# 002 003

004

000

001

006 007

008 009 010

011 012

013 014 015

016

018 019 021

023 024 025

026

027

028 029

031 032

033 034 035

037 040 041

042

043

044

046 047 048

050

051 052

# ON THE PRACTICALITY OF BOLTZMANN NEURAL SAMPLERS

#### Anonymous authors

Paper under double-blind review

#### **ABSTRACT**

We tackle a challenge at the heart of the missions of computational chemistry and biophysics—to sample a Boltzmann-type distribution

$$p(\mathbf{x}|\mathcal{G}) \propto e^{-U(\mathbf{x}|\mathcal{G})}$$
 (1)

on  $\mathbb{R}^{N\times 3}$  associated with some N-body system  $\mathcal{G}$ , where U is an energy function (termed force field) with orthogonal invariance and deep, isolated minima. Traditionally, this is sampled sequentially using Markov chain Monte Carlo methods, which can be so slow that one, for weeks of wall time, never breaks free from the local minima defined by the starting pose. Neural samplers have been designed to speed up this process by optimizing the dynamics, prescribed by a stochastic differential equation (SDE). Though sound and elegant in continuous time, they can be practically unstable and inefficient when discretized. In this paper, we attribute this phenomena to the limited expressiveness of the finite additive transition kernels, and their inability to bridge distant distributions. To remedy this, we design a new type of highly flexible prior by mixing orthogonally invariant densities (Mint), as well as a new discretized non-volume-preserving kernel, termed Jacobian-unpreserving Langevin with explicit projection (Julep). Together, MintJulep greatly improve the practical performance of neural samplers, while keeping the underlying SDE intact.

#### INTRODUCTION: BOLTZMANN DISTRIBUTION AND NEURAL SAMPLERS

Statistical mechanics, some [1] say, bridges the microscopic and the macroscopic world,

$$\bar{\mathcal{O}} = \int d\mathbf{x} \mathcal{O}(\mathbf{x}) p(\mathbf{x}), \tag{2}$$

with the probability distribution p, conditioned on some N-body system  $\mathcal{G}$ , adopting the Boltzmann [2] form (Equation 1), known up to a constant. On one end of the bridge are per-frame  $(\mathbf{x} \in \mathbb{R}^{N \times 3})$  computable quantities  $\mathcal{O}(\mathbf{x})$ ; on the other,  $\bar{\mathcal{O}}$ , some ensemble observable tangibly measurable in laboratories, such as the binding affinity of a newly designed therapeutics, or the physical properties of an innovative material. As such, to draw samples from Equation 1 in an efficient and unbiased manner to estimate Equation 2, will shed quantitative light on the understandings and discoveries spanning various domains, from chemistry, material science, to biophysics. Many machine learning pipelines in these disciplines can be seen as approximating (force field construction [3–10]) or minimizing (conformer generation [11], docking [12, 13], and protein folding [14– 16]) the Boltzmann distribution. Nevertheless, if one wishes to rigorously sample such distribution till convergence, Monte Carlo methods are typically needed, known as molecular dynamics (MD) simulations [17–19], which is slow and biased towards the starting pose, due to the sequential nature.

**Preliminaries.** The aforementioned sampling process typically involves integrating a SDE (from t=0 to 1 without loss of generality), using, for instance the overdamped Langevin dynamics,

$$dX = -\epsilon \nabla U_t dt + \sqrt{2\epsilon} dB, X_{t=0} \sim q_0$$
(3)

where  $\epsilon$  denotes the volatility (inverse friction).  $U_1 = U$  is required to target the correct Boltzmann distribution. A constant  $U_{t\in[0,1]}=U$  with a static starting point  $q_0=\delta(\mathbf{x})$  represents the traditional sampling process ubiquitously used by MD practitioners. Alternatively, a tractable starting point, such as a wide isotropic distribution  $q_0 = \mathcal{N}(0, \sigma I), \sigma >> 0$ , together with a linearly annealing potential  $U_t = tU + (1-t)(-\log \mathcal{N})$ , recovers the simulated annealing approach. When  $U_t$  is not constant, the quantity

$$W(X) = \int_0^1 \mathrm{d}t \partial_t U_t(X),\tag{4}$$

known as generalized work for physicists and path weights for statisticians, can be used as to correct for the bias introduced in this process when estimating arbitrary functions f (annealed importance sampling, AIS [20]), as well as to estimate the ratio of normalizing constants between  $p_0$  and  $p_1$  (Jarzynski's equality [21], JE):

$$\int dX f(X) p_1(X) = \frac{\mathbb{E}[e^W f(X_1)]}{\mathbb{E}[e^W]}; \frac{Z_1}{Z_0} = \frac{\int dX_1 p_1(X_1)}{\int dX_0 p_0(X_0)} = \mathbb{E}[e^W].$$
 (5)

**Problem statement & related works.** To improve the convergence of  $X_1$  to Equation 1, a non-equilibrium control [22]  $b_t$  can be added to the drift term in Equation 3, resulting in the *nonequilibrium annealing process*:

$$d\overrightarrow{X} = -\epsilon \nabla U_t dt + \sqrt{2\epsilon} dB + \frac{\mathbf{b_t}}{2} dt, \tag{6}$$

with the corresponding path weights:

$$W(X) = \int_0^1 \mathrm{d}t(-\nabla \cdot b_t(X) + \nabla U_t(X) \cdot b_t(X) + \partial_t U_t(X)). \tag{7}$$

A perfect control term [23] exists so that  $X_0$  can be transported to exactly match  $X_1$  to Equation 1, mitigating the need of reweighting, i.e.,  $W \equiv 0$ . We call the neural parametrization and optimization of  $b_t$  towards this goal neural samplers. For this purpose, the most obvious choice of objectives seeks to formulate this as a stochastic optimal control (SOC) problem [24–27] minimizing the control energy and the terminal reverse-KL divergence  $D_{\rm KL}[q_1||p_1]$ , where  $q_1$  is governed by the law of Equation 6. These online approaches require the differentiation through the SDE integration, and can therefore be expensive or unstable. When  $\epsilon=0$ , the deterministic counterpart of Equation 6 reduces to a (continuous) normalizing flow [28–30], referred to as Boltzmann generators [31, 32] in our context. To speed up convergence and prevent mode-collapsing ubiquitous in reverse-KL-based methods requiring only the energy function (energy-based training), these types of approaches typically requires samples from  $p_1$  (data-based training) to evaluate the forward-KL  $D_{\rm KL}[p_1||q_1]$ , incurring an overhead. Overall, offline methods relying on neither samples nor differentiable trajectories seem theoretically attractive for scalability. Specifically, consider a backward SDE with  $X_0 \sim p_1$ :

$$d\overleftarrow{X} = -\epsilon \nabla U_t dt + \sqrt{2\epsilon} d\overleftarrow{B} - b_t dt.$$
 (8)

With Equation 7, the *controlled* [33] Crook's fluctuation theorem reads exactly like (and recovers, when  $b \equiv 0$ ), the original Crook's fluctuation theorem (CLT [34]):

$$d\overrightarrow{\mathbb{P}}/d\overleftarrow{\mathbb{P}} = \exp(W - Z_0 + Z_1), \tag{9}$$

where the derivative is taken in the Radon-Nikodym (RND) sense, and  $\overrightarrow{\mathbb{P}}$ ,  $\overleftarrow{\mathbb{P}}$  are the path measure associated with Equation 6, 8, respectively. This furthermore generalizes JE (5) since  $\mathbb{E}[d\overrightarrow{\mathbb{P}}/d\overleftarrow{\mathbb{P}}] = 1$ . Although many divergences can be employed to find the perfect control term [35], such as the physics-inspired neural network (PINN [36]) or the action matching (AM [37]) loss, the most straightforward offline method [33, 38] minimizes the log-variance of the RND (9):

$$\mathcal{L} = \mathbb{V}[\log[d\overrightarrow{\mathbb{P}}/d\overleftarrow{\mathbb{P}}]],\tag{10}$$

with can be taken w.r.t. any measure (hence the offline nature), albeit usually w.r.t.  $\overrightarrow{\mathbb{P}}$  for minimal variance. When this approaches zero,  $b_t$  is perfect since  $W \equiv 0$  and  $\overleftarrow{\mathbb{P}}$  is exactly the time reversal [39] of  $\overrightarrow{\mathbb{P}}$ .

**Pathology:** Why neural samplers fail in practice? While the aforementioned formulation is simple and elegant in theory (continuous-time), practically, when discretized, it fails to perform when realistic physical systems are involved [26]. We postulate that this can be attributed to the limited expressiveness of the discretized kernel, and its inability to bridge drastically distant distributions. Slightly formally, we observe that:

Remark 1.1. To transport  $q_0$  to  $q_1$  with 1-Wasserstein distance  $W_1$ , using discrete transition kernel:

$$k_{+}(X_{t+\Delta t}|X_{t}) = \mathcal{N}(X_{t} + b_{t}(X_{t})\Delta t, \sqrt{2\epsilon}I)$$
(11)

where  $b_t \Delta t$  are L-Lipschitz continuous, at least  $\log_L W_1$  kernels are needed.

**Main contributions.** Motivated by the aforementioned pathology, while keeping the SDE (Equations 6, 8) and the objective (Equation 9) intact, we propose:

- A new prior called **Mint** (mixture of invariant densities, §2), where we achieve high parametrized flexibility while respecting the symmetry of physical systems.
- A new discretized kernel termed **Julep** (Jacobian-unpreserving Langevin with explicit projection, §3), which adds additional expressiveness to each step by allowing not only additive but also multiplicative transformations.

Further relating the discoveries to prior literature, we note that [40] also optimizes the prior of the sequential Monte Carlo process while leaving the actual annealing dynamics invariant, albeit using a much more detailed but expensive invertible flow model evaluated only once during the SDE integration. The Julep kernel, on the other hand, can be regarded as a continuous stochastic normalizing flow model [27] sandwiched by deterministic bijections built with matrix exponential [41], of which the time-discretized integration on a graph manifold is inspired by [42]. In §4, we show that these two innovations greatly enhance the feasibility of Boltzmann neural samplers, bringing us one step closer to the efficient and scalable sampling of Boltzmann distributions for physical systems.

#### 2 MINT PRIOR: MIXTURE OF ORTHOGONALLY INVARIANT DENSITIES

To model the highly irregular distributions defined by the force field u in Equation 1, which we know very little except that it is orthogonally invariant (w.r.t. internal rotation or reflection, definition above), we ask the question: whether it is possible to find a class of distribution that can approximate any arbitrary distributions on  $\mathbb{R}^{N \times n}$  up to the orthogonal symmetry group O(n)? Formally:

**Definition 2.1.** A function f is said to be orthogonally invariant on  $\mathbb{R}^{N\times n}$  if  $\forall X\in\mathbb{R}^{N\times n},Q\in\mathbb{R}^{d\times d},QQ^{\top}=Q^{\top}Q=I,$ 

$$f(X) = f(XQ). (12)$$

A distribution is said to be orthogonally invariant if its density function satisfies Equation 12.

It is also worth noting that, while most of the force fields used by computational physicists and chemists are actually E(n)-, rather than O(n)-invariant, we constrain the translational degrees of freedom here and work with a radial, internal coordinate system. Practically, this can be done by consistently placing one particle, termed *anchor atom* henceforth (See an illustration in Figure 1.), at the origin of the coordinate system. In addition, although we are more interested in n=3, all results shown in this paper can be generalized to all  $n\in Z^+$ , as we do not reply on any operators other than the dot product. In the following section, we firstly handle the angular degrees of freedom, before we proceed to the radial part and derive a parametric class of distribution capable of approximating any densities up to the defined symmetry.

**Pairwise von Mises-Fisher** (PvMF) **distribution:** O(n)-invariant on  $(\mathbb{S}^{n-1})^N$ . To start, we propose a new class of distribution termed the *pairwise von Mises-Fisher* distribution, on the manifold  $(\mathbb{S}^{n-1})^N$ , which is the N-product of  $\mathbb{S}^{n-1}$  spheres:

**Definition 2.2** (PvMF distribution).

$$PvMF(\Theta; \mu, \kappa) \propto \exp(\kappa \cos(\forall ec(\Theta\Theta^{\top}), vec(\mu\mu^{\top}))),$$
(13)

where  $\kappa \in \mathbb{R}^+$  and  $\Theta, \mu \in (\mathbb{S}^{n-1})^N$ . Simply put, the energy function measures the *cosine similarity* between the gram matrices defined respectively by the variable  $\Theta$  and the parameter  $\mu$ , both on the  $(\mathbb{S}^{n-1})^N$  manifold, and prescribes the density similar to the vanilla von Mises-Fisher (vMF) distribution [43]. Since such a density peaks at the perfect alignment of X and  $\mu$  up to an orthogonal transformation Q, with  $\text{PvMF}(\mu Q, \mu) \propto \exp(\kappa) < \infty$ , and therefore the integration over a finite-volume manifold  $Z = \int_{\Theta \in (\mathbb{S}^d)^N} \text{d}\Theta \text{PvMF}(\Theta, \mu) < \infty$  is normalizable. Besides, the gram operator is O(n)-invariant, so it naturally follows that:

*Remark* 2.3. PvMF( $\Theta$ ;  $\mu$ ,  $\kappa$ ) is a valid probability distribution.

Remark 2.4.  $PvMF(\Theta; \mu, \kappa) = PvMF(\Theta Q; \mu, \kappa), \forall Q^{\top}Q = I$  is orthogonally invariant.

Evidently, this probability density requires  $\mathcal{O}(N^2)$  runtime complexity to evaluate.

Mixture of  $\operatorname{PvMFLogNormal}$  products: universally approximative. Having the angular degrees of freedom  $(\Theta = X/\|X\| \in (\mathbb{S}^{n-1})^N)$  taken care of, we pair the  $\operatorname{PvMF}$  distribution with a simple  $\operatorname{LogNormal}$  distribution on the radial axis  $(r = \|X\| \in \mathbb{R}^N)$  and can now define a distribution on the entire  $X \in \mathbb{R}^{N \times d}$  space, called the  $\operatorname{PvMFLogNormal}$  product, which stays orthogonally invariant:

$$PvMFLogNormal(\Theta, r; \mu, \kappa, \rho, \sigma) = PvMF(\Theta, \kappa)LogNormal(r; \rho, \sigma), \tag{14}$$

where  $\Theta \in (\mathbb{S}^{n-1})^N, r \in N$  and  $\mu \in (\mathbb{S}^{n-1})^N, \kappa \in \mathbb{R}^+, \rho, \sigma \in \mathbb{R}^N$ . We now arrive at the complete form of the family distribution used henceforth—the mixture of PvMFLogNormal products, followed by the expressive characterization.

Definition 2.5 (Mixture definition).

$$q(X; \{\pi_i, \mu_i, \kappa_i, \rho_i, \sigma_i\}) = \sum_i \pi_i \text{PvMFLogNormal}(X/||X||, ||X||; \mu, \kappa, \rho, \sigma),$$
 (15)

**Theorem 2.6** (Universal approximator). *Mixture of* PvMFLogNormal distributions with sufficient components can approximate any arbitrary orthogonally invariant distributions on  $\mathbb{R}^{N\times d}$  with arbitrarily small error.

The proof, defferred to the appendix, follows the first fundamental theorem of the orthogonal group [44] and the style of the universal approximation theorem of the Gaussian mixture models [45].

Sampling and energy-based variational inference (VI). While the family of distribution is defined, we realize that sampling from this distribution (or the PvMF distribution itself) is highly non-trivial. Recall that, for regular vMF distributions, in the high concentration limit ( $\kappa \to \infty$ ), its behavior converges to that of a projected normal distribution, which is easy to sample. Following the same procedure ([43] §3.5.22 and 9.3.15) to Taylor-expand the density on the tangent space  $I - \mu \mu^{\top}$ , we see that our PvMF distribution, when concentrated, can also be approximated by a projected normal distribution rotated by an arbitrary angle Q (or reflection):

$$PvMF(\Theta; \mu, \kappa) \approx \mathcal{PN}(\Theta; \mu Q, 1/\sqrt{\kappa}), \forall Q, \kappa \to \infty.$$
 (16)

This approximation efficiently generates samples from the PvMF distribution, which can be further corrected by a brief Langevin dynamics integration. If we are working with a problem where the loss function is also orthogonally invariant, as is in this paper, the rotation can also be practically emitted Q = I. These samples are then used downstream to be multiplied by the radial components and blended into a mixture.

Samples from q (Equation 15) at hand, we can easily fit this highly flexible function to arbitrary target densities p by optimizing, for instance, the reverse KL divergence  $D_{\mathrm{KL}}[q||p] = \mathbb{E}_q[\log q - \log p]$ . Here, since we are dealing with particularly rugged energy landscape where the gradient of p can be numerically overwhelming, we adopt the trick from [40] and optimize the REINFORCE [46] policy gradient surrogate instead:

$$\mathcal{L}_{\text{Mint}} = -\mathbb{E}_q \log q [\overline{\text{SoftMax}}(\log p - \log q)] \tag{17}$$

This objective fills the energy landscape with elliptical probability masses. The mode-seeking behavior of the reverse KL-divergence is not problematic here as the multimodal nature of p can be

captured by explicit discrete mixtures. We can also add an additional Stein VI [47]-style repulsion kernel among the gram matrix of the mixture components to encourage the diversification of modes.

219

220 221

222

223

224

225

226 227

228

229 230

231

232

233

234

235 236

237

238

239

240

241 242

243

244

245 246

247

249 250

251

253 254

255 256

257

258

259

260 261

262

263

264

265

266

267

268

269

$$\mathcal{L}_{Repulsion}$$

 $\mathcal{L}_{\text{Repulsion}} = \sum_{i} \sum_{j \neq i} \cos(\mu_i \mu_i^T, \mu_j \mu_j^T)$ 

As such, we can view Mint as the *optimization* stage of neural sampler training, quickly and cheaply finding diverse minima on the energy landscape. In § 4, we see that Mint alone can achieve satisfactory results in terms of mode finding. Of course, despite of Theorem 2.6, in the finite limit of the number of mixtures, the elliptical density q cannot fill arbitrarily sophisticated shapes. This motivates the design of a highly expressive kernel in the following section.

## JULEP KERNEL: JACOBIAN-UNPRESERVING LANGEVIN WITH EXPLICIT **PROJECTION**

We design a novel forward kernel  $k^+$  to replace that in Remark 1.1, together with its corresponding backward kernel  $k^-$ :

$$k_{\pm}(X_{t\pm\Delta t}|X_{t}) = \mathcal{N}\Big(\exp\big(\pm A_{t}(X_{t})\Delta t\big)X_{t} \pm b_{t}(X_{t})\Delta t + \partial U_{t}/\partial X\Delta t, 2\epsilon\Delta t\Big),\tag{19}$$

where exp denotes matrix exponential. One can easily see that this is but a different discretization of Equation 6, now written as  $d\vec{X} = -\epsilon \nabla U_t dt + \sqrt{2\epsilon} dB + (b_t(X_t) + \exp W_t(A_t)) dt$ , with the last term omitted into  $b_t$  by considering the first-order Taylor expansion of the matrix exponential. Since the handling of the Brownian motion is unchanged, the convergence criteria is no difference than that of the vanilla Euler-Maruyama method. Intuitively, our method affords the traditionally additive kernel a multiplicative structure, thus greatly enhancing the expressiveness of each step, allowing it to bridge faraway distributions. We can see that the constraint in Remark 1.1 no longer holds.

Following [23, 27, 48, 49], the path weight (Equation 4) can be discretized as:

$$W \approx U_0(X_0) - U_1(X_1) + \sum_{t=0}^{1} \log k^+(X_t|X_{t-\Delta t}) - \sum_{t=0}^{1} \log k^-(X_{t-\Delta t}|X_t).$$
 (20)

This relies on the assumption of the Gaussianality [33] of the reverse kernel, which is exact when  $\Delta t \to 0$ . Plugging this into the log-variance objective (Equations 9, 10) and omitting the constant  $Z_0 - Z_1$ , we arrive at the loss for training the Julep kernel:

$$\mathcal{L}_{\text{Julep}} = \mathbb{V}[W] = \mathbb{V}[U_0(X_0) - U_1(X_1) + \sum \log(k^+/k^-)], \tag{21}$$

where the variance is evaluated over  $q_0$  and the law of the forward SDE.

**Flexible neural parametrization.** We stress that any arbitrary parametrization of A(X,t), b(X,t)are all fair game, and the parametrization of U(X,t) also does not break the mathematical framework as long as  $U_0$  and  $U_1$  stay invariant, which can be easily parametrized as  $U_t = (1-t)U_0 +$  $tU_1 + t(1-t)U_t$ , with a free-form  $U_t$ , which is significantly more flexible than pre-defined linear mixing schedule [36]. When it comes to the noise schedule, although it is possible to prescribe a state-heteroschedastic noise  $\epsilon(X,t)$ , doing so would require a divergence correction term  $\Delta_X \epsilon$  for both the forward and backward SDE. We therefore only optimize  $\epsilon$  as a function of t.

**Preserving the orthogonal symmetry.** With the amount of care taken to design an orthogonally invariant prior, we cannot afford to lose the O(n)-symmetry in the integration stage. Fortunately, this is almost trivial thanks to the rich literature about designing E(n)-equivariant force fields and generative models [3-10]—in a sense, we are merely building a time-dependent version of these models. Consider such a model  $f_{\theta}: \mathcal{X} \times \mathcal{H} \to \mathcal{X} \times \mathcal{H}$  that map from and to the joint spaces of (ndimensional) geometry  $\mathcal{X} \in \mathbb{R}^n$  and semantic embedding  $\mathcal{H} \in \mathbb{R}^C$  such that it is permutationally, rotationally, translationally, and reflectionally equivariant on  $\mathcal{X}$  and invariant on  $\mathcal{H}$ , i.e.,  $\mathbf{x} \in \mathcal{X}, h \in$  $\mathcal{H}$  and  $T: \mathcal{X} \to \mathcal{X}$  is rotation, translation, and reflection, we have:

$$\mathbf{x}_f, h_f = f_\theta(\mathbf{x}, h) \iff T(\mathbf{x}_f), h_f = f_\theta(T(\mathbf{x}), h).$$
 (22)

We can make this model O(n)-equivariant (and time-dependent) by constructing an embedding combining the radial component of X, the time representation: h = [t:||X||]. The output of this model is connected to Equation 19 as:

$$b_t = \mathbf{x}_f; U_t = \sum h_{f_0}; A = h_{f_1} h_{f_1}^{\top}, \tag{23}$$

where the control term reuses the equivariant output directly; the potential term aggregates invariant embeddings among the particles (ubiquitous in force field constructions); and the projection term are parametrized using low-rank form, where  $h_{f_{0,1}}$  are channels of the invariant output  $h_f$ . Note that, the resulting projection term A is on the space of  $\mathbb{R}^{N\times N}$ , similar to that in graph diffusion [42], so we only linearly combine the positions of the particles without introducing internal rotation, thus easily preserving orthogonal symmetry. Although most of these architectures can be reduced to linear complexity, they require pre-specified graph structure (edge connection), which is not possible here. So this backbone also incurs a  $\mathcal{O}(N^2)$  runtime complexity. In this paper, we used the simplest equivariant graph neural networks (EGNNs) [3] as the backbone  $f_{\theta}$ , and leave more sophisticated architecture for future studies.

#### 4 MINTJULEP RESULTS: SEPARATING OPTIMIZATION AND SAMPLING.

Having defined the two components, we now put them in a coherent framework and provide a straightforward recipe to optimize them sequentially:

```
Algorithm 1: MintJulep training.
Input: Energy function U.
Input: Randomly initialized Mint prior q and Julep kernel k^+.
Input: Hyperparameters: Integration steps T and sample size S.
Output: Samples from the Boltzmann distribution (Equation 1)
 ......
while \mathcal{L}_{\text{Mint}}(\cdot;q) not converging do
   for i \sim \{1, ..., S-1\} do
      sample X_i \sim q (Equation 15;
   end
   descent \mathcal{L}_{\mathrm{Mint}}(X;q) (Equation 17) to optimize q
                  while \mathcal{L}_{\mathrm{Julep}}(\cdot;k^+) not converging do
   for i \sim \{1, ..., S-1\} do
       sample X \sim q (Equation 15;
       sample t_i \sim U(0,1), i = 1, \dots, T and sort;
       for i \in \{1, ..., T-1\} do
          sample X_{t_{i+1}} \sim k^+(\cdot|X_{t_i}) (Equation 19)
       end
   end
   descent \mathcal{L}_{\text{Julep}}(\overline{\{X_t\}}) (Equation 21) to optimize k^+
```

Note that two stages are required in the training of the model. In the first stage, the Mint prior q is optimized to descend the surrogate KL divergence between q and  $p \propto \exp(-U)$ , so that it becomes already close to p. Next, the details of the distribution, which cannot be filled by elliptical probability mass, are refined in the second stage using the Julep kernel. Although it is also possible to optimize the parameter using the log-variance loss (Equation 21), we found that doing so harms the stability of the training process. Empirically, the training of the prior, due to the lack of SDE integration, only takes seconds on a GPU to converge.

Next, we test the performance of this formulation using synthesized and real-world energy land-scape. Again, it is worth emphasizing that we only have access to the target energy function, or probability density up to a constant, and do not have access to samples, which explains the seemingly slightly inferior numerical performance to methods which do require such access [50]—these are two different settings that are not comparable.

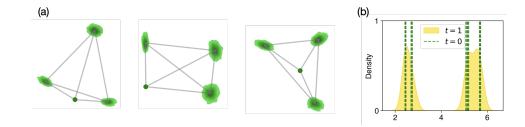


Figure 1: **Mint minimizes; Julep samples**—Illustration of the DW4 experiment. (a) Modes discovered by Mint on the 2-dimensional space. The dot represents the *anchor atom*. Kernel density estimation (KDE) plots from samples taken from the prior. The probability masses, though constrained to be elliptical in shape, are already placed at the minima of the energy surface. (b) KDE plots of the distances among particles from the posterior t=1, with modes of the prior t=0 marked by verticle lines.

	DW-4	LJ-13	LJ-55
PIS [24]	$46.2 \pm 8.1$	$1.2 \pm 1.1$	$0.1 \pm 0.0$
DDS [33]	$46.1 \pm 7.6$	$1.0\pm1.1$	$0.1 \pm 0.0$
MintJulep	$83.7 \pm 1.0$	$10.6 \pm 1.4$	$1.0 \pm 0.6$

Table 1: **MintJulep efficiently samples toy energy functions.** Effective sample size (ESS, %) normalized by the total sample size compared with state-of-the-art path-based models.

tions defined pairwise  $r_{ij} = ||\mathbf{x}_i - \mathbf{x}_j||$  among particles. The double wall (DW):

$$U_{\rm DW}(r) = \frac{1}{\tau} [\lambda_2 (r - r_0)^2 + \lambda_4 (r - r_0)^4], \tag{24}$$

with  $\lambda_2 = -4, \lambda_4 = 0.9, \tau = 1$ , and Leonard-Jones [51] potential:

$$U_{\rm LJ}(r) = \frac{\epsilon}{\tau} [(r/\sigma)^{12} - (r/\sigma)^6],$$
 (25)

with  $\sigma=1, \tau=1$ , and an additional harmonic potential constraining particles to the center-of-mass added [50] to prevent the dissolution of the system. Evidently, the LJ potential increases rapidly when  $r\to 0$ , which physically represents the strong repulsion among particles when they are about to collide, which contributes significantly to the ruggedness of the energy landscape. This poses a significant challenge for numerical optimization—when a linear path is used, the gradient soon causes overflow because of the 12-th power term. We therefore adopt a smooth annealing path:

$$\tilde{r} = r + \sigma(1-t); \tilde{U}_{LJ} = \epsilon/\tau [(\tilde{r}/\sigma)^{-6} - (2-t)^{-6}]^2,$$
 (26)

which preserves the minima but slowly anneal the minimal distance among particles from  $\sigma$  to 0 as  $t:0\to 1$ . We reuse this annealing path in the real-world experiment as well.

We compare with the path integral sampler (PIS) [24], which propagates the gradient across the SDE, and the denoising diffusion sampler (DDS) [25], which proposes the log-variance objective. Both of these methods use an isotropic Gaussian distribution and an additive kernel, which might explain the drastic difference in the normalized effective sample size (ESS):

$$ESS = \frac{1}{n} \mathbb{E}^{-1}[W] \approx \frac{1}{n} \frac{(\sum W_i)^2}{\sum W_i^2} = \frac{1}{n} \sum (SoftMax^2(W_i))^{-1}.$$
 (27)

In Figure 1, we show the compartimentalization and collaboration of the two parts of the model, with minima discovered firstly by Mint and refined by Julep. This trend is repeated in Figure 2 as well.

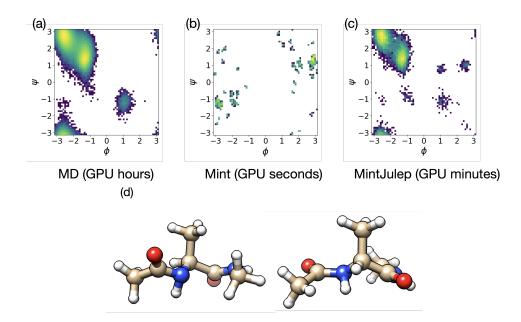


Figure 2: MintJulep can be used to produce convergent molecular dynamics (MD) simulation trajectories. Ramachandran plot (KDE plot of the dihedral angles of the molecule) of: (a) Reference equilibrium MD trajecotries of alanine dipeptide; (b) Samples generated from Mint, which captures the minima; (c) Samples generated from MintJulep, which contains finer detail. (d) Representative samples from MintJulep. Note that both chirality is possible since it is not specified in the energy function. The path ESS for MintJulep is  $18.0 \pm 1.2\%$ .

**Real-world energy landscape: alanine dipeptide.** Having established the satisfactory performance on synthetic sandboxes, we move on to test if our model can achieve real world utility by accelerating molecular dynamics (MD) simulation of biomolecular systems. In such case, the energy function comes from a molecular mechanics (MM) force field, typically expressed as:

$$\begin{split} U_{\text{MM}}(\mathbf{x}; \Phi_{\text{FF}}) &= \sum_{\text{bond}} &\frac{K_r}{2} (r_{ij} - r_0)^2 \\ &+ \sum_{\text{angle}} &\frac{K_{\theta}}{2} (\theta_{ijk} - \theta_0)^2 \\ &+ \sum_{\text{torsion}} &\sum_{n=1}^{n_{\text{max}}} K_{\phi,n} \left[ 1 + \cos(n\phi_{ijkl} - \phi_0) \right] \\ &+ \sum_{\text{Coulomb}} &\frac{1}{4\pi\epsilon_0} \frac{q_i \, q_j}{r_{ij}} \\ &+ \sum_{\text{LJ}} &4\epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right], \end{split}$$

where the total potential energy  $U_{\rm MM}$  as a function of the coordinates of the system  ${\bf x}$  and the collection of force field parameters  $\Phi_{\rm FF}=\{K_r,K_\theta,r_0,\theta_0,K_{\phi,n},\phi_0,q,\sigma,\epsilon\}_i$  is modeled as the sum of bond, angle, torsion, and nonbonded energy (For a machine learning community-friendly explanation, see [52]). In our setting, we adopt the topology from an alanine dipeptide and uses the collection of parameters from [53].

Reusing the annealing path (Equation 26), we notice a similar trend as the toy experiment—Mint captures the location of the minima quickly while Julep completes the fine detail of the energy landscape. It is worth noting that, since the chirality, which is an important trait of the biomolecules, is not encoded in the energy function, the model may generate samples that has different chirality than those abundant in nature (Figure 2 (d)), which might also explain the additional minima on the Ramachandran plot. In sum, a very brief Mint training can already capture the minima of the

energy landscape of alanine dipeptide, whereas MintJulep can recover most of the regions sampled by GPU-days-worth of MD trajectory in less than an hour of training.

#### 5 CONCLUSION

If we were able to sample the Boltzmann distribution associated with various physical systems efficiently and accurately, we would be able to build a more reliable bridge between the microscopic and the macroscopic, with which we can gain a deeper quantitative understanding of such systems, thereby rationally designing better pharmaceuticals, materials, and other physicochemical entities with microscopic structure and macroscopic functions. The approach presented here, MintJulep, represents a meaningful step towards this goal.

Concretely, the Boltzmann distributions associated with physical systems can oftentimes be described as rugged, i.e. with isolated minima. Another feature of such functions is that they are almost always O(n)—invariant. Starting from these two features of the realistic systems, as well as the failure mode of traditional neural samplers, we design a brand new class of distributions (Mint) and a powerful discretized kernel associated with the non-equilibrium annealing dynamics. These two practical improvements drastically increase the performance of the path-based neural samplers, allowing us to rapidly generate samples from the Boltzmann distributions associated with real systems given only the energy function. Informally, the Mint prior can be viewed as a minimization step (albeit still preserving the entropy structure), respecting the multimodality of the energy land-scape with the mixture of component design. As such, the prior is close to the desired target energy function by KL divergence, leaving the training of the already powerful Julep kernel a breeze.

**Limitations.** As discussed in §2, 3, both Mint and Julep incur  $\mathcal{O}(N^2)$  runtime complexity, and need further speed up before they can be efficiently used on realistic protein systems containing thousands of atoms. Furthermore, the reduction from the E(n) to O(n) group requires a careful choice of *anchor atom*, and the representation power of the internal coordinate system is sensitive to such choices.

**Future directions.** We plan to investigate further methods to simplify and accelerate the optimization of the Mint prior and the integration of the Julep kernel. This would allow us to model larger protein systems, and unify the pipelines of docking, sampling, and folding within one method. In addition, to make this model generalizable, in the style of [32], is a natural next step.

**Ethics statement.** We acknowledge and adhere to the Ethnics statement of the ICLR.

**Reproducibility statement.** The implementation of our method can be found at https://anonymous.4open.science/r/mint\_julep-22D7/, with core dependencies including JAX [54] and its eco-system. For hyperparameters, we use a 3-layer EGNN [3] with 64-units each and TanH activation everywhere in this paper. We also fix the integration steps to be 100 and the sample size to be 100, except when evaluating ESS for reporting, for which we use 1e5. The random seed is fixed as 2666 everywhere in this paper. Interestingly, our method only requires an energy function and is therefore a *data-free* method, requiring no datasets or data-processing.

#### REFERENCES

- [1] Mark E Tuckerman. *Statistical mechanics: theory and molecular simulation*. Oxford university press, 2023.
- [2] Ludwig Boltzmann. Studien uber das gleichgewicht der lebenden kraft. Wissenschafiliche Abhandlungen, 1:49–96, 1868.
- [3] Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E (n) equivariant graph neural networks. In *International conference on machine learning*, pages 9323–9332. PMLR, 2021.
- [4] Simon Batzner, Albert Musaelian, Lixin Sun, Mario Geiger, Jonathan P Mailoa, Mordechai Kornbluth, Nicola Molinari, Tess E Smidt, and Boris Kozinsky. E (3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nature communications*, 13(1): 2453, 2022.

[5] Ilyes Batatia, David P Kovacs, Gregor Simm, Christoph Ortner, and Gábor Csányi. Mace: Higher order equivariant message passing neural networks for fast and accurate force fields. *Advances in neural information processing systems*, 35:11423–11436, 2022.

- [6] Yuanqing Wang, Kenichiro Takaba, Michael S Chen, Marcus Wieder, Yuzhi Xu, Tong Zhu, John ZH Zhang, Arnav Nagle, Kuang Yu, Xinyan Wang, et al. On the design space between molecular mechanics and machine learning force fields. *Applied Physics Reviews*, 12(2), 2025.
- [7] Yuanqing Wang, Josh Fass, Benjamin Kaminow, John E Herr, Dominic Rufa, Ivy Zhang, Iván Pulido, Mike Henry, Hannah E Bruce Macdonald, Kenichiro Takaba, et al. End-to-end differentiable construction of molecular mechanics force fields. *Chemical Science*, 13(41):12016–12033, 2022.
- [8] Yuanqing Wang and John D Chodera. Spatial attention kinetic networks with e(n)-equivariance. In *ICLR* 2023, 2023.
- [9] Kenichiro Takaba, Anika J Friedman, Chapin E Cavender, Pavan Kumar Behara, Iván Pulido, Michael M Henry, Hugo MacDermott-Opeskin, Christopher R Iacovella, Arnav M Nagle, Alexander Matthew Payne, John D Chodera, and Yuanqing Wang. Machine-learned molecular mechanics force fields from large-scale quantum chemical data. *Chemical Science*, 2024.
- [10] Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Sauceda Felix, Stefan Chmiela, Alexandre Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *Advances in neural information processing systems*, 30, 2017.
- [11] Elman Mansimov, Omar Mahmood, Seokho Kang, and Kyunghyun Cho. Molecular geometry prediction using a deep generative graph neural network. *Scientific Reports*, 9(1), December 2019. ISSN 2045-2322. doi: 10.1038/s41598-019-56773-5. URL http://dx.doi.org/10.1038/s41598-019-56773-5.
- [12] Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock: Diffusion steps, twists, and turns for molecular docking, 2023. URL https://arxiv.org/abs/2210.01776.
- [13] Yuzhi Xu, Wei Xia, Chao Zhang, Xinxin Liu, Chengwei Ju, Xuhang Dai, Pujun Xie, Yuan-qing Wang, Guangyong Chen, and John Zhang. Quantifying protein-protein interaction with a spatial attention kinetic graph neural network. *bioRxiv*, pages 2025–06, 2025.
- [14] Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, 2024.
- [15] Gustaf Ahdritz, Nazim Bouatta, Christina Floristean, Sachin Kadyan, Qinghui Xia, William Gerecke, Timothy J O'Donnell, Daniel Berenberg, Ian Fisk, Niccolò Zanichelli, et al. Openfold: Retraining alphafold2 yields new insights into its learning mechanisms and capacity for generalization. *Nature methods*, 21(8):1514–1524, 2024.
- [16] Chai Discovery team, Jacques Boitreaud, Jack Dent, Matthew McPartlon, Joshua Meier, Vinicius Reis, Alex Rogozhonikov, and Kevin Wu. Chai-1: Decoding the molecular interactions of life. *BioRxiv*, pages 2024–10, 2024.
- [17] Michael Levitt and Shneior Lifson. Refinement of protein conformations using a macromolecular energy minimization procedure. *Journal of molecular biology*, 46(2):269–279, 1969.
- [18] AT Hagler, E Huler, and Shneior Lifson. Energy functions for peptides and proteins. i. derivation of a consistent force field including the hydrogen bond from amide crystals. *Journal of the American Chemical Society*, 96(17):5319–5327, 1974.
- [19] J Andrew McCammon, Bruce R Gelin, and Martin Karplus. Dynamics of folded proteins. *nature*, 267(5612):585–590, 1977.
- [20] Radford M. Neal. Annealed importance sampling, 1998. URL https://arxiv.org/abs/physics/9803008.
- [21] Christopher Jarzynski. Nonequilibrium equality for free energy differences. *Physical Review Letters*, 78(14):2690, 1997.
- [22] Suriyanarayanan Vaikuntanathan and Christopher Jarzynski. Escorted free energy simulations. *The Journal of Chemical Physics*, 134(5), February 2011. ISSN 1089-7690. doi: 10.1063/1. 3544679. URL http://dx.doi.org/10.1063/1.3544679.

[23] Michael S Albergo and Eric Vanden-Eijnden. Nets: A non-equilibrium transport sampler. arXiv preprint arXiv:2410.02711, 2024.

- [24] Qinsheng Zhang and Yongxin Chen. Path integral sampler: a stochastic control approach for sampling, 2022. URL https://arxiv.org/abs/2111.15141.
- [25] Francisco Vargas, Will Grathwohl, and Arnaud Doucet. Denoising diffusion samplers, 2023. URL https://arxiv.org/abs/2302.13834.
- [26] Aaron Havens, Benjamin Kurt Miller, Bing Yan, Carles Domingo-Enrich, Anuroop Sriram, Brandon Wood, Daniel Levine, Bin Hu, Brandon Amos, Brian Karrer, Xiang Fu, Guan-Horng Liu, and Ricky T. Q. Chen. Adjoint sampling: Highly scalable diffusion samplers via adjoint matching, 2025. URL https://arxiv.org/abs/2504.11713.
- [27] Hao Wu, Jonas Köhler, and Frank Noé. Stochastic normalizing flows. *Advances in neural information processing systems*, 33:5933–5944, 2020.
- [28] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. *CoRR*, abs/1806.07366, 2018. URL http://arxiv.org/abs/1806.07366.
- [29] Avishek Joey Bose, Tara Akhound-Sadegh, Guillaume Huguet, Kilian Fatras, Jarrid Rector-Brooks, Cheng-Hao Liu, Andrei Cristian Nica, Maksym Korablyov, Michael Bronstein, and Alexander Tong. Se (3)-stochastic flow matching for protein backbone generation. *arXiv* preprint arXiv:2310.02391, 2023.
- [30] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling, 2023. URL https://arxiv.org/abs/2210.02747.
- [31] Frank Noé, Simon Olsson, Jonas Köhler, and Hao Wu. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science*, 365(6457):eaaw1147, 2019.
- [32] Leon Klein and Frank Noé. Transferable boltzmann generators, 2025. URL https://arxiv.org/abs/2406.14426.
- [33] Francisco Vargas, Shreyas Padhy, Denis Blessing, and Nikolas Nüsken. Transport meets variational inference: Controlled monte carlo diffusions, 2025. URL https://arxiv.org/abs/2307.01050.
- [34] Gavin E Crooks. Entropy production fluctuation theorem and the nonequilibrium work relation for free energy differences. *Physical Review E*, 60(3):2721, 1999.
- [35] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2209.15571*, 2022.
- [36] Bálint Máté and François Fleuret. Learning interpolations between boltzmann densities. *arXiv* preprint arXiv:2301.07388, 2023.
- [37] Kirill Neklyudov, Rob Brekelmans, Daniel Severo, and Alireza Makhzani. Action matching: Learning stochastic dynamics from samples. In *International conference on machine learning*, pages 25858–25889. PMLR, 2023.
- [38] Lorenz Richter and Julius Berner. Improved sampling via learned diffusions, 2024. URL https://arxiv.org/abs/2307.01198.
- [39] Edward Nelson. *Dynamical theories of Brownian motion*, volume 3. Princeton university press, 1967.
- [40] Laurence Illing Midgley, Vincent Stimper, Gregor NC Simm, Bernhard Schölkopf, and José Miguel Hernández-Lobato. Flow annealed importance sampling bootstrap. *arXiv* preprint *arXiv*:2208.01893, 2022.
- [41] Changyi Xiao and Ligang Liu. Generative flows with matrix exponential. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 10452–10461. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/xiao20a.html.
- [42] Benjamin Paul Chamberlain, James Rowbottom, Maria Gorinova, Stefan Webb, Emanuele Rossi, and Michael M. Bronstein. Grand: Graph neural diffusion, 2021. URL https://arxiv.org/abs/2106.10934.

[43] Kanti V Mardia and Peter E Jupp. Directional Statistics. Wiley, 2000.

- [44] Soledad Villar, David W. Hogg, Kate Storey-Fisher, Weichi Yao, and Ben Blum-Smith. Scalars are universal: Gauge-equivariant machine learning, structured like classical physics. *CoRR*, abs/2106.06610, 2021. URL https://arxiv.org/abs/2106.06610.
- [45] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org.
- [46] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- [47] Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm, 2019. URL https://arxiv.org/abs/1608.04471.
- [48] Jiajun He, Yuanqi Du, Francisco Vargas, Yuanqing Wang, Carla P. Gomes, José Miguel Hernández-Lobato, and Eric Vanden-Eijnden. Feat: Free energy estimators with adaptive transport, 2025. URL https://arxiv.org/abs/2504.11516.
- [49] Jerome P. Nilmeier, Gavin E. Crooks, David D. L. Minh, and John D. Chodera. Nonequilibrium candidate monte carlo is an efficient tool for equilibrium simulation. *Proceedings of the National Academy of Sciences*, 108(45), October 2011. ISSN 1091-6490. doi: 10.1073/pnas. 1106094108. URL http://dx.doi.org/10.1073/pnas.1106094108.
- [50] Leon Klein, Andreas Krämer, and Frank Noé. Equivariant flow matching. *Advances in Neural Information Processing Systems*, 36:59886–59910, 2023.
- [51] John Edward Jones. On the determination of molecular fields.—i. from the variation of the viscosity of a gas with temperature. *Proceedings of the Royal Society of London. Series A, containing papers of a mathematical and physical character*, 106(738):441–462, 1924.
- [52] Yuanqing Wang, Kenichiro Takaba, Michael S. Chen, Marcus Wieder, Yuzhi Xu, John Z. H. Zhang, Kuang Yu, Xinyan Wang, Linfeng Zhang, Daniel J. Cole, Joshua A. Rackers, Joe G. Greener, Peter Eastman, Stefano Martiniani, and Mark E. Tuckerman. On the design space between molecular mechanics and machine learning force fields. *Applied Physics Review*, 2025. URL https://doi.org/10.1063/5.0237876.
- [53] David A Case, Thomas E Cheatham III, Tom Darden, Holger Gohlke, Ray Luo, Kenneth M Merz Jr, Alexey Onufriev, Carlos Simmerling, Bing Wang, and Robert J Woods. The amber biomolecular simulation programs. *Journal of computational chemistry*, 26(16):1668–1688, 2005.
- [54] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http://github.com/jax-ml/jax.

### **Proof of Theorem 2.6**

*Proof.* Suppose we have a probability density function p that is orthogonally invariant according to Definition 2.1. It can be written as:

$$p(X) = \int dY \delta(X - Y), \tag{28}$$

which, since p is piecewise continuous, can be approximated arbitrarily well by a Riemann sum:

$$p(X) = \frac{1}{k} \sum p_i(X|\xi_i), \tag{29}$$

where  $\xi_i$  is a region in which  $p_i$  stays constant. Due to the first fundamental theorem of the orthogonal group,  $\xi$  can be embedded in any coordinate system up to the orthogonal transformation. As such, the mixture component

$$p_i \text{PvMFLogNormal}(\cdot, \mu = ||\Xi||, \kappa \to \text{inf}, \rho = ||\Xi||, \sigma \to 0),$$
 (30)

where  $\Xi Q \in \xi, \forall QQ^{\top} = I$ , can approximate any region  $\xi$  arbitrarily well.