JAMUN: Bridging Smoothed Molecular Dynamics and Score-Based Learning for Conformational Ensemble Generation

Ameya Daigavane* MIT ameyad@mit.edu Bodhi P. Vani* Prescient Design, Genentech vanib@gene.com Darcy Davidson
Prescient Design, Genentech
davidsd5@gene.com

Saeed Saremi Prescient Design, Genentech saremis@gene.com **Joshua A. Rackers**[†] Prescient Design, Genentech

Joseph Kleinhenz[†] Prescient Design, Genentech kleinhej@gene.com

Abstract

Conformational ensembles of protein structures are immensely important both for understanding protein function and drug discovery in novel modalities such as cryptic pockets. Current techniques for sampling ensembles such as molecular dynamics (MD) are computationally inefficient, while many recent machine learning methods do not transfer to systems outside their training data. We propose JAMUN which performs MD in a smoothed, noised space of all-atom 3D conformations of molecules by utilizing the framework of walk-jump sampling. JAMUN enables ensemble generation for small peptides at rates of an order of magnitude faster than traditional molecular dynamics. The physical priors in JAMUN enables transferability to systems outside of its training data, even to peptides that are longer than those originally trained on. Our model, code and weights are available at https://github.com/prescient-design/jamun.

1 Introduction

Proteins are inherently dynamic entities constantly in motion, and these movements can be vitally important. They are not well characterized as single structures as has traditionally been the case, but rather as ensembles of structures drawn from the Boltzmann distribution [30]. Protein dynamics is required for the function of most proteins, for instance the global tertiary structure motions for myglobin to bind oxygen and move it around the body [57], or the beta-sheet transition to a disordered strand for insulin to dissociate and find and bind to its receptor [5]. Similarly, drug discovery on protein kinases depends on characterizing kinase conformational ensembles [26]. In general the search for druggable 'cryptic pockets' requires understanding protein dynamics [15], and antibody design is deeply affected by conformational ensembles [23]. However, while machine learning (ML) methods for molecular structure prediction have experienced enormous success recently, ML methods for dynamics have yet to have similar impact. ML models for generating molecular ensembles are widely considered the 'next frontier' [11, 57, 91]. In this work, we present JAMUN (Walk-Jump Accelerated Molecular ensembles with Universal Noise), a generative ML model which advances this frontier by demonstrating improvements in both speed and transferability over previous approaches.

While the importance of protein dynamics is well-established, it can be exceedingly difficult to sufficiently sample large biomolecular systems. The most common sampling method is molecular

^{*}Equal contribution.

[†]Equal contribution.

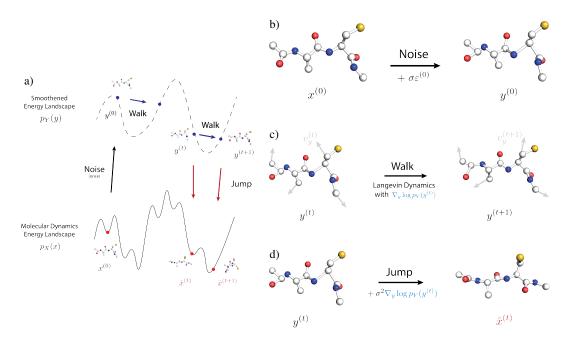


Figure 1: (a) Overview of walk-jump sampling in JAMUN. (b) Adding noise to an initial conformation $x^{(0)}$ to obtain $y^{(0)} \sim p_Y$. (c) One iteration of BAOAB-discretized Langevin dynamics (Equation 3 and Equation 68) starting from $y^{(t)} \sim p_Y$ leads to a new sample $y^{(t+1)} \sim p_Y$. (d) Denoising of $y^{(t)}$ according to Equation 7 gives us new samples $\hat{x}^{(t)}$.

dynamics (MD), which will sample the Boltzmann distribution in the limit of infinite sampling time, but is limited by the need for very short timesteps of 1-2 femtoseconds in the numerical integration scheme. Many important protein dynamic phenomena occur on the timescale of milliseconds. As described by Borhani and Shaw [10], simulating with this resolution is '... equivalent to tracking the advance and retreat of the glaciers of the last Ice Age – tens of thousands of years – by noting their locations each and every second.' Importantly, there is nothing fundamental about this small time-step limitation; it is an artifact of high-frequency motions, such as bond vibrations, that have little effect on protein ensembles [51]. Enhanced sampling methods have been developed in an attempt to accelerate sampling, but they often require domain knowledge about relevant collective variables, and, more importantly, do not address the underlying time-step problem [87].

A large number of generative models have been developed to address the sampling inefficiency problems of MD using machine learning, which we discuss in greater detail in Section 4. The key requirement is that of *transferability*: any model must be able to generate conformational ensembles for molecules that are significantly different from those in its training set. To benchmark this transferability, we focus on small peptides whose MD trajectories can be run to convergence in a reasonable amount of time, unlike those for much larger proteins [78].

Here, we propose a new method, JAMUN, that bridges molecular dynamics with score-based learning in a latent space. This physical prior enables JAMUN to transfer well – just like force fields for molecular dynamics can – to unseen systems. The key idea is to run Langevin molecular dynamics in a noisy 'latent' space $Y \in \mathbb{R}^{N \times 3}$, instead of the original space $X \in \mathbb{R}^{N \times 3}$ of all-atom 3D positions. Indeed, Y is constructed by adding a small amount of i.i.d Gaussian noise ε to X:

$$Y = X + \sigma \varepsilon \tag{1}$$

To run this Langevin MD over Y, the crucial component is the score function $\nabla_y \log p_Y(y)$, which needs to be modelled. (This is identical to how a force field needs to be parametrized in classical MD). Once the MD trajectory over Y is run, the resulting samples need to be mapped back to clean data via a denoising procedure.

This framework is actually mathematically described by Walk-Jump Sampling (WJS), as first introduced by Saremi and Hyvärinen [70]. WJS has been used in voxelized molecule generation [62, 63]

and protein sequence generation [24]. In particular, JAMUN corresponds to a SE(3)-equivariant Walk-Jump sampler of point clouds.

The WJS framework actually tells us that the score function $\nabla_y \log p_Y(y)$ can be used for denoising as well, removing the need to learn a separate model. Indeed, this mathematical connection is used to learn the score function, similar to the training of diffusion models. Unlike diffusion models such as DDPM [32], however, we only need to learn the score function at a single noise level. The choice of this noise level is important; we aim to simply smooth out the distribution enough to resolve sampling difficulties without fully destroying the information present in the data distribution.

In short, the score function $\nabla_y \log p_Y(y)$ is learned by adding noise to clean data x, and a denoising neural network is trained to recover the clean samples x from y. This denoiser defines the score function of the noisy manifold Y which we sample using Langevin dynamics (walk step) and allows us to periodically project back to the original data distribution (jump step). Crucially, the walk and jump steps are *decoupled* from each other.

Rather than starting over from an uninformative prior for each sample as is commonly done in diffusion [32, 79] and flow-matching [54, 47], JAMUN is able to simply denoise samples from the slightly noised distribution, enabling much greater sampling efficiency. Appendix A contains a comparison of JAMUN at a single noise level, against full diffusion which begins the sampling process at a high noise level.

We train JAMUN on a large dataset of MD simulations of small peptides. We demonstrate that this model can generalize to a holdout set of unseen peptides. In all of these cases, generation with JAMUN yields converged sampling of the conformational ensemble faster than MD with a standard force field, even outperforming several state-of-the-art baselines. These results suggest that this transferability is a consequence of retaining the physical priors inherent in MD data. Significantly, we find that JAMUN performs well even for peptides longer than the ones seen in the training set.

2 Methods

2.1 Representing Peptides as Point Clouds

Each point cloud of N atoms can be represented by the tuple (x,h) where $x \in \mathbb{R}^{N \times 3}$ represents the 3D coordinates of each of the N atoms and $h \in \mathbb{R}^D$ represents atom and covalent bonding information. h can be easily computed from the amino acid sequence for each peptide, and hence is not learned or sampled. For clarity of presentation, we omit the conditioning on h in the distributions and models below. We discuss how our model uses h in Section 2.4.

At sampling time, we assume access to an initial sample $x^{(0)} \in \mathbb{R}^{N \times 3}$ sampled from the clean data distribution p_X . Similarly to how MD simulations of small peptides are commonly seeded, we use the sequence command in the LEaP program packaged with the Amber force fields to procedurally generate $x^{(0)}$. In theory, $x^{(0)}$ could also be obtained from experimental data, such as crystallized structures from the Protein Data Bank [9]. We plan to explore this approach in the future to seed JAMUN simulations for larger proteins.

2.2 Walk-Jump Sampling

JAMUN operates by performing walk-jump sampling on molecular systems represented as 3D point clouds. A conceptual overview of the process is illustrated in Figure 1.

Given the initial sample $x^{(0)} \sim p_X$, walk-jump sampling performs the following steps:

1. **Noise** the initial structure $x^{(0)}$ to create the initial sample $y^{(0)}$ from the noisy data distribution p_Y (Figure 1a):

$$y^{(0)} = x^{(0)} + \sigma \varepsilon^{(0)}, \text{ where } \varepsilon^{(0)} \sim \mathcal{N}(0, \mathbb{I}_{N \times 3}).$$
 (2)

2. Walk to obtain samples $y^{(1)}, \dots, y^{(N)}$ from p_Y using Langevin dynamics which conists of numerically solving the following Stochastic Differential Equation (SDE) (Figure 1b):

$$dy = v_y dt, (3)$$

$$dv_y = \nabla_y \log p_Y(y)dt - \gamma v_y dt + M^{-\frac{1}{2}} \sqrt{2} dB_t, \tag{4}$$

where v_y represents the particle velocity, $\nabla_y \log p_Y(y)$ is the gradient of the log of the probability density function (called the score function) of p_Y , γ is friction, M is the mass, and B_t is the standard Wiener process in $N \times 3$ -dimensions: $B_t \sim \mathcal{N}(0, t\mathbb{I}_{N \times 3})$. In practice, we employ the BAOAB solver (Appendix I) to integrate Equation 3 numerically.

3. **Jump** back to p_X to obtain samples $\hat{x}_1, \dots, \hat{x}_N$ (Figure 1c):

$$\hat{x}_i = \hat{x}(y_i) = \mathbb{E}[X \mid Y = y_i],\tag{5}$$

where $\hat{x}(\cdot) \equiv \mathbb{E}[X \mid Y = \cdot]$ is called the denoiser. It corresponds to the minimizer (Section H.1) of the ℓ_2 -loss between clean samples X and samples denoised back from $Y = X + \sigma \varepsilon$.

$$\hat{x}(\cdot) = \underset{f}{\arg\min} \mathbb{E}_{X \sim p_X, \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})} [\|f(Y) - X\|^2], \tag{6}$$

where $f: \mathbb{R}^{N \times 3} \to \mathbb{R}^{N \times 3}$. As shown by Robbins [66], Miyasawa [58] (and Section H.2), the denoiser \hat{x} is closely linked to the score $\nabla_u \log p_Y$:

$$\hat{x}(y) = y + \sigma^2 \nabla_y \log p_Y(y). \tag{7}$$

Importantly, the score function $\nabla_y \log p_Y$ shows up in both the **walk** and **jump** steps, and we need to approximate this quantity.

2.3 Learning to Denoise

In order to run Walk-Jump Sampling as outlined above, we have the choice of modelling either the score $\nabla_y \log p_Y$ or the denoiser \hat{x} as they are equivalent by Equation 7. Following trends in diffusion models [41, 42], we model the denoiser as a neural network $\hat{x}_{\theta}(y, \sigma) \approx \hat{x}(y)$ parameterized by model parameters θ .

Importantly, we only need to learn a model at a **single, fixed noise level** σ . This is unlike training diffusion or flow-matching models where a wide range of noise levels are required for sampling. In particular, the choice of noise level σ for WJS is important because mode-mixing becomes faster as σ is increased, but the task asked of the denoiser becomes harder.

The denoiser \hat{x}_{θ} thus takes in noisy point clouds y formed by adding noise (at a fixed noise level σ) to clean point clouds x. The denoiser is tasked to reconstruct back x, given y. To be precise, training the denoiser \hat{x}_{θ} consists of solving the following optimization problem:

$$\theta^* = \underset{\theta}{\arg\min} \mathbb{E}_{X \sim p_X, \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})} \|\hat{x}_{\theta}(Y, \sigma) - X\|^2$$
(8)

to obtain θ^* , the optimal model parameters. As is standard in the empirical risk minimization (ERM) [84] setting, we approximate the expectation in Equation 8 by sampling $X \sim p_X$ and $\varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})$. We minimize the loss as a function of model parameters θ using the first-order optimizer Adam [44] in PyTorch 2.0 [4, 22].

2.4 Parametrization of the Denoiser Network

We summarize the key features of the denoiser network $\hat{x}_{\theta}(y,\sigma)$ which will approximate $\hat{x}(y)$ in this section. Note that σ is fixed in our setting, but we explicitly mention it in this section for clarity. A diagrammatic overview of our model along with specific hyperparameters are presented in Appendix F.

We utilize the same parametrization of the denoiser as originally proposed by Karras et al. [41, 42] (in the context of image generation):

$$\hat{x}_{\theta}(y,\sigma) = c_{\text{skin}}(\sigma)y + c_{\text{out}}(\sigma)F_{\theta}(c_{\text{in}}(\sigma)y, c_{\text{noise}}(\sigma)), \tag{9}$$

where F_{θ} represents a learned network parameterized by parameters θ . In particular, F_{θ} (Figure 14) is a geometric graph neural network (GNN) model similar to NequIP [8, 81]. Importantly, F_{θ} is chosen to be SE(3)-equivariant, in contrast to existing methods [33, 45, 46] that utilize the E(3)-equivariant EGNN model [71]. As rightly pointed out by Dumitrescu et al. [20] and Schreiner et al. [73], E(3)-equivariant models are equivariant under parity, which means that are forced to transform

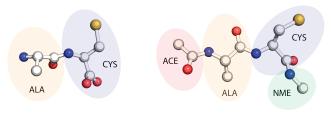


Figure 2: A side-by-side comparison of uncapped (left) compared to capped (right) ALA-CYS. The acetyl (ACE) and N-methyl (NME) capping groups provide steric hindrance and prevent local charge interactions on the N-terminal and C-terminal ends.

mirrored structures identically. When we experimented with E(3)-equivariant architectures, we found symmetric Ramachandran plots which arise from the unnecessary parity constraint of the denoising network. For this reason, TBG [45] and Timewarp [46] use a 'chirality checker' to post-hoc fix the generated structures from their model. For JAMUN, such post-processing is unnecessary because our model can distinguish between chiral structures.

The coefficients $c_{\text{skip}}(\sigma)$, $c_{\text{out}}(\sigma)$, $c_{\text{in}}(\sigma)$, $c_{\text{noise}}(\sigma)$ in Equation 9 are normalization functions (from \mathbb{R}^+ to \mathbb{R}) which adjust the effective inputs and outputs to F_{θ} . They are chosen to encourage re-use of the input y at low noise levels, but the opposite at high noise levels. Importantly, based on the insight that F_{θ} uses relative vectors in the message passing steps, we adjust the values of these coefficients instead of simply using the choices made in Karras et al. [41, 42], Wohlwend et al. [89], Abramson et al. [1], as discussed in Appendix G.

In F_{θ} , edges between atoms are computed using a radial cutoff of 10Å over the noisy positions in y. The edge features are a concatenation of a one-hot feature indicating bonded-ness and the radial distance embedded using Bessel functions. As obtained from h, atom-level features are computed using the embedding of the atomic number (eg. C and N), and the atom name following PDB notation (eg. CA, CB for alpha and beta carbons). Similarly, residue-level features are obtained using the embedding of the residue code (eg. ALA, CYS) and concatenated to each atom in the residue. Importantly, we *do not use the sequence index* of the residues (eg. $0, 1, \ldots$) as we found that it hurts generalization to longer peptide lengths.

3 Datasets

For development, demonstration, and benchmarking against existing models, we use three different datasets consisting of peptides from 2 to 5 amino acids (AA) long: TIMEWARP 2AA-LARGE and TIMEWARP 4AA-LARGE from Klein et al. [46], MDGEN 4AA EXPLICIT from Jing et al. [37], and our own MD data simulated with OpenMM [21]. A summary of these datasets is presented in Table 6. The differing simulation conditions across these datasets allows us to test the broad applicability of our approach. We describe these, as well as an internal dataset we use for more realistic setting analysis in C.

We also use the implicit MD code from Timewarp to generate trajectories for three randomly picked 5AA peptides with codes NRLCQ, VWSPF and KTYDI. We call this dataset UNCAPPED 5AA. Further details on MD simulation conditions can be found in Appendix C.

4 Related Work and Baseline Models

The goal of building machine learning models that can generate conformational ensembles of molecular systems is not new. While a full overview of this field is beyond the scope of this work – see [6] for a recent review – we note a few relevant previous efforts. Boltzmann Generators [60] introduced the idea that a neural network could be used to transform the underlying data distribution into an easier-to-sample Gaussian distribution. DiffMD [90] learns a diffusion model over conformations of small organic molecules from MD17 [14], showing some level of transferability across $C_7O_2H_{10}$ isomers [76]. Timewarp [46] uses a normalizing flow as a proposal distribution in MCMC sampling of the Boltzmann distribution to approximate the conditional distribution of future

conformational states $x^{(t+\Delta t)}$ conditional on the present state $x^{(t)}$. ITO [73] modelled this distribution using diffusion with a SE(3)-equivariant PaiNN architecture. Equijump [19] extended this idea with a protein-specific message-passing neural network with reweighting to sample rarer conformations of fast-folding proteins. Importantly, Timewarp was the first truly transferable Boltzmann generator, but was still too slow relative to molecular dynamics on unseen systems. Later, Transferable Boltzmann Generators (TBG) [45] built upon Timewarp by using flow-matching instead of maximum likelihood estimation and a more efficient continuous normalizing flow architecture.

We discuss a few more related methods and studies in E.

A related problem is that of protein structure prediction, as successfully tackled by the AlphaFold models [77, 38, 1]. Here, we compare to Boltz-1 [89], an open-source reproduction of AlphaFold3 [1]. Boltz-1 was trained exclusively on static crystalline structures of folded states, without any dynamics or conformational information, allowing us to evaluate how effectively the conformational landscape can be inferred from structural data alone.

By fixing covalent bond lengths along the backbone and side chains, AlphaFold2 introduced a 'frames' parametrization of protein structures consisting of a roto-translation together with 7 torsion angles for each residue. MDGen [37] cleverly builds on this parametrization by creating a SE(3)-invariant tokenization of the backbone torsion angles, relative to a known initial conformation (here, $x^{(0)}$). Then, they learn a stochastic interpolant [3] (a generalization of diffusion and flow-matching) over the trajectories of these tokens. While their overall objective is different from ours – MD trajectory generation as opposed to Boltzmann distributions – we can compare to their 'forward simulation' model. Similarly to AlphaFold2, their SE(3)-invariant tokenization is limited to single-chain proteins and peptides, but allows for more efficient architectures. On the other hand, JAMUN models all atoms explicitly with a SE(3)-equivariant network, which makes it far more flexible and easily extendable to molecules other than proteins, as discussed in the supplement.

Several models build upon AlphaFold2 to sample conformational ensembles, discussed further in E. One such method is BioEmu– technically only a backbone-only model, but their repository provides an additional side-chain reconstruction step using H-Packer [86], allowing comparison to the all-atom models. As seen in Figure 3, the side-chain reconstruction can be quite expensive, because of the lack of support for batched inference with H-Packer.

5 Metrics

We adopt the analysis methods of MDGen [37]. In particular, we compare models by projecting their sampled distributions of all-atom positions onto a variety of variables: pairwise distances, dihedral angles of backbone (known as Ramachandran plots) and sidechain torsion angles, TICA (time-lagged independent coordinate analysis) projections, and metastable state probabilities as computed by Markov State Models (MSMs) fit with PyEMMA [72]. TICA [59] is a popular dimensionality reduction method for larger molecules which aims to extract slow collective degrees of freedom from a trajectory [61, 75]. As is standard practice, all TICA projections and MSMs are estimated using the reference MD data.

6 Experiments

Here, we compare to several of these state-of-the art methods on our benchmark datasets: 1. TBG [45] on TIMEWARP 2AA-LARGE, 2. MDGen [37] on MDGEN 4AA-EXPLICIT, and 3. Boltz-1 [89] and BioEmu [52] on UNCAPPED 5AA.

Figure 3 contains a summary of the sampling efficiencies for different models, averaged over the corresponding test sets peptides. While the TBG model can technically produce reweighted samples, we run it in the un-reweighted mode, making it a Boltzmann emulator which is almost $10\times$ faster in practice. This allows for a fair comparison to JAMUN.

Due to the different simulation conditions across the datasets shown in Table 6, we train a different JAMUN model for each dataset. However, the same noise level of $\sigma = 0.4$ Å is applied for training and sampling on all datasets. In fact, all training hyperparameters *are kept identical* across datasets. We trained each JAMUN model for 3 days on 2 NVIDIA RTX A100 GPUs with 40 GB memory,

Model	Time per Sample	Number of Samples	Total Time
TBG	$720\mathrm{ms}$	5,000	$60\mathrm{min}$
MDGen	$6\mathrm{ms}$	10,000	$1\mathrm{min}$
Boltz-1	$360\mathrm{ms}$	10,000	$60\mathrm{min}$
BioEmu	$15\mathrm{ms}$	10,000	$2.5\mathrm{min}$
BioEmu + H-Packer	$4320\mathrm{ms}$	10,000	$720\mathrm{min}$
JAMUN (2AA)	$2\mathrm{ms}$	100,000	$3 \min$
JAMUN (4AA)	$3\mathrm{ms}$	100,000	$5\mathrm{min}$
JAMUN (5AA)	$8\mathrm{ms}$	100,000	$12.5\mathrm{min}$
CAPPED 2AA	$40\mathrm{ms}$	60,000	$40\mathrm{min}$
MDGEN 4AA-EXPLICIT	$11\mathrm{ms}$	1,000,000	$180\mathrm{min}$
UNCAPPED 5AA	$108\mathrm{ms}$	100,000	$180\mathrm{min}$

Figure 3: Comparison of (approximate and batched) sampling times per test peptide for baseline models (top), baseline MD simulations (middle) and JAMUN. All models and baselines were run on a single NVIDIA RTX A100 GPU, except for MDGEN 4AA-EXPLICIT which was simulated on a single NVIDIA T4 GPU, as mentioned in Jing et al. [37].

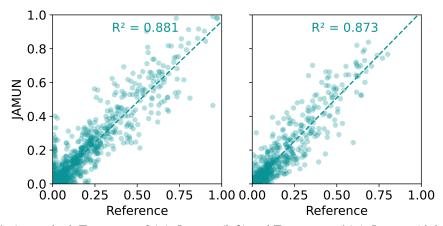


Figure 4: Across both TIMEWARP 2AA-LARGE (left) and TIMEWARP 4AA-LARGE (right), MSM state probabilities for JAMUN samples (on the y-axis) and those for the reference MD trajectories (on the x-axis) across all test peptides are strongly correlated. The perfect sampler will obtain an R^2 of 1.

although competitive results can be obtained by training for only 1 day. Our denoiser model is built with the e3nn library [25], and contains approximately 10.5M parameters.

This scale of noise ($\sigma=0.4\text{Å}$) is large enough to result in significant disruption of structure, leading to the smoothed Gaussian convolved 'walk' manifold. However, the scale is also small enough to avoid atoms 'swapping' position, for instance, or pairs of bonded atoms ending up very far from each other with reasonable probability. Indeed, as shown in Figure 8a, we find that higher noise levels (such as $\sigma=0.8\text{Å}$) result in samples with broken topologies, while lower noise levels (such as $\sigma=0.2\text{Å}$) require many more sampling steps to explore the entire conformational landscape.

6.1 Results on Timewarp 2AA-Large and Timewarp 4AA-Large

First, we show that JAMUN samples similar states as the reference MD data. Indeed, Figure 4 shows that the metastable state probabilities over JAMUN sampled trajectories match very well with those over the reference MD data on the TIMEWARP 2AA-LARGE and TIMEWARP 4AA-LARGE datasets.

Table 1 shows that on TIMEWARP 2AA-LARGE, JAMUN outperforms TBG when run for equal amounts of time (based on Figure 3), and is only slightly worse when TBG is run for $20 \times$ longer.

In Figure 5a and Figure 5b, we visualize the TICA-0,1 projections and Ramachandran plots for randomly chosen test peptides from TIMEWARP 2AA-LARGE, highlighting that TBG misses certain basins that JAMUN is able to sample.

Trajectory	Backbone Torsions	Sidechain Torsions	All Torsions	TICA-0	TICA-0,1	Metastable Probs
JAMUN TBG TBG (20× shorter)	$\begin{array}{c} 0.130 \pm 0.020 \\ 0.083 \pm 0.028 \\ 0.203 \pm 0.070 \end{array}$	$\begin{array}{c} 0.185 \pm 0.044 \\ 0.115 \pm 0.045 \\ 0.235 \pm 0.073 \end{array}$	$\begin{array}{c} 0.165 \pm 0.030 \\ 0.105 \pm 0.038 \\ 0.225 \pm 0.071 \end{array}$	$\begin{array}{c} 0.177 \pm 0.053 \\ 0.122 \pm 0.051 \\ 0.240 \pm 0.072 \end{array}$	$\begin{array}{c} 0.260 \pm 0.052 \\ 0.225 \pm 0.070 \\ 0.484 \pm 0.073 \end{array}$	$\begin{array}{c} 0.155 \pm 0.063 \\ 0.101 \pm 0.046 \\ 0.124 \pm 0.054 \end{array}$
JAMUN Reference (10× shorter) Reference (100× shorter) MDGen	$\begin{array}{c} 0.159 \pm 0.060 \\ 0.100 \pm 0.035 \\ 0.227 \pm 0.062 \\ 0.129 \pm 0.039 \end{array}$	$\begin{array}{c} 0.210 \pm 0.057 \\ 0.092 \pm 0.027 \\ 0.254 \pm 0.060 \\ 0.089 \pm 0.032 \end{array}$	$\begin{array}{c} 0.187 \pm 0.054 \\ 0.095 \pm 0.025 \\ 0.240 \pm 0.051 \\ 0.107 \pm 0.028 \end{array}$	$\begin{array}{c} 0.257 \pm 0.111 \\ 0.234 \pm 0.068 \\ 0.444 \pm 0.131 \\ 0.228 \pm 0.092 \end{array}$	$\begin{array}{c} 0.353 \pm 0.120 \\ 0.332 \pm 0.067 \\ 0.569 \pm 0.108 \\ 0.320 \pm 0.087 \end{array}$	$\begin{array}{c} 0.262 \pm 0.118 \\ 0.286 \pm 0.066 \\ 0.482 \pm 0.138 \\ 0.233 \pm 0.093 \end{array}$
JAMUN Ref. (10× shorter) Ref. (100× shorter) Boltz-1 BioEmu	$\begin{array}{c} 0.196 \pm 0.027 \\ 0.118 \pm 0.013 \\ 0.272 \pm 0.062 \\ 0.425 \pm 0.033 \\ 0.329 \pm 0.013 \end{array}$	$\begin{array}{c} 0.196 \pm 0.013 \\ 0.150 \pm 0.032 \\ 0.307 \pm 0.023 \\ 0.402 \pm 0.036 \\ 0.489 \pm 0.024 \end{array}$	$\begin{array}{c} 0.197 \pm 0.010 \\ 0.135 \pm 0.015 \\ 0.290 \pm 0.039 \\ 0.411 \pm 0.029 \\ 0.420 \pm 0.018 \end{array}$	$\begin{array}{c} 0.336 \pm 0.049 \\ 0.430 \pm 0.077 \\ 0.555 \pm 0.070 \\ 0.457 \pm 0.050 \\ 0.415 \pm 0.092 \end{array}$	$\begin{array}{c} 0.440 \pm 0.048 \\ 0.504 \pm 0.079 \\ 0.678 \pm 0.034 \\ 0.584 \pm 0.026 \\ 0.597 \pm 0.026 \end{array}$	$\begin{array}{c} 0.250 \pm 0.075 \\ 0.460 \pm 0.051 \\ 0.601 \pm 0.112 \\ 0.483 \pm 0.047 \\ 0.321 \pm 0.018 \end{array}$

Table 1: Comparison of Jenson-Shannon distances between JAMUN and baselines, both benchmark models and reference MD trajectories, averaged over test peptides in each dataset. Baselines are shortened to represent comparable sampling time to Jamun. First, compared with TBG for TIMEWARP 2AA-LARGE. Note that TBG ($20\times$ shorter) has a similar sampling time as JAMUN. Second, with MDgen and MD trajectories (shortened by a factor of 10 and 100) for the MDGEN 4AA-EXPLICIT. Third, for generalization, reference MD (shortened by a factor of 10 and 100), Boltz-1 and BioEmu, for three test peptides in UNCAPPED 5AA.

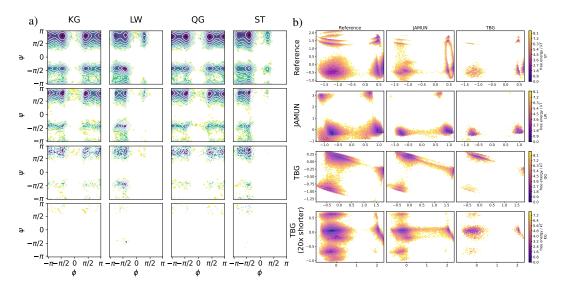


Figure 5: Results for 4 randomly chosen test peptides on TIMEWARP 2AA-LARGE: (a) Ramachandran plots, and (b) TICA-0,1 projections.

6.2 Results on MDGEN 4AA-EXPLICIT

Table 1 shows that JAMUN is very competitive with MDGen on the JSD metrics. In fact, the TICA projections in Figure 6 shows that MDGen is missing some basins that JAMUN is able to sample, while the Ramachandran plots show an example where JAMUN hallucinates a basin.

Figure 7 shows the significant speedups in backbone and sidechain torsion decorrelation times for JAMUN compared to the reference MD data. This matches our intuition about taking 'larger' integrator steps in the smoothed manifold of the noisy latent space.

We perform an analysis of the physical validity and the energy of JAMUN samples in Appendix B. In summary, we find that JAMUN samples pass all Posebusters [12] checks at a rate of $\approx 94.7\%$ and have energies close to the ground truth MD for randomly selected 20 unseen peptides.

6.3 Assessing Generalization over Peptide Lengths with UNCAPPED 5AA

JAMUN's graph neural network architecture enables it to operate on molecules of larger sizes than it was originally trained on. Thus, we test whether JAMUN can generalize to peptides of lengths

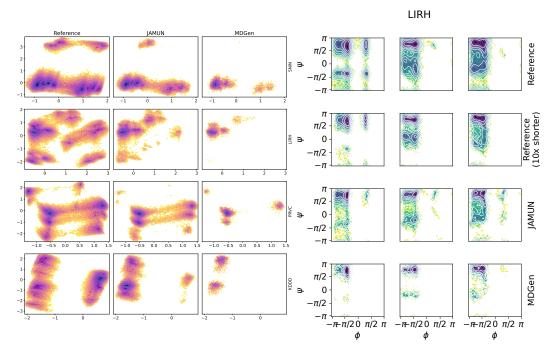


Figure 6: MDGEN 4AA-EXPLICIT results: (left) TICA-0,1 projections for 4 randomly chosen test peptides; (right) Ramachandran plots for JAMUN and MDGen on a randomly chosen test.

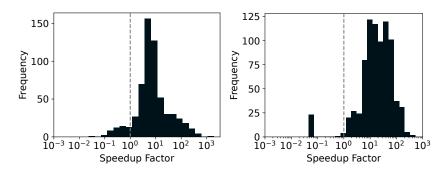


Figure 7: MDGEN 4AA-EXPLICIT results: Speedup defined as ratio of decorrelation time histogrammed for all test backbone (left) and sidechain (right) torsions.

beyond its training set. This is a challenging task, and one that we believe conformational generation models have not been adequately benchmarked on. Unfortunately, neither TBG nor MDGen transfer to UNCAPPED 5AA due to fixed-length positional embeddings, despite our best efforts. In particular, we initialized the positional embeddings to support longer peptides, but this resulted in broken topologies in the resulting samples. Instead, we choose Boltz-1 and BioEmu which support sampling on UNCAPPED 5AA to compare against JAMUN trained on TIMEWARP 4AA-LARGE.

Surprisingly, we find that the JAMUN model *trained only on 4AA peptides can accurately predict ensembles for 5AA peptides*. Figure 8b and Figure 11 show that JAMUN is able to recover most states and even reproduce relative probabilities. Interestingly, the same experiment does not work if we train on TIMEWARP 2AA-LARGE instead, suggesting that the 2AA reference MD data may not be informative enough to generalize from.

On the other hand, we find that Boltz-1 is unable to sample the diversity of peptide conformations. This is not entirely surprising as Boltz-1 was not trained on any MD data, as we noted before. Further, Boltz-1 also utilizes a common pair representation across all diffusion samples, as computed by its Pairformer stack. The pair representation intuitively represents the residue-wise distance matrix,

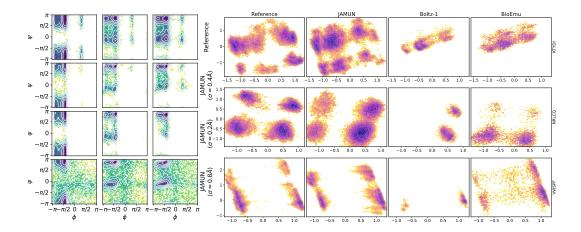


Figure 8: (left) Comparing noise sensitivity for an example test peptide GCSL for JAMUN, sampled identically, showing the tradeoff between slower mode mixing at $\sigma = 0.2\text{Å}$) and broken topologies at $\sigma = 0.8\text{Å}$. (right) TICA-0,1 projections for three test peptides in UNCAPPED 5AA.

and thus encodes a significant portion of the geometry. Keeping this representation fixed possibly prevents the sampling of large conformational changes.

Surprisingly, BioEmu also seems to struggle in this setting, even when considering distributions of backbone torsion angles only. This suggests that BioEmu cannot capture the relative flexibility of smaller peptides, even when trained on MD data for much larger proteins.

Quantitatively, Table 1 shows that JAMUN significantly outperforms Boltz-1 and BioEmu on the metrics from Section 5. As seen in Figure 3, JAMUN is roughly $5\times$ faster than Boltz-1, and is roughly $60\times$ faster than BioEmu when we perform side-chain reconstruction with H-Packer.

7 Conclusion

We present JAMUN, a walk-jump sampling model for generating ensembles of molecular conformations, outperforming the state-of-the-art TBG model, and competitive with the performance of MDGen with no protein-specific parametrization. This represents an important step toward the ultimate goal of a transferable generative model for protein conformational ensembles. Performing MD in the noised space gives the model a clear physics interpretation, and allows faster decorrelation (and hence, sampling) than classical MD.

The model has some limitations that motivate future work. While it is highly transferable in the space of two to five amino acid peptides, scaling up is likely to require more exploration and intensive data generation in future work. Additionally, while the current SE(3)-equivariant denoiser architecture works well, further development of the denoising network could speed up sampling. Alternative jump methods, such as multiple denoising steps (à la diffusion), could also serve to sharpen generation. Lastly, a promising direction that has not yet been explored is the application of classical enhanced sampling methods, such as metadynamics, for traversing the noisy space.

References

- [1] J. Abramson, J. Adler, J. Dunger, R. Evans, T. Green, A. Pritzel, O. Ronneberger, L. Willmore, A. J. Ballard, J. Bambrick, S. W. Bodenstein, D. A. Evans, C.-C. Hung, M. O'Neill, D. Reiman, K. Tunyasuvunakool, Z. Wu, A. Žemgulytė, E. Arvaniti, C. Beattie, O. Bertolli, A. Bridgland, A. Cherepanov, M. Congreve, A. I. Cowen-Rivers, A. Cowie, M. Figurnov, F. B. Fuchs, H. Gladman, R. Jain, Y. A. Khan, C. M. R. Low, K. Perlin, A. Potapenko, P. Savy, S. Singh, A. Stecula, A. Thillaisundaram, C. Tong, S. Yakneen, E. D. Zhong, M. Zielinski, A. Žídek, V. Bapst, P. Kohli, M. Jaderberg, D. Hassabis, and J. M. Jumper. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, Jun 2024. ISSN 1476-4687. doi: 10.1038/s41586-024-07487-w. URL https://doi.org/10.1038/s41586-024-07487-w.
- [2] V. Agarwal and A. C. McShan. The power and pitfalls of alphafold2 for structure prediction beyond rigid globular proteins. *Nature Chemical Biology*, 20(8):950–959, Aug 2024. ISSN 1552-4469. doi: 10.1038/s41589-024-01638-w. URL https://doi.org/10.1038/s41589-024-01638-w.
- [3] M. S. Albergo, N. M. Boffi, and E. Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions, 2023. URL https://arxiv.org/abs/2303.08797.
- [4] J. Ansel, E. Yang, H. He, N. Gimelshein, A. Jain, M. Voznesensky, B. Bao, P. Bell, D. Berard, E. Burovski, G. Chauhan, A. Chourdia, W. Constable, A. Desmaison, Z. DeVito, E. Ellison, W. Feng, J. Gong, M. Gschwind, B. Hirsh, S. Huang, K. Kalambarkar, L. Kirsch, M. Lazos, M. Lezcano, Y. Liang, J. Liang, Y. Lu, C. Luk, B. Maher, Y. Pan, C. Puhrsch, M. Reso, M. Saroufim, M. Y. Siraichi, H. Suk, M. Suo, P. Tillet, E. Wang, X. Wang, W. Wen, S. Zhang, X. Zhao, K. Zhou, R. Zou, A. Mathews, G. Chanan, P. Wu, and S. Chintala. PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation. In 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24). ACM, Apr. 2024. doi: 10.1145/3620665. 3640366. URL https://pytorch.org/assets/pytorch2-2.pdf.
- [5] A. Antoszewski, C.-J. Feng, B. P. Vani, E. H. Thiede, L. Hong, J. Weare, A. Tokmakoff, and A. R. Dinner. Insulin dissociates by diverse mechanisms of coupled unfolding and unbinding. *The Journal of Physical Chemistry B*, 124(27):5571–5587, 2020. doi: 10.1021/acs.jpcb.0c03521. URL https://doi.org/10.1021/acs.jpcb.0c03521. PMID: 32515958.
- [6] A. Aranganathan, X. Gu, D. Wang, B. Vani, and P. Tiwary. Modeling Boltzmann weighted structural ensembles of proteins using AI based methods. 2024.
- [7] M. Arts, V. Garcia Satorras, C.-W. Huang, D. Zugner, M. Federici, C. Clementi, F. Noé, R. Pinsler, and R. van den Berg. Two for one: Diffusion models and force fields for coarse-grained molecular dynamics. *Journal of Chemical Theory and Computation*, 19(18):6151–6159, 2023.
- [8] S. Batzner, A. Musaelian, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, and B. Kozinsky. E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nature Communications*, 13(1):2453, May 2022. ISSN 2041-1723. doi: 10.1038/s41467-022-29939-5. URL https://doi.org/10.1038/s41467-022-29939-5.
- [9] H. M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weissig, I. N. Shindyalov, and P. E. Bourne. The Protein Data Bank. *Nucleic Acids Research*, 28(1):235–242, 01 2000. ISSN 0305-1048. doi: 10.1093/nar/28.1.235. URL https://doi.org/10.1093/nar/28.1.235.
- [10] D. W. Borhani and D. E. Shaw. The future of molecular dynamics simulations in drug discovery. *Journal of computer-aided molecular design*, 26:15–26, 2012.
- [11] G. R. Bowman. AlphaFold and Protein Folding: Not Dead Yet! The Frontier Is Conformational Ensembles. *Annual Review of Biomedical Data Science*, 7, 2024.
- [12] M. Buttenschoen, G. M. Morris, and C. M. Deane. Posebusters: Ai-based docking methods fail to generate physically valid poses or generalise to novel sequences. *Chemical Science*, 15(9): 3130–3139, 2024. ISSN 2041-6539. doi: 10.1039/d3sc04185a. URL http://dx.doi.org/10.1039/D3SC04185A.
- [13] D. Chakravarty and L. L. Porter. AlphaFold2 fails to predict protein fold switching. *Protein Sci*, 31(6):e4353, June 2022.

- [14] S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller. Machine learning of accurate energy-conserving molecular force fields. *Science Advances*, 3 (5):e1603015, 2017. doi: 10.1126/sciadv.1603015. URL https://www.science.org/doi/abs/10.1126/sciadv.1603015.
- [15] G. Colombo. Computing allostery: from the understanding of biomolecular regulation and the discovery of cryptic sites to molecular design. *Current Opinion in Structural Biology*, 83: 102702, 2023.
- [16] T. Darden, D. York, and L. Pedersen. Particle Mesh Ewald: An N log(N) method for Ewald sums in large systems. *The Journal of Chemical Physics*, 98(12):10089–10092, June 1993. doi: 10.1063/1.464397. URL https://doi.org/10.1063/1.464397.
- [17] D. del Alamo, D. Sala, H. S. Mchaourab, and J. Meiler. Sampling alternative conformational states of transporters and receptors with alphafold2. *eLife*, 11:e75751, mar 2022. ISSN 2050-084X. doi: 10.7554/eLife.75751. URL https://doi.org/10.7554/eLife.75751.
- [18] A. R. Dinner, J. C. Mattingly, J. O. B. Tempkin, B. Van Koten, and J. Weare. Trajectory stratification of stochastic dynamics. SIAM Rev Soc Ind Appl Math, 60(4):909–938, Nov. 2018.
- [19] A. dos Santos Costa, I. Mitnikov, F. Pellegrini, A. Daigavane, M. Geiger, Z. Cao, K. Kreis, T. Smidt, E. Kucukbenli, and J. Jacobson. Equijump: Protein dynamics simulation via so(3)-equivariant stochastic interpolants, 2024. URL https://arxiv.org/abs/2410.09667.
- [20] A. Dumitrescu, D. Korpela, M. Heinonen, Y. Verma, V. Iakovlev, V. Garg, and H. Lähdesmäki. Field-based Molecule Generation, 2024. URL https://arxiv.org/abs/2402.15864.
- [21] P. Eastman, J. Swails, J. D. Chodera, R. T. McGibbon, Y. Zhao, K. A. Beauchamp, L.-P. Wang, A. C. Simmonett, M. P. Harrigan, C. D. Stern, R. P. Wiewiora, B. R. Brooks, and V. S. Pande. OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLOS Computational Biology*, 13(7):e1005659, July 2017. doi: 10.1371/journal.pcbi.1005659. URL https://doi.org/10.1371/journal.pcbi.1005659.
- [22] W. Falcon and The PyTorch Lightning team. PyTorch Lightning, Mar. 2019. URL https://github.com/Lightning-AI/lightning.
- [23] M. L. Fernández-Quintero, N. D. Pomarici, A.-L. M. Fischer, V. J. Hoerschinger, K. B. Kroell, J. R. Riccabona, A. S. Kamenik, J. R. Loeffler, J. A. Ferguson, H. R. Perrett, et al. Structure and dynamics guiding design of antibody therapeutics and vaccines. *Antibodies*, 12(4):67, 2023.
- [24] N. C. Frey, D. Berenberg, J. Kleinhenz, I. Hotzel, J. Lafrance-Vanasse, R. L. Kelly, Y. Wu, A. Rajpal, S. Ra, R. Bonneau, K. Cho, A. Loukas, V. Gligorijevic, and S. Saremi. Protein discovery with discrete walk-jump sampling. In *International Conference on Learning Representations*, 2024.
- [25] M. Geiger and T. Smidt. e3nn: Euclidean neural networks. *arXiv preprint arXiv:2207.09453*, 2022.
- [26] N. R. Gough and C. G. Kalodimos. Exploring the conformational landscape of protein kinases. *Current Opinion in Structural Biology*, 88:102890, 2024.
- [27] C. A. Grambow, H. Weir, C. N. Cunningham, T. Biancalani, and K. V. Chuang. Ringer: Rapid inference of macrocyclic peptide conformational ensembles. *arXiv preprint*, arXiv:2310.01234, 2023. URL https://arxiv.org/abs/2310.01234.
- [28] C. A. Grambow, H. Weir, C. N. Cunningham, T. Biancalani, and K. V. Chuang. Cremp: Conformer–rotamer ensembles of macrocyclic peptides for machine learning. *Scientific Data*, 11, 2024. doi: 10.1038/s41597-024-03698-y. URL https://doi.org/10.1038/s41597-024-03698-y.
- [29] X. Guan, Q.-Y. Tang, W. Ren, M. Chen, W. Wang, P. G. Wolynes, and W. Li. Predicting protein conformational motions using energetic frustration analysis and alphafold2. *Proceedings of the National Academy of Sciences*, 121(35):e2410662121, 2024. doi: 10.1073/pnas.2410662121. URL https://www.pnas.org/doi/abs/10.1073/pnas.2410662121.
- [30] K. Henzler-Wildman and D. Kern. Dynamic personalities of proteins. *Nature*, 450(7172): 964–972, Dec. 2007.
- [31] B. Hess, H. Bekker, H. J. Berendsen, and J. G. Fraaije. Lincs: a linear constraint solver for molecular simulations. *Journal of computational chemistry*, 18(12):1463–1472, 1997.

- [32] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models, 2020. URL https://arxiv.org/abs/2006.11239.
- [33] E. Hoogeboom, V. G. Satorras, C. Vignac, and M. Welling. Equivariant Diffusion for Molecule Generation in 3D, 2022. URL https://arxiv.org/abs/2203.17003.
- [34] T. Hsu, B. Sadigh, V. Bulatov, and F. Zhou. Score dynamics: Scaling molecular dynamics with picoseconds time steps via conditional diffusion model. *Journal of Chemical Theory and Computation*, 20(6):2335–2348, 2024.
- [35] X. Huang, R. Pearce, and Y. Zhang. FASPR: an open-source tool for fast and accurate protein side-chain packing. *Bioinformatics*, 36(12):3758–3765, June 2020.
- [36] B. Jing, B. Berger, and T. Jaakkola. AlphaFold meets flow matching for generating protein ensembles. *arXiv preprint arXiv:2402.04845*, 2024.
- [37] B. Jing, H. Stärk, T. Jaakkola, and B. Berger. Generative modeling of molecular dynamics trajectories. *arXiv preprint arXiv:2409.17808*, 2024.
- [38] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, A. Bridgland, C. Meyer, S. A. A. Kohl, A. J. Ballard, A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D. Reiman, E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein, D. Silver, O. Vinyals, A. W. Senior, K. Kavukcuoglu, P. Kohli, and D. Hassabis. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, Aug 2021. ISSN 1476-4687. doi: 10.1038/s41586-021-03819-2. URL https://doi.org/10.1038/s41586-021-03819-2.
- [39] A. Jussupow and V. R. I. Kaila. Effective molecular dynamics from neural network-based structure prediction models. *Journal of Chemical Theory and Computation*, 19(7):1965– 1975, 2023. doi: 10.1021/acs.jctc.2c01027. URL https://doi.org/10.1021/acs.jctc. 2c01027. PMID: 36961997.
- [40] W. Kabsch. A solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A*, 32(5):922–923, Sep 1976. doi: 10.1107/S0567739476001873. URL https://doi.org/10.1107/S0567739476001873.
- [41] T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the Design Space of Diffusion-Based Generative Models. In *Proc. NeurIPS*, 2022.
- [42] T. Karras, M. Aittala, J. Lehtinen, J. Hellsten, T. Aila, and S. Laine. Analyzing and Improving the Training Dynamics of Diffusion Models. In *Proc. CVPR*, 2024.
- [43] S. Kieninger and B. G. Keller. GROMACS Stochastic Dynamics and BAOAB Are Equivalent Configurational Sampling Algorithms. *Journal of Chemical Theory and Computation*, 18(10): 5792–5798, 2022.
- [44] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2017. URL https://arxiv.org/abs/1412.6980.
- [45] L. Klein and F. Noé. Transferable Boltzmann Generators. arXiv preprint arXiv:2406.14426, 2024.
- [46] L. Klein, A. Foong, T. Fjelde, B. Mlodozeniec, M. Brockschmidt, S. Nowozin, F. Noé, and R. Tomioka. Timewarp: Transferable acceleration of molecular dynamics by learning timecoarsened dynamics. *Advances in Neural Information Processing Systems*, 36, 2024.
- [47] L. Klein, A. Krämer, and F. Noé. Equivariant flow matching. *Advances in Neural Information Processing Systems*, 36, 2024.
- [48] M. Lazou, O. Khan, T. Nguyen, D. Padhorny, D. Kozakov, D. Joseph-McCarthy, and S. Vajda. Predicting multiple conformations of ligand binding sites in proteins suggests that alphafold2 may remember too much. *Proceedings of the National Academy of Sciences*, 121(48): e2412719121, 2024. doi: 10.1073/pnas.2412719121. URL https://www.pnas.org/doi/abs/10.1073/pnas.2412719121.
- [49] B. Leimkuhler and C. Matthews. Rational Construction of Stochastic Numerical Methods for Molecular Sampling. *Applied Mathematics Research eXpress*, 2013(1):34–56, June 2012. ISSN 1687-1200. doi: 10.1093/amrx/abs010. URL https://doi.org/10.1093/amrx/abs010. eprint: https://academic.oup.com/amrx/article-pdf/2013/1/34/397230/abs010.pdf.

- [50] B. Leimkuhler and C. Matthews. Robust and efficient configurational molecular sampling via langevin dynamics. *The Journal of Chemical Physics*, 138(17):174102, May 2013. doi: 10.1063/1.4802990. URL https://doi.org/10.1063/1.4802990.
- [51] B. Leimkuhler and C. Matthews. Molecular Dynamics: With Deterministic and Stochastic Numerical Methods. Number 39 in Interdisciplinary Applied Mathematics. Springer International Publishing: Imprint: Springer, Cham, 1st ed. 2015 edition, 2015. ISBN 978-3-319-16375-8.
- [52] S. Lewis, T. Hempel, J. Jiménez-Luna, M. Gastegger, Y. Xie, A. Y. K. Foong, V. G. Satorras, O. Abdin, B. S. Veeling, I. Zaporozhets, Y. Chen, S. Yang, A. Schneuing, J. Nigam, F. Barbero, V. Stimper, A. Campbell, J. Yim, M. Lienen, Y. Shi, S. Zheng, H. Schulz, U. Munir, C. Clementi, and F. Noé. Scalable emulation of protein equilibrium ensembles with generative deep learning. bioRxiv, 2024. doi: 10.1101/2024.12.05.626885.
- [53] Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, N. Smetanin, R. Verkuil, O. Kabeli, Y. Shmueli, A. dos Santos Costa, M. Fazel-Zarandi, T. Sercu, S. Candido, and A. Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/science.ade2574. URL https://www.science.org/doi/abs/10.1126/science.ade2574.
- [54] Y. Lipman, R. T. Q. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=PqvMRDCJT9t.
- [55] J. Lu, B. Zhong, Z. Zhang, and J. Tang. Str2str: A score-based framework for zero-shot protein conformation sampling. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=C4BikKsgmK.
- [56] S. J. Marrink, H. J. Risselada, S. Yefimov, D. P. Tieleman, and A. H. de Vries. The martini force field: Coarse grained model for biomolecular simulations. *The Journal of Physical Chemistry B*, 111(27):7812–7824, Jul 2007. ISSN 1520-6106. doi: 10.1021/jp071097f. URL https://doi.org/10.1021/jp071097f.
- [57] M. D. Miller and G. N. Phillips. Moving beyond static snapshots: Protein dynamics and the Protein Data Bank. *Journal of Biological Chemistry*, 296, 2021.
- [58] K. Miyasawa. An Empirial Bayes Estimator of the Mean of a Normal Population. *Bulletin de l'Institut international de statistique.*, 38(4):181–188, 1960.
- [59] L. Molgedey and H. G. Schuster. Separation of a mixture of independent signals using time delayed correlations. *Phys. Rev. Lett.*, 72:3634–3637, Jun 1994. doi: 10.1103/PhysRevLett.72. 3634. URL https://link.aps.org/doi/10.1103/PhysRevLett.72.3634.
- [60] F. Noé, S. Olsson, J. Köhler, and H. Wu. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science*, 365(6457):eaaw1147, 2019.
- [61] G. Pérez-Hernández, F. Paul, T. Giorgino, G. De Fabritiis, and F. Noé. Identification of slow molecular order parameters for markov model construction. *J Chem Phys*, 139(1):015102, July 2013.
- [62] P. O. Pinheiro, A. Jamasb, O. Mahmood, V. Sresht, and S. Saremi. Structure-based Drug Design by Denoising Voxel Grids. *arXiv preprint arXiv:2405.03961*, 2024.
- [63] P. O. Pinheiro, J. Rackers, J. Kleinhenz, M. Maser, O. Mahmood, A. Watkins, S. Ra, V. Sresht, and S. Saremi. 3D molecule generation by denoising voxel grids. *Advances in Neural Information Processing Systems*, 36, 2024.
- [64] J. W. Ponder and D. A. Case. Force fields for protein simulations. 66:27-85, 2003. ISSN 0065-3233. doi: https://doi.org/10.1016/S0065-3233(03)66002-X. URL https://www.sciencedirect.com/science/article/pii/S006532330366002X.
- [65] P. Pracht and S. Grimme. Automated exploration of the low-energy chemical space with fast quantum chemical methods. *Physical Chemistry Chemical Physics*, 22(14):7169–7192, 2020. doi: 10.1039/C9CP06869D. URL https://pubs.rsc.org/en/content/articlelanding/2020/cp/c9cp06869d.
- [66] H. Robbins. An Empirical Bayes Approach to Statistics. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, volume 3.1, 1956.

- [67] J. P. Roney and S. Ovchinnikov. State-of-the-art estimation of protein model accuracy using alphafold. *Phys. Rev. Lett.*, 129:238101, Nov 2022. doi: 10.1103/PhysRevLett.129.238101. URL https://link.aps.org/doi/10.1103/PhysRevLett.129.238101.
- [68] M. Sachs, B. Leimkuhler, and V. Danos. Langevin dynamics with variable coefficients and nonconservative forces: from stationary states to numerical methods. *Entropy*, 19(12):647, 2017.
- [69] T. Saldaño, N. Escobedo, J. Marchetti, D. J. Zea, J. Mac Donagh, A. J. Velez Rueda, E. Gonik, A. García Melani, J. Novomisky Nechcoff, M. N. Salas, T. Peters, N. Demitroff, S. Fernandez Alberti, N. Palopoli, M. S. Fornasari, and G. Parisi. Impact of protein conformational diversity on alphafold predictions. *Bioinformatics*, 38(10):2742–2748, 04 2022. ISSN 1367-4803. doi: 10.1093/bioinformatics/btac202. URL https://doi.org/10.1093/bioinformatics/btac202.
- [70] S. Saremi and A. Hyvärinen. Neural Empirical Bayes. *Journal of Machine Learning Research*, 20(181):1–23, 2019.
- [71] V. G. Satorras, E. Hoogeboom, and M. Welling. E(n) Equivariant Graph Neural Networks, 2022. URL https://arxiv.org/abs/2102.09844.
- [72] M. K. Scherer, B. Trendelkamp-Schroer, F. Paul, G. Pérez-Hernández, M. Hoffmann, N. Plattner, C. Wehmeyer, J.-H. Prinz, and F. Noé. Pyemma 2: A software package for estimation, validation, and analysis of markov models. *Journal of Chemical Theory and Computation*, 11(11):5525–5542, 2015. doi: 10.1021/acs.jctc.5b00743. URL https://doi.org/10.1021/acs.jctc.5b00743. PMID: 26574340.
- [73] M. Schreiner, O. Winther, and S. Olsson. Implicit transfer operator learning: Multiple timeresolution surrogates for molecular dynamics, 2023. URL https://arxiv.org/abs/2305. 18046.
- [74] M. Schreiner, O. Winther, and S. Olsson. Implicit transfer operator learning: multiple timeresolution surrogates for molecular dynamics. arXiv preprint arXiv:2305.18046, 2023.
- [75] C. R. Schwantes and V. S. Pande. Improvements in markov state model construction reveal many Non-Native interactions in the folding of NTL9. *J Chem Theory Comput*, 9(4):2000–2009, Apr. 2013.
- [76] K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, and A. Tkatchenko. Quantum-chemical insights from deep tensor neural networks. *Nature Communications*, 8(1), Jan. 2017. ISSN 2041-1723. doi: 10.1038/ncomms13890. URL http://dx.doi.org/10.1038/ncomms13890.
- [77] A. W. Senior, R. Evans, J. Jumper, J. Kirkpatrick, L. Sifre, T. Green, C. Qin, A. Žídek, A. W. R. Nelson, A. Bridgland, H. Penedones, S. Petersen, K. Simonyan, S. Crossan, P. Kohli, D. T. Jones, D. Silver, K. Kavukcuoglu, and D. Hassabis. Improved protein structure prediction using potentials from deep learning. *Nature*, 577(7792):706–710, Jan 2020. ISSN 1476-4687. doi: 10.1038/s41586-019-1923-7. URL https://doi.org/10.1038/s41586-019-1923-7.
- [78] D. E. Shaw, P. Maragakis, K. Lindorff-Larsen, S. Piana, R. O. Dror, M. P. Eastwood, J. A. Bank, J. M. Jumper, J. K. Salmon, Y. Shan, and W. Wriggers. Atomic-level characterization of the structural dynamics of proteins. *Science*, 330(6002):341–346, 2010. doi: 10.1126/science. 1187409. URL https://www.science.org/doi/abs/10.1126/science.1187409.
- [79] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models, 2022. URL https://arxiv.org/abs/2010.02502.
- [80] M. Steinegger and J. Söding. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature Biotechnology*, 35(11):1026–1028, Nov 2017. ISSN 1546-1696. doi: 10.1038/nbt.3988. URL https://doi.org/10.1038/nbt.3988.
- [81] N. Thomas, T. Smidt, S. Kearnes, L. Yang, L. Li, K. Kohlhoff, and P. Riley. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. arXiv preprint arXiv:1802.08219, 2018.
- [82] S. Umeyama. Least-squares estimation of transformation parameters between two point patterns. IEEE Transactions on Pattern Analysis and Machine Intelligence, 13(4):376–380, 1991. doi: 10.1109/34.88573.
- [83] B. P. Vani, J. Weare, and A. R. Dinner. Computing transition path theory quantities with trajectory stratification. *J Chem Phys*, 157(3):034106, July 2022.

- [84] V. Vapnik. Principles of risk minimization for learning theory. In J. Moody, S. Hanson, and R. Lippmann, editors, *Advances in Neural Information Processing Systems*, volume 4. Morgan-Kaufmann, 1991. URL https://proceedings.neurips.cc/paper_files/paper/1991/file/ff4d5fbbafdf976cfdc032e3bde78de5-Paper.pdf.
- [85] M. Varadi, D. Bertoni, P. Magana, U. Paramval, I. Pidruchna, M. Radhakrishnan, M. Tsenkov, S. Nair, M. Mirdita, J. Yeo, O. Kovalevskiy, K. Tunyasuvunakool, A. Laydon, A. Žídek, H. Tomlinson, D. Hariharan, J. Abrahamson, T. Green, J. Jumper, E. Birney, M. Steinegger, D. Hassabis, and S. Velankar. Alphafold protein structure database in 2024: providing structure coverage for over 214 million protein sequences. *Nucleic Acids Research*, 52(D1):D368–D375, 11 2023. ISSN 0305-1048. doi: 10.1093/nar/gkad1011. URL https://doi.org/10.1093/nar/gkad1011.
- [86] G. M. Visani, W. Galvin, M. Pun, and A. Nourmohammad. H-packer: Holographic rotationally equivariant convolutional neural network for protein side-chain packing. In D. A. Knowles and S. Mostafavi, editors, *Proceedings of the 18th Machine Learning in Computational Biology meeting*, volume 240 of *Proceedings of Machine Learning Research*, pages 230–249. PMLR, 30 Nov-01 Dec 2024. URL https://proceedings.mlr.press/v240/visani24a.html.
- [87] A. Vitalis and R. V. Pappu. Methods for monte carlo simulations of biomacromolecules. *Annual reports in computational chemistry*, 5:49–76, 2009.
- [88] H. K. Wayment-Steele, A. Ojoawo, R. Otten, J. M. Apitz, W. Pitsawong, M. Hömberger, S. Ovchinnikov, L. Colwell, and D. Kern. Predicting multiple conformations via sequence clustering and alphafold2. *Nature*, 625(7996):832–839, Jan 2024. ISSN 1476-4687. doi: 10.1038/s41586-023-06832-9. URL https://doi.org/10.1038/s41586-023-06832-9.
- [89] J. Wohlwend, G. Corso, S. Passaro, M. Reveiz, K. Leidal, W. Swiderski, T. Portnoi, I. Chinn, J. Silterra, T. Jaakkola, and R. Barzilay. Boltz-1 democratizing biomolecular interaction modeling. bioRxiv, 2024. doi: 10.1101/2024.11.19.624167. URL https://www.biorxiv.org/content/early/2024/11/20/2024.11.19.624167.
- [90] F. Wu and S. Z. Li. Diffmd: A geometric diffusion model for molecular dynamics simulations, 2023. URL https://arxiv.org/abs/2204.08672.
- [91] L.-E. Zheng, S. Barethiya, E. Nordquist, and J. Chen. Machine learning generation of dynamic protein conformational ensembles. *Molecules*, 28(10):4047, 2023.
- [92] S. Zheng, J. He, C. Liu, Y. Shi, Z. Lu, W. Feng, F. Ju, J. Wang, J. Zhu, Y. Min, et al. Predicting equilibrium distributions for molecular systems with deep learning. *Nature Machine Intelligence*, pages 1–10, 2024.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes] In the abstract, our claims are that JAMUN is able to generate ensembles faster than traditional methods and recent ML methods, and can generalize outside its training data. Figure 3 shows the comparison of sampling times for JAMUN and all baselines. Table 1 and Table 1 show that JAMUN is able to outperform TBG, Boltz-1 and BioEmu in terms of efficiency in sampling the conformational space in total wall-clock time. Table 1 shows that JAMUN is competitive with the recently proposed peptide-specific MDGen model in terms of sampling quality. Section 6.3 shows that JAMUN is able to generalize to 5AA peptides, even when trained only on 4AA peptides, which is out of the reach of models such as TBG and MDGen. For these 5AA peptides, JAMUN is able to outperform Boltz-1 and BioEmu which were trained on crystal structures and MD data of much larger protein structures respectively.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes] The key limitation of the paper is described in Section 7, that JAMUN is transferable across 2-5 amino acid long peptides, but the model has not been tested on longer peptides and full proteins. The SE(3)-equivariant denoiser works well, but further development of the denoising network could speed up sampling to match up to MDGen.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes] The main contribution of our work is not theoretical. We largely build upon the theoretical framework of Walk-Jump Sampling as developed by [70]. For completeness, we re-derive all key theoretical results in Appendix H and re-derive the Langevin dynamics update in Appendix J.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes] We provide details about the training and sampling of JAMUN in Section 2.2, with details of the denoiser in Section 2.4 and Appendix F. We also provide code to reproduce our results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes] The code was submitted in accordance with the (https://nips.cc/public/guides/CodeSubmissionPolicy) which includes scripts to reproduce the results in the paper. Section 3 includes the description of the data used within the paper. The TIMEWARP 2AA-LARGE, TIMEWARP 4AA-LARGE, and MDGEN 4AA-EXPLICIT datasets are available to the public under an MIT licence. We are in the legal process of sharing the CAPPED 2AA dataset for public use.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes] The data splits are described in Section 3, and were chosen to allow a comparison to existing baseline models. The parametrization and the hyperparameters of our model are described in Section 2.4 and Appendix F. We also provide code to reproduce our results.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes] Our metrics are averaged over the test dataset, where we report the mean and 1-sigma errors for all JSD metrics.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes] We provide the compute details for training JAMUN in Section 6 and sampling JAMUN and all baseline models in Figure 3.

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes] We have made our best effort to conform to NeurIPS Code of Ethics.

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA] This paper presents work whose goal is to advance the field of accelerated sampling of protein conformations. Due to the niche field of study, we do not believe that there are any negative societal impacts that arise directly from our work.

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA] The authors do not foresee any ways in which this model could be misused.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes] Yes, we credit all authors for creating the benchmark datasets: Klein et al. [46] for TIMEWARP 2AA-LARGE and TIMEWARP 4AA-LARGE, and Jing et al. [37] for MDGEN 4AA-EXPLICIT available under the MIT license. We have respected their terms of use.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes] The new CAPPED 2AA dataset is described in Section 3 and Appendix C. We are in the legal process of releasing this dataset for public use.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA] This paper did not involve crowdsourcing or research with human subjects...

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA] This paper does not involve crowdsorcing or research with humans. The data was all collected in silico.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA] The authors did not use LLMs in the core research methods.

A Comparison of Walk-Jump Sampling Against Full Diffusion

Section 6 contains a comparison of JAMUN with existing baselines, but does not directly address the comparison of the techinique of walk-jump sampling with standard diffusion. We perform this comparison here; we fix the JAMUN architecture and directly compare walk-jump sampling (at a single noise level) to full diffusion (which samples a range of noise levels from very large to very small).

We train a diffusion model with the same JAMUN architecture on TIMEWARP 2AA. We then sample from this model using both diffusion (specifically, the ODE sampler with a noise schedule of 64 steps from 0.01Å to 10Å using the Heun second order method as recommended by EDM [41]) and walk-jump sampling. We compare the Jensen-Shannon divergence (JSD) of the backbone torsions averaged over the TIMEWARP 2AA test set as a function of number of samples and number of function evaluations (NFE):

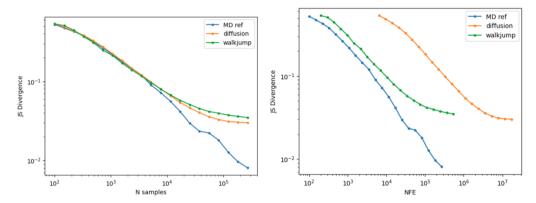


Figure 9: Comparison between diffusion, walk-jump and the ground truth MD in terms of JS divergence of backbone torsions to the full ground truth MD distributions, as a function of (left) number of samples and (right) number of function evaluations.

Sampler	Number of Samples	NFE	JSD-Backbone Torsions ↓
Walk-Jump	3149	6298	0.1501
Diffusion	3149	399923	0.1363
Walk-Jump	200000	400000	0.0496
Diffusion	200000	25400000	0.0460

Table 2: Comparison of walk-jump and diffusion sampling with the same model. Diffusion obtains a better JSD metric than walk-jump sampling, but at the cost of ≈ 30 more sampling time.

Table 2 contains a summary of the key comparison. Essentially, we find that walk-jump sampling is faster with only minor loss in fidelity; because it works in a partially noised space, instead of having to generate every sample from an uninformative Gaussian prior over many steps.

B Physical Validity and Energy Analysis

Here, we perform a physical validity analysis on the generated JAMUN samples, using the popular Posebusters [12] package. We randomly selected 20 unseen test peptides from the MDGEN 4AA-EXPLICIT dataset for this analysis.

We see that the bond lengths are correctly captured with high probability by JAMUN. Furthermore, the overall quality of the JAMUN samples is high. For a finer-grained view into the performance across the 20 test peptides, we report the empirical CDFs of the pass rates:

Next, we compute the force field energies for JAMUN samples. We add hydrogen atoms using OpenMM's PDBFixer [21], and compute energies using the amber 14 force field. We find that the energies of the JAMUN samples overlap well with those of the reference MD samples:

Posebusters Metric	Average Pass Rate
Valid Bond Lengths	97.0%
Valid Bond Angles	99.4%
Internal Steric Clash	100.0%
Internal Energy	97.5%
Overall	94.7%

Table 3: Average pass rates for Posebusters metrics on JAMUN samples.

Posebusters Pass Rate	> 90%	> 92%	> 94%	> 96%	> 98%	100%
Valid Bond Lengths	$\frac{20}{20}$	$\frac{18}{20}$	$\frac{17}{20}$	$\frac{14}{20}$	$\frac{10}{20}$	$\frac{5}{20}$
Valid Bond Angles	$\frac{20}{20}$	$\frac{20}{20}$	$\frac{20}{20}$	$\frac{20}{20}$	$\frac{18}{20}$	$\frac{14}{20}$
Internal Steric Clash	$\frac{20}{20}$	$\frac{20}{20}$	$\frac{20}{20}$	$\frac{20}{20}$	$\frac{20}{20}$	$\frac{20}{20}$
Internal Energy	$\frac{20}{20}$	$\frac{19}{20}$	$\frac{19}{20}$	$\frac{17}{20}$	$\frac{10}{20}$	$\frac{5}{20}$
Overall	$\frac{19}{20}$	$\frac{14}{20}$	$\frac{13}{20}$	$\frac{11}{20}$	$\frac{6}{20}$	$\frac{2}{20}$

Table 4: Empirical CDF of pass rates across 20 test peptides.

C Data Details

The TIMEWARP and MDGEN datasets consist of 'uncapped' peptides, whose termini are zwitterionic amino and carboxyl groups, as shown in the left panel of Figure 2. These are not ideal analogues of amino acids in proteins due to local charge interactions as well as lack of steric effects.

We also create a similar dataset called CAPPED 2AA of 2AA peptides by adding ACE (acetyl) and NME (N-methyl amide) caps, a common practice in molecular dynamics simulations of very small peptides. As illustrated in the right panel of Figure 2, these caps introduce additional peptide bonds with the first and last residues. These peptide bonds remove the need for the zwitterion, while the methyl group provides some steric interactions. These capping groups increase the complexity of the modelling task, but ensure a more realistic distribution of conformations. We choose the same splits as in TIMEWARP 2AA-LARGE. Since we simulated this data ourselves, we can also measure the wall-clock speed-ups of JAMUN relative to MD on this dataset.

We ensure that our unbiased molecular dynamics runs are converged or representative by comparing against biased molecular dynamics runs using Non-Equilibrium Umbrella Sampling (NEUS) [18, 83], a trajectory stratification based enhanced sampling algorithm. The protein is represented by the amber99sbildn force field [64]. The simulations are performed at 300 K with the BAOAB integrator [50] in OpenMM [21]; LINCS is used to constrain the lengths of bonds to hydrogen atoms [31]; Particle Mesh Ewald is used to calculate electrostatics [16]; the step size was 2 fs. The systems are solvated with TIP3P water models and equilibrated under NVT and NPT ensembles for 100ps each.

For the 2AA datasets, the training set consists of 50% of all possible 2AA peptides. For the 4AA datasets, the generalization task is much harder, because the number of 4AA peptides in the training sets is less than 1% and 2% respectively of the total number of possible 4AA peptides.

D Results for CAPPED 2AA

D.1 Further Results on internally simulated data

Here, we compute the ratio of the decorrelation times for the backbone and sidechain torsions in JAMUN and the reference MD data. Figure 10 highlights how sampling in the smoothed space Y compared to the original space X enables much faster decorrelation.

In Table 7, we compare JAMUN to the reference MD, shortened by a factor of 10 along the variables mentioned in Section 5, using the Jensen-Shannon distance to the full reference MD data. JAMUN

Sequence	MDGEN 4AA-EXPLICIT	JAMUN
FHSE	-675.6 ± 20.1	-633.2 ± 113.3
FKKL	-699.0 ± 24.0	-562.5 ± 192.1
FLRH	-1272.7 ± 18.9	-1198.5 ± 129.3
FSDP	-697.5 ± 23.7	-687.6 ± 91.5
FSRK	-1333.0 ± 21.9	-1349.9 ± 0.0
GCIC	-557.5 ± 21.4	-538.8 ± 56.4
GGHN	-905.3 ± 21.9	-821.3 ± 129.6
GLIL	-743.3 ± 20.7	-711.6 ± 71.8
HELI	-794.2 ± 25.5	-780.6 ± 73.5
HENV	-1156.9 ± 17.3	-1123.4 ± 148.7
HTIQ	-762.9 ± 17.0	-726.8 ± 105.8
IAMI	-428.0 ± 15.3	-426.8 ± 74.8
IDRH	-1416.8 ± 18.1	-722.8 ± 2608.1
IHNV	-845.4 ± 21.1	-864.4 ± 48.8
IMRY	-1230.7 ± 23.6	-1100.6 ± 207.0
INVH	-793.6 ± 22.9	-745.3 ± 128.6
IPGD	-611.5 ± 15.0	-582.5 ± 55.2

Table 5: Force field energies (in (kJ mol⁻¹) comparison between samples from MDGEN 4AA-EXPLICIT and JAMUN.

Dataset	Peptide Length	Capped?	Force Field	Solvent Model	Temperature	# Train	# Validation	# Test
TIMEWARP 2AA-LARGE	2	Х	amber14	Implicit Water	$310\mathrm{K}$	200	80	100
CAPPED 2AA	2	/	amber99sbildn	Explicit Water	$300\mathrm{K}$	200	80	100
TIMEWARP 4AA-LARGE	4	X	amber14	Implicit Water	$310\mathrm{K}$	1459	379	182
MDGEN 4AA-EXPLICIT	4	X	amber14	Explicit Water	$350\mathrm{K}$	3109	100	100
UNCAPPED 5AA	5	X	amber14	Implicit Water	$310\mathrm{K}$	_	_	3
CREMP 4AA	4	X	GFN2-xTB	Explicit Chloroform	$300\mathrm{K}$	15842	1000	1000

Table 6: A short description of the simulation conditions across the different datasets.

outperforms this shortened MD trajectory across all metrics, even though it takes approximately $2 \times$ longer to sample than JAMUN, from Figure 3.

Trajectory	Backbone Torsions	Sidechain Torsions	All Torsions	TICA-0	TICA-0,1	Metastable Probs
JAMUN	0.291 ± 0.119	0.320 ± 0.108	0.304 ± 0.112	0.351 ± 0.130	0.438 ± 0.117	0.264 ± 0.108
Reference ($10 \times$ shorter)	0.447 ± 0.057	0.406 ± 0.071	0.424 ± 0.056	0.557 ± 0.043	0.564 ± 0.041	0.543 ± 0.073

Table 7: Comparison of Jenson-Shannon distances between JAMUN and the reference MD (shortened by a factor of 10), averaged over the test peptides in CAPPED 2AA. Note that this shortened reference MD takes $2 \times$ longer to sample as JAMUN.

E Related Work

Protein structure prediction only requires prediction of a few folded states, and is usually trained on crystallized structures of proteins which can be quite different from their native states. The size of these proteins are significantly larger than the peptides we study here. On the other hand, conformational ensemble generation requires many samples from the Boltzmann distribution to capture the effects of solvent atoms and intramolecular interactions, even if dynamics is not modelled explicitly. Research has found that protein structure prediction models tend to struggle with capturing this diversity, especially when sampling the conformations of flexible domains [69, 13, 2] and rarer conformations [48] not found in the Protein Data Bank (PDB) [9]. Another issue is that the quality of these models' predictions can be dependent on multiple sequence alignment (MSA) information, a preprocessing step where sequence databases are queried for similar protein sequences which indicate evolutionary conservation patterns as a supplementary input to the model. In fact, some of the first attempts to sample alternative conformations of proteins with AlphaFold2 were performed by subsampling [17], clustering [88] or manipulating MSA information [29]. Sequence-to-structure

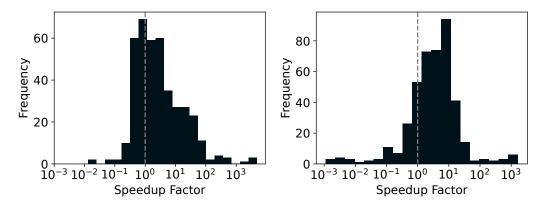


Figure 10: Speedups defined as the ratio between decorrelation times between the reference MD and JAMUN for backbone (left) and sidechain (right) torsions for all test peptides in CAPPED 2AA.

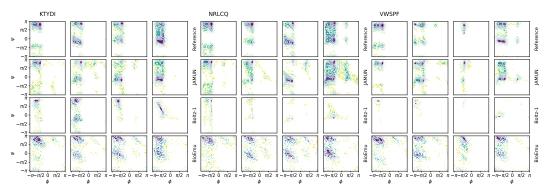


Figure 11: Ramachandran plots for test peptides NRLCQ, VWSPF and KTYDI in UNCAPPED 5AA.

models such as ESMFold [53] which do not require MSA information can be faster but tend to produce less physically accurate structures, as noted by Lu et al. [55]. Unfortunately, MSA information can be quite unreliable for small peptides due to the presence of many hits. In fact, popular MSA software suites such as MMseqs2 [80] querying will, by default, simply return an empty MSA for the peptides we study here. Nevertheless, previous research [67] has shown that AlpahFold2 can still learn an accurate energy function for protein structures without MSA co-evolution information.

Jussupow and Kaila [39] finds that AlphaFold2's predicted local distance difference test (pLDDT) and predicted aligned error (PAE) scores correlate with local protein dynamics and global conformational flexibility respectively. They use these scores to parametrize an additional harmonic potential for coarse-grained MD with the MARTINI [56] force field. AlphaFlow [36] develops flow-matching over the quotient space of 3D positions modulo rotations to learn a distribution over 3D positions of the β -carbon atoms, outperforming MSA subsampling with AlphaFold2 over the Protein Data Bank (PDB) [9]. Distributional Graphormer [92] parametrizes a diffusion model over α -carbon positions, also trained on the PDB. Str2Str [55] proposes a noising-denoising process along a range of noise levels for sampling backbone atom coordinates, followed by regression of side-chain coordinates with the rotamer-based FASPR [35] packgae. BioEmu [52] is a diffusion model built using the EvoFormer stack from AlphaFold2 [38]. BioEmu is pretrained on 200 million protein structures from the AlphaFold Protein Structure Database [85] and finetuned on over 200ms of MD data, which are orders of magnitude larger than the datasets we benchmark here.

A predecessor of many of the models discussed in 4, 'Two for One' [7] showed that the score learned by diffusion models can be used for running molecular dynamics simulations. However, as they choose the noise level for the score function close to 0, the molecular dynamics is effectively run in the original space X, not in the latent space Y as JAMUN does, which again limits the effective timestep of simulation.

There have also been efforts to build ML models for taking longer MD time-steps [46, 74, 34] and for approximating conformations of large proteins [92, 36]. These methods rely on hand-crafted featurizations (eg. backbone torsion angles). In practice, this has made generalization to unseen molecules challenging for these models as well. The above models are often classified either as Boltzmann Generators or Boltzmann Emulators. Boltzmann Generators are guaranteed to draw unbiased samples from the Boltzmann distribution, while Boltzmann Emulators do not have this guarantee. Strictly speaking, JAMUN is a Boltzmann Emulator.

JAMUN is a walk-jump sampling method which uses an SE(3)-equivariant neural network for denoising. WJS is built on the seminal work of Neural Empirical Bayes [70], and has been used in voxelized molecule generation [62, 63] and protein sequence generation [24]. Our work is the first to our knowledge to apply walk-jump sampling to point clouds.

F Overview of Denoiser

The denoiser is a SE(3)-equivariant graph neural network. The graph is defined by a radial cutoff of 10\AA in y. The overall computation performed by the denoiser is shown in Figure 12, with the initial embedding, message-passing and output head blocks shown in Figure 13, Figure 14 and Figure 15 respectively.

For all datasets, we train and sample with $\sigma = 0.4$ Å. For the Langevin dynamics (Equation 68), we set M = 1, friction of $\gamma = 1.0$ and a step size of $\Delta t = \sigma$.

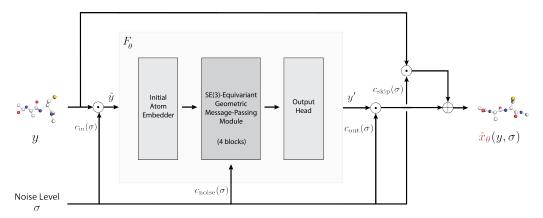


Figure 12: Overview of the denoiser network \hat{x}_{θ} . The submodule F_{θ} sees input atom coordinates $\tilde{y} = c_{\text{in}}(\sigma)y$ and outputs predicted atom coordinates y', which gets scaled and added to a noise-conditional skip connection to finally obtain $\hat{x}_{\theta}(y)$.

The hidden features $h^{(n)}$ for $n=0,\ldots,4$ contain 120 scalar and 32 vector features per atom. We use spherical harmonics up to l=1 for the tensor product.

We use the Adam optimizer with learning rate .002. Models are trained with a batch size of 32 over 2 NVIDIA RTX A100 GPUs.



Figure 13: Overview of the initial embedder in the denoiser network, creating initial features $h^{(0)}$ at each atom and edge.

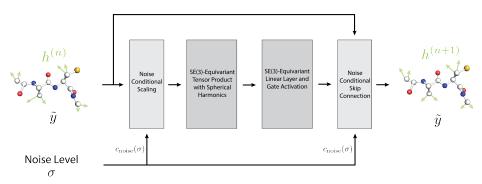


Figure 14: Overview of a single SE(3)-equivariant message-passing block (indexed by n) in the denoiser network. There are four such blocks iteratively updating the atom features from $h^{(0)}$ to $h^{(4)}$. The atom coordinates denoted by $\tilde{y}=c_{\rm in}(\sigma)y$ (and hence, the edge features) are unchanged throughout these blocks.

$$\hat{y} = \begin{bmatrix} h^{(4)} \\ h^{(4)} \\ h^{(4)} \\ h^{(4)} \end{bmatrix} \rightarrow \begin{bmatrix} SE(3) \text{ Equivariant} \\ \text{Linear Layer} \\ 1 \text{ vector per atom} \end{bmatrix} \rightarrow \begin{bmatrix} Se(3) \text{ Equivariant} \\ \text{Linear Layer} \\ \text{Vectors} \end{bmatrix} \rightarrow \begin{bmatrix} Shift Atom \\ \text{Coordinates by} \\ \text{Vectors} \end{bmatrix}$$

Figure 15: Overview of the output head, which predicts the coordinates $y' = F_{\theta}(c_{\text{in}}(\sigma)y, c_{\text{noise}}(\sigma))$.

G Normalization

As the noise level σ is increased, $y = x + \sigma \varepsilon$ where $\varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})$ expands in space. Let \tilde{y} represent the 'normalized' input y, as seen by the network F_{θ} :

$$\tilde{y} = c_{\rm in}(\sigma)y\tag{10}$$

To control the expansion of y, $c_{in}(\sigma)$ is chosen such that the following property holds:

$$\mathbb{E}_{\substack{(i,j) \sim \text{Uniform}(E) \\ \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})}} [\|\tilde{y}_i - \tilde{y}_j\|^2] = 1 \text{ at all noise levels } \sigma.$$
(11)

Note that here, we define E such that Note that this is distinct from the normalization chosen by [41, 42], which normalizes ||y|| directly. The intuition behind this normalization is that the GNN model F_{θ} does not operate on atom positions y directly, but instead uses the relative vectors $y_i - y_j$ to account for translation invariance, and controlling this object directly ensures that the topology of the graph does not change with varying noise level σ .

To achieve this, we compute:

$$c_{\rm in}(\sigma) = \frac{1}{\sqrt{C + 6\sigma^2}} \tag{12}$$

where $C = \mathbb{E}_{(i,j) \sim \text{Uniform}(E)} \|x_i - x_j\|^2$ can be easily estimated from the true data distribution. The full derivation can be found in Section G.1.

As the input is now appropriately normalized, the target output of the network F_{θ} should also be appropriately normalized. A full derivation, found in Section G.2, leads to:

$$c_{\rm skip}(\sigma) = \frac{C}{C + 6\sigma^2} \tag{13}$$

$$c_{\text{out}}(\sigma) = \sqrt{\frac{C \cdot 6\sigma^2}{C + 6\sigma^2}} \tag{14}$$

$$c_{\text{noise}}(\sigma) = \log_{10} \sigma \tag{15}$$

The noise normalization is a scaled version of the recommendation of $\frac{1}{4} \ln \sigma$ for images in Karras et al. [41, 42].

G.1 Input Normalization

Fix an $(i,j) \in E$ from Equation 11. As $\varepsilon_i, \varepsilon_j \stackrel{\text{iid}}{\sim} \mathcal{N}(0,\mathbb{I}_3)$, we have $\varepsilon_i - \varepsilon_j \sim \mathcal{N}(0,2\mathbb{I}_3)$ from the closure of the multivariate Gaussian under linear combinations. Thus, for each component d=1,2 and 3, we have: $(\varepsilon_i - \varepsilon_j)_{(d)} \sim \mathcal{N}(0,2)$ and hence:

$$\mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})}[(x_i - x_j)^T (\varepsilon_i - \varepsilon_j)] = \sum_{d=1}^3 (x_i - x_j)_{(d)} \mathbb{E}[(\varepsilon_i - \varepsilon_j)_{(d)}] = 0$$
 (16)

$$\mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})}[\|\varepsilon_i - \varepsilon_j\|^2] = \sum_{d=1}^3 \mathbb{E}[(\varepsilon_i - \varepsilon_j)_{(d)}^2] = 6$$
(17)

We can now compute:

$$\mathbb{E}_{z}[\|\tilde{y}_{i} - \tilde{y}_{j}\|^{2}]
= c_{\text{in}}(\sigma)^{2} \mathbb{E}_{\varepsilon}[\|y_{i} - y_{j}\|^{2}]
= c_{\text{in}}(\sigma)^{2} \mathbb{E}_{\varepsilon}[\|x_{i} - x_{j} + \sigma(\varepsilon_{i} - \varepsilon_{j})\|^{2}]
= c_{\text{in}}(\sigma)^{2} \left(\|x_{i} - x_{j}\|^{2} + 2\sigma \mathbb{E}_{\varepsilon}[(x_{i} - x_{j})^{T}(\varepsilon_{i} - \varepsilon_{j})] + \sigma^{2} \mathbb{E}_{\varepsilon}[\|\varepsilon_{i} - \varepsilon_{j}\|^{2}]\right)
= c_{\text{in}}(\sigma)^{2} \left(\|x_{i} - x_{j}\|^{2} + \sigma^{2} \mathbb{E}_{\varepsilon}[\|\varepsilon_{i} - \varepsilon_{j}\|^{2}]\right)
= c_{\text{in}}(\sigma)^{2} \left(\|x_{i} - x_{j}\|^{2} + 6\sigma^{2}\right).$$
(18)

Now, taking the expectation over all $(i, j) \in E$ uniformly:

$$\mathbb{E}_{\substack{(i,j) \sim \text{Uniform}(E) \\ \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})}} [\|\tilde{y}_i - \tilde{y}_j\|^2] = \mathbb{E}_{(i,j) \sim \text{Uniform}(E)} [\mathbb{E}_{\varepsilon} [\|\tilde{y}_i - \tilde{y}_j\|^2]]$$

$$= c_{\text{in}}(\sigma)^2 \left(\mathbb{E}_{(i,j) \sim \text{Uniform}(E)} \|x_i - x_j\|^2 + 6\sigma^2 \right)$$
(19)

Let $C = \mathbb{E}_{(i,j) \sim \text{Uniform}(E)} \|x_i - x_j\|^2$, which we estimate from the true data distribution. Then, from Equation 19 and our intended normalization given by Equation 11:

$$c_{\rm in}(\sigma) = \frac{1}{\sqrt{C + 6\sigma^2}} \tag{20}$$

G.2 Output Normalization

The derivation here is identical that of [41, 42], but with our normalization. The denoising loss at a single noise level is:

$$\mathcal{L}(\hat{x}_{\theta}, \sigma) = \mathbb{E}_{X \sim p_X} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})} [\|\hat{x}_{\theta}(X + \sigma \varepsilon, \sigma) - X\|^2]$$
(21)

which gets weighted across a distribution p_{σ} of noise levels by (unnormalized) weights $\lambda(\sigma)$:

$$\mathcal{L}(\hat{x}_{\theta}) = \mathbb{E}_{\sigma \sim p_{\sigma}} [\lambda(\sigma) \mathcal{L}(\hat{x}_{\theta}, \sigma)]
= \mathbb{E}_{\sigma \sim p_{\sigma}} \mathbb{E}_{X \sim p_{X}} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})} [\lambda(\sigma) \| \hat{x}_{\theta}(X + \sigma \varepsilon, \sigma) - X \|^{2}]
= \mathbb{E}_{\sigma \sim p_{\sigma}} \mathbb{E}_{X \sim p_{X}} \mathbb{E}_{Y \sim \mathcal{N}(X, \sigma^{2} \mathbb{I}_{N \times 3})} [\lambda(\sigma) \| \hat{x}_{\theta}(Y, \sigma) - X \|^{2}]
= \mathbb{E}_{\sigma \sim p_{\sigma}} \mathbb{E}_{X \sim p_{X}} \mathbb{E}_{Y \sim \mathcal{N}(X, \sigma^{2} \mathbb{I}_{N \times 3})} [\lambda(\sigma) \| c_{\text{skip}}(\sigma) Y + c_{\text{out}}(\sigma) F_{\theta}(c_{\text{in}}(\sigma) Y, c_{\text{noise}}(\sigma)) - x \|^{2}]
= \mathbb{E}_{\sigma \sim p_{\sigma}} \mathbb{E}_{X \sim p_{X}} \mathbb{E}_{Y \sim \mathcal{N}(X, \sigma^{2} \mathbb{I}_{N \times 3})} \left[\lambda(\sigma) c_{\text{out}}(\sigma)^{2} \| F_{\theta}(c_{\text{in}}(\sigma) Y, c_{\text{noise}}(\sigma)) - \frac{x - c_{\text{skip}}(\sigma) Y}{c_{\text{out}}(\sigma)} \|^{2} \right]
= \mathbb{E}_{\sigma \sim p_{\sigma}} \mathbb{E}_{X \sim p_{X}} \mathbb{E}_{Y \sim \mathcal{N}(X, \sigma^{2} \mathbb{I}_{N \times 3})} \left[\lambda(\sigma) c_{\text{out}}(\sigma)^{2} \| F_{\theta}(c_{\text{in}}(\sigma) Y, c_{\text{noise}}(\sigma)) - F \|^{2} \right] \tag{22}$$

where:

$$F(y,\sigma) = \frac{x - c_{\text{skip}}(\sigma)y}{c_{\text{out}}(\sigma)}$$
 (23)

is the effective training target for the network F_{θ} . We want to normalize F similarly as the network input:

$$\mathbb{E}_{\substack{(i,j) \sim \text{Uniform}(E) \\ \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})}} [\|F_i - F_j\|^2] = 1 \text{ at all noise levels } \sigma.$$
 (24)

Again, for a fixed $(i, j) \in E$, we have:

$$\mathbb{E}_{\varepsilon} \|F_{i} - F_{j}\|^{2} = \frac{\mathbb{E}_{\varepsilon} \|(x_{i} - x_{j}) - c_{\text{skip}}(\sigma)(y_{i} - y_{j})\|^{2}}{c_{\text{out}}(\sigma)^{2}}$$

$$= \frac{\mathbb{E}_{\varepsilon} \|(1 - c_{\text{skip}}(\sigma))(x_{i} - x_{j}) - c_{\text{skip}}(\sigma)\sigma \cdot (\varepsilon_{i} - \varepsilon_{j})\|^{2}}{c_{\text{out}}(\sigma)^{2}}$$

$$= \frac{(1 - c_{\text{skip}}(\sigma))^{2} \|x_{i} - x_{j}\|^{2} + c_{\text{skip}}(\sigma)^{2} \cdot 6\sigma^{2}}{c_{\text{out}}(\sigma)^{2}}$$
(25)

and hence:

$$\mathbb{E}_{(i,j) \sim \text{Uniform}(E)}[\|F_i - F_j\|^2] = 1$$

$$\Longrightarrow \frac{(1 - c_{\text{skip}}(\sigma))^2 \cdot C + c_{\text{skip}}(\sigma)^2 \cdot 6\sigma^2}{c_{\text{out}}(\sigma)^2} = 1$$

$$\Longrightarrow c_{\text{out}}(\sigma)^2 = (1 - c_{\text{skip}}(\sigma))^2 \cdot C + c_{\text{skip}}(\sigma)^2 \cdot 6\sigma^2$$
(26)

where C was defined above. Now, to minimize $c_{\text{out}}(\sigma)$ to maximize reuse and avoid amplifying network errors, as recommended by Karras et al. [41, 42]:

$$\frac{d}{dc_{\text{skip}}(\sigma)}c_{\text{out}}(\sigma)^{2} = 0$$

$$\implies -2(1 - c_{\text{skip}}(\sigma)) \cdot C + 2c_{\text{skip}}(\sigma) \cdot 6\sigma^{2} = 0$$

$$\implies c_{\text{skip}}(\sigma) = \frac{C}{C + 6\sigma^{2}}$$
(27)

Substituting into Equation 26, we get after some routine simplification:

$$c_{\text{out}}(\sigma) = \sqrt{\frac{C \cdot 6\sigma^2}{C + 6\sigma^2}}$$
 (28)

The noise normalization is chosen as $c_{\text{noise}}(\sigma) = \log_{10} \sigma$, a scaled version of the recommendation of $\frac{1}{4} \ln \sigma$ for images in Karras et al. [41, 42].

From Equation 22, we set $\lambda(\sigma) = \frac{1}{c_{\text{out}}(\sigma)^2}$ to normalize the loss at at all noise levels, as in Karras et al. [41, 42].

G.3 Rotational Alignment

As described in Algorithm 1, we use the Kabsch-Umeyama algorithm [40, 82] to rotationally align y to x before calling the denoiser.

Algorithm 1 Rotational Alignment with the Kabsch-Umeyama Algorithm

Note that both y and x are mean-centered to respect translational equivariance:

$$\sum_{i=1}^{N} y_i = \vec{0} \in \mathbb{R}^3 \tag{29}$$

$$\sum_{i=1}^{N} x_i = \vec{0} \in \mathbb{R}^3 \tag{30}$$

so there is no net translation.

H Proofs of Theoretical Results

For completeness, we prove the main theoretical results here, as first established by Robbins [66], Miyasawa [58], Saremi and Hyvärinen [70].

H.1 The Denoiser Minimizes the Expected Loss

Here, we prove Equation 6, rewritten here for clarity:

$$\hat{x}(\cdot) \equiv \mathbb{E}[X \mid Y = \cdot] = \operatorname*{arg\,min}_{f:\mathbb{R}^{N\times3} \to \mathbb{R}^{N\times3}} \mathbb{E}_{X \sim p_X, \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N\times3})} [\|f(Y) - X\|^2]$$
(31)

First, we can decompose the loss over the domain $\mathbb{R}^{N\times 3}$ of Y:

$$\mathbb{E}_{X \sim p_X, \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})}[\|f(Y) - X\|^2] = \mathbb{E}_{X \sim p_X, Y \sim p_Y}[\|f(Y) - X\|^2]$$

$$(32)$$

$$= \int_{\mathbb{R}^{N\times 3}} \int_{\mathbb{R}^{N\times 3}} \|f(y) - x\|^2 p_{X,Y}(x,y) dx dy$$
 (33)

$$= \int_{\mathbb{R}^{N\times 3}} \underbrace{\int_{\mathbb{R}^{N\times 3}} \|f(y) - x\|^2 p_{Y|X}(y \mid x) p_X(x) dx dy}_{l(f,y)}$$
(34)

$$= \int_{\mathbb{R}^{N\times 3}} l(f, y) dy \tag{35}$$

where $l(f,y) \geq 0$ for all functions f and inputs y. Hence, any minimizer f^* must minimize the local denoising loss $l(f^*,y)$ at each point $y \in \mathbb{R}^{N \times 3}$. For a fixed $y \in \mathbb{R}^{N \times 3}$, the loss l(f,y) is convex as a function of f(y). Hence, the global minimizer can be found by finding the critical points of l(f,y) as a function of f(y):

$$\nabla_{f(y)}l(f,y) = 0 \tag{36}$$

$$\implies \nabla_{f(y)} \int_{\mathbb{R}^{N \times 3}} \|f(y) - x\|^2 p_{Y|X}(y \mid x) p_X(x) dx = 0$$
 (37)

$$\Longrightarrow \int_{\mathbb{R}^{N\times 3}} 2(f^*(y) - x) p_{Y|X}(y \mid x) p_X(x) dx = 0 \tag{38}$$

Rearranging:

$$f^*(y) = \frac{\int_{\mathbb{R}^{N \times 3}} x \, p_{Y|X}(y \mid x) p_X(x) dx}{\int_{\mathbb{R}^{N \times 3}} p_{Y|X}(y \mid x) p_X(x) dx}$$
(39)

$$= \frac{\int_{\mathbb{R}^{N\times 3}} x \, p_{Y|X}(y \mid x) p_X(x) dx}{p_Y(y)} \tag{40}$$

$$= \int_{\mathbb{R}^{N\times 3}} x \, \frac{p_{Y|X}(y \mid x)p_X(x)}{p_Y(y)} dx \tag{41}$$

$$= \int_{\mathbb{R}^{N\times 3}} x \, p_{X|Y}(x \mid y) dx \tag{42}$$

$$= \mathbb{E}[X \mid Y = y] \tag{43}$$

$$=\hat{x}(y) \tag{44}$$

by Bayes' rule. Hence, the denoiser as defined by Equation 5 is indeed the minimzer of the denoising loss:

$$\hat{x}(\cdot) \equiv \mathbb{E}[X \mid Y = \cdot] = \underset{f:\mathbb{R}^{N \times 3} \to \mathbb{R}^{N \times 3}}{\arg \min} \mathbb{E}_{X \sim p_X, \varepsilon \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})} [\|f(Y) - X\|^2]$$
(45)

as claimed.

H.2 Relating the Score and the Denoiser

Here, we rederive Equation 7, relating the score function $\nabla \log p_Y$ and the denoiser \hat{x} .

Let $X \sim p_X$ defined over $\mathbb{R}^{N \times 3}$ and $\eta \sim \mathcal{N}(0, \mathbb{I}_{N \times 3})$. Let $Y = X + \sigma \eta$, which means:

$$p_{Y|X}(y \mid x) = \mathcal{N}(y; x, \mathbb{I}_{N \times 3}) = \frac{1}{(2\pi\sigma^2)^{\frac{3N}{2}}} \exp\left(-\frac{\|y - x\|^2}{2\sigma^2}\right)$$
 (46)

Then:

$$\mathbb{E}[X \mid Y = y] = y + \sigma^2 \nabla_y \log p_Y(y)$$
(47)

To prove this:

$$\nabla_{y} p_{Y|X}(y \mid x) = -\frac{y - x}{\sigma^{2}} p_{Y|X}(y \mid x)$$
 (48)

$$\implies (x - y)p_{Y|X}(y \mid x) = \sigma^2 \nabla_y p_{Y|X}(y \mid x) \tag{49}$$

$$\implies \int_{\mathbb{R}^{N\times 3}} (x-y) p_{Y|X}(y\mid x) \ p_X(x) dx = \int_{\mathbb{R}^{N\times 3}} \sigma^2 \nabla_y p_{Y|X}(y\mid x) \ p_X(x) dx \tag{50}$$

By Bayes' rule:

$$p_{Y|X}(y \mid x)p_X(x) = p_{X,Y}(x,y) = p_{X|Y}(x|y)p_Y(y)$$
(51)

and, by definition of the marginals:

$$\int_{\mathbb{R}^{N\times 3}} p_{X,Y}(x,y)dx = p_Y(y) \tag{52}$$

For the left-hand side, we have:

$$\int_{\mathbb{R}^{N\times 3}} (x-y) p_{Y|X}(y \mid x) p_X(x) dx = \int_{\mathbb{R}^{N\times 3}} (x-y) p_{X,Y}(x,y) dx$$
 (53)

$$= \int_{\mathbb{R}^{N\times3}} x p_{X,Y}(x,y) dx - \int_{\mathbb{R}^{N\times3}} y p_{X,Y}(x,y) dx \tag{54}$$

$$= p_Y(y) \left(\int_{\mathbb{R}^{N \times 3}} x p_{X|Y}(x \mid y) dx - y \int_{\mathbb{R}^{N \times 3}} p_{X|Y}(x \mid y) dx \right)$$

$$\tag{55}$$

$$= p_Y(y) \left(\mathbb{E}[X \mid Y = y] - y \right) \tag{56}$$

For the right-hand side, we have:

$$\sigma^2 \int_{\mathbb{R}^{N\times 3}} \nabla_y p_{Y\mid X}(y\mid x) p_X(x) dx = \sigma^2 \nabla_y \int_{\mathbb{R}^{N\times 3}} p_{Y\mid X}(y\mid x) p_X(x) dx \tag{57}$$

$$= \sigma^2 \nabla_y \int_{\mathbb{D}_{N\times 3}} p_{X,Y}(x,y) dx \tag{58}$$

$$= \sigma^2 \nabla_y p_Y(y) \tag{59}$$

Thus.

$$p_Y(y) \left(\mathbb{E}[X \mid Y = y] - y \right) = \sigma^2 \nabla_y p_Y(y) \tag{60}$$

$$\implies \mathbb{E}[X \mid Y = y] = y + \sigma^2 \frac{\nabla_y p_Y(y)}{p_Y(y)} \tag{61}$$

$$= y + \sigma^2 \nabla_y \log p_Y(y) \tag{62}$$

as claimed.

I Numerical Solvers for Langevin Dynamics

As mentioned in Section 2.2, solving the Stochastic Differential Equation corresponding to Langevin dynamics is often performed numerically. In particular, BAOAB [49, 51, 68] refers to a 'splitting method' that solves the Langevin dynamics SDE by splitting it into three different components labelled by \mathcal{A} , \mathcal{B} and \mathcal{O} below:

$$dy = \underbrace{v_y dt} \tag{63}$$

$$dv_y = \underbrace{M^{-1}\nabla_y \log p_Y(y)dt}_{\mathcal{B}} - \underbrace{\gamma v_y dt + \sqrt{2\gamma}M^{-\frac{1}{2}}dB_t}_{\mathcal{O}}$$
(64)

where both $y, v_y \in \mathbb{R}^d$. This leads to the following update operators:

$$\mathcal{A}_{\Delta t} \begin{bmatrix} y \\ v_y \end{bmatrix} = \begin{bmatrix} y + v_y \Delta t \\ v_y \end{bmatrix} \tag{65}$$

$$\mathcal{B}_{\Delta t} \begin{bmatrix} y \\ v_y \end{bmatrix} = \begin{bmatrix} y \\ v_y + M^{-1} \nabla_y \log p_Y(y) \Delta t \end{bmatrix}$$
 (66)

$$\mathcal{O}_{\Delta t} \begin{bmatrix} y \\ v_y \end{bmatrix} = \begin{bmatrix} y \\ e^{-\gamma \Delta t} v_y + M^{-\frac{1}{2}} \sqrt{1 - e^{-2\gamma \Delta t}} B \end{bmatrix}$$
 (67)

where $B \sim \mathcal{N}(0, \mathbb{I}_d)$ is resampled every iteration. As highlighted by Kieninger and Keller [43], the \mathcal{A} and \mathcal{B} updates are obtained by simply discretizing the updates highlighted in Equation 63 by the Euler method. The \mathcal{O} update refers to a explicit solution of the Ornstein-Uhlenbeck process, which we rederive for completeness in Appendix J.

Finally, the iterates of the BAOAB algorithm are given by a composition of these update steps, matching the name of the method:

$$\begin{bmatrix} y^{(t+1)} \\ v_{y}^{(t+1)} \end{bmatrix} = \mathcal{B}_{\frac{\Delta t}{2}} \mathcal{A}_{\frac{\Delta t}{2}} \mathcal{O}_{\Delta t} \mathcal{A}_{\frac{\Delta t}{2}} \mathcal{B}_{\frac{\Delta t}{2}} \begin{bmatrix} y^{(t)} \\ v_{y}^{(t)} \end{bmatrix}$$
(68)

J The Ornstein-Uhlenbeck Process

For completeness, we discuss the distributional solution of the Ornstein-Uhlenbeck process, taken directly from the excellent Leimkuhler and Matthews [51]. In one dimension, the Ornstein-Uhlenbeck Process corresponds to the following Stochastic Differential Equation (SDE):

$$dv_y = -\gamma v_y dt + \sqrt{2\gamma} M^{-\frac{1}{2}} dB_t \tag{69}$$

Multiplying both sides by the integrating factor $e^{\gamma t}$:

$$e^{\gamma t}dv_y = -\gamma e^{\gamma t} (v_y dt + e^{\gamma t} \sqrt{2\gamma} M^{-\frac{1}{2}} dB_t$$
 (70)

$$\implies e^{\gamma t}(dv_y + \gamma v_y dt) = e^{\gamma t} \sqrt{2\gamma} M^{-\frac{1}{2}} dB_t \tag{71}$$

and identifying:

$$e^{\gamma t}(dv_y + \gamma v_y dt) = d(e^{\gamma t}v_y) \tag{72}$$

We get after integrating from t_1 to t_2 , two adjacent time steps of our integration grid:

$$d(e^{\gamma t}v_y) = e^{\gamma t}\sqrt{2\gamma}M^{-\frac{1}{2}}dB_t \tag{73}$$

$$\implies \int_{t_1}^{t_2} d(e^{\gamma t} v_y) = \int_{t_1}^{t_2} e^{\gamma t} \sqrt{2\gamma} M^{-\frac{1}{2}} dB_t \tag{74}$$

$$\implies e^{\gamma t_2} v_y(t_2) - e^{\gamma t_1} v_y(t_1) = \sqrt{2\gamma} M^{-\frac{1}{2}} \int_{t_1}^{t_2} e^{\gamma t} dB_t$$
 (75)

Now, for a Wiener process B_t , if g(t) is a deterministic function, $\int_{t_1}^{t_2} g(t) dB_t$ is distributed as $\mathcal{N}\left(0, \int_{t_1}^{t_2} g(t)^2 dt\right)$ by Itô's integral. Thus, applying this result to $g(t) = e^{\gamma t}$, we get:

$$e^{\gamma t_2} v_y(t_2) - e^{\gamma t_1} v_y(t_1) = \sqrt{2\gamma} M^{-\frac{1}{2}} \mathcal{N} \left(0, \frac{e^{2\gamma t_2} - e^{2\gamma t_1}}{2\gamma} \right)$$
 (76)

$$\implies v_y(t_2) = e^{-\gamma(t_2 - t_1)} v_y(t_1) + \sqrt{2\gamma} M^{-\frac{1}{2}} e^{-\gamma t_2} \mathcal{N}\left(0, \frac{e^{2\gamma t_2} - e^{2\gamma t_1}}{2\gamma}\right) \tag{77}$$

$$= e^{-\gamma(t_2 - t_1)} v_y(t_1) + \sqrt{2\gamma} M^{-\frac{1}{2}} \sqrt{\frac{1 - e^{2\gamma(t_1 - t_2)}}{2\gamma}} \mathcal{N}(0, 1)$$
 (78)

$$=e^{-\gamma(t_2-t_1)}v_y(t_1)+M^{-\frac{1}{2}}\sqrt{1-e^{2\gamma(t_1-t_2)}}\mathcal{N}\left(0,1\right) \tag{79}$$

In the $N \times 3$ dimensional case, as the Wiener processes are all independent of each other, we directly get:

$$v_y(t_2) = e^{-\gamma(t_2 - t_1)} v_y(t_1) + M^{-\frac{1}{2}} \sqrt{1 - e^{2\gamma(t_1 - t_2)}} \mathcal{N}(0, \mathbb{I}_{N \times 3})$$
(80)

Setting $\Delta t = t_2 - t_1$, we get the form of the \mathcal{O} operator (Equation 65) of the BAOAB integrator in Appendix I.

K Parallelizing Sampling with Multiple Independent Chains

Our sampling strategy batches peptides in order to increase throughput. Another potential method to increase throughput (but which we did not employ for the results in this paper) is to sample multiple chains in parallel.

This can be done by initializing multiple chains: $y_1^{(0)}, \dots, y_{N_{\mathrm{ch}}}^{(0)},$ where:

$$y_{\rm ch}^{(0)} = x^{(0)} + \sigma \varepsilon_{\rm ch}^{(0)}$$
 (81)

where $\varepsilon_{\mathrm{ch}}^{(0)} \stackrel{\mathrm{iid}}{\sim} \mathcal{N}(0,\mathbb{I}_{N \times 3})$ for $\mathrm{ch} = 1,\ldots,N_{\mathrm{ch}}$ are all independent of each other. Then, the chains can be evolved independently withindependent walk steps (Equation 3) and denoised with independent jump steps (Equation 7). This independence allows batching over the $y_{\mathrm{ch}}^{(t)}$ over all chains ch at each iteration t.

Note that at t=0, the chains are correlated as they are all initialized from the same $x^{(0)}$. However, if the number of samples per chain is large enough, the chains are no longer correlated, as they have now mixed into the stationary distribution.

L Model transferability exploration using Macrocycles

One of the most important aspects of JAMUN is its highly general input, a point cloud, unlike other protein ensemble models that are frequently more bespoke, using dihedral or frame representations. While this may be a slight disadvantage for us in the protein space, it also makes us extremely flexible and easily transferable to other modalities. Here, we demonstrate that on macrocyclic peptides with several non-canonical residues.

The exploration of conformational ensembles in macrocyclic peptides is crucial due to their emerging role as therapeutic modalities. These molecules present significant challenges in computational modeling because of their conformational diversity and inherent geometric constraints. In fact, molecular dynamics trajectories for these molecules are particularly slow as good classical forcefields are unavailable and it is necessary to use quantum mechanical calculations to compute forces. Macrocyclic peptides are also extremely unwieldy in their "open", most common conformations, forming hydrogen bond networks with water. However, those macrocycles that are able to occupy smaller "crumpled" conformations are greasy and able to permeate through biological membranes, making them more suitable for biodelivery. Here, as an example, we use macrocycles from the CREMP dataset [28], generate ensembles with the CREST protocol [65], and benchmark the resulting conformers against the RINGER model [27]. Figure 16 illustrates the transferability of JAMUN to macrocyclic peptides.

It is clear that we are able to recover most basins sampled, even though it does seem like there are new basins uncovered. We note that our outputs look significantly more diffusive than the ground truth. The main reason for this is that the CREST data is clustered and filtered to represent the local minima, whereas JAMUN is designed to sample entire distributions. In some sense, the data is strictly not complete for the task JAMUN is designed for. It is impressive that in spite of this JAMUN learns enough from the denoising "jump" step with very local data to still perform the Langevin dynamics "walk" step and sample multiple basins.

We find that for 4-mers, which is what we trained our model for, we are able to recover all the sampled basins. We also attempt to run inference for 5 and 6-mers with the same model to test generalizability.

While this is a preliminary study, it points to the potential of JAMUN being used universally, not just for proteins, and in particular shows that it is effective even in a low-data regime.

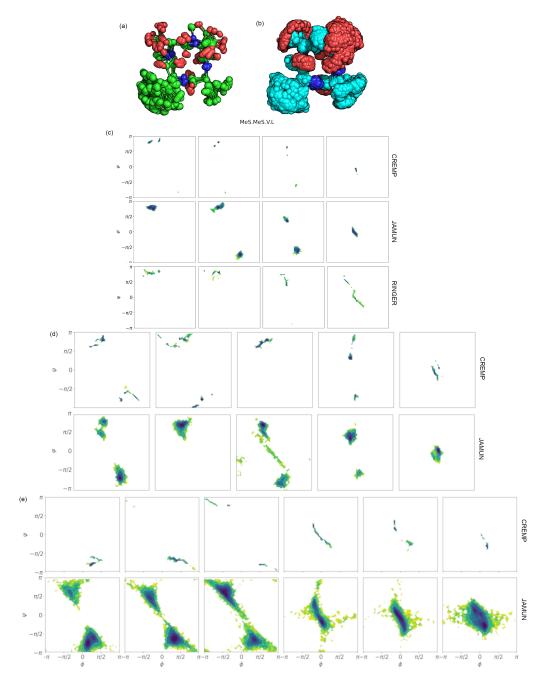


Figure 16: Macrocycle results trained on 4AA: 3D image of the (a) CREMP (green) and (b) JAMUN (cyan) MeS.MeS.V.L macrocycle. (c) Ramachandran plot for CREMP, JAMUN, and RINGER samples of the 4AA MeS.MeS.V.L macrocycle. (d) Ramachandran plot for CREMP, JAMUN, and RINGER samples of the 5AA F.Q.L.G.Met macrocycle.(e) Ramachandran plot for CREMP and JAMUN the 6AA Mes.T.Q.Mei.V.W macrocycle.