

# Digital Discovery

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: D. Reidenbach, F. Nikitin, O. Isayev and S. G. Paliwal, *Digital Discovery*, 2025, DOI: 10.1039/D5DD00380F.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

# Applications of Modular Co-Design for De Novo 3D Molecule Generation

Danny Reidenbach,<sup>\*,†,||</sup> Filipp Nikitin,<sup>\*,‡,¶,§,||</sup> Olexandr Isayev,<sup>¶,‡</sup> and Saeed Paliwal<sup>†</sup>

<sup>†</sup>*NVIDIA, Santa Clara, CA, USA*

<sup>‡</sup>*Computational Biology Department, Carnegie Mellon University, Pittsburgh, PA, USA*

<sup>¶</sup>*Department of Chemistry, Carnegie Mellon University, Pittsburgh, PA, USA*

<sup>§</sup>*Work performed during internship at NVIDIA.*

<sup>||</sup>*These authors contributed equally.*

E-mail: dreidenbach@nvidia.com; fnikitin@andrew.cmu.edu

## Abstract

De novo 3D molecule generation is a pivotal task in drug discovery. However, many recent geometric generative models struggle to produce high-quality geometries, even if they are able to generate valid molecular graphs. To tackle this issue and enhance the learning of effective molecular generation dynamics, we present Megalodon—a family of scalable transformer models. These models are enhanced with basic equivariant layers and trained using a joint continuous and discrete denoising co-design objective. We assess Megalodon’s performance on established molecule generation benchmarks and introduce new 3D structure benchmarks that evaluate a model’s capability to generate realistic molecular structures, particularly focusing on geometry precision. We show that Megalodon achieves state-of-the-art results in 3D molecule generation, conditional structure generation, and structure energy benchmarks using diffusion and flow matching. Furthermore, we demonstrate that scaling Megalodon produces up to 49x more



valid molecules at large sizes and 2-10x lower energy compared to the prior best generative models. The code and the model are available at <https://github.com/NVIDIA-Digital-Bio/megalodon>.

## 1 Introduction

Molecular Generative models have been heavily explored due to the allure of enabling efficient virtual screening and targeted drug design<sup>1</sup>. Similar to the rise in their application to computer vision (CV)<sup>2,3</sup>, Diffusion and Flow Matching models have been applied for tasks including molecule design, molecular docking, and protein folding<sup>4-6</sup>. Across CV and chemical design, the scaling of model architectures and training data have seen significant accuracy improvements but questions surrounding how to scale effectively still persist<sup>7</sup>.

Specifically for 3D molecule generation (3DMG), where the task is to unconditionally generate valid and diverse 3D molecules, diffusion models have shown great promise in enabling accurate generation starting from pure noise<sup>8</sup>. The iterative nature of diffusion models allows them to explore a diverse range of molecular configurations, ideally providing valuable insights into potential drug candidates and facilitating the discovery of novel compounds. However, unlike in CV, which has seen systematic evaluations of training data and scaling, with tangible benchmark results<sup>9</sup>, measuring success in de novo molecule generation is quite difficult. As a result, there is a nonlinear path to determining what truly is making an impact if, in each model, the data, architecture, training objective, and benchmarks differ. Furthermore, the commonly shared 3DMG benchmarks that do exist only evaluate molecular topologies, ignoring geometry, conformational energy, and model generalization to large molecule sizes—all quantities that are imperative for real-world use. In this work, we explore the above in the context of 3DMG and its interpretable benchmarks to directly target larger molecules.

Our main contributions are as follows:

- We present Megalodon, a scalable transformer-based architecture for multi-modal



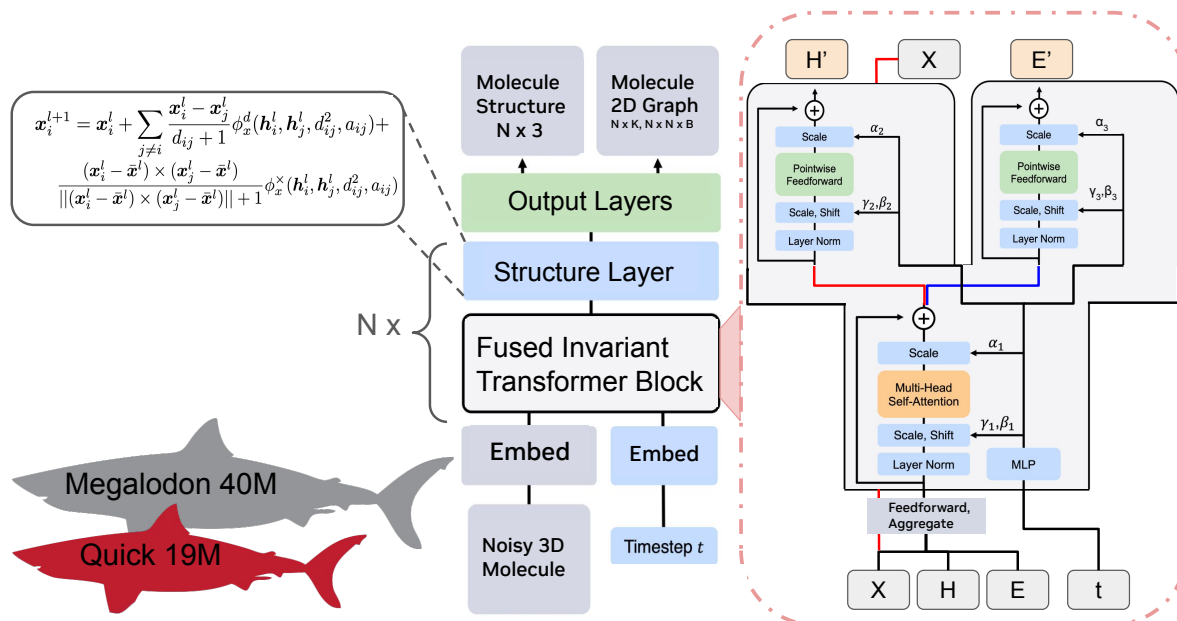


Figure 1: Megalodon Architecture: molecules are separated into 3D structures and discrete atom types, bond types, and atom charges features. All features are embedded separately, passed through a feed-forward neural network layer, and aggregated to produce the input tokens for the fused Invariant Transformer blocks. The embedded structural features and transformer outputs for the discrete features are passed to a single equivariant graph neural network (EGNN) layer for structure updates. The output heads consist of standard MLPs and an EGNN layer for bond refinement.

molecule diffusion and flow matching. This is the first 3DMG model to be tested with both objectives, with both obtaining state-of-the-art results. We show that our diffusion model excels at structure and energy benchmarks, whereas our flow matching model yields better 2D stability and the ability to use 25x fewer inference steps than its diffusion counterpart.

- Megalodon is the first model capable of unconditional molecule generation and conditional structure generation without retraining or finetuning.



## 2 Background

### 2.1 3D Molecule Generation

In de novo 3D molecule generation (3DMG), a molecule’s 3D structure and 2D topology are simultaneously generated. We define a molecule  $\mathbf{M} = (X, H, E, C)$  with  $N$  atoms where  $X \in \mathbf{R}^{N \times 3}$ ,  $H \in \{0, 1\}^{N \times A}$ ,  $E \in \{0, 1\}^{N \times N \times B}$ , and  $C \in \{0, 1\}^{N \times D}$  represents the atom coordinates, element types, bond types adjacency matrix, and formal charges respectively. Here,  $A$  denotes the number of atom types,  $B$  the number of bond types, and  $D$  the number of formal charge states.  $X$  is modeled as a continuous variable, whereas  $H$ ,  $E$ , and  $C$  are discrete one-hot variables.  $X$  is modeled as a continuous variable whereas  $H$ ,  $E$ , and  $C$  are discrete one-hot variables.

### 2.2 Important Qualities of 3D Molecules

The GEOM dataset<sup>10</sup> is widely used for 3D molecular structure (conformer) generation tasks, containing 3D conformations from both the QM9 and drug-like molecule (DRUGS) databases, with the latter presenting more complex and realistic molecules. Conformers in the dataset were generated using CREST<sup>11</sup>, which performs extensive conformational sampling based on the semi-empirical extended tight-binding method (GFN2-xTB)<sup>12</sup>. This ensures that each conformation represents a local minimum in the GFN2-xTB potential energy surface (PES).

A key requirement for generative models is their ability to implicitly learn PES of the training data and produce molecules that are local minima of the PES. However, since GFN2-xTB is itself a model rather than a universal energy function, comparing energies across different potentials (e.g., using GFN2-xTB optimized structures but computing energies with MMFF<sup>13</sup>) can introduce systematic errors. Differences in potential models, such as optimal bond lengths, may lead to unreliable results. Overall, the goal of 3DMG is to generate valid molecules mimicking the energy landscape of the GEOM dataset.



## 2.3 Related Work

Hoogetboom et al.<sup>8</sup> first introduced continuous diffusion modeling for coordinates and atom types using a standard equivariant graph neural network (EGNN) architecture<sup>14</sup>. Following this, many models have been produced that make slight changes to the architecture and diffusion interpolant schedule to generate atom coordinates and types<sup>15</sup>. While initially effective, they rely on OpenBabel<sup>16</sup> software to infer chemical bonds, which is a standalone hard problem and introduces additional sources of error into the pipeline. So, methods began to generate the bond locations and types in the generative process<sup>17</sup>. Vignac et al.<sup>18</sup> was the first to use continuous diffusion for coordinates and discrete diffusion for the atom and bond types, removing the OpenBabel requirement. Le et al.<sup>19</sup> used the same training objective but introduced a more effective equivariant architecture. Recently Irwin et al.<sup>20</sup> uses continuous and discrete flow matching with a latent equivariant graph message passing architecture to show improved performance.

Xu et al.<sup>21</sup> introduces GeoLDM a geometric latent diffusion model for 3DMG. GeoLDM applies its diffusion process over a learned latent representation. So rather than updating the atom position and types in euclidean space everything is done inside the model. Similar to EDM, GeoLDM uses OpenBabel for bond prediction. Pinheiro et al.<sup>22</sup> takes a different approach than majority of prior work in representing molecules as 3D voxels rather than graphs. This is akin to 3D image processing rather than point cloud processing. This however requires a recovery process as the voxels are not a natural molecule representation. Voxels however provide a better link to the applications of vision models which majority of the diffusion framework was created for. Lastly, Song et al.<sup>23</sup> introduces GeoBFN a Geometric Bayesian Flow Network, that unlike diffusion models operate in the parameter space rather than product space. While the integration of 3D voxels would not work for Megalodon, latent diffusion and BFN extensions are something relevant to future work.



## 2.4 Stochastic Interpolants

**Continuous Gaussian Interpolation** Following<sup>24,25</sup>, in the generative modeling setting, we construct interpolated states between an empirical data and a Gaussian noise distribution  $\mathcal{N}(\mathbf{x}_t; \beta(t)\mathbf{x}_1, \alpha(t)^2\mathbf{I})$ , this is,

$$\mathbf{x}_t = \alpha(t)\boldsymbol{\epsilon} + \beta(t)\mathbf{x}_1, \quad (1a)$$

$$\mathbf{x}_1 = \frac{\mathbf{x}_t - \alpha(t)\boldsymbol{\epsilon}}{\beta(t)} \quad (1b)$$

where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I})$  and  $\mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)$ . Common choices for the interpolation include (assuming  $t \in [0, 1]$ ), with  $t = 1$  corresponding to data and  $t = 0$  to noise:

- Variance-preserving SDE-like from the diffusion model literature<sup>26</sup>:  $\alpha(t) = \sqrt{1 - \gamma_t^2}$  and  $\beta(t) = \sqrt{\gamma_t^2}$  with some specific “noise schedule”  $\gamma_t$  which is commonly written as  $\sqrt{\bar{\alpha}_t}$  from Ho et al.<sup>27</sup>.
- Conditional linear vector field<sup>24</sup>:  $\alpha(t) = 1 - (1 - \sigma_{\min})t$  and  $\beta(t) = t$  with some smoothening of the data distribution  $\sigma_{\min}$ .

**Continuous Diffusion** Continuous Denoising Diffusion Probabilistic Models (DDPM) integrate a gradient-free forward noising process based on a predefined discrete-time variance schedule (Eq. 1a) and a gradient-based reverse or denoising process<sup>27</sup>. The denoising model can be parameterized by data or noise prediction as they can be equilibrated via Eq. 1b. Following Le et al.<sup>19</sup>, we use the following training objective and update rule:

$$\mathcal{L}_{\text{DDPM}}(\theta) = \mathbb{E}_{t, \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} \|\mathbf{x}_\theta(t, \mathbf{x}_t) - \mathbf{x}_1\|^2 \quad (2)$$

$$\mu_\theta(t, \mathbf{x}_t) = \mathbf{f}(\alpha(t), \beta(t)) * \mathbf{x}_\theta(t, \mathbf{x}_t) + \mathbf{g}(\alpha(t), \beta(t)) * \mathbf{x}_t \quad (3)$$

$$\mathbf{x}_{t+1} = \mu_\theta(t, \mathbf{x}_t) + \sigma((\alpha(t), \beta(t)) * \boldsymbol{\epsilon})$$

where functions  $\mathbf{f}$ ,  $\mathbf{g}$ , and  $\boldsymbol{\sigma}$  are defined for any noise schedule such as the cosine noise schedule used in Vignac et al.<sup>18</sup>.



**Continuous Flow Matching** Flow matching (FM) models are trained using the conditional flow matching (CFM) objective to learn a time-dependent vector field  $\mathbf{v}_\theta(t, \mathbf{x}_t)$  derived from a simple ordinary differential equation (ODE) that pushes samples from an easy-to-obtain noise distribution to a complex data distribution.

$$\begin{aligned}\mathcal{L}_{\text{CFM}}(\theta) &= \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} \left\| \mathbf{v}_\theta(t, \mathbf{x}_t) - \frac{d}{dt} \mathbf{x}_t \right\|^2 \\ &= \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} \left\| \mathbf{v}_\theta(t, \mathbf{x}_t) - \dot{\alpha}(t)\epsilon - \dot{\beta}(t)\mathbf{x}_1 \right\|^2,\end{aligned}\quad (4)$$

The time-differentiable interpolation seen in Eq. 1a gives rise to a probability path that can be easily sampled. For more details on how to relate the Gaussian diffusion and CFM objectives with the underlying score function of the data distribution, please see Appendix A.

In practice, many methods use a "data prediction" objective to simplify training, which gives rise to the following loss function and inference Euler ODE update step following the conditional linear vector field<sup>20,24</sup>.

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} \left\| \mathbf{x}_\theta(t, \mathbf{x}_t) - \mathbf{x}_1 \right\|^2 \quad (5)$$

$$\mathbf{v}_\theta(t, \mathbf{x}_t) = \frac{\mathbf{x}_\theta(t, \mathbf{x}_t) - \mathbf{x}_t}{1 - t}, \quad (6)$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \mathbf{v}_\theta(t, \mathbf{x}_t)dt$$

**Discrete Diffusion** Following Austin et al.<sup>28</sup>, Discrete Denoising Diffusion Probabilistic Models (D3PMs) apply the same concept as continuous diffusion but over a discrete state space. Like the continuous counterpart that relies on a predefined schedule to move mass from the data to prior distribution, D3PM uses a predefined transition matrix that controls how the model transitions from one discrete state to another.

For scalar discrete random variables with  $K$  categories  $a_t, a_{t-1} \in 1, \dots, K$  the forward transition probabilities can be represented by matrices:  $[Q_t]_{ij} = q(a_t = j | a_{t+1} = i)$ . Starting from our data  $a_1$  or  $a_T$  (where  $T$  is the total number of discrete time steps)<sup>1</sup>, we obtain the

<sup>1</sup>We adjust the direction of time for diffusion to match the FM equations such that  $T=1$  is data.





following  $T - t + 1$  step marginal and posterior at time  $t$ :

$$\begin{aligned}\bar{Q}_t &= Q_t Q_{t+1} \dots Q_T \\ q(a_t|a_{t+1}) &= \text{Cat}(a_t; p = a_{t+1}Q_t), \quad q(a_t|a_T) = \text{Cat}(a_t; p = a_T\bar{Q}_t), \\ q(a_{t+1}|a_t, a_T) &= \frac{q(a_t|a_{t+1}, a_T)q(a_{t+1}|a_T)}{q(a_t|a_T)} \\ &= \text{Cat}\left(a_{t+1}; p = \frac{a_tQ_t^\top \odot a_T\bar{Q}_{t+1}}{a_T\bar{Q}_ta_t^\top}\right)\end{aligned}\tag{7}$$

Here  $Q$  is defined as a function of the same cosine noise schedule used in continuous DDPM such that the discrete distribution converges to the desired terminal distribution (*i.e.* uniform prior) in  $T$  discrete steps. Similar to the use of mean squared error loss for DDPM, D3PM uses a discrete cross-entropy objective.

**Discrete Flow Matching** Following Campbell et al.<sup>29</sup>, we use the Discrete Flow Matching (DFM) framework to learn conditional flows for the discrete components of molecule generation (atom types, bond types, and atom charges). We use the following DFM interpolation in continuous time, where  $S$  is the size of the discrete state space:

$$p_{t|1}^{\text{unif}}(a_t|a_1) = q(a_t|a_1) = \text{Cat}(t\delta\{a_1, a_t\} + (1-t)\frac{1}{S}),\tag{8}$$

Similar to discrete diffusion, we use the cross-entropy objective for training. Please see Campbell et al.<sup>29</sup> for sampling procedure details.

**Diffusion vs. Flow Matching** We see that for both Diffusion and CFM, the loss functions used in practice are identical. Differences arise in how we build the interpolation, how we sample from these models, and their theoretical constraints. Diffusion models rely on complex interpolation schedules that are tuned to heavily weight the data distribution using a uniform time distribution. In contrast, FM commonly uses a simple linear interpolation but can achieve that same data distribution weighting by sampling from more complex time distributions. The choices of time distributions and interpolation schedules can be chosen appropriately to make FM and Diffusion equivalent in the Gaussian setting (see Appendix. A).



We show in Fig. 2 the interpolation and time distribution differences that mimic the same weighting of  $p_{\text{data}}$  at  $T=1$  that are currently used in recent 3DMG models<sup>19,20</sup>.

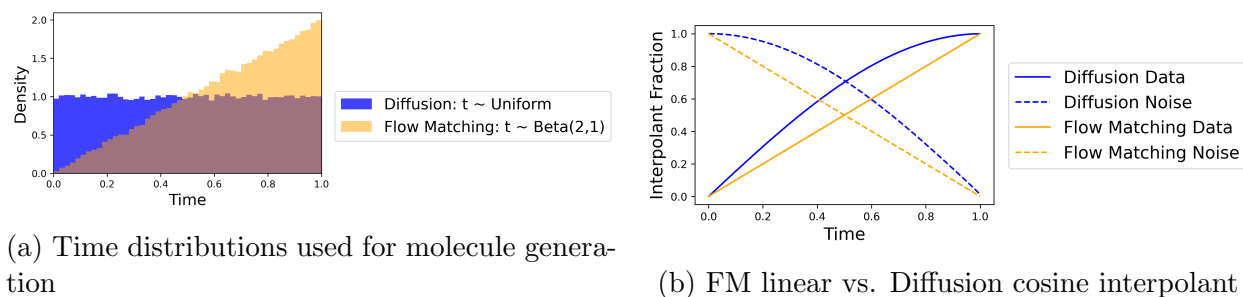


Figure 2: Time and interpolation comparison between Megalodon and Megalodon-flow

Diffusion models inherently rely on simulating Gaussian stochastic processes. In the forward process, data points are progressively noised, converging towards a Gaussian distribution. This process, derived from score-based generative models, aims to learn the score function (the gradient of the data distribution’s log density) to reverse the diffusion process. The generative model effectively solves a Stochastic Differential Equation (SDE) that describes how data diffuses towards noise and how it can be denoised in reverse. The reverse process requires SDE simulation at every step, which involves sampling from a learned probabilistic model that estimates how to remove noise. This involves simulating random variables at each time step, making diffusion models highly dependent on repeated stochastic simulation.

Flow Matching, on the other hand, learns a continuous vector field that deterministically ”flows” one distribution to another. The model learns this flow by matching the velocity field that pushes samples from a source distribution to a target distribution. Once the vector field is learned, generating samples involves solving an ODE that defines a continuous and deterministic trajectory from the source to the target distribution. Unlike diffusion models, which require simulating a series of stochastic transitions (noising and denoising) over many steps, flow matching learns a single, continuous flow. Sampling involves solving an ODE (or, in some cases, a deterministic SDE with noise) to move from the base distribution to the target in a smooth, deterministic fashion.



For DDPM, the equations only hold for the Gaussian path with access to a well-formed score function. This is why techniques like mini-batch Optimal Transport (OT) can be applied to FM but not Diffusion to align  $p_{\text{data}}$  and  $p_{\text{ref}}$ <sup>30</sup>. In FM, the vector field is learned, which, in the absence of OT, can be derived as a function of the score function, but having access to the score function is not a requirement to sample deterministically (simulation-free).

### 3 Methods

**Megalodon Architecture** Since 3DMG allows for the simultaneous generation of a discrete 2D molecular graph and its 3D structure, we intentionally designed our architecture with a core transformer trunk to better model discrete data<sup>31,32</sup>. Fig. 1 illustrates the model architecture, which is comprised of N blocks made up of fused invariant transformer blocks and simple structure update layers, followed by linear layers for discrete data projection.

In the fused invariant transformer block, the embedded structure, atom types, and bond types are fused and aggregated to create a single molecule feature. This is passed into a standard multi-head attention module with adaptive layernorm. The scaled output is then passed into separate adaptive layernorm feedforward blocks for the atom types and bond types. The transformer also produces an unchanged molecule structure via a residual connection to the input. The updated atom and bond types are then passed into a simple structure layer. The structure layer only updates the predicted structure via a standard distance-based EGNN update with a cross-product term<sup>4,14</sup>. At a high level, the transformer block updates our discrete invariant data, and our equivariant layer updates our structure. For more details, please see Appendix. B.

We introduce a generative scaling benchmark, and as we show, the performance of 3DMG models is correlated with the size of the generated molecules. We note that our large model is, in fact, not that large compared to recent biological models<sup>33,34</sup> and can be further scaled beyond 40M params if further benchmarks are developed.



**Training Objective** We explore Megalodon in the context of diffusion and flow matching. For our diffusion flavored model, following Vignac et al.<sup>18</sup>, Le et al.<sup>19</sup> we use the same weighted cosine noise schedules, DDPM, and discrete D3PM objective. When using conditional flow matching, we apply the same training objective and hyperparameters as Irwin et al.<sup>20</sup>, including equivariant optimal transport. In this way, for diffusion and flow matching, we train and evaluate our model in an *identical way* including hyperparameters to prior models of same types.

In our experiments with EQGAT-diff, we found that the diffusion objective with data-like priors possesses an interesting but potentially harmful behavior. Although the noise sample from the data-prior and the true data sample have bonds, the model consistently generates no bonds for all time  $\leq 0.5$ , which corresponds to an interpolation with  $\leq 70\%$  of the data as seen in Fig. 2(b). Therefore there is no useful information for the edge features in half the training and inference samples. As a result, only when the structure error is low, as the model starts with 70% data in the interpolation, does the bond prediction accuracy jump to near-perfect accuracy. Thus, only when the structure is accurate was the 2D graph accurate, which is counterintuitive to the independent and simultaneous objective. In other words, the 2D graph does not inform the 3D structure as one would expect to happen, and we would want equal importance on the 2D topology and 3D structure.

To address this inefficiency, as the structure, atom type, and bond type prediction inform each other to improve molecule generation, we introduce a subtle change to the training procedure similar to Campbell et al.<sup>29</sup>. Keeping each data type having its own independent noise schedule, we enable a concrete connection between the discrete and continuous data that it is modeling. Explicitly, rather than sampling a single time variable, we introduce a second noise variable to create  $t_{continuous}$  and  $t_{discrete}$ , both sampled from the same time distribution. Now discrete and continuous data are interpolated with their respective time variable *and* maintain the independent weighted noise schedules. We note that the MiDi weighted cosine schedules were already adding different levels of noise for the same time



value. Now, we take that one step further and allow the model to fill in the structure given the 2D graph and learn to handle more diverse data interpolations.

**Self Conditioning** Following Chen et al.<sup>35</sup>, we train Megalodon with self-conditioning similar to prior biological generative models<sup>20,36,37</sup>. We found that constructing self-conditioning as an outer model wrapper with a residual connection led to faster training convergence:

$$\begin{aligned}h_{\text{sc}} &= \text{model}(h_t) \\ h_t &= \text{MLP}([h_{\text{sc}}, h_t]) + h_t \\ h_{\text{pred}} &= \text{model}(h_t)\end{aligned}\tag{9}$$

where  $h_t$  represent one of the molecule component.

Specifically for 3DMG, self-conditioning is applied independently to each molecule component  $\mathbf{M} = (X, H, E, C)$ , where the structure component uses linear layers without bias and all discrete components operate over the raw logits rather than the one-hot predictions.

## 4 Experiments

**Data** GEOM Drugs is a dataset of drug-like molecules with an average size of around 44 atoms<sup>10</sup>. Following standard practice in prior work<sup>18–20</sup>, we train on the five lowest-energy conformers per molecule, using the same splits as these baselines. We emphasize that traditional metrics are calculated by first sampling molecule sizes from the dataset( Fig. 5) and then generating molecules with the sampled number of atoms, including explicit hydrogens. We show in Sec. 4.1 that this does not illustrate the full generative capacity, as in many real-world instances, people want to generate molecules with greater than 100 atoms<sup>38</sup>.

### 4.1 Unconditional De Novo Generation

**Problem Setup** Following Le et al.<sup>19</sup> we generate 5000 molecules (randomly sampling the number of atoms from the train distribution see Fig. 5), and report (1) Atom Stability:



Table 1: Measuring Unconditional Molecule Generation: 2D and 3D benchmarks. \* Denotes taken from EQGAT-Diff.

Model	Steps	2D Topological ( $\uparrow$ )			3D Distributional ( $\downarrow$ )	
		Atom Stab.	Mol Stab.	Validity	Bond Angle	Dihedral
EDM+OpenBabel*	1000	0.978	0.403	0.363	–	–
MiDi*	500	0.997	0.897	0.705	–	–
EQGAT-diff <sub>disc</sub> <sup>x0</sup>	500	0.998	0.935	0.830	0.858	2.860
EGNN + cross product	500	0.982	0.713	0.223	14.778	17.003
Megalodon-quick	500	0.998	0.961	0.900	0.689	2.383
Megalodon	500	<b>0.999</b>	<b>0.977</b>	<b>0.927</b>	<b>0.461</b>	<b>1.231</b>
SemlaFlow	100	0.998	0.979	0.920	<b>1.274</b>	<b>1.934</b>
Megalodon-flow	100	<b>0.999</b>	<b>0.988</b>	<b>0.944</b>	1.286	2.379

the percentage of individual atoms that have the correct valency according to its electronic configuration that was predefined in a lookup table, (2) Molecule Stability: percentage of molecules in which all atoms are stable, (3) Connected Validity: fraction of molecules with a single connected component which can be sanitized with RDKit. We also introduce two structural distributional metrics for the generated data: (4) bond angles and (5) dihedral angles, calculated as the weighted sum of the Wasserstein distance between the true and generated angle distributions, with weights based on the central atom type for bond angles and the central bond type for dihedral angles, respectively.

**Baselines** EQGAT-diff has 12.3M parameters and leverages continuous and discrete diffusion<sup>19</sup>. SemlaFlow has 23.3M params<sup>2</sup> and is trained with conditional flow matching with equivariant optimal transport<sup>20</sup>. We report two Megalodon sizes, small (19M) and large (40.6M). We train with identical objectives and settings to both EQGAT-diff and SemlaFlow. We also compare to older diffusion models, including MiDi and EDM, as they introduce imperative techniques from which the more recent models are built.

**Analysis** Both the diffusion and flow matching versions of Megalodon achieve state-of-the-art results. With the FM version obtaining better topological accuracy and the diffusion

<sup>2</sup>Checkpoint from public code has 2 sets of 23.2M params, one for the last gradient step and EMA weights



version seeing significantly improved structure accuracy. This experiment shows that the underlying augmented transformer is useful for the discrete and continuous data requirements of 3DMG, regardless of the interpolant and sampling methodology. We also see that the transformer part is crucial for Megalodon’s success as just using the EGNN with cross-product updates with standard edge and feature updates for the non-equivariant quantities performs quite poorly. We also note that all methods obtain 100% *uniqueness*, 88-90% *diversity*, and 99% *novelty* following<sup>19</sup> definitions with no meaningful performance differences. For additional model comparisons and ablations related to reducing the number of inference steps, please refer to Appendix Table 6 and Appendix Sec. C.2. To illustrate the generalizability of Megalodon, Appendix Table 5 and Appendix Sec. C.1 report its performance on the QM9 dataset.

**Impact of molecule size on performance** As Table 1 shows average results over 5000 molecules of relatively small and similar sizes, it is hard to understand if the models are learning how to generate molecules or just regurgitating training-like data. We design an experiment to directly evaluate this question and see how models perform as they are tasked to generate molecules outside the support region of the train set. We see in Fig. 3 that the topological model performance is a function of length (for full-size distribution, see Fig. 5). Here for each length [30, 125] we generate 100 molecules and report the percentage of stable and valid molecules.

We emphasize that Table 1 illustrates only a slice of the performance via the average of 5K molecules sampled from the train set size distribution. We note that although molecules with greater than 72 atoms make up  $\leq 1\%$  of the train set, Megalodon demonstrates roughly 2-49x better performance than EQGAT-diff for the larger half of the generated molecule sizes. We hypothesize that since molecule stability is a discrete 2D measurement, the transformer blocks in Megalodon allow it to better generalize even if seeing similar molecules in less than 0.1% of the training data. In other words, the ability of transformers to excel at modeling discrete sequential data improves our generative performance. We want to point out that all



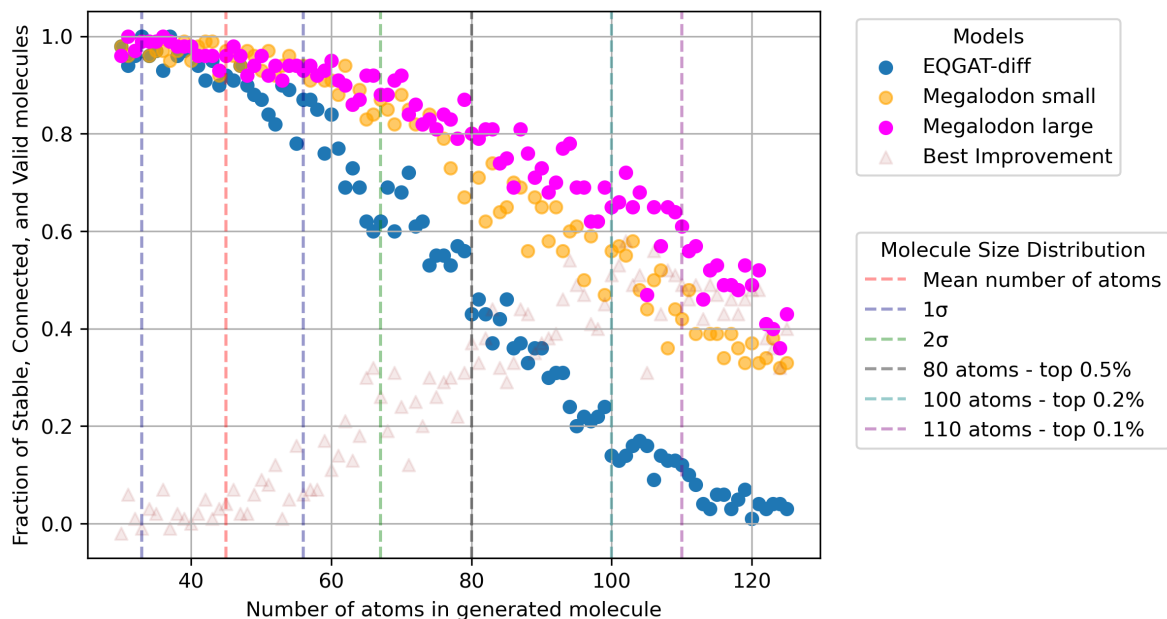


Figure 3: Diffusion model performance as a function of molecule size. Note the ability for Megalodon to generate valid and stable molecules with little training data support.

tested models are trained with identical datasets, hyperparameters, diffusion schedules, and training objectives. The only difference is the architecture. We also see that the ability to scale our simple architecture allows the model to even better generate molecules outside the region of data support. Lastly, we chose to focus on only the diffusion models here as they exhibit the best structure benchmark performance.

## 4.2 Conditional Structure Generation

Similar to the 3D molecule generation task, we use the GEOM-Drugs dataset to evaluate the conditional structure generation capabilities of our model. Given all unconditional 3DMG models are trained with independent noising of coordinates, atoms, and, in some cases, bonds, we want to evaluate how accurate the structural component is. We note this is something that is lacking from the existing prior benchmarks, as when generating de novo molecules, there is no ground truth structure to compare against. In the task of conditional structure generation, all models are given the molecule 2D graph (atom types, bonds) and asked to generate the 3D structure in which ground truth data exists. Given Vignac et al.<sup>18</sup>





and Jing et al.<sup>39</sup> use different train/test splits, we evaluate all methods on the overlap of 200 held-out molecules, with all methods generating 43634 structures in total. Due to the similarities with the baselines and its superior unconditional structure accuracy, we compare Megalodon trained with diffusion against recent methods with public reproducible code.

**Problem setup.** We report the average minimum RMSD (AMR) between ground truth and generated structures, and Coverage for Recall and Precision. Coverage is defined as the percentage of conformers with a minimum error under a specified AMR threshold. Recall matches each ground truth structure to its closest generated structure, and Precision measures the overall spatial accuracy of each generated structure. Following Jing et al.<sup>39</sup>, we generate two times the number of ground truth structures for each molecule. More formally the precision metrics are defined, for  $K = 2L$ , let  $\{C_l^*\}_{l \in [1..L]}$  and  $\{C_k\}_{k \in [1..K]}$  respectively be the sets of ground truth and generated structures:

$$\begin{aligned} \text{COV-Prec.} &:= \frac{1}{K} \left| \{k \in [1..K] : \min_{l \in [1..L]} \text{RMSD}(C_k, C_l^*) < \delta\} \right| \\ \text{AMR-Prec.} &:= \frac{1}{K} \sum_{k \in [1..K]} \min_{l \in [1..L]} \text{RMSD}(C_k, C_l^*) \end{aligned} \quad (10)$$

where  $\delta$  is the coverage threshold. The recall metrics are obtained by swapping ground truth and generated conformers.

**Baselines** We compare Megalodon with EQGAT-Diff, GeoDiff<sup>40</sup>, and TorsionalDiffusion<sup>39</sup>. For the unconditional 3DMG models, including Megalodon, we prompt them with the ground truth atom types and bond types to guide the generation of the structure along the diffusion process. This is done by replacing the input and output with the fixed conditional data. We do this to assess what the model is actually learning across the multiple data domains. The central question being, is the model learning how to generate molecules over the spatial and discrete manifolds, or is it just learning how to copy snapshots of training-like data?

**Analysis** We see in 2 that EQGAT-diff is unable to generate any remotely valid structures. Even though all modalities are being denoised independently at different rates, the model



cannot generate the structure given ground truth 2D molecule graphs. This is also seen during the sampling process, where diffusion models trained with similar denoising objectives as EQGAT-diff generate no bonds until the structure has seemingly converged. Therefore during most of the sampling process, the edge features which make up a large portion of the computational cost hold no value.

In comparison, Megalodon generates structures with competitive precision and recall by building a relationship between the discrete and continuous data directly in the training process described in Sec. 3. Half the time all data types are independently noised as normal with their respective time variables and schedules, the other half we only add noise to the structure. Therefore, our model learns to build a relationship between true 2D graphs and their 3D structure, as well as any interpolation between the three data tracks that are interpolated independently with different schedulers.

Table 2: Quality of ML generated conformer ensembles for GEOM-DRUGS ( $\delta = 0.75\text{\AA}$ ) test set in terms of Coverage (%) and Average RMSD ( $\text{\AA}$ ). Bolded best, underlined second best.

Method	Recall				Precision			
	Coverage $\uparrow$		AMR $\downarrow$		Coverage $\uparrow$		AMR $\downarrow$	
	Mean	Med	Mean	Med	Mean	Med	Mean	Med
GeoDiff	42.1	37.8	0.835	0.809	24.9	14.5	1.136	1.090
Tor. Diff.	<b>75.3</b>	<b>82.3</b>	<b>0.569</b>	<b>0.532</b>	<u>56.5</u>	<u>57.9</u>	<u>0.778</u>	<u>0.731</u>
EQGAT	0.8	0.0	2.790	2.847	0.1	0.0	3.754	3.771
Megalodon	<u>71.4</u>	<u>75.0</u>	<u>0.573</u>	<u>0.557</u>	<b>61.2</b>	<b>63.1</b>	<b>0.719</b>	<b>0.696</b>

Megalodon demonstrates that its unconditional discrete diffusion objective is crucial for its conditional performance. In other words, the discrete diffusion training objective improves the conditional continuous generative performance. This is evident in the comparison between GeoDiff and Megalodon. GeoDiff is trained on the same conditional Euclidean structure objective as Megalodon (with similar EGNN-based architecture) with 10x more diffusion steps, with both models taking in identical inputs. We see that since Megalodon is able to generate molecules from pure noise, it better learns structure and as a result can be prompted to generate accurate structures.



Interestingly, compared to Torsional Diffusion, which initializes the 3D structure with an RDKit approximation to establish all bond lengths and angles and then only modifies the dihedral angles, we see quite competitive performance. Before, it was understood that by restricting the degrees of freedom with good RDKit structures, the performance jump from GeoDiff to Torsional Diffusion was observed. Now we see that with the same euclidean diffusion process, similar accuracy improvements can be gained by learning how to generate accurate discrete molecule topology via discrete diffusion. We want to note that there have been recent advances on top of Torsional Diffusion<sup>41</sup> and other conformer-focused models that are not public<sup>42</sup>. We use this benchmark more to analyze the underlying multi-modal diffusion objective and focus on the underlying model comparisons. Megalodon is not a conformer generation model but a molecule generation model capable of de novo and conditional design. Overall, Megalodon shows how independent time interpolation and discrete diffusion create the ability for the model to be prompted or guided with a desired 2D topology to generate accurate 3D structures.

### 4.3 Unconditional Structure-based Energy Benchmarks

**Problem setup** Each ground truth structure in GEOM dataset represents a low-energy conformer within its ensemble, highlighting two key aspects. First, these molecules are local minima on the GFN2-xTB potential energy surface. Second, their energies are lower compared to other conformations sampled in the ensemble. Previously, these quantities have not been thoroughly evaluated for generated molecules. To address this gap, we directly measure how closely a generated molecule approximates its nearest local minimum (i.e., its relaxed structure). We measure the energy difference between the initial generated structure and its relaxed counterpart, as well as structural changes in bond lengths, bond angles, and dihedral (torsion) angles. This approach allows us to evaluate the ability of generative models to produce molecules that are true local minima, facilitating faster ranking of generated structures without additional minimization steps. For a more rigorous treatment of this benchmark



framework, including additional analyses and methodological considerations, we refer the reader to our accompanying work on benchmarking generative models on the GEOM-Drugs dataset<sup>43</sup>.

Table 3: xTB Relaxation Error: Length Å, angles degrees, energy kcal/mol. These metrics are taken over the valid molecules from Table 1. Methods are grouped by model type: diffusion (500 steps) and flow matching (100 steps)

Model	Bond Length	Bond Angles	Dihedral	Median $\Delta E_{\text{relax}}$	Mean $\Delta E_{\text{relax}}$
GEOM-Drugs	0.0000	0.00	7.2e-3	0.00	1.0e-3
EQGAT-diff	0.0076	0.95	7.98	6.36	11.06
Megalodon-quick	0.0085	0.88	7.28	5.78	9.74
Megalodon	<b>0.0061</b>	<b>0.66</b>	<b>5.42</b>	<b>3.17</b>	<b>5.71</b>
SemlaFlow	0.0309	2.03	6.01	32.96	93.13
Megalodon-flow <sup>a</sup>	0.0112	0.930	5.63	5.90	8.61
Megalodon-flow <sup>b</sup>	<b>0.0101</b>	<b>0.79</b>	<b>4.07</b>	<b>4.60</b>	<b>6.10</b>

<sup>a</sup> 100-step evaluation (SemlaFlow-compatible).

<sup>b</sup> 500-step evaluation (diffusion-aligned).

**Analysis** We see that for both diffusion and flow matching, Megalodon is better than its prior counterparts. Overall, Megalodon trained with diffusion performs best with roughly 2-10x lower median energy when compared to prior generative models. Notably, our model’s median relaxation energy difference  $\Delta E_{\text{relax}}$  is around 3 kcal/mol, which approaches the thermally relevant interval of 2.5 kcal/mol<sup>10</sup>. Megalodon is the first method to achieve such proximity to this thermodynamic threshold, marking a significant milestone in 3D molecular generation.

We note that while the loss function between FM and diffusion is identical in this instance, we see both flow models have an considerably larger bond length error, which translates to a similar energy performance gap. The xTB energy function is highly sensitive to bond lengths; small deviations in bond lengths can lead to significant increases in energy due to the steepness of the PES in these dimensions. A precise representation of bond lengths is crucial because inaccuracies directly impact the calculated energy, making bond length errors a primary contributor to higher relaxation energies in flow models.



When increasing the number of integration steps to 500, the gap between diffusion and flow matching narrows substantially, although flow matching still yields slightly higher relaxation energies. This behavior is fully aligned with prior analyses in SiT paper<sup>3</sup>, which attribute the remaining gap to the inherent advantages of stochastic interpolants in capturing fine-grained geometric structure. Importantly, because a flow-matching model is trained only once and can be evaluated at arbitrary numbers of integration steps, practitioners can directly trade off computation for geometric precision, making FM particularly flexible for downstream applications that demand tunable accuracy.

We also include a SOAP-based<sup>44</sup> comparison in the Appendix Sec. C.4, which shows that Megalodon’s generated structures closely track the geometric manifold of GEOM. Appendix Fig. 7 provides several examples of generated molecules, and we additionally include two SDF files containing structures produced by the Megalodon and Megalodon-flow models.

## 5 Conclusions

Megalodon enables the accurate generation of de novo 3D molecules with both diffusion and flow matching. We show with a scalable augmented transformer architecture that significant improvements are gained, especially when generating outside the region of support for the training distribution as it pertains to molecule sizes. Megalodon demonstrates the ability to achieve great accuracy in conditional structure generation due to being trained to generate complete molecules from scratch. We also introduce more interpretable quantum mechanical energy benchmarks that are grounded in the original creation of the GEOM Drugs dataset.

While unconditional generation alone is rarely directly actionable in drug discovery, the proposed framework provides a strong foundation for a wide range of structure-based tasks. The architectural components introduced in Megalodon, together with its ability to generate full 3D geometries from scratch, can be leveraged for structure-based drug design (SBDD), conformer generation, fragment or pharmacophore-guided molecule construction, and other spatially conditioned design problems. In addition, the flexibility of our implementation



allows the codebase to be rapidly adapted to any of these settings, enabling systematic comparisons of diffusion and flow-matching approaches on task-specific objectives. In this sense, Megalodon functions not only as an unconditional generator but as a practical foundation model whose learned spatial priors can be transferred across downstream molecular design workflows.

While we show that Megalodon performs well across a variety of 3D de novo molecules tasks there are still some limitations that are worthy of discussion.

Megalodon like Le et al.<sup>19</sup> and the prior edge prediction generative models before it relies on maintaining  $N^2$  edge features, which is quite expensive. Recently<sup>20</sup> was able to avoid this issue for a majority of the model architecture by fusing the edge and atom features, but this creates a trade-off between model speed and accuracy. Our ablations show that the larger edge features are critical for strong energy performance, so it is still an open question for how to best deal with discrete edge types as each atom can have a maximum of 6 bonds at a time, so is needing to model all  $N$  potential pairings at all times really necessary? We leave future work to explore this in greater depth.

As discussed herein, the existing 3D molecule generation benchmarks are quite limited. A common theme that has been discussed in prior work<sup>19,20</sup>. While we make strides in expanding the field of view of de novo design and energy-based benchmarks. More work needs to be done to measure important qualities, as even for common conditional design benchmarks, metrics such as QED are not meaningful in practice, and even more complex properties like protein-ligand binding affinity can be directly optimized for with non-3D structure-based methods<sup>45</sup>. For these reasons, we looked to explore conditional structure generation, but across the board, small molecule benchmarking is a current field-wide limitation when compared to the current drug discovery practices.

A general limitation of current 3D molecular generative models, including Megalodon, is that they inherit both the representational constraints of their architectures and the chemical coverage of the training data. Because our model uses a one-hot atom-type representation,



it cannot generalize to elements not present in the training set, and this interacts with the limitations of GEOM-Drugs, a gas-phase semiempirical dataset with a relatively narrow medicinal-chemistry bias. Unlike 2D or SMILES-based language models, which benefit from vast and chemically diverse datasets spanning a wide range of elements, the amount of available high-quality optimized 3D conformer data is far more limited. As a result, 3D generators tend to remain close to the GEOM-Drugs distribution. To improve downstream applicability, future work may involve training on more chemically diverse 3D datasets with broader element coverage and developing larger, more representative 3D molecular datasets that extend beyond the chemical space currently available.

Overall, we explore the similarities and differences between flow matching and diffusion while improving 3D molecule design.

## 6 Acknowledgement

O.I. acknowledges support by the NSF grant CHE-2154447. This work used Expanse at SDSC and Delta at NCSA through allocation CHE200122 from the Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support (ACCESS) program, which is supported by NSF grants #2138259, #2138286, #2138307, #2137603, and #2138296.

## 7 Data Availability Statement

The implementation of the Megalodon model, along with pre-trained weights and data processing scripts, is available at the GitHub repository: <https://github.com/NVIDIADigitalBio/megalodon>. A persistent, citable snapshot of the codebase is archived on Zenodo under the DOI: <https://doi.org/10.5281/zenodo.17945981>. The provided scripts enable both reproducible model training from scratch and sampling from existing models.



## References

- (1) Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; Aspuru-Guzik, A. Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Central Science* **2018**, *4*, 268–276, PMID: 29532027.
- (2) Peebles, W.; Xie, S. Scalable Diffusion Models with Transformers. *arXiv preprint arXiv:2212.09748* **2022**,
- (3) Ma, N.; Goldstein, M.; Albergo, M. S.; Boffi, N. M.; Vanden-Eijnden, E.; Xie, S. SiT: Exploring Flow and Diffusion-based Generative Models with Scalable Interpolant Transformers. 2024; <https://arxiv.org/abs/2401.08740>.
- (4) Schneuing, A.; Du, Y.; Harris, C.; Jamasb, A.; Igashov, I.; Du, W.; Blundell, T.; Lió, P.; Gomes, C.; Welling, M.; Bronstein, M.; Correia, B. Structure-based Drug Design with Equivariant Diffusion Models. *arXiv preprint arXiv:2210.13695* **2022**,
- (5) Corso, G.; Stärk, H.; Jing, B.; Barzilay, R.; Jaakkola, T. DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking. International Conference on Learning Representations (ICLR). 2023.
- (6) Abramson, J. et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature* **2024**, *630*, 493–500.
- (7) Durairaj, J. et al. PLINDER: The protein-ligand interactions dataset and evaluation resource. *bioRxiv* **2024**,
- (8) Hoogeboom, E.; Satorras, V. G.; Vignac, C.; Welling, M. Equivariant diffusion for molecule generation in 3d. International conference on machine learning. 2022; pp 8867–8887.





- (9) Esser, P.; Kulal, S.; Blattmann, A.; Entezari, R.; Müller, J.; Saini, H.; Levi, Y.; Lorenz, D.; Sauer, A.; Boesel, F.; Podell, D.; Dockhorn, T.; English, Z.; Rombach, R. Scaling Rectified Flow Transformers for High-Resolution Image Synthesis. Forty-first International Conference on Machine Learning. 2024.
- (10) Axelrod, S.; Gómez-Bombarelli, R. GEOM, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data* **2022**, *9*, 185.
- (11) Pracht, P.; Grimme, S.; Bannwarth, C.; Bohle, F.; Ehlert, S.; Feldmann, G.; Gorges, J.; Müller, M.; Neudecker, T.; Plett, C.; others CREST—A program for the exploration of low-energy molecular chemical space. *The Journal of Chemical Physics* **2024**, *160*.
- (12) Bannwarth, C.; Ehlert, S.; Grimme, S. GFN2-xTB—An Accurate and Broadly Parametrized Self-Consistent Tight-Binding Quantum Chemical Method with Multipole Electrostatics and Density-Dependent Dispersion Contributions. *Journal of Chemical Theory and Computation* **2019**, *15*, 1652–1671.
- (13) Halgren, T. A. Merck molecular force field. I. Basis, form, scope, parameterization, and performance of MMFF94. *Journal of Computational Chemistry* **1996**, *17*, 490–519.
- (14) Satorras, V. G.; Hoogeboom, E.; Welling, M. E (n) equivariant graph neural networks. International conference on machine learning. 2021; pp 9323–9332.
- (15) Song, Y.; Gong, J.; Xu, M.; Cao, Z.; Lan, Y.; Ermon, S.; Zhou, H.; Ma, W.-Y. Equivariant Flow Matching with Hybrid Probability Transport for 3D Molecule Generation. Thirty-seventh Conference on Neural Information Processing Systems. 2023.
- (16) O’Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R. Open Babel: An open chemical toolbox. *Journal of Cheminformatics* **2011**, *3*, 33.



- (17) Walters, P. Generative Molecular Design Isn't As Easy As People Make It Look. 2024; <https://practicalcheminformatics.blogspot.com/2024/05/generative-molecular-design-isnt-as.html>.
- (18) Vignac, C.; Osman, N.; Toni, L.; Frossard, P. MiDi: Mixed Graph and 3D Denoising Diffusion for Molecule Generation. *arXiv preprint arXiv:2302.09048* **2023**,
- (19) Le, T.; Cremer, J.; Noe, F.; Clevert, D.-A.; Schütt, K. T. Navigating the Design Space of Equivariant Diffusion-Based Generative Models for De Novo 3D Molecule Generation. The Twelfth International Conference on Learning Representations. 2024.
- (20) Irwin, R.; Tibo, A.; Janet, J. P.; Olsson, S. Efficient 3D Molecular Generation with Flow Matching and Scale Optimal Transport. 2024.
- (21) Xu, M.; Powers, A.; Dror, R.; Ermon, S.; Leskovec, J. Geometric Latent Diffusion Models for 3D Molecule Generation. 2023; <https://arxiv.org/abs/2305.01140>.
- (22) Pinheiro, P. O.; Rackers, J.; Kleinhenz, J.; Maser, M.; Mahmood, O.; Watkins, A. M.; Ra, S.; Sresht, V.; Saremi, S. 3D molecule generation by denoising voxel grids. 2024; <https://arxiv.org/abs/2306.07473>.
- (23) Song, Y.; Gong, J.; Qu, Y.; Zhou, H.; Zheng, M.; Liu, J.; Ma, W.-Y. Unified Generative Modeling of 3D Molecules via Bayesian Flow Networks. 2024; <https://arxiv.org/abs/2403.15441>.
- (24) Lipman, Y.; Chen, R. T. Q.; Ben-Hamu, H.; Nickel, M.; Le, M. Flow Matching for Generative Modeling. The Eleventh International Conference on Learning Representations. 2023.
- (25) Albergo, M. S.; Boffi, N. M.; Vanden-Eijnden, E. Stochastic Interpolants: A Unifying Framework for Flows and Diffusions. 2023; <https://arxiv.org/abs/2303.08797>.



- (26) Song, Y.; Sohl-Dickstein, J.; Kingma, D. P.; Kumar, A.; Ermon, S.; Poole, B. Score-Based Generative Modeling through Stochastic Differential Equations. *International Conference on Learning Representations (ICLR)*. 2021.
- (27) Ho, J.; Jain, A.; Abbeel, P. Denoising Diffusion Probabilistic Models. *arXiv preprint arxiv:2006.11239* **2020**,
- (28) Austin, J.; Johnson, D. D.; Ho, J.; Tarlow, D.; Van Den Berg, R. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems* **2021**, *34*, 17981–17993.
- (29) Campbell, A.; Yim, J.; Barzilay, R.; Rainforth, T.; Jaakkola, T. Generative Flows on Discrete State-Spaces: Enabling Multimodal Flows with Applications to Protein Co-Design. *arXiv preprint arXiv:2402.04997* **2024**,
- (30) Tong, A.; Malkin, N.; Huguet, G.; Zhang, Y.; Rector-Brooks, J.; Fatras, K.; Wolf, G.; Bengio, Y. Improving and generalizing flow-based generative models with minibatch optimal transport. *arXiv preprint arXiv:2302.00482* **2023**,
- (31) Vaswani, A. Attention is all you need. *Advances in Neural Information Processing Systems* **2017**,
- (32) Brown, T. B. et al. Language Models are Few-Shot Learners. 2020.
- (33) Lin, Z.; Akin, H.; Rao, R.; Hie, B.; Zhu, Z.; Lu, W.; Smetanin, N.; Verkuil, R.; Kabeli, O.; Shmueli, Y.; others Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **2023**, *379*, 1123–1130.
- (34) Hayes, T. et al. Simulating 500 million years of evolution with a language model. *bioRxiv* **2024**,
- (35) Chen, T.; Zhang, R.; Hinton, G. Analog bits: Generating discrete data using diffusion models with self-conditioning. *arXiv preprint arXiv:2208.04202* **2022**,



- (36) Yim, J.; Trippe, B. L.; De Bortoli, V.; Mathieu, E.; Doucet, A.; Barzilay, R.; Jaakkola, T. SE (3) diffusion model with application to protein backbone generation. *arXiv preprint arXiv:2302.02277* **2023**,
- (37) Stärk, H.; Jing, B.; Barzilay, R.; Jaakkola, T. Harmonic Self-Conditioned Flow Matching for Multi-Ligand Docking and Binding Site Design. 2024; <https://arxiv.org/abs/2310.05764>.
- (38) Békés, M.; Langley, D. R.; Crews, C. M. PROTAC targeted protein degraders: the past is prologue. *Nature Reviews Drug Discovery* **2022**, *21*, 181–200.
- (39) Jing, B.; Corso, G.; Chang, J.; Barzilay, R.; Jaakkola, T. Torsional Diffusion for Molecular Conformer Generation. *arXiv preprint arXiv:2206.01729* **2022**,
- (40) Xu, M.; Yu, L.; Song, Y.; Shi, C.; Ermon, S.; Tang, J. GeoDiff: A Geometric Diffusion Model for Molecular Conformation Generation. International Conference on Learning Representations. 2022.
- (41) Corso, G.; Xu, Y.; De Bortoli, V.; Barzilay, R.; Jaakkola, T. Particle guidance: non-iid diverse sampling with diffusion models. *arXiv preprint arXiv:2310.13102* **2023**,
- (42) Wang, Y.; Elhag, A. A.; Jaitly, N.; Susskind, J. M.; Bautista, M. A. Swallowing the Bitter Pill: Simplified Scalable Conformer Generation. 2023.
- (43) Nikitin, F.; Dunn, I.; Koes, D. R.; Isayev, O. GEOM-drugs revisited: toward more chemically accurate benchmarks for 3D molecule generation. *Digital Discovery* **2025**, *4*, 3282–3291.
- (44) Bartók, A. P.; Kondor, R.; Csányi, G. On representing chemical environments. *Physical Review B—Condensed Matter and Materials Physics* **2013**, *87*, 184115.
- (45) Reidenbach, D. EvoSBDD: Latent Evolution for Accurate and Efficient Structure-Based



Drug Design. ICLR 2024 Workshop on Machine Learning for Genomics Explorations. 2024.

- (46) Vincent, P. A Connection Between Score Matching and Denoising Autoencoders. *Neural Computation* **2011**, *23*, 1661–1674.
- (47) Henry, A.; Dachapally, P. R.; Pawar, S. S.; Chen, Y. Query-Key Normalization for Transformers. Findings of the Association for Computational Linguistics: EMNLP 2020. 2020; p 4246–4253.
- (48) Peng, X.; Guan, J.; Liu, Q.; Ma, J. MolDiff: Addressing the Atom-Bond Inconsistency Problem in 3D Molecule Diffusion Generation. Proceedings of the 40th International Conference on Machine Learning. 2023; pp 27611–27629.



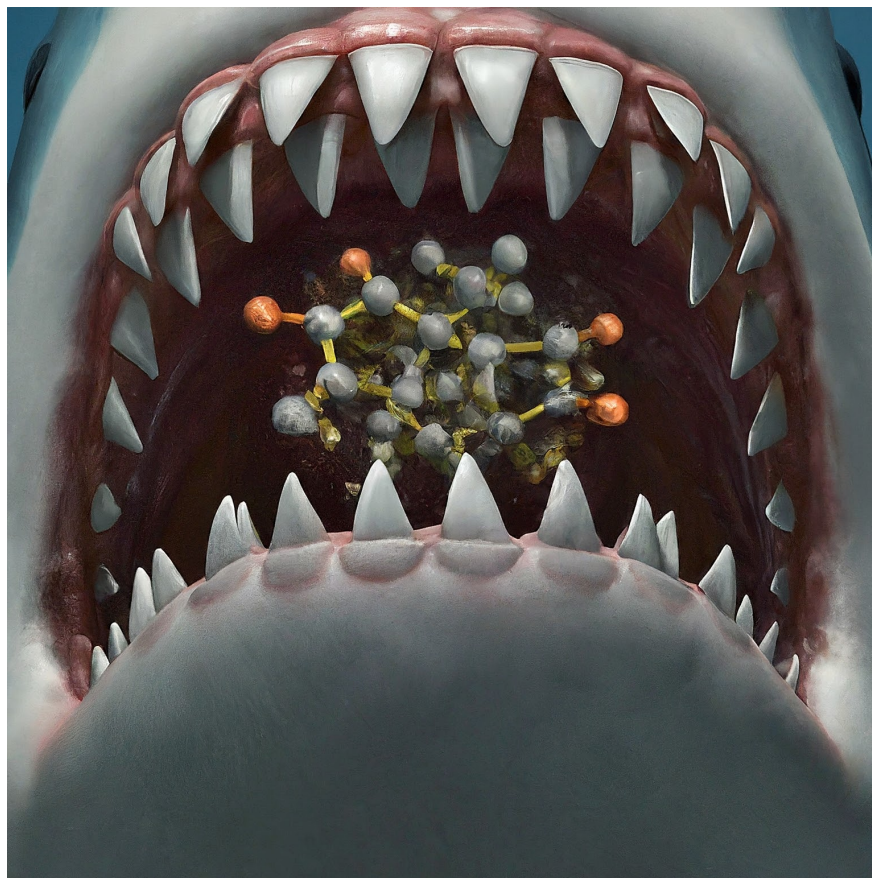


Figure 4: Megalodon molecule generation dynamics generated with Imagen 2

## A Equating Continuous Gaussian Diffusion and Flow Matching

A part of our work was to explore when to use diffusion versus flow matching and what the empirical differences are. We show below that from a training perspective in the continuous domain, they can be made equivalent.

It can be shown that this objective under the Gaussian setting is a time-dependent scalar multiple of the standard denoising objective explored in Ho et al.<sup>27</sup>. Let's insert equation 1b



into the flow matching objective

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} \left\| \mathbf{v}_\theta(t, \mathbf{x}_t) - \dot{\alpha}(t)\epsilon - \frac{\dot{\beta}(t)}{\beta(t)}(\mathbf{x}_t - \alpha(t)\epsilon) \right\|^2. \quad (11)$$

where the dot notation denotes the partial time derivative.

Now we see that we can construct an objective that is similar to the “noise prediction” objective that is used in diffusion models:

$$\begin{aligned} \mathcal{L}_{\text{CFM}}(\theta) &= \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} \left\| \mathbf{v}_\theta(t, \mathbf{x}_t) - \dot{\alpha}(t)\epsilon - \frac{\dot{\beta}(t)}{\beta(t)}(\mathbf{x}_t - \alpha(t)\epsilon) \right\|^2 \\ &= \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} \left\| \mathbf{v}_\theta(t, \mathbf{x}_t) - \frac{\dot{\beta}(t)}{\beta(t)}\mathbf{x}_t - \underbrace{\left( \dot{\alpha}(t) - \frac{\dot{\beta}(t)}{\beta(t)}\alpha(t) \right)}_{=:s(t)} \epsilon \right\|^2 \\ &= \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} s^2(t) \left\| \underbrace{\frac{1}{s(t)} \left( \mathbf{v}_\theta(t, \mathbf{x}_t) - \frac{\dot{\beta}(t)}{\beta(t)}\mathbf{x}_t \right)}_{=: \epsilon_\theta(t, \mathbf{x}_t)} - \epsilon \right\|^2 \\ &= \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I}), \mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x}_1)} s^2(t) \left\| \epsilon_\theta(t, \mathbf{x}_t) - \epsilon \right\|^2. \end{aligned} \quad (12)$$

We see that the resulting mean squared error of noise prediction is the original core loss derived in Ho et al.<sup>27</sup>. This allows us to choose time-dependent scalars via the time distribution itself or the noise or variance schedule to equate the CFM and Diffusion objectives.

In the generative modeling case, we interpolate between a data distribution and a Gaussian density, meaning all data-conditional paths are Gaussian. In that special case, we can, in fact, easily extract the score function from the regular flow matching objective, and we get stochastic sampling for free. We know that  $\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_1)$  follows Gaussian probability paths. Based on equation 1, we know that

$$\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_1) = \mathcal{N}(\mathbf{x}_t; \beta(t)\mathbf{x}_1, \alpha^2(t)\mathbf{I}). \quad (13)$$





Let's calculate the score:

$$\begin{aligned}
 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_1) &= -\nabla_{\mathbf{x}_t} \frac{(\mathbf{x}_t - \beta(t)\mathbf{x}_1)^2}{2\alpha^2(t)} \\
 &= -\frac{\mathbf{x}_t - \beta(t)\mathbf{x}_1}{\alpha^2(t)} \\
 &= -\frac{\boldsymbol{\epsilon}}{\alpha(t)},
 \end{aligned} \tag{14}$$

where we used equation 1 in the last step. We can solve this for  $\boldsymbol{\epsilon}$  and insert into the reparametrized  $\mathcal{L}_{\text{CFM}}$  in equation 12 and see that we obtain denoising score matching,<sup>46</sup> which implies that  $\boldsymbol{\epsilon}_\theta(t, \mathbf{x}_t)$ , or analogously  $\mathbf{v}_\theta(t, \mathbf{x}_t)$  via their connection, learn a model of the marginal score  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ .

Specifically, we have alternatively

$$\boldsymbol{\epsilon}_\theta(t, \mathbf{x}_t) = -\alpha(t) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t), \tag{15}$$

$$\mathbf{v}_\theta(t, \mathbf{x}_t) = -\alpha(t) \frac{\beta(t)\dot{\alpha}(t) - \dot{\beta}(t)\alpha(t)}{\beta(t)} \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \frac{\dot{\beta}(t)}{\beta(t)} \mathbf{x}_t. \tag{16}$$

We note that these equations only hold for a Gaussian prior without optimal transport.

## B Megalodon Architecture

### B.1 Architecture

As described in Fig. 1, Megalodon consists of  $N$  augmented transformer blocks that consist of a Fused Invariant Transformer (FiT) block and a structure layer. We refer to it as Megalodon and Megalodon Quick, as we maintain the same number of layers but weaken the representation size to achieve 2x sampling speeds compared to the base model.





Table 4: Comparison of Megalodon Quick and Megalodon hyperparameter configurations.

Parameter	Megalodon Quick	Megalodon
Invariant Edge Feature Dimension	64	256
Invariant Node Feature Dimension	256	256
Number of Vector Features	64	128
Number of Layers	10	10
Number of FiT Attention Heads	4	4
Distance Feature Size	16	128

### B.1.1 Input/Output Layers

Megalodon takes the input molecules structures and projects them into a  $N \times D$  tensor where  $D$  is the number of vector features. After all augmented transformer blocks, the predicted structure is projected back down to  $N \times 3$ .

Similarly, the input discrete components are projected from their one hot variable to a hidden dimension size. The bonds leverage the edge feature size, and the atom types and charges use the node feature size. After all augmented transformer blocks, final prediction heads are applied to project the values back into their respective vocabulary size for discrete prediction.

### B.1.2 Fused Invariant Transformer Block

Our Fused Invariant Transformer (FiT) block has several key differences compared to other diffusion transformers<sup>2</sup>.

- Rather than just operating over the discrete atom type features  $H$ , we operate over a fused feature  $m = \frac{1}{N} \sum_{i,j \in N} f(h_{\text{norm},i,j}, h_{\text{norm},i,j}, e_{\text{norm},i,j}, \text{distance}_{i,j})$  where  $h_{\text{norm}}$  and  $e_{\text{norm}}$  are the outputs of the time conditioned adaptive layer norm for the atom type and edge type features. The distance features are the concatenation of scalar distances and dot products. We note that this fusing step is important to ground the simple equivariant structure update layer to the transformer trunk.
- We employ query key normalization<sup>34,47</sup>.



- The multi-head attention is applied to  $m$  to produce  $\text{mha\_out}$  and then used directly in the standard feed-forward to produce  $H_{out}$ . To create  $E_{out}$  we mimic the same steps but use  $f(\text{mha\_out}_i + \text{mha\_out}_j)$  for all edges between nodes  $i$  and  $j$ . Our feed-forward is the standard SWiGLU layer with a feature projection of 4. We note that this feed-forward for edge features is the most expensive component of the model, which is why Megalodon-quick is designed the way it is.

### B.1.3 Structure Layer

Following Schneuing et al.<sup>4</sup>, the structure layer of Megalodon consists of a single EGNN layer with a positional and cross-product update component. Before this operation, all inputs are normalized to prevent value and gradient explosion, a common problem faced when using EGNNs<sup>14</sup>. The invariant features use standard layer norm, whereas the equivariant features use an E3Norm<sup>18</sup>.

$$\mathbf{x}_i^{l+1} = \mathbf{x}_i^l + \sum_{j \neq i} \frac{\mathbf{x}_i^l - \mathbf{x}_j^l}{d_{ij} + 1} \phi_x^d(\mathbf{h}_i^l, \mathbf{h}_j^l, d_{ij}^2, a_{ij}) + \frac{(\mathbf{x}_i^l - \bar{\mathbf{x}}^l) \times (\mathbf{x}_j^l - \bar{\mathbf{x}}^l)}{\|(\mathbf{x}_i^l - \bar{\mathbf{x}}^l) \times (\mathbf{x}_j^l - \bar{\mathbf{x}}^l)\| + 1} \phi_x^\times(\mathbf{h}_i^l, \mathbf{h}_j^l, d_{ij}^2, a_{ij}), \quad (17)$$

## B.2 Compute and Data Requirements

Similar to Le et al.<sup>19</sup>, we use MiDi’s adaptive dataloader for GEOM DRUGS with a batch cost of 200. We note that the adaptive logic randomly selects one molecule and fills in the batch with similar-sized molecules, tossing any molecules selected that do not fit the adaptive criteria out of the current epoch’s available molecules. As a result, an epoch in this setting does not hold the standard connotation as time for the model to see each training data point. We use this dataloader as it was used by prior methods and we felt it important to standardize the data to best create a fair comparison. Megalodon-quick is trained on 4 NVIDIA A100 GPUs for 250 epochs. Megalodon was trained on 8 A100 GPUs for 250



epochs, taking roughly 2 days.

Megalodon-flow was trained using the data splits and adaptive data loader from Irwin et al.<sup>20</sup>, which does not discard molecules though was prefiltered to only include molecules with  $\leq 72$  atoms. It was trained for 200 epochs on 8 A100 NVIDIA GPUs.

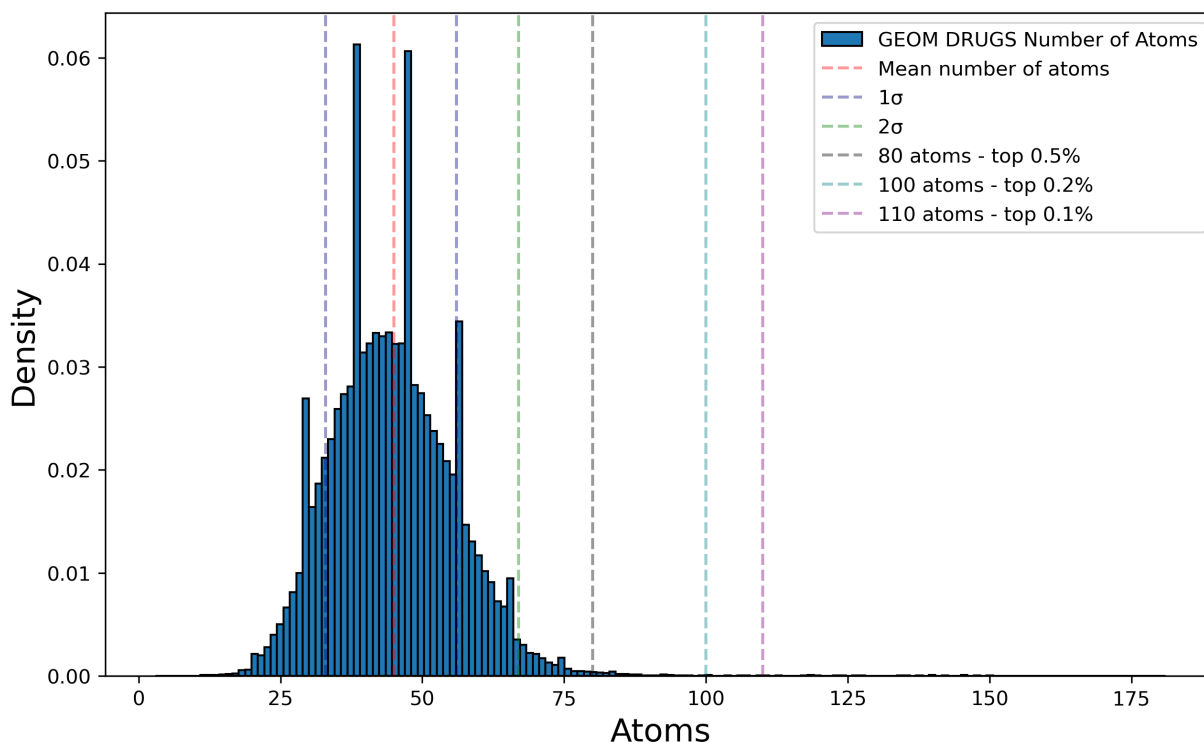


Figure 5: Distribution of molecule sizes

## C Extended Unconditional Generation

### C.1 Performance on QM9

There are three popular datasets of 3D molecular structures commonly used to benchmark generative models: **QM9**, **PubChem3D**, and **GEOM Drugs**. In this work, we primarily focus on GEOM Drugs because PubChem3D provides relatively low-quality 3D structures that do not necessarily reflect low-energy conformations. Nevertheless, QM9 remains a well-established and frequently used small-scale benchmark, despite the fact that its me-



Table 5: Measuring Unconditional Molecule Generation: 2D and 3D benchmarks on QM9 dataset. \* Denotes taken from MiDi.

Model	Steps	2D Topological ( $\uparrow$ )			3D Distributional ( $\downarrow$ )	
		Atom Stab.	Mol Stab.	Validity	Bond Angle	Dihedral
MiDi*	500	0.998	0.975	0.979	0.670	—
EQGAT-diff <sub>disc</sub> <sup>r0</sup>	500	0.998	0.977	0.979	0.365	0.815
Megalodon-quick	500	<b>0.999</b>	<b>0.986</b>	<b>0.988</b>	<b>0.241</b>	0.662
Megalodon	500	<b>0.999</b>	<b>0.986</b>	0.987	0.422	<b>0.637</b>
SemlaFlow	100	<b>0.999</b>	<b>0.986</b>	0.986	0.775	1.194
Megalodon-flow	100	0.998	0.973	0.976	0.804	0.970

dian molecular size is unrealistically small (approximately 20 atoms). For completeness, we therefore report results on QM9 as shown in Table 5.

We trained SemlaFlow and EQGAT-Diff from scratch while using MiDi results from the original paper. In their original codebases, both MiDi and EQGAT-Diff were trained on molecular representations with three bond types: single, double, and triple, whereas Semla included aromatic bonds. We found that the inclusion of aromatic bonds negatively impacted molecular stability metrics, even though any molecule can be represented without them in a molecular graph using the Kekulized form. To ensure comparability, we trained all models using only single, double, and triple bonds.

In Table 5 we observe that Megalodon-quick achieves the best overall performance on QM9 in terms of both 2D (topological) and 3D (distributional) metrics. This outcome makes intuitive sense because Megalodon-quick is a more lightweight variant, and smaller models often suit datasets of reduced scale, such as QM9, more effectively. In contrast, Megalodon-flow appears to be too large for this dataset, leading to slightly weaker performance; however, we include it here for completeness and consistency with results on GEOM Drugs.

While flow matching benefits from fewer integration steps (here 100 steps), the diffusion-based objective with more denoising steps (in our case, 500) ultimately achieves stronger 3D quality metrics. Thus, the trade-off between computational efficiency (fewer steps) and generative fidelity (more steps) is highlighted once again in this smaller-scale setting.



## C.2 Unconditional Ablations

Table 6: Measuring Unconditional Molecule Generation: 2D topological and 3D distributional benchmarks.

Model	Steps	2D Topological ( $\uparrow$ )			3D Distributional ( $\downarrow$ )	
		Atom Stab.	Mol Stab.	Connected Validity	Bond Angle	Dihedral
EDM + OpenBabel*	1000	0.978	0.403	0.363	–	–
MolDiff taken from Peng et al. <sup>48</sup>	1000	–	–	0.739	–	–
GeoBFN taken from Song et al. <sup>23</sup>	2000	0.862	0.917	–	–	–
MiDi*	500	0.997	0.897	0.705	–	–
EQGAT-diff <sub>disc</sub> <sup>x0</sup>	100	0.996	0.891	0.768	1.772	3.514
EQGAT-diff <sub>disc</sub> <sup>x0</sup>	500	0.998	0.935	0.830	0.858	2.860
EGNN + cross product	500	0.982	0.713	0.223	14.778	17.003
Megalodon-quick	500	0.998	0.961	0.900	0.689	2.383
Megalodon	100	0.998	0.939	0.817	0.871	3.367
Megalodon	500	<b>0.999</b>	0.977	0.927	0.461	<b>1.231</b>
SemlaFlow	20	0.997	0.962	0.875	2.188	3.173
SemlaFlow	100	0.998	0.979	0.920	1.274	1.934
Megalodon-flow	20	0.998	0.937	0.852	2.695	3.892
Megalodon-flow	100	<b>0.999</b>	0.98	0.944	1.286	2.379
Megalodon-flow <sup>†</sup>	100	0.997	0.990	0.948	0.976	2.085
Megalodon-flow	500	<b>0.999</b>	<b>0.991</b>	<b>0.965</b>	<b>0.438</b>	1.646

<sup>†</sup> Uses the original Semla-Flow preprocessing (variance-1 scaling; > 72-atom molecules removed).

\* Denotes taken from EQGAT-Diff.

We include each primary model in its base form as well as with 5x fewer inference steps. The flow models do not have to be retrained as they were trained to learn a continuous vector field, whereas the diffusion models must be retrained due to the change in variance discretization in the forward diffusion process.

We also include EGNN + cross product which is similar to Megalodon except the transformer layers were replaced by the standard invariant and edge feature updates in Satorras et al.<sup>14</sup>. Prior methods exist that improve upon EDM + Open Babel and maintain that bonds are generated external to the model via Open Babel<sup>15</sup>. We do not include such methods in our comparison as, for the most part, public code with weights is not available, and Open Babel introduces significant bias and errors, which make evaluating the model difficult<sup>15,17</sup>.

Open Babel, while a powerful tool for molecular manipulation and conversion, can intro-



duce several potential errors, particularly in the context of bond assignment and 3D structure generation. Some common errors include:

- **Incorrect bond orders:** Open Babel often assigns bond orders based on geometric heuristics or atom types, which can lead to inaccuracies, especially in complex or exotic molecules where bond orders are not trivial.
- **Geometric distortions:** When converting between different formats or generating 3D coordinates, Open Babel may generate suboptimal or distorted geometries, especially if the input structure is incomplete or poorly defined.
- **Protonation state assumptions:** Open Babel may incorrectly infer or standardize protonation states, which can lead to chemical inaccuracies, especially in sensitive systems such as drug-like molecules or biologically active compounds.
- **Ambiguous aromaticity:** Open Babel can sometimes misinterpret or incorrectly assign aromaticity, which can lead to an incorrect representation of the molecular structure.
- **Missing stereochemistry:** While converting or generating structures, stereochemistry can be incorrectly assigned or lost altogether, affecting the overall molecular properties.

### C.3 3D Distributional Metrics

To evaluate the geometric fidelity of the generated molecules, we compute the Wasserstein-1 distance between the generated and target distributions of bond angles, following the methodology of<sup>19</sup>. The overall bond angle metric is defined as:

$$W_{\text{angles}} = \sum_{y \in \text{atom types}} p(y) \cdot W_1(\hat{D}_{\text{angle}}(y), D_{\text{angle}}(y)),$$

where  $p(y)$  is the probability of atom type  $y$ ,  $W_1$  denotes the Wasserstein-1 distance,  $\hat{D}_{\text{angle}}(y)$  is the bond angle distribution for atom type  $y$  in the generated data, and  $D_{\text{angle}}(y)$  is the corresponding distribution in of test set.



Similarly, for torsion angles, the metric is calculated as:

$$W_{\text{torsions}} = \sum_{y \in \text{bond types}} p(y) \cdot W_1(\hat{D}_{\text{torsion}}(y), D_{\text{torsion}}(y)),$$

where  $p(y)$  is the probability of bond type  $y$ ,  $\hat{D}_{\text{torsion}}(y)$  is the torsion angle distribution for bond type  $y$  in the generated data, and  $D_{\text{torsion}}(y)$  is the corresponding distribution in the test set. Since we utilized RDKit to identify torsions, the torsional distribution difference was computed only for valid molecules.

## C.4 SOAP-Based Structural Similarity Between Generated and GEOM Conformers

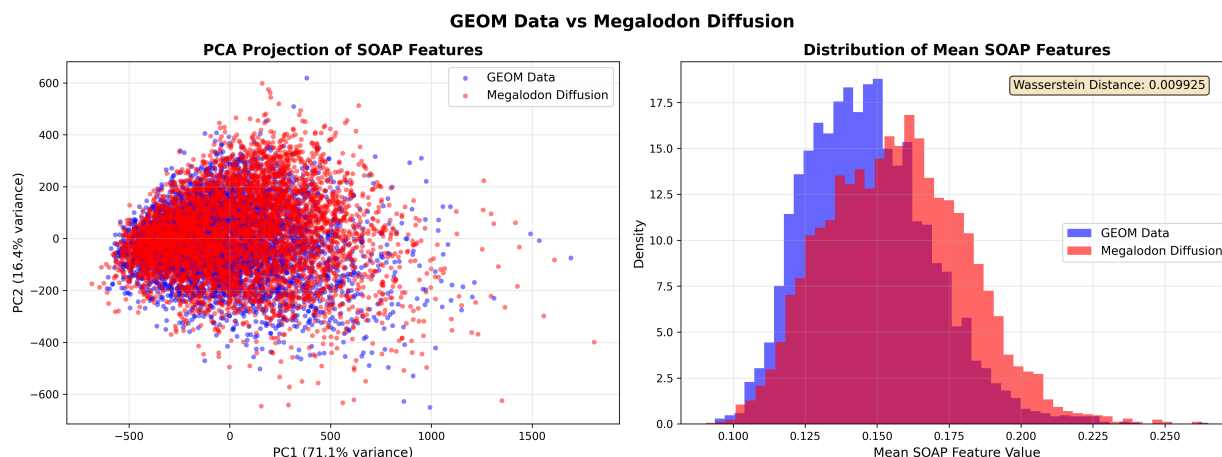


Figure 6: Comparison of SOAP feature distributions between GEOM conformers and Megalodon-generated molecules. Left: PCA projection of SOAP embeddings shows substantial overlap between the two structural distributions. Right: Histogram of mean SOAP feature values with the corresponding Wasserstein distance.

We also examined how closely Megalodon-generated structures resemble the broader GEOM conformer space by comparing SOAP descriptors for a 5k-molecule subset of GEOM and 5k molecules generated by our model. SOAP encodes local atomic environments in a geometry-aware manner, independent of bonding assignments. As shown in Appendix Fig. 6, a PCA projection of the SOAP embeddings reveals substantial overlap between the



two distributions. The Wasserstein distance between their mean SOAP feature distributions is small, indicating that Megalodon captures the global geometric statistics of the dataset. This analysis confirms that our generated molecules occupy a similar region of 3D conformational space as the GEOM reference set.

## D Megalodon Molecule Visualization

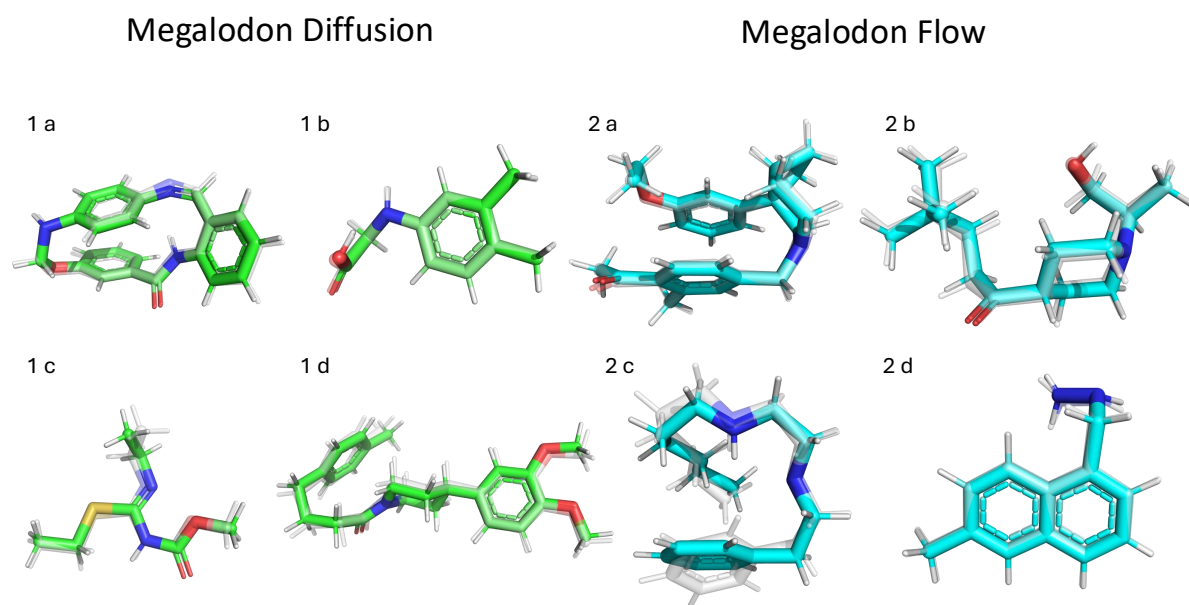


Figure 7: Examples of generated molecules using Megalodon: (1) Diffusion and (2) Flow Matching. Each generated molecule is displayed alongside its corresponding optimized structure (shown in transparent grey). The examples include small aromatic molecules (1b, 2d), molecules exhibiting pi-stacking interactions (1a, 2a), non-aromatic molecules (1c, 2b), and a molecule with a macrocycle (1a).





## Data Availability Statement

The implementation of the Megalodon model, along with pre-trained weights and data processing scripts, is available at the GitHub repository: <https://github.com/NVIDIADigital-Bio/megalodon>. A persistent, citable snapshot of the codebase is archived on Zenodo under the DOI: <https://doi.org/10.5281/zenodo.17945981>. The provided scripts enable both reproducible model training from scratch and sampling from existing models.

