

---

# Rethinking Diffusion Models with Symmetries through Canonicalization with Applications to Molecular Graph Generation

---

Cai Zhou<sup>\*1</sup> Zijie Chen<sup>\*2</sup> Zian Li<sup>3</sup> Jike Wang<sup>2</sup> Pan Li<sup>4</sup> Kaiyi Jiang<sup>1,5</sup>  
Rose Yu<sup>6</sup> Muhan Zhang<sup>3</sup> Stephen Bates<sup>1</sup> Tommi Jaakkola<sup>1</sup>

## Abstract

Many generative tasks in chemistry and science involve distributions invariant to group symmetries (e.g., permutation and rotation). A common strategy enforces invariance and equivariance through architectural constraints such as equivariant denoisers and invariant priors. In this paper, we challenge this tradition through the alternative canonicalization perspective: first map each sample to an orbit representative with a canonical pose or order, train an unconstrained (non-equivariant) diffusion or flow model on the canonical slice, and finally recover the invariant distribution by sampling a random symmetry transform at generation time. Building on a formal quotient-space perspective, our work provides a comprehensive theory of canonical diffusion by proving: (i) the correctness, universality and superior expressivity of canonical generative models over invariant targets; (ii) canonicalization accelerates training by removing diffusion score complexity induced by group mixtures and reducing conditional variance in flow matching. We then show that aligned priors and optimal transport act complementarily with canonicalization and further improves training efficiency. We instantiate the framework for molecular graph generation under  $S_N \times SO(3)$  symmetries. By leveraging geometric spectral-based canonicalization and mild positional encodings, canonical diffusion significantly outperforms equivariant baselines in 3D molecule generation tasks, with similar or even less computation. Moreover, with a novel architecture *Canon*, CanonFlow achieves state-of-the-art performance on the challenging GEOM-DRUG dataset.

<sup>1</sup>Massachusetts Institute of Technology <sup>2</sup>Zhejiang University  
<sup>3</sup>Peking University <sup>4</sup>Georgia Institute of Technology <sup>5</sup>Princeton University <sup>6</sup>University of California, San Diego. Correspondence to: Cai Zhou <caiz428@mit.edu>.

Proceedings of the 43<sup>rd</sup> International Conference on Machine Learning, Seoul, South Korea. PMLR 306, 2026. Copyright 2026 by the author(s).

## 1. Introduction

Generative modeling is fundamentally grounded in the geometric structure of data. In domains such as computer vision and natural language processing (NLP), data exhibits specific symmetries—for instance, objects in images possess translation invariance, while semantic meaning in text relies on sequential order. Recent diffusion (Song & Ermon, 2019; Ho et al., 2020; Song et al., 2021) and flow-based (Liu et al., 2022; Lipman et al., 2022; Albergo et al., 2023) generative models have achieved great success across images (Dhariwal & Nichol, 2021; Rombach et al., 2022), videos (Ho et al., 2022; Singer et al., 2022), text (Li et al., 2022; Gong et al., 2022), and biomolecules (Watson et al., 2023; Corso et al., 2023). Crucially, however, modalities like images are not fully invariant to all transformations: an upside-down landscape or a reversed sentence typically loses its semantic validity. Consequently, state-of-the-art diffusion models in these fields explicitly break symmetries: they utilize fixed grid topologies or inject positional encodings (PEs) to anchor the generation process to a canonical orientation.

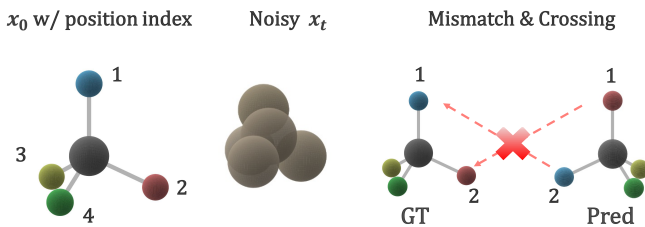


Figure 1. Motivation of CanonFlow: position-wise supervision induces mismatch loss under symmetries for equivalent molecules in diffusion models.

In contrast, molecular generation operates on a fundamentally different geometric space defined by the direct product of permutation and Euclidean symmetries,  $S_N \times SE(3)$ . Unlike natural images, a rotated molecule or a re-indexed molecular graph represents the exact same physical entity. Molecular generative models typically enforce the constraints by building equivariant architectures and/or invariant priors (Hoogeboom et al., 2022; Xu et al., 2022; Tian et al., 2024), so that the learned vector field (score/velocity) respects the group action and the generated distribution is

symmetry-consistent. While principled, this approach often incurs substantial architectural and computational overhead (e.g., equivariant layers, tensor algebra), and it can obscure an additional challenge: symmetry creates latent “gauge” ambiguity, so intermediate noisy states may correspond to multiple equivalent group-transformed configurations. This complex mixture-like nature results in “trajectory crossing” and conflicting gradients (Figure 1), as the model struggles to determine *which* valid orientation to generate from a symmetric noise vector.

Motivated by this observation, we argue that the efficiency of image generation stems precisely from its lack of total invariance, and we propose to transfer this advantage to molecular modeling via **canonicalization**. In this work, we challenge the necessity of equivariant generative models and introduce a novel canonical diffusion framework, with theoretical analysis of how this symmetry breaking procedure accelerates training while ensuring validity over invariant or equivariant targets. By mapping the quotient space of molecules to a unique canonical slice, we explicitly break the symmetry, aligning the noise and data distributions. This effectively addresses the trajectory mismatch problem, transforming the complex equivariant generation task into a simplified problem on the canonical gauge.

Our starting point is the flow-matching identity that the irreducible regression error equals an expected conditional variance of endpoint displacements given an intermediate state. In Section 3, under a group-aligned lifting construction, we show this conditional variance decomposes into two nonnegative components: a within-slice term that reflects genuine transport difficulty on a canonical slice, and a symmetry-ambiguity term that arises purely from marginalizing over the latent group element. Canonicalization eliminates the symmetry-ambiguity term by fixing a gauge (working on a canonical representative per orbit), whereas optimal-transport-like couplings target the within-slice term by making the coupling more Monge-like and the trajectory closer to straight-line transport.

Based on this analysis, in Section 4 we propose a practical symmetry-aware pipeline for molecular generation that combines a canonicalizer with flexible non-equivariant backbones. We canonicalize molecules by jointly fixing permutation and rotation gauges (e.g., a Fiedler-vector rank for  $S_N$  and a rank-anchored frame for  $SO(3)$ ), then train diffusion/flow models directly in canonical space with a gap-free prior and canonical conditions. This yields a simple yet powerful alternative to fully equivariant generative modeling: it reduces symmetry-induced ambiguity, improves few-step accuracy with the straighter learned transport, and unlocks the expressive and computational advantages of generic Transformers/GNNs for large-scale molecular generation. Furthermore, we develop a novel architecture termed

**Canon** (Figure 7) which explicitly incorporates and refines canonical information in an additional atom hidden state, enabling canonicity-aware denoising.

Experimentally, canonicalization is compatible with any (equivariant or non-equivariant) base models. Section 5 shows that the canonicalized counterparts of baselines such as SemlaFlow (Irwin et al., 2024) yield significantly better results on popular unconditional 3D molecule generation benchmarks, including QM9 and GEOM-DRUG, with negligible computation overheads. Furthermore, leveraging the “Canon” architecture, *CanonFlow* leads to state-of-the-art molecule stability and validity metrics, outperforming all baselines with large margins. Remarkably, canonicalized models reveal faster training convergence and better few-step sampling quality.

## 2. Preliminary

**Canonicalization in symmetric data.** We work in a measurable space  $\mathcal{M}$  of structured objects. In molecular generation, we use the following state throughout the paper:  $\mathbf{Z} = (\mathbf{X}, \mathbf{H}, \mathbf{A}) \in \mathcal{M}$ , where  $\mathbf{X} \in \mathbb{R}^{N \times 3}$  are coordinates;  $\mathbf{H} \in \mathbb{R}^{N \times d_h}$  are atom features, each element consisting of  $d_h$  scalars, e.g., atom types  $\{1, \dots, K\}^N$  and formal charges;  $\mathbf{A} \in \mathbb{R}^{N \times N \times d_e}$  encode bonds or edge features. Let a group action  $\mathcal{G}$  act on  $\mathcal{M}$  by

$$g \cdot \mathbf{Z} = (R, t, \pi) \cdot (\mathbf{X}, \mathbf{H}, \mathbf{A}) \\ := (\pi(\mathbf{X})R^\top + \mathbf{1}t^\top, \pi(\mathbf{H}), \pi(\mathbf{A})), \quad (1)$$

where  $(R, t) \in SE(3)$  and  $\pi \in S_N$  (simultaneous permutation of nodes and adjacency).

**Definition 2.1** (Orbit and quotient). The orbit of  $\mathbf{Z}$  is  $\mathcal{O}(\mathbf{Z}) = \{g \cdot \mathbf{Z} : g \in \mathcal{G}\}$ . The quotient space is  $\mathcal{M}/\mathcal{G} = \{\mathcal{O}(\mathbf{Z}) : \mathbf{Z} \in \mathcal{M}\}$ .

**Definition 2.2** (Invariance). A probability measure  $\mu$  on  $\mathcal{M}$  is  $\mathcal{G}$ -invariant if  $g\#\mu = \mu$  for all  $g \in \mathcal{G}$ .

We now introduce *canonicalization* as orbit selection.

**Definition 2.3** (Canonicalization map and slice). A measurable map  $\Psi : \mathcal{M} \rightarrow \mathcal{M}$  is a *canonicalization map* if (i)  $\Psi(\mathbf{Z}) \in \mathcal{O}(\mathbf{Z})$  and (ii)  $\Psi(g \cdot \mathbf{Z}) = \Psi(\mathbf{Z})$  for all  $g \in \mathcal{G}$ . The image  $S := \Psi(\mathcal{M})$  is a *canonical slice*.

**Remark 2.4** (Stabilizers and non-uniqueness). If  $\mathbf{Z}$  has non-trivial stabilizer  $\text{Stab}(\mathbf{Z}) = \{g : g \cdot \mathbf{Z} = \mathbf{Z}\}$  (e.g., graph automorphisms or symmetric geometries), then canonical representatives may be non-unique. This inherently causes discontinuities or multi-valued choices; weighted/probabilistic frames are one remedy (Dym et al., 2024).

Fortunately, we assume the following holds throughout the paper, which is reasonable for (noisy) real-world 3D molecules (exact symmetries are measure-zero).

**Assumption 2.5** (Free action a.s.). Under  $p_0$ , the stabilizer is trivial almost surely.

There exists a finite, translation-invariant measure (Haar measure  $\lambda$ ) that can be normalized to have total mass 1 on a compact topological group. For permutations  $S_N$ , Haar is uniform; for  $SO(3)$ , Haar is the uniform rotation measure; for translations, however, Haar is not finite. Thus following previous work (Hoogeboom et al., 2022; Li et al., 2025), throughout the paper we remove translations by centering (e.g., center-of-mass), leaving a compact symmetry  $SO(3) \times S_N$ .

**Assumption 2.6** (Centered representation). We assume  $\mathbf{X}$  is zero-centered so global translations are removed; the remaining group is compact:  $\mathcal{G} = SO(3) \times S_N$  with Haar probability  $\lambda$ .

**Diffusion and flow-matching models.** We present a continuous-time formulation, while discrete-time DDPM (Ho et al., 2020) and CTMC variants follow similarly. Let  $\mathbf{Z}_0 \sim p_0$  be data. We first present score-based diffusion (Song et al., 2021). Consider a forward noising SDE on an ambient Euclidean embedding of  $\mathcal{M}$  (continuous parts such as  $\mathbf{X}$ ), a standard VP-SDE is

$$d\mathbf{Z}_t = -\frac{1}{2}\beta(t)\mathbf{Z}_t dt + \sqrt{\beta(t)} dW_t, \quad (2)$$

with marginals  $p_t$  and score  $s_*(\mathbf{Z}, t) = \nabla_{\mathbf{Z}} \log p_t(\mathbf{Z})$  learned by a parameterized score network.

Some recent generative models are built from the equivalent flow matching viewpoint, learning a time-dependent vector field  $v_t(\mathbf{Z})$  that transports a prior  $p_1$  to data  $p_0$  via an ODE:

$$\frac{d\mathbf{Z}_t}{dt} = v_t(\mathbf{Z}_t), \quad \mathbf{Z}_1 \sim p_1, \quad \mathbf{Z}_0 \sim p_0. \quad (3)$$

In conditional flow matching, one specifies a coupling  $\gamma(\mathbf{Z}_0, \mathbf{Z}_1)$  and a *microscopic* conditional path  $\mathbf{Z}_t = \Phi_t(\mathbf{Z}_0, \mathbf{Z}_1)$  with conditional vector field  $u_t(\cdot | \mathbf{Z}_0, \mathbf{Z}_1)$ . The optimal marginal field is the conditional expectation:  $v_t(\mathbf{Z}) = \mathbb{E}_\gamma[u_t(\mathbf{Z} | \mathbf{Z}_0, \mathbf{Z}_1) | \mathbf{Z}_t = \mathbf{Z}]$ . The categorical parts (such as  $\mathbf{A}, \mathbf{H}$ ), in comparison, are usually modeled by an appropriate discrete diffusion, typically with uniform/absorbing/masked states as the prior.

### 3. Understanding Canonicalization in Diffusion Models with Symmetries

Building on a formal quotient-space perspective, we provide a comprehensive theory of canonical diffusion in this section. We prove (i) a factorization theorem for invariant measures via slice distributions and Haar randomization, and the universality and superior expressivity of canonicalized generative models over invariant targets (Section 3.1);

(ii) an explicit expression for the diffusion score complexity induced by group mixtures while being removable through canonicalization, and a flow-matching conditional variance decomposition showing that canonicalization reduces conditional variance, thereby accelerating diffusion/flow matching training (Section 3.2). Complete proof with mathematical notation and more discussions are available in Appendix C.

#### 3.1. Canonical Generative Models Induce Universal Invariant Distributions

This subsection formalizes why learning on the canonical slice while sampling with Haar randomization is not only sufficient in representing any invariant target (Section 3.1.1), but also more expressive with non-equivariant denoising backbones (Section 3.1.2).

##### 3.1.1. UNIVERSALITY OF CANONICAL PARAMETERIZATIONS OVER INVARIANT AND EQUIVARIANT TARGETS

First, we show that invariant distributions factor through the canonical slice. Recall that Haar probability  $\lambda$  is always uniform for  $S_N \times SO(3)$ .

**Theorem 3.1** (Factorization of invariant measures; known). *Suppose Assumptions 2.5 and 2.6 hold. Let  $\mu$  be any  $\mathcal{G}$ -invariant probability measure on  $\mathcal{M}$ . Let  $\Psi$  be an orbit representative map defined  $\mu$ -a.s., and let  $\nu = \Psi\#\mu$  be the slice distribution on  $S = \Psi(\mathcal{M})$ . Then*

$$\mu = \int_S \left( \int_{\mathcal{G}} \delta_{g \cdot \mathbf{z}} d\lambda(g) \right) d\nu(\mathbf{Z}). \quad (4)$$

*Equivalently, if  $\tilde{\mathbf{Z}} \sim \nu$ ,  $g \sim \lambda$  independent, then  $g \cdot \tilde{\mathbf{Z}} \sim \mu$ .*

**Corollary 3.2** (Sufficiency of slice modeling). *To model any invariant target  $\mu$ , it suffices to model the slice distribution  $\nu$ ; invariance is recovered by Haar randomization.*

We propose that canonicalization is a general technique for constructing equivariant/invariant functions *without equivariant backbones*, generalizing (Kaba et al., 2023):

**Proposition 3.3** (Universality of canonicalized parameterizations). *Let  $\mathcal{G}$  acts continuously and orthogonally on a compact set  $K \subset \mathbb{R}^d$ . Suppose we have a (measurable) canonicalization map  $\Psi : K \rightarrow K$  as Theorem 2.3, and a (measurable) gauge map  $\kappa : K \rightarrow \mathcal{G}$  s.t.*

$$\Psi(g \cdot x) = \Psi(x), \quad \kappa(g \cdot x) = g\kappa(x), \quad x = \kappa(x) \cdot \Psi(x) \quad (5)$$

*Consider the parametrization*

$$\phi(x) = \kappa(x) \cdot f(\Psi(x)), \quad (6)$$

*where  $f$  is a universal approximator on  $\Psi(K)$ . Then  $\phi$  is a universal approximator of continuous  $\mathcal{G}$ -equivariant functions on  $K$ , and  $f \circ \Psi$  is universal for continuous  $\mathcal{G}$ -invariant functions on  $K$ .*

### 3.1.2. STRONGER EXPRESSIVITY VIA SYMMETRY BREAKING

We now argue that symmetry breaking with the help of non-equivariant models can be practically more expressive. Even though the true score of an invariant distribution is equivariant, enforcing equivariance inside the denoiser can **restrict architectural choices**, *limiting expressivity while increasing unnecessary computation costs* (Yan et al., 2023). Remarkably, non-equivariant models can be more expressive than their equivariant counterparts (Zhang et al., 2021; Zhou et al., 2023a). However, frame averaging or relational pooling (Puny et al., 2021; Murphy et al., 2019) are needed to recover invariant outputs, leading to a complexity of the order of the group. This motivates canonicalized diffusion: it uses a symmetry-breaking gauge (canonical order/pose) and leverages a stronger non-equivariant backbone on the slice, then restores invariance. Without the necessity to traverse and average, canonicalized non-equivariant models yield stronger expressivity with similar or less computation. More details are available in Appendix C.1.2.

### 3.2. Canonicalization Accelerates Diffusion Training

In this subsection, we provide two complementary analyses, score mixture complexity (Section 3.2.1) and flow-matching conditional variance (Section 3.2.2), rigorously formalizing how canonicalization decreases symmetry-induced variance and accelerates diffusion model training. We further show that aligned canonical prior and optimal transport further reduces within-slice cost (Section 3.2.3).

#### 3.2.1. MIXTURE STRUCTURE IN DIFFUSION SCORE

Consider a finite group  $\mathcal{G}$  with  $M$  elements acting orthogonally on  $\mathbb{R}^d$  and  $g_m \in G$  the  $m$ -th element. Let  $q$  be a slice density supported on a canonical region  $S$ . Define the invariant mixture:

$$p(z_t) = \frac{1}{M} \sum_{m=1}^M q(g_m^{-1} \cdot z_t). \quad (7)$$

where  $z_t$  denotes noisy states, and the summation can be substituted by integration for infinite groups. When  $q$  is differentiable and positive, the score follows given the chain rule and orthogonality,

$$\begin{aligned} \nabla \log p(z_t) &= \sum_{m=1}^M w_m(z_t) g_m \cdot \nabla \log q(g_m^{-1} \cdot z_t), \\ w_m(z_t) &:= \frac{q(g_m^{-1} \cdot z_t)}{\sum_{j=1}^M q(g_j^{-1} \cdot z_t)}. \end{aligned} \quad (8)$$

Where multiple group copies overlap, the responsibilities  $w_m(z_t)$  vary rapidly. The induced sharp score fields require more capacity and smaller step sizes for stable reverse integration. As noted in (Yan et al., 2023), invariant training

leads to increased modes and mixture-like optimal denoising scores (GMM analogy). Instead, canonicalization removes the mixture: there is a single component on the slice ( $M = 1$ ), which smooths the score landscape, simplifying denoising score network training.

#### 3.2.2. CONDITIONAL FLOW VARIANCE REDUCTION VIA SYMMETRY AMBIGUITY ELIMINATION

Here we formalize our central result: *canonicalization reduces conditional variance in flow matching training* by eliminating symmetry ambiguity. Without loss of generality, consider the linear-path setting, the conditional vector field is constant:

$$\begin{aligned} \mathbf{Z}_t &= (1-t)\mathbf{Z}_0 + t\mathbf{Z}_1, \quad (\mathbf{Z}_0, \mathbf{Z}_1) \sim \gamma, \\ u_t(\mathbf{Z}_t | \mathbf{Z}_0, \mathbf{Z}_1) &= \frac{d}{dt} \mathbf{Z}_t = \mathbf{Z}_1 - \mathbf{Z}_0. \end{aligned} \quad (9)$$

The optimal marginal field is  $v_t(\mathbf{Z}) = \mathbb{E}[\mathbf{Z}_1 - \mathbf{Z}_0 | \mathbf{Z}_t = \mathbf{Z}]$ . Let  $\hat{v}$  be any measurable predictor of  $v_t(\mathbf{Z}_t)$  from  $(t, \mathbf{Z}_t)$ . Then the minimum achievable MSE (equivalently, the *irreducible flow-matching error*) is the conditional variance:

$$\inf_{\hat{v}} \mathbb{E} [\|\hat{v}(t, \mathbf{Z}_t) - (\mathbf{Z}_1 - \mathbf{Z}_0)\|^2] = \mathbb{E} [\text{Var}(\mathbf{Z}_1 - \mathbf{Z}_0 | t, \mathbf{Z}_t)]. \quad (10)$$

We aim to understand when  $\text{Var}(\mathbf{Z}_1 - \mathbf{Z}_0 | t, \mathbf{Z}_t)$  is large, especially in the presence of symmetry. Consider the following *symmetry-mixture data model*. Let a compact group  $\mathcal{G}$  act linearly and orthogonally on  $\mathbb{R}^d$  (so  $\|g \cdot x\| = \|x\|$ ). Assume the data distribution  $p_0$  is  $\mathcal{G}$ -invariant:

$$\begin{aligned} Z_0 &\stackrel{d}{=} G \cdot \tilde{Z}_0 \sim p_0, \quad G \sim \lambda(\text{Haar}), \\ \tilde{Z}_0 &\sim q_0, \quad \tilde{Z}_0 = \Psi(Z_0), \quad q_0 = \Psi_{\#} p_0, \end{aligned} \quad (11)$$

where  $\Psi$  is a (measurable) canonicalizer mapping each orbit to a representative on a slice. Intuitively, without canonicalization the posterior  $p(Z_0 | Z_t)$  can be multi-modal over the latent symmetry element  $G$ , inflating conditional variance.

Compare the following two training paradigms: (i) **canonical slice training** (quotient space), train on canonicalized data  $\tilde{Z}_0 := \Psi(Z_0) \sim q_0$  with a slice prior  $\tilde{Z}_1 \sim q_1$  and a coupling  $\tilde{\gamma}(\tilde{Z}_0, \tilde{Z}_1)$ ; (ii) **invariant-mixture training** (ambient space), train directly on the invariant data  $Z_0 \sim p_0$  with an ambient prior  $Z_1 \sim p_1$  and a coupling  $\gamma(Z_0, Z_1)$ . To relate slice coupling to invariant ambient-space coupling, we define the *group-aligned lift* of a given  $\tilde{\gamma}$ :

$$\begin{aligned} G &\sim \lambda, \quad (\tilde{Z}_0, \tilde{Z}_1) \sim \tilde{\gamma} \text{ independent}, \\ (Z_0, Z_1) &:= (G \cdot \tilde{Z}_0, G \cdot \tilde{Z}_1) \sim \gamma^{\text{lift}}. \end{aligned} \quad (12)$$

We emphasize that the coupling correspondence (namely the invariant training uses the induced group-aligned lifting as the ambient coupling) is practically reasonable. When  $\tilde{Z}_1 \sim$

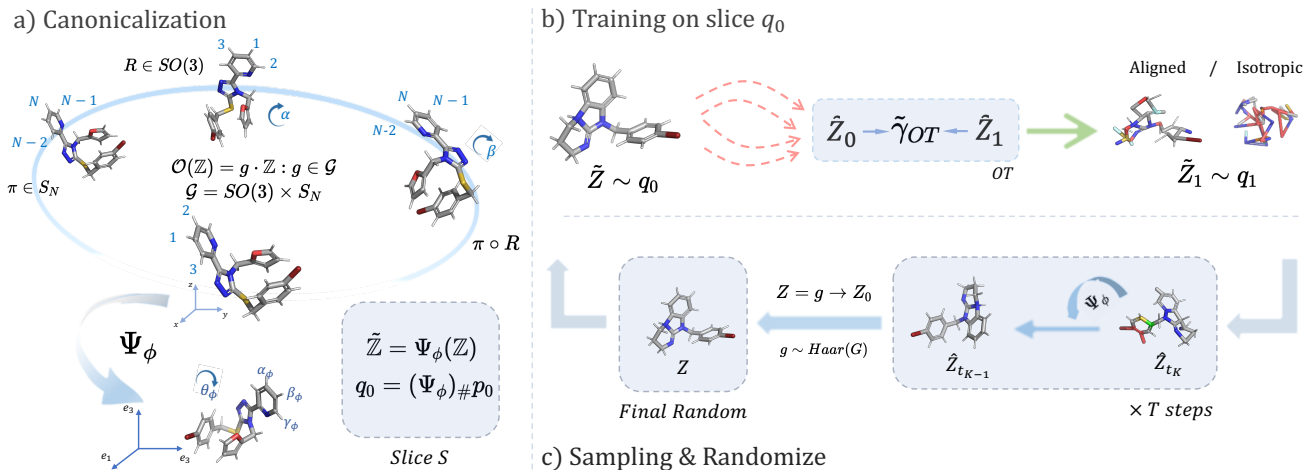


Figure 2. Overview of our canonicalized generation pipeline. (a) Canonicalization: map a molecule  $Z$  to a slice representative  $\tilde{Z} = \Psi_\phi(Z)$  under  $\mathcal{G} = SO(3) \times S_N$ , inducing  $q_0 = (\Psi_\phi)_\# p_0$ . (b) Training: learn a diffusion/flow model on  $\tilde{Z}_0 \sim q_0$  with slice prior  $q_1$  (optionally using an OT coupling for aligned pairings). (c) Sampling: generate on the slice (optionally with projected canonical sampling) and then apply  $g \sim \text{Haar}(\mathcal{G})$  to obtain an invariant sample  $Z = g \cdot \tilde{Z}_0$ .

$\mathcal{N}(0, I)$  and the slice coupling is the independent product  $\tilde{\gamma} := q_0 \otimes q_1$ , then the ambient prior  $Z_1 \sim p_1 = \mathcal{N}(0, I)$  and the invariant coupling recovers the widely adopted product coupling (w/o optimal transport)  $\gamma_{mix} := p_0 \otimes p_1$ ; refer to Proposition C.12 for more details. This generally does not hold, however, if  $\tilde{Z}_1$  is group-aware, or  $\tilde{Z}_1$  and  $\tilde{Z}_0$  are not independent; more details available in Appendix C.2.2.

**Remark 3.4.** All variance decompositions below compare under the lifted coupling  $\gamma^{\text{lift}}$  (i.e., they isolate the effect of marginalizing over  $G$  versus conditioning on  $G$ ).

We have the following sharp conditional variance comparison under the lifted coupling.

**Theorem 3.5** (Variance decomposition under group-aligned lift). *Assume  $\mathcal{G}$  acts orthogonally. Under the group-aligned lifted coupling  $G \sim \lambda$  independent of  $(\tilde{Z}_0, \tilde{Z}_1)$  and  $(Z_0, Z_1) = (G \cdot \tilde{Z}_0, G \cdot \tilde{Z}_1)$ , let  $Z_t = (1-t)Z_0 + tZ_1$  and  $\tilde{Z}_t = (1-t)\tilde{Z}_0 + t\tilde{Z}_1$ . With  $U := Z_1 - Z_0 = G \cdot \Delta$  and  $\Delta := \tilde{Z}_1 - \tilde{Z}_0$ ,*

$$\begin{aligned} \text{Var}(U | Z_t) &= \mathbb{E} \left[ \underbrace{\text{Var}(\Delta | \tilde{Z}_t | Z_t)}_{\text{within-slice difficulty}} \right] \\ &\quad + \underbrace{\text{Var}(\mathbb{E}[U | Z_t, G] | Z_t)}_{\text{symmetry ambiguity} \geq 0}. \end{aligned} \quad (13)$$

Consequently, invariant training without canonicalization admits larger conditional variance:

$$\mathbb{E}[\text{Var}(U | Z_t)] \geq \mathbb{E}[\text{Var}(\Delta | \tilde{Z}_t)]. \quad (14)$$

In particular, canonicalization eliminates the second term caused by symmetry ambiguity, and strict inequality typically holds whenever  $G$  remains ambiguous given  $Z_t$  and

the conditional mean drift depends on  $G$ ; see concrete posterior collision bounds in Appendix C.2.2.

**Remark 3.6** (Only non-equivariant models benefit from canonicalization). Notably, Theorem 3.5 only provides the lower bound on the irreducible error, but does not indicate all models are able to achieve this lower bound in practice. We emphasize that only non-equivariant models benefit from symmetry ambiguity elimination practically since canonicalization and canonical conditions expands the effective hypothesis class, echoing our analysis in Section 3.1.2. Equivariant models, by contrast, adhere the same ambient Bayes risk after canonicalization as in invariant training. More details are referred to Proposition C.15 in the Appendix.

### 3.2.3. WITHIN-SLICE SIMPLIFICATION WITH ALIGNED CANONICAL PRIOR AND MONGE COUPLING

Theorem 3.5 isolates two distinct sources of irreducible error: within-slice difficulty and symmetry ambiguity. Canonicalization reduces the second term, but the first term still depends on the slice coupling  $\tilde{\gamma}(\tilde{Z}_0, \tilde{Z}_1)$  and can remain large when the slice coupling is far from straight transport. We hereby present techniques targeting the first term: *aligned canonical priors* and *optimal transport*.

**Principles of aligned prior for canonical slice.** We start from the choice of prior  $\tilde{Z}_1 \sim q_1$ . Consider the case if one keeps a *misaligned* simple prior (e.g.  $q_1 = \mathcal{N}(0, I)$  while  $q_0$  on the canonical slice is far from centered/isotropic or is low-dimensional), then the slice coupling is still difficult and the first term in Theorem 3.5 can remain large; see Proposition C.19 for an example with closed-form within-slice conditional variance. Therefore, we propose two simple yet effective approaches to derive **aligned canonical pri-**

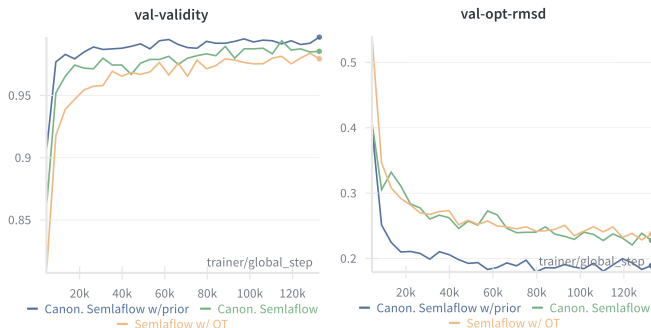


Figure 3. Training trajectories of our canonicalized model (visualization of the learning dynamics).

**ors:** (i) *learnable priors* within certain function classes (e.g. a Gaussian with learnable mean and variance); (ii) *KL projections of  $q_0$*  onto certain function classes (e.g. moment-matched Gaussian). For example, among Gaussian canonical priors, the moment matched Gaussian is  $q_1^* \sim \mathcal{N}(\mathbb{E}_{q_0}[\tilde{Z}_0], \text{Cov}_{q_0}(\tilde{Z}_0))$  (Proposition C.20).

**Optimal transport as complementary.** Given the canonical data  $q_0$  and the prior  $q_1$ , the slice conditional variance  $\mathbb{E}[\text{Var}(\tilde{Z}_1 - \tilde{Z}_0 \mid \tilde{Z}_t)]$  still depends on the coupling  $\tilde{\gamma}$ . In favorable regimes, e.g. Monge-like optimal transport (OT) under standard regularity conditions, this term can be small. Remarkably, OT has been applied in some previous work (Irwin et al., 2024), and we now formally clarify its effects on symmetry space. While OT generally reduces trajectory intersections and conditional variances regardless of whether canonicalization is applied, the solutions are not unique or stable in the presence of symmetry. Canonicalization stabilizes the solution by providing a gauge, and further benefits non-equivariant models with canonical conditions, as illustrated in Appendix C.2.2. Alternatively, canonicalization serves as an “implicit OT” by aligning the data-noise pair into the same gauge.

*Remark 3.7* (Canonicalization and optimal transport act complementarily). In the presence of symmetry, canonicalization eliminates group ambiguity, while optimal transport reduces within-slice difficulty. OT solutions can be closer to unique/Monge-like with canonicalization (Figure 5).

**Training-sampling consistency.** Suppose the network is conditioned on an auxiliary variable  $C$  (e.g. “canonical rank” positional encoding, a frame ID, or any deterministic or stochastic output of a canonicalizer). Training then implicitly targets a conditional field  $\tilde{v}(t, z, c)$  under a joint law  $\tilde{\pi}_t(c, z)$  induced by the training pipeline. We say there is *no condition mismatch at the start* if the inference-time joint law  $\tilde{\pi}_1^{\text{inf}}(C, \tilde{Z}_1)$  equals the training-time joint law  $\tilde{\pi}_1^{\text{tr}}(C, \tilde{Z}_1)$ . With either (i) aligned canonical priors, or (ii) isotropic priors and training-time OT, one keeps the training–sampling consistency. The training trajectories in Figure 3

strongly support our claim that canonicalization, especially when combined with aligned priors, significantly accelerates training.

## 4. Canonical Diffusion for Molecular Graph Generation

Building on the theoretical foundations established in Section 3, we now instantiate our canonical molecular graph diffusion framework under the joint symmetry group  $\mathcal{G} = S_N \times SO(3)$  (atom index permutations and global rotations).

### 4.1. Canonicalization of Molecular Graphs

The key design choice lies in constructing a canonicalizer that is both *geometrically principled* and *computationally tractable*. We decompose the full symmetry group, addressing permutation and rotation symmetries sequentially:

$$\Psi_\phi(\mathbf{Z}) = \Psi_\phi^{(\text{rot})} \left( \Psi_\phi^{(\text{perm})}(\mathbf{Z}) \right), \quad (15)$$

where  $\Psi_\phi^{(\text{perm})}$  establishes a canonical atom ordering and  $\Psi_\phi^{(\text{rot})}$  fixes a global reference frame.

We first consider **canonicalization via geometric spectra**, which breaks permutation symmetry of atom indices in a principled manner. The Fiedler vector, i.e. the eigenvector corresponding to the second smallest eigenvalue (the algebraic connectivity) of a graph’s Laplacian matrix, determines a canonical atom ordering in a natural way. Concretely, given a molecular, we build a symmetric kernel matrix  $\mathbf{W}(\mathbf{X}, \mathbf{A}) \in \mathbb{R}^{N \times N}$  from pairwise distances (e.g.,  $W_{ij} = k(\|\mathbf{X}_i - \mathbf{X}_j\|)$  for a radial kernel  $k$ ) to handle the 3D geometry, forming the normalized geometric Laplacian

$$\mathbf{L}(\mathbf{X}, \mathbf{A}) = \mathbf{D}^{-1}(\mathbf{D} - \mathbf{W}), \quad \mathbf{D} = \text{diag}(\mathbf{W}\mathbf{1}). \quad (16)$$

$\mathbf{L}$  depends only on distances and is rotation-invariant. We then use the Fiedler vector  $\mathbf{u}_2 \in \mathbb{R}^N$  to define a canonical ordering

$$\pi^*(\mathbf{Z}) := \text{argsort}(\mathbf{u}_2) \in S_N. \quad (17)$$

Under the standard mass–spring interpretation of the Laplacian, this corresponds to the lowest-frequency non-rigid vibration mode (Knyazev, 2018). This rank yields the classical spectral bisection and provides a principled seriation keeping strongly connected substructures contiguous (Pothén et al., 1990). In the molecular context, the spectral ordering serves as a *geometric-aware linearization*.

Finally, we optionally define a global  $SO(3)$  gauge for global rotations in a rank-consistent way after applying  $\pi^*$ . We choose a small set of anchor atoms in the canonical order (e.g., the two extremes define a primary axis and a third anchor determines a plane normal), and rotate

the centered coordinates into the resulting right-handed orthonormal frame. This yields a deterministic representative under the joint  $S_N \times SO(3)$  action; the full algorithm is deferred to Algorithm 2.

There are also other existing canonicalization methods, mostly for permutation in abstract graphs (Zhao et al., 2024; Ma et al., 2023; Dong et al., 2024) instead of geometric graphs, while existing OT methods mostly target coordinates; our method consistently consider the joint group in a stable manner. Different canonicalization methods may lead to various stability and practical performance, yet we experimentally find that our novel geometric spectra canonicalization outperforms others (Appendix D.1).

## 4.2. Canonical Diffusion and CanonFlow

**Canonicity-aware denoising architectures.** We elaborate on two types of architectural changes to enable non-equivariant denoising model in canonical diffusion: (i) *positional encoding only*, where we keep the original model architectures except that the model is conditioned on canonical rank embeddings  $C = \varphi(r_i)$  with normalized rank  $r_i = (i - 1)/(N - 1)$  to handle various molecular sizes; (ii) *Canon architecture*, our newly designed architecture that explicitly keeps atom canonical embeddings as a learnable hidden state  $C \in \mathbb{R}^{N \times d_c}$ , which interacts with the embeddings of  $\mathbf{X}, \mathbf{H}, \mathbf{A}$  in each layer (Figure 7); see more details in Appendix D.2. We use the prefix “Canon.” to indicate the first type, where everything is consistent with the base model except the PE; “CanonFlow” specifically refers to the flow matching model with *Canon* architecture.

**Training on the canonical slice.** Training proceeds via standard diffusion or flow matching procedures on the canonical slice  $\tilde{\mathbf{Z}} = \Psi_\phi(\mathbf{Z})$  and the slice distribution  $q_0 = (\Psi_\phi)\#p_0$ . For convenience, we illustrate with the flow matching framework, and diffusion models follow analogously. Given a coupling  $\tilde{\gamma}(\tilde{\mathbf{Z}}_0, \tilde{\mathbf{Z}}_1)$  between canonicalized data  $\tilde{\mathbf{Z}}_0 \sim q_0$  and the (optionally aligned) prior  $\tilde{\mathbf{Z}}_1 \sim q_1$ , we construct interpolated paths  $\tilde{\mathbf{Z}}_t = \Phi_t(\tilde{\mathbf{Z}}_0, \tilde{\mathbf{Z}}_1)$  for  $t \sim \text{Unif}[0, 1]$ . The model  $m_\theta$  learns to predict the conditional velocity field:

$$\min_{\theta} \mathbb{E} \left[ \left\| m_{\theta}(\tilde{\mathbf{Z}}_t, t; C) - u_t(\tilde{\mathbf{Z}}_t \mid \tilde{\mathbf{Z}}_0, \tilde{\mathbf{Z}}_1) \right\|^2 \right], \quad (18)$$

where  $u_t = \frac{d}{dt} \Phi_t$  including velocities for both continuous and discrete features. As explained in Section 3.2.3, we also optionally adopt an aligned prior, or enable OT coupling early in training and anneal its usage over epochs (termed as *OT anneal*) to stabilize optimization while avoiding over-specialization to aligned pairings.

**Sampling and restoring invariance.** Starting from the slice prior  $\tilde{\mathbf{Z}}_1 \sim q_1$  (the same marginal prior used in training), we integrate the learned dynamics to obtain a canonical

sample  $\hat{\tilde{\mathbf{Z}}}_0$  on the slice  $S$ ; no OT matching is required at inference. Denote the sampled canonical distribution as  $\tilde{\mu}$ , we can then restore the invariance by sampling  $g \sim \lambda$ ,  $\tilde{\mathbf{Z}} \sim \tilde{\mu}$  and outputting  $g \cdot \tilde{\mathbf{Z}}$ , thus the randomized distribution  $\mu := \int_G (g \cdot) \# \tilde{\mu} d\lambda(g)$  is invariant (Proposition D.1).

**Improved techniques with canonicity.** A subtle discrepancy arises between training and inference: canonical ranks are computed exactly from  $\Psi_\phi(\mathbf{Z})$  during training, whereas at inference the model generates from noise without access to ground-truth ranks. Moreover, in the spirit of (Dym et al., 2024), continuous canonicalization is impossible for  $S_N$  and  $SO(d)$ . To address this issue and to mitigate the train-test gap, we relaxed the ranks into continuous values and inject small perturbations to rank features during training, with an optional auxiliary head that self-estimates ranks from intermediate representations. In sampling, models either utilize fixed canonical conditions, or adaptively re-estimate the ranks to align the noisy data and the canonical conditions, which we termed as *projected canonical sampling* (PCS). Inspired by (Li et al., 2025), we also randomly drop canonical conditions with probability  $p_{\text{drop}}$  during training, and enable classifier-free guidance (CFG) (Ho & Salimans, 2022) extrapolating conditional and unconditional generation during inference. More training and sampling techniques are deferred to Appendix D.

## 5. Experiments

**Experimental Setup.** We evaluate our canonicalization approach on two widely adopted benchmarks for unconditional 3D molecule generation: QM9 (Ramakrishnan et al., 2014) and GEOM-DRUG (Axelrod & Gomez-Bombarelli, 2022). While QM9 consists of relatively small molecules, GEOM-DRUG provides a more challenging testbed that better differentiates the capabilities of various generative models. We adopt identical train/validation/test splits as prior work (Vignac et al., 2023; Le et al., 2023). Following (Irwin et al., 2024), we exclude molecules containing more than 72 atoms from the GEOM-DRUG training set for computational efficiency, which accounts for approximately 1% of the data; the validation and test sets remain intact. During evaluation, we sample molecule sizes from the empirical distribution in the test set and generate the corresponding number of atoms by numerically integrating the learned flow dynamics. All reported metrics are computed over 10k generated molecules with 3 independent runs. Without specification, we adopt the settings and hyperparameters in SemlaFlow (Irwin et al., 2024), which serves as our base model. By default, canonicalization is conducted on the  $S_N$  group; results of canonicalizing  $S_N \times SO(3)$  as well as *more ablation studies* are deferred to Appendix E.

**Main results.** Applying the canonical diffusion framework to SemlaFlow and CanonFlow, we achieve state-of-

Table 1. Performance of canonicalized SemlaFlow on QM9 dataset. Atom stability, molecule stability, and validity are reported as percentages; Opt-RMSD is reported in Å.

Methods	Atom Stab $\uparrow$	Mol Stab $\uparrow$	Valid $\uparrow$	Opt-RMSD $\downarrow$	NFE $\downarrow$
EDM	98.7	82.0	91.9	–	1000
GCDM	98.7	85.7	94.8	–	1000
MUDiff	98.3	89.9	95.3	–	1000
FlowMol	99.7	96.2	97.3	–	100
MiDi	99.8	97.5	97.9	–	500
EQGAT-diff	99.9 $\pm$ 0.0	98.7 $\pm$ 0.18	99.0 $\pm$ 0.16	–	500
SemlaFlow w/ OT	99.9 $\pm$ 0.0	99.56 $\pm$ 0.07	99.31 $\pm$ 0.08	0.23 $\pm$ 0.00	100
Canon. SemlaFlow w/ OT anneal + PCS (ours)	99.9 $\pm$ 0.0	99.62 $\pm$ 0.08	99.49 $\pm$ 0.06	0.21 $\pm$ 0.00	100
Canon. SemlaFlow w/ Prior + PCS (ours)	99.9 $\pm$ 0.0	<b>99.64</b> $\pm$ 0.04	<b>99.58</b> $\pm$ 0.06	<b>0.17</b> $\pm$ 0.00	100

Table 2. Performance of canonicalized SemlaFlow and CanonFlow on GEOM-DRUG dataset. Atom stability, molecule stability, validity, as well as uniqueness and novelty are reported as percentages.

Methods	Atom Stab $\uparrow$	Mol Stab $\uparrow$	Valid $\uparrow$	Unique $\uparrow$	Novel $\uparrow$	NFE $\downarrow$
FlowMol	99.0	67.5	51.2	–	–	100
MiDi	99.8	91.6	77.8	100.0	100.0	500
EQGAT-diff	99.8 $\pm$ 0.0	93.4 $\pm$ 0.21	94.6 $\pm$ 0.24	100.0 $\pm$ 0.0	99.9 $\pm$ 0.07	500
SemlaFlow w/ OT	99.8 $\pm$ 0.0	97.3 $\pm$ 0.08	93.9 $\pm$ 0.19	100.0 $\pm$ 0.0	99.6 $\pm$ 0.03	100
Canon. SemlaFlow w/ OT anneal (ours)	99.8 $\pm$ 0.0	98.1 $\pm$ 0.03	95.0 $\pm$ 0.20	100.0 $\pm$ 0.0	99.7 $\pm$ 0.02	100
CanonFlow w/ OT anneal (ours)	<b>99.9</b> $\pm$ 0.0	<b>98.4</b> $\pm$ 0.02	<b>95.9</b> $\pm$ 0.08	100.0 $\pm$ 0.0	99.7 $\pm$ 0.01	100

the-art results on most metrics across both datasets.

Table 1 shows results on QM9. While all methods achieve near-saturated atom stability, our canonicalized variants improve upon SemlaFlow in both molecule stability and validity. Notably, Canon. SemlaFlow (Prior + PCS) achieves the lowest Opt-RMSD (a 26% reduction compared to the baseline), indicating that canonicalization guides the model toward energy-minimized conformations. The training dynamics on QM9 illustrated in Figure 3 also demonstrates that canonicalization significantly accelerates training convergence.

Table 2 presents results on the more challenging GEOM-DRUG dataset, which contains larger and structurally diverse drug-like molecules that better differentiate model capabilities. On this benchmark, our canonicalized approach demonstrates clear improvements over baselines. In particular, Canon. SemlaFlow (trained w/ CFG) consistently improves validity and molecule stability over SemlaFlow baseline. Remarkably, our CanonFlow achieves **SOTA** performance across atom stability, molecule stability and validity – the most important metrics, outperforming all previous models with a large margin (surpassing SemlaFlow baseline 1.1% in molecule stability and 2.0% in validity).

We also report *few-step generation* performance in Table 6, Table 7 and Table 8. The advantages of canonicalized models are still significant with even only 50 steps, which is only ten percent of most previous models. Notably, our method almost induces no overheads in computation or sample time, validating the efficient and effective guidance

signal of canonical conditioning in sampling.

## 6. Conclusion

We studied symmetry in diffusion and flow-based generative modeling through the lens of canonicalization. For molecular data, where representations are invariant to permutations and rigid motions, we showed that symmetry can induce latent *group ambiguity* at intermediate noise levels, inflating conditional variance and making few-step solvers less reliable. Canonicalization provides a principled gauge fixing, removing the symmetry-induced ambiguity term and enabling the use of expressive non-equivariant backbones. The remaining within-slice difficulty can be mitigated by better aligned priors and near-Monge couplings during training. Canonical diffusion yields state-of-the-art results across several metrics on popular geometric molecule generation benchmarks, with significant advantages in training convergence. While the current experimental scope of this paper lies in unconditional small molecule generation, we hope to extend to other biomolecule (e.g., protein) structure prediction and conditional generation (e.g., binder design) tasks with more potential symmetry groups in the future.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning, with applications in molecular generative modeling for chemistry and biology. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630 (8016):493–500, 2024.
- Albergo, M. S., Boffi, N. M., and Vanden-Eijnden, E. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- Axelrod, S. and Gomez-Bombarelli, R. Geom, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data*, 9(1):185, 2022.
- Bo, D., Shi, C., Wang, L., and Liao, R. Specformer: Spectral graph neural networks meet transformers. *ArXiv*, abs/2303.01028, 2023.
- Cai, C., Hy, T. S., Yu, R., and Wang, Y. On the connection between mpnn and graph transformer. In *International conference on machine learning*, pp. 3408–3430. PMLR, 2023.
- Campbell, A., Benton, J., Bortoli, V. D., Rainforth, T., Deligiannidis, G., and Doucet, A. A continuous time framework for discrete denoising models. 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=DmT862YAieY>.
- Cao, N. D. and Kipf, T. Molgan: An implicit generative model for small molecular graphs, 2022. URL <https://arxiv.org/abs/1805.11973>.
- Chen, X., He, J., Han, X., and Liu, L.-P. Efficient and degree-guided graph generation via discrete diffusion modeling. *arXiv preprint arXiv:2305.04111*, 2023.
- Corso, G., St’ark, H., Jing, B., Barzilay, R., and Jaakkola, T. Diffdock: Diffusion steps, twists, and turns for molecular docking. In *International Conference on Learning Representations*, 2023.
- Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Ding, Y. and Hofmann, T. Scalable non-equivariant 3d molecule generation via rotational alignment. *arXiv preprint arXiv:2506.10186*, 2025.
- Dong, Z., Zhang, M., Payne, P., Province, M., Cruchaga, C., Zhao, T., Li, F., and Chen, Y. Rethinking the power of graph canonization in graph representation learning with stability. In *International Conference on Learning Representations*, volume 2024, pp. 30797–30817, 2024.
- Duval, A. A., Schmidt, V., Hernández-García, A., Miret, S., Malliaros, F. D., Bengio, Y., and Rolnick, D. Faenet: Frame averaging equivariant gnn for materials modeling. In *International Conference on Machine Learning*, pp. 9013–9033. PMLR, 2023.
- Dwivedi, V. P., Luu, A. T., Laurent, T., Bengio, Y., and Bresson, X. Graph neural networks with learnable structural and positional representations. *ArXiv*, abs/2110.07875, 2021.
- Dym, N., Lawrence, H., and Siegel, J. W. Equivariant frames and the impossibility of continuous canonicalization. *arXiv preprint arXiv:2402.16077*, 2024.
- Feng, W., Wang, L., Lin, Z., Zhu, Y., Wang, H., Dong, J., Bai, R., Wang, H., Zhou, J., Peng, W., Huang, B., and Zhou, W. Generation of 3d molecules in pockets via language model, 2023. URL <https://arxiv.org/abs/2305.10133>.
- 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.
- Gong, S., Li, M., Feng, J., Wu, Z., and Kong, L. Diffuseq: Sequence to sequence text generation with diffusion models. *arXiv preprint arXiv:2210.08933*, 2022.
- Hassan, M., Shenoy, N., Lee, J., Stark, H., Thaler, S., and Beaini, D. Et-flow: Equivariant flow-matching for molecular conformer generation, 2024. URL <https://arxiv.org/abs/2410.22388>.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *ArXiv*, abs/2006.11239, 2020.
- Ho, J., Salimans, T., Gritsenko, A., Chan, W., Norouzi, M., and Fleet, D. J. Video diffusion models. *Advances in neural information processing systems*, 35:8633–8646, 2022.
- Hong, H., Lin, W., and Tan, K. C. Accelerating 3d molecule generation via jointly geometric optimal transport, 2025. URL <https://arxiv.org/abs/2405.15252>.
- 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.

- Hou, X., Zhu, T., Ren, M., Bu, D., Gao, X., Zhang, C., and Sun, S. Improving molecular graph generation with flow matching and optimal transport, 2024. URL <https://arxiv.org/abs/2411.05676>.
- Hua, C., Luan, S., Xu, M., Ying, R., Fu, J., Ermon, S., and Precup, D. Mudiff: Unified diffusion for complete molecule generation, 2024. URL <https://arxiv.org/abs/2304.14621>.
- Huang, H., Sun, L., Du, B., and Lv, W. Learning joint 2d & 3d diffusion models for complete molecule generation, 2023. URL <https://arxiv.org/abs/2305.12347>.
- Huang, Y., Peng, X., Ma, J., and Zhang, M. Boosting the Cycle Counting Power of Graph Neural Networks with  $l^2$ -GNNs. *arXiv e-prints*, art. arXiv:2210.13978, October 2022. doi: 10.48550/arXiv.2210.13978.
- Irwin, R., Tibo, A., Janet, J. P., and Olsson, S. Semlaflow—efficient 3d molecular generation with latent attention and equivariant flow matching. *arXiv preprint arXiv:2406.07266*, 2024.
- Jang, Y., Kim, D., and Ahn, S. Hierarchical graph generation with  $k^2$ -trees. In *ICML 2023 Workshop on Structured Probabilistic Inference* {\&} *Generative Modeling*, 2023.
- Jiang, K., Cui, J., Dong, X., and Toni, L. Bures-wasserstein flow matching for graph generation, 2025. URL <https://arxiv.org/abs/2506.14020>.
- Jin, W., Barzilay, R., and Jaakkola, T. Junction tree variational autoencoder for molecular graph generation, 2019. URL <https://arxiv.org/abs/1802.04364>.
- Jing, B., Berger, B., and Jaakkola, T. Alphafold meets flow matching for generating protein ensembles. *arXiv preprint arXiv:2402.04845*, 2024.
- Kaba, S.-O., Mondal, A. K., Zhang, Y., Bengio, Y., and Ravanbakhsh, S. Equivariance with learned canonicalization functions. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 15546–15566. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/kaba23a.html>.
- Kim, J., Nguyen, T. D., Min, S., Cho, S., Lee, M., Lee, H., and Hong, S. Pure transformers are powerful graph learners. *ArXiv*, abs/2207.02505, 2022.
- Knyazev, A. V. On spectral partitioning of signed graphs, 2018. URL <https://arxiv.org/abs/1701.01394>.
- Kornilov, N., Mokrov, P., Gasnikov, A., and Korotin, A. Optimal flow matching: Learning straight trajectories in just one step, 2024. URL <https://arxiv.org/abs/2403.13117>.
- Lawrence, H., Hofgard, E., Chen, Y., Smidt, T., and Walters, R. Detecting symmetry-breaking in molecular data distributions. In *AI for Accelerated Materials Design-ICLR 2025*, 2025a.
- Lawrence, H., Portilheiro, V., Zhang, Y., and Kaba, S.-O. Improving equivariant networks with probabilistic symmetry breaking. *arXiv preprint arXiv:2503.21985*, 2025b.
- Le, T., Cremer, J., Noé, F., Clevert, D.-A., and Schütt, K. Navigating the design space of equivariant diffusion-based generative models for de novo 3d molecule generation. *arXiv preprint arXiv:2309.17296*, 2023.
- Lee, J., Kim, S., Moon, S., Kim, H., and Kim, W. Y. Fragfm: Hierarchical framework for efficient molecule generation via fragment-level discrete flow matching, 2025. URL <https://arxiv.org/abs/2502.15805>.
- Li, X., Thickstun, J., Gulrajani, I., Liang, P. S., and Hashimoto, T. B. Diffusion-lm improves controllable text generation. *Advances in neural information processing systems*, 35:4328–4343, 2022.
- Li, Z., Wang, X., Huang, Y., and Zhang, M. Is distance matrix enough for geometric deep learning? In *Advances in Neural Information Processing Systems*, volume 36, pp. 37413–37447, 2023.
- Li, Z., Wang, X., Kang, S., and Zhang, M. On the completeness of invariant geometric deep learning models. *arXiv preprint arXiv:2402.04836*, 2024.
- Li, Z., Zhou, C., Wang, X., Peng, X., and Zhang, M. Geometric representation condition improves equivariant molecule generation. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 36921–36953. PMLR, 13–19 Jul 2025. URL <https://proceedings.mlr.press/v267/li25dz.html>.
- Lim, D., Robinson, J., Zhao, L., Smidt, T., Sra, S., Maron, H., and Jegelka, S. Sign and basis invariant networks for spectral graph representation learning. *arXiv preprint arXiv:2202.13013*, 2022.
- Lin, Y., Helwig, J., Gui, S., and Ji, S. Equivariance via minimal frame averaging for more symmetries and efficiency. *arXiv preprint arXiv:2406.07598*, 2024.

- Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Lippmann, P., Gerhartz, G., Remme, R., and Hamprecht, F. A. Beyond canonicalization: How tensorial messages improve equivariant message passing. *arXiv preprint arXiv:2405.15389*, 2024.
- Liu, X., Gong, C., and Liu, Q. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- Luo, T., Mo, Z., and Pan, S. J. Fast graph generation via spectral diffusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Luo, Y. and Ji, S. An autoregressive flow model for 3d molecular geometry generation from scratch. In *International conference on learning representations (ICLR)*, 2022.
- Luo, Y., Yan, K., and Ji, S. Graphdf: A discrete flow model for molecular graph generation, 2021. URL <https://arxiv.org/abs/2102.01189>.
- Ma, G., Wang, Y., Lim, D., Jegelka, S., and Wang, Y. A canonicalization perspective on invariant and equivariant learning. *Advances in Neural Information Processing Systems*, 37:60936–60979, 2024.
- Ma, J., Wang, Y., and Wang, Y. Laplacian canonization: A minimalist approach to sign and basis invariant spectral embedding. *arXiv preprint arXiv:2310.18716*, 2023.
- Maron, H., Ben-Hamu, H., Serviansky, H., and Lipman, Y. Provably powerful graph networks. *ArXiv*, abs/1905.11136, 2019.
- Morris, C., Ritzert, M., Fey, M., Hamilton, W. L., Lenssen, J. E., Rattan, G., and Grohe, M. Weisfeiler and leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, pp. 4602–4609, 2019.
- Morris, C., Rattan, G., and Mutzel, P. Weisfeiler and leman go sparse: Towards scalable higher-order graph embeddings. *Advances in Neural Information Processing Systems*, 33:21824–21840, 2020.
- Murphy, R., Srinivasan, B., Rao, V., and Ribeiro, B. Relational pooling for graph representations. In *International Conference on Machine Learning*, pp. 4663–4673. PMLR, 2019.
- Pothen, A., Simon, H. D., and Liou, K.-P. Partitioning sparse matrices with eigenvectors of graphs. *SIAM Journal on Matrix Analysis and Applications*, 11(3):430–452, 1990. doi: 10.1137/0611030.
- Puny, O., Atzmon, M., Ben-Hamu, H., Misra, I., Grover, A., Smith, E. J., and Lipman, Y. Frame averaging for invariant and equivariant network design. *arXiv preprint arXiv:2110.03336*, 2021.
- Qin, Y., Madeira, M., Thanou, D., and Frossard, P. Defog: Discrete flow matching for graph generation. *arXiv preprint arXiv:2410.04263*, 2024.
- 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.
- Rampasek, L., Galkin, M., Dwivedi, V. P., Luu, A. T., Wolf, G., and Beaini, D. Recipe for a general, powerful, scalable graph transformer. *ArXiv*, abs/2205.12454, 2022.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Satorras, V. G., Hoogeboom, E., and Welling, M. E (n) equivariant graph neural networks. In *International conference on machine learning*, pp. 9323–9332. PMLR, 2021.
- Shi, C., Xu, M., Zhu, Z., Zhang, W., Zhang, M., and Tang, J. Graphaf: a flow-based autoregressive model for molecular graph generation, 2020. URL <https://arxiv.org/abs/2001.09382>.
- Shirzad, H., Velingker, A., Venkatachalam, B., Sutherland, D. J., and Sinop, A. K. Exphormer: Sparse transformers for graphs. *ArXiv*, abs/2303.06147, 2023.
- Simonovsky, M. and Komodakis, N. Graphvae: Towards generation of small graphs using variational autoencoders, 2018. URL <https://arxiv.org/abs/1802.03480>.
- Singer, U., Polyak, A., Hayes, T., Yin, X., An, J., Zhang, S., Hu, Q., Yang, H., Ashual, O., Gafni, O., et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. *ArXiv*, abs/1907.05600, 2019.
- 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*, 2021. URL <https://openreview.net/forum?id=PxTIG12RRHS>.

- 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. *Advances in Neural Information Processing Systems*, 36, 2023.
- Tahmasebi, B. and Jegelka, S. Generalization bounds for canonicalization: A comparative study with group averaging. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=n01Xaskyk5>.
- Tahmasebi, B. and Jegelka, S. Regularity in canonicalized models: A theoretical perspective. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025b.
- Tian, Q., Xu, Y., Yang, Y., Wang, Z., Liu, Z., Yan, P., and Li, X. Equiflow: Equivariant conditional flow matching with optimal transport for 3d molecular conformation prediction, 2024. URL <https://arxiv.org/abs/2412.11082>.
- Tong, A., Fatras, K., Malkin, N., Huguet, G., Zhang, Y., Rector-Brooks, J., Wolf, G., and Bengio, Y. Improving and generalizing flow-based generative models with minibatch optimal transport, 2024. URL <https://arxiv.org/abs/2302.00482>.
- Vignac, C., Krawczuk, I., Siraudin, A., Wang, B., Cevher, V., and Frossard, P. Digress: Discrete denoising diffusion for graph generation. *arXiv preprint arXiv:2209.14734*, 2022.
- Vignac, C., Osman, N., Toni, L., and Frossard, P. Midi: Mixed graph and 3d denoising diffusion for molecule generation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 560–576. Springer, 2023.
- Wang, C., Zhou, C., Gupta, S., Lin, Z., Jegelka, S., Bates, S., and Jaakkola, T. Learning diffusion models with flexible representation guidance. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=cIGfKdfy3N>.
- Wang, J., Luo, H., Qin, R., Wang, M., Fang, M., Zhang, O., Gou, Q., Su, Q., Shen, C., You, Z., Wan, X., Liu, L., Hsieh, C.-Y., Hou, T., and Kang, Y. 3d molecular pocket-based generation with token-only large language model. *ChemRxiv*, 2024(0819), 2024. doi: 10.26434/chemrxiv-2024-0ckgt-v2. URL <https://chemrxiv.org/doi/abs/10.26434/chemrxiv-2024-0ckgt-v2>.
- Watson, J. L., Juergens, D., Bennett, N. R., Trippe, B. L., Yim, J., Eisenach, H. E., Ahern, W., Borst, A. J., Ragotte, R. J., Milles, L. F., et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 620:1089–1100, 2023.
- Wu, L., Gong, C., Liu, X., Ye, M., and Liu, Q. Diffusion-based molecule generation with informative prior bridges. *Advances in Neural Information Processing Systems*, 35: 36533–36545, 2022.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? *ArXiv*, abs/1810.00826, 2018.
- Xu, M., Yu, L., Song, Y., Shi, C., Ermon, S., and Tang, J. Geodiff: A geometric diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*, 2022.
- 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.
- Xu, Y., Wang, Y., Luo, S., Gao, K., He, T., Liu, C., and He, D. Quotient-space diffusion model. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=3JPAkwSVc4>.
- Yan, Q., Liang, Z., Song, Y., Liao, R., and Wang, L. Swingnn: Rethinking permutation invariance in diffusion models for graph generation. *arXiv preprint arXiv:2307.01646*, 2023.
- You, J., Gomes-Selman, J. M., Ying, R., and Leskovec, J. Identity-aware graph neural networks. In *AAAI Conference on Artificial Intelligence*, 2021.
- You, Y., Zhou, R., Park, J., Xu, H., Tian, C., Wang, Z., and Shen, Y. Latent 3d graph diffusion. In *The Twelfth International Conference on Learning Representations*, 2023.
- Zhang, L., Ashouritaklimi, K., Teh, Y. W., and Cornish, R. Symdiff: Equivariant diffusion via stochastic symmetrisation. *arXiv preprint arXiv:2410.06262*, 2024.
- Zhang, M., Li, P., Xia, Y., Wang, K., and Jin, L. Labeling trick: A theory of using graph neural networks for multi-node representation learning. *Advances in Neural Information Processing Systems*, 34:9061–9073, 2021.
- Zhao, L., Ding, X., and Akoglu, L. Pard: Permutation-invariant autoregressive diffusion for graph generation. *Advances in Neural Information Processing Systems*, 37: 7156–7184, 2024.

Zhou, C., Wang, X., and Zhang, M. From relational pooling to subgraph GNNs: A universal framework for more expressive graph neural networks. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 42742–42768. PMLR, 2023a.

Zhou, C., Wang, X., and Zhang, M. Facilitating graph neural networks with random walk on simplicial complexes. In *Advances in Neural Information Processing Systems*, volume 36, pp. 16172–16206, 2023b.

Zhou, C., Wang, X., and Zhang, M. Unifying generation and prediction on graphs with latent graph diffusion. *Advances in Neural Information Processing Systems*, 37, 2024a.

Zhou, C., Yu, R., and Wang, Y. On the theoretical expressive power and the design space of higher-order graph transformers. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pp. 2179–2187. PMLR, 02–04 May 2024b.

## A. Notation Guide

This section summarizes the notation used in the theoretical development. We use tildes to denote quantities after canonicalization, i.e., variables living on the canonical slice. We do not repeat standard diffusion-model notation unless it is needed to distinguish ambient-space and slice-space objects.

Table 3. Spaces, group actions, and canonicalization notation.

Notation	Meaning
$\mathcal{M}$	Measurable space of structured objects. In molecular generation, an element is $\mathbf{Z} = (\mathbf{X}, \mathbf{H}, \mathbf{A})$ .
$\mathcal{G}$	Symmetry group acting on $\mathcal{M}$ ; after centering, typically $\mathcal{G} = S_N \times SO(3)$ .
$g \cdot \mathbf{Z}$	Action of a group element $g \in \mathcal{G}$ on an object $\mathbf{Z}$ .
$g\#\mu$	Pushforward of a measure $\mu$ by the map $\mathbf{Z} \mapsto g \cdot \mathbf{Z}$ .
$\mathcal{O}(\mathbf{Z})$	Orbit of $\mathbf{Z}$ : $\mathcal{O}(\mathbf{Z}) = \{g \cdot \mathbf{Z} : g \in \mathcal{G}\}$ .
$\mathcal{M}/\mathcal{G}$	Quotient space whose elements are group orbits.
$\text{Stab}(\mathbf{Z})$	Stabilizer of $\mathbf{Z}$ : $\{g \in \mathcal{G} : g \cdot \mathbf{Z} = \mathbf{Z}\}$ . A free action means the stabilizer is trivial.
$\lambda$	Haar probability measure on the compact group $\mathcal{G}$ ; uniform on $S_N$ and uniform over rotations on $SO(3)$ .
$\Psi$	Canonicalization map selecting an orbit representative, with $\Psi(\mathbf{Z}) \in \mathcal{O}(\mathbf{Z})$ and $\Psi(g \cdot \mathbf{Z}) = \Psi(\mathbf{Z})$ .
$S = \Psi(\mathcal{M})$	Canonical slice, i.e., the image of the canonicalization map.
$\kappa$	Gauge map or canonical frame satisfying $x = \kappa(x) \cdot \Psi(x)$ and $\kappa(g \cdot x) = g\kappa(x)$ when defined.
$\delta_{\mathbf{Z}}$	Dirac measure at $\mathbf{Z}$ , used in the factorization of invariant measures.

Table 4. Distributional and transport notation used in the theory.

Notation	Meaning
$\mu$	Target distribution on the ambient object space $\mathcal{M}$ , usually assumed $\mathcal{G}$ -invariant.
$\nu = \Psi\#\mu$	Slice distribution induced by pushing $\mu$ forward through the canonicalizer $\Psi$ .
$p_0, p_1, p_t$	Ambient data, prior, and intermediate/noised distributions, respectively.
$q_0, q_1$	Canonical-slice data distribution and slice prior distribution, respectively.
$q_1^*$	Aligned slice prior, e.g., the moment-matched Gaussian approximation to $q_0$ in the aligned-prior discussion of Appendix C.2.
$\mathbf{Z}, \tilde{\mathbf{Z}}$	Ambient variable and its canonicalized slice representative.
$\mathbf{Z}_t, \tilde{\mathbf{Z}}_t$	Intermediate ambient state and intermediate slice state along a diffusion or flow path.
$G \sim \lambda$	Latent random symmetry element used to lift slice variables back to the ambient space.
$\gamma$	Coupling between ambient endpoints $(\mathbf{Z}_0, \mathbf{Z}_1)$ .
$\tilde{\gamma}$	Coupling between slice endpoints $(\tilde{\mathbf{Z}}_0, \tilde{\mathbf{Z}}_1)$ .
$\gamma^{\text{lift}}$	Group-aligned lift of a slice coupling, defined by $(\mathbf{Z}_0, \mathbf{Z}_1) = (G \cdot \tilde{\mathbf{Z}}_0, G \cdot \tilde{\mathbf{Z}}_1)$ with $G \sim \lambda$ .
$\gamma_{\text{mix}}$	Product ambient coupling, typically $p_0 \otimes p_1$ , recovered by the lifted product slice coupling under an isotropic prior.
$\otimes$	Product coupling or product measure.
$\Phi_t$	Conditional interpolation path mapping endpoints to an intermediate state.
$u_t(\cdot   \mathbf{Z}_0, \mathbf{Z}_1)$	Microscopic conditional vector field associated with a chosen endpoint pair.
$v_t$	Marginal vector field, given by the conditional expectation of the microscopic field.

## B. Related Work

**Symmetries in generative models.** Many datasets of interest in generative modeling are naturally defined only up to group actions—e.g., rotations and translations in 3D, or permutations in sets and graphs, so a principled generator should respect these symmetries to avoid spurious modes and overcounting equivalent configurations. A dominant approach is to build symmetry into the architecture by enforcing group equivariance (and thus invariance of the induced distribution), as in E(n)/SE(3)-equivariant message passing networks and their use in equivariant flows and diffusion/score-based models for geometric data (Satorras et al., 2021; Garcia Satorras et al., 2021; Hoogeboom et al., 2022). Systematic studies further clarify the design space and practical trade-offs of equivariant generative modeling under such geometric symmetries (Le et al., 2023; Lawrence et al., 2025b;a). An alternative architecture-agnostic perspective achieves symmetry by stochastic symmetrization (Zhang et al., 2024), which can enable the use of expressive and scalable non-equivariant backbones; learned canonicalization has been shown to provide a general mechanism for constructing equivariant functions from

Table 5. Proof-level operators and auxiliary quantities.

Notation	Meaning
$U = \mathbf{Z}_1 - \mathbf{Z}_0$	Ambient displacement in the linear-path flow-matching analysis.
$\Delta = \bar{\mathbf{Z}}_1 - \bar{\mathbf{Z}}_0$	Slice displacement in the canonicalized linear-path analysis.
$\text{Var}(\cdot   \cdot)$	Conditional variance; the irreducible flow-matching regression error is an expected conditional variance.
$\text{Cov}(\cdot   \cdot)$	Conditional covariance, used in the Gaussian supplement for within-slice variance.
$w_m(x)$	Responsibility weight of the $m$ -th group copy in the symmetry-mixture score formula.
$\varrho(g   z)$	Posterior distribution over the latent group element given the intermediate ambient state.
$m_g(z), m(z)$	Group-conditioned and marginal conditional mean drifts used to quantify symmetry ambiguity.
$\mathcal{L}_t(v)$	Population squared loss for a candidate vector field at time $t$ .
$v^*$	Bayes regressor, i.e., the conditional mean vector field minimizing $\mathcal{L}_t(v)$ .
$(g, g)_{\#} \gamma = \gamma$	Diagonal invariance of a coupling under simultaneous group action on both endpoints.
$\mathcal{H}_{\text{eq}}$	Class of equivariant functions or vector fields, used when comparing equivariant and canonicalized model classes.
$\mathcal{H}_{\text{lift}}(\mathcal{A})$	Equivariant class induced by lifting functions from a slice function class $\mathcal{A}$ .
$\text{TV}(\cdot, \cdot)$	Total variation distance, used to compare induced model and target distributions.
$\text{KL}(\cdot    \cdot)$	Kullback-Leibler divergence, used to define the Gaussian projection for aligned priors.

canonical representatives (Kaba et al., 2023). However, canonicalization is not “free”: for common groups there are fundamental continuity/stability obstructions, motivating softened or averaged constructions such as weighted frames to mitigate discontinuities near symmetric configurations (Dym et al., 2024; Lin et al., 2024; Ma et al., 2024). Related theory analyzes when canonicalization or group averaging offers better statistical behavior, highlighting distinct generalization regimes (Tahmasebi & Jegelka, 2025a; Puny et al., 2021; Lippmann et al., 2024; Tahmasebi & Jegelka, 2025b; Duval et al., 2023). Complementing these perspectives, recent work in graph diffusion argues that enforcing strict invariance can make learning harder by inducing mixture-like objectives over symmetry transformations, and proposes post-hoc group randomization at sampling time to recover invariance without constraining the training model (Yan et al., 2023). These works motivate the broader question of whether symmetry should be imposed architecturally or removed before learning.

**Molecular generative models.** Early work largely focused on generating discrete molecular structures via string or graph-based parameterizations, including structured latent-variable models that explicitly construct chemically valid graphs (Jin et al., 2019; Simonovsky & Komodakis, 2018; Jang et al., 2023). A major recent line uses diffusion processes on discrete graph attributes. Discrete denoising diffusion models operate directly on categorical attributes and have demonstrated strong performance at scale on molecular benchmarks (Vignac et al., 2022). In parallel, flow-based methods provide alternative likelihood-based or transport-inspired formulations for molecular graph generation (Luo et al., 2021; Shi et al., 2020), and discrete flow matching further improves sampling flexibility and efficiency while retaining strong generation quality (Qin et al., 2024; Hou et al., 2024; Lee et al., 2025; Chen et al., 2023; Luo et al., 2023).

Generating molecules directly in 3D space has been explored with sequential and autoregressive schemes that place atoms step-by-step while maintaining geometric consistency (Gebauer et al., 2019; Luo & Ji, 2022; Feng et al., 2023; Wang et al., 2024). Diffusion models have become a dominant paradigm for 3D molecular generation by learning to denoise corrupted coordinates with architectures designed to respect Euclidean symmetries (Hoogeboom et al., 2022). Closely related conditional settings, such as conformer generation given a fixed molecular graph, have also benefited from geometric diffusion and flow matching formulations (Xu et al., 2022; Hassan et al., 2024; Hong et al., 2025; Wu et al., 2022).

A growing body of work aims to jointly generate discrete molecular graphs and continuous 3D geometries to avoid brittle post-hoc bond inference and to better couple chemical structure with spatial arrangement. MiDi (Vignac et al., 2023) proposes a mixed discrete–continuous diffusion approach that generates molecular graphs together with conformers in an end-to-end differentiable manner. Subsequent efforts further explore joint diffusion over comprehensive molecular representations (Huang et al., 2023; Hua et al., 2024; Irwin et al., 2024).

There are also latent generative models developed for 2D, 3D, or joint representations of molecules, including GeoLDM (Xu et al., 2023), LGD (Zhou et al., 2024a), and LDM-3DG (You et al., 2023). More recently, GeoRCG (Li et al., 2025) further utilizes a two-stage generation: first generating molecule representations, then using geometric representations to guide the molecule generation. By contrast, REED (Wang et al., 2025) leverages the pretrained molecular representation as the guidance, aligning internal diffusion model features with these high-quality representations to accelerate training and

improve generation performance. We emphasize that our analysis and proposed methods are also applicable and effective for these latent/representation space generative models under mild conditions.

**Non-equivariant, alignment-based, and quotient-space biochemical generators.** A related line of work in molecular and biomolecular generation relaxes fully equivariant architectures and instead handles rigid-motion symmetries through alignment, augmentation, symmetrization, or quotient-space modeling. In 3D molecule and conformer generation, geometric diffusion models often align coordinates or targets when defining losses or sampling procedures to reduce redundant rotational degrees of freedom (Xu et al., 2022), and recent work explicitly studies scalable non-equivariant 3D molecule generation via rotational alignment (Ding & Hofmann, 2025). Similar ideas also appear in biomolecular structure modeling and design, where protein and complex generation systems rely on geometric alignment, frame choices, or symmetry-aware losses to stabilize learning over equivalent structures (Watson et al., 2023; Corso et al., 2023; Abramson et al., 2024; Jing et al., 2024). Architecture-agnostic stochastic symmetrization SymDiff further shows that equivariance can be recovered without hard-coding it into the backbone (Zhang et al., 2024). However, it still requires an  $S_N$ -equivariant backbone, only relaxing  $O(3)$  – our approach relaxed both  $S_N$  and  $SE(3)$ , enabling fully non-equivariant architectures with positional encodings. Most closely related to our motivation, Quotient-Space Diffusion Models formulate diffusion directly on the quotient space to avoid learning redundant symmetry directions and to provide a principled sampler, in contrast to heuristic alignment strategies (Xu et al., 2026). Our work shares this quotient-space view but takes a different route: rather than defining the diffusion intrinsically on the quotient, we choose a measurable canonical slice, train a flexible non-equivariant model on the slice, and recover the invariant distribution by Haar randomization. This canonicalization perspective gives an explicit slice distribution, a conditional-variance explanation for training acceleration, and a practical mechanism for combining non-equivariant backbones with aligned priors and optimal transport.

**Optimal transport.** In continuous-time generative modeling, optimal transport (OT) is often used as a principle for selecting cost-effective and more structured probability paths between a simple base distribution and the data distribution. A widely used choice in flow matching is the OT displacement interpolation, which yields trajectories closer to minimal-cost transport and empirically leads to faster convergence and fewer sampling steps (Lipman et al., 2022). Related “trajectory straightening” viewpoints, such as Rectified Flow, can be interpreted as progressively transforming the learned dynamics toward straighter transport-like paths, further improving numerical stability and reducing the number of integration steps required at inference (Liu et al., 2022; Tong et al., 2024; Kornilov et al., 2024).

In molecular generation, OT-inspired objectives have been adopted more explicitly to enable high-quality few-step generation under geometric symmetries. For 2D molecular graph generation, MolGAN (Cao & Kipf, 2022) was among the first to successfully use the Wasserstein-1 distance to stabilize the training of molecular graph generators. Subsequent work such as BWFlow (Jiang et al., 2025) also demonstrates the effectiveness of OT in 2D molecule design. For 3D molecule generation, various methods have demonstrated that OT can serve not only as a theoretical lens but also as a practical design tool for fast, high-fidelity molecule synthesis (Song et al., 2023; Tian et al., 2024; Hong et al., 2025).

## C. Omitted Proof

This section provides all omitted proof and more detailed discussions on our theoretical results, presented mainly in Section 3.

### C.1. Proof for Section 3.1

This subsection proves that canonicalized generative models can induce universal invariant distributions and outperform those equivariant baselines trained directly on invariant data.

#### C.1.1. PROOF FOR SECTION 3.1.1

We first show the correctness and sufficiency of canonical slice modeling.

**Theorem C.1** (Factorization of invariant measures; Theorem 3.1 in main text). *Suppose Assumptions 2.5 and 2.6 hold. Let  $\mu$  be any  $\mathcal{G}$ -invariant probability measure on  $\mathcal{M}$ . Let  $\Psi$  be an orbit representative map defined  $\mu$ -a.s., and let  $\nu = \Psi\#\mu$  be the slice distribution on  $S = \Psi(\mathcal{M})$ . Then*

$$\mu = \int_S \left( \int_{\mathcal{G}} \delta_{g \cdot \mathbf{z}} d\lambda(g) \right) d\nu(\mathbf{Z}). \quad (19)$$

Equivalently, if  $\tilde{\mathbf{Z}} \sim \nu$ ,  $g \sim \lambda$  independent, then  $g \cdot \tilde{\mathbf{Z}} \sim \mu$ .

*Proof.* This is a known result, included for completeness. Let  $\mathbf{Z} \sim \mu$  and define  $\tilde{\mathbf{Z}} = \Psi(\mathbf{Z}) \in S$ . By definition of orbit representative, there exists  $g(\mathbf{Z}) \in \mathcal{G}$  with  $\mathbf{Z} = g(\mathbf{Z}) \cdot \tilde{\mathbf{Z}}$ . Under 2.5,  $g(\mathbf{Z})$  is unique a.s. Fix any measurable  $A \subseteq \mathcal{G}$ . For any  $h \in \mathcal{G}$ , invariance of  $\mu$  implies

$$\mathbb{P}(g(\mathbf{Z}) \in A \mid \tilde{\mathbf{Z}}) = \mathbb{P}(h g(\mathbf{Z}) \in A \mid \tilde{\mathbf{Z}}),$$

so the conditional law of  $g(\mathbf{Z})$  given  $\tilde{\mathbf{Z}}$  is left-invariant. By uniqueness of Haar probability on compact  $\mathcal{G}$ , this conditional law is  $\lambda$ . Therefore,

$$\begin{aligned} \mu(B) &= \mathbb{E}[\mathbf{1}\{g(\mathbf{Z}) \cdot \tilde{\mathbf{Z}} \in B\}] = \mathbb{E}\left[\int_{\mathcal{G}} \mathbf{1}\{g \cdot \tilde{\mathbf{Z}} \in B\} d\lambda(g)\right] \\ &= \int_S \left(\int_{\mathcal{G}} \delta_{g \cdot \mathbf{z}}(B) d\lambda(g)\right) d\nu(\mathbf{z}), \end{aligned} \quad (20)$$

which is (4). □

**Corollary C.2** (Sufficiency of slice modeling; Corollary 3.2 in main text). *To model any invariant target  $\mu$ , it suffices to model the slice distribution  $\nu$ ; invariance is recovered by Haar randomization.*

In the next part, we show the universality of canonicalized parameterizations over invariant and equivariant functions.

**Proposition C.3** (Universality over canonicalized parameterizations; Proposition 3.3 in the main text). *Let  $\mathcal{G}$  act continuously and orthogonally on a compact set  $K \subset \mathbb{R}^d$ . Suppose we have a (measurable) canonicalization map  $\Psi : K \rightarrow K$  as Theorem 2.3, and a (measurable) gauge map  $\kappa : K \rightarrow \mathcal{G}$  s.t.*

$$\Psi(g \cdot x) = \Psi(x), \quad \kappa(g \cdot x) = g\kappa(x), \quad x = \kappa(x) \cdot \Psi(x) \quad (21)$$

(defined on a full-measure set; under free actions,  $\kappa$  is unique a.s.). Consider the parametrization

$$\phi(x) = \kappa(x) \cdot f(\Psi(x)), \quad (22)$$

where  $f$  is a universal approximator on  $\Psi(K)$ . Then  $\phi$  is a universal approximator of continuous  $\mathcal{G}$ -equivariant functions on  $K$ , and  $f \circ \Psi$  is universal for continuous  $\mathcal{G}$ -invariant functions on  $K$ .

*Proof.* Let  $F : K \rightarrow \mathbb{R}^d$  be continuous and  $\mathcal{G}$ -equivariant:  $F(g \cdot x) = g \cdot F(x)$ . Define its restriction to the slice  $S := \Psi(K)$  by  $F_S(z) := F(z)$  for  $z \in S$ . For any  $x \in K$ , write  $x = \kappa(x) \cdot \Psi(x)$ . By equivariance,

$$F(x) = F(\kappa(x) \cdot \Psi(x)) = \kappa(x) \cdot F(\Psi(x)) = \kappa(x) \cdot F_S(\Psi(x)). \quad (23)$$

Since  $f$  is universal on  $S$ , for any  $\varepsilon > 0$  there exists  $f$  such that  $\sup_{z \in S} \|f(z) - F_S(z)\| < \varepsilon$ . Then

$$\sup_{x \in K} \|\phi(x) - F(x)\| \leq \sup_{z \in S} \|f(z) - F_S(z)\| < \varepsilon, \quad (24)$$

using that the group action is norm-preserving. □

### C.1.2. PROOF FOR SECTION 3.1.2

We now show the superiority of canonicalized generative models with expressive (possibly) non-equivariant networks over standard generative models with equivariant backbones.

**Expressivity gain in  $S_N$  from canonical ordering.** It is known that the expressivity of message-passing based GNNs is bounded by the first-order Weisfeiler–Lehman (1-WL) algorithm on abstract graphs in terms of distinguishing isomorphic graphs (Xu et al., 2018). This immediately implies their inability in counting substructures such as cycles (Huang et al., 2022) and even link prediction (Zhang et al., 2021). Analogous issues exist for transformers due to their equivalence with MPNNs (Cai et al., 2023), and equivariant GNNs including EGNNs (Satorras et al., 2021) when processing 3D geometric graphs (Li et al., 2023; 2024). There are some mainstream methods to overcome the permutation-related expressivity pitfalls:

1. *High-order architectures*: high-order GNNs (Morris et al., 2019; 2020; Maron et al., 2019) and high-order transformers (Kim et al., 2022; Zhou et al., 2024b) have provably more powerful expressivity within and beyond the  $k$ -WL framework, at the cost of exponentially higher computation; for example, SwinGNN (Yan et al., 2023) adopts 3-WL equivalent PPGN (Maron et al., 2019) as the denoising backbone. We argue, however, stronger expressivity does not need to emerge from high-order networks with exponential complexity: symmetry breaking also improves expressivity as detailed as follows.
2. *Symmetry breaking*: by introducing asymmetry intentionally, the expressivity of underlying network can be provably enhanced. A widely adopted technique is “labeling” the nodes with extra identity features (You et al., 2021; Zhang et al., 2021; Huang et al., 2022; Zhou et al., 2023a), and the expressivity monotonically increases with the number of IDs. Unfortunately, these methods need the relational pooling or averaging to produce invariant outputs, which again introduce complexity at the scale of order of the group - canonicalization offers a way to produce appropriate and stable outputs (Dong et al., 2024) without the costly traverse and averaging procedure.
3. *Positional encodings and structural encodings*: graph structure-based PEs significantly improve the expressivity and practical performance of neural backbones (Rampasek et al., 2022). Popular PEs include spatial structure-based (e.g. RWSE (Dwivedi et al., 2021) and RRWP (Shirzad et al., 2023)) and spectra-based (using eigenvalues and eigenvectors of graph Laplacians (Ma et al., 2023; Lim et al., 2022; Bo et al., 2023) or Hodge Laplacians (Zhou et al., 2023b)). However, different from the fixed positional encodings in images or texts based on absolute indices of tokens, the calculations of above PEs on graphs rely on ground truth structures, which is not available in diffusion generation. Some generative models like DeFoG (Qin et al., 2024) utilize an estimated PE from the current noisy graph, which we argue might not be the optimal choice. In comparison, our canonicalization rank can also be viewed as spatial and spectral PEs, yet yields almost no computation overhead and no train-test gap.

The three observations above are specific to the  $S_N$  side of the symmetry: 1-WL limitations, node IDs, and graph positional encodings all target the loss of information caused by quotienting over vertex permutations. In this sense, canonical ordering provides a stable symmetry-breaking mechanism that strengthens the effective hypothesis class on the permutation quotient.

**Expressivity gain in  $SO(3)$  from measurable gauge fixing.** The rotational side gives a different, complementary separation. Practical equivariant architectures built from continuous layers realize continuous equivariant maps, whereas canonicalization may use a measurable gauge that is discontinuous at unavoidable frame-switching sets. This allows canonicalized non-equivariant models to represent equivariant targets outside the continuous equivariant class.

**Proposition C.4** (Gauge lifting and strict separation for  $SO(3)$ ). *Let  $\mathcal{X} \subset \mathbb{R}^{N \times 3}$  be a space of centered point clouds with the  $SO(3)$  action  $R \cdot X := XR^\top$ . Let  $\mathcal{X}_0 \subseteq \mathcal{X}$  be a full-measure subset on which the action is free. Suppose there exist measurable maps  $\Psi : \mathcal{X}_0 \rightarrow S$  and  $\kappa : \mathcal{X}_0 \rightarrow SO(3)$  such that, for all  $R \in SO(3)$  and  $X \in \mathcal{X}_0$ ,*

$$\Psi(R \cdot X) = \Psi(X), \quad \kappa(R \cdot X) = R\kappa(X), \quad X = \kappa(X) \cdot \Psi(X). \quad (25)$$

For any measurable  $a : S \rightarrow SO(3)$ , define

$$F_a(X) := \kappa(X)a(\Psi(X)). \quad (26)$$

Then  $F_a$  is  $SO(3)$ -equivariant on  $\mathcal{X}_0$ . Moreover, every continuous  $SO(3)$ -equivariant map  $f : \mathcal{X} \rightarrow SO(3)$  is represented by this form with  $a = f|_S$ . If  $\kappa$  is not equal on  $\mathcal{X}_0$  to the restriction of any continuous  $SO(3)$ -equivariant map on  $\mathcal{X}$ , then the canonicalized class strictly contains the continuous equivariant class.

*Proof.* Using (25), for any  $R \in SO(3)$ ,

$$F_a(R \cdot X) = \kappa(R \cdot X)a(\Psi(R \cdot X)) = R\kappa(X)a(\Psi(X)) = RF_a(X), \quad (27)$$

so  $F_a$  is equivariant. For any continuous equivariant  $f$ , choose  $a = f|_S$ . Since  $X = \kappa(X) \cdot \Psi(X)$ ,

$$f(X) = f(\kappa(X) \cdot \Psi(X)) = \kappa(X)f(\Psi(X)) = F_a(X). \quad (28)$$

For strictness, choose  $a \equiv I$ , so  $F_a = \kappa$ . If  $\kappa$  is not the restriction of any continuous equivariant map, then  $F_a$  lies in the canonicalized class but not in the continuous equivariant class. This obstruction is consistent with the impossibility of globally continuous canonical frames for groups such as  $SO(d)$  (Dym et al., 2024).  $\square$

*Remark C.5.* Theorem C.4 is a function-class statement, not a continuity claim about the learned neural backbone. The non-equivariant backbone only needs to model functions on the canonical slice, while the possibly discontinuous gauge is supplied by the canonicalizer. Stabilizers or degenerate frames can be handled by restricting to the full-measure free-action set, as in Assumption 2.5.

**TV guarantee for canonicalized non-equivariant backbones.** The preceding  $S_N$  and  $SO(3)$  discussions justify why the slice model class induced by canonicalization can be at least as expressive as an equivariant ambient baseline. We next show that any such slice-level advantage transfers to the invariant target distribution after Haar randomization.

**Theorem C.6** (Canonicalized models match (or improve) compact-group equivariant baselines in TV). *Under Assumptions 2.5 and 2.6, let  $\mathcal{G}$  be the compact symmetry group after centering, e.g.,  $S_N$ ,  $SO(3)$ , or  $S_N \times SO(3)$ . Let  $\mu$  be any  $\mathcal{G}$ -invariant target distribution on  $\mathcal{M}$  and let  $\nu := \Psi_{\#}\mu$  be its canonical-slice distribution on  $S = \Psi(\mathcal{M})$ . Assume there exists an equivariant baseline (defined in ambient space) producing an invariant model distribution  $\hat{\mu}_{\theta(\text{eq})}$  such that*

$$\text{TV}(\hat{\mu}_{\theta(\text{eq})}, \mu) < \varepsilon. \quad (29)$$

Let  $\hat{\nu}_{\theta(\text{eq})} := \Psi_{\#}\hat{\mu}_{\theta(\text{eq})}$  be the induced slice distribution. If the slice model class  $\{\hat{\nu}_{\theta(\text{free})}\}$  is at least as expressive on  $S$  in the sense that

$$\inf_{\theta(\text{free})} \text{TV}(\hat{\nu}_{\theta(\text{free})}, \nu) \leq \text{TV}(\hat{\nu}_{\theta(\text{eq})}, \nu), \quad (30)$$

then there exists  $\theta(\text{free})$  such that the Haar-randomized model

$$\hat{\mu}_{\theta(\text{free})} := \int_S \left( \int_{\mathcal{G}} \delta_{g \cdot \mathbf{z}} d\lambda(g) \right) d\hat{\nu}_{\theta(\text{free})}(\mathbf{Z}) \quad (31)$$

satisfies  $\text{TV}(\hat{\mu}_{\theta(\text{free})}, \mu) < \varepsilon$ .

*Proof.* First, pushforward cannot increase TV (data processing inequality). Hence

$$\text{TV}(\hat{\nu}_{\theta(\text{eq})}, \nu) = \text{TV}(\Psi_{\#}\hat{\mu}_{\theta(\text{eq})}, \Psi_{\#}\mu) \leq \text{TV}(\hat{\mu}_{\theta(\text{eq})}, \mu) < \varepsilon. \quad (32)$$

By (30) (note that this usually holds according to our discussions above), choose  $\theta(\text{free})$  such that  $\text{TV}(\hat{\nu}_{\theta(\text{free})}, \nu) \leq \text{TV}(\hat{\nu}_{\theta(\text{eq})}, \nu) < \varepsilon$ . Finally, Haar randomization is a Markov kernel  $T$  (cf. Theorem 3.1), and thus is TV-nonexpansive:

$$\text{TV}(T(\hat{\nu}_{\theta(\text{free})}), T(\nu)) \leq \text{TV}(\hat{\nu}_{\theta(\text{free})}, \nu). \quad (33)$$

Since  $T(\nu) = \mu$  (Theorem 3.1), we conclude  $\text{TV}(\hat{\mu}_{\theta(\text{free})}, \mu) < \varepsilon$ .  $\square$

## C.2. Proof for Section 3.2

In this subsection, we provide more detailed analysis and complete proof for our results regarding accelerating diffusion and flow model training through canonicalization.

### C.2.1. PROOF FOR SECTION 3.2.1

This part mainly considers score-based diffusion modeling, showing that the existence of group symmetry in invariant data induces the mixture structure of the score function. A toy visualization of the corresponding intermediate-time score field is provided in Figure 4, which illustrates the mixture-induced ambiguity near the symmetry axis and the simplification brought by canonicalization.

**Proposition C.7** (Score of a symmetry mixture). *Assume  $q$  is differentiable and positive where needed. Then*

$$\nabla \log p(x) = \sum_{m=1}^M w_m(x) g_m \cdot \nabla \log q(g_m^{-1} \cdot x), \quad w_m(x) := \frac{q(g_m^{-1} \cdot x)}{\sum_{j=1}^M q(g_j^{-1} \cdot x)}. \quad (34)$$

*Proof.* Differentiate  $\log p(x) = \log\left(\frac{1}{M} \sum_m q(g_m^{-1} \cdot x)\right)$  and apply the chain rule. Orthogonality implies  $\nabla_x q(g_m^{-1} \cdot x) = g_m \cdot \nabla q(g_m^{-1} \cdot x)$ .  $\square$

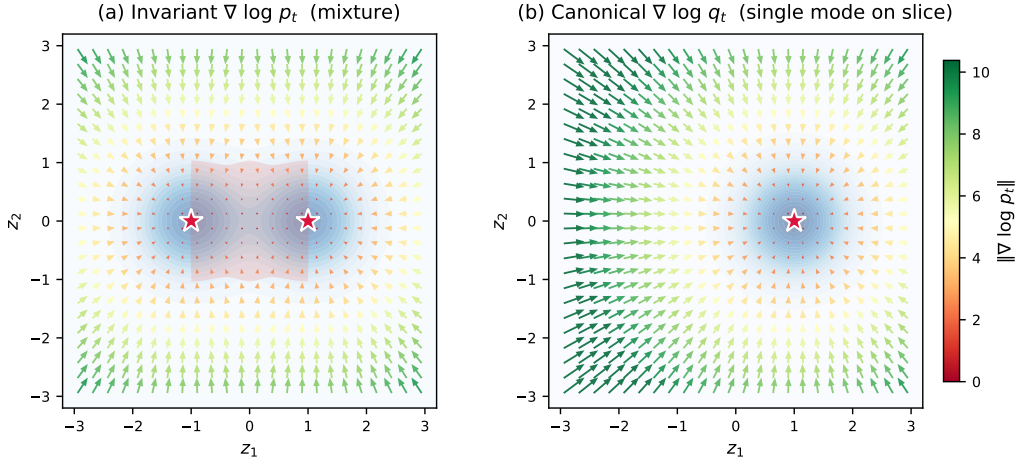
Score fields at  $t=0.6$  for data with  $\mathbb{Z}_2$  reflection symmetry


Figure 4. Toy illustration of score fields at intermediate diffusion time  $t = 0.6$  for a  $\mathbb{Z}_2$ -symmetric Gaussian mixture in 2D. (a) The invariant score  $\nabla \log p_t$  is a responsibility-weighted mixture of per-mode scores (Eq. 9). Near the symmetry axis, the two modes' contributions cancel, creating a low-magnitude "dead zone" (red shading) where the score provides little learning signal and conflicting gradient directions. (b) On the canonical slice, only a single mode remains. The score  $\nabla \log q_t$  is smooth, unambiguous, and consistently points toward the mode center, eliminating the mixture complexity entirely. Arrow length is proportional to score magnitude; color encodes  $\|\nabla \log p_t\|$  (green = high, red = low).

### C.2.2. PROOF FOR SECTION 3.2.2

This part highlights our central result: canonical training reduces conditional variance in flow-matching.

**Irreducible error in flow matching.** First we prove a standard lemma, stating that the irreducible flow-matching error is the conditional variance.

**Lemma C.8** (Irreducible flow-matching error = conditional variance). *Let  $\hat{v}$  be any measurable predictor of  $v_t(\mathbf{Z}_t)$  from  $(t, \mathbf{Z}_t)$ . Then the minimum achievable MSE satisfies*

$$\inf_{\hat{v}} \mathbb{E}[\|\hat{v}(t, \mathbf{Z}_t) - (\mathbf{Z}_1 - \mathbf{Z}_0)\|^2] = \mathbb{E}[\text{Var}(\mathbf{Z}_1 - \mathbf{Z}_0 \mid t, \mathbf{Z}_t)]. \quad (35)$$

*Proof.* This is the standard  $L^2$  regression identity: the conditional expectation  $\mathbb{E}[\mathbf{Z}_1 - \mathbf{Z}_0 \mid t, \mathbf{Z}_t]$  is the unique minimizer, and the minimum risk equals the expected conditional variance (law of total variance).  $\square$

**Corollary C.9** (Deterministic couplings yield zero irreducible error). *If  $\mathbf{Z}_1 - \mathbf{Z}_0$  is (a.s.) a deterministic function of  $(t, \mathbf{Z}_t)$ , then the right-hand side is 0, and in principle the flow can be learned with arbitrarily small error (up to approximation and optimization).*

Below we give a concrete example of this irreducible error, showing that the minimum achievable MSE loss is strictly positive in the presence of multi-group ambiguity.

**Quantifying symmetry ambiguity through posterior collision bounds** We now strengthen the variance-decomposition viewpoint by (i) giving an exact identity for the *symmetry-ambiguity* term in Theorem 3.5, (ii) deriving lower bounds purely in terms of the posterior  $\varrho(g \mid Z_t)$  over the latent symmetry.

Under the group-aligned lifted coupling, write the posterior over the latent symmetry as  $\varrho(g \mid z) := \mathbb{P}(G = g \mid Z_t = z)$  (finite group case; for compact groups interpret  $\varrho(\cdot \mid z)$  as a posterior density w.r.t. Haar).

We first show that ambiguity term can be written as a pairwise mean-separation identity. The following identity is standard (*variance via i.i.d. copies*), but its specialization here yields an explicit geometric handle on the symmetry-ambiguity term.

**Lemma C.10** (Symmetry ambiguity as pairwise conditional-mean separation). *Define the group-conditioned conditional mean drift*

$$m_g(z) := \mathbb{E}[U \mid Z_t = z, G = g], \quad m(z) := \mathbb{E}[U \mid Z_t = z]. \quad (36)$$

Then the ambiguity term in Theorem 3.5 satisfies, for each  $z$ ,

$$\text{Var}(\mathbb{E}[U \mid Z_t, G] \mid Z_t = z) = \mathbb{E}_{g \sim \varrho(\cdot|z)} [\|m_g(z) - m(z)\|^2] \quad (37)$$

$$= \frac{1}{2} \mathbb{E}_{g, g' \stackrel{i.i.d.}{\sim} \varrho(\cdot|z)} [\|m_g(z) - m_{g'}(z)\|^2]. \quad (38)$$

*Proof.* Eq. (37) is the definition of conditional variance of a discrete random variable taking values  $m_G(z)$ . Eq. (38) follows from the standard identity  $\text{Var}(X) = \frac{1}{2} \mathbb{E}\|X - X'\|^2$  for i.i.d. copies  $(X, X')$  (applied conditionally on  $Z_t = z$ ).  $\square$

We hereby provide the lower bounds from posterior uncertainty (collision probability). Lemma C.10 implies that ambiguity is large whenever the posterior  $\varrho(g \mid Z_t)$  spreads mass over group elements with meaningfully different conditional mean drifts.

**Proposition C.11** (Posterior-collision lower bound). *Fix  $z$  and let  $\mathcal{G}_\Delta(z) \subseteq \mathcal{G}$  be any subset such that for all  $g \in \mathcal{G}_\Delta(z)$ ,*

$$\|m_g(z) - m(z)\| \geq \Delta(z). \quad (39)$$

Then

$$\text{Var}(\mathbb{E}[U \mid Z_t, G] \mid Z_t = z) \geq \Delta(z)^2 \mathbb{P}(G \in \mathcal{G}_\Delta(z) \mid Z_t = z). \quad (40)$$

In particular, if there exist two disjoint subsets  $\mathcal{A}(z), \mathcal{B}(z)$  such that  $\|m_g(z) - m_{g'}(z)\| \geq \Delta(z)$  for all  $g \in \mathcal{A}(z), g' \in \mathcal{B}(z), g \neq g'$ , then

$$\text{Var}(\mathbb{E}[U \mid Z_t, G] \mid Z_t = z) \geq \frac{\Delta(z)^2}{2} \left(1 - \sum_g \varrho(g \mid z)^2\right), \quad (41)$$

where  $1 - \sum_g \varrho(g \mid z)^2$  is the complement of the posterior collision probability.

*Proof.* The first bound is immediate from (37) by keeping only  $g \in \mathcal{G}_\Delta(z)$  and using  $\|m_g(z) - m(z)\|^2 \geq \Delta(z)^2$  there. For the second bound, use (38) and lower bound the pairwise distance by  $\Delta(z)$  whenever  $(g, g') \in \mathcal{A}(z) \times \mathcal{B}(z) \cup \mathcal{B}(z) \times \mathcal{A}(z)$ . Then

$$\mathbb{P}(g \neq g' \mid z) = 1 - \sum_g \varrho(g \mid z)^2, \quad (42)$$

and the stated inequality follows (up to the factor 1/2 from (38)).  $\square$

The collision term  $1 - \sum_g \varrho(g \mid Z_t)^2$  is 0 iff  $G$  is a.s. determined by  $Z_t$ . Thus, unless the conditional mean drifts  $m_g(Z_t)$  coincide across  $g$ , symmetry ambiguity creates a strictly positive variance component whenever  $G$  remains uncertain.

**Discussions on lifted coupling.** A subtle but important point is that the group-aligned lift in (12) shares the same latent  $G$  across  $(Z_0, Z_1)$ , which *a priori* could introduce dependence. However, for *standard isotropic Gaussian* priors and orthogonal group actions (including  $S_N$  permutations and  $SO(3)$  rotations), this dependence disappears.

**Proposition C.12** (Lift of a product slice coupling equals the ambient product coupling). *Assume  $\mathcal{G}$  acts orthogonally on  $\mathbb{R}^d$ . Let  $G \sim \lambda$  and  $(\tilde{Z}_0, \tilde{Z}_1) \sim q_0 \otimes \mathcal{N}(0, I)$  be independent, and define  $(Z_0, Z_1) = (G \cdot \tilde{Z}_0, G \cdot \tilde{Z}_1)$ . Then (i)  $Z_1 \sim \mathcal{N}(0, I)$  and  $Z_1 \perp G$ , and (ii)  $Z_1 \perp Z_0$ . Moreover, if the invariant data distribution satisfies the disintegration  $p_0 = \int (g \cdot)_{\#} q_0 d\lambda(g)$ , then the lifted ambient coupling satisfies*

$$(Z_0, Z_1) \sim p_0 \otimes \mathcal{N}(0, I). \quad (43)$$

*Proof.* For any measurable  $B$  and any  $g$ , orthogonal invariance of  $\mathcal{N}(0, I)$  gives  $\mathbb{P}(Z_1 \in B \mid G = g) = \mathbb{P}(g \cdot \tilde{Z}_1 \in B) = \mathbb{P}(\tilde{Z}_1 \in B)$ , which does not depend on  $g$ , hence  $Z_1 \perp G$  and  $Z_1 \sim \mathcal{N}(0, I)$ . For any measurable  $A, B$ ,

$$\mathbb{P}(Z_0 \in A, Z_1 \in B) = \mathbb{E} \left[ \mathbf{1}_{\{Z_0 \in A\}} \mathbb{P}(Z_1 \in B \mid G, \tilde{Z}_0) \right] = \mathbb{P}(\tilde{Z}_1 \in B) \mathbb{P}(Z_0 \in A), \quad (44)$$

since  $\tilde{Z}_1$  is independent of  $(G, \tilde{Z}_0)$  and  $\mathbb{P}(g \cdot \tilde{Z}_1 \in B)$  is constant in  $g$ . Thus  $Z_1 \perp Z_0$ . If additionally  $Z_0 \sim p_0$  via the stated disintegration, the joint law is  $p_0 \otimes \mathcal{N}(0, I)$ .  $\square$

Therefore, the commonly used ambient product coupling  $\gamma_{\text{mix}} = p_0 \otimes \mathcal{N}(0, I)$  can be viewed as a special case of slice training with product coupling and isotropic Gaussian slice prior, thus *conditional variances of two paradigms become directly comparable*. Canonicalization does *not* change the coupling in this case; it changes the *representation* so that the group ambiguity term in Theorem 3.5 becomes explicit and avoidable on the slice. Hence our theory covers most practical training regimes without optimal transport.

However, this generally does not hold if  $\tilde{Z}_1$  is group-aware (e.g., learned or group-statistics related prior in our implementation), or  $\tilde{Z}_1$  and  $\tilde{Z}_0$  are not independent (e.g., optimal transport on the canonicalized slice). In particular, if the slice prior is *not* isotropic (e.g., an aligned Gaussian  $\mathcal{N}(\mu, \Sigma)$  with  $g\Sigma g^\top \neq \Sigma$  for some  $g$ ), then  $g \cdot \tilde{Z}_1$  depends on  $g$  in distribution, so  $Z_1$  is no longer independent of  $G$ , and the lifted coupling generally is *not* a product. Likewise, if the slice coupling  $\tilde{\gamma}$  is non-product (e.g., OT/Monge), the lifted coupling remains non-product.

**Canonicalization reduces conditional variance via symmetry ambiguity elimination.** We now prove our central result on the advantages of canonicalized flow matching.

**Theorem C.13** (Variance decomposition under group-aligned lift; Theorem 3.5 in main text). *Assume  $\mathcal{G}$  acts orthogonally. Under the group-aligned lifted coupling  $G \sim \lambda$  independent of  $(\tilde{Z}_0, \tilde{Z}_1)$  and  $(Z_0, Z_1) = (G \cdot \tilde{Z}_0, G \cdot \tilde{Z}_1)$ , let  $Z_t = (1-t)Z_0 + tZ_1$  and  $\tilde{Z}_t = (1-t)\tilde{Z}_0 + t\tilde{Z}_1$ . With  $U := Z_1 - Z_0 = G \cdot \Delta$  and  $\Delta := \tilde{Z}_1 - \tilde{Z}_0$ , we have*

$$\text{Var}(U \mid Z_t) = \underbrace{\mathbb{E}[\text{Var}(\Delta \mid \tilde{Z}_t) \mid Z_t]}_{\text{within-slice difficulty}} + \underbrace{\text{Var}(\mathbb{E}[U \mid Z_t, G] \mid Z_t)}_{\text{symmetry ambiguity} \geq 0}. \quad (45)$$

Consequently,

$$\mathbb{E}[\text{Var}(U \mid Z_t)] \geq \mathbb{E}[\text{Var}(\Delta \mid \tilde{Z}_t)]. \quad (46)$$

*Proof.* Apply the standard law of total variance with respect to the latent symmetry variable  $G$ :

$$\text{Var}(U \mid Z_t) = \mathbb{E}[\text{Var}(U \mid Z_t, G) \mid Z_t] + \text{Var}(\mathbb{E}[U \mid Z_t, G] \mid Z_t). \quad (47)$$

The second term is nonnegative. For the first term, conditioning on  $(Z_t, G)$  determines  $\tilde{Z}_t = G^{-1} \cdot Z_t$ . Under the lifted coupling,  $(\Delta, \tilde{Z}_t)$  is independent of  $G$ , and orthogonality implies  $\text{Var}(G \cdot X \mid \cdot) = \text{Var}(X \mid \cdot)$ . Therefore,

$$\text{Var}(U \mid Z_t, G) = \text{Var}(G \cdot \Delta \mid \tilde{Z}_t, G) = \text{Var}(\Delta \mid \tilde{Z}_t). \quad (48)$$

Substituting into (47) yields (13). Taking expectations over  $Z_t$  and using the tower property gives (14).  $\square$

The lower bound in Theorem 3.5 is the “within-slice” conditional variance; the gap is exactly the additional uncertainty induced by not knowing which symmetry element generated the observation. This formalizes the intuition that symmetry induces posterior multi-modality over group elements, which inflates the conditional variance and creates an irreducible error floor for coarse solvers.

*Remark C.14* (What this does *and does not* compare). The lemma above compares “with vs. without marginalizing  $G$ ” under the same lifted coupling. It does *not* imply that a canonicalized training paradigm universally dominates any other paradigm under arbitrary choices of couplings. Its value is to isolate a nonnegative variance component that arises purely from symmetry ambiguity and is absent on the slice. However, as discussed in Proposition C.12, the lifted coupling can easily recover the standard product with Gaussian noises used in most practical flow matching models (w/o OT).

In the next part we will further explain Remark 3.6. Fix a coupling  $\gamma$  on  $(Z_0, Z_1)$  and define the linear interpolation

$$Z_t = (1-t)Z_0 + tZ_1, \quad U := Z_1 - Z_0. \quad (49)$$

For any predictor  $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , the population squared loss at time  $t$  is

$$\mathcal{L}_t(v) := \mathbb{E}[\|v(Z_t) - U\|^2]. \quad (50)$$

It is standard that the Bayes regressor is  $v^*(z) = \mathbb{E}[U \mid Z_t = z]$  and the Bayes risk equals the conditional variance:

$$\inf_v \mathcal{L}_t(v) = \mathbb{E}[\text{Var}(U \mid Z_t)]. \quad (51)$$

**Canonicalization does not improve the ambient Bayes risk of equivariant models.** Let a compact group  $\mathcal{G}$  act orthogonally on  $\mathbb{R}^d$  (e.g. permutations and rotations). Let  $\gamma$  be *diagonal-invariant*:

$$(g, g)_{\#} \gamma = \gamma \quad \forall g \in \mathcal{G}, \quad (52)$$

which holds in particular for the ambient product coupling  $p_0 \otimes \mathcal{N}(0, I)$  when  $p_0$  is  $\mathcal{G}$ -invariant, and also for symmetrized OT couplings.

**Proposition C.15** (Equivariance of the Bayes regressor under diagonal invariance). *Assume  $\mathcal{G}$  acts orthogonally and  $\gamma$  is diagonal-invariant. Then the Bayes regressor  $v^*(z) = \mathbb{E}[U \mid Z_t = z]$  is  $\mathcal{G}$ -equivariant:*

$$v^*(g \cdot z) = g \cdot v^*(z) \quad \forall g \in \mathcal{G}. \quad (53)$$

Consequently, restricting to the equivariant function class  $\mathcal{H}_{\text{eq}} := \{v : v(g \cdot z) = g \cdot v(z)\}$  does not increase the minimal population loss:

$$\inf_{v \in \mathcal{H}_{\text{eq}}} \mathcal{L}_t(v) = \inf_v \mathcal{L}_t(v) = \mathbb{E}[\text{Var}(U \mid Z_t)]. \quad (54)$$

*Proof.* Diagonal invariance implies  $(Z_t, U) \stackrel{d}{=} (g \cdot Z_t, g \cdot U)$ . Hence for any measurable set  $A$ ,

$$\mathbb{E}[U \mathbf{1}_{\{Z_t \in A\}}] = \mathbb{E}[g \cdot U \mathbf{1}_{\{g \cdot Z_t \in A\}}] = g \cdot \mathbb{E}[U \mathbf{1}_{\{Z_t \in g^{-1}A\}}], \quad (55)$$

which yields  $v^*(g \cdot z) = g \cdot v^*(z)$  by Radon–Nikodym characterization of conditional expectation. The equality of infima follows because  $v^* \in \mathcal{H}_{\text{eq}}$ .  $\square$

*Remark C.16.* If one keeps the *same ambient regression problem* (51), then canonicalizing the inputs does not improve the population optimum for an equivariant model class, because the Bayes regressor is already equivariant (Proposition C.15).

**Canonical conditions expand the effective hypothesis class of non-equivariant models.** Canonicalization is most useful when the model is *not* architecturally equivariant. Intuitively, a non-equivariant model in ambient space must learn symmetry averaging from data. Canonicalization and canonical conditions provide a *gauge* (e.g. a canonicalizer  $\Psi$  and an associated group element  $\kappa(\cdot)$ ) that allows a non-equivariant model to implement equivariant behavior via an explicit formula.

We assume an exact canonicalizer model, i.e., a canonicalizer provides  $\Psi : \mathbb{R}^d \rightarrow S$  (slice) and  $\kappa : \mathbb{R}^d \rightarrow \mathcal{G}$  such that  $\Psi(x) = \kappa(x)^{-1} \cdot x$  on a full-measure set (up to stabilizers).

**Proposition C.17** (Canonical-condition lifting realizes equivariant functions). *Let  $\mathcal{A}$  be any function class on the slice  $S$  (not necessarily equivariant). Define the induced ambient class using the canonical condition:*

$$\mathcal{H}_{\text{lift}}(\mathcal{A}) := \{h(x) = \kappa(x) \cdot a(\Psi(x)) : a \in \mathcal{A}\}. \quad (56)$$

Then every  $h \in \mathcal{H}_{\text{lift}}(\mathcal{A})$  is  $\mathcal{G}$ -equivariant (on the set where  $\Psi, \kappa$  are consistent). Moreover, if  $\mathcal{A}$  is universal on  $S$ , then  $\mathcal{H}_{\text{lift}}(\mathcal{A})$  can approximate any equivariant target function on  $\mathbb{R}^d$  up to the usual stabilizer caveats.

*Proof.* Equivariance follows from  $\Psi(g \cdot x) = \Psi(x)$  and  $\kappa(g \cdot x) = g\kappa(x)$  (in the free-action regime):

$$h(g \cdot x) = \kappa(g \cdot x) \cdot a(\Psi(g \cdot x)) = g\kappa(x) \cdot a(\Psi(x)) = g \cdot h(x). \quad (57)$$

Universality is inherited because any equivariant function is determined by its restriction to the slice.  $\square$

*Remark C.18* (Why this is a real “shortcut” for non-equivariant networks). A generic non-equivariant network that only sees  $x$  must *implicitly* learn the equivariant structure. Providing  $(\Psi(x), \kappa(x))$  (or a discrete proxy such as canonical rank/frame ID) allows implementing equivariance by construction as in Proposition C.17. This can reduce approximation complexity and improve optimization.

To summarize, slice training does *not* claim to reduce the ambient Bayes risk (51). The empirical Bayes risk that can be actually achieved still depends on the architectures and could not be improved for equivariant models (Proposition C.15). Instead, it learns a *different* regression problem on the slice (predicting  $\Delta$  from  $\tilde{Z}_t$ ), whose irreducible variance is the first term in Theorem 3.5 and is provably no larger than the ambient one. This is where canonicalization can enable easier training and fewer steps sampling for non-equivariant models.

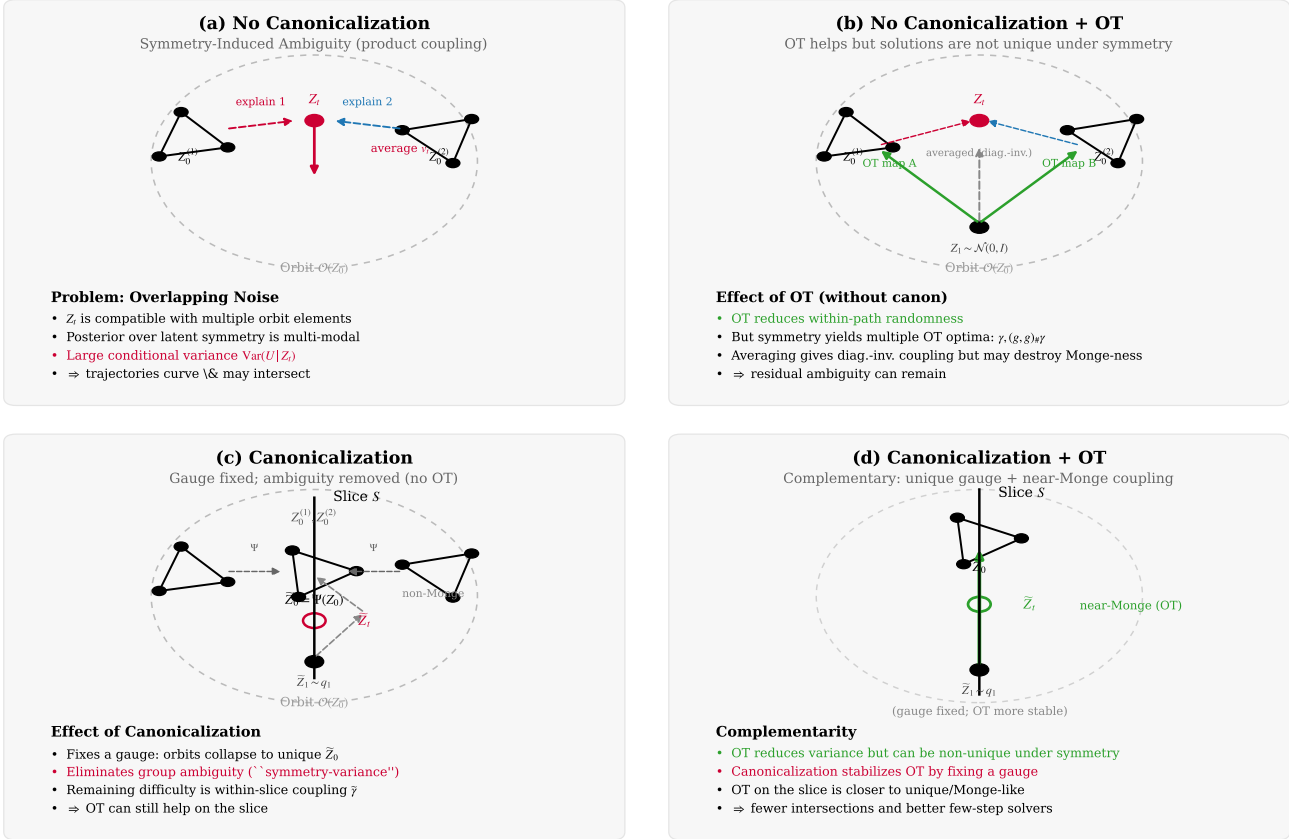


Figure 5. OT and canonicalization act complementarily. OT can reduce conditional variances with or without canonicalization, but OT solutions are generally non-unique in the presence of symmetry (if  $\gamma$  is optimal then  $(g, g) \neq \gamma$  is also optimal); averaging yields a diagonal-invariant optimum but may destroy Monge-ness. Canonicalization fixes a gauge (collapsing each orbit to a unique slice representative), eliminating group ambiguity and stabilizing OT on the slice.

### C.2.3. PROOF FOR SECTION 3.2.3

This part provides supplementary to aligned slice priors, OT couplings, and train–test consistency.

Canonicalization eliminates the symmetry-ambiguity term, but the remaining within-slice conditional variance  $\mathbb{E}[\text{Var}(\Delta | \tilde{Z}_t)]$  depends on the *slice prior*  $q_1$  and the *slice coupling*  $\tilde{\gamma}$ . Section 3.2.3 explains that after canonicalization, the remaining irreducible term is the *within-slice difficulty*  $\mathbb{E}[\text{Var}(\Delta | \tilde{Z}_t)]$ , which depends on (i) the slice prior  $q_1$  and (ii) the slice coupling  $\tilde{\gamma}$  (Theorem 3.5). This appendix subsection collects three short technical supplements: (i) a closed-form expression showing how *misaligned* Gaussian priors can inflate within-slice variance under product coupling, (ii) why OT/near-Monge couplings are complementary but can be non-unique under symmetry, and (iii) a minimal consistency statement clarifying when training-time OT is compatible with inference and when conditioning causes mismatch.

**Gaussian misalignment under product coupling: an exact within-slice variance formula.** Canonicalization removes the symmetry-ambiguity term, but even on the slice the conditional variance can remain large if the prior is poorly aligned with the slice geometry. The simplest (and common) baseline is the product coupling  $\tilde{\gamma} = q_0 \otimes q_1$ . The following proposition gives an exact expression for the within-slice conditional covariance in the linear-path (rectified-flow) setting, illustrating how mismatch between  $\Sigma_0$  and  $\Sigma_1$  controls the irreducible error.

**Proposition C.19** (Closed-form within-slice conditional variance for independent Gaussians). *Assume the slice data and slice prior are Gaussian and independent:*

$$\tilde{Z}_0 \sim \mathcal{N}(\mu_0, \Sigma_0), \quad \tilde{Z}_1 \sim \mathcal{N}(\mu_1, \Sigma_1), \quad \tilde{Z}_0 \perp \tilde{Z}_1.$$

Let  $\tilde{Z}_t = (1-t)\tilde{Z}_0 + t\tilde{Z}_1$  and  $\Delta = \tilde{Z}_1 - \tilde{Z}_0$ . Then

$$\text{Cov}(\Delta \mid \tilde{Z}_t) = \frac{1}{t^2} \text{Cov}(\tilde{Z}_0 \mid \tilde{Z}_t), \quad (58)$$

and the conditional covariance admits the closed form

$$\text{Cov}(\tilde{Z}_0 \mid \tilde{Z}_t) = \Sigma_0 - (1-t)^2 \Sigma_0 \left( (1-t)^2 \Sigma_0 + t^2 \Sigma_1 \right)^{-1} \Sigma_0. \quad (59)$$

Consequently, the within-slice irreducible error in flow matching is

$$\mathbb{E}[\text{Var}(\Delta \mid \tilde{Z}_t)] = \text{tr} \left( \frac{1}{t^2} \Sigma_0 - \frac{(1-t)^2}{t^2} \Sigma_0 \left( (1-t)^2 \Sigma_0 + t^2 \Sigma_1 \right)^{-1} \Sigma_0 \right). \quad (60)$$

*Proof.* The identity  $\Delta = (\tilde{Z}_1 - \tilde{Z}_0)/t$  is algebraic, hence  $\text{Cov}(\Delta \mid \tilde{Z}_t) = \frac{1}{t^2} \text{Cov}(\tilde{Z}_0 \mid \tilde{Z}_t)$ . Since  $(\tilde{Z}_0, \tilde{Z}_t)$  is jointly Gaussian with  $\text{Cov}(\tilde{Z}_0, \tilde{Z}_t) = (1-t)\Sigma_0$  and  $\text{Cov}(\tilde{Z}_t) = (1-t)^2 \Sigma_0 + t^2 \Sigma_1$ , the standard Gaussian conditioning formula yields (59). Taking traces gives the final expression.  $\square$

Therefore, if one uses an isotropic prior  $q_1 = \mathcal{N}(0, I)$  while  $q_0$  is strongly anisotropic (large condition number in  $\Sigma_0$ ), then the trace above remains large along directions where  $\Sigma_0$  dominates  $I$ , formalizing why a *misaligned* simple prior can degrade few-step accuracy even after canonicalization. A practical remedy is to choose an *aligned* prior within a tractable family, e.g. the moment-matched Gaussian  $q_1^* = \mathcal{N}(\mathbb{E}_{q_0}[\tilde{Z}_0], \text{Cov}_{q_0}(\tilde{Z}_0))$  (KL projection onto Gaussians), or a learned prior. While it is possible to learn the canonical prior  $q_1$ , using KL projections of  $q_0$  is a simple yet effective way. The computation only occurs in pre-processing and does not induce any overhead in training. For instance, among Gaussians, the moment-matched Gaussian is the KL projection of  $q_0$ :

**Proposition C.20** (Moment-matched Gaussian is the KL-optimal Gaussian approximation). *Let  $q_0$  be any distribution on  $\mathbb{R}^d$  with finite second moments. Consider the family  $\mathcal{N}(\mu, \Sigma)$  with  $\Sigma \succ 0$ . The minimizer of  $\text{KL}(q_0 \parallel \mathcal{N}(\mu, \Sigma))$  is*

$$\mu^* = \mathbb{E}_{q_0}[\tilde{Z}_0], \quad \Sigma^* = \text{Cov}_{q_0}(\tilde{Z}_0). \quad (61)$$

*Proof.* This is a standard exponential-family projection: the Gaussian log-density is an affine function of the sufficient statistics  $(x, xx^\top)$ , so minimizing  $\text{KL}(q_0 \parallel \mathcal{N}(\mu, \Sigma))$  yields moment matching.  $\square$

**OT/near-Monge couplings: complementary to canonicalization, but symmetry can destroy uniqueness.** Given  $(q_0, q_1)$ , choosing  $\tilde{\gamma}$  close to an OT coupling can make transport more deterministic (Monge-like), shrinking  $\mathbb{E}[\text{Var}(\Delta \mid \tilde{Z}_t)]$  and straightening trajectories (hence benefiting few-step solvers). However, when working in ambient space with symmetric marginals and symmetric costs, OT solutions are generally *not unique*: if  $\gamma$  is optimal then  $(g, g)_{\#} \gamma$  is also optimal. Averaging over  $g$  yields a diagonal-invariant optimum, but can destroy Monge-ness by turning a deterministic map into a mixture. Canonicalization fixes a gauge and thus *stabilizes* OT: on the slice, the symmetry-induced degeneracy is reduced, and OT solutions are empirically closer to unique/Monge-like. This is the precise sense in which canonicalization and OT act *complementarily*: canonicalization removes group ambiguity (Theorem 3.5), while OT targets the remaining within-slice difficulty. Figure 5 provides an illustration.

**Training-time OT is compatible with inference; mismatch comes from conditioning variables.** A frequent confusion is whether using OT in training forces OT at inference. In flow matching / rectified flow, OT changes the *training coupling*  $\tilde{\gamma}$  (and thus the supervision signal), but inference only requires sampling from the *marginal prior*  $q_1$  and integrating the learned dynamics. In particular, if the population vector field is well-defined and sufficiently regular, sampling is *coupling-free* given the correct field: one draws  $\tilde{Z}_1 \sim q_1$  and integrates the learned ODE/SDE to obtain a sample from  $q_0$ ; no paired endpoints are needed at inference.

The actual train–test mismatch risk arises when the model is conditioned on an auxiliary variable  $C$  (e.g. “canonical rank” PE, frame ID, or any deterministic/stochastic output of a canonicalizer). Let  $\tilde{\pi}_1^{\text{tr}}(C, \tilde{Z}_1)$  be the training-time joint law at the start, and  $\tilde{\pi}_1^{\text{inf}}(C, \tilde{Z}_1)$  the inference-time one. A sufficient condition for *no mismatch* is:

$$\tilde{\pi}_1^{\text{tr}}(C, \tilde{Z}_1) = \tilde{\pi}_1^{\text{inf}}(C, \tilde{Z}_1). \quad (62)$$

In particular, if  $C$  is a deterministic index inherent to the canonical coordinate system (e.g. the coordinate index  $i$  is the canonical rank), then this equality holds for any choice of  $q_1$ ; conversely, if inference computes  $C$  as a function of sampled noise in a way that differs from training, mismatch is unavoidable. Finally, note that *canonicalizing noise at inference* (e.g. setting  $\tilde{Z}_1 := \Psi(\varepsilon)$  with  $\varepsilon \sim \mathcal{N}(0, I)$ ) generally induces a complicated slice prior  $\Psi_{\#}\mathcal{N}(0, I)$ ; unless training uses the same induced prior, this introduces an avoidable prior mismatch and can increase within-slice variance.

## D. Implementation Details

### D.1. Details of Canonicalization Methods

#### D.1.1. GEOMETRIC SPECTRAL ORDERING

We define a deterministic canonicalization function  $\kappa : \mathcal{M} \rightarrow \mathcal{M}$  that maps each molecular graph to a unique representative within its orbit under  $S_N \times SE(3)$ . The key insight is to leverage the geometric Laplacian constructed from 3D atomic coordinates, whose Fiedler vector encodes the molecular connectivity structure in a continuous, rotation-invariant manner.

**Geometric Laplacian construction and canonical permutation.** The construction of the geometric Laplacian and the calculation of the (signed) Fiedler vector can be cleanly summarized as the pseudo-algorithm in Algorithm 1.

---

#### Algorithm 1 Geometric Laplacian and signed Fiedler vector

---

**Input:** coordinates  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^{N \times 3}$ , bond set  $E$ .  
 Compute mean bond length  $\bar{d}_{\text{bond}} \leftarrow \frac{1}{|E|} \sum_{(i,j) \in E} \|\mathbf{x}_i - \mathbf{x}_j\|$ .  
 Set bandwidth  $\sigma^2 \leftarrow 4\bar{d}_{\text{bond}}^2$ .  
**for**  $i = 1, \dots, N$  **do**  
   **for**  $j = 1, \dots, N$  **do**  
     **if**  $i = j$  **then**  
        $W_{ij} \leftarrow 0$ .  
     **else**  
        $W_{ij} \leftarrow \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (2\sigma^2))$ .  
     **end if**  
   **end for**  
**end for**  
 Degree matrix  $\mathbf{D} \leftarrow \text{diag}(\sum_j W_{ij})$ .  
 Random-walk Laplacian  $L_{\text{rw}} \leftarrow \mathbf{D}^{-1}(\mathbf{D} - \mathbf{W})$ .  
 Compute eigenpairs of  $L_{\text{rw}}$ :  $0 = \lambda_1 \leq \lambda_2 \leq \dots$  with eigenvectors  $\{u_k\}$ .  
 Fiedler vector  $u_2 \leftarrow$  eigenvector of  $\lambda_2$ .  
 Centroid  $\bar{\mathbf{x}} \leftarrow \frac{1}{N} \sum_i \mathbf{x}_i$  and mean radius  $\bar{d} \leftarrow \frac{1}{N} \sum_i \|\mathbf{x}_i - \bar{\mathbf{x}}\|$ .  
 Fix sign by centroid-direction convention:  
    $u_2 \leftarrow \text{sign}(\sum_i u_{2,i}(\|\mathbf{x}_i - \bar{\mathbf{x}}\| - \bar{d})) \cdot u_2$ .  
**Output:** signed Fiedler vector  $u_2$ .

---

The canonical ordering is then:  $\pi^* = \text{argsort}(u_2)$  and the canonicalized molecule is  $\kappa(\mathcal{M}) = \pi^*(\mathcal{M})$ . The Fiedler vector approximately captures the graph Laplacian’s fundamental mode of oscillation, providing a “core-to-periphery” ordering. Atoms with similar Fiedler values tend to be spatially close and structurally equivalent, which induces a natural generation order from molecular core to functional groups; see Figure 6 for illustrative examples.

**Canonical  $SO(3)$  frame.** The above procedure determines a canonical order on  $S_N$ , based on which we further (optionally) process the  $SO(3)$  symmetry using Algorithm 2.

#### D.1.2. ALTERNATIVE ORDERINGS

Graph canonicalization is a widely studied topic and there are also other existing canonicalization methods. However, most of them are defined for permutation in abstract graphs (Zhao et al., 2024; Ma et al., 2023; Dong et al., 2024) instead of geometric graphs.

**Algorithm 2** Spectral  $SO(3)$  Canonicalization

**Require:** Atomic coordinates  $\mathbf{X} \in \mathbb{R}^{N \times 3}$  and signed Fiedler vector  $v$  (computed by Algorithm 1)

**Ensure:** Canonicalized coordinates  $\mathbf{X}'$

- 1: Apply spectral ordering  $\pi^* = \text{argsort}(v)$  and enforce  $\text{sign} \sum_i v_i^3 > 0$
- 2: Let head index  $h$  be rank 0 and tail index  $t$  be rank  $N - 1$
- 3: Choose anchor index

$$a = \text{argmax}_{k \in \{[N/3], [2N/3]\}} \|(\mathbf{x}_k - \mathbf{x}_h) \times (\mathbf{x}_t - \mathbf{x}_h)\|$$

- 4: Compute longitudinal axis  $\mathbf{e}_1 \leftarrow \frac{\mathbf{x}_t - \mathbf{x}_h}{\|\mathbf{x}_t - \mathbf{x}_h\|}$
- 5: Compute plane normal  $\mathbf{n} \leftarrow \mathbf{e}_1 \times (\mathbf{x}_a - \mathbf{x}_h)$
- 6: Normalize normal  $\mathbf{e}_3 \leftarrow \frac{\mathbf{n}}{\|\mathbf{n}\|}$
- 7: Complete right-handed basis  $\mathbf{e}_2 \leftarrow \mathbf{e}_3 \times \mathbf{e}_1$
- 8: Assemble rotation matrix  $\mathbf{R} \leftarrow [\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3]^\top$
- 9: Center and rotate  $\mathbf{X}' \leftarrow (\mathbf{X} - \bar{\mathbf{X}})\mathbf{R}^\top$

We also implement two simpler orderings for ablation:

1. Structural Multihop (inspired by PARD (Zhao et al., 2024)): Iterative degree peeling using weighted multihop degrees  $w_K(v) = \sum_{k=1}^K d_k(v) \cdot N^{K-k}$ , where  $d_k(v)$  counts nodes at exactly  $k$  hops.
2. Atomic Numbering: Priority-based ordering by atomic number (e.g., heavy/rare atoms first, hydrogen last).

However, we experimentally find that our structure-aware geometric spectral ordering outperforms all these methods, validating the effectiveness.

## D.2. Canonical Diffusion and CanonFlow

We present additional algorithmic and architectural design for the canonical diffusion and flow matching.

### D.2.1. DETAILS OF CANONICAL DIFFUSION AND FLOW MATCHING

We already describe the overall framework of canonical diffusion or flow matching in Section 4.2. We now provide additional details in this part.

**Canonical positional encoding.** To explicitly break permutation equivariance, we inject positional information derived from the canonical rank into the model architecture. For each atom  $i$  with normalized rank  $r_i = \text{rank}_i/N \in [0, 1)$ , we compute a sinusoidal positional encoding:

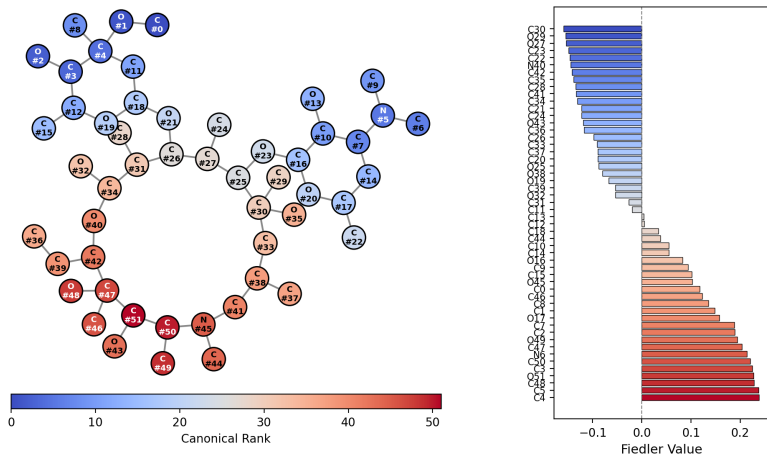
$$\text{PE}(r, 2k) = \sin\left(\frac{r \cdot M}{10000^{2k/d_{\text{pe}}}}\right), \quad \text{PE}(r, 2k+1) = \cos\left(\frac{r \cdot M}{10000^{2k/d_{\text{pe}}}}\right) \quad (63)$$

where  $M = 10000$  is the max position scale and  $d_{\text{pe}}$  is the encoding dimension. This encoding is concatenated with atom features before the initial projection. This completely breaks permutation equivariance at the architecture level, allowing the model to distinguish atoms by their structural role. Note that this PE can be applied to either our novel Canon architecture, or other existing architectures to make a (non-equivariant) canonicity-aware counterpart. The *PE-drop* mechanism replaces  $\mathbf{R}$  with a learned “fake PE” embedding with some probability  $p_{\text{drop}}$  during training, which is used for classifier-free guidance w.r.t. canonical conditions.

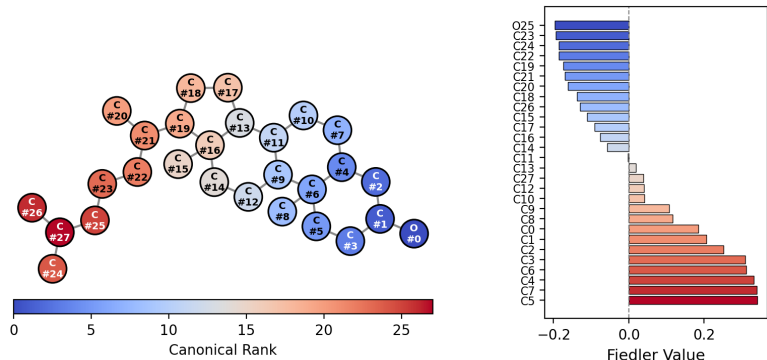
**Adaptive position-dependent prior.** Standard flow matching samples noise from a uniform prior  $p_0$ . We replace this with a position-dependent prior that encodes empirical correlations between canonical position and atom type:

$$p_0^{\text{eff}}(c | r) = \beta_r \cdot p_{\text{prior}}(c) + (1 - \beta_r) \cdot p_r(c | r) \quad (64)$$

where  $r = r_i/N$  is the relative position,  $\beta_r$  is a mixing coefficient (typically 0.1 for categorical),  $p_{\text{prior}}$  is the prior distribution such as  $p_{\text{uniform}}$ , and  $p_r(\cdot | r)$  is the empirical distribution estimated from the training set.



(a) Azithromycin (high symmetry / macrocycle): the spectral ordering is core-centric, concentrating the macrocycle near the spectral center and mapping peripheral sugar moieties toward the sequence terminals.



(b) Cholesterol (low symmetry / anisotropic): the spectral ordering unfolds along the principal axis from the polar head through the rigid tetracyclic core to the flexible hydrocarbon tail, while preserving local connectivity.

Figure 6. Spectral canonicalization via the (signed) Fiedler vector produces stable, locality-preserving atom orderings for both highly symmetric and weakly symmetric molecules (blue: start, red: end), demonstrating robust performance across diverse symmetry regimes.

For *prior estimation*, we discretize  $r \in [0, 1]$  into  $K$  bins and compute:

$$P(c | B_k) = \frac{\text{Count}(c, k) + \epsilon}{\sum_{c'} (\text{Count}(c', k) + \epsilon)} \quad (65)$$

for categorical data; we use a Gaussian with sufficient statistics mean and variance in the dataset for continuous coordinates. At runtime, we use linear interpolation between adjacent bins for continuous  $r$ :

$$p_r(c|r) = (1 - \delta) \cdot P(c|B_{k_0}) + \delta \cdot P(c|B_{k_1}) \quad (66)$$

where  $k_0 = \lfloor r \cdot K \rfloor$ ,  $k_1 = \min(k_0 + 1, K - 1)$ , and  $\delta = r \cdot K - k_0$ .

### D.2.2. CANON ARCHITECTURE

The Semla architecture proposed by (Irwin et al., 2024) maintains two hidden states corresponding to  $\mathbf{X}$ ,  $\mathbf{H}$  in each layer. In our *Canon* architecture, we additionally incorporate a third hidden state to update and refine canonical information in each layer by interacting with other atom and bond features. Figure 7 summarizes the Canon architecture.

**Input processing.** We implement a three-stream molecular dynamics network, **Canon**, which augments a Semla-like equivariant message passing trunk with an explicit *canonical-rank stream* and optional canonical positional encodings (PE).

Concretely, at diffusion/flow time  $t$ , the model takes as input

$$\mathbf{X}_t \in \mathbb{R}^{B \times N \times 3}, \quad \mathbf{H}_t \in \mathbb{R}^{B \times N \times d_h}, \quad \mathbf{R} \in \mathbb{R}^{B \times N},$$

where  $\mathbf{X}_t$  are coordinates,  $\mathbf{H}_t$  are node-wise invariant features (e.g. atom features, optional conditioning), and  $\mathbf{R}$  is the canonical rank (or its normalized variant). The positional encoding processing is identical to the general description above. Optional self-conditioning provides a previous estimate  $\hat{\mathbf{X}}_0$ , and optionally a previous rank estimate  $\hat{\mathbf{R}}$ . In the implementation, self-conditioning is realized by stacking  $(\mathbf{X}_t, \hat{\mathbf{X}}_0)$  and projecting into  $K$  coordinate sets via a linear map  $\Pi_{CS}$ :

$$\mathbf{CS}_t = \Pi_{CS}([\mathbf{X}_t, \hat{\mathbf{X}}_0]) \in \mathbb{R}^{B \times K \times N \times 3}, \quad (67)$$

with masking applied on padded atoms.

We append a size embedding  $\text{Emb}(|V|)$  and optional canonical PE to the invariant features and project them into the trunk width  $d$ :

$$\tilde{\mathbf{H}}_t = \phi_{\text{in}}([\mathbf{H}_t, t, \text{Emb}(|V|), \text{PE}(\mathbf{R})]) \in \mathbb{R}^{B \times N \times d}, \quad (68)$$

where  $\phi_{\text{in}}$  is a two-layer MLP with SiLU nonlinearity.

**CanonDynamics: message-passing layers with node/coord/rank streams.** The trunk (`CanonDynamics`) stacks  $L$  copies of a message passing layer, optionally with edge features. Each layer updates three coupled states:

$$(\mathbf{CS}, \mathbf{H}, \mathbf{R}) \mapsto (\mathbf{CS}', \mathbf{H}', \mathbf{R}'),$$

In each *edge message layer*, we construct pairwise messages using (i) projected node features, (ii) projected rank features, and (iii) geometric features derived from coord-set dot-products. Specifically, after normalization we compute

$$\mathbf{G}_{ij}^{(k)} = \langle \mathbf{CS}_i^{(k)}, \mathbf{CS}_j^{(k)} \rangle, \quad k = 1, \dots, K, \quad (69)$$

and concatenate

$$\mathbf{m}_{ij} = \text{MLP}([\mathbf{W}_h \mathbf{h}_i, \mathbf{W}_h \mathbf{h}_j, \mathbf{W}_r \mathbf{r}_i, \mathbf{W}_r \mathbf{r}_j, \mathbf{G}_{ij}^{(1:K)}, e_{ij}]), \quad (70)$$

where  $e_{ij}$  is an optional edge feature input, and  $\mathbf{W}_h, \mathbf{W}_r$  are learnable projection weights.

Next we adopt attention-based updates for node/coord/rank features. The message tensor is split into three channels (plus optional edge-out):

$$\mathbf{m}_{ij} = (\mathbf{m}_{ij}^{\text{node}}, \mathbf{m}_{ij}^{\text{coord}}, \mathbf{m}_{ij}^{\text{rank}}, \mathbf{m}_{ij}^{\text{edge}}), \quad (71)$$

which drive three attention modules:

$$\mathbf{H} \leftarrow \mathbf{H} + \text{Attn}_{\text{node}}(\mathbf{H}, \mathbf{m}^{\text{node}}), \quad \mathbf{CS} \leftarrow \mathbf{CS} + \text{Attn}_{\text{coord}}(\mathbf{CS}, \mathbf{m}^{\text{coord}}), \quad \mathbf{R} \leftarrow \mathbf{R} + \text{Attn}_{\text{rank}}(\mathbf{R}, \mathbf{m}^{\text{rank}}), \quad (72)$$

interleaved with feed-forward residual blocks consisting of equivariant or invariant MLPs for  $(\mathbf{H}, \mathbf{CS})$  and  $\mathbf{R}$ .

**Output heads.** After  $L$  layers, Canon produces coordinate updates through a normalization + linear head:

$$\hat{\mathbf{X}} = \text{Head}_X(\text{Norm}(\mathbf{CS})) \in \mathbb{R}^{B \times N \times 3}, \quad (73)$$

and predicts atom type and charge logits via linear classifiers on the final node features. If enabled, it also predicts a rank score per node via a linear head and min-max normalization to  $[0, 1]$ . Optional edge features are refined by a bond-refinement module and projected to edge-type logits.

### D.3. Training

**Rank noise for robustness.** To reduce the train-inference gap (since ground-truth ranks are unavailable during sampling), we add noise to canonical ranks during training:

$$r_{\text{noised}} = r + \mathcal{N}(0, \sigma_r^2) \cdot (t) \quad (74)$$

The noise decays with  $t$  to avoid large perturbations when the molecule is nearly reconstructed.

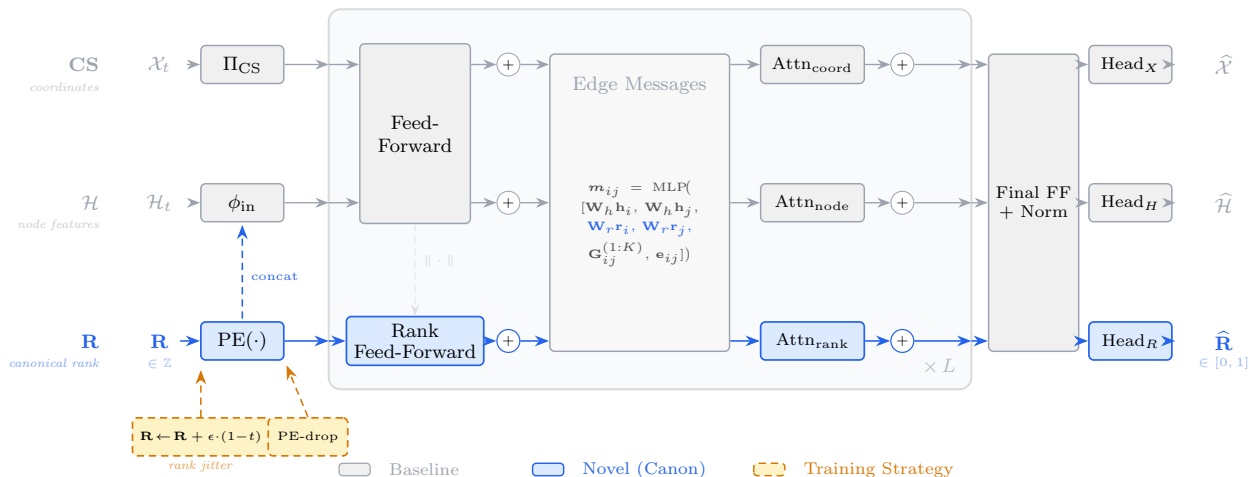


Figure 7. Overview of the Canon architecture. Three parallel streams—coordinates (**CS**), node features (**H**), and canonical rank (**R**)—are updated through  $L$  message-passing layers consisting of feed-forward blocks, pairwise edge messages, and attention-based aggregation. Gray blocks denote components inherited from Semla; blue blocks denote novel canonical-rank components introduced by Canon. Rank information is jittered prior to positional encoding and dropped afterward as a training strategy.

**Auxiliary rank prediction.** Optionally, the model predicts canonical rank as an auxiliary task:

$$\hat{r} = \sigma(\text{MLP}(h_i^{(L)})) \in [0, 1] \quad (75)$$

with min-max normalization across atoms and training loss:

$$\mathcal{L}_{\text{rank}} = \lambda_r \cdot \mathbb{E} [\|r^* - \hat{r}\|^2] \quad (76)$$

This enables dynamic rank estimation during inference.

---

**Algorithm 3** Gap-free canonical training with optional OT coupling (no OT needed at inference)

---

**Require:** Dataset  $(G, X)$ ; symmetry group  $G \in \{S_N, S_N \times SO(3)\}$ ; canonicalizer  $\Psi_G$  producing rank-ordered canonical coordinates  $\tilde{X}$ .

- 1: **Group setup:** choose  $\Psi_G$  and the canonical-space prior  $q_1$  accordingly.
  - 2: Canonicalize each training example once:  $\tilde{Z}_0 = \text{vec}(\Psi(X))$ ; treat rank as the fixed coordinate index.
  - 3: Choose a marginal prior  $q_1$  in canonical space (Aligned or Isotropic).
  - 4: **for** each training iteration **do**
  - 5:   Sample minibatch  $\{\tilde{z}_0^i\} \sim q_0$  and noise minibatch  $\{\tilde{z}_1^i\} \sim q_1$ .
  - 6:   Construct a coupling:
    - Product: pair  $(\tilde{z}_0^i, \tilde{z}_1^i)$ ; or
    - OT/Sinkhorn: compute plan  $\Pi$  minimizing  $\sum_{i,j} \Pi_{i,j} \|\tilde{z}_0^i - \tilde{z}_1^j\|^2$ , then sample pairs from  $\Pi$ .
  - 7:   Train the flow/diffusion model on canonical coordinates using fixed/predicted rank embeddings  $\text{PE}(i)$ .
  - 8: **end for**
- 

**Optimal transport annealing.** We observe that optimal transport (OT) matching between noise and target molecules reduces flow variance but may limit generalization. We implement **OT annealing**: the probability of using OT decreases linearly during training:

$$p_{\text{OT}}(\text{epoch}) = \max\left(0, 1 - \frac{\text{epoch}}{\text{max epochs}}\right) \quad (77)$$

This allows the model to benefit from OT alignment early in training while learning to handle arbitrary noise-target pairings later.

**Training objective.** The overall training loss combines coordinate regression, categorical cross-entropy for discrete features, and optional rank prediction:

$$\mathcal{L}_{\text{coord}} = \frac{1}{N} \sum_{i=1}^N \|\hat{\mathcal{X}}_{0,i} - \mathcal{X}_{0,i}\|^2, \quad \mathcal{L}_{\text{type}} = \frac{1}{N} \sum_{i=1}^N \text{CE}(\hat{\mathbf{h}}_{0,i}, \mathbf{h}_{0,i}) \quad (78)$$

$$\mathcal{L}_{\text{bond}} = \frac{1}{N^2} \sum_{i,j} \text{CE}(\hat{\mathbf{E}}_{0,ij}, \mathbf{E}_{0,ij}), \quad \mathcal{L}_{\text{rank}} = \frac{1}{N} \sum_{i=1}^N (\hat{r}_i - r_i^*)^2 \quad (79)$$

$$\mathcal{L} = \mathcal{L}_{\text{coord}} + \lambda_{\text{type}} \mathcal{L}_{\text{type}} + \lambda_{\text{bond}} \mathcal{L}_{\text{bond}} + \lambda_{\text{charge}} \mathcal{L}_{\text{charge}} + \lambda_{\text{rank}} \mathcal{L}_{\text{rank}} \quad (80)$$

where default weights are  $\lambda_{\text{type}} = 0.2$ ,  $\lambda_{\text{bond}} = \lambda_{\text{charge}} = 1.0$  and  $\lambda_{\text{rank}} = 0.1$ . Following (Irwin et al., 2024), we set  $\lambda_{\text{bond}} = 0.5$  for QM9 dataset.

The complete training procedure can be summarized into Algorithm 3.

#### D.4. Sampling

---

**Algorithm 4** Unified few-step canonical sampling (Regime A/B + options)

---

**Require:** Symmetry group  $G \in \{S_N, S_N \times SO(3)\}$ ; trained model  $v_\theta$ ; baseline prior  $q_1$ ; step times  $1 = t_K > \dots > t_0 = 0$ .

**Require:** Sampling regime flag  $\text{Regime} \in \{A, B\}$ .

**Require:** Options: (i) aligned prior; (ii) PCS (if enabled, apply  $\Psi$  after each step; requires  $\text{Regime} = B$ ).

- 1: **Note (rank conditioning):** if rank is used as a condition/PE, treat it as the *fixed coordinate index* and *do not* recompute it from noise.
  - 2: **(Optional prior choice):** optionally replace  $q_1$  with an alternative prior.
  - 3: Sample  $\tilde{Z}_{t_K} \sim q_1$  in the fixed canonical coordinate system.
  - 4: **for**  $k = K, \dots, 1$  **do**
  - 5:      $\tilde{Z}_{t_{k-1}} \leftarrow \text{Step}(\tilde{Z}_{t_k}, v_\theta, t_k \rightarrow t_{k-1})$  {ODE/SDE solver}
  - 6:     **if**  $\text{Regime} = B$  **then**
  - 7:          $\tilde{Z}_{t_{k-1}} \leftarrow \Psi(\tilde{Z}_{t_{k-1}})$
  - 8:     **else**
  - 9:         **Note:** do *not* re-canonicalize / re-rank / project during sampling in Regime A; keep rank indices fixed.
  - 10:    **end if**
  - 11: **end for**
  - 12: Output  $\tilde{Z}_{t_0}$ .
  - 13: **Optional invariance restoration:** sample  $g \sim \lambda(G)$  and output  $g \cdot \tilde{Z}_{t_0}$ .
- 

**Recovering invariance through Haar randomization.** The critical observation is that canonical samples, while not themselves drawn from an invariant distribution, can be *lifted* to an invariant distribution through a simple post-processing step. As an immediate result of Theorem 3.1, the following theorem, which generalizes the observation in (Yan et al., 2023) to arbitrary compact groups, provides the theoretical guarantee in sampling invariant distributions:

**Proposition D.1** (Post-hoc randomization yields invariant sampling). *Let  $\tilde{\mu}$  be any distribution on  $\mathcal{M}$ . Define the randomized distribution*

$$\mu := \int_{\mathcal{G}} (g \cdot \cdot) \# \tilde{\mu} \, d\lambda(g). \quad (81)$$

*Then  $\mu$  is  $\mathcal{G}$ -invariant. Moreover, sampling  $g \sim \lambda$  and  $\tilde{\mathbf{Z}} \sim \tilde{\mu}$  and outputting  $g \cdot \tilde{\mathbf{Z}}$  produces  $\mu$ .*

*Proof.* For any  $h \in \mathcal{G}$ ,

$$(h \cdot \cdot) \# \mu = \int_{\mathcal{G}} (hg \cdot \cdot) \# \tilde{\mu} \, d\lambda(g) = \int_{\mathcal{G}} (g' \cdot \cdot) \# \tilde{\mu} \, d\lambda(g') = \mu, \quad (82)$$

where we changed variables  $g' = hg$  and used left-invariance of Haar. □

*Remark D.2* (Canonicalization vs. randomization). Randomization alone enforces invariance but does not simplify the learning problem. Canonicalization enforces a *gauge choice* that can substantially reduce multimodality/variance induced by symmetry and thereby accelerate training convergence, while randomization then restores invariance at the end. The combination of canonicalized training with post-hoc randomization thus achieves the best of both worlds: efficient learning and provably invariant generation.

**Dynamic rank estimation.** Since sampling can be done using a fixed canonical rank as the condition, the dynamic updates of rank estimation is completely optional.

A practical few-step stabilizer is to project intermediate states back to the slice (Regime B in Algorithm 4). In particular, to keep rank conditions and partially noisy samples consistent, we allow projection-to-slice sampling, which we termed as Projected Canonical Sampling (PCS) as detailed in Algorithm 5. Note that projecting the data to the canonical slice while keeping the PE ranks fixed is conceptually equivalent to keeping the data in the original order but update the canonical rank to align the two.

We support two strategies to estimate canonical rank in the absence of ground truth in Algorithm 5:

1. **Predict Mode:** Use the model’s rank prediction output  $\hat{r}$  directly for warping.
2. **Canonicalize Mode:** Periodically recompute canonical rank from intermediate predictions: Discretize current atom type predictions:  $\hat{a}_i = \text{argmax}(\text{atomics}_i)$ . Then apply the same canonicalization algorithm to the predicted molecule. Update rank every  $K$  steps when  $t \geq T$ .

We experimentally find that the predict mode is always better, likely due to the instability of canonicalization algorithms for intermediate noisy states; instead, self-prediction integrates this procedure into the model learning, facilitating robustness and expressivity.

Remarkably, if the dynamic rank estimation is disabled (Regime A in Algorithm 4), the model always conducts a “conditional generation” task given the canonical rank conditions, which is also a valid task with no train-test gap.

The complete sampling procedure can be summarized into Algorithm 4.

---

#### Algorithm 5 Projected Canonical Sampling (PCS)

---

**Require:** Canonicalizer  $\Psi_\phi$ ; learned slice model (score or vector field)  $m_\theta$ ; time grid  $1 = t_K > \dots > t_0 = 0$ .

- 1: Initialize  $\tilde{\mathbf{Z}}_{t_K} \sim q_1$ .
  - 2: **for**  $k = K, \dots, 1$  **do**
  - 3:   Take one reverse step on the slice:  $\hat{\mathbf{Z}} \leftarrow \text{Step}(\tilde{\mathbf{Z}}_{t_k}, m_\theta, t_k \rightarrow t_{k-1})$ .
  - 4:   Project:  $\tilde{\mathbf{Z}}_{t_{k-1}} \leftarrow \Psi_\phi(\hat{\mathbf{Z}})$ .
  - 5: **end for**
  - 6: Output canonical sample  $\hat{\mathbf{Z}}_0 = \tilde{\mathbf{Z}}_{t_0}$  and optionally randomize by  $g \sim \lambda$ .
- 

## E. Experimental Details and Additional Results

### E.1. Experimental Details

#### E.1.1. EVALUATION METRICS

We evaluate generated molecules using the following metrics. Let  $\mathcal{M} = \{m_i\}_{i=1}^N$  denote the set of generated molecules, and let  $\tilde{m}_i$  denote the RDKit force-field optimized conformation of  $m_i$ .

- **Validity.** The fraction of generated molecules that pass RDKit chemical validity checks.
- **Atom Stability.** The proportion of all generated atoms whose explicit valence lies within a predefined allowed range.
- **Molecule Stability.** The fraction of generated molecules in which *every* atom is stable.
- **Uniqueness.** Among molecules that can be converted to canonical SMILES, the fraction of distinct SMILES strings.

- **Novelty.** Among molecules with valid canonical SMILES, the fraction absent from the training set.
- **Opt-RMSD.** The average root-mean-square deviation between generated and optimized conformations:

$$\text{Opt-RMSD} = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \text{RMSD}(m_i, \tilde{m}_i),$$

where  $\mathcal{I} = \{i : \text{optimization succeeds and RMSD is finite}\}$ . This quantifies how far the generated 3D geometry deviates from a local energy minimum.

- **NFE.** The Number of Function Evaluations, i.e. the total number of neural network forward passes required by the ODE solver during sampling. NFE serves as a hardware-agnostic measure of generation cost.
- **Sampling Time.** The wall-clock time to generate the full evaluation set ( $N=1,000$  molecules).

### E.1.2. HYPER-PARAMETERS.

Our most hyper-parameters completely follow (Irwin et al., 2024) for both Canonical SemlaFlow and our CanonFlow. All of our backbones consist of  $L=12$  equivariant message-passing layers ( $d_{\text{model}}=384$ ,  $d_{\text{msg}}=128$ ,  $d_{\text{edge}}=128$ ,  $H=32$  attention heads,  $S=64$  coordinate sets with per-set length normalization), with two edge-aware layers at the input and output. A learnable molecule-size embedding of dimension 64 is concatenated to each atom’s input features. All output heads (atom type, charge, bond type, and the optional canonical-rank head) are two-layer MLPs with SiLU activation. Self-conditioning is employed by default: during training, with probability 0.5 the model first produces a preliminary prediction with zeroed conditioning inputs, which is then fed back as additional context; during inference it is applied at every integration step.

Models are trained for 200 epochs on GEOM-DRUG and 300 epochs on QM9 with Adam ( $\text{lr}=3 \times 10^{-4}$ , `amsgrad`, no weight decay), a linear warm-up (10 000 steps for GEOM-DRUG, 2 000 for QM9) followed by a constant learning rate, gradient clipping at global norm 1.0, and bucket-based dynamic batching with a cost budget of 4096 and linear cost scaling. An exponential moving average (EMA, decay 0.999) of model parameters is maintained and used for all inference. All training is conducted in FP32 precision. The total training time is approximately 10 hours for QM9 and 28 hours for GEOM-DRUG on a single NVIDIA A100 80GB GPU.

Coordinates are interpolated linearly,  $\mathbf{x}_t = (1-t)\mathbf{x}_0 + t\mathbf{x}_1$ , with additive Gaussian noise ( $\sigma=0.2$ ). Atom types and bond types use the uniform-sample categorical interpolation (Campbell et al., 2022). The interpolation time is drawn from  $\text{Beta}(2, 1)$ . Equivariant optimal transport (Kabsch alignment + Hungarian matching) is optionally applied to coordinate pairs. When classifier-free guidance (CFG) is enabled, the positional-encoding dropout rate defaults to  $p_{\text{drop}}=0.1$ .

Table 6. CanonFlow w/ our novel Canon architecture on GEOM-DRUG. Canonicalization is conducted on  $S_N$ .

Model	Mol Stab $\uparrow$	Valid $\uparrow$	NFE
EQGAT-diff	93.4 $\pm$ 0.21	94.6 $\pm$ 0.24	500
SemlaFlow <sub>50</sub> w/ OT	97.0 $\pm$ 0.21	93.9 $\pm$ 0.12	50
<b>CanonFlow<sub>50</sub></b>	<b>98.0</b> $\pm$ 0.08	<b>95.4</b> $\pm$ 0.09	50
SemlaFlow <sub>100</sub> w/ OT	97.3 $\pm$ 0.08	93.9 $\pm$ 0.19	100
<b>CanonFlow<sub>100</sub></b>	<b>98.4</b> $\pm$ 0.02	<b>95.9</b> $\pm$ 0.08	100
Data	100.0	100.0	–

## E.2. Additional Results

In addition to our state-of-the-art main results presented in Section 5, we provide more comprehensive ablation studies here.

**Superior few-step generation with CanonFlow.** Table 6 presents few-step generation results of our CanonFlow (trained with OT annealing), which consistently outperforms the SemlaFlow baseline w/o canonicalization, and the advantage is still significant even with only 50 steps. This further supports our claim that canonicalization not only improves the upper bound of generative models, but also provides strong guidance signal through the conditions to facilitate few-step generation.

Table 7. Canonical SemlaFlow on GEOM-DRUG dataset w/ classifier-free guidance (CFG). Canonicalization is conducted on  $S_N$ .

Model	CFG scale	Mol Stab $\uparrow$	Valid $\uparrow$	NFE
EQGAT-diff	–	93.4 $\pm$ 0.21	94.6 $\pm$ 0.24	500
SemlaFlow <sub>50</sub>	–	97.0 $\pm$ 0.21	93.9 $\pm$ 0.12	50
Canon. SemlaFlow <sub>50</sub>	1.0	<b>97.6</b> $\pm$ 0.13	94.6 $\pm$ 0.32	50
Canon. SemlaFlow <sub>50</sub>	1.5	97.4 $\pm$ 0.07	<b>94.7</b> $\pm$ 0.12	50
Canon. SemlaFlow <sub>50</sub>	2.0	96.7 $\pm$ 0.04	94.5 $\pm$ 0.03	50
SemlaFlow <sub>100</sub>	–	97.3 $\pm$ 0.08	93.9 $\pm$ 0.19	100
Canon. SemlaFlow <sub>100</sub>	1.0	<b>98.1</b> $\pm$ 0.03	95.0 $\pm$ 0.20	100
Canon. SemlaFlow <sub>100</sub>	1.5	97.8 $\pm$ 0.03	95.0 $\pm$ 0.09	100
Canon. SemlaFlow <sub>100</sub>	2.0	97.5 $\pm$ 0.08	<b>95.1</b> $\pm$ 0.07	100
Data	–	100.0	100.0	–

Table 8. Canonical SemlaFlow on GEOM-DRUG dataset w/ classifier-free guidance (CFG). Canonicalization is conducted on  $S_N \times SO(3)$ .

Model	CFG scale	Mol Stab $\uparrow$	Valid $\uparrow$	NFE
EQGAT-diff	–	93.4 $\pm$ 0.21	94.6 $\pm$ 0.24	500
SemlaFlow <sub>100</sub>	–	97.3 $\pm$ 0.08	93.9 $\pm$ 0.19	100
Canon. SemlaFlow <sub>100</sub>	1.0	<b>97.9</b> $\pm$ 0.09	93.8 $\pm$ 0.30	100
Canon. SemlaFlow <sub>100</sub>	1.5	97.7 $\pm$ 0.06	94.2 $\pm$ 0.23	100
Canon. SemlaFlow <sub>100</sub>	2.0	97.6 $\pm$ 0.09	<b>94.4</b> $\pm$ 0.05	100
Data	–	100.0	100.0	–

**Effects of classifier-free guidance.** We further ablate the sampling quality with CFG. To comprehensively evaluate the effects of canonicalized symmetry groups, we report the performance of canonical SemlaFlow on both  $S_N$  (trained w/ OT annealing) in Table 7 and  $S_N \times SO(3)$  (trained w/o OT) in Table 8. Within an appropriate range, a larger CFG scale tends to improve the validity metric, at the cost of slightly decreasing the molecule stability. Interestingly, the improvement is more significant on the  $S_N \times SO(3)$  canonicalization, since the prediction relies more on the canonical information, which benefits from extrapolating the conditional shortcut and unconditional baseline.

**Ablation on aligned prior and OT.** In Table 9 we study the interaction between aligned priors and OT for both equivariant and non-equivariant backbones under  $S_N$  canonicalization. Without PE, the model is not conditioned on an additional symmetry-breaking variable, so training-time equivariant OT only changes the endpoint coupling while preserving the inference prior marginal; in this case, OT is compatible with the aligned slice prior and can be viewed as an aligned coupling on the canonical space. The situation changes once a non-equivariant conditioning signal such as canonical-rank PE is introduced. In this case, the model learns under a joint law  $\tilde{\pi}_1^{\text{tr}}(C, \tilde{Z}_1)$  over the condition and the initial noise, and training–sampling consistency requires this law to match the inference law  $\tilde{\pi}_1^{\text{inf}}(C, \tilde{Z}_1)$ . A naive combination of PE, aligned prior, and OT does not explicitly preserve this conditional joint law: OT may improve endpoint matching, but it can simultaneously destroy the alignment between  $C$  and  $\tilde{Z}_1$  that the sampler relies on. We therefore include this non-aligned variant as a negative ablation, highlighted in red in Table 9. Its weaker few-step performance supports our theoretical claim that OT remains beneficial only when the induced coupling is compatible with the conditioning joint law.

Table 9. Canonical SemlaFlow ablations with aligned prior on GEOM-DRUG under  $S_N$  canonicalization. We compare equivariant and non-equivariant backbones with aligned priors and OT. The red rows report a deliberately non-aligned variant that combines PE with OT without preserving the conditional start-time joint law  $\tilde{\pi}_1^{\text{tr}}(C, \tilde{Z}_1) = \tilde{\pi}_1^{\text{inf}}(C, \tilde{Z}_1)$ ; its weaker few-step performance is consistent with the training–sampling consistency analysis.

Group	Model	PE	OT	Mol Stab $\uparrow$	Valid $\uparrow$	NFE
	EQGAT-diff	-	-	93.4 $\pm$ 0.21	94.6 $\pm$ 0.24	500
	SemlaFlow <sub>20</sub>	-	equivariant	95.3 $\pm$ 0.14	93.0 $\pm$ 0.10	20
	SemlaFlow <sub>50</sub>	-	equivariant	97.0 $\pm$ 0.21	93.9 $\pm$ 0.12	50
	SemlaFlow <sub>100</sub>	-	equivariant	97.3 $\pm$ 0.08	93.9 $\pm$ 0.19	100
$S_N$	Canon. SemlaFlow <sub>20</sub>	-	equivariant	95.1 $\pm$ 0.09	93.5 $\pm$ 0.17	20
	Canon. SemlaFlow <sub>50</sub>	-	equivariant	97.3 $\pm$ 0.08	94.4 $\pm$ 0.02	50
	Canon. SemlaFlow <sub>100</sub>	-	equivariant	97.5 $\pm$ 0.10	94.2 $\pm$ 0.11	100
$S_N$	Canon. SemlaFlow <sub>20</sub>	True	equivariant	92.6 $\pm$ 0.18	91.7 $\pm$ 0.20	20
	Canon. SemlaFlow <sub>50</sub>	True	equivariant	96.3 $\pm$ 0.03	93.9 $\pm$ 0.07	50
	Canon. SemlaFlow <sub>100</sub>	True	equivariant	97.5 $\pm$ 0.10	94.2 $\pm$ 0.20	100
	Data			100.0	100.0	-

## F. Samples from Canonical Diffusion

This section presents samples from Canonicalized SemlaFlow trained on GEOM-DRUG. The samples were generated randomly, but we have rotated them where necessary to aid visualization.

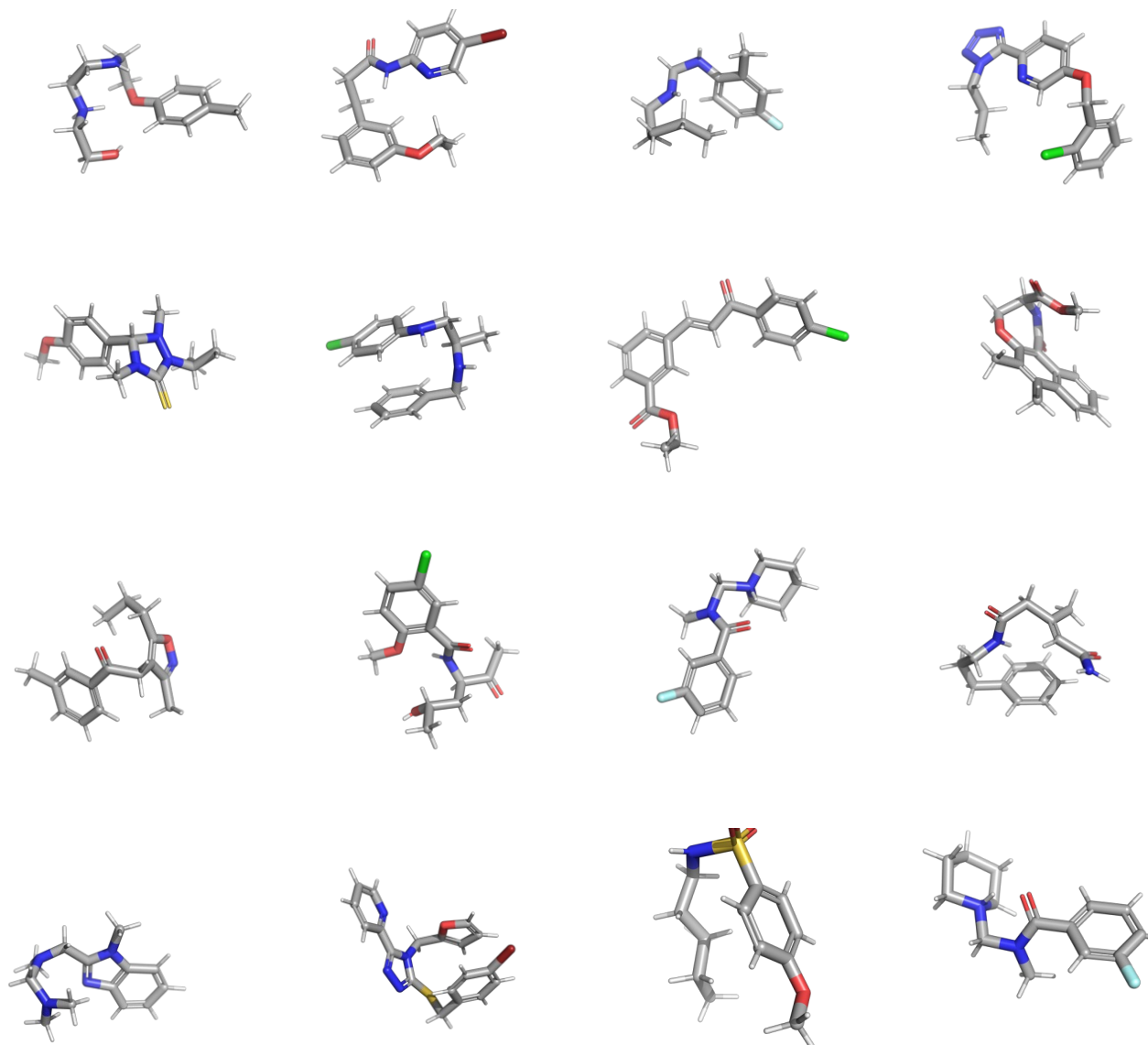


Figure 8. Random samples from Canonicalized Semlaflow model trained on GEOM DRUG. These samples were generated using 100 ODE integration steps.