
Transferable Boltzmann Generators

Anonymous Author(s)

Affiliation

Address

email

Abstract

The generation of equilibrium samples of molecular systems has been a long-standing problem in statistical physics. Boltzmann Generators are a generative machine learning method that addresses this issue by learning a transformation via a normalizing flow from a simple prior distribution to the target Boltzmann distribution of interest. Recently, flow matching has been employed to train Boltzmann Generators for small molecular systems in Cartesian coordinates. We extend this work and propose a first framework for Boltzmann Generators that are transferable across chemical space, such that they predict zero-shot Boltzmann distributions for test molecules without being retraining for these systems. These transferable Boltzmann Generators allow approximate sampling from the target distribution of unseen systems, as well as efficient reweighting to the target Boltzmann distribution. The transferability of the proposed framework is evaluated on dipeptides, where we show that it generalizes efficiently to unseen systems. Furthermore, we demonstrate that our proposed architecture enhances the efficiency of Boltzmann Generators trained on single molecular systems.

1 Introduction

Generative models have demonstrated remarkable success in the physical sciences, including protein structure prediction [1, 2, 3], generation of de novo molecules [4, 5, 6], and efficiently generating samples from the Boltzmann distribution [7, 8, 9]. In this work, we will focus on the latter for molecular systems, which represents a promising avenue for addressing the sampling problem. The sampling problem refers to the long-standing challenge in statistical physics to generate samples from equilibrium Boltzmann distributions $\mu(x) \propto \exp(-U(x)/k_B T)$, where $U(x)$ is the potential energy of the system, k_B the Boltzmann constant, and T the temperature. Traditionally, samples are generated with sequential sampling algorithms such as Markov Chain Monte Carlo and Molecular Dynamics (MD) simulations. However, these algorithms require a significant amount of time to generate uncorrelated samples from the target distribution. This is due to the necessity of performing small update steps, in the order of femtoseconds, for stability. This is especially challenging in the presence of well-separated metastable states, where transitions are unlikely due to high energy barriers. In recent years, numerous machine learning methods have emerged to address this challenge [7, 10]. One such method is the Boltzmann Generators (BG) [7]. In this work, we refer to BGs as a model that allows for the approximate sampling of the Boltzmann distribution of interest and the subsequent reweighting to the unbiased target distribution. If the model is only capable of generating approximate samples, which may stem from a subset of the Boltzmann distribution, we refer to them as Boltzmann Emulators¹. Boltzmann Generators transform an often simple, prior distribution via a normalizing flow [11, 12, 13, 14] to an approximation of the target Boltzmann distribution. Subsequently, generated samples can be reweighted to the unbiased target distribution. The effectiveness of the reweighting depends on how close the generated distribution matches the

¹To the best of our knowledge Bowen Jing introduced the name first.

target distribution. Hence, it is possible to generate uncorrelated and unbiased samples from the target Boltzmann distribution, potentially generating significant speed-up over classical MD simulations. There are many ways to build a BG, because of the various available realizations of normalizing flows. In this work, we will focus on continuous normalizing flows (CNFs) [15, 16], rather than coupling flows [17]. Recently, flow matching [18, 19, 20, 21] emerged as an alternative training method for CNFs, which is simulation free, allowing for more efficient training of CNFs.

Thus far, Boltzmann Generators have been found to be limited by the necessity of training them on the system of interest. This training process, which requires a significant amount of time, makes it challenging to achieve any significant speed-up over classical MD simulations. Furthermore, the training time must be taken into account, which may even necessitate the execution of MD simulations of the system of interest on its own. It is therefore desirable to have a transferable Boltzmann Generator that can be trained on one set of molecules and generalize to another set, where Boltzmann samples can be efficiently generated at inference time without retraining. Only recently, reliable Boltzmann Generators in Cartesian coordinates for molecules were introduced [22, 23], which paved the way for transferable Boltzmann Generators, as they do not depend on the molecule specific internal coordinate representation, which make it difficult to construct transferable models.

In this work, we introduce a framework for *transferable* Boltzmann Generators based on CNFs, allowing effective sample generation from unseen Boltzmann distributions. Transferable Boltzmann Generators are desirable, as they do not require retraining for similar systems and can be trained on short training trajectories that miss metastable states. The Boltzmann Generator can still learn these from other similar trajectories of other systems.

We make the following main contributions:

1. We introduce, to the best of our knowledge, the first *transferable* Boltzmann Generator. We demonstrate the transferability on dipeptides, where we are able to generate unbiased samples from Boltzmann distributions of unseen dipeptides.
2. We describe a general framework for training and sampling with transferable Boltzmann Generators based on continuous normalizing flows. This includes also the post-processing of generated samples.
3. We perform several ablation studies to investigate the effect of different architectures, training set sizes, as well as biasing the training data. The results demonstrate that small training sets can be sufficient to train transferable Boltzmann Generators.

2 Related work

The initial work on Boltzmann Generators [7] has led to a great deal of subsequent research. The most common application of BGs is to generate samples from Boltzmann distributions of molecules [24, 25, 26, 27, 28, 29], as well as lattice systems [25, 30, 31]. Most BGs for molecular systems require system-specific featurizations such as internal coordinates [7, 32, 25, 26, 33, 34, 29]. Only recently, BGs for small molecular systems in Cartesian coordinates were introduced [22, 23], using CNFs and coupling flows, respectively. Equivariant normalizing flows [35, 36, 37, 27, 22, 38] played a pivotal role in the success of Boltzmann Generators in Cartesian coordinates, not only for molecular systems. The majority of BGs employ a Gaussian prior distribution, but it is also possible to start from prior distributions close to the target distribution [39, 34, 40], which makes the learning task simpler. However, all previous Boltzmann Generators are not transferable. Arguably, the work of [6] represents an exception, as they are able to generate samples from unseen conditional (Boltzmann) distributions in torsion space. However, the distribution is conditioned on a single local structure for each molecule, namely fixed bonds and angles. Consequently, in contrast to our work, they are unable to generate samples from the full Boltzmann distribution in Euclidean space. The first transferable deep generative model that was able to generate asymptotically unbiased samples from the Boltzmann distribution is [10]. Instead of generating independent samples, they learn large time steps and combine these with Metropolis-Hastings acceptance steps, to ensure asymptotically unbiased samples. However, in contrast to our work, they do not generate uncorrelated samples.

Boltzmann Emulators are analogous to Boltzmann Generators, yet they are not designed to generate unbiased equilibrium samples from the target Boltzmann distribution. Instead, they are intended to generate approximate samples that do not undergo reweighting. Furthermore, the generation of all metastable states may not be a necessary requirement, depending on the system. Boltzmann

Emulators do not need to be based on flow models, as they do not aim to do reweighing to the target distribution. They are often similar to Boltzmann Generators and use normalizing flows or diffusion models for the architecture, but due to removing the constraint of sampling the unbiased Boltzmann distribution, they can target significantly larger systems or are transferable. One example is [41], who propose a three stage transferable CNF model to learn peptide ensembles. [42] use flow matching to learn distributions of proteins, while [43] utilize diffusion models. [44] build a transferable Boltzmann Emulator for small molecules. Others aim to additionally also capture the correct dynamics of the molecular systems, such as [45], who use a diffusion model to predict transition probabilities. Scaling to larger systems often requires coarse graining [46, 47, 42], e.g. describing amino acids by a single bead rather than the individual atoms. However, this approach precludes the possibility of reweighing to the Boltzmann distribution.

A distinct, though related, learning objective is to generate novel molecular conformations. However, approximations from the Boltzmann distribution are not necessary; it is sufficient to generate a few (or even a single) conformation per molecule. The utilized architectures are once again analogous, as flow and diffusion models are employed [4, 5, 48, 49, 50, 6].

3 Boltzmann Generators and Normalizing Flows

Here, we describe Boltzmann Generators and normalizing flows, which are a central part of our proposed transferable Boltzmann Generator framework. We follow the notation of [22].

3.1 Boltzmann Generators

Boltzmann Generators (BGs) [7] combine an exact likelihood deep generative model and a reweighting algorithm to reweight the generated distribution to the target Boltzmann distribution. The exact likelihood deep generative model is trained to generate samples from a distribution $\tilde{p}(x)$ that is close to the target Boltzmann distribution $\mu(x)$. A common choice for the exact likelihood model are normalizing flows.

The Boltzmann Generator can be used to generate unbiased samples by first sampling $x \sim \tilde{p}(x)$ with the exact likelihood model and then computing corresponding importance weights $w(x) = \mu(x)/\tilde{p}(x)$ for each sample. These allow to reweight generated samples to the target Boltzmann distribution $\mu(x)$. It is possible to estimate observables of interest (asymptotically unbiased) using the weights $w(x)$ with importance sampling via

$$\langle O \rangle_\mu = \frac{\mathbb{E}_{x \sim \tilde{p}(x)}[w(x)O(x)]}{\mathbb{E}_{x \sim \tilde{p}(x)}[w(x)]}. \quad (1)$$

Furthermore, these reweighting weights can be employed to assess the efficiency of trained BGs by computing the effective sample size (ESS) with Kish’s equation [51]. In this work, we will compute the relative ESS, rather than the absolute one, and refer to it as ESS.

3.2 Continuous Normalizing Flows (CNFs)

Normalizing flows [14, 52] are a type of deep generative model used to learn complex probability densities $\mu(x)$ by transforming a prior distribution $q(x)$ through an invertible transformation $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^n$, resulting in the push-forward distribution $\tilde{p}(x)$.

Continuous Normalizing Flows (CNFs) [15, 16] are a specific kind of normalizing flow. In CNFs, the invertible transformation $f_\theta^t(x)$ is defined by the ordinary differential equation

$$\frac{df_\theta^t(x)}{dt} = v_\theta(t, f_\theta^t(x)), \quad f_\theta^0(x) = x_0, \quad (2)$$

where $v_\theta(t, x) : \mathbb{R}^n \times [0, 1] \rightarrow \mathbb{R}^n$ is a time-dependent vector field. The solution to this initial value problem provides the transformation equation

$$f_\theta^t(x) = x_0 + \int_0^t dt' v_\theta(t', f_\theta^{t'}(x)), \quad (3)$$

with $f_\theta^1(x) = \tilde{p}^t(x)$. The corresponding change in log density from the prior to the push-forward distribution is described by the continuous change of variable equation

$$\log \tilde{p}(x) = \log q(x) - \int_0^1 dt \nabla \cdot v_\theta(t, f_\theta^t(x)). \quad (4)$$

Equivariant flows The energy of molecular systems is typically invariant under permutations of interchangeable particles and global rotations and translations. Consequently, it is advantageous for the push-forward distribution of a Boltzmann generator to possess the same symmetries as the target system. In [35, 36] the authors demonstrate that the push-forward distribution $\tilde{p}(x)$ of a permutation and rotation equivariant normalizing flow with a permutation and rotation invariant prior distribution, is again rotation and permutation invariant. Furthermore, [35] present a method to construct such equivariant CNFs by using an equivariant vector field v_θ .

3.3 Flow matching

Flow matching [18, 19, 20, 21] enables efficient, simulation-free training of CNFs. The conditional flow matching training objective allows for the direct training of the vector field $v_\theta(t, x)$ through

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t \sim [0,1], x \sim p_t(x|z)} \|v_\theta(t, x) - u_t(x|z)\|_2^2. \quad (5)$$

There are many possible parametrizations for the conditional vector field $u_t(x|z)$ and the conditional probability path $p_t(x|z)$. One of the most simple, but powerful possible parametrization is

$$z = (x_0, x_1) \quad \text{and} \quad p(z) = q(x_0)\mu(x_1) \quad (6)$$

$$u_t(x|z) = x_1 - x_0 \quad \text{and} \quad p_t(x|z) = \mathcal{N}(x|t \cdot x_1 + (1-t) \cdot x_0, \sigma^2), \quad (7)$$

which we use in this work to train our models. For a more detailed description refer to [18, 21, 22, 48].

4 Transferable Boltzmann Generators

This section presents our proposed framework for transferable Boltzmann Generators (TBGs).

4.1 Architecture

Our proposed transferable Boltzmann Generator is based on a CNF. The corresponding vector field $v_\theta(t, x)$ is parametrized by an $O(D)$ - and $S(N)$ -equivariant graph neural network (EGNN) [37, 53, 41], as commonly used in prior work, e.g. [22, 37]. Although, less expressive than other equivariant networks such as [54, 55, 56, 57], it is faster to evaluate, which is important for CNFs as there can be hundreds of vector field calls during inference.

The vector field $v_\theta(t, x)$ consists of L consecutive layers. The position of the i -th particle x_i is updated according to the following equations:

$$h_i^0 = (t, a_i, b_i, c_i), \quad m_{ij}^l = \phi_e(h_i^l, h_j^l, d_{ij}^2), \quad (8)$$

$$x_i^{l+1} = x_i^l + \sum_{j \neq i} \frac{(x_i^l - x_j^l)}{d_{ij} + 1} \phi_d(m_{ij}^l), \quad (9)$$

$$h_i^{l+1} = \phi_h(h_i^l, m_i^l), \quad m_i^l = \sum_{j \neq i} \phi_m(m_{ij}^l) m_{ij}^l, \quad (10)$$

$$v_\theta(t, x^0)_i = x_i^L - x_i^0 - \frac{1}{N} \sum_j (x_i^L - x_i^0), \quad (11)$$

where ϕ_α represents different neural networks, d_{ij} is the Euclidean distance between particle i and j , t is the time, a_i is an embedding for the particle type, b_i for the amino acid, and c_i or the amino acid position in the peptide. In the final step, the geometric center is subtracted to ensure that the center of positions is conserved. When combined with a symmetric mean-free prior distribution, the push-forward distribution of the CNF will be $O(D)$ - and $S(N)$ -invariant, as demonstrated in [58].

162 The embedding of each atom is constructed from three parts. The first part is the atom type a_i , which
163 is a one-hot vector of 54 classes. The classes are defined by the atom type in the topology for a
164 classical force field. Therefore, only a few atoms are indistinguishable, such as hydrogen atoms that
165 are bound to the same carbon or nitrogen atom. The second part is the amino acid to which the atom
166 belongs, which is divided into 20 classes. The third part is the position of the amino acid in the peptide
167 sequence. This embedding is similar to the embedding used in [41] for the rotamer embeddings.
168 The amino acid and positional embeddings are only used for the transferable experiments. For more
169 details see Appendix B.5. In this study, we refer to this transferable Boltzmann Generator architecture
170 as *TBG + full*, and we use this name even when we apply it to a non-transferable setting.

171 The proposed architecture in [22] employs distinct encodings for all backbone atoms and the atom
172 types for the remainder. This represents a special case of our architecture, wherein b_i and c_i are
173 omitted and a_i is simply the atom types or the backbone encoding. We refer to this architecture
174 as *TBG + backbone*. Furthermore, we refer to the specific architecture employed in [22] as BG +
175 backbone for the alanine dipeptide experiments.

176 Moreover, we employ a model that utilizes the atom type as the sole encoding (there are only five
177 distinct atom types). This model is referred to as simply *TBG*.

178 4.2 Training transferable Boltzmann Generators

179 All transferable Boltzmann Generators are trained with flow matching. As there are different peptides
180 in each batch, the individual flow matching loss is divided by the number of atoms in each peptide.
181 For the 2AA dataset, all training peptides in each batch are used. All different architectures are
182 trained in the same way; for more details, see Appendix B.

183 4.3 Inference with transferable Boltzmann Generators

184 Sampling with a transferable Boltzmann Generator, especially on unseen peptides, poses multiple
185 challenges: (i) Some generated samples may not correspond to the molecule of interest, but rather to
186 a molecule that contains the same atom number and types but has a different bonding graph. This is
187 because the model has never encountered such a configuration during training. These configurations
188 might even have much lower quantum Chemical potential energies than the molecule of interest. For
189 some examples see Appendix A.4. However, as we are in this work interested in sampling from
190 the equilibrium Boltzmann distribution for a given molecular bonding graph, rather than sampling
191 distinct molecules, we would like to avoid these cases. Nevertheless, this effect can be largely
192 mitigated by our proposed TBG + full architecture. (ii) When working with classical force fields,
193 the correct ordering with respect to topology is crucial for evaluating energies. This is not a concern
194 for semi-empirical force fields, as they respect the permutation symmetry of particles of the same
195 type. As we typically use Gaussian prior distribution, it is common that the generated samples are
196 not arranged according to the topology. Consequently, in order to evaluate the energy, it is necessary
197 to reorder the generated samples according to the topology. (iii) It is possible that the chirality of
198 generated samples may differ from that of the peptide of interest. This can be rectified by simple
199 mirroring if all chirality centers of a peptide are flipped. Otherwise, these samples are assigned high
200 energies, as they are not from the target Boltzmann distribution of interest.

201 We resolve (i) and (ii) by generating a bond graph for the generated samples, based on empirical
202 bond distances and atom types. This graph is then compared with a reference bond graph. If the
203 two graphs are isomorphic, we can conclude that the configuration is correct. For more details, see
204 Appendix B.1. For (iii), we employ the code of [10] to check all chiral centers. It should be noted that
205 only if the generated samples possess the correct configuration and chirality will they be considered
206 valid samples from the Boltzmann distribution of interest.

207 5 Experiments

208 In this section, we compare our model with similar previous work [22] on equivariant Boltzmann
209 Generators for alanine dipeptide. Moreover, we show the transferability of our model on dipeptides.
210 More experimental details, such as dataset details, the specifics of the employed models, and the
211 utilized computing infrastructure can be found in Appendix B.

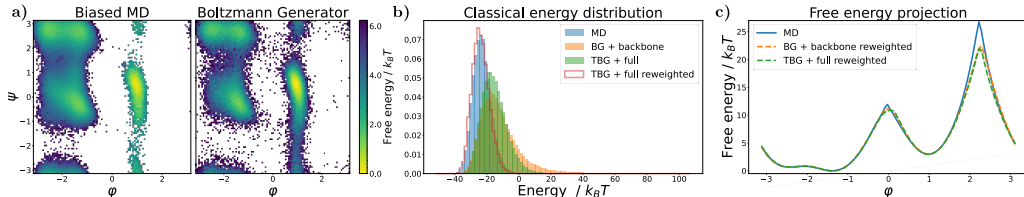


Figure 1: Results for the alanine dipeptide system simulated with a classical force field (a) Ramachandran plots for the biased MD distribution (left) and for samples generated with the TGB + full model (right). (b) Energies of samples generated with different methods. (c) Free energy projection along the slowest transition (φ angle), computed with different methods.

Table 1: Comparison of Boltzmann Generators with different architectures for the single molecular system alanine dipeptide. Errors are computed over five runs. The results for Boltzmann Generator and backbone encoding (BG + backbone) for the semi-empirical force field are taken from [22].

Model	NLL (\downarrow)	ESS (\uparrow)
Alanine dipeptide - semi-empirical force field		
BG + backbone [22]	-107.56 ± 0.09	$0.50 \pm 0.13\%$
TGB + full (ours)	-124.71 ± 0.08	$1.03 \pm 0.17\%$
Alanine dipeptide - classical force field		
BG + backbone [22]	-109.02 ± 0.01	$1.56 \pm 0.30\%$
TGB + full (ours)	-127.06 ± 0.12	$6.03 \pm 1.34\%$

212 5.1 Alanine dipeptide

213 In our first experiment, we investigate the single molecule alanine dipeptide in implicit solvent,
 214 described in Cartesian coordinates. The dataset was introduced in [22], for more details see Ap-
 215 pendix B.2. The training trajectory was generated by sampling with respect to a classical force field,
 216 and subsequently, 10^5 random samples were relaxed with respect to the semi-empirical *GFN2-xTB*
 217 force-field [59] for 100fs each. The objective is to train a Boltzmann Generator capable of sampling
 218 from the equilibrium Boltzmann distribution defined by the semi-empirical *GFN2-xTB* force-field
 219 efficiently and to recover the free energy surface along the slowest transition, i.e. the φ angle.
 220 Following the methodology outlined in [22], the training data is biased towards the less probable
 221 (positive) φ state. It is evident that any trained model on this set will be biased in comparison to
 222 the true Boltzmann distribution defined by the semi-empirical energy. However, the reweighting
 223 technique allows for the debiasing of the samples. The model is trained in the same way as described
 224 in [22]. Overall, the likelihoods and ESS values observed for the TGB + full model are superior to
 225 those reported in [22] (Table 1). This is achieved with nearly the same amount of parameters and
 226 maintaining comparable training and inference times (see Appendix B.3). Furthermore, the correct
 227 free energy difference is recovered, as demonstrated in Appendix A.1.

228 In [22] the authors used a semi-empirical potential to avoid the required ordering of the atoms to the
 229 topology for classical force fields. As the prior distribution of the Boltzmann Generator is usually
 230 a multivariate standard Gaussian distribution, generated samples will almost certainly not have the
 231 correct ordering. As we have introduced an efficient way to reorder samples in Section 4.3, we can
 232 now also evaluate alanine dipeptide for a classical force field. Therefore, we retrain the model in [22]
 233 on the classical MD trajectory and compare with our TGB + full architecture. We bias the training
 234 data as before towards the unlikely φ state. As expected, the likelihood and ESS for the classical
 235 force field are much better than for the semi-empirical one, as the training data stems from the target
 236 distribution. Our proposed architecture again performs significantly better, as shown in Section 5.1
 237 and Figure 1. The majority of generated samples with the TGB + full model and the BG + backbone
 238 sample nearly exclusively correct configurations, i.e. configurations with the correct bond graph,
 239 namely nearly 100% and about 98%, respectively. As presented in Figure 1, both models recover the
 240 free energy landscape correctly.

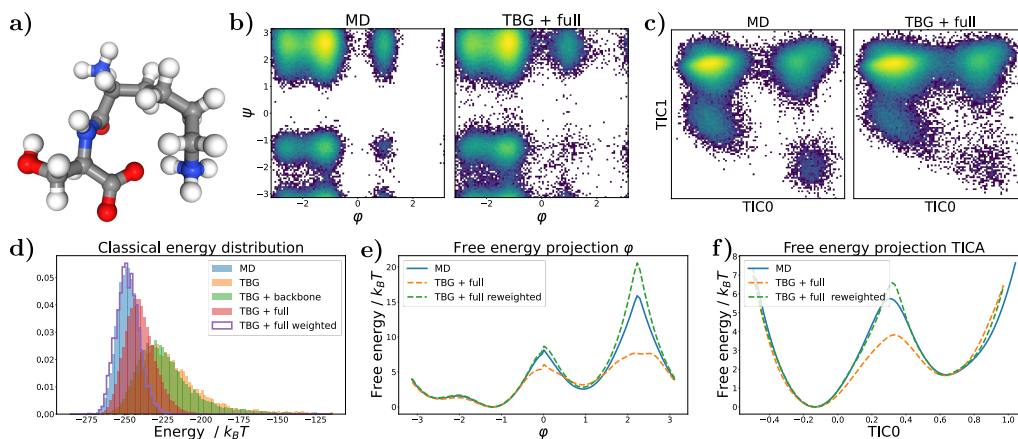


Figure 2: Results for the KS dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of samples generated with different methods and architectures. (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

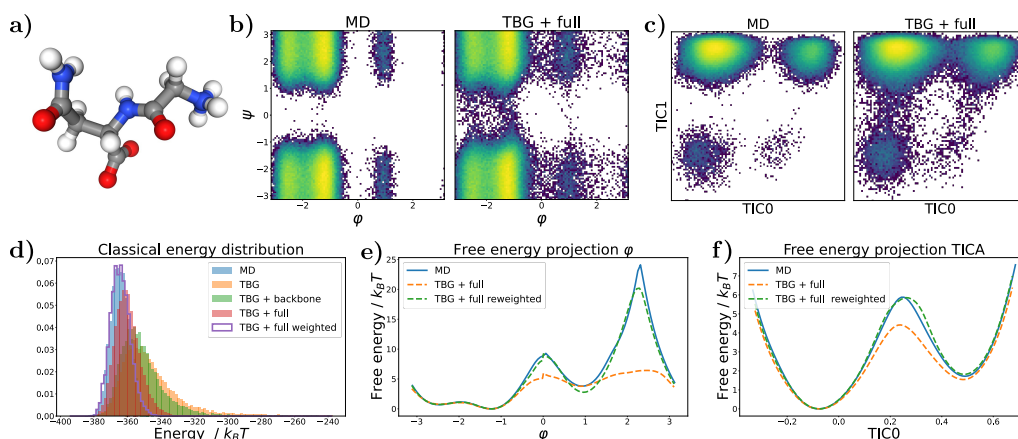


Figure 3: Results for the GN dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of samples generated with different methods and architectures. (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

241 5.2 Dipeptides (2AA)

242 In our second experiment, we evaluate our model on dipeptides and show transferability. The dataset
 243 was introduced in [10]. The training set consists of 200 dipeptides, which were simulated each with a
 244 classical force field for 50 ns and, therefore, may not have reached convergence. Nevertheless, as
 245 previously demonstrated, it is not necessary to train on unbiased data in order to obtain unbiased
 246 samples with a Boltzmann Generator.

247 We compare the three different transferable architectures described in Section 4.1 and use the same
 248 training procedure for all of them. Similar to the alanine dipeptide experiments, we obtain significantly
 249 better results for the TBG + full model in terms of ESS (Table 2 and Figure 4a), energies (Figure 2d),
 250 the ratio of correct configurations (Table 2), and likelihoods of test set samples (Appendix A.5). In
 251 particular, the extremely low number of correct configurations for numerous test peptides for the

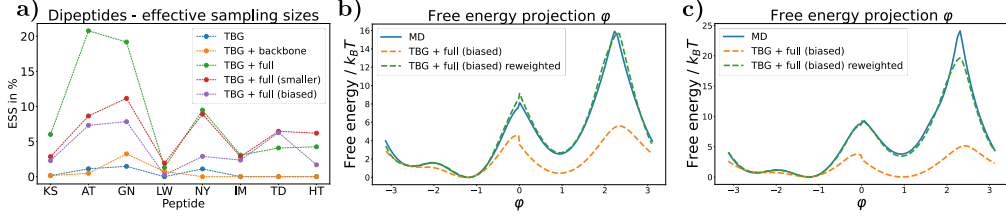


Figure 4: (a) Effective samples sizes (ESS) for the first 8 test peptides for different transferable architectures and training sets. (b) Free energy projection along the φ angle for the TBG + full model trained on the *biased* dataset for the KS dipeptide. The weighted free energy projection demonstrates a superior fit compared to the TBG + full model (see Figure 2e). (c) Free energy projection along the φ angle for the TBG + full model trained on the *biased* dataset for the GN dipeptide. The weighted free energy projection demonstrates a superior fit compared to the TBG + full model (see Figure 3e).

TBG and TBG + backbone models renders them unsuitable as Boltzmann Generators for this setting (Table 2 and Appendix A.5). Furthermore, the TBG + full model always find all metastable states for unseen test peptides (see also Appendix A.5).

The results for the well-performing TBG + full model are presented for two exemplary peptides from the test set in Figure 2 and Figure 3. They were chosen as all architectures sample relevant amounts of correct configurations. Detailed results for other evaluated test peptides are shown in Appendix A.5.

The TBG + full model is an exemplary Boltzmann Emulator, as it is capable of capturing all metastable states of the target Boltzmann distribution (Figure 2b,c and Figure 3b,c). However, it is furthermore also a capable Boltzmann Generator, as it allows for efficient reweighting (Figure 2d,e,f and Figure 3d,e,f). To identify different metastable states, we employ time-lagged independent component analysis (TICA) [60], a dimensionality reduction technique that separates metastable states. We show this analysis in addition to the Ramachandran plots for the dihedral angles.

Moreover, we investigate the influence of the training set in two ablation studies.

Training on a biased training set Our alanine dipeptide results as well as [22] indicate that it can be advantageous to bias the training data towards states that are less probable, such as positive φ states, to recover free energy landscapes. Therefore, we bias the training data by weighting positive φ states for each training peptide, such that they have nearly equal weight to the negative states (see also Appendix B.4). We show that a TBG + full model trained on this dataset (TBG + full (biased)) produces even more accurate free energy landscapes for both the Ramachandran and TICA projections (Figure 4bc). Notably, the unweighted projection shows a clear bias, as expected. However, as the training data is now biased, the effective sample size (ESS) is generally lower (Table 2 and Figure 4a).

Training on a smaller training set Additionally, we examine the impact of smaller training sets on the generalization results. To this end, we train the TBG + full model on two smaller datasets with shorter simulation times: (i) 5ns for each training simulation and (ii) only 500ps of each training simulation. Consequently, the training trajectories are 10 times and 100 times smaller than before. As we utilize only the initial portion of each trajectory, a greater number of metastable states are missed during the brief simulations, as illustrated in Appendix A.2. While the training on the tenfold smaller trainings set, we refer to the model as TBG + full (smaller), shows similar results to training on the whole trainings set (Table 2 and Appendix A.5), the even smaller trainings set leads to inferior results, with several metastable states being missed as presented in Appendix A.3. Nevertheless, we demonstrated that TBGs can be trained with very small datasets, with trajectories that individually miss many metastable states.

6 Discussion

For the first time, we demonstrated the feasibility of training *transferable* Boltzmann Generators. We introduced a general framework for training and evaluating transferable Boltzmann Generators based on continuous normalizing flows. Furthermore, we developed a transferable architecture based on equivariant graph neural networks and demonstrated the importance of including topology

Table 2: Effective samples size and correct configuration rate for unseen dipeptides across different transferable Boltzmann Generator (TBG) architectures.

Model	ESS (\uparrow)		Correct configurations (\uparrow)	
	Mean	Range	Mean	Range
TBG	$0.48 \pm 0.59\%$	(0.0%, 1.47%)	$13 \pm 18\%$	(1%, 48%)
TBG + backbone	$0.58 \pm 1.04\%$	(0.0%, 3.24%)	$17 \pm 21\%$	(1%, 52%)
TBG + full	$8.53 \pm 6.99\%$	(1.31%, 20.79%)	$95 \pm 4\%$	(86%, 100%)
TBG + full (smaller)	$6.13 \pm 3.13\%$	(1.93%, 11.16%)	$96 \pm 3\%$	(88%, 100%)
TBG + full (biased)	$3.86 \pm 2.67\%$	(0.24%, 7.84%)	$96 \pm 4\%$	(87%, 100%)

information in the architecture to enable efficient generalization to unseen, but similar systems. The transferable Boltzmann Generator was evaluated on dipeptides, where significant effective sample sizes were demonstrated on unseen test peptides and accurate sampling of physical properties, such as the free energy difference between metastable states, was achieved. Moreover, we have shown in ablation studies that transferable Boltzmann Generators can be extremely data efficient, with even small training trajectories being sufficient. Future research will determine whether and how transferable Boltzmann Generators can be scaled to larger systems.

7 Limitations / Future work

We leave the scaling to larger system for future work. Notably, this usually requires large amounts of computational resources, as e.g. shown in [10], where they are able to train their transferable model on tetrapeptides, but use more than 100 times more parameters than us.

Instead of flow matching, one could use optimal transport flow matching [21] or equivariant optimal transport flow matching [22] for training, but as indicated in [22] the effect for molecular systems, especially in the presence of many distinguishable particles, are small.

Throughout our work, we utilize a standard Gaussian prior distribution. However, as recently introduced, an alternative is to use a Harmonic prior distribution [61, 62], where atoms that are close in the bond graph are sampled in the vicinity of each other. Notably, we experimented with this different prior distribution, but did not find relevant improvements for our transferable model. This finding is in alignment with the results of [63] that chemical informed prior distributions do not enhance performance significantly compared to simpler uninformed prior distributions for flow matching for molecules. Instead, the network architecture and inductive bias are more important.

Despite conducting a series of ablation studies, we did not pursue the impact of a training set comprising a smaller number of peptides. Instead, we opted for a shorter trajectory approach. Furthermore, we could consider relaxing the 2AA dataset with the semi-empirical force field and training on this modified version, analogous to the alanine dipeptide experiment.

The EGNN architecture was employed for the vector field, as it permits fast evaluation. However, a promising avenue for future research is to explore alternative architectures for the vector field, such as [54, 55, 56, 57, 64, 65, 62], to ascertain whether this enhances performance, which may be necessary to enable scaling to larger systems than those considered. We hope that our provided framework will enable the scaling of transferable Boltzmann Generators to larger systems in future research.

8 Broader Impact

This work represents foundational research with no immediate societal impact. However, if our method is scalable to larger, more relevant systems, it could facilitate the acceleration of drug and material discovery by replacing or enhancing MD simulations, which often play a crucial part in the process. A potential risk is that it could then be used to identify new diseases or biological weapons. Another potential risk associated with this method is that, at present, no convergence criterion is known. This implies that it is not possible to be certain that all potential configurations have been identified, even if an infinite number of samples are taken. This could result in false claims regarding the results, which could have an impact on subsequent applications.

References

- [1] 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*, pages 1–3, 2024.
- [2] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [3] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *BioRxiv*, 2022:500902, 2022.
- [4] Emiel Hoogeboom, Víctor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion for molecule generation in 3D. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 8867–8887. PMLR, 17–23 Jul 2022.
- [5] Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric diffusion model for molecular conformation generation. In *International Conference on Learning Representations*, 2022.
- [6] Bowen Jing, Gabriele Corso, Jeffrey Chang, Regina Barzilay, and Tommi Jaakkola. Torsional diffusion for molecular conformer generation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 24240–24253. Curran Associates, Inc., 2022.
- [7] Frank Noé, Simon Olsson, Jonas Köhler, and Hao Wu. Boltzmann generators — sampling equilibrium states of many-body systems with deep learning. *Science*, 365:eaaw1147, 2019.
- [8] Alessandro Coretti, Sebastian Falkner, Jan Weinreich, Christoph Dellago, and O Anatole von Lilienfeld. Boltzmann generators and the new frontier of computational sampling in many-body systems. *arXiv preprint arXiv:2404.16566*, 2024.
- [9] Grant M Rotskoff. Sampling thermodynamic ensembles of molecular systems with generative neural networks: Will integrating physics-based models close the generalization gap? *Current Opinion in Solid State and Materials Science*, 30:101158, 2024.
- [10] Leon Klein, Andrew Y. K. Foong, Tor Erlend Fjelde, Bruno Kacper Mlodozieniec, Marc Brockschmidt, Sebastian Nowozin, Frank Noe, and Ryota Tomioka. Timewarp: Transferable acceleration of molecular dynamics by learning time-coarsened dynamics. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [11] Esteban G Tabak, Eric Vanden-Eijnden, et al. Density estimation by dual ascent of the log-likelihood. *Communications in Mathematical Sciences*, 8(1):217–233, 2010.
- [12] Esteban G Tabak and Cristina V Turner. A family of nonparametric density estimation algorithms. *Communications on Pure and Applied Mathematics*, 66(2):145–164, 2013.
- [13] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *CoRR*, 2019.
- [14] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- [15] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in neural information processing systems*, pages 6571–6583, 2018.
- [16] Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, and David Duvenaud. Scalable reversible generative models with free-form continuous dynamics. In *International Conference on Learning Representations*, 2019.

- [17] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *CoRR*, 2014.
- [18] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- [19] Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *The Eleventh International Conference on Learning Representations*, 2023.
- [20] Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023.
- [21] Alexander Tong, Nikolay Malkin, Guillaume Hugué, Yanlei Zhang, Jarrod Rector-Brooks, Kilian Fatras, Guy Wolf, and Yoshua Bengio. Conditional flow matching: Simulation-free dynamic optimal transport. *arXiv preprint arXiv:2302.00482*, 2023.
- [22] Leon Klein, Andreas Krämer, and Frank Noé. Equivariant flow matching. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [23] Laurence Illing Midgley, Vincent Stimper, Javier Antoran, Emile Mathieu, Bernhard Schölkopf, and José Miguel Hernández-Lobato. SE(3) equivariant augmented coupling flows. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [24] Peter Wirnsberger, Andrew J Ballard, George Papamakarios, Stuart Abercrombie, Sébastien Racanière, Alexander Pritzel, Danilo Jimenez Rezende, and Charles Blundell. Targeted free energy estimation via learned mappings. *The Journal of Chemical Physics*, 153(14):144112, 2020.
- [25] Manuel Dibak, Leon Klein, Andreas Krämer, and Frank Noé. Temperature steerable flows and Boltzmann generators. *Phys. Rev. Res.*, 4:L042005, Oct 2022.
- [26] Jonas Köhler, Andreas Krämer, and Frank Noé. Smooth normalizing flows. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 2796–2809. Curran Associates, Inc., 2021.
- [27] Laurence Illing Midgley, Vincent Stimper, Gregor N. C. Simm, Bernhard Schölkopf, and José Miguel Hernández-Lobato. Flow annealed importance sampling bootstrap. In *The Eleventh International Conference on Learning Representations*, 2023.
- [28] Xinqiang Ding and Bin Zhang. Deepbar: A fast and exact method for binding free energy computation. *Journal of Physical Chemistry Letters*, 12:2509–2515, 3 2021.
- [29] Joseph C. Kim, David Bloore, Karan Kapoor, Jun Feng, Ming-Hong Hao, and Mengdi Wang. Scalable normalizing flows enable boltzmann generators for macromolecules. In *International Conference on Learning Representations (ICLR)*, 2024.
- [30] Rasool Ahmad and Wei Cai. Free energy calculation of crystalline solids using normalizing flows. *Modelling and Simulation in Materials Science and Engineering*, 30(6):065007, September 2022.
- [31] Kim A Nicoli, Christopher J Anders, Lena Funcke, Tobias Hartung, Karl Jansen, Pan Kessel, Shinichi Nakajima, and Paolo Stornati. Estimation of thermodynamic observables in lattice field theories with deep generative models. *Physical review letters*, 126(3):032001, 2021.
- [32] Hao Wu, Jonas Köhler, and Frank Noé. Stochastic normalizing flows. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 5933–5944. Curran Associates, Inc., 2020.
- [33] Xinqiang Ding and Bin Zhang. Computing absolute free energy with deep generative models. *Biophysical Journal*, 120(3):195a, 2021.
- [34] Andrea Rizzi, Paolo Carloni, and Michele Parrinello. Multimap targeted free energy estimation, 2023.

- [35] Jonas Köhler, Leon Klein, and Frank Noé. Equivariant flows: exact likelihood generative learning for symmetric densities. In *International conference on machine learning*, pages 5361–5370. PMLR, 2020.
- [36] Danilo Jimenez Rezende, Sébastien Racanière, Irina Higgins, and Peter Toth. Equivariant Hamiltonian flows. *arXiv preprint arXiv:1909.13739*, 2019.
- [37] Victor Garcia Satorras, Emiel Hoogetboom, Fabian Fuchs, Ingmar Posner, and Max Welling. E(n) equivariant normalizing flows. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 4181–4192. Curran Associates, Inc., 2021.
- [38] Jonas Köhler, Michele Invernizzi, Pim de Haan, and Frank Noé. Rigid body flows for sampling molecular crystal structures. In *International Conference on Machine Learning, ICML 2023*, volume 202 of *Proceedings of Machine Learning Research*, pages 17301–17326. PMLR, 2023.
- [39] Andrea Rizzi, Paolo Carloni, and Michele Parrinello. Targeted free energy perturbation revisited: Accurate free energies from mapped reference potentials. *Journal of Physical Chemistry Letters*, 12:9449–9454, 2021.
- [40] Michele Invernizzi, Andreas Krämer, Cecilia Clementi, and Frank Noé. Skipping the replica exchange ladder with normalizing flows. *The Journal of Physical Chemistry Letters*, 13:11643–11649, 2022.
- [41] Osama Abidin and Philip M Kim. Pepflow: direct conformational sampling from peptide energy landscapes through hypernetwork-conditioned diffusion. *bioRxiv*, pages 2023–06, 2023.
- [42] Bowen Jing, Bonnie Berger, and Tommi Jaakkola. Alphafold meets flow matching for generating protein ensembles. *arXiv preprint arXiv:2402.04845*, 2024.
- [43] Shuxin Zheng, Jiyan He, Chang Liu, Yu Shi, Ziheng Lu, Weitao Feng, Fusong Ju, Jiayi Wang, Jianwei Zhu, Yaosen Min, et al. Predicting equilibrium distributions for molecular systems with deep learning. *Nature Machine Intelligence*, pages 1–10, 2024.
- [44] Juan Viguera Diez, Sara Romeo Atance, Ola Engkvist, and Simon Olsson. Generation of conformational ensembles of small molecules via surrogate model-assisted molecular dynamics. *Machine Learning: Science and Technology*, 5(2):025010, 2024.
- [45] Mathias Schreiner, Ole Winther, and Simon Olsson. Implicit transfer operator learning: Multiple time-resolution models for molecular dynamics. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [46] Nicholas E Charron, Felix Musil, Andrea Guljas, Yaoyi Chen, Klara Bonneau, Aldo S Pasos-Trejo, Jacopo Venturin, Daria Gusew, Iryna Zaporozhets, Andreas Krämer, et al. Navigating protein landscapes with a machine-learned transferable coarse-grained model. *arXiv preprint arXiv:2310.18278*, 2023.
- [47] Jonas Köhler, Yaoyi Chen, Andreas Krämer, Cecilia Clementi, and Frank Noé. Flow-matching: Efficient coarse-graining of molecular dynamics without forces. *Journal of Chemical Theory and Computation*, 19(3):942–952, 2023.
- [48] Yuxuan Song, Jingjing Gong, Minkai Xu, Ziyao Cao, Yanyan Lan, Stefano Ermon, Hao Zhou, and Wei-Ying Ma. Equivariant flow matching with hybrid probability transport for 3d molecule generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [49] Yan Wang, Lihao Wang, Yuning Shen, Yiqun Wang, Huizhuo Yuan, Yue Wu, and Quanquan Gu. Protein conformation generation via force-guided se (3) diffusion models. *arXiv preprint arXiv:2403.14088*, 2024.
- [50] Felix Draxler, Peter Sorrenson, Lea Zimmermann, Armand Rousselot, and Ullrich Köthe. Free-form flows: Make any architecture a normalizing flow. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 2197–2205. PMLR, 02–04 May 2024.

- [51] H Wiegand. Kish, I.: Survey sampling. John Wiley & Sons, Inc., New York, London 1965, ix+643 s., 31 abb., 56 tab., Preis 83 s., 1968.
- [52] George Papamakarios, Eric T Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *J. Mach. Learn. Res.*, 22(57):1–64, 2021.
- [53] 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.
- [54] Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael John Lamarre Townshend, and Ron Dror. Learning from protein structure with geometric vector perceptrons. In *International Conference on Learning Representations*, 2021.
- [55] Kristof Schütt, Oliver Unke, and Michael Gastegger. Equivariant message passing for the prediction of tensorial properties and molecular spectra. In *International Conference on Machine Learning*, pages 9377–9388. PMLR, 2021.
- [56] Yi-Lun Liao and Tess Smidt. Equiformer: Equivariant graph attention transformer for 3d atomistic graphs. In *The Eleventh International Conference on Learning Representations*, 2023.
- [57] Johannes Gasteiger, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional graph neural networks for molecules. *Advances in Neural Information Processing Systems*, 34:6790–6802, 2021.
- [58] Jonas Köhler, Leon Klein, and Frank Noé. Equivariant flows: exact likelihood generative learning for symmetric densities. In *International Conference on Machine Learning*, pages 5361–5370. PMLR, 2020.
- [59] Christoph Bannwarth, Sebastian Ehlert, and Stefan Grimme. Gfn2-xtb—an accurate and broadly parametrized self-consistent tight-binding quantum chemical method with multipole electrostatics and density-dependent dispersion contributions. *Journal of Chemical Theory and Computation*, 15(3):1652–1671, 2019. PMID: 30741547.
- [60] Guillermo Pérez-Hernández, Fabian Paul, Toni Giorgino, Gianni De Fabritiis, and Frank Noé. Identification of slow molecular order parameters for Markov model construction. *The Journal of chemical physics*, 139(1):07B604_1, 2013.
- [61] Bowen Jing, Ezra Erives, Peter Pao-Huang, Gabriele Corso, Bonnie Berger, and Tommi Jaakkola. Eigenfold: Generative protein structure prediction with diffusion models. *arXiv preprint arXiv:2304.02198*, 2023.
- [62] Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Harmonic self-conditioned flow matching for multi-ligand docking and binding site design. *arXiv preprint arXiv:2310.05764*, 2023.
- [63] Dina A Sharon, Yining Huang, Motolani Oyewole, and Sammy Mustafa. How to go with the flow: an analysis of flow matching molecular docking performance with priors of varying information content. In *ICLR 2024 Workshop on Generative and Experimental Perspectives for Biomolecular Design*, 2024.
- [64] Weitao Du, Yuanqi Du, Limei Wang, Dieqiao Feng, Guifeng Wang, Shuiwang Ji, Carla Gomes, and Zhi-Ming Ma. A new perspective on building efficient and expressive 3d equivariant graph neural networks. *arXiv preprint arXiv:2304.04757*, 2023.
- [65] Han Yang, Chenxi Hu, Yichi Zhou, Xixian Liu, Yu Shi, Jielan Li, Guanzhi Li, Zekun Chen, Shuizhou Chen, Claudio Zeni, et al. Mattersim: A deep learning atomistic model across elements, temperatures and pressures. *arXiv preprint arXiv:2405.04967*, 2024.
- [66] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems* 32, pages 8024–8035. Curran Associates, Inc., 2019.

- 524 [67] Michael Poli, Stefano Massaroli, Atsushi Yamashita, Hajime Asama, Jinkyoo Park, and Stefano
525 Ermon. Torchdyn: Implicit models and neural numerical methods in pytorch. In *Neural*
526 *Information Processing Systems, Workshop on Physical Reasoning and Inductive Biases for the*
527 *Real World*, volume 2, 2021.
- 528 [68] Aric Hagberg, Pieter Swart, and Daniel S Chult. Exploring network structure, dynamics, and
529 function using networkx. Technical report, Los Alamos National Lab.(LANL), Los Alamos,
530 NM (United States), 2008.
- 531 [69] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
532 *arXiv:1412.6980*, 2014.

Table 3: Dimensionless free energy differences for the slowest transition of alanine dipeptide estimated with various methods. Umbrella sampling yields a converged reference solution. Errors are calculated over five runs. Values for BG + backbone and Umbrella sampling are taken from [22].

	Umbrella sampling	BG + backbone [22]	TBG + full (ours)
Free energy difference / $k_B T$	4.10 ± 0.26	4.10 ± 0.08	4.09 ± 0.05

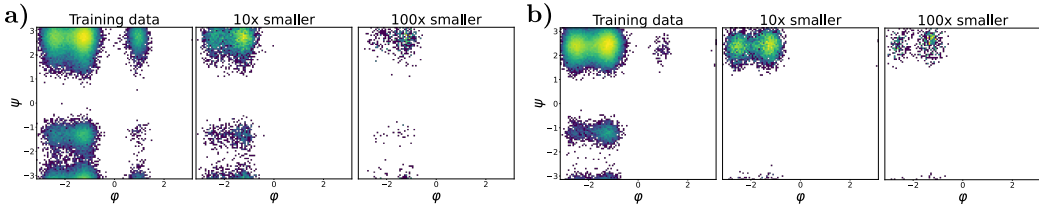


Figure 5: Example Ramachandran plots for different trajectory lengths for the training data. It can be observed that as the trajectory length decreases, the number of metastable states that are missed increases, thereby making the learning task more challenging. (a) AY dipeptide (b) IH dipeptide.

Appendix

A Additional results and experiments

A.1 Semi-empirical force field for alanine dipeptide

We report the free energy differences for the slowest transitions of alanine dipeptide for a semi-empirical force field in Table 3. See Section 5 for more details.

A.2 Dipeptide training data

When training on smaller training sets, i.e. with shorter trajectories, additional metastable states will not be visited during the short simulation times. We show this for two example training peptides in Figure 5. Nevertheless, the TBG + full (smaller) model trained on 10 times shorter trajectories, is nearly as good as the model trained on the full trajectories, see Section 5. However, for the 100 times smaller trajectories, the TBG + full model perform significantly worse, see Appendix A.3.

A.3 Smaller dataset

We investigate the effect of 100 times smaller trainings trajectories, i.e. simulation time of only 500ps. As shown in Appendix A.2, these trajectories miss many metastable states. This can be also observed for the so trained models, which we refer as TBG + full (smaller500), as they do not capture especially unlikely metastable states well as presented in Figure 6 and Figure 7. In contrast, models trained on larger trajectories find all metastable states and allow for efficient reweighting, as discussed in Section 5 and Appendix A.5.

A.4 Sampled dipeptide configurations

For some amino acid combinations, both the TBG and TBG + backbone models sample only a small number of correct configurations. Although the generated configurations are potentially valid molecular configurations, they are not the one of the target dipeptide as shown in Figure 8. Only the various TBG + full architectures samples nearly exclusively correct configurations.

A.5 Additional results for dipeptides

Inference is a costly process, and extensive sampling is necessary to obtain reliable estimates for the expected sample size (ESS). Therefore, we only evaluate the transferable models on a subset of the test set. The dipeptides are randomly selected, but it is ensured that all amino acids are represented at least once. However, we evaluate the best-performing model, namely TBG + full, for twice as many

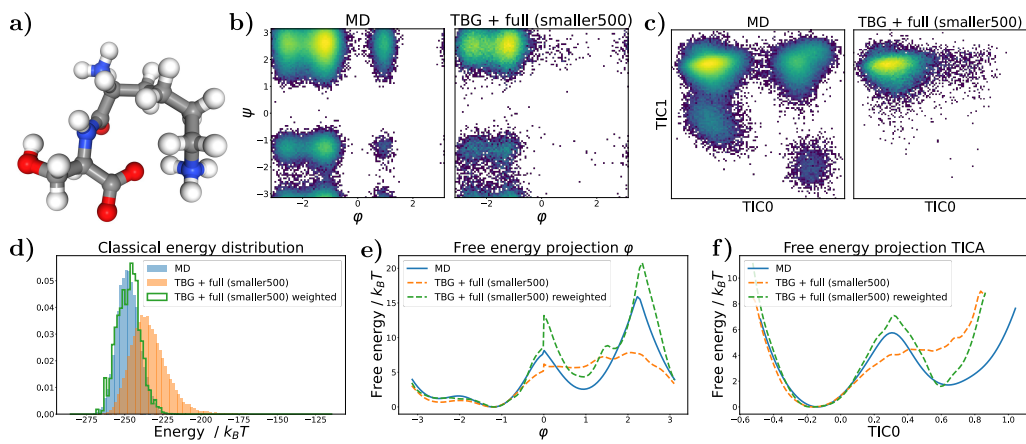


Figure 6: Results for the KS dipeptide for TBG + full model trained on 100 times smaller training trajectories. As can be seen in Figure 2, the results for the TBG + full model trained on the whole trajectories are much better. (a) KS dipeptide (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the model (right). (d) Energies of samples generated with the model. (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

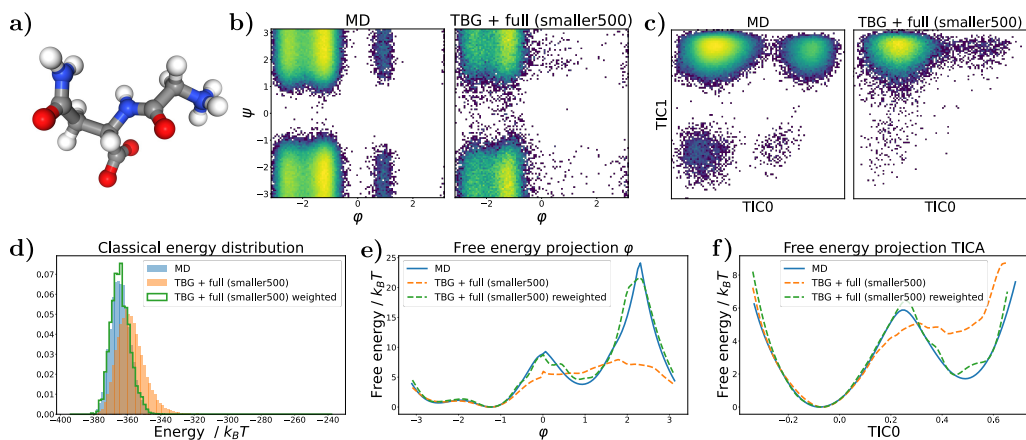


Figure 7: Results for the GN dipeptide for TBG + full model trained on 100 times smaller training trajectories. As can be seen in Figure 3, the results for the TBG + full model trained on the whole trajectories are much better. (a) GN dipeptide (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the model (right). (d) Energies of samples generated with the model. (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

561 test peptides as the rest. The results for the additional test peptides are in good agreement with the
 562 first half as presented in Table 4 and Figure 9.

563 We report individual results for the different architectures in Figure 9.

564 To illustrate the performance of the TBG + full model, we present additional examples of dipeptides
 565 from the test set in Figure 10a-f and Figure 11a-f. Furthermore, we also again show results for
 566 training on the biased dataset in the same figures (Figure 10g,h,i and Figure 11g,h,i). As observed
 567 previously, the TBG + full (biased) model recovers the free energy landscape better than the TBG +
 568 full model, especially for the φ projections. We present additional examples of dipeptides from the
 569 test set for the TBG + full model in Figure 12, Figure 13, Figure 14, Figure 15, and Figure 16.

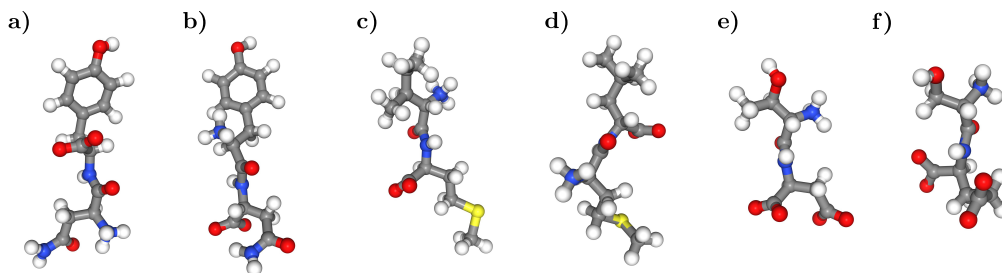


Figure 8: Sampled molecules with the TBG and TBG + backbone models, which do not have the correct topology. (a) NY dipeptide reference (b) Generated molecule with NY atoms by the TBG model. (c) IM dipeptide reference (d) Generated molecule with IM atoms by the TBG model. (e) TD dipeptide reference (f) Generated molecule with TD atoms by the TBG + backbone model.

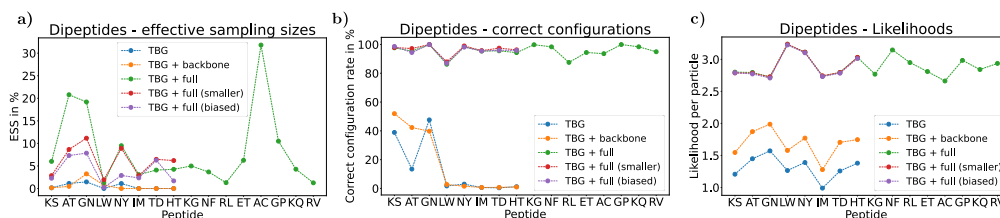


Figure 9: Performance comparison for different transferable architectures and training sets on the test set (a) Effective samples sizes (ESS) (b) Correct configuration rate (c) Likelihood per particle.

Furthermore, we present results for two example peptides from the test set for the TBG + full (smaller) model, which is trained on trajectories that are tenfold smaller than those used for the TBG + full model. These results are shown in Figure 17 and Figure 18.

A.6 Transferable Boltzmann Generators as Boltzmann Emulators

Given the high cost of sampling with CNFs, which necessitates integrating the Jacobian trace along the positions, we did not evaluate all available test peptides (see Appendix B.2). However, since sampling without the Jacobian trace is less expensive and we do not require as many samples as for estimating the ESS, we also employ the TBG + full (smaller) model as a Boltzmann Emulator to ascertain whether we have identified all metastable states, despite the fact that it was only trained on the 10 times smaller training set. The Boltzmann Emulator is evaluated on a diverse set of test peptides, and nearly always finds all metastable states within less than one hour of wall clock time. This is a notable improvement over MD simulations, which often take longer to explore due to the iterative nature of MD. Some examples are shown in Figure 19. This experiment shares similarities with the exploration mode of [10], where they employ their model without the acceptance step and therefore also explore a potentially biased distribution rather than the unbiased Boltzmann distribution.

Table 4: Effective samples size and correct configuration rate for unseen dipeptides for the TBG + full architecture for different number of test peptides.

Model	ESS (\uparrow)		Correct configurations (\uparrow)	
	Mean	Range	Mean	Range
TBG + full (8 test peptides)	$8.53 \pm 6.99\%$	(1.31%, 20.79%)	$95 \pm 4\%$	(86%, 100%)
TBG + full (16 test peptides)	$8.27 \pm 8.29\%$	(1.26%, 31.80%)	$96 \pm 4\%$	(86%, 100%)

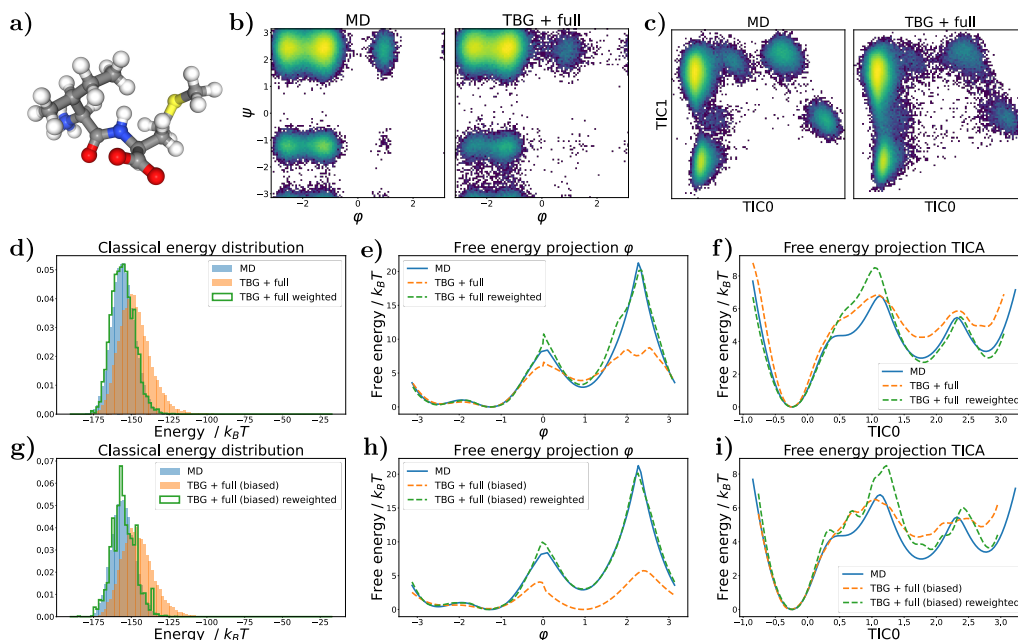


Figure 10: Results for the IM dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of samples generated with the TBG + full model. (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0). (g) Energies of samples generated with the TBG + full (biased) model. (h) Free energy projection along the φ angle for the TBG + full (biased) model. (i) Free energy projection along the slowest transition (TIC0) for the TBG + full (biased) model.

B Technical details

B.1 Code libraries

We primarily use the following code libraries: *PyTorch* (BSD-3) [66], *bgflow* (MIT license) [7, 35], *torchdyn* (Apache License 2.0) [67], and *NetworkX* (BSD-3) [68] for validating graph isomorphisms. Additionally, we use the code from [37] (MIT license) for EGNNs, as well as the code from [10] (MIT license) and [22] (MIT license) for datasets and related evaluation methods.

Our code is available here: https://osf.io/n8vz3/?view_only=1052300a21bd43c08f700016728aa96e. We will make our code public upon publication.

B.2 Benchmark systems

The investigated benchmark systems were created in prior studies [22, 10].

Alanine dipeptide The alanine dipeptide datasets were created in [22] (CC BY 4.0), we refer to them for detailed simulation details. The classical trajectory was created at $T = 300\text{K}$ with the classical *Amber ff99SBildn* force-field. The subsequent relaxation was performed with the semi-empirical *GFN2-xTB* force-field [59].

Dipeptides (2AA dataset) The original dipeptide dataset as introduced in [10] (MIT License) is available here: <https://huggingface.co/datasets/microsoft/timewarp>. As this includes a lot of intermediate saved states and quantities, like energies, we create a smaller version with is available here: https://osf.io/n8vz3/?view_only=1052300a21bd43c08f700016728aa96e. For a comprehensive overview of the simulation details, refer to [10]. All dipeptides were simulated

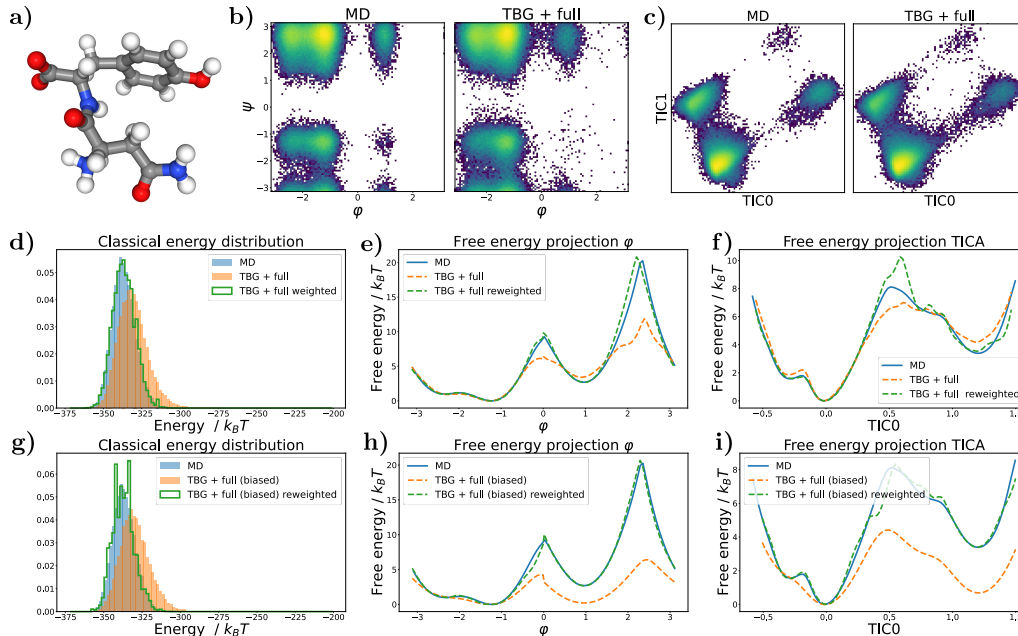


Figure 11: Results for the NY dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of samples generated with the TBG + full model. (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0). (g) Energies of samples generated with the TBG + full (biased) model. (h) Free energy projection along the φ angle for the TBG + full (biased) model. (i) Free energy projection along the slowest transition (TIC0) for the TBG + full (biased) model.

with a classical *amber-14* force-field at $T = 310\text{K}$. The simulation of the training peptides were run for 50ns, while the test set peptides were run for $1\mu\text{s}$.

Choice of test set peptides Inference is a costly process with CNFs (see Appendix B.7), and extensive sampling is necessary to obtain reliable estimates for the relative effective sample size (ESS). Therefore, we only evaluate the transferable models on a subset of the test set. The dipeptides are randomly selected, but it is ensured that all amino acids are represented at least once.

B.3 Hyperparameters

We report the model hyperparameters for the different model architectures as describes in Section 4.1 in Table 5. As in [22] all neural networks ϕ_α have one hidden layer with n_{hidden} neurons and *SiLU* activation functions. The input size of the embedding $n_{\text{embedding}}$ depends on the model architecture.

We report training hyperparameters for the different model architectures in Table 6. It should be noted that all TBG models are trained in an identical manner if the training set is identical. We use the ADAM optimizer for all experiments [69]. For the dipeptide training, each batch consists of three samples for each peptide.

B.4 Biasing target samples

As introduced in [22], it can be beneficial to bias the training data in such a way that unlikely states are more prominent. For alanine dipeptide and many dipeptides, the positive φ states at $\varphi = 1$ are often the unlikely ones and transition between the positive and negative φ states are slow. For the alanine dipeptide dataset, the biasing methodology proposed in [22] is employed. Similarly, we bias the dipeptides based on the von Mises distribution f_{vM} . The weights ω are computed along the φ

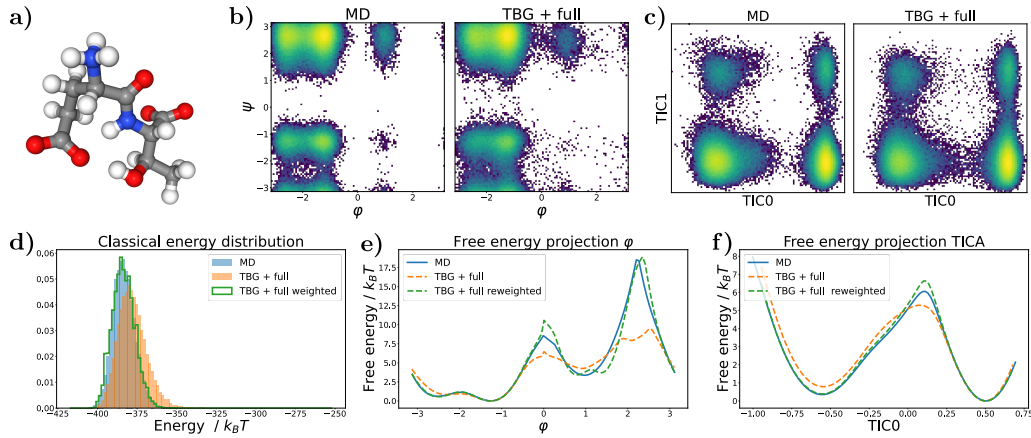


Figure 12: Results for the ET dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of generated samples (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

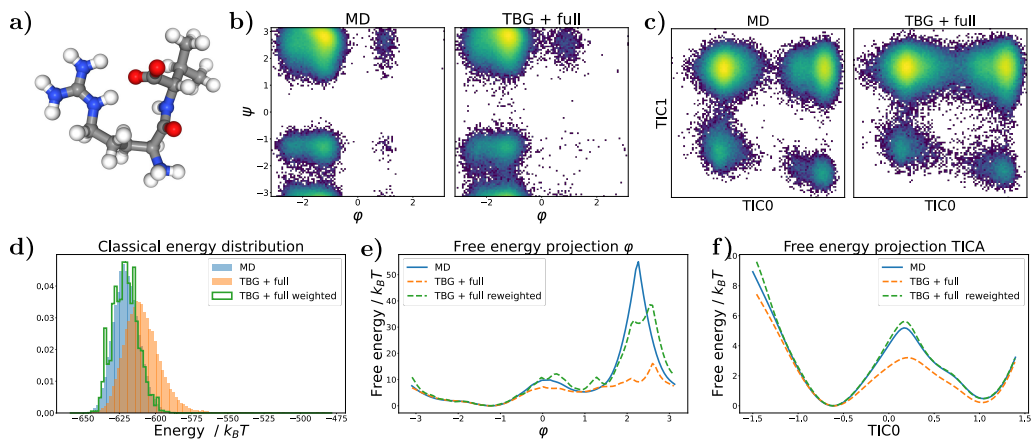


Figure 13: Results for the RV dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of generated samples (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

625 dihedral angle as

$$\omega(\varphi) = r \cdot f_{\text{VM}}(\varphi | \mu = 1, \kappa = 10) + 1, \quad (12)$$

626 where r is computed based on the ratio of positive and negative φ states, such that both have nearly
627 the same weight after the biasing.

628 B.5 Encoding of atom types

629 The atom type embedding a_i is a one-hot vector of 54 classes. The classes are mostly defined by the
630 atom type in the topology for a classical force field. Therefore, only a few atoms are indistinguishable,
631 such as hydrogen atoms that are bound to the same carbon or nitrogen atom. Moreover, we also treat
632 oxygen atoms bound to the same carbon atom as indistinguishable, unless they are in the carboxyl
633 group. Notably, we never treat particle groups as indistinguishable, such as two CH3 groups bound to
634 the same carbon atom.

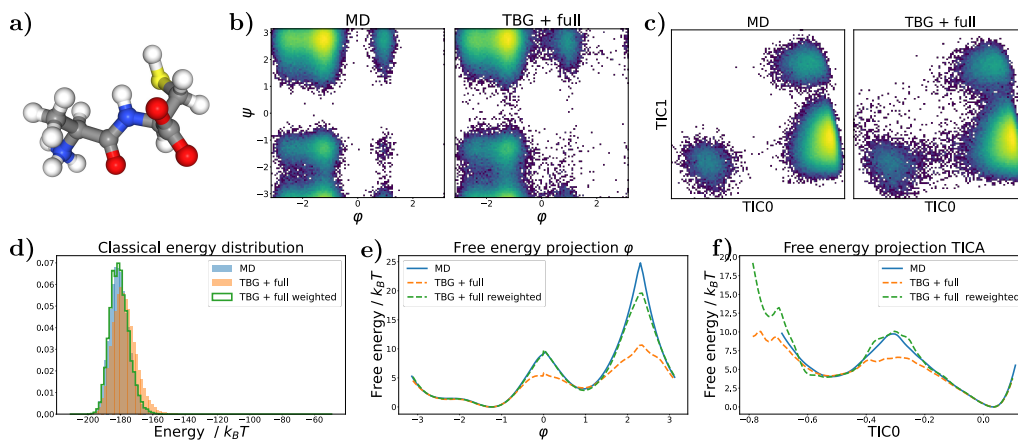


Figure 14: Results for the AC dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of generated samples (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

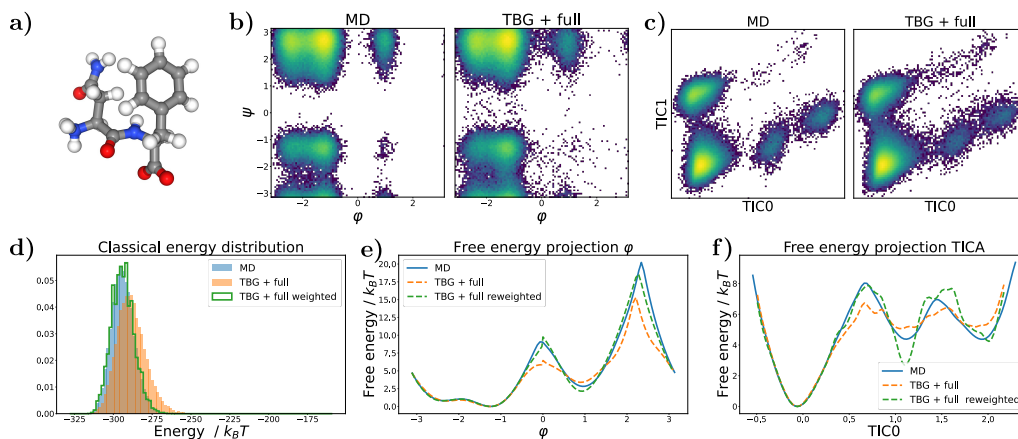


Figure 15: Results for the NF dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of generated samples (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

635 B.6 Effective samples sizes

636 The relative effective sample sizes (ESS) are computed with Kish’s equation [51] as in prior work.
 637 For the alanine dipeptide experiments we use 2×10^5 samples for the forward ESS and 1×10^4 for the
 638 negative log likelihood computation. A total of 3×10^4 samples were used for each dipeptide in the
 639 forward ESS, while 4.5×10^3 samples were employed for the negative log likelihood computation.

640 B.7 Computing resources

641 All training and inference was performed on single *NVIDIA A100 GPUs* with 80GB of RAM.

642 The training time for the models is reported in Appendix B.3, although it should be noted that a
 643 significant amount of time was required for hyperparameter tuning. It is estimated that at least ten
 644 times the compute time reported in Appendix B.3 was necessary to identify suitable hyperparam-
 645 eters. Furthermore, inference with CNFs is expensive, especially if one requires the reweighting

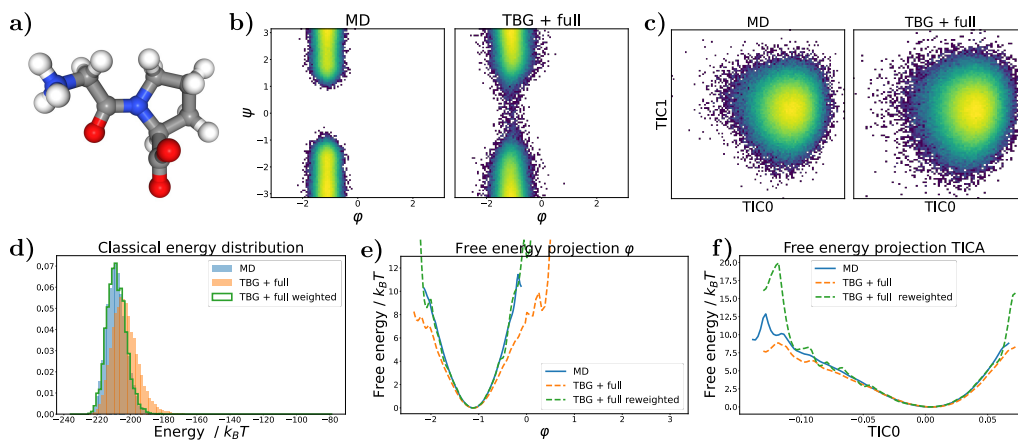


Figure 16: Results for the GP dipeptide (a) Sample generated with the TBG + full model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full model (right). (d) Energies of generated samples (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

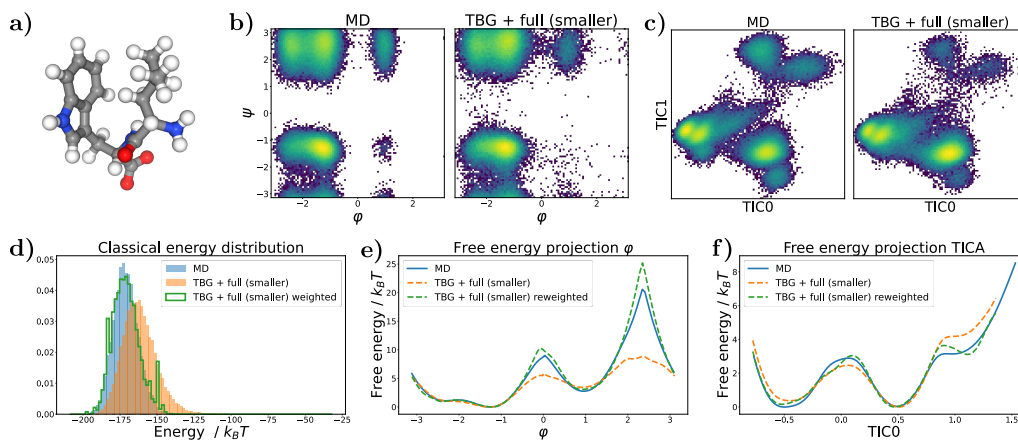


Figure 17: Results for the LW dipeptide for the TBG + full (smaller) model, which is trained on tenfold smaller trajectories than the TBG + full model. (a) Sample generated with the TBG + full (smaller) model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full (smaller) model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full (smaller) model (right). (d) Energies of generated samples. (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

646 weights. Generating 3×10^4 samples with the large transferable models for the dipeptides requires
 647 approximately four days, whereas generating 2×10^5 samples for the alanine dipeptide experiments
 648 takes less than one day. However, generating samples without corresponding weights significantly
 649 accelerates the sampling process. In the case of the dipeptides, the generation of 2×10^5 samples can
 650 be completed in less than one day. However, it should be noted that sampling can be done fully in
 651 parallel, as Boltzmann Generators generate independent samples.

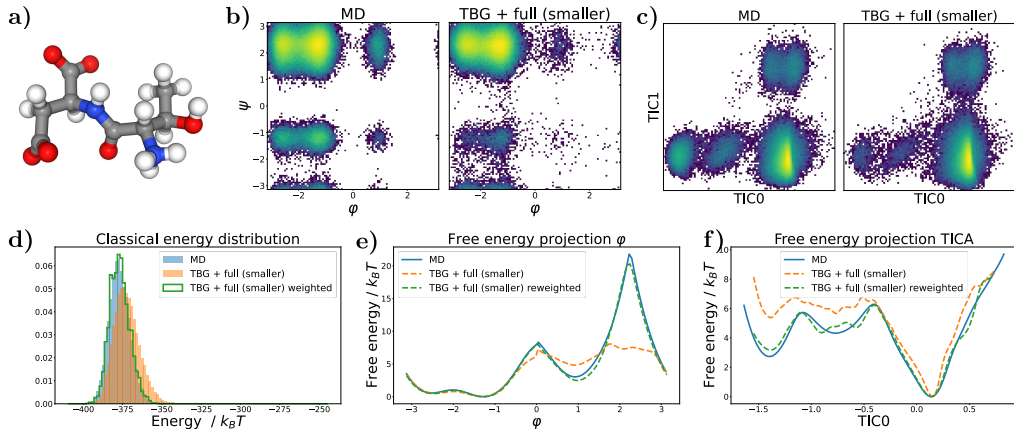


Figure 18: Results for the TD dipeptide for the TBG + full (smaller) model, which is trained on tenfold smaller trajectories than the TBG + full model. (a) Sample generated with the TBG + full (smaller) model (b) Ramachandran plot for the weighted MD distribution (left) and for samples generate with the TBG + full (smaller) model (right). (c) TICA plot for the weighted MD distribution (left) and for samples generate with the TBG + full (smaller) model (right). (d) Energies of generated samples (e) Free energy projection along the φ angle. (f) Free energy projection along the slowest transition (TIC0).

Table 5: Model hyperparameters

Model	L	n_{hidden}	$n_{\text{embedding}}$	Num. of parameters
alanine dipeptide				
BG + backbone	5	64	8	147599
TBG + full encoding	5	64	15	149147
Dipeptides (2AA)				
TBG	9	128	5	1044239
TBG + backbone	9	128	13	1046295
TBG + full encoding	9	128	76	1062486

Table 6: Training hyperparameters

Mdoel	Batch size	Learning rate	Epochs	Training time
Alanine dipeptide				
BG + backbone	256	$5e-4/5e-5$	500/500	3.5h
TBG + full	256	$5e-4/5e-5$	500/500	3.5h
Dipeptides (2AA)				
TBG	600	$5e-4/5e-5/5e-6$	4/4/4	3d
TBG + backbone	600	$5e-4/5e-5/5e-6$	4/4/4	3d
TBG + full	600	$5e-4/5e-5/5e-6$	4/4/4	3d
TBG + full (smaller)	600	$5e-4/5e-5/5e-6$	30/30/30	2.5d

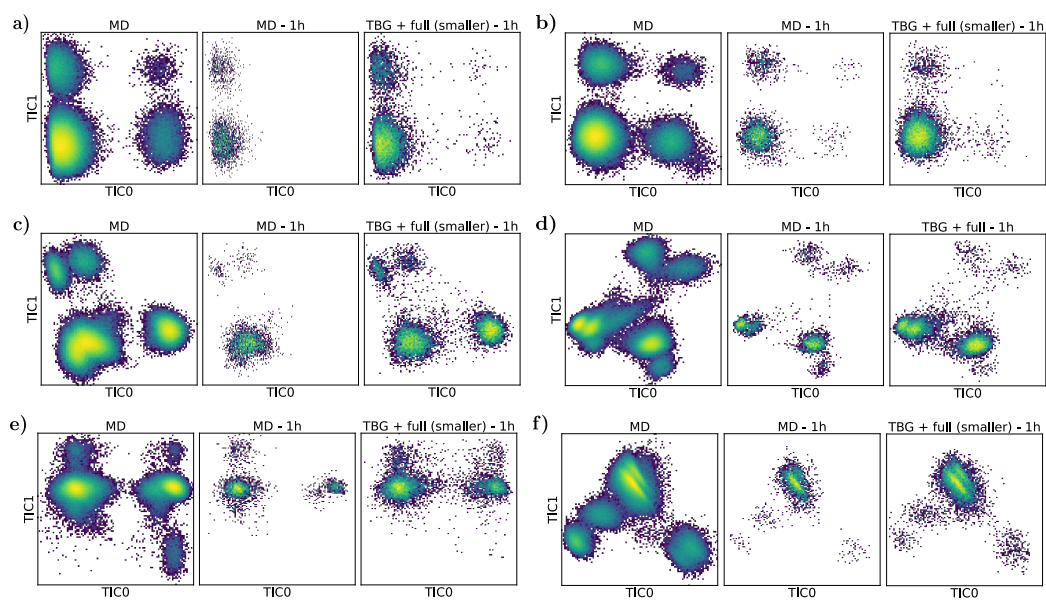


Figure 19: Comparison of classical MD runs for 1 hour (MD - 1h) and the sampling with the TBG + full (smaller) model without weight computation for 1 hour (TBG + full (smaller) - 1h). The TICA plots of different peptides from the test set are shown. It is important to note that the TICA projection is always computed with respect to the long MD trajectory (MD). All peptides stem from the test set. (a) CS dipeptide (b) EK dipeptide (c) KI dipeptide (d) LW dipeptide (e) RL dipeptide (f) TF dipeptide.

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: [\[Yes\]](#)

Justification: The claims made in the abstract and introduction are all based on the results of our work, as shown in Section 5 and Appendix A .

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

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

Answer: [\[Yes\]](#)

Justification: We discuss the limitations in Section 7.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

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: [\[NA\]](#)

Justification: We only utilize theorems of prior work, which we reference accordingly.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

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: [\[Yes\]](#)

Justification: We explain how we performed our experiments in Section 5 and Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the 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]

Justification: Code, data, model checkpoints as well as detailed instructions are available here: https://osf.io/n8vz3/?view_only=1052300a21bd43c08f700016728aa96e, as stated in Appendix B.1. Nevertheless, high level training, evaluation and implementation details are described in Appendix B and Section 4.1.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

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]

Justification: We discuss the dataset details in Appendix B and refer to related work for datasets introduced in prior work.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

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]

Justification: We provide error bars for the alanine dipeptide experiments. In contrast, for the much more expensive transferable experiments, we utilize our computational resources to sample a multitude of different peptides from the test set, rather than training and sampling distinct instances of the same architecture for the same peptide. Consequently, we obtain error bounds by averaging results over different test dipeptides rather than runs.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

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]

Justification: We discuss the required computational resources for this work for the training and inference in Appendix B.3 and Appendix B.7.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

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]

Justification: The research conducted conforms with the NeurIPS Code of Ethics. All authors have read the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

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

Answer: [Yes]

Justification: We discuss the broader impact of our work in Section 8.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

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]

Justification: Our models do not pose such a risk, as they are for molecular data.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

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]

Justification: We give credit to used assets in this work in Appendix B.1 and Appendix B.2. Our assets will be available under the MIT / CC BY 4.0 license.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

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

Answer: [\[Yes\]](#)

Justification: We document our models in Appendix B.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

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: [\[NA\]](#)

Justification: No crowdsourcing nor human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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: [\[NA\]](#)

Justification: No crowdsourcing nor human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- 965 • We recognize that the procedures for this may vary significantly between institutions
966 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
967 guidelines for their institution.
- 968 • For initial submissions, do not include any information that would break anonymity (if
969 applicable), such as the institution conducting the review.