ATROPDIFF: DATA-SCARCE ATROPISOMER GENER-ATION VIA MULTI-TASK PRETRAINED CLASSIFIER-GUIDED DIFFUSION

Letian Chen^{1,2}, Xi Wang³, Gufeng Yu¹, Caihua Shan^{4,*}, Yang Yang^{1,*} ¹Shanghai Jiao Tong University, ²Shanghai Innovation Institute, ³New York University, ⁴Microsoft Research Asia

clt2001@sjtu.edu.cn, xw3763@nyu.edu, jm5820zz@sjtu.edu.cn, caihuashan@microsoft.com,yangyang@cs.sjtu.edu.cn

Abstract

Customized molecule generation remains a challenging task, especially for datascarce categories such as atropisomers. Their structural complexity and scarcity of labeled data often lead to unstable optimization and poor controllability in traditional diffusion models. To address this gap, we propose AtropDiff, a diffusion framework guided by a multi-task pretrained classifier that implicitly encodes stereochemical knowledge. Key innovations include: 1) multi-task pretraining of a chirality-classifier for gradient-guided chirality control; 2) dynamically weighted gradient guidance to balance chirality accuracy and functional group validity during denoising; and 3) progressive integration of scaffold into the generation process to enable optimization on known atropisomers. Our results demonstrate that AtropDiff successfully overcomes data limitations, advancing AI-driven scientific discovery in chemistry.

1 INTRODUCTION

Atropisomers, arising from hindered rotation around stereogenic axes, play a pivotal role in asymmetric catalysis and drug discovery, with over 30% FDA-approved small-molecule drugs containing at least one atropisomeric axis (ST & JL, 2018). While deep generative models show promise for molecular design (Bagal et al., 2021; Luo et al., 2021; Hoogeboom et al., 2022; Xu et al., 2023), generating atropisomeric compounds faces twin challenges: (1) data scarcity (experimental determination requires resource-intensive chiroptical spectroscopy or quantum chemistry calculations) and (2) inadequate inductive biases (general molecular models fail to capture atropisomeric determinants).

Current solutions exacerbate these issues: fine-tuning foundation models suffers from mode collapse, while structure-free generation lacks chirality control. This creates a paradoxical situation where generation guidance requires precisely the stereochemical knowledge unlearnable from minimal data.

We address this through embedding domain knowledge into diffusion models via pretrained classifiers and scaffold information. Our AtropDiff framework introduces: 1) Multi-task pretrained atropisomer classifier: A classifier is pretrained using synthetic data augmentation on multiple chiralityrelated determinants, including chiral axis presence, rotational energy barriers, and relative energy. This enables robust performance during fine-tuning, even under real-data scarcity. 2) Dynamically weighted gradient guidance: The framework adaptively balances chiral fidelity and structural validity during the denoising process by periodically adjusting the strength of gradient guidance. 3) Progressive integration of scaffold into generation process: Chiral scaffold optimization is decoupled from substituent sampling, resulting in significantly improved chiral fidelity (97%). This framework represents a novel paradigm for addressing molecular generation challenges in data-scarce scenarios.



Figure 1: Framework overview of AtropDiff

2 Methods

The framework of our methods is shown in Figure 1. The method consists of two parts: (1) Prediction part: A classifier is trained on a dataset of computationally generated atropisomer candidates with three objectives: chiral axis identification, rotational energy barrier prediction, and relative energy prediction. (2) Generation part: An Equivariant Diffusion Model (EDM) is trained on the same dataset for molecular generation. Two generation strategies are used: gradient-guided generation and scaffold-guided generation. The gradient-guided generation uses the trained classifier to guide the EDM generation process, while the scaffold-guided generation uses a pre-defined scaffold to guide the EDM generation process.

2.1 DATASET CONSTRUCTION

Pre-training dataset. We generate a large-scale atropisomeric candidates dataset by: 1) Connecting 20 heterocyclic units into atropisomeric scaffolds, 2) Introducing two specialized substituent libraries and randomly attaching them to the scaffolds (10 common functional groups + 10 chirality-relevant groups), yielding 1.7M unique structures. For each molecule, we perform a relaxed scan of the dihedral angle around the inter-ring single bond (2.5° increments) using *xtb* to calculate rotational energy barriers and relative conformer energies. After DBSCAN clustering for redundancy reduction, we obtain 12M annotated conformers labeled with ChiralFinder (Shi et al., 2025), a package for finding chiral axes. The GEOM dataset (Axelrod & Gómez-Bombarelli, 2022) is incorporated with chiral axis annotations to enhance model generalization.

Fine-tuning dataset. We construct a challenging benchmark through: 1) PubChem similarity search based on known atropisomers (Yu et al., 2024), 2) Strict filtering using ChiralFinder and CSM software (Inbal et al., 2023), 3) Scaffold-based 7:3 train-test split on 2,000 curated molecules. For the fine-tuning dataset, the negative samples share structural similarity with positive samples to increase classification difficulty.

2.2 ATROPISOMER PREDICTION

Based on the Uni-Mol architecture (Zhou et al., 2023), we introduce three additional atropisomerspecific pre-training tasks:

- 1. Chiral axis identification (binary classification task)
- 2. Rotational energy barrier prediction (regression task)
- 3. Relative energy prediction (regression task)

Additional tasks enable the model to learn fundamental atropisomeric patterns from computationally generated labels before fine-tuning on limited experimental data.

2.3 Atropisomer generation

Our framework adopts the E(3)-equivariant Diffusion Model (EDM) (Hoogeboom et al., 2022) as the generative backbone, which operates on both continuous atomic coordinates $x \in \mathbb{R}^{M \times 3}$ and discrete atom types $h \in \mathbb{R}^{M \times K}$. The diffusion process is defined by gradually adding noise to data through a variance-preserving scheme:

$$q(\boldsymbol{z}_t | \boldsymbol{x}, \boldsymbol{h}) = \mathcal{N}_x(\boldsymbol{z}_t^{(x)} | \alpha_t \boldsymbol{x}, \sigma_t^2 \mathbf{I}) \cdot \mathcal{N}(\boldsymbol{z}_t^{(h)} | \alpha_t \boldsymbol{h}, \sigma_t^2 \mathbf{I}),$$

where $\alpha_t^2 + \sigma_t^2 = 1$ governs the signal-to-noise ratio. The reverse process learns to denoise $z_t = [z_t^{(x)}, z_t^{(h)}]$ via an equivariant graph network (EGNN):

$$\hat{\boldsymbol{\epsilon}}_t = \operatorname{EGNN}(\boldsymbol{z}_t^{(x)}, [\boldsymbol{z}_t^{(h)}, t/T]) - [\boldsymbol{z}_t^{(x)}, \boldsymbol{0}],$$

with the denoised estimate \hat{x} , $\hat{h} = z_t/\alpha_t - \sigma_t \hat{\epsilon}_t/\alpha_t$. The training objective minimizes the weighted noise prediction error:

$$\mathcal{L} = \mathbb{E}_{t,\boldsymbol{\epsilon}_t} \left[\frac{1}{2} w(t) \| \boldsymbol{\epsilon}_t - \hat{\boldsymbol{\epsilon}}_t \|^2 \right], \quad w(t) = 1 - \frac{\mathrm{SNR}(t-1)}{\mathrm{SNR}(t)}.$$

This formulation ensures E(3)-equivariant generation while jointly modeling geometric and categorical features, making it ideal for atropisomer-aware molecular design.

Gradient-guided generation To enable generation under property constraints, we adapt a training-free classifier-guidance strategy using our pretrained predictor (Han et al., 2024). By modifying Uni-Mol's first transformer layer to accept differentiable probability matrices, we create a gradient-guided diffusion process. During generation, the diffusion model navigates the chemical space through gradient signals from the chiral axis classifier, effectively leveraging prior knowledge without requiring additional training.

Scaffold-guided generation We adapt image inpainting techniques to molecular generation by fixing core scaffolds during diffusion. For a scaffold with n atoms, the model initializes coordinates and one-hot features for the scaffold atoms. At each denoising time step t, the model combines denoised new atoms (backward process) with noised scaffold atoms (forward process) to form the result as the input for the next time step. Finally, at time step 0, the output preserves the original scaffold while generating novel substituents around the scaffold. This approach enables lead optimization by providing a promising scaffold that is likely to form atropisomers.

3 EXPERIMENTS

3.1 ATROPISOMER PREDICTION

We assess the impact of incremental pretraining tasks on atropisomer prediction, as summarized in Table 1. The baseline Uni-Mol model achieves moderate performance with the F1 score of 0.785, while progressively adding pretraining tasks yields consistent improvements. The model with all additional pretraining tasks achieves the best performance, with an F1 score of 0.824.

To evaluate model generalizability, we use a hold-out test set comprising seven newly identified atropisomeric ligands. All pretraining configurations successfully identified atropisomerism in these molecules, demonstrating robust performance despite positive samples constituting only one-fourth of the fine-tuning dataset (563 out of 2000).

3.2 Atropisomer generation

We quantitatively compare the generation performance of scaffold-guided and gradient-guided approaches against baseline unconditional generation, as detailed in Table 2. To address the prohibitive

Table 1: Performance of Uni-Mol on atropisomer prediction tasks with different pretraining se	ttings.
t1 refers to the chiral axis identification task. t2 refers to the rotational energy barrier regress	ion. t3
refers to the relative energy regression.	

Pretraining Setting	Precision	Recall	F1 Score
Uni-Mol	0.779	0.793	0.785
Uni-Mol + t1	<u>0.832</u>	0.797	0.814
Uni-Mol + t1 + t2	0.825	<u>0.813</u>	0.819
Uni-Mol + t1 + t2 + t3	0.834	0.815	0.824

cost of experimental validation, we employ Chiralfinder to assess the presence of at least one chiral axis as our success rate metric. Our results demonstrate that both guided generation strategies achieve significantly higher success rates compared to unconditional generation, while maintaining comparable performance in standard molecular generation metrics including uniqueness and novelty. The relatively lower validity observed in scaffold-guided generation is because the scaffold size is approaching the upper limit of atom count supported by the model architecture. A similar scaffold does not appear in the training set. Nevertheless, the generative model still gives many valuable results. Several newly designed molecules are shown in Figure 2.

Table 2: Performance of molecular generation approaches



Figure 2: a) Gradient-guided generated molecules b) Scaffold-guided generated molecules

4 CONCLUSION

We present AtropDiff, the first deep learning-based framework capable of generating atropisomers. This work establishes a new paradigm for atropisomer design, which generates chemically plausible candidates computationally rather than through experimental trial-and-error. The scaffold-guided generation technique is particularly advantageous for optimizing lead compounds with a limited number of known examples. Our findings underscore the potential of leveraging physics-informed synthetic data and guided generation strategies to address data scarcity in AI-driven scientific discovery, particularly for complex molecular properties such as atropisomerism.

REFERENCES

- Simon Axelrod and Rafael Gómez-Bombarelli. Geom, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data*, 9(1):185, 2022. doi: 10.1038/s41597-022-01288-4. URL https://doi.org/10.1038/s41597-022-01288-4.
- Viraj Bagal, Rishal Aggarwal, PK Vinod, and U Deva Priyakumar. Molgpt: molecular generation using a transformer-decoder model. *Journal of Chemical Information and Modeling*, 62(9):2064– 2076, 2021.
- Xu Han, Caihua Shan, Yifei Shen, Can Xu, Han Yang, Xiang Li, and Dongsheng Li. Trainingfree multi-objective diffusion model for 3d molecule generation. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum? id=X41c4uB4k0.
- Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion for molecule generation in 3d. In *International conference on machine learning*, pp. 8867–8887. PMLR, 2022.
- Tuvi-Arad Inbal, Alon Gil, and Anir David. Csm, 2023. URL https://csm.ouproj.org. il/.
- Shitong Luo, Jiaqi Guan, Jianzhu Ma, and Jian Peng. A 3d generative model for structure-based drug design. Advances in Neural Information Processing Systems, 34:6229–6239, 2021.
- Runhan Shi, Chi Zhang, Gufeng Yu, Xiaohong Huo, and Yang Yang. Chiralfinder: Automated detection of stereogenic elements and discrimination of stereoisomers in complex molecules. *Chem-Rxiv*, 2025. doi: 10.26434/chemrxiv-2025-wz7kh. This content is a preprint and has not been peer-reviewed.
- Toenjes ST and Gustafson JL. Atropisomerism in medicinal chemistry: challenges and opportunities. *Future Med Chem*, 10, 2018. doi: 10.4155/fmc-2017-0152.
- Minkai Xu, Alexander S Powers, Ron O Dror, Stefano Ermon, and Jure Leskovec. Geometric latent diffusion models for 3d molecule generation. In *International Conference on Machine Learning*, pp. 38592–38610. PMLR, 2023.
- Gufeng Yu, Kaiwen Yu, Xi Wang, Xiaohong Huo, and Yang Yang. Clc-db: an online open-source database of chiral ligands and catalysts. 2024.
- Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Linfeng Zhang, and Guolin Ke. Uni-mol: A universal 3d molecular representation learning framework. In *The Eleventh International Conference on Learning Representations*, 2023. URL https: //openreview.net/forum?id=6K2RM6wVqKu.