\mathbf{x}) + \varepsilon \sigma_t$, where $E_\varphi(\mathbf{x})$ corresponds to the encoder of the VAE. Recently, several works have shown that leveraging 2D connectivity information can lead to improved results (Peng et al., 2023; Vignac et al., 2023; Le et al., 2024). While not incompatible with END, we perform our experiments without modeling that auxiliary information, and therefore do not compare to these approaches directly. Other generative frameworks have also been tailored to molecule generation, such as Flow Matching (Lipman et al., 2022; Song et al., 2023) or Bayesian Flow Networks (Graves et al., 2024; Song et al., 2024).

6. Conclusion

In this work, we have presented Equivariant Neural Diffusion (END), a novel diffusion model that is equivariant to Euclidean transformations. The key innovation in END lies in the forward process that is specified by a learnable data- and time-dependent transformation. Experimental results demonstrate the benefits of our method. In the unconditional setting, we show that END yields competitive generative performance across two different benchmarks. In the conditional setting, END offers better controllability, when conditioning on composition and substructure. Avenues for future work are numerous. In particular, leveraging the flexible framework of NFDM (Bartosh et al., 2024) to constrain the generative trajectories, e.g. to be straight and enable faster sampling, modelling bond information, or extending the conditional setting to other types of conditioning information, e.g. other point cloud or target property, are all promising directions.

References

- Anderson, B. D. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- Anstine, D. M. and Isayev, O. Generative models as an emerging paradigm in the chemical sciences. *Journal of the American Chemical Society*, 145(16):8736–8750, 2023.
- Axelrod, S. and Gomez-Bombarelli, R. Geom, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data*, 9(1):185, 2022.
- Bao, F., Zhao, M., Hao, Z., Li, P., Li, C., and Zhu, J. Equivariant energy-guided SDE for inverse molecular design. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=r0otLtOwYW>.
- Bartosh, G., Vetrov, D., and Naesseth, C. A. Neural diffusion models. *arXiv preprint arXiv:2310.08337*, 2023.
- Bartosh, G., Vetrov, D., and Naesseth, C. A. Neural flow diffusion models: Learnable forward process for improved diffusion modelling. *arXiv preprint arXiv:2404.12940*, 2024.
- Batatia, I., Kovacs, D. P., Simm, G., Ortner, C., and Csányi, G. Mace: Higher order equivariant message passing neural networks for fast and accurate force fields. *Advances in Neural Information Processing Systems*, 35: 11423–11436, 2022.
- Batzner, S., Musaelian, A., Sun, L., Geiger, M., Mailoa, J. P., Kornbluth, M., Molinari, N., Smidt, T. E., and Kozinsky, B. E (3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nature communications*, 13(1):2453, 2022.
- Daigavane, A., Kim, S. E., Geiger, M., and Smidt, T. Symphony: Symmetry-equivariant point-centered spherical harmonics for 3d molecule generation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=MIEnYtlGyv>.
- Garcia Satorras, V., Hoogeboom, E., Fuchs, F., Posner, I., and Welling, M. E (n) equivariant normalizing flows. *Advances in Neural Information Processing Systems*, 34: 4181–4192, 2021.
- Gebauer, N., Gastegger, M., and Schütt, K. Symmetry-adapted generation of 3d point sets for the targeted discovery of molecules. *Advances in neural information processing systems*, 32, 2019.
- Gebauer, N. W., Gastegger, M., Hessmann, S. S., Müller, K.-R., and Schütt, K. T. Inverse design of 3d molecular structures with conditional generative neural networks. *Nature communications*, 13(1):973, 2022.
- Graves, A., Srivastava, R. K., Atkinson, T., and Gomez, F. Bayesian flow networks, 2024.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Hoogeboom, E., Satorras, V. G., Vignac, C., and Welling, M. Equivariant diffusion for molecule generation in 3d. In *International conference on machine learning*, pp. 8867–8887. PMLR, 2022.
- Igashov, I., Stärk, H., Vignac, C., Schneuing, A., Satorras, V. G., Frossard, P., Welling, M., Bronstein, M., and Correia, B. Equivariant 3d-conditional diffusion model for molecular linker design. *Nature Machine Intelligence*, pp. 1–11, 2024.
- Kingma, D., Salimans, T., Poole, B., and Ho, J. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- Landrum, G. et al. Rdkit: A software suite for cheminformatics, computational chemistry, and predictive modeling. *Greg Landrum*, 8:31, 2013.
- Le, T., Noe, F., and Clevert, D.-A. Representation learning on biomolecular structures using equivariant graph attention. In *The First Learning on Graphs Conference*, 2022. URL <https://openreview.net/forum?id=kv4xUo5Pu6>.
- Le, T., Cremer, J., Noe, F., Clevert, D.-A., and Schütt, K. T. Navigating the design space of equivariant diffusion-based generative models for de novo 3d molecule generation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=kzGuiRXZrQ>.
- Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2022.
- Luo, Y. and Ji, S. An autoregressive flow model for 3d molecular geometry generation from scratch. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=C03Ajc-NS5W>.

- Merchant, A., Batzner, S., Schoenholz, S. S., Aykol, M., Cheon, G., and Cubuk, E. D. Scaling deep learning for materials discovery. *Nature*, 624(7990):80–85, 2023.
- Nielsen, B. M. G., Christensen, A., Dittadi, A., and Winther, O. Diffenc: Variational diffusion with a learned encoder. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=8nxy1bQWTG>.
- O Pinheiro, P. O., Rackers, J., Kleinhenz, J., Maser, M., Mahmood, O., Watkins, A., Ra, S., Sresht, V., and Saremi, S. 3d molecule generation by denoising voxel grids. *Advances in Neural Information Processing Systems*, 36, 2024.
- O’Boyle, N. M., Banck, M., James, C. A., Morley, C., Vandermeersch, T., and Hutchison, G. R. Open babel: An open chemical toolbox. *Journal of cheminformatics*, 3: 1–14, 2011.
- Peng, X., Guan, J., Liu, Q., and Ma, J. Moldiff: Addressing the atom-bond inconsistency problem in 3d molecule diffusion generation. In *International Conference on Machine Learning*, pp. 27611–27629. PMLR, 2023.
- Ragoza, M., Masuda, T., and Koes, D. R. Generating 3d molecules conditional on receptor binding sites with deep generative models. *Chemical science*, 13(9):2701–2713, 2022.
- Ramakrishnan, R., Dral, P. O., Rupp, M., and Von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- Ruddigkeit, L., Van Deursen, R., Blum, L. C., and Reymond, J.-L. Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17. *Journal of chemical information and modeling*, 52(11):2864–2875, 2012.
- Schütt, K., Kindermans, P.-J., Saucedo Felix, H. E., Chmiela, S., Tkatchenko, A., and Müller, K.-R. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *Advances in neural information processing systems*, 30, 2017.
- Schütt, K., Unke, O., and Gastegger, M. Equivariant message passing for the prediction of tensorial properties and molecular spectra. In *International Conference on Machine Learning*, pp. 9377–9388. PMLR, 2021.
- Simm, G., Pinsler, R., and Hernández-Lobato, J. M. Reinforcement learning for molecular design guided by quantum mechanics. In *International Conference on Machine Learning*, pp. 8959–8969. PMLR, 2020.
- Skalic, M., Jiménez, J., Sabbadin, D., and De Fabritiis, G. Shape-based generative modeling for de novo drug design. *Journal of chemical information and modeling*, 59(3):1205–1214, 2019.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.
- Song, Y., Gong, J., Xu, M., Cao, Z., Lan, Y., Ermon, S., Zhou, H., and Ma, W.-Y. Equivariant flow matching with hybrid probability transport for 3d molecule generation. In Oh, A., Neumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 549–568. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/01d64478381c33e29ed611f1719f5a37-Paper-Conference.pdf.
- Song, Y., Gong, J., Zhou, H., Zheng, M., Liu, J., and Ma, W.-Y. Unified generative modeling of 3d molecules with bayesian flow networks. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=NSVtmmzeRB>.
- Unke, O. T., Chmiela, S., Saucedo, H. E., Gastegger, M., Poltavsky, I., Schuett, K. T., Tkatchenko, A., and Mueller, K.-R. Machine learning force fields. *Chemical Reviews*, 121(16):10142–10186, 2021.
- Vignac, C., Osman, N., Toni, L., and Frossard, P. Midi: Mixed graph and 3d denoising diffusion for molecule generation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 560–576. Springer, 2023.
- Weininger, D. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.
- Wu, L., Gong, C., Liu, X., Ye, M., and qiang liu. Diffusion-based molecule generation with informative prior bridges. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho,

K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=TJUNtiZiTKKE>.

Wu, L., Trippe, B., Naesseth, C., Blei, D., and Cunningham, J. P. Practical and asymptotically exact conditional sampling in diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.

Xu, M., Yu, L., Song, Y., Shi, C., Ermon, S., and Tang, J. Geodiff: A geometric diffusion model for molecular conformation generation. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=PzcvxEMzvQC>.

Xu, M., Powers, A. S., Dror, R. O., Ermon, S., and Leskovec, J. Geometric latent diffusion models for 3d molecule generation. In *International Conference on Machine Learning*, pp. 38592–38610. PMLR, 2023.

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

A. Appendix / supplemental material

A.1. Time-derivative of an $O(3)$ -equivariant function

Let $f : \mathcal{X} \times [0, 1] \rightarrow \mathcal{Y}$ be a function that is equivariant to actions of the group $O(3)$, such that $\mathbf{R} \cdot f(\mathbf{x}, t) = f(\mathbf{R} \cdot \mathbf{x}, t), \forall \mathbf{R} \in O(3)$.

Proof sketch We need to show that $\frac{\partial}{\partial t}(f(\mathbf{R} \cdot \mathbf{x}, t)) = \mathbf{R} \cdot \frac{\partial}{\partial t}(f(\mathbf{x}, t)), \forall \mathbf{R} \in O(3)$ and $\forall \mathbf{x} \in \mathcal{X}$.

$$\frac{\partial}{\partial t}(f(\mathbf{R} \cdot \mathbf{x}, t)) = \frac{\partial}{\partial t}(\mathbf{R} \cdot f(\mathbf{x}, t)), \quad (10)$$

$$= \mathbf{R} \cdot \frac{\partial}{\partial t}(f(\mathbf{x}, t)) \quad (11)$$

where the last equality follows by linearity.

A.2. Inverse of an $O(3)$ -equivariant function

Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a function that (1) is equivariant to the action of the group $O(3)$, and (2) admits an inverse $f^{-1} : \mathcal{Y} \rightarrow \mathcal{X}$, then f^{-1} is also equivariant to the action of $O(3)$.

Proof sketch We need to show that $f^{-1}(\mathbf{R} \cdot \mathbf{y}) = \mathbf{R} \cdot f^{-1}(\mathbf{y}), \forall \mathbf{R} \in O(3)$ and $\forall \mathbf{y} \in \mathcal{Y}$.

Since f is invertible, we have that to any $\mathbf{y} \in \mathcal{Y}$ corresponds a unique $\mathbf{x} \in \mathcal{X}$, such that $\mathbf{y} = f(\mathbf{x})$, and that $f^{-1}(\mathbf{y}) = f^{-1}(f(\mathbf{x})) = \mathbf{x}$. As f is equivariant to the action of $O(3)$, we have that $\mathbf{R} \cdot f(\mathbf{x}) = f(\mathbf{R} \cdot \mathbf{x}), \forall \mathbf{R} \in O(3)$:

$$f^{-1}(\mathbf{R} \cdot \mathbf{y}) = f^{-1}(\mathbf{R} \cdot f(\mathbf{x})), \quad (12)$$

$$= f^{-1}(f(\mathbf{R} \cdot \mathbf{x})), \quad (13)$$

$$= \mathbf{R} \cdot \mathbf{x}, \quad (14)$$

$$= \mathbf{R} \cdot f^{-1}(\mathbf{y}). \quad (15)$$

A.3. O(3)-invariance of the objective function

In this section, we show that the objective function is invariant under the action of $O(3)$: $\mathcal{L}_{\text{END}}(\mathbf{R}\mathbf{x}) = \mathcal{L}_{\text{END}}(\mathbf{x})$, $\forall \mathbf{R} \in O(3)$, provided that F_φ and $\hat{\mathbf{x}}_\theta$.

$$\begin{aligned}
\mathcal{L}_{\text{END}}(\mathbf{R}\mathbf{x}) &= \mathbb{E}_{u(t), q_\varphi(\mathbf{z}_t|\mathbf{R}\mathbf{x})q(\mathbf{R}\mathbf{x})} \left[\frac{1}{2g_\varphi^2(t)} \left\| f_\varphi(\mathbf{z}_t, t, \mathbf{R}\mathbf{x}) - \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{z}_t} \log q_\varphi(\mathbf{z}_t|\mathbf{R}\mathbf{x}) \right. \right. \\
&\quad \left. \left. - f_\varphi(\mathbf{z}_t, t, \hat{\mathbf{x}}_\theta(\mathbf{z}_t, t)) + \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{z}_t} \log q_\varphi(\mathbf{z}_t|\hat{\mathbf{x}}_\theta(\mathbf{z}_t, t)) \right\|_2^2 \right], \\
&= \mathbb{E}_{u(t), q_\varphi(\mathbf{z}_t|\mathbf{R}\mathbf{x})q(\mathbf{R}\mathbf{x})} \left[\frac{1}{2g_\varphi^2(t)} \left\| f_\varphi(\mathbf{R}\mathbf{R}^{-1}\mathbf{z}_t, t, \mathbf{R}\mathbf{x}) - \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{z}_t} \log q_\varphi(\mathbf{R}\mathbf{R}^{-1}\mathbf{z}_t|\mathbf{R}\mathbf{x}) \right. \right. \\
&\quad \left. \left. - f_\varphi(\mathbf{R}\mathbf{R}^{-1}\mathbf{z}_t, t, \hat{\mathbf{x}}_\theta(\mathbf{R}\mathbf{R}^{-1}\mathbf{z}_t, t)) + \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{z}_t} \log q_\varphi(\mathbf{R}\mathbf{R}^{-1}\mathbf{z}_t|\hat{\mathbf{x}}_\theta(\mathbf{R}\mathbf{R}^{-1}\mathbf{z}_t, t)) \right\|_2^2 \right], \\
&= \mathbb{E}_{u(t), q_\varphi(\mathbf{z}_t|\mathbf{R}\mathbf{x})q(\mathbf{R}\mathbf{x})} \left[\frac{1}{2g_\varphi^2(t)} \left\| \mathbf{R}f_\varphi(\mathbf{R}^{-1}\mathbf{z}_t, t, \mathbf{x}) - \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{z}_t} \log q_\varphi(\mathbf{R}^{-1}\mathbf{z}_t|\mathbf{x}) \right. \right. \\
&\quad \left. \left. - \mathbf{R}f_\varphi(\mathbf{R}^{-1}\mathbf{z}_t, t, \hat{\mathbf{x}}_\theta(\mathbf{R}^{-1}\mathbf{z}_t, t)) + \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{z}_t} \log q_\varphi(\mathbf{R}^{-1}\mathbf{z}_t|\hat{\mathbf{x}}_\theta(\mathbf{R}^{-1}\mathbf{z}_t, t)) \right\|_2^2 \right], \\
&= \mathbb{E}_{u(t), q_\varphi(\mathbf{R}\mathbf{y}_t|\mathbf{R}\mathbf{x})q(\mathbf{R}\mathbf{x})} \left[\frac{1}{2g_\varphi^2(t)} \left\| \mathbf{R}f_\varphi(\mathbf{y}_t, t, \mathbf{x}) - \frac{g_\varphi^2(t)}{2} \mathbf{R} \nabla_{\mathbf{y}_t} \log q_\varphi(\mathbf{y}_t|\mathbf{x}) \right. \right. \\
&\quad \left. \left. - \mathbf{R}f_\varphi(\mathbf{y}_t, t, \hat{\mathbf{x}}_\theta(\mathbf{y}_t, t)) + \frac{g_\varphi^2(t)}{2} \mathbf{R} \nabla_{\mathbf{y}_t} \log q_\varphi(\mathbf{y}_t|\hat{\mathbf{x}}_\theta(\mathbf{y}_t, t)) \right\|_2^2 \right], \\
&= \mathbb{E}_{u(t), q_\varphi(\mathbf{y}_t|\mathbf{x})q(\mathbf{x})} \left[\frac{1}{2g_\varphi^2(t)} \left\| f_\varphi(\mathbf{y}_t, t, \mathbf{x}) - \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{y}_t} \log q_\varphi(\mathbf{y}_t|\mathbf{x}) \right. \right. \\
&\quad \left. \left. - f_\varphi(\mathbf{y}_t, t, \hat{\mathbf{x}}_\theta(\mathbf{y}_t, t)) + \frac{g_\varphi^2(t)}{2} \nabla_{\mathbf{y}_t} \log q_\varphi(\mathbf{y}_t|\hat{\mathbf{x}}_\theta(\mathbf{y}_t, t)) \right\|_2^2 \right], \\
&= \mathcal{L}_{\text{END}}(\mathbf{x})
\end{aligned}$$

The first equality is obtained by replacing \mathbf{x} by $\mathbf{R}\mathbf{x}$ in the definition of the objective function Equation (5). The second is obtained by multiplying by $\mathbf{R}\mathbf{R}^{-1} = \mathbb{I}$. The third equality by leveraging that f_φ , q_φ and $\hat{\mathbf{x}}_\theta$ are equivariant. We then perform a change of variable $\mathbf{y}_t = \mathbf{R}^{-1}\mathbf{z}_t$. As rotation does preserve distances, we obtain the last equality.

A.4. Details about F_φ

A.4.1. WORKING IN AMBIENT SPACE

In practice, we would like F_φ and $\hat{\mathbf{x}}_\theta$ to operate in ambient space, i.e. on $\mathbf{z}_t^{(r)} \in \mathbb{R}^{M \times 3}$.

Garcia Satorras et al. (2021) showed that the Jacobian of the transformation can be directly computed in ambient space for $\mathbf{z}_t^{(r)}$ and $\boldsymbol{\varepsilon}^{(r)}$ that live in the linear subspace R , provided that the transformation F_φ leaves the center of mass unchanged. Additionally, F_φ is required to be invertible with respect to $\boldsymbol{\varepsilon}$.

If we consider a flat representation of $\mathbf{z}_t^{(r)} \in \mathbb{R}^{M \cdot d}$, such transformation can be written as,

$$U_\varphi(\mathbf{x}, t) = (\mathbb{I}_{M \cdot d} - \frac{1}{M} T T^\top) \tilde{U}_\varphi(\mathbf{x}, t) + \frac{1}{M} T T^\top, \quad (16)$$

where $T \in \mathbb{R}^{M \cdot d \times d} = [\mathbb{I}_d, \mathbb{I}_d, \dots]^\top$, and $\frac{1}{M} T T^\top$ corresponds to the linear operator computing the center of mass. Intuitively, the first term corresponds to a linear transformation followed by a projection to the 0-CoM subspace, while the second term translates the system back to the initial center of mass.

Computing $\mathbf{z}_t^{(r)}$ from $\boldsymbol{\varepsilon}^{(r)}$ Given that $\boldsymbol{\varepsilon}^{(r)} \in R$, obtaining $\mathbf{z}_t^{(r)} \in R$ from $\boldsymbol{\varepsilon}^{(r)}$ simply amounts to (1) computing $\tilde{\mathbf{z}}_t^{(r)} = \tilde{U}_\varphi(\mathbf{x}, t) \boldsymbol{\varepsilon}^{(r)}$, and then (2) removing the center of mass from $\tilde{\mathbf{z}}_t^{(r)}$.

Computing $|J_F|$ and F^{-1} In what follows, we shorten the notation, and denote $U_\varphi(\mathbf{x}, t)$ by U and $\tilde{U}_\varphi(\mathbf{x}, t)$ by \tilde{U} . We

start by reorganizing Equation (16), as

$$U = \tilde{U} + \frac{1}{M}TT^\top(\mathbb{I}_{M \cdot d} - \tilde{U}). \quad (17)$$

Computing $\det U$ can be done by leveraging the Matrix Determinant Lemma,

$$\det U = \det \tilde{U} \cdot \det \left(\mathbb{I}_d + \frac{1}{M}T^\top(\mathbb{I}_{M \cdot d} - \tilde{U})\tilde{U}^{-1}T \right), \quad (18)$$

$$= \det \tilde{U} \cdot \det \left(\mathbb{I}_d + \frac{1}{M}T^\top(\tilde{U}^{-1} - \mathbb{I}_{M \cdot d})T \right), \quad (19)$$

$$= \det \tilde{U} \cdot \det \left(\frac{1}{M} \sum_{m=1}^M (\tilde{U}^m)^{-1} \right), \quad (20)$$

$$= \det \tilde{U} \cdot \det V, \quad (21)$$

$$= \prod_{m=1}^M \det \tilde{U}^m \cdot \det V \quad (22)$$

where $V = \frac{1}{M} \sum_{m=1}^M (\tilde{U}^m)^{-1}$ is a $d \times d$ -matrix, and \tilde{U}^m denotes the m -th $d \times d$ -block in \tilde{U} .

The inverse can be obtained by leveraging the Woodbury matrix identity,

$$U^{-1} = \left(\tilde{U} + \frac{1}{M}TT^\top(\mathbb{I}_{M \cdot d} - \tilde{U}) \right)^{-1} \quad (23)$$

$$= \tilde{U}^{-1} - \frac{1}{M}\tilde{U}^{-1}TV^{-1}T^\top(\tilde{U}^{-1} - \mathbb{I}_{M \cdot d}), \quad (24)$$

$$= \tilde{U}^{-1}(\mathbb{I} - C) \quad (25)$$

where $V = \frac{1}{M} \sum_{m=1}^M (\tilde{U}^m)^{-1}$, as previously defined, and $C = \frac{1}{M}TV^{-1}T^\top(\tilde{U}^{-1} - \mathbb{I}_{M \cdot d})$.

In practice, we do not need the inverse itself, but rather $\varepsilon^{(r)}$ given $\mathbf{z}_t^{(r)}$,

$$\varepsilon^{(r)} = U^{-1}\mathbf{z}_t^{(r)}, \quad (26)$$

$$= \tilde{U}^{-1}(\mathbb{I} - C)\mathbf{z}_t^{(r)}, \quad (27)$$

$$= \tilde{U}^{-1}[\mathbf{z}_t^{(r)} - \mathbf{c}], \quad (28)$$

where \mathbf{c} is a translation vector that can be obtained without constructing the full matrix C .

A.4.2. INVARIANT FEATURES

For simplicity, we omitted in Section 3.1 that molecules are described as tuples: $\mathbf{x} = (\mathbf{r}, \mathbf{h})$, as only \mathbf{r} transform under Euclidean transformations. For the invariant features \mathbf{h} , we use the following parameterization

$$\mu_\varphi^{(h)}(\mathbf{x}, t) = (1 - t)\mathbf{h} + t(1 - t)\bar{\mu}_\varphi^{(h)}(\mathbf{x}, t), \quad (29)$$

which ensures that $\mu_\varphi^{(h)}(\mathbf{x}, 0) = \mathbf{h}$, and $\mu_\varphi^{(h)}(\mathbf{x}, 1) = \mathbf{0}$; whereas for $\sigma_\varphi^{(h)}(\mathbf{x}, t)$, we use the following

$$\sigma_\varphi^{(h)}(\mathbf{x}, t) = \delta^{1-t}\bar{\sigma}_\varphi(\mathbf{x}, t)^{t(1-t)}. \quad (30)$$

Implementation F_φ is implemented as a neural network with an architecture similar to that of the data point predictor $\hat{\mathbf{x}}_\theta(\mathbf{z}_t, t)$, but with a specific readout layer that produces the outputs related to \mathbf{r} ($[\bar{\mu}_\varphi(\mathbf{x}, t), \bar{\sigma}_\varphi, \bar{U}_\varphi(\mathbf{x}, t)]$), as described in the main text in Section 3.1. Additionally, it produces $\bar{\mu}_\varphi^{(h)}(\mathbf{x}, t)$ and $\bar{\sigma}_\varphi^{(h)}(\mathbf{x}, t)$ as invariant outputs.

Inverse transformation The logarithm of the determinant of the inverse transformation $\log |J_F^{-1}|$ writes

$$\log |J_F^{-1}| = -\log |J_F| = - \underbrace{\sum_{i=1}^{M \times D} \sigma_\varphi^{(h),i}(\mathbf{x}, t)}_{\text{invariant features}} - \underbrace{\sum_{i=m}^M \log |\det(U_\varphi^m(\mathbf{x}, t))|}_{\text{vectorial features}} - \log \det |V|, \quad (31)$$

where V is defined as in Equation (22).

A.5. Algorithms

Algorithm 1 Training algorithm of END

Require: $q(\mathbf{x}), F_\varphi, \hat{\mathbf{x}}_\theta$
for training iterations **do**
 $\mathbf{x} \sim q(\mathbf{x}), t \sim u(t), \varepsilon \sim p(\varepsilon)$
 $\mathbf{z}_t \leftarrow \mu_\varphi(\mathbf{x}, t) + U_\varphi(\mathbf{x}, t)\varepsilon$
 $\mathcal{L} = \frac{1}{2g_\varphi^2(t)} \|f_\varphi^B(\mathbf{z}_t, t, \mathbf{x}) - \hat{f}_{\theta, \varphi}(\mathbf{z}_t, t)\|_2^2$
 Gradient step on θ and φ
end for

Algorithm 2 Stochastic sampling from END

Require: $F_\varphi, \hat{\mathbf{x}}_\theta$, integration steps T , empirical distribution of number of atoms $p(N)$
 $\Delta t = \frac{1}{T}$
 $N \sim p(N)$
 $\mathbf{z}_1 \sim p(\mathbf{z}_1)$
for $t = 1, \dots, \frac{1}{T}$ **do**
 $\bar{\mathbf{w}} \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$
 $\mathbf{z}_{t-\Delta t} \leftarrow \mathbf{z}_t - \hat{f}_{\theta, \varphi}(\mathbf{z}_t, t)\Delta t + g_\varphi(t)\bar{\mathbf{w}}\sqrt{\Delta t}$
end for
 $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z}_0)$

A.6. Experimental details

A.6.1. EVALUATION METRICS

Stability An atom is deemed stable if it has a charge of 0, whereas a molecule is stable if all its atoms have 0 charge. We reuse the lookup table from (Hoogeboom et al., 2022) to infer bond types from pairwise distances.

Validity Validity corresponds to the percentage of samples that can be parsed and sanitized by RDKit (Landrum et al., 2013), after inference of the bonds using a lookup table (Hoogeboom et al., 2022). It should be noted that the metric does not penalize fragmented samples as long as each individual fragment appears valid. This can be problematic when running evaluation on larger compounds, as models tend to generate disconnected structures.

Uniqueness It is the proportion of samples that are valid and have a unique SMILES string (Weininger, 1988).

Total variation The total variation is computed as the MAE between the (discrete) marginal obtained on the training data and on the generated samples. For bond types, we compute the ground truth and generated distributions using the lookup table mechanism.

MMD We follow the procedure of (Daigavane et al., 2024), and compute the MMD between true and generated pairwise distances distributions for the 10 most common bonds in QM9: ["C-H:1.0", "C-C:1.0", "C-O:1.0", "C-N:1.0", "H-N:1.0", "C-O:2.0", "C-N:1.5", "H-O:1.0", "C-C:1.5", "C-N:2.0"]. Here, we infer bonds using RDKit’s `rdDetermineBonds.DetermineBonds` with `charge=0`.

A.6.2. ARCHITECTURE

Our forward transformation F_φ and data point predictor $\hat{\mathbf{x}}_\theta$ share a common neural network architecture that we detail here. The architecture is similar to that of EQCAT (Le et al., 2022), and updates a collection of invariant and equivariant features for each node in the graph. We chose that architecture because it allows for an easy construction of $\bar{U}_\varphi(\mathbf{x}, t)$ by linear projection of the final equivariant layer.

We follow previous work (Hoogeboom et al., 2022) and consider fully-connected graphs. We initially featurize pairwise distances through Gaussian Radial Basis functions, with dataset-specific cutoff taken large enough to ensure full connectivity. In opposition to (Hoogeboom et al., 2022), we do not update positions in the message-passing phase, but instead obtain

the positions prediction through a linear projection of the final equivariant hidden states. The predictions for the invariant features are obtained by reading out the final invariant hidden states.

Optimization For all model variants, we employ Adam with a learning rate of 10^{-4} . We perform gradient clipping (norm) with a value of 10 on QM9, and a value of 1 on GEOM-DRUGS.

A.6.3. UNCONDITIONAL GENERATION

We reuse the data setup from previous work (Hoogeboom et al., 2022; Xu et al., 2023).

QM9 On QM9, we use 10 layers of message passing for EDM*, while the variants of END feature 5 layers of message-passing in F_φ and 5 layers in \hat{x}_θ . For all models, we use 256 invariant and 256 equivariant hidden features, along with an RBF expansion of dimension 64 with a cutoff of 10\AA for pairwise distances. This ensures that the compared models have the same number of learnable parameters, i.e. 9.4M each. We train all models for 1000 epochs with a batch size of 64.

GEOM-Drugs On GEOM-DRUGS, we use 10 layers of message passing for EDM*, while the variants of END feature 5 layers of message-passing in F_φ and 5 in \hat{x}_θ . The hidden size of the invariant and equivariant features is set to 192, along with an RBF expansion of dimension 64 with a cutoff of 30\AA for pairwise distances. Each model features 5.4M learnable parameters. We train all models for 10 epochs with an effective batch size of 64.

A.6.4. CONDITIONAL GENERATION

We use 10 layers of message passing for EDM*, while the variants of END feature 5 layers of message-passing in F_φ and 5 in \hat{x}_θ . The hidden size of the invariant and equivariant features is set to 192, along with an RBF expansion of dimension 64 with a cutoff of 10\AA for pairwise distances. We train all models for 1000 epochs with a batch size of 64.

After an initial encoding, the conditional information is introduced at the end of each message passing step, and alters the scalar hidden states through a one-layer MLP, that shares the same dimension as the hidden scalar state.

Composition-conditioned generation The encoding of the condition follows that of (Gebauer et al., 2022). Each atom type gets its own embedding, weighted by the proportion it represents in the provided formula. The weighted embeddings of all atom types are then concatenated and flattened, and the obtained vector is processed through a 2-layer MLP with 64 hidden units. The composition used at sampling time are extracted from the validation and test set. For each unique formula, the model gets to generate 10 samples.

Substructure-conditioned generation The encoding of the condition follows that of (Bao et al., 2023). The 1024-dimensional fingerprint is simply processed by a 2-layer MLP with hidden dimensions [512, 256], and a final linear projection to 192, i.e. the hidden size of the invariant features.

A.7. Compared models

In Table 5, we detail the compared all models in terms of their transformation F_φ .

Table 5: Compared models.

	$F_\varphi(\varepsilon, t, \mathbf{x})$	Comment
EDM (Hoogeboom et al., 2022) / EDM*	$\alpha_t \mathbf{x} + \sigma_t \varepsilon$	$\alpha_t = \exp\left(-\frac{1}{2} \int_0^t \beta(s) ds\right)$ $\sigma_t = 1 - \exp\left(-\frac{1}{2} \int_0^t \beta(s) ds\right)$
GEOLDM (Xu et al., 2023)	$\alpha_t E_\varphi(\mathbf{x}) + \sigma_t \varepsilon$	α_t and σ_t similar to EDM $p(\mathbf{x} \mathbf{z}_0) = \mathcal{N}(\mathbf{x} D_\varphi(\mathbf{z}_0), \delta^2 \mathbb{I})$
EDM* + γ_φ	$\alpha_\varphi(t) \mathbf{x} + \sigma_\varphi(t) \varepsilon$	learned γ_φ with 2 separate outputs (for \mathbf{r} and \mathbf{h})
END (μ_φ only)	$\mu_\varphi(\mathbf{x}, t) + \sigma_t \varepsilon$	σ_t similar to EDM
END	$\mu_\varphi(\mathbf{x}, t) + U_\varphi(\mathbf{x}, t) \varepsilon$	as introduced in Equation (6)

A.8. Compute resources

All experiments were run on a single GPU. The experiments on QM9 were run on a NVIDIA SM3090 with 24 GB of memory. The experiments on GEOM-DRUGS were run on NVIDIA A100 with 40 GB. Training took up to 7 days.