

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

Elign: Equivariant Diffusion Model Alignment from Foundational Machine Learning Force Fields

Anonymous Authors¹

Abstract

Generative models for 3D molecular conformations must respect Euclidean symmetries and concentrate probability mass on thermodynamically favorable, mechanically stable structures. However, E(3)-equivariant diffusion models often reproduce biases from semi-empirical training data rather than capturing the equilibrium distribution of a high-fidelity Hamiltonian. While physics-based guidance can correct this, it faces two computational bottlenecks: expensive quantum-chemical evaluations (e.g., DFT) and the need to repeat such queries at every sampling step. We present Elign, a post-training framework that amortizes both costs. First, we replace expensive DFT evaluations with a faster, pretrained foundational machine-learning force field (MLFF) to provide physical signals. Second, we eliminate repeated run-time queries by shifting physical steering to the training phase. To achieve the second amortization, we formulate reverse diffusion as a reinforcement learning problem and introduce Force–Energy Disentangled Group Relative Policy Optimization (FED-GRPO) to fine-tune the denoising policy. FED-GRPO includes a potential-based energy reward and a force-based stability reward, which are optimized and group-normalized independently. Experiments show that Elign generates conformations with lower gold-standard DFT energies and forces, while improving stability. Crucially, inference remains as fast as unguided sampling, since no energy evaluations are required during generation.

1. Introduction

The generation of realistic three-dimensional molecular conformations is a central problem in computational chemistry, materials science, and drug discovery (Xu et al., 2023; Hoogeboom et al., 2022). For practical use, a generative model must satisfy two requirements. First, it must respect the symmetries of physics, most notably invariance to rigid body translations and rotations. Second, it must generate samples that correspond to low energy and physically stable configurations.

Score-based diffusion models with E(3)-equivariant architectures have shown notable performance in molecular generation (Xu et al., 2023; Hoogeboom et al., 2022; Cornet et al., 2025). In the idealized case where the training data are drawn from a true thermodynamic equilibrium, the data distribution itself would follow the Boltzmann law (Car & Parrinello, 1985). Under this condition, maximizing likelihood naturally coincides with favoring mechanically plausible, low-energy configurations. In practice, however, this rarely holds. Standard datasets are typically constructed via semi-empirical relaxations (Bannwarth et al., 2019) or heuristic enumeration (Ramakrishnan et al., 2014; Isert et al., 2022), approximating equilibrium only coarsely. Consequently, the resulting model replicates the biases of the dataset generation rather than the equilibrium distribution of a high-fidelity Hamiltonian (Hohenberg & Kohn, 1964; Kohn & Sham, 1965). This gap highlights a limitation of diffusion models: the likelihood objective does not, by design, inherently enforce physical stability.

A direct way to mitigate this is to introduce physical guidance or rewards during the generation process (Wu et al., 2022). In this setting, an oracle approximating the potential energy surface penalizes mechanically unstable configurations during sampling. While prior work (Zhou et al., 2025) has attempted to use first-principles oracles like density functional theory (DFT) (Hohenberg & Kohn, 1964; Kohn & Sham, 1965), the computational cost is often high. As a result, most approaches are forced to rely on tractable surrogates and compromises, such as applying rewards at terminal sampling stages (Zhou et al., 2025), relying on post-processing (Wu et al., 2022), or utilizing semi-empirical potential methods (Shen et al., 2024; Zhou et al., 2025). No-

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

tably, many guidance-based approaches (Shen et al., 2024) require *run-time alignment*, meaning that energy or force evaluations must be performed during sampling, incurring additional computational overhead at inference. Another line of work reframes reward-guided diffusion as stochastic optimal control, leading to algorithms such as Adjoint Matching (Havens et al., 2025) and Adjoint Schrödinger Bridge Sampler (Liu et al., 2025a). These methods rely on adjoint equations that assume reward differentiability with respect to the state. In molecular diffusion models with mixed discrete and continuous variables, this assumption could fail, requiring gradients to be approximated using zeroth-order or surrogate estimators, e.g., Simultaneous Perturbation Stochastic Approximation (SPSA) (Spall, 1992; Shen et al., 2024).

MLFFs offer a natural mechanism for injecting physical constraints into generative models (Chmiela et al., 2017; Schütt et al., 2017). By amortizing the cost of expensive first-principles calculations (Car & Parrinello, 1985) into the training phase, inference reduces to efficient neural network evaluations of energies and forces (Behler, 2021; Unke et al., 2021). We see this as the *first level of amortization* in our framework. Moreover, a suggestive connection exists between MLFFs and diffusion models: the denoising objective used to train diffusion models (Hoogeboom et al., 2022; Xu et al., 2023) has also proven effective for MLFF pretraining (Zaidi et al., 2022; Feng et al., 2023). It is also worth noting that historically, MLFFs are typically trained on narrow, system-specific datasets (Chmiela et al., 2017; 2018) such as single molecules, small chemical families, or homogeneous materials. Thus, they do not define a unified potential across the vast chemical space spanned by modern diffusion models. Recent “foundation” MLFFs go a long way toward addressing these limitations (Amin et al., 2025; Mannan et al., 2025; Unke et al., 2021; Wood et al., 2025): trained on large, chemically diverse quantum-mechanical datasets (Chanussot et al., 2021; Tran et al., 2023; Eastman et al., 2023; Unke et al., 2024), they approximate potential energy surfaces across wide molecular classes while retaining amortized efficiency.

MLFFs could enable physically grounded simulation pipelines, such as molecular dynamics or equilibrium sampling. It can also provide guidance that is more accurate than purely data-driven generators. However, those methods remain expensive when many diverse conformations must be generated, as it requires long trajectories and repeated force evaluations. This motivates a *second level of amortization*: compiling physical constraints into the diffusion model itself, shifting computation from inference to a post-training stage. By *post-training*, we mean an additional training stage applied after the likelihood-based pretraining. Together, MLFFs amortize quantum-mechanical calculations into reusable potentials, while post-training amortizes

repetitive inference into the denoiser, eliminating run-time overhead. Based on these, we introduce Elign, a framework that fine-tunes E(3)-equivariant diffusion models using MLFF-derived rewards. Elign employs potential-based reward shaping derived from MLFF energies. Both energy and force signals are combined under a disentangled group-relative policy optimization scheme (FED-GRPO), analogous to MLFF training. Experiments on QM9 and GEOM-drugs show that Elign generates low-energy, mechanically stable conformations, outperforming runtime-guided methods while matching the inference speed of unguided samplers.

2. Preliminaries

Denoising diffusion models. Diffusion models define a generative distribution $p_\theta(z_0)$ over data $p_{\text{data}}(z_0)$ on a space \mathcal{Z} by learning to reverse an iterative forward noising process. The forward process is defined by the SDE $dz_t = b_t(z_t) dt + \sigma_t d\mathbf{B}_t$ with initial condition $z_0 \sim p_{\text{data}}$. Let $p_t = \text{Law}(z_t)$. For the general forward diffusion $dz_t = -\frac{1}{2}\beta_t z_t dt + \sqrt{\beta_t} d\mathbf{B}_t$, the transition distribution is $p(z_t|z_0) = \mathcal{N}(z_t; \alpha_t z_0, \bar{\sigma}_t^2 \mathbf{I})$ with $\alpha_t = \exp(-\frac{1}{2} \int_0^t \beta_s ds)$ and $\bar{\sigma}_t^2 = 1 - \exp(-\int_0^t \beta_s ds)$. With a constant noise schedule $\beta_t = 2$, we recover the Ornstein–Uhlenbeck forward process $dz_t = -z_t dt + \sqrt{2} d\mathbf{B}_t$ with $\alpha_t = e^{-t}$ and $\bar{\sigma}_t^2 = 1 - e^{-2t}$. This gives a signal-to-noise ratio $\text{SNR}(t) := \alpha_t^2/\bar{\sigma}_t^2$ that decreases monotonically in t , ensuring that as $t \rightarrow T$, $p_T \approx \mathcal{N}(0, \mathbf{I})$. To define the reverse process, let $z_t^\leftarrow = z_{T-t}$ denote the time-reversed trajectory with law $p_t^\leftarrow = p_{T-t}$. By Nelson’s theorem (Nelson, 1967) with mild regularity conditions, the reverse-time process satisfies the drift relation $b_t(z) + b_{T-t}^\leftarrow(z) = a_t \nabla \log p_t(z)$ where $a_t = \sigma_t \sigma_t^\top$, yielding the reverse SDE $dz_t^\leftarrow = (-b_{T-t}(z_t^\leftarrow) + a_{T-t} \nabla \log p_{T-t}(z_t^\leftarrow)) dt + \sigma_{T-t} d\mathbf{B}_t$. In practice, the score function is approximated by a neural network $s_\theta : \mathcal{Z} \times [0, T] \rightarrow \mathcal{Z}$ trained via denoising score matching by minimizing $\mathbb{E} [\|s_\theta(z_t, t) - \nabla_{z_t} \log p(z_t|z_0)\|^2]$.

Equivariant diffusion models. Let \mathcal{G} be a group acting on a space \mathcal{Z} . A function $f : \mathcal{Z} \rightarrow \mathcal{Z}$ is \mathcal{G} -equivariant if $f(g \cdot z) = g \cdot f(z)$ for all $g \in \mathcal{G}, z \in \mathcal{Z}$, and a distribution p on \mathcal{Z} is \mathcal{G} -invariant if its density satisfies $p(g \cdot z) = p(z)$. To construct a diffusion model that generates samples from a \mathcal{G} -invariant law, both the forward and reverse SDEs must respect this symmetry. For the forward process, the Fokker–Planck equation $\partial_t p_t = \frac{1}{2} \langle \sigma_t \sigma_t^\top, \nabla^2 p_t \rangle - \nabla \cdot (b_t p_t)$ preserves invariance when b_t is equivariant and $\sigma_t \sigma_t^\top$ commutes with the group action. By defining a \mathcal{G} -equivariant forward process, the true reverse-time drift is automatically \mathcal{G} -equivariant, since it is composed of an equivariant drift b_t and an equivariant score function $\nabla \log p_t$ (for a \mathcal{G} -invariant distribution p_t , $\nabla \log p_t(g \cdot z) = g \cdot \nabla \log p_t(z)$). In prac-

tice, the unknown true score is estimated by a parameterized model $s_\theta(z_t, t)$. To preserve this symmetry, the score model must be parameterized as a \mathcal{G} -equivariant neural network. The reverse process must also be initiated by drawing samples from a \mathcal{G} -invariant terminal distribution p_T .

***N*-body data and subspace diffusion.** An N -body system is defined by state $\mathbf{z} = [\mathbf{x}, \mathbf{h}] \in \mathbb{R}^{N \times 3} \times \mathbb{R}^{N \times d_h}$, comprising positions $\mathbf{x} \in \mathcal{X}$ and invariant features $\mathbf{h} \in \mathcal{H}$ (e.g., atom types). The physical distribution $p(\mathbf{z})$ is invariant under the Euclidean group $E(3)$ acting on \mathcal{X} . Translation invariance implies that $p(\mathbf{z})$ is not normalizable on the full space \mathbb{R}^{3N} . We therefore restrict the diffusion to the linear subspace of center-of-mass (CoM) zero configurations (Hoogeboom et al., 2022), $\mathcal{M} = \{\mathbf{x} \in \mathbb{R}^{N \times 3} \mid \sum_i \mathbf{x}_i = 0\}$. We define the projection operator \mathbf{P}_{CoM} as a block-diagonal matrix: on positions it is the $3N \times 3N$ centering projector, while on features it acts as identity. Explicitly, $\mathbf{P}_{\text{CoM}} = \text{diag}(\mathbf{P}_{\mathcal{M}}, \mathbf{I}_{Nd_h})$ where $\mathbf{P}_{\mathcal{M}}$ projects onto \mathcal{M} . The forward SDE is modified to diffuse only within this subspace: $d\mathbf{z}_t = -\frac{1}{2}\beta_t \mathbf{z}_t dt + \sqrt{\beta_t} \mathbf{P}_{\text{CoM}} d\mathbf{B}_t$. Because $\mathbf{P}_{\mathcal{M}}$ is idempotent (that is, $\mathbf{P}_{\mathcal{M}} \circ \mathbf{P}_{\mathcal{M}} = \mathbf{P}_{\mathcal{M}}$) and commutes with the $E(3)$ action restricted to rotations and reflections, the marginal distributions p_t remain supported on \mathcal{M} and retain $E(3)$ -invariance. The score network s_θ is parameterized to be $E(3)$ -equivariant and to output vectors in the tangent space of \mathcal{M} (i.e., satisfying the zero-CoM constraint). This could be realized by projecting coordinate updates with $\mathbf{P}_{\mathcal{M}}$ at each denoising step.

Machine-learning force fields. Machine-learning force fields (MLFFs) serve as efficient surrogates for quantum mechanical (QM) calculations, enabling dynamics simulations at scales inaccessible to *ab initio* methods. We assume access to a *quantum chemistry oracle* which, for any given molecular configuration \mathbf{z} , computes a reference energy E^{ref} and reference atomic forces \mathbf{F}^{ref} . As direct queries to such oracles are computationally expensive, often scaling cubically with the number of atoms, MLFFs amortize this cost by training a neural network to approximate the potential energy function $E_\phi(\mathbf{z})$. Atomic forces are typically obtained via automatic differentiation as the negative gradient of this potential, $\mathbf{F}_\phi^{(i)} = -\nabla_{\mathbf{x}^{(i)}} E_\phi(\mathbf{z})$, though they may also be parameterized directly via a dedicated force regression head. In practice (e.g., (Wood et al., 2025)), one can first train a direct force-regression head and then fine-tune via automatic differentiation. During training, the network parameters ϕ are optimized by minimizing a combined energy-and-force matching loss over a dataset of reference calculations:

$$\mathcal{L} = \lambda_E \left(E_\phi(\mathbf{z}) - E^{\text{ref}} \right)^2 + \lambda_F \sum_{i=1}^N \left\| \mathbf{F}_\phi^{(i)}(\mathbf{z}) - \mathbf{F}^{\text{ref},(i)}(\mathbf{z}) \right\|^2, \quad (1)$$

where λ_E and λ_F are weighting coefficients. Physical consistency imposes strict symmetry constraints: E_ϕ must be E(3)-invariant (unaffected by global rotation, translation, and inversion), while the derived forces must be E(3)-

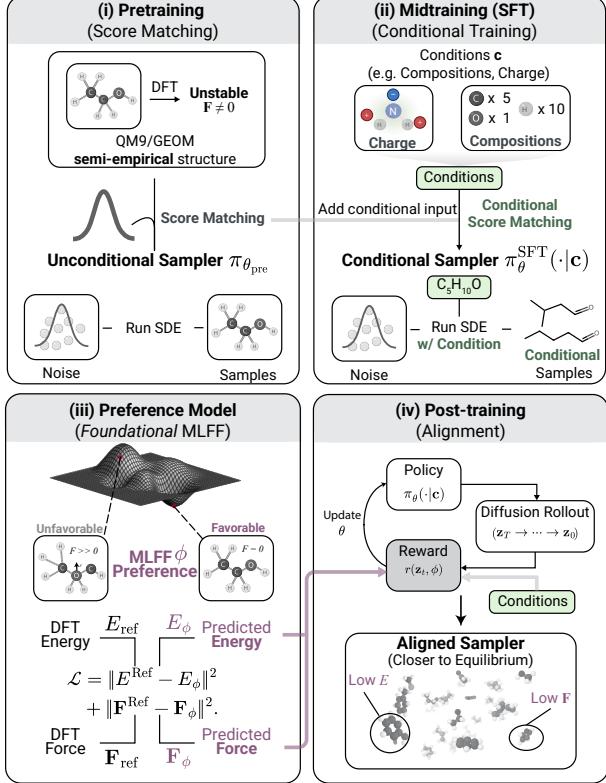


Figure 1. Staged pipeline for equilibrium molecular generation. (i) **Pretraining:** score-matching trains an E(3)-equivariant diffusion model π_θ^{pre} on approximate structures. (ii) **Optional SFT:** conditional fine-tuning improves adherence to discrete specifications. (iii) **Preference model:** a foundation MLFF ϕ provides energy/force signals. (iv) **Post-training:** RL fine-tunes the sampler to improve stability without run-time oracle calls.

equivariant. These constraints are enforced via specialized invariant or equivariant network architectures. Recently, foundational MLFFs have scaled this paradigm to massive heterogeneous datasets, learning universal potentials that are *transferable* across diverse chemical systems.

3. Diffusion post-training with a force-field alignment model

Setup. As shown in Figure 1, we organize equilibrium molecular generation into a staged pipeline analogous to LLM training: base model training, optional supervised conditioning, preference modeling, and post-training alignment.

- (i) **Pretraining (score matching):** We train an equivariant diffusion model $\pi_{\theta_{\text{pre}}}$ via score matching on large collections of approximate equilibrium structures (e.g., QM9-style semi-empirical pipelines). This stage learns a broad generative distribution but does not enforce physical stability.
- (ii) **Optional SFT (conditional training):** When explicit

constraints are required, one could optionally fine-tune the model to condition on variables such as by training a conditional diffusion model $\pi_{\theta_{\text{soft}}}(\cdot | \mathbf{c})$ on paired data (\mathbf{c}, \mathbf{z}) using the same denoising objective. This step corresponds to in-

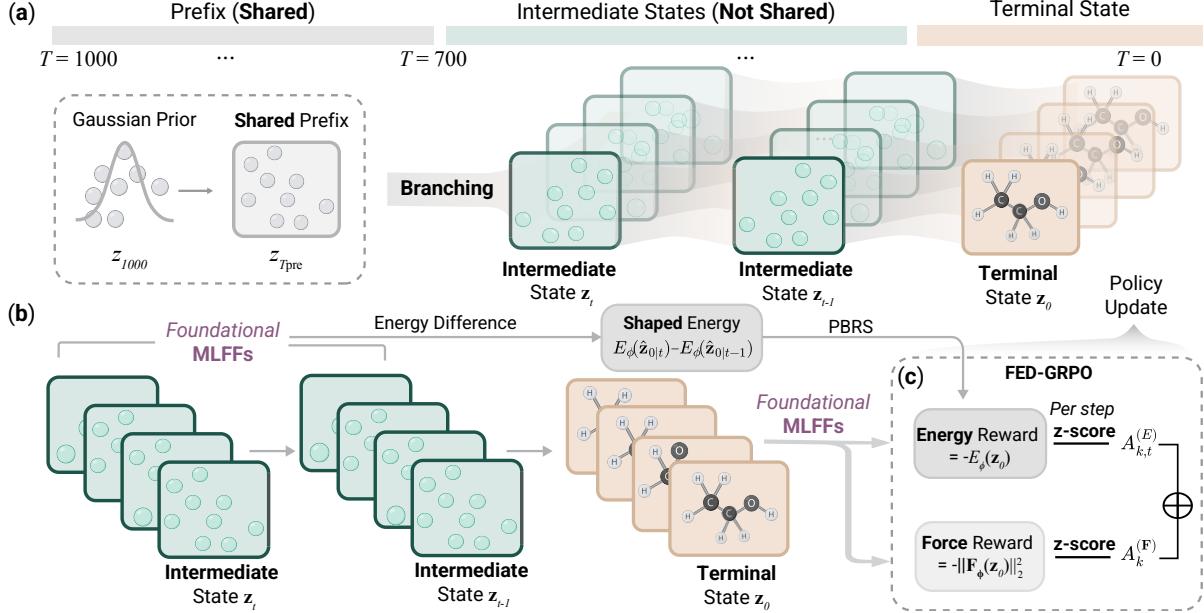


Figure 2. Overview of Elgin. (a) Rollout branching with shared prefix: starting from a CoM-free Gaussian prior at $t = T$, an EGNN policy denoises to $t = 0$. At $t = T_{\text{prefix}}$, we cache a shared prefix state and branch into K rollouts with independent noise. Each trajectory is propagated to its terminal state \mathbf{z}_0 and scored by a foundational MLFF. (b) Energy-based reward shaping: intermediate predicted clean geometries $\hat{\mathbf{z}}_{0|t}$ are evaluated by the MLFF, and local energy differences $E_\phi(\hat{\mathbf{z}}_{0|t}) - E_\phi(\hat{\mathbf{z}}_{0|t-1})$ provide dense shaping signals that bias the policy toward lower-energy conformations. (c) FED-GRPO: terminal energy and force rewards are z-score normalized separately per timestep, then combined to compute the final advantage for policy updates.

struction SFT in LLMs and improves constraint adherence without introducing preference optimization¹. (iii) **Preference model (foundation MLFF):** We use a pretrained foundation MLFF ϕ as a preference model, providing energies and forces that quantify thermodynamic stability and mechanical equilibrium. (iv) **Post-training (alignment):** We treat reverse diffusion as a trajectory-generating process and apply reinforcement learning to align π_θ with the MLFF-defined preferences, while limiting deviation from the pretrained or SFT-initialized model via trust-region regularization.

Diffusion as an MDP. We formulate the iterative denoising process as a finite-horizon Markov Decision Process (MDP). We index reverse time with t decreasing from T to 0. We initialize our method with a pretrained equivariant diffusion model, $\pi_{\theta_{\text{pre}}}$, which serves as the reference base policy. The state is defined by $S_t = (\mathbf{z}_t, t)$, comprising the latent molecular geometry \mathbf{z}_t and the discrete reverse-diffusion time index t . The *policy* $\pi_\theta(\mathbf{z}_{t-1} \mid \mathbf{z}_t, t)$ executes a *one-step discretization* of the reverse SDE update (e.g., Euler-Maruyama. For simplicity, we assume first-order solvers so the unaugmented state S_t is Markov.). This induces a Gaussian policy parameterized by the score network: $\pi_\theta(\mathbf{z}_{t-1} \mid \mathbf{z}_t, t) = \mathcal{N}(\mathbf{z}_{t-1}; \mu_\theta(\mathbf{z}_t, t), \Sigma_t)$, where μ_θ is the learned reverse mean and Σ_t is given. The *action* is defined as a realization sampled from the Gaussian policy.

¹Due to the page limit, we include results on conditional generation in Appendix E.2.

It is important to distinguish the per-step policy, denoted $\pi(\mathbf{z}_{t-1} \mid \mathbf{z}_t, t)$, from the resulting distribution ρ , which arises from the sequential chaining of these policy steps. All stochasticity in the trajectory arises from the Gaussian perturbation inside the policy itself. There are two apparent sources of randomness: (i) the noise term of the discretized reverse-time SDE and (ii) sampling from a stochastic policy. In our case, these two sources *coincide* because the policy itself is the SDE kernel. Here, we offer an alternative interpretation of this MDP through the lens of stochastic optimal control (of which maximum-entropy RL is a special case). In this view, the goal of RL is to learn a policy that acts as a “control knob,” modulating the drift term $\mu_\theta(\cdot)$ of the diffusion process to maximize the expected reward. This is also analogous to a *learnable guidance function* in diffusion guidance.

3.1. Alignment Objective

The objective of RL post-training is to ensure the final generated structures are physically valid. Validity requires two properties: thermodynamic stability (low potential energy) and mechanical equilibrium (zero net force). To design the reward, we draw inspiration from training MLFFs, which utilize supervision from both energy and force labels to learn the potential energy surface (Eq. 1). Because force correlates with the gradient of energy, these objectives are coupled but distinct; a structure can have low energy but high unstable forces if it sits on a steep slope of the potential energy surface. We therefore define the *terminal alignment*

reward, assigned only at the final step when $t = 0$. Specifically:

$$r_0^{(E)}(\mathbf{z}_0) := -E_\phi(\mathbf{z}_0); r_0^{(\mathbf{F})}(\mathbf{z}_0) := -\|\mathbf{F}_\phi(\mathbf{z}_0)\|_F^2. \quad (2)$$

Here, E_ϕ and \mathbf{F}_ϕ are the energy and force predictions from the pretrained MLFF. For implementation, we normalize the energy and force rewards using formation energy per atom and RMS force.

Energy-based potential shaping. Optimizing solely for a terminal alignment reward yields a sparse learning signal over long reverse diffusion horizons. To facilitate credit assignment while preserving the optimal policy of the original terminal-reward MDP, we add an intermediate signal using *potential-based reward shaping* (PBRS). Traditionally, PBRS relies on a potential function $\Psi(S)$ that estimates the “closeness” of a current state to the goal, much like an admissible heuristic in A* search (Hart et al., 1968). In our setting, this role is naturally played by physics: because our objective is thermodynamic stability, the physical potential energy E_ϕ directly quantifies the distance from equilibrium.

The MLFF energy $E_\phi(\cdot)$ is physically meaningful only on clean atomic coordinates, not on the noisy intermediates \mathbf{z}_t . At each reverse step t we therefore reconstruct a *predicted clean geometry* $\hat{\mathbf{z}}_{0|t}$ from \mathbf{z}_t using the diffusion posterior mean estimate: $\hat{\mathbf{z}}_{0|t} = \frac{1}{\alpha_t}(\mathbf{z}_t + \bar{\sigma}_t^2 s_\theta(\mathbf{z}_t, t))$. We define a shaping potential on this reconstructed estimate: $\Psi(S_t) := -E_\phi(\hat{\mathbf{z}}_{0|t})$, $S_t = (\mathbf{z}_t, t)$. The intermediate shaping reward for a transition $S_t \rightarrow S_{t-1}$ is then

$$r_t^{\text{shape}} := \gamma \Psi(S_{t-1}) - \Psi(S_t), \quad t = T_{\text{prefix}}, \dots, 1. \quad (3)$$

This is the canonical PBRS form and preserves the set of optimal policies under the discounted return. In particular, if we define the (discounted) shaped return-to-go from state S_t as $G_t^{\text{shape}} := \sum_{u=1}^t \gamma^{t-u} r_u^{\text{shape}}$, then the shaping contribution telescopes:

$$G_t^{\text{shape}} = \sum_{u=1}^t \gamma^{t-u} (\gamma \Psi(S_{u-1}) - \Psi(S_u)) = \gamma^t \Psi(S_0) - \Psi(S_t). \quad (4)$$

Intuitively, the return-to-go is the discounted difference between terminal and current state potentials. This shaping provides informative per-step feedback while preserving policy optimality under the discounted return (Ng et al., 1999). We also considered a force-based potential, e.g., $\Psi_F(S_t) = -\|\mathbf{F}_\phi(\hat{\mathbf{z}}_{0|t})\|_F^2$, and apply PBRS analogously. We found it unstable in practice, possibly because forces (first-order derivatives) are more sensitive than energies (zeroth-order) to small geometric errors in $\hat{\mathbf{z}}_{0|t}$.

3.2. Theoretical View of the Alignment Objective

The rewards above specify what properties we desire in the final sample. We now adopt a distributional perspective to

characterize the *terminal law* that post-training encourages. Because the diffusion policy is parameterized by an E(3)-equivariant architecture, post-training preserves invariance of the terminal distribution. For analytical clarity, we focus on the energy objective and ignore PBRS.

Let \mathcal{Z} denote the (CoM-free) configuration space of clean molecular structures, and identify the terminal denoised sample $\mathbf{z}_0 \in \mathcal{Z}$. Let $\rho_{\theta_{\text{pre}}} \in \mathcal{P}(\mathcal{Z})$ denote the terminal distribution induced by the pretrained diffusion model $\pi_{\theta_{\text{pre}}}$, i.e., $\mathbf{z}_0 \sim \rho_{\theta_{\text{pre}}}$ when sampling from $\pi_{\theta_{\text{pre}}}$. Any fine-tuned policy π_θ induces its own terminal law ρ_θ over \mathcal{Z} .

Theorem 1 (Energy-aligned terminal distribution). *Assume the policy class is rich enough that any $\rho \ll \rho_{\theta_{\text{pre}}}$ can be realized as the terminal law of some admissible reverse diffusion policy, and for a trust-region regularized reward maximization objective: $\mathcal{J}(\rho) := \mathbb{E}_{\mathbf{z}_0 \sim \rho}[-E_\phi(\mathbf{z}_0)] - w_{\text{KL}} \text{KL}(\rho \| \rho_{\theta_{\text{pre}}})$, the maximum is attained. Then the maximizer is unique and equals*

$$\rho^*(\mathbf{z}) = \frac{1}{Z_\phi} \rho_{\theta_{\text{pre}}}(\mathbf{z}) \exp(-\beta_{\text{eff}} E_\phi(\mathbf{z})), \quad \beta_{\text{eff}} := \frac{1}{w_{\text{KL}}}, \quad (5)$$

with normalizer $Z_\phi = \int_{\mathcal{Z}} \rho_{\theta_{\text{pre}}}(\mathbf{u}) \exp(-\beta_{\text{eff}} E_\phi(\mathbf{u})) d\mathbf{u}$.

Proof sketch. This is the standard variational form of KL-regularized optimization: enforce normalization with a Lagrange multiplier, take first-order optimality conditions in ρ , and normalize. Strict concavity in ρ (from the negative entropy term) gives uniqueness.

This result admits an energy decomposition view: defining a “prior energy” $E_{\text{prior}}(\mathbf{z}) := -\beta_{\text{eff}}^{-1} \log \rho_{\theta_{\text{pre}}}(\mathbf{z})$, the optimal density satisfies

$$\rho^*(\mathbf{z}) \propto \exp(-\beta_{\text{eff}}[E_\phi(\mathbf{z}) + E_{\text{prior}}(\mathbf{z})]).$$

The terminal distribution is therefore determined by two components: the physical energy $E_\phi(\mathbf{z})$ from the MLFF and an implicit prior induced by the pretrained model. Note that we do not compute E_{prior} explicitly, this is a conceptual decomposition. In this view, alignment sharpens thermodynamic plausibility while remaining confined to the learned data manifold, with the pretrained model acting as both a regularizer and a support constraint (via the absolute-continuity condition). This highlights the importance of high-quality pretraining. Moreover, post-training can be viewed as amortized sampling: instead of applying energy guidance at inference time via repeated E_ϕ evaluations, RL compiles the corresponding Gibbs reweighting into the reverse diffusion policy. We next relate MLFF approximation error to the true physical Boltzmann law:

MLFF approximation and distribution error. Let $E_\star : \mathcal{Z} \rightarrow \mathbb{R}$ be a target potential energy and E_ϕ its MLFF approximation. Fix a reference terminal law $\rho_{\theta_{\text{pre}}}$ and

define $\rho_\star^*(\mathbf{z}) \propto \rho_{\theta_{\text{pre}}}(\mathbf{z}) e^{-\beta_{\text{eff}} E_\star(\mathbf{z})}$, $\rho_\phi^*(\mathbf{z}) \propto \rho_{\theta_{\text{pre}}}(\mathbf{z}) e^{-\beta_{\text{eff}} E_\phi(\mathbf{z})}$. We relate the MLFF error to the distributional error via the total variation distance:

Theorem 2 (Energy error implies distribution error). *Assume the MLFF uniformly approximates the target energy, $\sup_{\mathbf{z} \in \mathcal{Z}} |E_\phi(\mathbf{z}) - E_\star(\mathbf{z})| \leq \delta$ for some $\delta \geq 0$. Then the aligned terminal laws in Eq. (3.2) satisfy*

$$\|\rho_\phi^* - \rho_\star^*\|_{\text{TV}} \leq \tanh(\beta_{\text{eff}} \delta). \quad (6)$$

The bound in Eq. (6) implies that the terminal distribution error is controlled by the product $\beta_{\text{eff}} \delta$. In the small-error or moderate-temperature regime, $\tanh(\beta_{\text{eff}} \delta) \approx \beta_{\text{eff}} \delta$, so the mismatch scales linearly with the uniform energy error. In contrast, when β_{eff} is large, even small MLFF biases can induce substantial deviations in the terminal law. This highlights a practical trade-off: aggressive alignment toward low effective temperatures requires either higher MLFF accuracy in the relevant region or stronger trust-region regularization (smaller β_{eff}) to avoid amplifying model miscalibration. We would also like to highlight that this uniform error is conservative, providing a worst-case guarantee; in practice we expect smaller errors in high-density regions. In the next section, we describe how this theoretical objective is implemented via reinforcement learning.

3.3. RL Training

Force–energy disentangled GRPO (FED-GRPO). To fine-tune the diffusion model, we require a policy gradient algorithm that is both sample-efficient and computationally inexpensive. Standard actor–critic methods are costly, as they would require training a separate value function alongside the score model. Moreover, in our setting the value function must be $E(3)$ -invariant and defined over high-dimensional molecular states; learning such a critic via bootstrapping can introduce instability, in addition to extra compute. We instead adopt *Group Relative Policy Optimization* (GRPO) (Shao et al., 2024), a critic-free algorithm, and introduce a *disentangled* advantage formulation that separately accounts for energy and force signals; we refer to the resulting method as Force–Energy Disentangled GRPO (**FED-GRPO**).

Shared-prefix grouping. We first define a *group* over a diffusion trajectory. In large language models, all members of a group share the same prompt; for example, given a fixed math question, multiple responses are sampled in parallel and compared to identify higher-quality answers. Analogously, we extend this idea to diffusion by rolling out a common reverse-diffusion prefix from $T \rightarrow T_{\text{prefix}}$ (for a chosen prefix time $T_{\text{prefix}} \in \{1, \dots, T\}$) and caching $\mathbf{z}_{T_{\text{prefix}}}$. From this shared state we branch into K stochastic continuations with fresh noise. This shared initialization has also been used in concurrent work (e.g., FlowGRPO (Liu

et al., 2025b)) and lowers variance by making early denoising dynamics comparable.

Two-channel advantages. For each branched rollout k , we distinguish between dense energy feedback and sparse force feedback: (i) **Step-wise energy advantage:** For energy, we utilize the dense PBRS signal. We first compute the raw *return-to-go* at each step t : $G_{k,t}^{(E)} = \sum_{u=1}^t \gamma^{t-u} r_{k,u}^{E,\text{shape}}$. Crucially, we do **not** normalize $r_{k,u}$ per step; we first sum raw shaped rewards to preserve telescoping (Eq. (4)), then normalize only at the advantage level. (ii) **Trajectory-level force advantage:** For force, we assign a single scalar return to the entire trajectory: $G_k^{(\mathbf{F})} = -\left\| \mathbf{F}_\phi(\mathbf{z}_0^{(k)}) \right\|_F^2$.

Disentangled group-relative normalization. Directly summing the raw rewards is undesirable, as their differing scales and temporal structures can lead to interference, with one signal dominating the optimization dynamics. We therefore standardize the two channels differently to mitigate these effects and to disentangle their contributions according to granularity. For **force**, we compute a single group-level mean μ_F and standard deviation σ_F . For **energy**, we compute statistics **per time-step** t , denoted $(\mu_{E,t}, \sigma_{E,t})$, across the K rollouts. This ensures that the advantage at step t is relative to the group’s performance *at that specific step*. The total advantage for rollout k at step t is:

$$\hat{A}_{k,t} = w_E \underbrace{\frac{G_{k,t}^E - \mu_{E,t}}{\sigma_{E,t} + \eta}}_{\text{Step-wise Energy Adv.}} + w_{\mathbf{F}} \underbrace{\frac{G_k^F - \mu_F}{\sigma_F + \eta}}_{\text{Global Force Adv.}}. \quad (7)$$

Here, $\eta > 0$ is a small constant added for numerical stability, and the force term is broadcast to the entire trajectory.

Clipped FED-GRPO objective. Let the probability ratio be $\xi_{k,t}(\theta) = \frac{\pi_\theta(\mathbf{z}_{t-1}^{(k)} | \mathbf{z}_t^{(k)}, t)}{\pi_{\theta_{\text{old}}}(\mathbf{z}_{t-1}^{(k)} | \mathbf{z}_t^{(k)}, t)}$, which determines how much more (or less) likely the current policy is to take a given denoising step relative to the old policy. The clipped surrogate objective aggregates over all steps using the per-step advantage $\hat{A}_{k,t}$: $L(\theta) = \mathbb{E}_k [\sum_{t=1}^{T_{\text{prefix}}} \min(\xi_{k,t}(\theta) \hat{A}_{k,t}, \text{clip}(\xi_{k,t}(\theta), 1-\varepsilon, 1+\varepsilon) \hat{A}_{k,t})]$. This objective maximizes (or minimizes) the likelihood of denoising trajectories proportional to their physics-based advantage, with clipping enforcing a *local* (per-step) trust region that stabilizes training; a *global* trust region can also be added via a KL penalty against the pretrained policy $\pi_{\theta_{\text{pre}}}$. Concretely, this is given as: $\mathcal{L}_{\text{reg}}(\theta) := \mathcal{L}(\theta) + w_{\text{KL}} \text{KL}(\pi_\theta \| \pi_{\theta_{\text{pre}}})$, where $w_{\text{KL}} > 0$ controls the strength of the global trust region. Note that this is a tractable proxy for terminal-law trust region defined in Sec.3.2.

Finally, we summarize the FED-GRPO procedure in Alg. 1–2.

Algorithm 1 FED-GRPO post-training

```

330 1: Input:  $\pi_{\theta_{\text{old}}}$ , MLFF  $\phi$ ,  $T_{\text{prefix}}$ ,  $K$ ,  $(w_E, w_F)$ ,  $\gamma$ , clip  $\varepsilon$ .
331 2: for each iteration do
332 3:   Rollout  $T \rightarrow T_{\text{prefix}}$  with  $\pi_{\theta_{\text{old}}}$ ; cache  $\mathbf{z}_{\text{start}}$ .  $\triangleright$  Shared.
333 4:   for  $k = 1, \dots, K$  do
334 5:      $\tau_k \leftarrow \text{ROLLOUT}(\mathbf{z}_{\text{start}}, \pi_{\theta_{\text{old}}})$ 
335 6:      $(\{G_{k,t}^{(E)}\}_{t=1}^{T_{\text{prefix}}}, G_k^{(F)}) \leftarrow \text{REWARD}(\tau_k, \phi, \pi_{\theta_{\text{old}}})$ .
336 7:      $\hat{A}_k^{(F)} \leftarrow (G_k^{(F)} - \mu_F) / (\sigma_F + \eta)$  over  $k$ .
337 8:      $\hat{A}_{k,t}^{(E)} \leftarrow (G_{k,t}^{(E)} - \mu_{E,t}) / (\sigma_{E,t} + \eta)$  over  $k$  for each  $t$ .
338 9:      $\hat{A}_{k,t} \leftarrow w_E \hat{A}_{k,t}^{(E)} + w_F \hat{A}_k^{(F)}$ ;
339 10:     $\xi_{k,t} \leftarrow \pi_{\theta}(\mathbf{z}_{t-1}^{(k)} | \mathbf{z}_t^{(k)}, t) / \pi_{\theta_{\text{old}}}(\mathbf{z}_{t-1}^{(k)} | \mathbf{z}_t^{(k)}, t)$ .
340 11:    Update  $\theta$  with the GRPO objective using  $(\xi_{k,t}, \hat{A}_{k,t})$ ;  $\triangleright$ 
341        Optionally add a KL Term.
342 12:    set  $\theta_{\text{old}} \leftarrow \theta$ .

```

Algorithm 2 REWARD: energy PBRS return-to-go + terminal force

```

347 1: function REWARD( $\tau, \phi, \pi_{\theta_{\text{old}}}$ )
348 2:   Extract  $\{\mathbf{z}_t\}_{t=0}^{T_{\text{prefix}}}$  from  $\tau$ .
349 3:   for  $t = T_{\text{prefix}}, \dots, 0$  do
350 4:     Compute  $\hat{\mathbf{z}}_{0|t}$  from  $\mathbf{z}_t$   $\triangleright$  Posterior from  $\pi_{\theta_{\text{old}}}$ .
351 5:      $\Psi_t \leftarrow -E_{\phi}(\hat{\mathbf{z}}_{0|t})$ .
352 6:   for  $t = 1, \dots, T_{\text{prefix}}$  do
353 7:      $G_t^{(E)} \leftarrow \gamma^t \Psi_0 - \Psi_t$ .  $\triangleright$  Eq. (4)
354 8:      $G^{(F)} \leftarrow -\|\mathbf{F}_{\phi}(\mathbf{z}_0)\|_F^2$ .
355 9:   return  $\{G_t^{(E)}\}_{t=1}^{T_{\text{prefix}}}, G^{(F)}$ .

```

4. Results

Implementation. We use the public EDM (Hoogeboom et al., 2022) checkpoint as our pretrained base model and UMA (Wood et al., 2025) as the preference model. UMA is a machine-learning force field trained on OMol25 (Levine et al., 2025), OC20 (Chanussot et al., 2021), ODAC23 (Sriram et al., 2024), OMat24 (Barroso-Luque et al., 2024), etc., that predicts per-molecule energies and per-atom forces from atom types and 3D coordinates. UMA has two variants: UMA-1p1-S and UMA-1p1-M. Unless stated otherwise, Elgin uses UMA-1p1-M with potential-based reward shaping (PBRS). Although UMA’s training data may include molecules from QM9 or GEOM-Drugs, the specific 3D conformations generated by EDM are novel and not seen during UMA training. Moreover, we use UMA only as a fixed reward oracle without fine-tuning, and our key evaluations (RDKit validity, DFT energies/forces) are computed independently of UMA. We evaluate on QM9 (Ramakrishnan et al., 2014) and GEOM-Drugs (Axelrod & Gómez-Bombarelli, 2022). QM9 contains 130k small molecules with up to 9 heavy atoms (29 atoms including hydrogens). GEOM-Drugs consists of larger organic compounds with up to 181 atoms (44.2 on average) across 37 million conformations for around 450k molecules. Following Xu et al. (2023), we report atom stability (A), molecule stability (M), RDKit

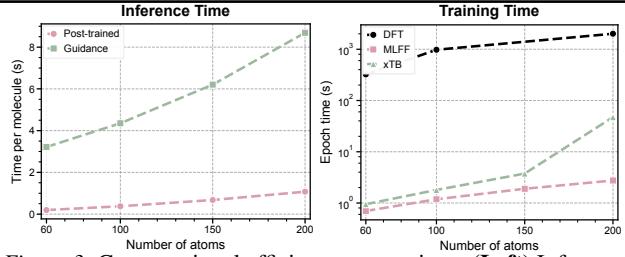


Figure 3. Computational efficiency comparison. **(Left)** Inference time per molecule for post-trained models versus guidance-based methods. **(Right)** Training epoch time comparing different reward oracles on a log scale. For a fair comparison, all methods use terminal-only rewards (no PBRS). Evaluation was conducted on a single NVIDIA H100 GPU.

validity (V), and Validity \times Uniqueness (V \times U). QM9 results average three runs of 10,000 samples; GEOM-Drugs uses 1,024 samples. For GEOM-Drugs, we omit M and V \times U due to limitations in ground-truth bond inference (bond inference is unreliable for larger molecules with 60–200 atoms) and near-saturated uniqueness ($\approx 100\%$). This is in line with community standards (Xu et al., 2023). Baseline results are taken from published work whenever possible.

QM9. Table 1 reports the benchmark results on QM9. Elgin improves molecule stability from 82.00% (EDM) to 93.70% and increases V \times U from 90.70% to 95.31%. Compared with DFT-based RL (EDM+DFT), Elgin matches atom stability and achieves higher V \times U (95.31% vs. 92.87%) without DFT queries during training or sampling.

GEOM-drug. Table 3 evaluates larger drug-like molecules. Elgin reaches 89.43% atom stability and 99.40% validity, improving over EDM (81.3% / 91.9%) and over RLPF (xTB rewards) on atom stability (87.94% vs. 87.52%).

Ablation analysis. We analyze the contribution of individual components in Elgin through controlled ablations on QM9 (Table 2). Our default configuration achieves the strongest overall validity \times uniqueness tradeoff across metrics, with validity \times uniqueness 95.31% (V \times U). Replacing dense PBRS with sparse terminal rewards reduces V \times U from 95.31% to 93.85%, suggesting that intermediate reward signals improve optimization effectiveness over long diffusion horizons. Substituting the larger MLFF (UMA-1p1-M) with its smaller counterpart (UMA-1p1-S) exposes a stability–diversity tradeoff: although the smaller model slightly improves molecule stability (94.83% vs. 93.70%), it leads to a marked decrease in uniqueness, resulting in a 4.5 percentage point reduction in V \times U. A plausible explanation is that the smaller model induces a less smooth energy landscape, yielding higher-variance gradients that concentrate samples around a limited set of local minima. Further decomposing the reward highlights the complementary roles of its components. Energy-only rewards result in low molecule stability (86.00%), while force-only rewards achieve higher stability but reduced diversity (V \times U of 90.21%). Together, these results indicate that jointly optimizing energy and

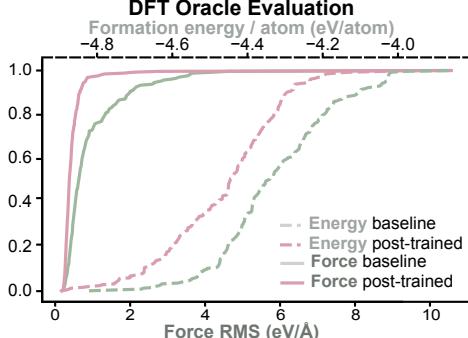


Figure 4. DFT oracle evaluation of generated molecules from QM9. Cumulative distribution functions comparing baseline and post-trained models on force RMS (left) and formation energy per atom (right). The results are obtained from 256 samples.

force objectives, combined with dense reward shaping, is important for balancing stability and diversity.

DFT Oracle Evaluation. To verify that Elign’s improvements reflect genuine physical quality rather than exploiting the MLFF reward, we evaluate generated molecules using independent DFT calculations. Specifically, we use PySCF (Sun et al., 2020) with B3LYP/6-31G(d) to compute total energies and nuclear gradients, then report formation energy per atom (total energy minus isolated atom energies, divided by atom count). Figure 4 shows cumulative distributions for baseline EDM and Elign-trained models. Post-training shifts both distributions toward improved values: reduced force RMS indicates conformations closer to stationary points, while lower formation energy reflects increased thermodynamic stability. Because these metrics come from a DFT oracle not used during training, the results confirm that Elign improves physical fidelity beyond standard stability metrics.

Speed. Figure 3 reports wall-clock scaling. DFT evaluation is orders of magnitude slower than MLFF/UMA and scales steeply with system size; xTB also becomes expensive beyond \sim 150–200 atoms. This observation supports the first level of amortization in our framework. During inference, run-time guidance is slower because it requires oracle gradients at each denoising step, whereas Elign uses the unguided EDM sampling procedure; across 60–200 atoms, post-trained sampling is about 8–16 \times faster in our benchmark. This supports the second level of amortization.

5. Related Work.

Diffusion Models for Molecules. For molecular generation, equivariant diffusion have achieved strong geometric fidelity while respecting symmetries (Hoogeboom et al., 2022; Xu et al., 2023; Cornet et al., 2025; Peng et al., 2023; Song et al., 2023; 2024; Feng et al., 2024; Qiang et al., 2023). Extensions incorporating physical guidance or control include energy-informed priors (Wu et al., 2022; Shen et al., 2024), reinforcement learning with DFT (Zhou et al., 2025), and stochastic control formulations such as Adjoint Matching

Table 1. QM9 3D generation. Atom stability (A), molecule stability (M), RDKit validity (V), and Validity \times Uniqueness (V \times U).

Model	A [%] \uparrow	M [%] \uparrow	V [%] \uparrow	V \times U [%] \uparrow
EDM (Hoogeboom et al., 2022)	98.70	82.00	91.90	90.70
EDM + DPM Solver++ (Lu et al., 2022)	95.72	72.23	87.82	87.82
EDM-BRIDGE (Wu et al., 2022)	98.80	84.60	92.00	90.70
GeoLDM (Xu et al., 2023)	98.90	89.40	93.80	92.70
EDN (Cornet et al., 2025)	98.90	89.10	94.80	92.60
UniGEM (Feng et al., 2024)	99.00	89.80	95.00	93.20
GeoBFN (Song et al., 2024)	99.08	90.87	95.31	92.96
RLPF (EDM + DFT PPO) (Zhou et al., 2025)	99.08	93.37	98.22	92.87
EDM + Rejection Sampling (Zhou et al., 2025)	98.99	89.47	93.20	92.60
GeoLDM + xTB Guidance (Shen et al., 2024)	99.02	90.60	91.40	91.40
EDM + Soft Metropolis-Hastings (Feng et al., 2025)	99.56	91.70	98.70	83.20
GeoLDM + Soft Metropolis-Hastings (Feng et al., 2025)	99.09	91.30	95.10	94.90
EDM + UMA-1p1-S Guidance	98.90	87.00	92.54	92.34
EDM + DPM Solver++ + UMA-1p1-S Guidance	97.72	76.23	89.84	89.84
Elign (EDM + UMA-1p1-M)	99.33	93.70	98.32	95.31
Data (Ground Truth)	99.00	95.20	97.70	97.70

Table 2. QM9 ablations. We vary reward composition (E vs F), reward density (Sparse terminal vs Dense PBRS shaping), and MLFF capacity (UMA-1p1-S vs UMA-1p1-M).

Reward	MLFF				A [%] \uparrow	M [%] \uparrow	V [%] \uparrow	V \times U [%] \uparrow
	E	F	PBRS	S	M			
✓	✓	✓	✓	✓	99.33	93.70	98.32	95.31
✓	✓	✓	✓	✓	99.42	94.83	97.53	90.81
✓	✓			✓	99.02	93.75	96.54	93.85
✓				✓	98.70	86.00	91.70	91.70
	✓			✓	99.40	94.92	96.84	90.21

and Adjoint Schrödinger Bridge Samplers (Havens et al., 2025; Liu et al., 2025a). While effective in biasing samples toward desired properties, these methods typically rely on sparse rewards, differentiable objectives, or run-time oracle evaluations, potentially limiting their scalability. We also include an expanded related-work appendix.

Table 3. GEOM-drug 3D generation. Atom stability (A) and RDKit validity (V).

Model	A [%] \uparrow	V [%] \uparrow
EDM (Hoogeboom et al., 2022)	81.3	91.9
EDM-BRIDGE (Wu et al., 2022)	82.4	91.9
GeoLDM (Xu et al., 2023)	84.4	99.3
EDN (Cornet et al., 2025)	87.0	92.9
UniGEM (Feng et al., 2024)	85.1	98.4
GeoBFN (Song et al., 2024)	85.6	92.08
RLPF (EDM+xBT (Zhou et al., 2025) PPO)	87.52	99.20
Elign (EDM+UMA-1p1-M)	87.94	99.40
Data (Ground Truth)	—	86.5

6. Conclusion & Outlook

We presented Elign, a framework that post-trains equivariant molecular diffusion models with physical preferences using a pretrained foundational MLFF. By casting reverse diffusion as an RL problem and optimizing a force–energy disentangled GRPO objective, Elign improves thermodynamic stability and mechanical equilibrium while preserving fast, unguided inference. In some settings, we observe improved reward values without corresponding gains in chemical validity, suggesting reward hacking; developing reliable diagnostics and mitigation strategies is an important direction for future work.

440 Impact Statement

441 This paper presents work whose goal is to advance the field
 442 of Machine Learning. There are many potential societal
 443 consequences of our work, none which we feel must be
 444 specifically highlighted here.

445 References

446 Amin, I., Raja, S., and Krishnapriyan, A. Towards fast,
 447 specialized machine learning force fields: Distilling
 448 foundation models via energy Hessians. *arXiv preprint*
 449 *arXiv:2501.09009*, 2025.

450 Axelrod, S. and Gómez-Bombarelli, R. Geom, energy-
 451 annotated molecular conformations for property predic-
 452 tion and molecular generation. *Scientific Data*, 9(1):185,
 453 2022. doi: 10.1038/s41597-022-01288-4. URL <https://doi.org/10.1038/s41597-022-01288-4>.

454 Aykent, S. and Xia, T. GotenNet: Rethinking Efficient 3D
 455 Equivariant Graph Neural Networks. In *The Thirteenth*
 456 *International Conference on Learning Representations*,
 457 2025. URL <https://openreview.net/forum?id=5wxCQDtMo>.

458 Bannwarth, C., Ehlert, S., and Grimme, S. Gfn2-xtb—an
 459 accurate and broadly parametrized self-consistent tight-
 460 binding quantum chemical method with multipole elec-
 461 trostatics and density-dependent dispersion contributions.
 462 *Journal of Chemical Theory and Computation*, 15(3):
 463 1652–1671, 2019.

464 Barroso-Luque, L., Shuaibi, M., Fu, X., Wood, B. M.,
 465 Dzamba, M., Gao, M., Rizvi, A., Zitnick, C. L., and
 466 Ulissi, Z. W. Open materials 2024 (omat24) inor-
 467 ganic materials dataset and models. *arXiv preprint*
 468 *arXiv:2410.12771*, 2024.

469 Batatia, I., Kovacs, D. P., Simm, G. N. C., Ortner, C.,
 470 and Csanyi, G. MACE: Higher order equivariant mes-
 471 sage passing neural networks for fast and accurate force
 472 fields. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho,
 473 K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=YPPsngE-ZU>.

474 Batzner, S., Musaelian, A., Sun, L., Geiger, M., Mailoa,
 475 J. P., Kornbluth, M., Molinari, N., Smidt, T. E., and
 476 Kozinsky, B. E(3)-equivariant graph neural networks for
 477 data-efficient and accurate interatomic potentials. *Nature*
 478 *Communications*, 13(1):2453, 2022.

479 Behler, J. Four generations of high-dimensional neural
 480 network potentials. *Chemical Reviews*, 121(16):10037–
 481 10072, 2021.

482 Black, K., Janner, M., Du, Y., Kostrikov, I., and Levine, S.
 483 Training diffusion models with reinforcement learning.
 484 *arXiv preprint arXiv:2305.13301*, 2024.

485 Car, R. and Parrinello, M. Unified approach for molecular
 486 dynamics and density-functional theory. *Physical Review*
 487 *Lettters*, 55(22):2471, 1985.

488 Chanussot, L., Das, A., Goyal, S., Lavril, T., Shuaibi, M.,
 489 Riviere, M., Tran, K., Heras-Domingo, J., Ho, C., Hu, W.,
 490 et al. Open catalyst 2020 (oc20) dataset and community
 491 challenges. *Acs Catalysis*, 11(10):6059–6072, 2021.

492 Chen, C., Ye, W., Zuo, Y., Zheng, C., and Ong, S. P. Graph
 493 networks as a universal machine learning framework for
 494 molecules and crystals. *Chemistry of Materials*, 31(9):
 495 3564–3572, 2019.

496 Chmiela, S., Tkatchenko, A., Sauceda, H. E., Poltavsky, I.,
 497 Schütt, K. T., and Müller, K.-R. Machine learning of accu-
 498 rate energy-conserving molecular force fields. *Science*
 499 *Advances*, 3(5):e1603015, 2017.

500 Chmiela, S., Sauceda, H. E., Müller, K.-R., and Tkatchenko,
 501 A. Towards exact molecular dynamics simulations with
 502 machine-learned force fields. *Nature Communications*, 9
 503 (1):1–10, 2018.

504 Cornet, F., Bartosh, G., Schmidt, M. N., and Naesseth,
 505 C. A. Equivariant neural diffusion for molecule gen-
 506 eration. *arXiv preprint arXiv:2506.10532*, 2025.

507 Du, W., Zhang, H., Du, Y., Meng, Q., Chen, W., Zheng,
 508 N., Shao, B., and Liu, T.-Y. SE(3)-equivariant graph
 509 neural networks with complete local frames. In *Inter-
 510 national Conference on Machine Learning*, pp. 5583–5608.
 511 PMLR, 2022.

512 Eastman, P., Behara, P. K., Dotson, D. L., Galvelis, R., Herr,
 513 J. E., Horton, J. T., Mao, Y., Chodera, J. D., Pritchard,
 514 B. P., Wang, Y., et al. Spice, a dataset of drug-like
 515 molecules and peptides for training machine learning
 516 potentials. *Scientific Data*, 10(1):11, 2023.

517 Feng, H., Qiu, P., Zhang, M.-C., Fan, Y., Tao, Y., and Poczos,
 518 B. Soft metropolis-hastings correction for generative
 519 model sampling. *bioRxiv preprint*, 2025.

520 Feng, R., Zhu, Q., Tran, H., Chen, B., Toland, A., Ram-
 521 prasad, R., and Zhang, C. May the force be with you:
 522 Unified force-centric pre-training for 3D molecular con-
 523 formations. *arXiv preprint arXiv:2308.14759*, 2023.

524 Feng, S., Zhou, Y., Li, Y., Wang, Y., Quan, Y., Ju, P., Chen,
 525 W., Wang, Y., Zheng, M., Liu, J., et al. Unigem: A
 526 unified approach to generation and property prediction
 527 for molecules. *arXiv preprint arXiv:2410.10516*, 2024.

- 495 Fu, X., Wood, B. M., Barroso-Luque, L., Levine, D. S., Gao,
 496 M., Dzamba, M., and Zitnick, C. L. Learning smooth and
 497 expressive interatomic potentials for physical property
 498 prediction. In *Forty-second International Conference on*
 499 *Machine Learning*, 2025.
- 500 Fuchs, F., Worrall, D., Fischer, V., and Welling, M. SE(3)-
 501 transformers: 3D roto-translation equivariant attention
 502 networks. *Advances in Neural Information Processing Systems*,
 503 33:1970–1981, 2020.
- 504 Gasteiger, J., Groß, J., and Günnemann, S. Directional
 505 message passing for molecular graphs. *arXiv preprint arXiv:2003.03123*, 2020.
- 506 Gasteiger, J., Becker, F., and Günnemann, S. Gemnet: Uni-
 507 versal directional graph neural networks for molecules.
 508 *Advances in Neural Information Processing Systems*, 34:
 509 6790–6802, 2021.
- 510 Hart, P. E., Nilsson, N. J., and Raphael, B. A formal basis
 511 for the heuristic determination of minimum cost paths.
 512 *IEEE Transactions on Systems Science and Cybernetics*,
 513 4(2):100–107, 1968.
- 514 Havens, A. J., Miller, B. K., Yan, B., Domingo-Enrich,
 515 C., Sriram, A., Levine, D. S., Wood, B. M., Hu, B.,
 516 Amos, B., Karrer, B., Fu, X., Liu, G.-H., and Chen, R.
 517 T. Q. Adjoint sampling: Highly scalable diffusion sam-
 518 plers via adjoint matching. In *Forty-second International*
 519 *Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=6Eg1OrHmg2>.
- 520 Hohenberg, P. and Kohn, W. Inhomogeneous electron gas.
 521 *Physical Review*, 136(3B):B864, 1964.
- 522 Hoogeboom, E., Garcia Satorras, V., Vignac, C., and
 523 Welling, M. Equivariant diffusion for molecule gener-
 524 ation in 3D. In *International Conference on Machine*
 525 *Learning*, pp. 8867–8887. PMLR, 2022.
- 526 Isert, C., Atz, K., Jiménez-Luna, J., and Schneider, G.
 527 Qmugs, quantum mechanical properties of drug-like
 528 molecules. *Scientific Data*, 9(1):273, 2022.
- 529 Kohn, W. and Sham, L. J. Self-consistent equations includ-
 530 ing exchange and correlation effects. *Physical Review*,
 531 140(4A):A1133, 1965.
- 532 Levine, D. S., Shuaibi, M., Spotte-Smith, E. W. C., Taylor,
 533 M. G., Hasim, M. R., Michel, K., Batatia, I., Csányi,
 534 G., Dzamba, M., Eastman, P., et al. The open molecules
 535 2025 (omol25) dataset, evaluations, and models. *arXiv*
 536 *preprint arXiv:2505.08762*, 2025.
- 537 Li, Y., Wang, Y., Huang, L., Yang, H., Wei, X., Zhang, J.,
 538 Wang, T., Wang, Z., Shao, B., and Liu, T.-Y. Long-short-
 539 range message-passing: A physics-informed framework
 540 to capture non-local interaction for scalable molecular
 541 dynamics simulation. *arXiv preprint arXiv:2304.13542*,
 542 2023.
- 543 Li, Y., Xia, Z., Huang, L., Wei, X., Harshe, S., Yang, H.,
 544 Luo, E., Wang, Z., Zhang, J., Liu, C., et al. Enhancing the
 545 scalability and applicability of Kohn-Sham Hamiltonians
 546 for molecular systems. In *The Thirteenth International*
 547 *Conference on Learning Representations*, 2024.
- 548 Li, Y., Huang, L., Ding, Z., Wei, X., Wang, C., Yang, H.,
 549 Wang, Z., Liu, C., Shi, Y., Jin, P., et al. E2former:
 550 An efficient and equivariant transformer with linear-
 551 scaling tensor products. In *The Thirty-ninth Annual*
 552 *Conference on Neural Information Processing Systems*,
 553 2025. URL <https://openreview.net/pdf?id=ls5L4IMEwt>.
- 554 Liao, Y.-L. and Smidt, T. Equiformer: Equivariant graph
 555 attention transformer for 3D atomistic graphs. In *Inter-
 556 national Conference on Learning Representations*,
 557 2023. URL <https://openreview.net/forum?id=KwmPfARgOTD>.
- 558 Liu, G.-H., Choi, J., Chen, Y., Miller, B. K., and Chen, R.
 559 T. Q. Adjoint schrödinger bridge sampler. In *The Thirty-
 560 ninth Annual Conference on Neural Information Process-
 561 ing Systems*, 2025a. URL <https://openreview.net/forum?id=rMhQB1hh4c>.
- 562 Liu, J., Liu, G., Liang, J., Li, Y., Liu, J., Wang, X., Wan,
 563 P., Zhang, D., and Ouyang, W. Flow-GRPO: Training
 564 flow matching models via online RL. In *The Thirty-ninth*
 565 *Annual Conference on Neural Information Processing*
 566 *Systems*, 2025b. URL <https://openreview.net/forum?id=oCBKGw5HNF>.
- 567 Liu, Y., Wang, L., Liu, M., Lin, Y., Zhang, X., Oztekin,
 568 B., and Ji, S. Spherical message passing for 3D molec-
 569 ular graphs. In *International Conference on Learning*
 570 *Representations (ICLR)*, 2022.
- 571 Lu, C., Zhou, Y., Bao, F., Chen, J., Li, C., and Zhu, J. Dpm-
 572 solver++: Fast solver for guided sampling of diffusion
 573 probabilistic models. *arXiv preprint arXiv:2211.01095*,
 574 2022.
- 575 Mannan, S., Bihani, V., Gonzales, C., Lee, K. L. K., Gos-
 576 vami, N. N., Ranu, S., Miret, S., and Krishnan, N. Evaluat-
 577 ing universal machine learning force fields against exper-
 578 imental measurements. *arXiv preprint arXiv:2508.05762*,
 579 2025.
- 580 Musaelian, A., Batzner, S., Johansson, A., Sun, L., Owen,
 581 C. J., Kornbluth, M., and Kozinsky, B. Learning local
 582 equivariant representations for large-scale atomistic dy-
 583 namics. *Nature Communications*, 14(1):579, 2023.

- 550 Nelson, E. *Dynamical Theories of Brownian Motion*. Princeton
 551 University Press, Princeton, NJ, 1967.
- 552 Ng, A. Y., Harada, D., and Russell, S. Policy invariance
 553 under reward transformations: Theory and application
 554 to reward shaping. In *Proceedings of the International
 555 Conference on Machine Learning*, pp. 278–287, 1999.
- 556 Passaro, S. and Zitnick, C. L. Reducing SO(3) convolu-
 557 tions to SO(2) for efficient equivariant GNNs. In *Inter-
 558 national Conference on Machine Learning*, pp. 27420–
 559 27438. PMLR, 2023.
- 560 Peng, X., Guan, J., Liu, Q., and Ma, J. Moldiff: Address-
 561 ing the atom-bond inconsistency problem in 3d molecule
 562 diffusion generation. In Krause, A., Brunskill, E., Cho,
 563 K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.),
 564 *Proceedings of the 40th International Conference on Ma-
 565 chine Learning*, volume 202 of *Proceedings of Machine
 566 Learning Research*, pp. 27611–27629. PMLR, 23–29 Jul
 567 2023. URL <https://proceedings.mlr.press/v202/peng23b.html>.
- 568 Qiang, B., Song, Y., Xu, M., Gong, J., Gao, B., Zhou, H.,
 569 Ma, W.-Y., and Lan, Y. Coarse-to-fine: a hierarchical
 570 diffusion model for molecule generation in 3d. In *Inter-
 571 national Conference on Machine Learning*, pp. 28277–
 572 28299. PMLR, 2023.
- 573 Qu, E. and Krishnapriyan, A. The importance of being
 574 scalable: Improving the speed and accuracy of neural
 575 network interatomic potentials across chemical domains.
 576 *Advances in Neural Information Processing Systems*, 37:
 577 139030–139053, 2024.
- 578 Ramakrishnan, R., Dral, P. O., Rupp, M., and Von Lilienfeld,
 579 O. A. Quantum chemistry structures and properties of
 580 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- 581 Satorras, V. G., Hoogeboom, E., and Welling, M. E(n)-
 582 equivariant graph neural networks. In *International Con-
 583 ference on Machine Learning*, pp. 9323–9332. PMLR,
 584 2021.
- 585 Schütt, K. T., Arbabzadah, F., Chmiela, S., Müller, K. R.,
 586 and Tkatchenko, A. Quantum-chemical insights from
 587 deep tensor neural networks. *Nature Communications*, 8
 588 (1):1–8, 2017.
- 589 Schütt, K. T., Sauceda, H. E., Kindermans, P.-J.,
 590 Tkatchenko, A., and Müller, K.-R. Schnet—a deep learn-
 591 ing architecture for molecules and materials. *The Journal
 592 of Chemical Physics*, 148(24):241722, 2018.
- 593 Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang,
 594 H., Zhang, M., Li, Y. K., Wu, Y., and Guo, D. DeepSeek-
 595 Math: Pushing the limits of mathematical reasoning in
 596 open language models. *arXiv preprint arXiv:2402.03300*,
 597 2024.
- 598 Shen, Y., Zhang, C., Fu, S., Zhou, C., Washburn, N.,
 599 and Poczos, B. Chemistry-inspired diffusion with non-
 600 differentiable guidance. In *The Thirteenth International
 601 Conference on Learning Representations*, 2024.
- 602 Simeon, G. and De Fabritiis, G. Tensornet: Cartesian ten-
 603 sor representations for efficient learning of molecular
 604 potentials. *Advances in Neural Information Processing
 605 Systems*, 36, 2024.
- 606 Song, Y., Gong, J., Xu, M., Cao, Z., Lan, Y., Ermon, S.,
 607 Zhou, H., and Ma, W.-Y. Equivariant flow matching with
 608 hybrid probability transport for 3d molecule generation.
 609 *Advances in Neural Information Processing Systems*, 36:
 610 549–568, 2023.
- 611 Song, Y., Gong, J., Zhou, H., Zheng, M., Liu, J., and Ma,
 612 W.-Y. Unified generative modeling of 3D molecules via
 613 Bayesian Flow Networks. In *International Conference
 614 on Learning Representations*, 2024.
- 615 Spall, J. C. Multivariate stochastic approximation using a
 616 simultaneous perturbation gradient approximation. *IEEE
 617 Transactions on Automatic Control*, 37(3):332–341, 1992.
- 618 Sriram, A., Choi, S., Yu, X., Brabson, L. M., Das, A., Ulissi,
 619 Z., Uyttendaele, M., Medford, A. J., and Sholl, D. S.
 620 The open dac 2023 dataset and challenges for sorbent
 621 discovery in direct air capture, 2024.
- 622 Sun, Q., Zhang, X., Banerjee, S., Bao, P., Barbry, M., Blunt,
 623 N. S., Bogdanov, N. A., Booth, G. H., Chen, J., Cui,
 624 Z.-H., et al. Recent developments in the pyscf program
 625 package. *The Journal of chemical physics*, 153(2), 2020.
- 626 Tran, R., Lan, J., Shuaibi, M., Wood, B. M., Goyal, S., Das,
 627 A., Heras-Domingo, J., Kolluru, A., Rizvi, A., Shoghi,
 628 N., et al. The open catalyst 2022 (oc22) dataset and
 629 challenges for oxide electrocatalysts. *ACS Catalysis*, 13
 630 (5):3066–3084, 2023.
- 631 Unke, O. T., Chmiela, S., Sauceda, H. E., Gastegger, M.,
 632 Poltavsky, I., Schütt, K. T., Tkatchenko, A., and Müller,
 633 K.-R. Machine learning force fields. *Chemical Reviews*,
 634 121(16):10142–10186, 2021.
- 635 Unke, O. T., Stöhr, M., Ganscha, S., Unterthiner, T.,
 636 Maennel, H., Kashubin, S., Ahlin, D., Gastegger, M.,
 637 Medrano Sandonas, L., Berryman, J. T., et al. Biomolecul-
 638 lar dynamics with machine-learned quantum-mechanical
 639 force fields trained on diverse chemical fragments. *Sci-
 640 ence Advances*, 10(14):eadn4397, 2024.
- 641 Wang, L., Liu, Y., Lin, Y., Liu, H., and Ji, S. ComENet: To-
 642 wards complete and efficient message passing for 3D
 643 molecular graphs. In Oh, A. H., Agarwal, A., Bel-
 644 grave, D., and Cho, K. (eds.), *Advances in Neural In-
 645 formation Processing Systems*, 2022. URL <https://openreview.net/forum?id=mCzMqeWSFJ>.

605 Wood, B. M., Dzamba, M., Fu, X., Gao, M., Shuaibi, M.,
606 Barroso-Luque, L., Abdelmaqsoud, K., Gharakhanyan,
607 V., Kitchin, J. R., Levine, D. S., et al. Uma: A
608 family of universal models for atoms. *arXiv preprint*
609 *arXiv:2506.23971*, 2025.

610 Wu, L., Gong, C., Liu, X., Ye, M., and Liu, Q. Diffusion-
611 based molecule generation with informative prior bridges.
612 In *Advances in Neural Information Processing Systems*,
613 volume 35, pp. 36533–36545, 2022.

614 Xu, M., Powers, A., Dror, R., Ermon, S., and Leskovec, J.
615 Geometric latent diffusion models for 3D molecule gener-
616 ation. In *International Conference on Machine Learning*,
617 pp. 38592–38610. PMLR, 2023.

618 Xue, Z., Wu, J., Gao, Y., Kong, F., Zhu, L., Chen, M., Liu,
619 Z., Liu, W., Guo, Q., Huang, W., and Luo, P. Dancegrpo:
620 Unleashing grp on visual generation. *arXiv preprint*
621 *arXiv:2505.07818*, 2025.

622 Zaidi, S., Schaarschmidt, M., Martens, J., Kim, H., Teh,
623 Y. W., Sanchez-Gonzalez, A., Battaglia, P., Pascanu, R.,
624 and Godwin, J. Pre-training via denoising for molecu-
625 lar property prediction. In *The Eleventh International*
626 *Conference on Learning Representations*, 2022.

627 Zhou, Z., An, J., Liu, Z., Shi, Y., Zhang, X., Cao, F., Qu, C.,
628 and Qi, Y. Guiding diffusion models with reinforcement
629 learning for stable molecule generation. *arXiv preprint*
630 *arXiv:2508.16521*, 2025.

Appendix Table of Contents

660		
661		
662		
663		
664		
665	A Hyperparameters	14
666	A.1 RL post-training	14
667	A.2 QM9	15
668	A.3 GEOM-Drugs	17
669		
670		
671		
672	B Notations	17
673		
674	C Related Works	18
675		
676	D Proofs for Section 3.2	19
677		
678		
679	E Supplementary Experiments	21
680	E.1 Qualitative analysis.	21
681	E.2 Property controllability and distributional fidelity.	23
682		
683	E.3 More Ablation Studies	24
684		
685		
686	F Supplementary Theory: PBRS as an Alchemical Force	24
687	F.1 Alchemical force analysis	25
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		
702		
703		
704		
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		

A. Hyperparameters

A.1. RL post-training

Table 4. Hyperparameters for RL postraining

Notation	Description	Value
learning_rate	Learning rate for policy optimization	2e-7 to 4e-6
clip_range	PPO clipping threshold	2e-3
train_micro_batch_size	Micro-batch size for policy updates	4
epoch_per_rollout	Number of optimization epochs per rollout	{1, 3}
kl_penalty_weight	Weight of KL regularization that constrains the policy update to stay close to the reference model	0.04
discount	Discount factor used in calculating intermediate shaping rewards	1
weight_energy	Weighting factor for energy reward	{0.05 0.1}
weight_force	Weighting factor for force reward	1
sample_group_size	Number of sample groups generated per prompt	4
each_prompt_sample	Number of trajectories generated per prompt	24
time_step	Number of diffusion timesteps used during sampling and rollout	1000
mlff_model	Identifier of the machine-learned force field (MLFF) model	uma-m-1p1 or uma-s-1p1
mlff_batch_size	Batch size for reward evaluation	8
force_aggregation	Aggregation method for force errors	rms or max
skip_prefix	Number of initial steps skipped when applying reward shaping	{600,700}
scheduler_name	Learning rate scheduler type	cosine
scheduler_warmup_steps	Number of warm-up steps during which the learning rate increases from zero to the initial value	60
scheduler_total_steps	Total number of scheduler steps before learning rate decay completes	1500
scheduler_min_lr_ratio	Final learning rate expressed as a fraction of the initial learning rate after decay	0.3
KL_weight	KL divergence with respect to a reference pretrained policy	{0, 0.02 }

A.2. QM9

Table 5. Hyperparameters used for the pretrained model on QM9 molecular dataset

Notation	Description	Value
model	Dynamics model type	egnn_dynamics
probabilistic_model	Probabilistic model type	diffusion
diffusion_steps	Number of diffusion steps	500
diffusion_noise_schedule	Noise schedule used in diffusion	polynomial_2
diffusion_loss_type	Loss function for diffusion training	vlb, l2
diffusion_noise_precision	Minimum diffusion noise precision used for numerical stability	1e-5
n_epochs	Total number of training epochs	200
batch_size	Batch size used for training	128
lr	Learning rate for the optimizer	0.0002
n_layers	Number of EGNN layers	6
inv_sublayers	Number of invariant sublayers per EGNN layer	1
nf	Hidden feature dimension (node feature size)	128
tanh	Whether to use tanh activation in coordinate MLPs	True
attention	Whether to use attention mechanisms in EGNN layers	True
norm_constant	Normalization constant in coordinate updates	1
sin_embedding	Whether to use sinusoidal time embeddings	False
ode_regularization	Regularization strength for ODE-based formulations	1e-3
dataset	Dataset to use	qm9 or qm9_second_half
filter_n_atoms	Restrict dataset to molecules with a fixed number of atoms	None
dequantization	Dequantization strategy	argmax_variational
ema_decay	Amount of Exponential Moving Average (EMA) decay, 0 means off, a reasonable value is 0.999	0.999
n_stability_samples	Number of samples to compute the stability	500
normalize_factors	Normalization factors for x, categorical, and integer features	1, 4, 1
include_charges	Whether to include atom charge	True
normalization_factor	Normalize the sum aggregation of EGNN	1
aggregation_method	Aggregation for the graph network	sum or mean

Table 6. Hyperparameters used for the pretrained model on the GEOM Drugs dataset

Notation	Description	Value
model	Dynamics model type	egnn_dynamics
probabilistic_model	Probabilistic model type	diffusion
diffusion_steps	Number of diffusion steps	500
diffusion_noise_schedule	Noise schedule used in diffusion	polynomial_2
diffusion_loss_type	Loss function for diffusion training	l2
diffusion_noise_precision	Minimum diffusion noise precision used for numerical stability	1e-5
n_epochs	Total number of training epochs	10000
batch_size	Batch size used for training	64
lr	Learning rate for the optimizer	5e-5
n_layers	Number of EGNN layers	6
inv_sublayers	Number of invariant sublayers per EGNN layer	1
nf	Hidden feature dimension (node feature size)	192
tanh	Whether to use tanh activation in coordinate MLPs	True
attention	Whether to use attention mechanisms in EGNN layers	True
norm_constant	Normalization constant in coordinate updates	1
sin_embedding	Whether to use sinusoidal time embeddings	False
ode_regularization	Regularization strength for ODE-based formulations	1e-3
dataset	Dataset to use	geom
filter_n_atoms	Restrict dataset to molecules with a fixed number of atoms	None
dequantization	Dequantization strategy	argmax-variational
ema_decay	Amount of Exponential Moving Average (EMA) decay, 0 means off, a reasonable value is 0.999	0
n_stability_samples	Number of samples to compute the stability	20
normalize_factors	Normalization factors for x, categorical, and integer features	1, 4, 10
include_charges	Whether to include atom charge	False
normalization_factor	Normalize the sum aggregation of EGNN	100
aggregation_method	Aggregation for the graph network	sum or mean

A.3. GEOM-Drugs
B. Notations

Symbol	Type	Description
Geometric states and molecular structure		
\mathbf{x}_0	$\mathbb{R}^{N \times 3}$	Clean atomic coordinates
\mathbf{x}_t	$\mathbb{R}^{N \times 3}$	Noisy conformation at diffusion time t
N	\mathbb{N}	Number of atoms
\mathcal{X}	space	Configuration space of clean (CoM free) molecular structures
z	\mathcal{X}	Generic clean structure (integration variable in appendix proofs)
\mathbf{z}_t	$\mathbb{R}^{N \times 3} \times \mathbb{R}^{N \times d_h}$	Diffusion state (coordinates plus hidden features)
Diffusion process and generative model		
$p_t(\mathbf{z})$	distribution	Marginal distribution at time t
$p_\theta(\mathbf{z}_{t-1} \mathbf{z}_t)$	distribution	Reverse diffusion transition
$s_\theta(\mathbf{z}_t, t)$	$\mathbb{R}^{N \times 3}$	Score function $\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t)$
ϵ	$\mathbb{R}^{N \times 3}$	Standard Gaussian noise
α_t	\mathbb{R}	Signal coefficient
σ_t	$\mathbb{R}_{>0}$	Noise scale
Trajectories and expectations		
τ	sequence	Reverse diffusion trajectory $(\mathbf{z}_T, \dots, \mathbf{z}_0)$
\mathbb{E}	operator	Expectation over samples or trajectories
\mathcal{T}	set	Set of sampled trajectories
Energy, forces, and physical models		
$E(z)$	\mathbb{R}	Potential energy of a clean structure z
$\mathbf{F}(z)$	$\mathbb{R}^{N \times 3}$	Force field $-\nabla_z E(z)$
$E_*(z)$	\mathbb{R}	Target energy (gold standard reference)
$E_\phi(z)$	\mathbb{R}	MLFF energy approximation
$\mathbf{F}_*(z)$	$\mathbb{R}^{N \times 3}$	Target forces $-\nabla_z E_*(z)$
$\mathbf{F}_\phi(z)$	$\mathbb{R}^{N \times 3}$	MLFF forces $-\nabla_z E_\phi(z)$
Reinforcement learning and FED GRPO		
π_θ	policy	Reverse diffusion policy
r_t^E	\mathbb{R}	Energy based reward
r_t^F	\mathbb{R}	Force based stability reward
r_t	\mathbb{R}	Total reward
\hat{A}_t	\mathbb{R}	Advantage estimate
$\mathcal{L}_{\text{GRPO}}$	objective	Group Relative Policy Optimization loss
$\mathcal{L}_{\text{FED-GRPO}}$	objective	Force Energy Disentangled GRPO loss
Measure theoretic and variational notation (appendix proofs)		
$\mathcal{P}(\mathcal{X})$	space	Probability measures on \mathcal{X}
μ	$\mathcal{P}(\mathcal{X})$	Pretrained terminal law, $\mu := \rho_{\theta_{\text{pre}}}$
ρ	$\mathcal{P}(\mathcal{X})$	Candidate terminal law
$\rho \ll \mu$	relation	Absolute continuity of ρ with respect to μ
$\text{KL}(\rho \ \mu)$	$\mathbb{R}_{\geq 0}$	Kullback Leibler divergence
$\mathcal{J}(\rho)$	functional	Variational objective over terminal laws
w_E	$\mathbb{R}_{>0}$	Energy weight in advantage calculation
w_{KL}	$\mathbb{R}_{>0}$	KL weight in \mathcal{J}
β_{eff}	$\mathbb{R}_{>0}$	Effective inverse temperature, $\beta_{\text{eff}} := w_E / w_{\text{KL}}$
ρ^*	$\mathcal{P}(\mathcal{X})$	Maximizer of \mathcal{J}
Z_ϕ	$\mathbb{R}_{>0}$	Partition function for ρ^* under E_ϕ
Z_*	$\mathbb{R}_{>0}$	Partition function for the tilt under E_*
f	density	Radon Nikodym derivative $f := d\rho/d\mu$
λ	\mathbb{R}	Lagrange multiplier for normalization of f
$\ \cdot\ _{\text{TV}}$	$\mathbb{R}_{\geq 0}$	Total variation distance
δ	$\mathbb{R}_{\geq 0}$	Uniform energy error bound, $\sup_{z \in \mathcal{X}} E_\phi(z) - E_*(z) \leq \delta$
Symmetry groups and spaces		
$E(3)$	group	Euclidean group of rotations and translations
$\text{SO}(3)$	group	Rotation group
$\mathbb{R}^{N \times 3}$	space	Cartesian coordinate space

Table 7. Summary of symbols and notation used throughout the paper and appendix.

935 C. Related Works

936 **Diffusion Models for Molecules.** E(3)-equivariant diffusion models have quickly become a leading paradigm for 3D
 937 molecular generation, as they inherently enforce rotational and translational symmetries and produce geometries with
 938 high fidelity. (Hoogeboom et al., 2022) introduced an SE(3)-equivariant diffusion process that generates molecules as
 939 sets of atoms with remarkable realism. Building on this foundation, (Xu et al., 2023) developed GeoLDM, a latent 3D
 940 diffusion approach that operates in a learned point-cloud latent space composed of both invariant scalars and equivariant
 941 tensors, improving sample efficiency and controllable generation of molecular geometries. More recently, (Cornet et al.,
 942 2025) proposed an Equivariant Neural Diffusion (END) model that also learns the forward-noising process, achieving
 943 state-of-the-art results on standard molecular generation benchmarks. Beyond unconditional generation, many methods aim
 944 to guide or bias diffusion models toward molecules with desired properties or constraints. One strategy is to inject domain
 945 knowledge via physically-informed priors. (Wu et al., 2022) introduces informative prior bridges that steer the diffusion
 946 trajectory toward low-energy conformations by constructing a biased intermediate distribution for the sampler. (Zhou et al.,
 947 2025) formulates 3D molecule generation as a Markov decision process and uses proximal policy optimization to fine-tune a
 948 pre-trained equivariant diffusion model’s sampling policy by providing a reward based on physics evaluations. An alternative
 949 formalism is to cast guided generation as a stochastic optimal control problem. (Havens et al., 2025) introduces Adjoint
 950 Sampling, which optimizes a matching objective for the diffusion’s probability flow by learning an optimal diffusion drift
 951 that samples from an unnormalized target density, such as a Boltzmann distribution. (Liu et al., 2025a) generalizes this idea
 952 by relaxing previous assumptions on the prior, solving a Schrödinger Bridge between a simple base distribution and the
 953 desired Boltzmann-like distribution.

955 **Machine-learning Force-Field.** Machine learning force fields (MLFFs) use geometric deep learning to approximate
 956 potential energy surfaces and interatomic forces with high accuracy at far lower cost than quantum chemistry. Early MLFFs
 957 relied on invariant message passing networks that predict an energy scalar invariant to rotations and translations, with
 958 forces obtained via differentiation. SchNet established this paradigm using continuous filter convolutions over interatomic
 959 distances (Schütt et al., 2018). More generally, graph network formalisms demonstrated that message passing can unify
 960 molecular and crystalline property prediction under a single framework (Chen et al., 2019). A major line of progress then
 961 focused on injecting richer geometric structure into invariant architectures. DimeNet introduced directional message passing
 962 through angular features and specialized basis expansions (Gasteiger et al., 2020), and GemNet improved higher order
 963 interaction modeling to better capture many-body effects relevant for energies and forces (Gasteiger et al., 2021). SphereNet
 964 further leveraged a spherical coordinate parameterization for local neighborhoods to encode distances and angles in a
 965 unified way (Liu et al., 2022), while ComENet emphasized complete yet efficient geometric message passing to improve
 966 accuracy without sacrificing scaling (Wang et al., 2022). Recent analyses underscore that architecture and training data
 967 choices critically determine transfer across chemical domains and system sizes, motivating careful study of inductive
 968 biases and data coverage (Qu & Krishnapriyan, 2024). In parallel, E(3) equivariant MLFFs enforce symmetry directly in
 969 intermediate representations, enabling vector and tensor features that transform correctly under rotations and translations.
 970 EGNN provides a lightweight recipe for equivariant message passing with coordinate updates (Satorras et al., 2021), and
 971 complete local frame constructions strengthen SE(3) equivariance and stability by anchoring messages in local geometric
 972 frames (Du et al., 2022). Tensor-based designs such as TensorNet offer an efficient Cartesian tensor alternative to spherical
 973 harmonic pipelines while preserving equivariant structure (Simeon & De Fabritiis, 2024), and newer architectures focus
 974 explicitly on practical efficiency, such as GotenNet (Aykent & Xia, 2025). Equivariant graph networks have also been used
 975 as surrogates for more demanding electronic structure components, improving the scalability of learned approximations
 976 to Kohn-Sham DFT quantities while maintaining symmetry constraints (Li et al., 2024). A particularly influential family
 977 of MLFFs builds on SO(3) representation theory with spherical harmonics. SE(3) Transformer introduced equivariant
 978 attention with irreducible representations (Fuchs et al., 2020), and Equiformer extended this transformer style approach with
 979 strong results on quantum chemistry benchmarks (Liao & Smidt, 2023). NequIP demonstrated striking data efficiency for
 980 learning interatomic potentials with spherical tensor message passing (Batzner et al., 2022), while MACE advanced higher
 981 order equivariant interactions for accurate molecular dynamics potentials (Batatia et al., 2022). Scalability has been further
 982 improved with locally deployed equivariant potentials such as Allegro (Musaelian et al., 2023), extensions that model both
 983 long range and short range interactions (Li et al., 2023), and methods that reduce the cost of equivariant operations (Passaro
 984 & Zitnick, 2023; Li et al., 2025). Finally, recent work targets improved physical consistency and broader transfer, including
 985 learning objectives tuned for stable molecular dynamics and phonon properties (Fu et al., 2025) and universal pretrained
 986 atomistic models that aim to generalize across elements and compounds (Wood et al., 2025).

990 **Reinforcement Learning for Diffusion Guidance.** Integrating RL with diffusion models is a recent trend aimed at
 991 steering generative models toward complex objectives that are hard to encode in the training likelihood. One representative
 992 approach is Denoising Diffusion Policy Optimization (DDPO) (Black et al., 2024). DDPO reframes the multi-step denoising
 993 process as a Markov decision process, where each diffusion step is an action, and a reward is obtained at the final sample.
 994 This allows applying policy gradient algorithms to fine-tune a pretrained diffusion model for higher rewards. (Xue et al.,
 995 2025) adapts GRPO to diffusion and continuous flow models with *DanceGRPO*, achieving robust RL-based fine-tuning
 996 on large-scale text-to-image and text-to-video generation tasks. (Liu et al., 2025b) introduces *Flow-GRPO*, an online RL
 997 algorithm that trains flow matching models using reward signals by aligning generated trajectories with desired outcomes
 998 through policy gradient updates. (Shen et al., 2024) proposes a chemistry-inspired guided diffusion that uses external
 999 force-field evaluations as a form of reward to bias the diffusion sampler toward low-energy structures, without requiring
 1000 those evaluations to be differentiable. (Zhou et al., 2025) applies Reinforcement Learning with Physical Feedback (RLPF)
 1001 to an equivariant diffusion model, significantly improving stability of generated conformations. Crucially, both works
 1002 demonstrate that RL can incorporate domain-specific criteria (like force-field energies) that lie outside the scope of the
 1003 original generative model’s training. Another perspective comes from control theory: *Adjoint methods* avoid explicit RL by
 1004 solving optimal control formulations for diffusion. (Havens et al., 2025) derives an adjoint formulation to directly match a
 1005 diffusion sampler to an optimal policy in a single backward pass, while (Liu et al., 2025a) develops an adjoint Schrödinger
 1006 bridge that computes an optimal transport between the prior and target distributions. These methods achieve goals similar
 1007 to RL-guided diffusion by biasing generation toward certain distributions or rewards, but do so by analytically aligning
 1008 sampling trajectories, thereby improving efficiency and scalability.

D. Proofs for Section 3.2

1012 **Notation.** Let \mathcal{X} denote the configuration space of clean (CoM-free) molecular structures. Let $\mu := \rho_{\theta_{\text{pre}}} \in \mathcal{P}(\mathcal{X})$ be the
 1013 pretrained terminal law. For any $\rho \in \mathcal{P}(\mathcal{X})$, write $\rho \ll \mu$ for absolute continuity. We use $\text{KL}(\rho \| \mu)$ for Kullback–Leibler
 1014 divergence. Throughout, we fix the energy weight $w_E = 1$ and define

$$\beta_{\text{eff}} := \frac{1}{w_{\text{KL}}}.$$

1018 (Equivalently, any non-unit w_E can be absorbed into w_{KL} by rescaling.)

1019 **Theorem 3** (Energy-aligned terminal distribution (restated)). Assume $w_{\text{KL}} > 0$ and consider the functional over $\rho \in \mathcal{P}(\mathcal{X})$,
 1020

$$\mathcal{J}(\rho) := \mathbb{E}_{z \sim \rho}[-E_\phi(z)] - w_{\text{KL}} \text{KL}(\rho \| \mu), \quad \text{with the convention } \mathcal{J}(\rho) = -\infty \text{ if } \rho \not\ll \mu.$$

1022 Assume the supremum of \mathcal{J} over $\mathcal{P}(\mathcal{X})$ is attained. Then the maximizer is unique and equals

$$\rho^*(dz) = \frac{1}{Z_\phi} \exp(-\beta_{\text{eff}} E_\phi(z)) \mu(dz), \quad Z_\phi := \int_{\mathcal{X}} \exp(-\beta_{\text{eff}} E_\phi(z)) \mu(dz).$$

1024 *Proof.* Let $\rho \in \mathcal{P}(\mathcal{X})$. If $\rho \not\ll \mu$, then $\text{KL}(\rho \| \mu) = +\infty$ and $\mathcal{J}(\rho) = -\infty$, so we restrict to $\rho \ll \mu$. Write the
 1025 Radon–Nikodym derivative $f := d\rho/d\mu$. Then $f \geq 0$ μ -a.e. and $\int f d\mu = 1$. Moreover,

$$\text{KL}(\rho \| \mu) = \int_{\mathcal{X}} f(z) \log f(z) \mu(dz), \quad \mathbb{E}_{z \sim \rho}[E_\phi(z)] = \int_{\mathcal{X}} E_\phi(z) f(z) \mu(dz).$$

1026 Thus maximizing $\mathcal{J}(\rho)$ over $\rho \ll \mu$ is equivalent to maximizing over densities f :

$$\mathcal{J}(f) = - \int_{\mathcal{X}} E_\phi(z) f(z) \mu(dz) - w_{\text{KL}} \int_{\mathcal{X}} f(z) \log f(z) \mu(dz), \quad \text{s.t. } \int_{\mathcal{X}} f(z) \mu(dz) = 1, \quad f \geq 0.$$

1027 Introduce a Lagrange multiplier $\lambda \in \mathbb{R}$ for the normalization constraint and define

$$\mathcal{L}(f, \lambda) = - \int E_\phi(z) f(z) \mu(dz) - w_{\text{KL}} \int f(z) \log f(z) \mu(dz) + \lambda \left(\int f(z) \mu(dz) - 1 \right).$$

1028 Let h be any bounded perturbation with $\int h d\mu = 0$ and consider $f_\varepsilon = f + \varepsilon h$. A first-order stationarity condition at an
 1029 optimum f^* implies $\frac{d}{d\varepsilon} \mathcal{L}(f^* + \varepsilon h, \lambda)|_{\varepsilon=0} = 0$ for all such h . Using $\frac{d}{d\varepsilon} [(f^* + \varepsilon h) \log(f^* + \varepsilon h)]|_{\varepsilon=0} = h(1 + \log f^*)$, we
 1030 obtain

$$0 = \int_{\mathcal{X}} h(z) \left(-E_\phi(z) - w_{\text{KL}}(1 + \log f^*(z)) + \lambda \right) \mu(dz).$$

1045 Since this holds for all h with zero μ -mean, the bracketed term must be μ -a.e. constant, i.e.

$$1046 \quad -E_\phi(z) - w_{\text{KL}}(1 + \log f^*(z)) + \lambda = 0 \quad \text{for } \mu\text{-a.e. } z.$$

1048 Rearranging gives

$$1050 \quad \log f^*(z) = -\frac{1}{w_{\text{KL}}} E_\phi(z) + c = -\beta_{\text{eff}} E_\phi(z) + c$$

1052 for some constant $c \in \mathbb{R}$, hence

$$1053 \quad f^*(z) = \exp(c) \exp(-\beta_{\text{eff}} E_\phi(z)).$$

1054 Imposing $\int f^* d\mu = 1$ yields $\exp(c) = 1/Z_\phi$ with

$$1056 \quad Z_\phi = \int_{\mathcal{X}} \exp(-\beta_{\text{eff}} E_\phi(z)) \mu(dz),$$

1059 so $\rho^*(dz) = f^*(z)\mu(dz)$ as claimed.

1060 Uniqueness follows because $f \mapsto \int f \log f d\mu$ is strictly convex on $\{f \geq 0 : \int f d\mu = 1\}$, hence $\mathcal{J}(f)$ is strictly concave
1061 when $w_{\text{KL}} > 0$. \square

1063 **MLFF approximation and distribution error.** Let $E_\star : \mathcal{X} \rightarrow \mathbb{R}$ be a target energy and E_ϕ its MLFF approximation.
1064 Define the two tilted laws (with the same reference μ):

$$1066 \quad \rho_\star^*(dz) = \frac{1}{Z_\star} \exp(-\beta_{\text{eff}} E_\star(z)) \mu(dz), \quad \rho_\phi^*(dz) = \frac{1}{Z_\phi} \exp(-\beta_{\text{eff}} E_\phi(z)) \mu(dz),$$

1069 where

$$1070 \quad Z_\star := \int_{\mathcal{X}} \exp(-\beta_{\text{eff}} E_\star(z)) \mu(dz), \quad Z_\phi := \int_{\mathcal{X}} \exp(-\beta_{\text{eff}} E_\phi(z)) \mu(dz).$$

1072 **Theorem 4** (Energy error implies distribution error (restated)). Assume $\sup_{z \in \mathcal{X}} |E_\phi(z) - E_\star(z)| \leq \delta$ for some $\delta \geq 0$. Then

$$1074 \quad \|\rho_\phi^* - \rho_\star^*\|_{\text{TV}} \leq \tanh(\beta_{\text{eff}} \delta).$$

1076 **Lemma 1** (Likelihood-ratio bound implies TV bound). Let P, Q be probability measures with $P \ll Q$ and suppose there
1077 exists $\varepsilon \geq 0$ such that

$$1078 \quad e^{-\varepsilon} \leq \frac{dP}{dQ} \leq e^\varepsilon \quad Q\text{-a.e.}$$

1080 Then $\|P - Q\|_{\text{TV}} \leq \tanh(\varepsilon/2)$.

1082 *Proof.* Let $r := dP/dQ$, so $r \in [e^{-\varepsilon}, e^\varepsilon]$ Q -a.e. and $\mathbb{E}_Q[r] = 1$. Let $A := \{r \geq 1\}$. Then

$$1084 \quad \|P - Q\|_{\text{TV}} = \frac{1}{2} \int |dP - dQ| = \frac{1}{2} \int |r - 1| dQ = \int_A (r - 1) dQ,$$

1087 where the last equality uses $\int_A (r - 1) dQ = \int_{A^c} (1 - r) dQ$ from $\mathbb{E}_Q[r] = 1$.

1088 Write $q := Q(A)$ and the conditional means $m_+ := \mathbb{E}_Q[r \mid A] \in [1, e^\varepsilon]$, $m_- := \mathbb{E}_Q[r \mid A^c] \in [e^{-\varepsilon}, 1]$. The constraint
1089 $\mathbb{E}_Q[r] = 1$ becomes $qm_+ + (1 - q)m_- = 1$, and the TV becomes $q(m_+ - 1)$. The maximum TV under the bounds is
1090 achieved at the extreme values $m_+ = e^\varepsilon$ and $m_- = e^{-\varepsilon}$. Solving

$$1092 \quad qe^\varepsilon + (1 - q)e^{-\varepsilon} = 1 \implies q = \frac{1}{e^\varepsilon + 1},$$

1094 we get

$$1096 \quad \|P - Q\|_{\text{TV}} \leq q(e^\varepsilon - 1) = \frac{e^\varepsilon - 1}{e^\varepsilon + 1} = \tanh(\varepsilon/2).$$

\square

1100 *Proof of Theorem 4.* Let $\beta := \beta_{\text{eff}}$. Define $w_\star(z) := \exp(-\beta E_\star(z))$ and $w_\phi(z) := \exp(-\beta E_\phi(z))$. By $|E_\phi(z) - E_\star(z)| \leq$
 1101 δ , for all $z \in \mathcal{X}$,

$$e^{-\beta\delta} w_\star(z) \leq w_\phi(z) \leq e^{\beta\delta} w_\star(z).$$

1103 Integrating w.r.t. μ yields the partition-function sandwich:

$$e^{-\beta\delta} Z_\star \leq Z_\phi \leq e^{\beta\delta} Z_\star.$$

1107 Therefore, for ρ_\star^* -a.e. z ,

$$\frac{d\rho_\phi^*}{d\rho_\star^*}(z) = \frac{Z_\star}{Z_\phi} \cdot \frac{w_\phi(z)}{w_\star(z)} \in [e^{-\beta\delta}, e^{\beta\delta}] \cdot [e^{-\beta\delta}, e^{\beta\delta}] = [e^{-2\beta\delta}, e^{2\beta\delta}].$$

1112 Applying Lemma 1 with $P = \rho_\phi^*$, $Q = \rho_\star^*$ and $\varepsilon = 2\beta\delta$ gives

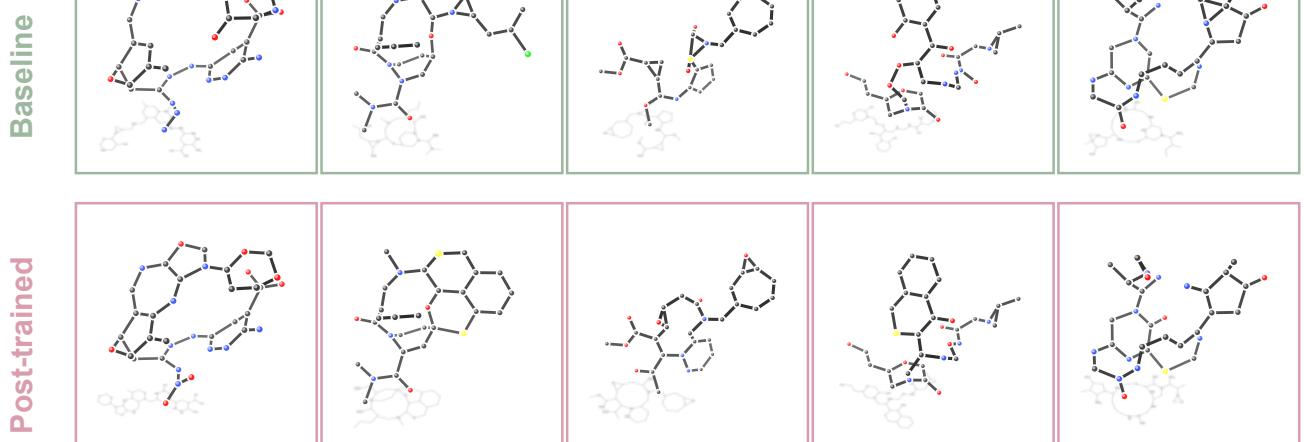
$$\|\rho_\phi^* - \rho_\star^*\|_{\text{TV}} \leq \tanh((2\beta\delta)/2) = \tanh(\beta\delta) = \tanh(\beta_{\text{eff}}\delta).$$

□

E. Supplementary Experiments

E.1. Qualitative analysis.

1121 Figures 5–8 provide a qualitative view of how post-training changes the generated 3D geometries on GEOM-Drugs. In
 1122 Figure 5, we overlay generated conformers (colored) with the reference structure (gray), highlighting that the post-trained
 1123 sampler more consistently produces a coherent, well-connected geometry that visually matches the target conformer.
 1124 Figure ?? complements this visualization by annotating the same examples with MLFF-predicted energies and RMS
 1125 force norms: the post-trained samples typically shift toward lower energies and smaller residual forces, consistent with
 1126 configurations closer to local equilibrium. Finally, Figure 8 compares matched-noise reverse trajectories (same initial
 1127 noise) and shows that the post-trained model corrects failure modes that persist throughout denoising in the baseline (e.g.,
 1128 distorted rings/strained geometries), yielding a chemically valid final structure at $t = 0$. Together, these figures illustrate that
 1129 alignment not only improves terminal metrics but also changes the denoising dynamics in a way that stabilizes intermediate
 1130 geometries.



1149 *Figure 5.* Qualitative comparison of molecular conformers generated by the baseline diffusion model (top row, green) and the post-trained
 1150 model (bottom row, pink). Each panel displays a generated conformer (dark atoms and bonds) overlaid on the ground-truth reference
 1151 structure (light gray). The post-trained model produces conformers with improved structural fidelity to the reference geometries.

1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179

Conformer Quality Comparison

Baseline
Post-trained

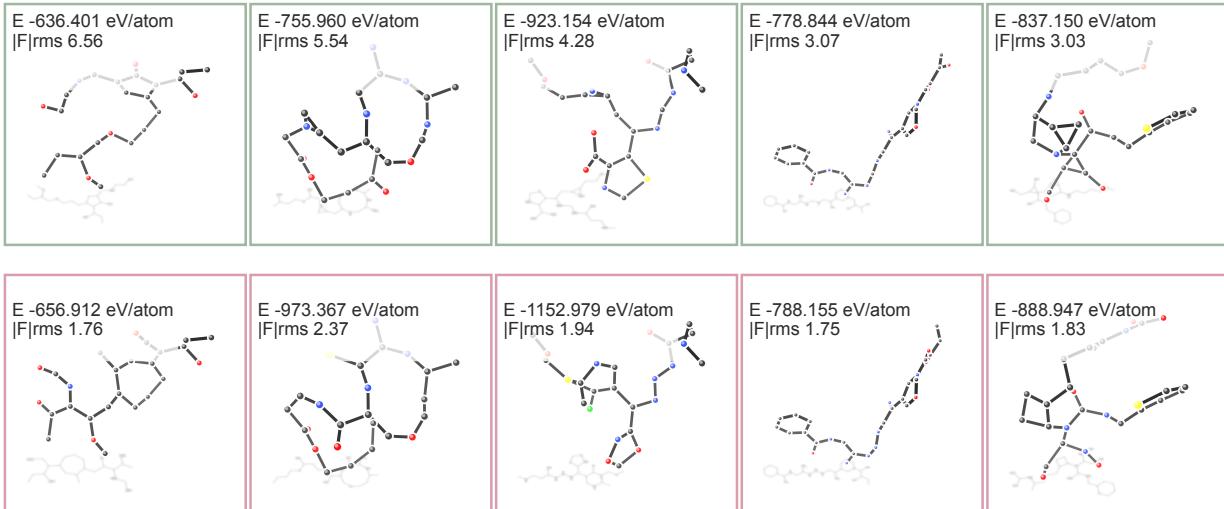


Figure 6. Conformer quality comparison showing energy and force metrics. Top row: baseline model generates structures with higher energies (E) and larger RMS forces, indicating geometries farther from equilibrium. Bottom row: post-trained model consistently achieves lower energies and reduced force magnitudes across all examples, demonstrating improved physical plausibility. Generated conformers (dark) are overlaid on reference structures (light gray).

1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209

Conformer quality Evaluation II

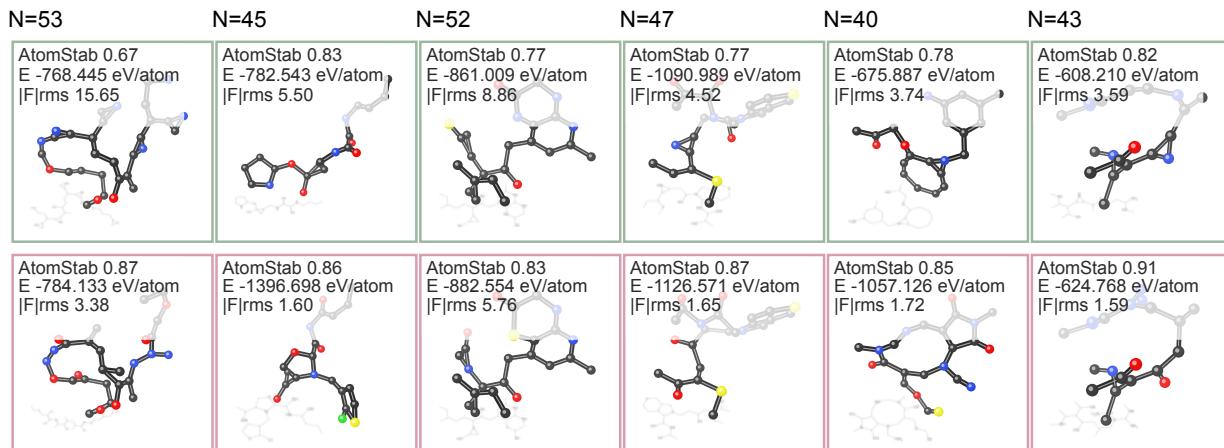


Figure 7. Conformer quality comparison showing energy and force metrics. Top row: baseline model generates structures with higher energies (E) and larger RMS forces, indicating geometries farther from equilibrium. Bottom row: post-trained model consistently achieves lower energies and reduced force magnitudes across all examples, demonstrating improved physical plausibility. Generated conformers (dark) are overlaid on reference structures (light gray).

1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224

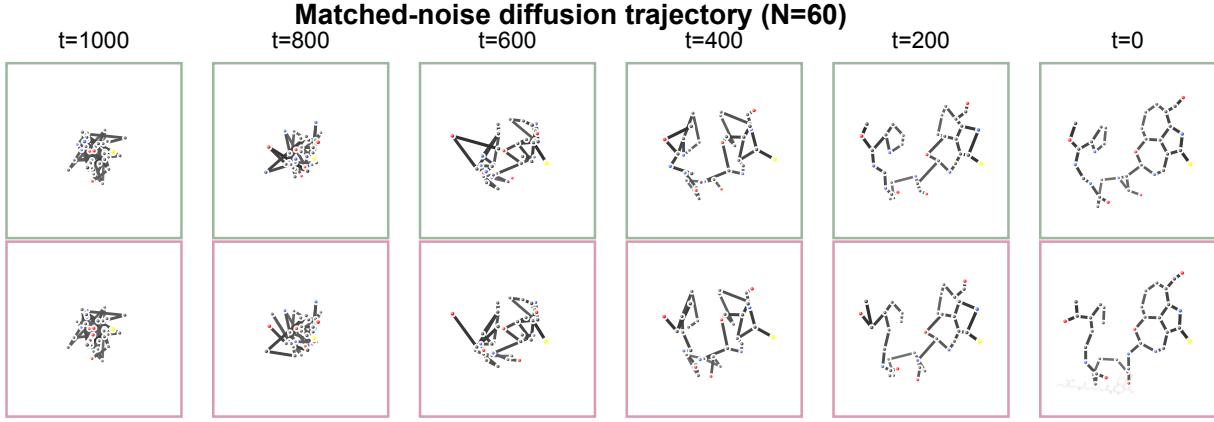


Figure 8. Matched-noise diffusion trajectories comparing the baseline (top row, green) and post-trained (bottom row, pink) models starting from identical initial noise ($N = 60$ atoms, seed=89). Snapshots are shown at diffusion timesteps $t \in \{1000, 800, 600, 400, 200, 0\}$. The baseline model fails to generate a chemically valid molecule, producing a final structure ($t = 0$) with distorted ring geometries and strained bond angles that violate standard valence constraints. In contrast, the post-trained model successfully generates a valid conformer with proper ring planarity and realistic bond geometries, demonstrating improved structural validity through physics-informed fine-tuning.

E.2. Property controllability and distributional fidelity.

We split the QM9 training set into two disjoint halves with 50k samples each. We train the property prediction network ω on the first half, and train conditional generative models on the second half. For a range of target property values \mathbf{c} , we draw conditional samples $\mathbf{z}_0 \sim \rho_\theta(\cdot | \mathbf{c})$ and compute the corresponding predicted properties via $\hat{\mathbf{c}}(\mathbf{z}_0) := \omega(\mathbf{z}_0)$. We then report the MAE between $\hat{\mathbf{c}}(\mathbf{z}_0)$ and the ground-truth QM9 property values.

Property alignment reward. For *conditional generation* with a target property vector \mathbf{y} , we add an auxiliary *property alignment reward* at the terminal state:

$$r_0^{(P)}(\mathbf{z}_0; \mathbf{y}) := -\|\omega(\mathbf{z}_0) - \mathbf{y}\|_1. \quad (8)$$

The total terminal reward combines $r_0^{(P)}$ with the energy/force rewards used by Elgin. The property reward is treated as an additional disentangled channel (analogous to energy and force) and is *group-normalized separately* before being weighted and added to the total advantage.

Results. As shown in Table 8, Elgin substantially reduces MAE compared to unguided diffusion baselines, narrowing the gap toward the QM9 oracle. Random samples define an upper bound on MAE, while QM9 ground-truth molecules provide a lower bound. Lower MAE indicates that generated samples reside closer to the true data distribution and admit more reliable property prediction. We report results for isotropic polarizability α , HOMO–LUMO gap $\Delta\varepsilon$, and LUMO energy $\varepsilon_{\text{LUMO}}$. Notably, Elgin improves all properties simultaneously, indicating that alignment with MLFF-derived energy and force feedback yields samples that are not only structurally stable but also chemically predictive.

Table 8. Mean Absolute Error (MAE) for molecular property prediction. Lower values indicate better controllable generation. Properties are predicted by a pretrained E(3)-equivariant EGNN regressor ω on molecular samples generated by each method. QM9 and Random serve as lower and upper bounds, respectively.

Property Units	α Bohr ³	$\Delta\varepsilon$ meV	$\varepsilon_{\text{LUMO}}$ meV
QM9	0.10	64	36
Random	9.01	1470	1457
N_{atoms}	3.86	866	813
EDM (Hoogeboom et al., 2022)	2.76	655	584
GeoLDM (Xu et al., 2023)	2.37	587	522
GeoBFN (Song et al., 2024)	2.34	577	516
Elgin (w/ Additional Property Alignment Reward)	2.32	564	512

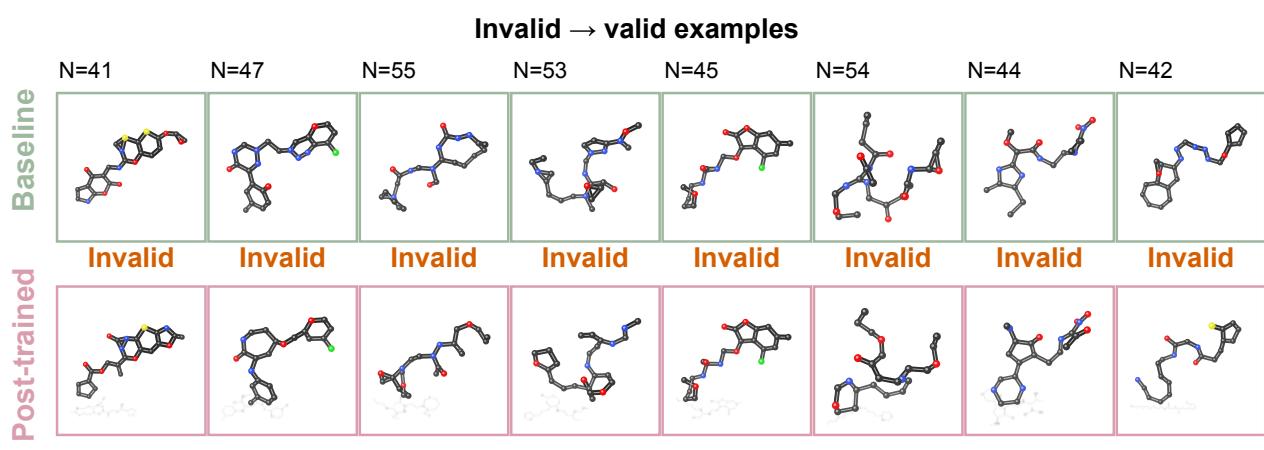


Figure 9. Qualitative comparison of molecular conformers generated by baseline versus post-trained models. Each column shows paired conformers for the same molecule (N indicates the number of atoms). The baseline model (top row, green borders) produces chemically invalid structures, while the post-trained model (bottom row, pink borders) generates valid conformers for the same molecular graphs. Atom colors: carbon (gray), nitrogen (blue), oxygen (red), sulfur (yellow), chlorine (green). Hydrogen atoms are shown as small gray spheres.

Table 9. FED-GRPO component ablations.

Setting	Disentangle?	Shared prefix?	PBRS?	A↑	M↑	V×U↑
Full FED-GRPO (Elgin)	✓	✓	✓	99.33	93.70	95.31
No disentanglement	✗	✓	✓	99.12	92.60	93.80
No shared prefix grouping	✓	✗	✓	98.85	83.40	91.20
Naive per-step energy (no PBRS)	✓	✓	✗	99.20	92.50	94.60

E.3. More Ablation Studies

FED-GRPO component ablations. Table 9 ablates key ingredients of FED-GRPO. **Disentangled normalization** stabilizes multi-objective optimization by preventing one reward channel (energy, force, or property) from dominating due to scale differences. **Shared-prefix grouping** reduces variance by comparing rollouts that share the same partially denoised prefix, making per-step advantages more meaningful. Finally, **dense PBRS shaping** improves credit assignment over long diffusion horizons, typically improving the stability–diversity trade-off relative to terminal-only rewards.

F. Supplementary Theory: PBRS as an Alchemical Force

While potential-based reward shaping (PBRS) is commonly introduced as a credit-assignment heuristic, in our setting it admits a concrete physical interpretation. We show that PBRS induces an exponential tilting of the one-step reverse diffusion kernel under a KL-regularized local control view, and that in the small-step limit this tilt manifests as an additional drift term proportional to the gradient of the shaping potential. When the potential is chosen as the negative MLFF energy evaluated on the predicted clean geometry, this drift corresponds to a force-like correction acting on the reverse dynamics. Crucially, because the diffusion state includes both coordinates and relaxed feature channels, the induced term decomposes into a physical force on atomic positions and an auxiliary “alchemical” force on non-geometric variables. This result formalizes PBRS as a principled approximation to energy-guided diffusion, clarifying how dense shaping rewards compile physical guidance into the learned reverse policy at training time, rather than injecting oracle forces during inference

Theorem 5 (PBRS induces an approximate (alchemical) force in the reverse drift). *Fix a reverse step t and denote the shaping potential by $\Psi_t(\mathbf{z}) := -E_\phi(\hat{\mathbf{z}}_{0|t}(\mathbf{z}))$ (Eq. (4)). Let the base one-step reverse kernel be Gaussian,*

$$q_t^{\Delta t}(\mathbf{z}_{t-\Delta t} \mid \mathbf{z}_t) = \mathcal{N}(\mathbf{z}_{t-\Delta t}; \mathbf{z}_t + b_t(\mathbf{z}_t)\Delta t, a_t\Delta t),$$

(e.g. Euler–Maruyama for a reverse SDE; a_t may include \mathbf{P}_{CoM}). Consider the per-step KL-regularized local update with

1320 PBRS reward $r_t^{\text{shape}} = \gamma \Psi_{t-\Delta t}(\mathbf{z}_{t-\Delta t}) - \Psi_t(\mathbf{z}_t)$ (Eq. (5)):

$$1321 \pi_t^{\star, \Delta t}(\cdot | \mathbf{z}_t) \in \arg \max_{\pi(\cdot | \mathbf{z}_t)} \left\{ \mathbb{E}_{\pi}[r_t^{\text{shape}}] - w_{\text{KL}} \text{KL}(\pi(\cdot | \mathbf{z}_t) \| q_t^{\Delta t}(\cdot | \mathbf{z}_t)) \right\}. \\ 1322$$

1323 Assume Ψ_t is C^1 in \mathbf{z} with locally Lipschitz $\nabla \Psi_t$. Then (i) the optimizer is the exponential tilt

$$1324 \pi_t^{\star, \Delta t}(\mathbf{z}_{t-\Delta t} | \mathbf{z}_t) \propto q_t^{\Delta t}(\mathbf{z}_{t-\Delta t} | \mathbf{z}_t) \exp\left(\frac{\gamma}{w_{\text{KL}}} \Psi_{t-\Delta t}(\mathbf{z}_{t-\Delta t})\right), \\ 1325$$

1326 and (ii) its conditional mean admits the small-step expansion

$$1327 \mathbb{E}_{\pi_t^{\star, \Delta t}}[\mathbf{z}_{t-\Delta t} | \mathbf{z}_t] = \mathbf{z}_t + \left(b_t(\mathbf{z}_t) + \frac{\gamma}{w_{\text{KL}}} a_t \nabla \Psi_t(\mathbf{z}_t) \right) \Delta t + o(\Delta t). \\ 1328$$

1329 Equivalently, in the $\Delta t \rightarrow 0$ limit the induced controlled reverse drift is

$$1330 \tilde{b}_t(\mathbf{z}) = b_t(\mathbf{z}) + \frac{\gamma}{w_{\text{KL}}} a_t \nabla \Psi_t(\mathbf{z}) = b_t(\mathbf{z}) - \frac{\gamma}{w_{\text{KL}}} a_t \nabla_{\mathbf{z}} E_{\phi}(\hat{\mathbf{z}}_{0|t}(\mathbf{z})). \\ 1331$$

1332 Writing $\mathbf{z} = [\mathbf{x}, \mathbf{h}]$, the \mathbf{x} -block recovers a force-like term $(-\nabla_{\mathbf{x}} E_{\phi})$ while the \mathbf{h} -block yields an “alchemical” preference direction $(-\nabla_{\mathbf{h}} E_{\phi})$ on relaxed types/features.

1333 *Proof.* Fix \mathbf{z}_t . The term $-\Psi_t(\mathbf{z}_t)$ is constant in $\mathbf{z}_{t-\Delta t}$, so the objective is

$$1334 \max_{\pi} \left\{ \gamma \mathbb{E}_{\pi}[\Psi_{t-\Delta t}(\mathbf{z}_{t-\Delta t})] - w_{\text{KL}} \text{KL}(\pi \| q_t^{\Delta t}) \right\}. \\ 1335$$

1336 A one-line variational calculation (Lagrange multiplier for $\int \pi = 1$) gives the unique maximizer

$$1337 \pi_t^{\star, \Delta t}(u | \mathbf{z}_t) \propto q_t^{\Delta t}(u | \mathbf{z}_t) \exp(\eta \Psi_{t-\Delta t}(u)), \quad \eta := \gamma / w_{\text{KL}}. \\ 1338$$

1339 Now use Gaussian integration by parts (Stein’s identity): if $U \sim \mathcal{N}(\mu, \Sigma)$ and g is smooth, $\mathbb{E}[(U - \mu)g(U)] = \Sigma \mathbb{E}[\nabla g(U)]$.
1340 Apply this with $U \sim q_t^{\Delta t}(\cdot | \mathbf{z}_t)$ and $g(u) = \exp(\eta \Psi_{t-\Delta t}(u))$ to obtain

$$1341 \mathbb{E}_{\pi_t^{\star, \Delta t}}[U | \mathbf{z}_t] = \mu + \Sigma \eta \mathbb{E}_{\pi_t^{\star, \Delta t}}[\nabla \Psi_{t-\Delta t}(U) | \mathbf{z}_t], \\ 1342$$

1343 where $\mu = \mathbf{z}_t + b_t(\mathbf{z}_t)\Delta t$ and $\Sigma = a_t \Delta t$. Since $U = \mathbf{z}_t + O_{\mathbb{P}}(\sqrt{\Delta t})$ under the Gaussian kernel and $\nabla \Psi$ is locally Lipschitz,

$$1344 \mathbb{E}_{\pi_t^{\star, \Delta t}}[\nabla \Psi_{t-\Delta t}(U) | \mathbf{z}_t] = \nabla \Psi_t(\mathbf{z}_t) + o(1). \\ 1345$$

1346 Substituting and collecting terms yields

$$1347 \mathbb{E}_{\pi_t^{\star, \Delta t}}[\mathbf{z}_{t-\Delta t} | \mathbf{z}_t] = \mathbf{z}_t + (b_t(\mathbf{z}_t) + \eta a_t \nabla \Psi_t(\mathbf{z}_t)) \Delta t + o(\Delta t), \\ 1348$$

1349 and replacing $\nabla \Psi_t = -\nabla E_{\phi}(\hat{\mathbf{z}}_{0|t})$ gives the stated “force” form. \square

1350 **From theory to measurement.** Theorem 5 provides a local-control view in which energy PBRS induces a force-like
1351 correction to the reverse drift. We next empirically probe this effect by measuring the score change induced by post-training
1352 and comparing its position component to the MLFF force direction.

1353 F.1. Alchemical force analysis

1354 To understand how post-training compiles physical preferences into the sampler, we measure an *alchemical force* induced by
1355 alignment and compare it to the actual force field. Write $\mathbf{z} = [\mathbf{x}, \mathbf{h}]$ for positions and atom-type features. Let $s_{\theta_{\text{pre}}}^{(\mathbf{x})}(\mathbf{z}_t, t)$
1356 and $s_{\theta}^{(\mathbf{x})}(\mathbf{z}_t, t)$ denote the *position* components of the pretrained and post-trained score networks, respectively. We define
1357 the alchemical force proxy (positions only) as

$$1358 \mathbf{F}_{\text{alc}}(\mathbf{z}_t, t) := s_{\theta}^{(\mathbf{x})}(\mathbf{z}_t, t) - s_{\theta_{\text{pre}}}^{(\mathbf{x})}(\mathbf{z}_t, t), \quad (9) \\ 1359$$

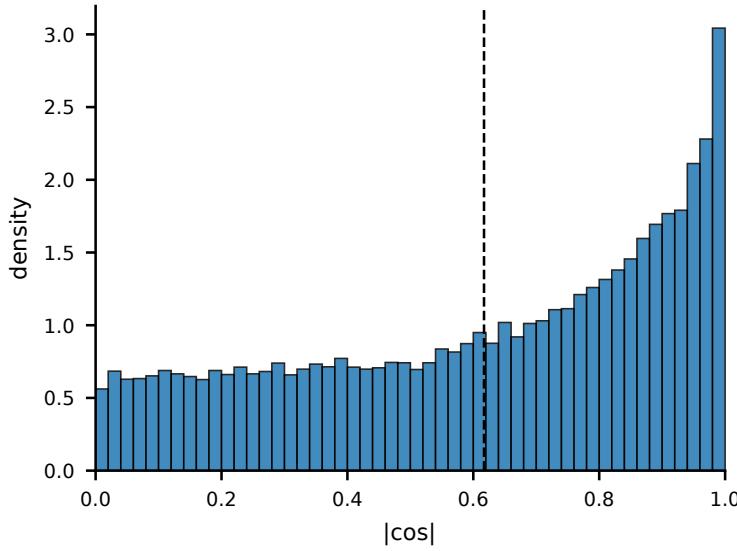


Figure 10. Distribution of $|\cos(\mathbf{F}_{\text{alc}}, \mathbf{F}_\phi)|$ between the position-space alchemical force \mathbf{F}_{alc} (post-trained minus pretrained score, Eq. 9) and the MLFF force $\mathbf{F}_\phi(\hat{\mathbf{z}}_{0|t})$ (Eq. 10). Larger values indicate stronger directional agreement.

so that $\mathbf{F}_{\text{alc}} \in \mathbb{R}^{N \times 3}$ is comparable to a physical force field. Given a force oracle (e.g., the MLFF) $\mathbf{F}_\phi(\hat{\mathbf{z}}_{0|t})$, we quantify agreement via cosine similarity

$$\cos(\mathbf{F}_{\text{alc}}, \mathbf{F}_\phi) := \frac{\langle \mathbf{F}_{\text{alc}}(\mathbf{z}_t, t), \mathbf{F}_\phi(\hat{\mathbf{z}}_{0|t}) \rangle}{\|\mathbf{F}_{\text{alc}}(\mathbf{z}_t, t)\| \|\mathbf{F}_\phi(\hat{\mathbf{z}}_{0|t})\|}. \quad (10)$$

Here $\hat{\mathbf{z}}_{0|t}$ is the predicted clean geometry from the diffusion posterior mean. Empirically, we find that \mathbf{F}_{alc} is strongly aligned with \mathbf{F}_ϕ . Figure 10 plots a histogram of $|\cos(\mathbf{F}_{\text{alc}}, \mathbf{F}_\phi)|$ computed over sampled diffusion states and molecules; higher values indicate that the post-training update (post-trained minus pretrained score, position component) points in a direction similar to the MLFF force. The distribution concentrates toward large cosine similarity, supporting the interpretation that alignment learns a force-like correction in the reverse dynamics.