
Latent Space Simulator for Unveiling Molecular Free Energy Landscapes and Predicting Transition Dynamics

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Free Energy Surfaces (FES) and metastable transition rates are key elements in
2 understanding the behaviour of molecules within a system. However, the typi-
3 cal approaches require computing force-fields across billions of time-steps in a
4 molecular dynamics (MD) simulation, which is often considered intractable when
5 dealing with large systems or databases. In this work we propose LAMODY, a
6 latent-space MD simulator to effectively tackle the intractability with around 20-
7 fold speed improvements compared to classical MD's. The model leverages a
8 chirality aware $SE(3)$ -invariant encoder-decoder architecture to generate a latent
9 space, coupled with a recurrent neural network to run the time-wise dynamics. We
10 show that LAMODY effectively recovers realistic trajectories and FES more accu-
11 rately and faster than existing methods, while capturing their major dynamical and
12 conformational properties. Furthermore, the proposed approach can generalize to
13 molecules outside the training distribution.

14 1 Introduction

15 Fundamental quantities of interest towards understand-
16 ing a molecule's dynamics and properties are its Free
17 Energy Surface (FES) and metastable states, along-
18 side its transition rates between metastable states. Ac-
19 cessing them enables many real-world applications in
20 drug discovery or material sciences (Peng et al., 2014;
21 Bochevarov et al., 2013). Each 3D conformation of a
22 molecule is associated with a potential energy that de-
23 termines its probability of occurring (via a Boltzmann
24 distribution).

25 The FES is a lower-dimensional representation of this
26 energy landscape, providing insights into stable states
27 (energy minima), transition pathways, and free energy
28 differences. Additionally, a molecule's kinetics are of
29 interest, such as the transition rates between metastable
30 states/modes of the Boltzmann distribution.

31 The usual approach to compute these properties is to
32 run long micro-second molecular dynamics (MD) simulations. Considering that each MD step is

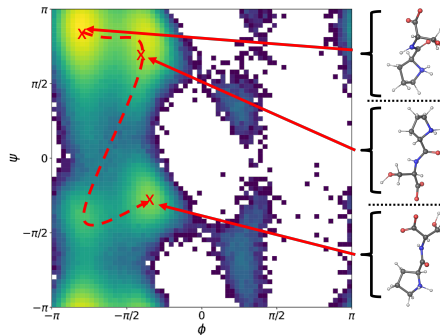


Figure 1: Free Energy Surface (FES) with minima corresponding to different conformations and an example MD trajectory as dotted arrow.

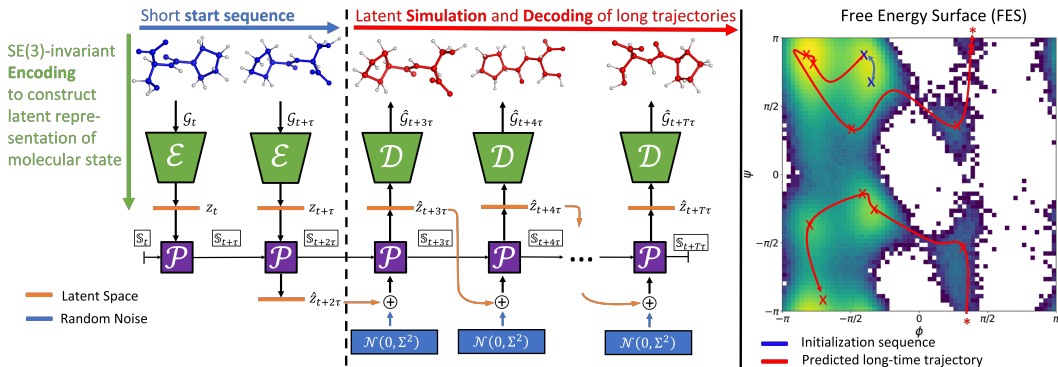


Figure 2: Overview of LAMODY. An encoder \mathcal{E} computes $SE(3)$ -invariant latent embeddings of a short initialization sequence, the dynamical propagator \mathcal{P} iteratively predicts the next states to produce a long-time trajectory in latent space from which molecular conformers can be reconstructed by the decoder \mathcal{D} . The warm-up sequence and predicted trajectory are visualized in the FES. Here, $\mathcal{N}(0, \Sigma)$ denotes random noise, \oplus is vector addition, \mathcal{G}_t denotes the 3D graph representation of a molecule at time t , z is a latent space state, τ is the time lag between states in a trajectory, and * denotes the point where the MD trajectory crosses the plane.

33 in the scale of femto-seconds, the simulation comes with a high computational cost. To accelerate
 34 the recovery of these properties, it is essential to develop a method that (1) can operate at time steps
 35 beyond the femtosecond level; (2) captures the key reaction coordinates; (3) does not suffer from
 36 instabilities (unphysical states) for long-time simulations.

37 Learned simulators operating in a latent space suit these requirements if the latent space captures
 38 reaction coordinates (a molecule’s most important degrees of freedom) since they allow for larger
 39 time steps (Sidky et al., 2020; Vlachas et al., 2022). However, existing architectures restrict the
 40 simulator to only work on a single molecule at a time, meaning that they cannot generalize to new
 41 molecules (Sidky et al., 2020; Vlachas et al., 2022). Furthermore, LED (Vlachas et al., 2022) fails
 42 to recover rare metastable states and lacks practical relevance as it has only been shown to work with
 43 multiple re-initializations from Boltzmann distributed states, meaning that a long MD simulation is
 44 still required to define the starting states.

45 Other approaches, such as Boltzmann generators (Noé et al., 2019) or Distributional Graphormer
 46 (Zheng et al., 2023) can predict the equilibrium distribution of unseen molecules but do not have
 47 a notion of time, i.e., no dynamical properties such as the transition rates can be extracted. In this
 48 regard, machine learning (ML) force fields (Unke et al., 2021; Batzner et al., 2022; Hu et al., 2021)
 49 have made significant progress for ab-initio simulations but are still slower for long simulations and
 50 larger molecules where classical force fields are applied (Fu et al., 2023).

51 To tackle these limitations, we propose a learned Latent Molecular Dynamics LAMODY, model. We
 52 employ an $SE(3)$ -invariant encoder-propagator-decoder scheme based on message-passing neural
 53 networks (MPNN) (Gilmer et al., 2017) that can be trained end-to-end on MD data and can general-
 54 ize to unseen molecules. For the tasks of FES recovery, past studies used different sampling and
 55 evaluation protocols, making it difficult to compare methods. We define scientifically meaningful
 56 tasks and metrics that allow that reflect a model’s practical relevance in probing the free energy
 57 surface of molecules. In summary, our contributions are:

- 58 • 20-fold speed improvements compared to classical MD, thanks to a long operating time
 59 step of 100 fs.
- 60 • Generalization to unseen molecules thanks to our chirality-aware $SE(3)$ -invariant encoder-
 61 decoder.
- 62 • Defining a systematic evaluation scheme to assess the performance of simulation methods
 63 against scientifically meaningful tasks for FES recovery.

64 2 Related work

65 **Enhanced sampling** methods inject bias to the potential energy function to facilitate fast sampling
66 of transitions between local energy minima that are separated by high energy barriers. Popular
67 methods include simulated annealing (Bernardi et al., 2015; Tsallis & Stariolo, 1996), metadynamics
68 (Laio & Gervasio, 2008), replica exchange (Bernardi et al., 2015), umbrella sampling (Torrie &
69 Valleau, 1977), and parallel tempering Yang et al. (2019). A major limitation of enhanced sampling
70 methods lies in the fact that they typically require determining collective variables (CVs) in advance,
71 which can be challenging for complex systems Wang et al. (2021). Furthermore, enhanced sampling
72 methods do not have an explicit notion of “time”, meaning that no extraction of dynamical properties
73 is possible (Stelzl & Hummer, 2017).

74 **Latent Space Simulators** enable to accelerate MD simulations in the 3D configuration space, by
75 updating a latent state generated by a learned encoder, instead of moving each atom according to its
76 velocity and computed force. The updates are performed by a dynamical propagator, and the all-
77 atom representation can be constructed with a decoder. Time-lagged autoencoders with propagators
78 (Otto & Rowley, 2019; Lusch et al., 2018) learn a linear propagator whereas Sidky et al. (2020) use
79 a mixture density network (Bishop, 1994) as a propagator. However, the above methods do not obey
80 the $SE(3)$ -invariance of molecules (they could, e.g., arbitrarily flip a chirality each step). Vlachas
81 et al. (2022) train an LSTM network as propagator and account use a mixture density network as
82 autoencoder. However, this method requires multiple re-initializations from Boltzmann distributed
83 states and it remains unclear if the method stays stable for longer simulations. Additionally, all
84 previously mentioned methods only work on a single molecule they have been trained on - they are
85 not able to generalize unlike LAMODY.

86 3 Method

87 3.1 Model Architecture

88 **Encoder** To make the encoder architecture generalizable to other molecules, we use a graph
89 representation of internal coordinates and employ a Graph Neural Network (GNN) archi-
90 tecture. Concretely, a molecular state is represented by a graph $\mathcal{G} \in (\mathcal{V}, \mathcal{B}, \mathcal{X}, \mathcal{C})$ with
91 each node representing a bond in the original molecule, and edges representing bond angles
92 and torsion angles defined by triplets and quadruplets of bonds respectively, hence $|\mathcal{V}| =$
93 $|\mathbb{B}|$ and $|\mathcal{B}| = |\mathbb{A}| + |\mathbb{T}|$. Nodes are featurized with information about the atoms forming the
94 bond and the bond length and edges are featurized with the respective bond or torsion angle
95 and a categorical feature indicating whether the edge defines a bond angle or a torsion angle
96 ¹. We then employ L message-passing layers akin to Shi et al. (2021), pool the nodes using
97 a learnable set-to-set mapping (Vinyals et al., 2016), and predict the final latent vector using
98 a linear layer.

109 **Decoder** To reconstruct the internal coordinates of a molecular state given a latent
110 representation, we use a second GNN similar to Winter et al. (2021). The decoder takes as input
111 a two-dimensional molecular graph with nodes representing atoms and edges representing bonds
112 and a latent vector describing the molecular state in the latent space. First node level embeddings
113 are computed by iteratively applying a sequence of message-passing layers similar to the encoder.
114 Then, bond lengths are predicted by applying a three-layer MLP onto the concatenated pairs of
115

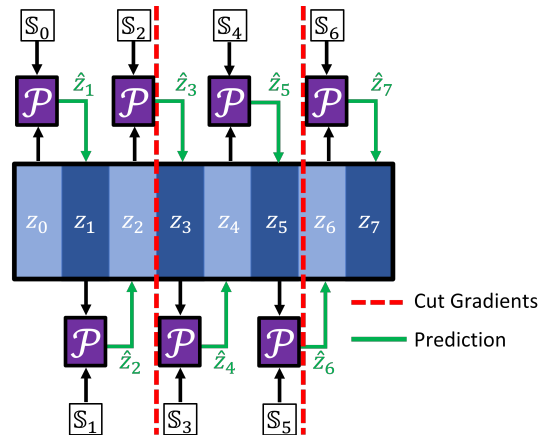


Figure 3: Training scheme for long sequences: The propagator \mathcal{P} takes in a latent state z_t and cell state S_t to predict the latent state at time $t + 1$. The cell states are not re-initialized and gradients are detached after a fixed-length interval.

¹for a detailed description see subsection C.1

116 nodes and the latent embedding, i.e. $d_i = \Pi_{bond}([h_a, h_b, z])$ with h_* being the node embeddings,
 117 z the latent vector and Π_{bond} the MLP. The same approach is taken for bond angles and torsion
 118 angles with triplets/quadruplets of node embeddings and $\Pi_{ang} \cdot \Pi_{tor}$ respectively.
 119

120 **Dynamical Propagator** As suggested by Vlachas et al. (2022), sequences of MD states are not
 121 necessarily Markovian since complex systems can exhibit long-term correlations in their behavior,
 122 meaning that future states can depend on past states, violating the assumption of independence
 123 between time steps. To account for this, we use an LSTM (Hochreiter & Schmidhuber, 1997) as the
 124 dynamical model that is trained to predict the next latent state given a short history. Concretely, we
 125 use

$$\begin{aligned} (\mathbf{h}_{t+\tau}, \mathbf{c}_{t+\tau}) &= LSTM(\mathbf{z}_t, \mathbf{h}_t, \mathbf{c}_t) \\ \mathbf{z}_{t+\tau} &= \Xi(\mathbf{h}_{t+\tau}) \end{aligned} \quad (1)$$

126 where $\mathbf{h}_t, \mathbf{c}_t$ denote the LSTM hidden state and cell state at time t , \mathbf{z}_t is the latent state at time t and
 127 Ξ is a two-layer MLP.

128 3.2 Training

129 We train our model end-to-end on MD data. To do so, we randomly sample a batch of starting points
 130 from the dataset from which we consider the consecutive k states with a time lag τ between states.
 131 Hence, we end up with a batch of sub-sequences of the full trajectory of length $k+1$ states. Starting
 132 with an initial LSTM state of $\mathbb{S}_0 = (\mathbf{h}_0, \mathbf{c}_0) = (\vec{0}, \vec{0})$, we iteratively unfold the LSTM to predict the
 133 next time step, while the LSTM cell states are passed through time. More specifically, we encode
 134 \mathcal{G}_0 into latent space by $z_0 = \mathcal{E}(\mathcal{G}_0)$, from which together with \mathbb{S}_0 the next time step latent state \hat{z}_1
 135 is predicted. Then \mathbb{S}_1 and $z_1 = \mathcal{E}(\mathcal{G}_1)$ are used to predict \hat{z}_2 , which can all be decoded back to
 136 molecular states.

137 To optimize the parameters of the model with backpropagation, we define an end-to-end propagation
 138 loss that is additionally regularized by a reconstruction loss and a latent loss :

$$\begin{aligned} \mathcal{L} &= \delta_{e2e} \frac{1}{k} \sum_{i=1}^k \mathcal{L}_{rec} [\mathcal{G}_i, \mathcal{D} \circ \mathcal{P} \circ \mathcal{E}(\mathcal{G}_{i-1})] \\ &+ \delta_{lat} \frac{1}{k} \sum_{i=1}^k \|\mathbf{z}_i - \hat{\mathbf{z}}_i\|^2 + \delta_{rec} \frac{1}{k+1} \sum_{i=0}^k \mathcal{L}_{rec} [\mathcal{G}_i, \mathcal{D} \circ \mathcal{E}(\mathcal{G}_i)] \end{aligned} \quad (2)$$

139 here $\delta_{rec}, \delta_{lat}, \delta_{e2e}$ are hyperparameters and \mathcal{L}_{rec} is defined as in Equation 11. Note that $\mathbf{z}_i =$
 140 $\mathcal{E}(\mathcal{G}_i)$, $\hat{\mathbf{z}}_i = \mathcal{P} \circ \mathcal{E}(\mathcal{G}_{i-1})$. Although the end-to-end part of our loss function theoretically encap-
 141 sulates the latent and the reconstruction loss, we found the explicit presence of both as additional
 142 regularization to be crucial for the training process to succeed.

143 **Training on long sequences** As we aim to predict long-timescale trajectories at inference time with
 144 $N_{steps} \gg k$, we require training on long sequences without suffering from vanishing or exploding
 145 gradients. To do so, we sample sub-trajectories of length $c * k$ with c being a hyperparameter and
 146 iteratively train on sequences of length k where we keep the LSTM states but detach the gradients
 147 as suggested by Vlachas et al. (2022).

148 3.3 Inference

149 At inference time, we "warm up" the LSTM with a sequence of k MD states from which we iter-
 150 atively unfold the propagator to predict latent trajectories. Additionally, we infuse artificial noise
 151 to the latent states before feeding them into the propagator. We found this to be crucial because
 152 otherwise, the dynamical model was prone to become stuck at a local energy minimum. Concretely,
 153 we predict the next latent state by :

$$\hat{\mathbf{z}}_{t+\tau} = \begin{cases} \mathcal{P}(\hat{\mathbf{z}}_t + \mathcal{N}(0, \Sigma)), & \text{if } x \sim U(0, 1) \leq \beta \\ \mathcal{P}(\hat{\mathbf{z}}_t), & \text{else} \end{cases} \quad (3)$$

154 where $\beta \in [0, 1]$ is a hyperparameter, $x \sim U(0, 1)$ indicates a sample from the uniform distribution
 155 and $\Sigma = \mathbf{I} * \sigma^2, \sigma^2 \in \mathbb{R}^+$ is computed from the warmup trajectory.

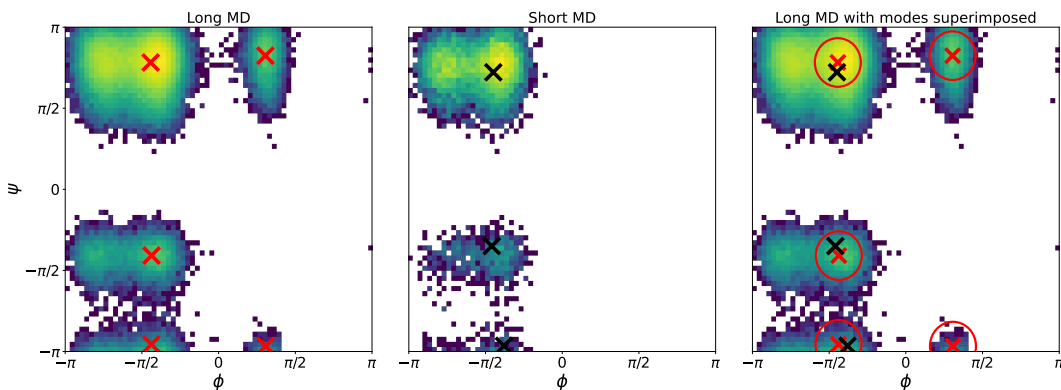


Figure 4: MSPR: Metastable State Precision/Recall; Ramachandran plot of a long and a short MD simulation for a peptide where identified metastable states are indicated by crosses. The third figure shows the long MD trajectory with modes identified by the short MD simulation superimposed and the circles denote the area where a mode is considered to be correct. This allows to compute metastable state precision and recall (MSPR).

156 4 Evaluation Protocol for FES recovery

157 This section aims to provide an evaluation protocol that is both robust and scalable. After identifying
 158 the issues with prior metrics, we propose a method of identifying metastable states and measuring
 159 the agreement between the model and the ground truth.

160 **Deficiencies of Past Metrics** Past studies have used different tasks and metrics for evaluation,
 161 making it difficult to compare methods. The metastable states of the free energy surface are fre-
 162 quently used for evaluation as they allow to reason about dominant conformations and transition
 163 rates. However, previous evaluation protocols are often not applicable to multiple systems but only
 164 allow qualitative inspection of single molecules at a time. To overcome these challenges, we propose
 165 a systematic evaluation protocol to reliably assess the quality of predicted trajectories for multiple
 166 systems.

167 A common practice to evaluate the quality of predicted FES is to use Kullback-Leibler (KL) diver-
 168 gences, either between one-dimensional marginals or the two-dimensional histogram (Klein et al.,
 169 2023). However, this method is heavily dependent on the chosen bin size of the histogram and
 170 ignores the fact that variations in the estimated density are negligible for multiple practical applica-
 171 tions, where the correct identification of modes and transition rates is the desired goal.

172 Work on conformation generation (Jing et al., 2022; Zhu et al., 2023) is typically evaluated by
 173 computing the coverage of predicted structures (in terms of RMSD) and reporting precision and
 174 recall, i.e. the fraction of correctly predicted structures and the fraction of identified structures
 175 compared to MD. Similar to the KL-based metrics, this protocol does not capture whether modes
 176 and transition rates are correctly identified.

177 **Identifying metastable states** Identifying modes in a two-dimensional FES is highly non-trivial.
 178 While previous works used K-MEANS clustering to identify metastable states (Pandey et al., 2023;
 179 Jain & Stock, 2012), we found that K-MEANS frequently converges to incorrect minima. There-
 180 fore, we use the method of Novelli et al. (2022) where the FES is first smoothed using a Gaussian
 181 kernel and local minima are identified via running multiple BFGS solvers from random starting
 182 points. For a detailed explanation, we refer to subsection B.3. Lastly, the identification of reac-
 183 tion coordinates varies across past methods where multiple methods a sophisticated scheme such as
 184 Time-Independent-Component-Analysis (TICA) (Pérez-Hernández et al., 2013) to define the reac-
 185 tion coordinates from which the FES is constructed (Sidky et al., 2020; Klein et al., 2023). While
 186 TICA is useful for a variety of applications, it requires a Chapman-Kolmogorov test and manual
 187 inspection of the lag time to guarantee high-quality dimensionality reduction. Therefore, we use the
 188 two dihedral angles ϕ, ψ as they are known to capture the conformation space of peptides (Choud-
 189 huri, 2014).

190 **Metrics** With the above-described procedure, we can identify
 191 metastable states without the need of manual specification.
 192 This allows to compute precision and recall in terms of found
 193 metastable states, i.e. the fraction of correctly predicted modes
 194 and the percentage of modes found where a mode is considered
 195 correct if it lies within close proximity to the ground-truth MD
 196 mode ².

197 Furthermore, the transition rates between these identified
 198 metastable states are relevant for many applications, such as
 199 inferring relaxation times or reaction rates, and can be studied
 200 using a Markov State Model (MSM) (Bowman et al., 2014).
 201 Hence, an MSM can be fitted to predicted and MD trajectories,
 202 allowing to compare transition rates. Specifically, the
 203 Mean First Passage Times (MFPTs) (Hoel et al., 1986) can be
 204 computed which represent the expected times for a transition
 205 to happen from a predefined origin state to a target state. The
 206 relative error across the MFPTs for multiple molecules compared
 207 to MD then gives insight about the practical use of the
 208 predicted dynamical properties.

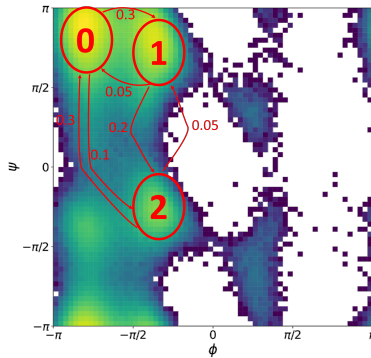


Figure 5: Example MSM with three states fitted to MD trajectory with transition probabilities.

209 5 Experimental Results

210 In this section, we first show LAMODY’s ability to recover the dynamics and transition states of
 211 alanine dipeptide, then show that it effectively generalizes across peptides. We further demonstrate
 212 the large benefits of LAMODY in terms of simulation speed in Appendix B. Finally, we do ablation
 213 studies on some of the architectural choices.

214 5.1 Alanine Dipeptide

215 Before we evaluate the generalization capabilities to unseen molecules, we test our method on a
 216 single molecule, namely alanine dipeptide (ALDP), which is a widely used benchmark for MD
 217 simulators and has been the subject of evaluation in previous works. In the case of ALDP, the
 218 primary degrees of freedom under consideration are the two backbone dihedral angles ϕ and ψ .
 219 Despite the model being trained on this exact molecule, it’s important to note that recovering long-
 220 time FES and transition rates remains highly nontrivial, as dynamical models are typically designed
 221 to predict single or a limited number of steps. Specifically, we train on 100ns of MD data of ALDP
 222 in implicit solvent to assess whether the model can qualitatively reproduce the free energy surface
 223 in terms of the backbone dihedral angles. Additionally, we analyze the model’s ability to predict
 224 transition rates between the identified metastable states, comparing them to MD results.

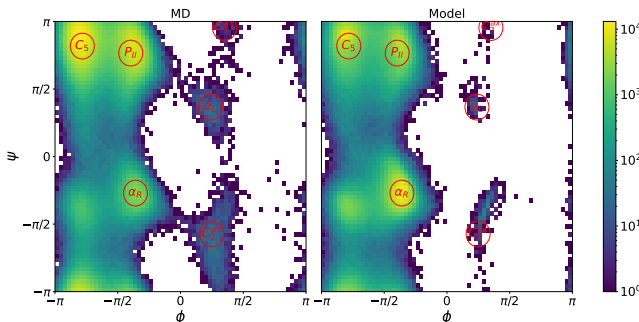


Figure 6: Ramachandran plots of trajectories from MD data and predictions of our model for alanine dipeptide with corresponding metastable states as defined by Vlachas et al. (2022).

225 **FES recovery** To use the trained model for simulating MD trajectories, we use the procedure de-
 226 scribed above. Starting from an initialization sequence of five states, we simulate a trajectory of

²See subsection B.3 and Figure 4

227 length $100ns$ without re-initialization. The Ramachandran plots of the predicted trajectory along-
 228 side the MD simulation are visualized in Figure 6. Figure 6 shows that our model is able to capture
 229 all metastable states without becoming unstable, i.e. no unphysical states are visited throughout the
 230 entire simulation. Notably, the model is able to explore the rare states C_7^{ax} , α_L , which previous
 231 latent space simulators (Vlachas et al., 2022) failed to achieve. The Ramachandran plots also show
 that our model slightly overestimates the density of α_R .

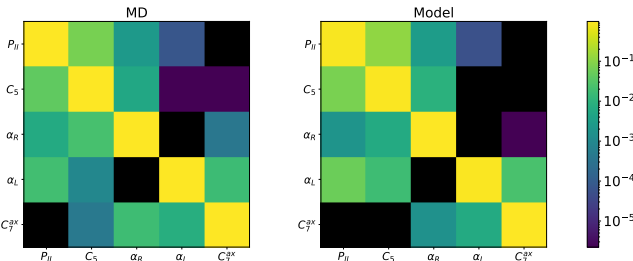


Figure 7: Transition probabilities of MSMs for alanine dipeptide estimated from MD data and predictions of our model. Black squares are transitions that were never observed.

232

233 **Transition dynamics** To examine whether the overestimation of α_R leads to unrealistic dynamical
 234 properties, we can compare the transition rates extracted from MSMs fitted to MD data as well as
 235 the predicted trajectory, which are shown in Figure 7. The transition probabilities clearly show that
 236 the dynamical properties that can be inferred from the model predictions closely match the true
 237 dynamics. Even for the highly unlikely states, our model approximates the correct transition rates.
 238 We found the training scheme for long trajectories as described above to be crucial for this.

239 5.2 Generalization across Molecules

240 After this first sanity check, we assess the capability of our approach to generalize to unseen
 241 molecules. To do so, we constructed a dataset of 216 dipeptides³ with a length of $12ns$ each of
 242 which 200 are used for training and 16 are held out for evaluation. We use the systematic evaluation
 243 protocol introduced in section 4.

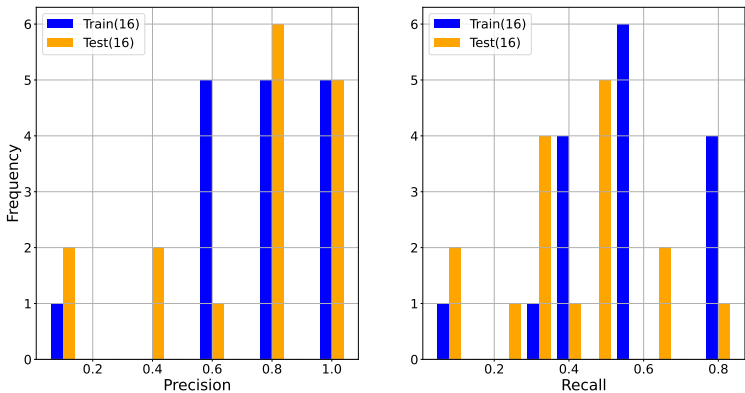


Figure 8: Metastable state precision and recall (MSPR) for train and test samples of the dipeptide model.

244 **FES recovery** In contrast to prior work on latent space simulators (Sidky et al., 2020; Vlachas
 245 et al., 2022) where the model can only be evaluated on the same molecule it has been trained on,
 246 our architecture is not restricted to single molecules. We evaluate the peptide model on 16 unseen
 247 molecules and randomly choose 16 peptides from the training set as a comparison. Figure 8 shows
 248 the precision and recall values the dipeptide model achieved. We can observe, that the model is better
 249 in terms of precision than recall. This suggests, that the learned simulator is more "conservative"

³Peptides with two amino acids

250 and avoids predicting unphysical modes rather than exploring the full state space which is desirable.
251 However, Figure 8 also shows that the model fails to recover the correct metastable states for a subset
of the peptides.

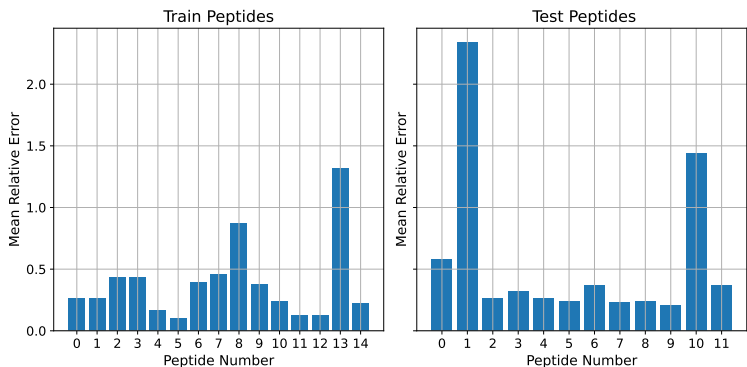


Figure 9: Mean relative error of MFPTs for MSMs fitted to predicted trajectories compared to MD for train and test set. Correctly extracted metastable states from the predicted trajectory are used to construct MSMs on MD and predicted data. Peptides where only one metastable state exists and therefore the MFPT error would always be zero are held out.

252

253 **Predicting transition dynamics** To gain more insight into the predicted trajectories, we evaluate
254 the relative error between predicted and MD MFPTs for MSMs constructed from correctly identified
255 states as defined in section 4. The results of this analysis are shown in Figure 9 where peptides that
256 only contain one mode are excluded, as the MFPT error would be 0 in this case (only one state in
257 the MSM, so no transitions). Figure 9 shows that the mean relative error is below 0.5 except for
258 two peptides from the training set and two peptides from the test set. This confirms the previous
259 results, i.e. that the model can approximate the majority of peptides very well, but misses a small
260 subset. Furthermore, this metric shows that the modes which are found by the model are captured
261 accurately and the transitions between the modes are captured within a relative error that existing
262 latent space simulators (Vlachas et al., 2022) achieve for a single molecule they have been trained
263 on. Furthermore, this shows the practical use of this method, as it can quickly and efficiently recover
264 the leading states of unseen molecules from which accurate transition rates can be extracted making
265 this model especially useful for screening large chemical spaces.

266 6 Discussion

267 We present MSPR, a reliable evaluation metric for FES that tackles the necessity of comparable
268 evaluation schemes for learned simulators. Additionally, we introduce LAMODY, a learned sim-
269 ulator operating in a latent space to efficiently recover free energy surfaces and transition rates.
270 LAMODY is trained end-to-end on MD data constructing its own latent space. The model employs
271 an $SE(3)$ -invariant encoder-propagator-decoder scheme. We show that our method can operate at
272 integration time steps that are two orders of magnitude larger than for MD while still being able
273 to conduct stable long-timescale simulations required for recovering properties such as FES and
274 transition rates.

275 In contrast to prior works, LAMODY does not require re-initialization throughout the simulation,
276 removing the need for prior MD simulations. We demonstrate that the predicted trajectories closely
277 match the results of MD and correct dynamical properties can be recovered even for rare metastable
278 states. Furthermore, our model is generalizable to molecules outside its training distribution and can
279 capture their leading structural and dynamical properties. Overall, our approach is approximately
280 20 times faster at recovering FES and transition rates than classical MD and can additionally easily
281 be parallelized for up to 128 peptides on a single GPU.

282 References

- 283 Simon Batzner, Albert Musaelian, Lixin Sun, Mario Geiger, Jonathan P. Mailoa, Mordechai Ko-
284 rnblyth, Nicola Molinari, Tess E. Smidt, and Boris Kozinsky. E(3)-equivariant graph neural
285 networks for data-efficient and accurate interatomic potentials. *Nature Communications*, 13
286 (1), May 2022. doi: 10.1038/s41467-022-29939-5. URL [https://doi.org/10.1038/
287 s41467-022-29939-5](https://doi.org/10.1038/s41467-022-29939-5).
- 288 Rafael C. Bernardi, Marcelo C.R. Melo, and Klaus Schulten. Enhanced sampling techniques in
289 molecular dynamics simulations of biological systems. *Biochimica et Biophysica Acta (BBA)*
290 - *General Subjects*, 1850(5):872–877, May 2015. doi: 10.1016/j.bbagen.2014.10.019. URL
291 <https://doi.org/10.1016/j.bbagen.2014.10.019>.
- 292 Christopher M Bishop. *Mixture density networks*. Aston University, 1994.
- 293 Art D. Bochevarov, Edward Harder, Thomas F. Hughes, Jeremy R. Greenwood, Dale A. Braden,
294 Dean M. Philipp, David Rinaldo, Mathew D. Halls, Jing Zhang, and Richard A. Friesner. Jaguar:
295 A high-performance quantum chemistry software program with strengths in life and materials
296 sciences. *International Journal of Quantum Chemistry*, 113(18):2110–2142, July 2013. doi:
297 10.1002/qua.24481. URL <https://doi.org/10.1002/qua.24481>.
- 298 Gregory R. Bowman, Vijay S. Pande, and Frank Noé. *An Introduction to Markov State Models and*
299 *Their Application to Long Timescale Molecular Simulation*, volume 797. Springer Netherlands,
300 Dordrecht, 2014. doi: 10.1007/978-94-007-7606-7.
- 301 Johannes Brandstetter, Daniel E. Worrall, and Max Welling. Message passing neural PDE solvers. In
302 *International Conference on Learning Representations, 2022*. URL [https://openreview.
303 net/forum?id=vSix3HPYKSU](https://openreview.net/forum?id=vSix3HPYKSU).
- 304 Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties
305 of neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth*
306 *Workshop on Syntax, Semantics and Structure in Statistical Translation*, pp. 103–111, Doha,
307 Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/W14-4012.
308 URL <https://aclanthology.org/W14-4012>.
- 309 Supratim Choudhuri. Additional bioinformatic analyses involving protein sequences. In *Bioinfor-*
310 *matics for Beginners*, pp. 183–207. Elsevier, 2014. doi: 10.1016/b978-0-12-410471-6.00008-6.
311 URL <https://doi.org/10.1016/b978-0-12-410471-6.00008-6>.
- 312 Peter Eastman, Jason Swails, John D. Chodera, Robert T. McGibbon, Yutong Zhao, Kyle A.
313 Beauchamp, Lee-Ping Wang, Andrew C. Simmonett, Matthew P. Harrigan, Chaya D. Stern,
314 Rafal P. Wiewiora, Bernard R. Brooks, and Vijay S. Pande. OpenMM 7: Rapid development
315 of high performance algorithms for molecular dynamics. *PLOS Computational Biology*, 13(7):
316 e1005659, July 2017. doi: 10.1371/journal.pcbi.1005659. URL [https://doi.org/10.
317 1371/journal.pcbi.1005659](https://doi.org/10.1371/journal.pcbi.1005659).
- 318 Matthias Fey and Jan Eric Lenssen. Fast Graph Representation Learning with PyTorch Geometric,
319 May 2019. URL https://github.com/pyg-team/pytorch_geometric.
- 320 Xiang Fu, Zhenghao Wu, Wujie Wang, Tian Xie, Sinan Ketten, Rafael Gomez-Bombarelli, and
321 Tommi S. Jaakkola. Forces are not enough: Benchmark and critical evaluation for machine
322 learning force fields with molecular simulations, 2023. URL [https://openreview.net/
323 forum?id=_V-nKeWvs7p](https://openreview.net/forum?id=_V-nKeWvs7p).
- 324 Johannes Gasteiger, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional
325 graph neural networks for molecules. In *Conference on Neural Information Processing Systems*
326 (*NeurIPS*), 2021.
- 327 Mario Geiger and Tess Smidt. e3nn: Euclidean neural networks, 2022. URL [https://arxiv.
328 org/abs/2207.09453](https://arxiv.org/abs/2207.09453).

- 329 Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neu-
330 ral message passing for quantum chemistry. In Doina Precup and Yee Whye Teh (eds.), *Pro-
331 ceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceed-
332 ings of Machine Learning Research*, pp. 1263–1272. PMLR, 06–11 Aug 2017. URL [https:
333 //proceedings.mlr.press/v70/gilmer17a.html](https://proceedings.mlr.press/v70/gilmer17a.html).
- 334 M.A. González. Force fields and molecular dynamics simulations. *École thématique de la Société
335 Française de la Neutronique*, 12:169–200, 2011. doi: 10.1051/sfn/201112009. URL [https:
336 //doi.org/10.1051/sfn/201112009](https://doi.org/10.1051/sfn/201112009).
- 337 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):
338 1735–1780, November 1997. doi: 10.1162/neco.1997.9.8.1735. URL [https://doi.org/
339 10.1162/neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735).
- 340 Paul G Hoel, Sidney C Port, and Charles J Stone. *Introduction to stochastic processes*. Waveland
341 Press, 1986.
- 342 Weihua Hu, Muhammed Shuaibi, Abhishek Das, Siddharth Goyal, Anuroop Sriram, Jure Leskovec,
343 Devi Parikh, and C. Lawrence Zitnick. Forcenet: A graph neural network for large-scale quantum
344 calculations, 2021. URL <https://arxiv.org/abs/2103.01436>.
- 345 Abhinav Jain and Gerhard Stock. Identifying metastable states of folding proteins. *Journal of
346 Chemical Theory and Computation*, 8(10):3810–3819, April 2012. doi: 10.1021/ct300077q. URL
347 <https://doi.org/10.1021/ct300077q>.
- 348 Bowen Jing, Gabriele Corso, Regina Barzilay, and Tommi S. Jaakkola. Torsional diffusion for
349 molecular conformer generation. In *ICLR2022 Machine Learning for Drug Discovery, 2022*.
350 URL <https://openreview.net/forum?id=D9IxPlXPJJS>.
- 351 Leon Klein, Andrew Y. K. Foong, Tor Erlend Fjelde, Bruno Mlodozieniec, Marc Brockschmidt,
352 Sebastian Nowozin, Frank Noé, and Ryota Tomioka. Timewarp: Transferable acceleration of
353 molecular dynamics by learning time-coarsened dynamics, 2023. URL [https://arxiv.
354 org/abs/2302.01170](https://arxiv.org/abs/2302.01170).
- 355 Alessandro Laio and Francesco L Gervasio. Metadynamics: a method to simulate rare events and
356 reconstruct the free energy in biophysics, chemistry and material science. *Rep. Prog. Phys.*, 71
357 (12):126601, December 2008.
- 358 Don S. Lemons and Anthony Gythiel. Paul langevin’s 1908 paper “on the theory of brownian mo-
359 tion” [“sur la théorie du mouvement brownien,” c. r. acad. sci. (paris) b146/b, 530–533 (1908)].
360 *American Journal of Physics*, 65(11):1079–1081, November 1997. doi: 10.1119/1.18725. URL
361 <https://doi.org/10.1119/1.18725>.
- 362 Bethany Lusch, J. Nathan Kutz, and Steven L. Brunton. Deep learning for universal linear em-
363 beddings of nonlinear dynamics. *Nature Communications*, 9(1), November 2018. doi: 10.1038/
364 s41467-018-07210-0. URL <https://doi.org/10.1038/s41467-018-07210-0>.
- 365 Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer New York, 2006. doi: 10.
366 1007/978-0-387-40065-5. URL <https://doi.org/10.1007/978-0-387-40065-5>.
- 367 Frank Noé, Simon Olsson, Jonas Köhler, and Hao Wu. Boltzmann generators: Sampling equilibrium
368 states of many-body systems with deep learning. *Science (New York, N.Y.)*, 365(6457), 2019. doi:
369 10.1126/science.aaw1147.
- 370 Pietro Novelli, Luigi Bonati, Massimiliano Pontil, and Michele Parrinello. Characterizing
371 metastable states with the help of machine learning. *Journal of Chemical Theory and Com-
372 putation*, 18(9):5195–5202, August 2022. doi: 10.1021/acs.jctc.2c00393. URL [https:
373 //doi.org/10.1021/acs.jctc.2c00393](https://doi.org/10.1021/acs.jctc.2c00393).
- 374 Samuel E. Otto and Clarence W. Rowley. Linearly recurrent autoencoder networks for learning
375 dynamics. *SIAM Journal on Applied Dynamical Systems*, 18(1):558–593, January 2019. doi:
376 10.1137/18m1177846. URL <https://doi.org/10.1137/18m1177846>.

377 Bhawna Pandey, Krishnendu Sinha, Aditya Dev, Himal K. Ganguly, Smarajit Polley, Suman
378 Chakrabarty, and Gautam Basu. Phosphorylation-competent metastable state of escherichia coli
379 toxin hipa. *Biochemistry*, 62(5):989–999, 2023. doi: 10.1021/acs.biochem.2c00614. URL
380 <https://doi.org/10.1021/acs.biochem.2c00614>. PMID: 36802529.

381 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
382 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf,
383 Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit
384 Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-
385 Performance Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché
386 Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp.
387 8024–8035. Curran Associates, Inc., 2019. URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
388 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
389 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).

390 Lili X. Peng, Monica T. Hsu, Massimiliano Bonomi, David A. Agard, and Matthew P. Jacobson.
391 The free energy profile of tubulin straight-bent conformational changes, with implications for mi-
392 crotubule assembly and drug discovery. *PLoS Computational Biology*, 10(2):e1003464, February
393 2014. doi: 10.1371/journal.pcbi.1003464. URL [https://doi.org/10.1371/journal.](https://doi.org/10.1371/journal.pcbi.1003464)
394 [pcbi.1003464](https://doi.org/10.1371/journal.pcbi.1003464).

395 Guillermo Pérez-Hernández, Fabian Paul, Toni Giorgino, Gianni De Fabritiis, and Frank Noé.
396 Identification of slow molecular order parameters for markov model construction. *The Jour-*
397 *nal of Chemical Physics*, 139(1):015102, July 2013. doi: 10.1063/1.4811489. URL [https:](https://doi.org/10.1063/1.4811489)
398 [//doi.org/10.1063/1.4811489](https://doi.org/10.1063/1.4811489).

399 Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Saucedo Felix, Stefan Chmiela, Alexan-
400 dre Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolutional
401 neural network for modeling quantum interactions. In I. Guyon, U. Von Luxburg,
402 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-*
403 *vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
404 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/](https://proceedings.neurips.cc/paper_files/paper/2017/file/303ed4c69846ab36c2904d3ba8573050-Paper.pdf)
405 [file/303ed4c69846ab36c2904d3ba8573050-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/303ed4c69846ab36c2904d3ba8573050-Paper.pdf).

406 David W. Scott. *Multivariate Density Estimation*. Wiley, 1992. doi: 10.1002/9780470316849. URL
407 <https://doi.org/10.1002/9780470316849>.

408 Yunsheng Shi, Zhengjie Huang, shikun feng, Hui Zhong, Wenjin Wang, and Yu Sun. Masked
409 label prediction: Unified message passing model for semi-supervised classification, 2021. URL
410 <https://openreview.net/forum?id=B9t708KMr9d>.

411 Hythem Sidky, Wei Chen, and Andrew L. Ferguson. Molecular latent space simulators. *Chemical*
412 *Science*, 11(35):9459–9467, 2020. doi: 10.1039/d0sc03635h. URL [https://doi.org/10.](https://doi.org/10.1039/d0sc03635h)
413 [1039/d0sc03635h](https://doi.org/10.1039/d0sc03635h).

414 Lukas S. Stelzl and Gerhard Hummer. Kinetics from replica exchange molecular dynamics sim-
415 ulations. *Journal of Chemical Theory and Computation*, 13(8):3927–3935, July 2017. doi:
416 [10.1021/acs.jctc.7b00372](https://doi.org/10.1021/acs.jctc.7b00372). URL <https://doi.org/10.1021/acs.jctc.7b00372>.

417 G.M. Torrie and J.P. Valleau. Nonphysical sampling distributions in monte carlo free-energy esti-
418 mation: Umbrella sampling. *Journal of Computational Physics*, 23(2):187–199, feb 1977. doi:
419 [10.1016/0021-9991\(77\)90121-8](https://doi.org/10.1016/0021-9991(77)90121-8). URL [https://doi.org/10.1016/0021-9991\(77\)90121-8](https://doi.org/10.1016/0021-9991(77)90121-8).
420 [2877%2990121-8](https://doi.org/10.1016/0021-9991(77)90121-8).

421 Constantino Tsallis and Daniel A. Stariolo. Generalized simulated annealing. *Physica A:*
422 *Statistical Mechanics and its Applications*, 233(1-2):395–406, November 1996. doi: 10.
423 [1016/s0378-4371\(96\)00271-3](https://doi.org/10.1016/s0378-4371(96)00271-3). URL [https://doi.org/10.1016/s0378-4371\(96\)](https://doi.org/10.1016/s0378-4371(96)00271-3)
424 [00271-3](https://doi.org/10.1016/s0378-4371(96)00271-3).

425 Oliver T. Unke, Stefan Chmiela, Huziel E. Saucedo, Michael Gastegger, Igor Poltavsky, Kristof T.
426 Schütt, Alexandre Tkatchenko, and Klaus-Robert Müller. Machine learning force fields. *Chem-*
427 *ical Reviews*, 121(16):10142–10186, March 2021. doi: 10.1021/acs.chemrev.0c01111. URL
428 <https://doi.org/10.1021/acs.chemrev.0c01111>.

- 429 Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for sets,
430 2016.
- 431 Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau,
432 Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der
433 Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson,
434 Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore,
435 Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero,
436 Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt,
437 and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing
438 in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- 439 Pantelis R. Vlachas, Julija Zavadlav, Matej Praprotnik, and Petros Koumoutsakos. Accelerated
440 simulations of molecular systems through learning of effective dynamics. *Journal of Chemical
441 Theory and Computation*, 18(1):538–549, 2022. doi: 10.1021/acs.jctc.1c00809. URL <https://doi.org/10.1021/acs.jctc.1c00809>. PMID: 34890204.
- 443 Dongdong Wang, Yanze Wang, Junhan Chang, Linfeng Zhang, Han Wang, and Weinan E. Efficient
444 sampling of high-dimensional free energy landscapes using adaptive reinforced dynamics. *Nature
445 Computational Science*, 2(1):20–29, December 2021. doi: 10.1038/s43588-021-00173-1. URL
446 <https://doi.org/10.1038/s43588-021-00173-1>.
- 447 Robin Winter, Frank Noé, and Djork-Arné Clevert. Auto-encoding molecular conformations, 2021.
448 URL <https://arxiv.org/abs/2101.01618>.
- 449 Neo Wu, Bradley Green, Xue Ben, and Shawn O’Banion. Deep transformer models for time series
450 forecasting: The influenza prevalence case, 2020. URL <https://arxiv.org/abs/2001.08317>.
- 452 Yi Isaac Yang, Qiang Shao, Jun Zhang, Lijiang Yang, and Yi Qin Gao. Enhanced sampling in
453 molecular dynamics. *The Journal of Chemical Physics*, 151(7):070902, August 2019. doi: 10.
454 1063/1.5109531. URL <https://doi.org/10.1063/1.5109531>.
- 455 Shuxin Zheng, Jiyan He, Chang Liu, Yu Shi, Ziheng Lu, Weitao Feng, Fusong Ju, Jiayi Wang,
456 Jianwei Zhu, Yaosen Min, He Zhang, Shidi Tang, Hongxia Hao, Peiran Jin, Chi Chen, Frank
457 Noé, Haiguang Liu, and Tie-Yan Liu. Towards predicting equilibrium distributions for molecular
458 systems with deep learning, 2023. URL <https://arxiv.org/abs/2306.05445>.
- 459 Jun-Jie Zhu, Ning-Jie Zhang, Ting Wei, and Hai-Feng Chen. Enhancing conformational sampling
460 for intrinsically disordered and ordered proteins by variational autoencoder. *International Journal
461 of Molecular Sciences*, 24(8):6896, April 2023. doi: 10.3390/ijms24086896. URL <https://doi.org/10.3390/ijms24086896>.

463 A Additional Explanations

464 A.1 Molecular Dynamics Simulation

465 Molecular Dynamics (MD) simulations are a computational tool that can be utilized to study the
466 behavior of molecules over time at an atomistic resolution. To do so, a popular method is Langevin
467 Dynamics (Lemons & Gythiel, 1997), which evolves the positions and velocities of the system under
468 study by the following stochastic differential equation:

$$m_i \frac{d^2 \mathbf{x}_i}{dt^2} = -\nabla_i U(\mathbf{x}_1, \dots, \mathbf{x}_N) - \gamma m_i \frac{d\mathbf{x}_i}{dt} + \sqrt{2m_i \gamma k_B T} dB_t \quad (4)$$

469 where \mathbf{x}_i denotes the position of atom i , U is the potential energy, γ is a friction constant, m_i is the
470 mass of atom i , T is the temperature of the system, k_B is the Boltzmann constant, and dB_t is standard
471 Brownian motion. To ensure the stability of the simulation, the integration time step size is typically
472 chosen to be in the range of a few femtoseconds. The potential energy of the molecule based on
473 the coordinates of the particles $U(\mathbf{x}_1, \dots, \mathbf{x}_N)$ is usually parameterized by a force field⁴. Machine

⁴see González (2011) for a detailed definition.

474 learning methods that aim to simulate molecular systems are normally evaluated by their ability to
 475 recover conformational modes, free energy surfaces, and dynamical properties in comparison to a
 476 classical MD simulation (Vlachas et al., 2022; Sidky et al., 2020; Klein et al., 2023).

477 A.2 Internal Coordinate Graph

478 Figure 10 shows the a visualization of the internal coordinate graph used by the encoder as defined
 479 in subsection 3.1.

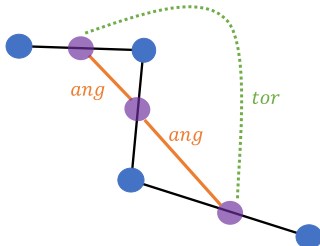


Figure 10: Graph of internal coordinates superimposed onto the molecular graph. Blue vertices and black edges show the corresponding molecular graph. The internal graph is superimposed with bond vertices in purple, bond angle edges in orange, and torsion angle edges in green.

480 B Additional Results

481 B.1 Simulation Speed

482 As high computational complexity/ slow simulation speed is the major limitation of MD simulations
 483 Table 1 shows the propagation speed of our method and MD in terms of iterations per second and
 484 the total wallclock time the respective simulation requires⁵. Table 1 clearly shows the advantage
 485 of our method that realizes a speedup of approximately 20, improving upon the results of Vlachas
 486 et al. (2022), who reported an acceleration by a factor of 3. Furthermore, in contrast to prior work,
 487 our model does not require re-initialization paired with short timescale predictions but can instead
 488 simulate long timescale trajectories starting from a five-state sequence without becoming unstable.
 489 Note that the predictions of our model can also be run in parallel with up to 128 peptides on a single
 490 GPU.

Table 1: Simulation Speed of MD and LAMODY given as averaged iterations per second and total wallclock times.

Molecule	iteration/second		wallclock time [minute]	
	MD	LAMODY	MD	LAMODY
ALDP	189	3788	88	4.6
Peptides	117	2239	34.2	1.8

491 B.2 model variations and ablations

492 **Cartesian Encoders** As the natural choice for an input representation seems to be representing a
 493 state by the two-dimensional molecular graph and associated cartesian positions, we also employed
 494 an $SE(3)$ -invariant encoder operating on cartesian coordinates based on Euclidean graph neural
 495 networks (Geiger & Smidt, 2022). Additionally, we also used the popular GEMNET (Gasteiger
 496 et al., 2021) as our encoder network since GEMNET operates on cartesian coordinates and uses the
 497 internal coordinates of a molecule as features during message passing. However, we unexpectedly
 498 encountered that the cartesian encoder as well as GemNet failed to identify rare metastable states.
 499 The results of these simulations are shown in Figure 11 and Figure 12. We suspect this to be the

⁵Hardware specifications are reported in Appendix F

500 case as both models are more memory intense than the internal encoder and we, therefore, had to
 501 reduce the length over which we unroll the propagator states during training ⁶.

502 B.3 Identification of metastable states

503 Following Novelli et al. (2022), we use a standard Gaussian kernel density estimator (Scott, 1992)
 504 to approximate the free energy surface in the space of the two dihedral angles ϕ, ψ that are known
 505 to capture the conformational space for peptides (Choudhuri, 2014). Then we aim to identify the
 506 local minima of the FES as these will represent the metastable states. To do so, 100 BFGS solvers
 507 (Nocedal & Wright, 2006) are initialized at random points and run until convergence from which
 508 we recover the unique local minima. By doing so, we are able to reliably identify metastable states
 509 without the need for manual specification ⁷.

510 To assess the quality of our predictions, we apply this procedure to the trajectories produced by
 511 our model as well as the MD data. This allows to compute precision and recall of the metastable
 512 states extracted from the predicted trajectories where we consider a metastable state to be correctly
 513 identified if $|\mu_{pred} - \mu_{MD}| \leq 0.15$. This allows us to judge the models' ability to recover correct
 514 FES for multiple peptides. Additionally, we use the set of correctly identified metastable states (from
 515 our model predictions) to construct an MSM for which we can compare the mean first passage times
 516 (MFPT) (Hoel et al., 1986) between MD and our model. The MFPTs are the expected time for a
 517 transition to happen from a predefined origin state to a target state. In practical applications this
 518 property is of great interest and can, for instance, be used to estimate the time it takes for a molecule
 519 to bind to a receptor. With this evaluation metric, we can judge the quality of the predicted dynamics
 520 and the practical use of the model, even if the model did not find all metastable states.

521 B.4 Model Variations and Ablations

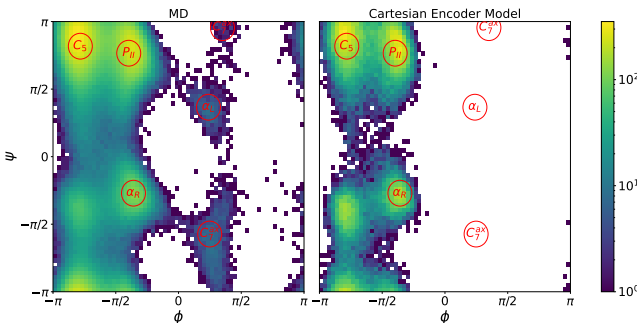


Figure 11: Ramachandran plots of trajectories from MD data and predictions of the model with cartesian encoder based on tensor product convolutions (Geiger & Smidt, 2022).

522 Figure 11 and Figure 12 show the inference results for the models with a cartesian/GEMNET en-
 523 coder respectively. The figures show that both models miss the rare metastable states, which we
 524 suspect to be caused by the shorter training sequences due to memory limitations as described in
 525 subsection B.2.

526 C Architecture Details

527 C.1 Encoder

528 The internal encoder operates on the internal coordinate graph as described in subsection 3.1, which
 529 is $SE(3)$ -invariant by construction. The internal coordinates are normalized to lie in $[0, 1]$.
 530 Nodes v_i are featurized with: Atomic number of the first atom in the bond, atomic number of the
 531 second atom in the bond, bond length, mass of the first atom, and mass of the second atom. Edges

⁶Unrolling propagator states for long trajectories with detaching gradients, see subsection 3.2 for details.

⁷An example is shown in Figure 4

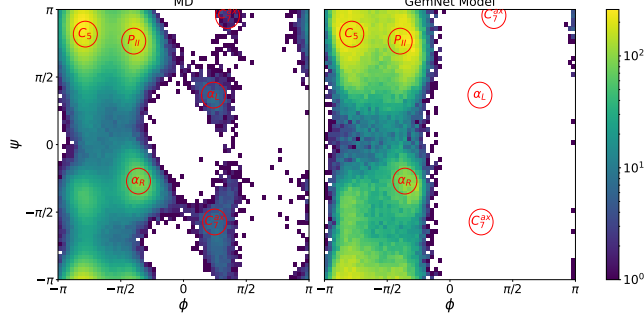


Figure 12: Ramachandran plots of trajectories from MD data and predictions of the model with GEMNET (Gasteiger et al., 2021) encoder.

532 between all pairs of bonds that form a bond angle are featurized with the bond angle and an addi-
 533 tional categorical feature indicating the edge type. Torsional edges are featurized with the torsion
 534 angle and the categorical feature accordingly. These scalar features are transformed by a set of
 535 learnable MLPs (one for each feature), to compute an initial feature embedding \mathbf{h}^0 for each node.
 536 After computing the initial embeddings \mathbf{h}_i^0 , we iteratively apply L message passing layers that addi-
 537 tionally employ a (multi-head) dot product attention mechanism to scale messages according to
 538 their importance, akin to Shi et al. (2021). More specifically, node embeddings for a node a at layer
 539 l get updated by:

$$\mathbf{h}_a^{l+1} = \beta_a \mathbf{W}_1 \mathbf{h}_a^l + (1 - \beta_a) \underbrace{\left(\sum_{b \in \mathcal{N}(a)} \alpha_{ab} (\mathbf{W}_2 \mathbf{h}_b^l + \mathbf{W}_6 \mathbf{c}_{ab}) \right)}_{m_a} \quad (5)$$

540 with

$$\alpha_{ab} = \text{softmax} \left(\frac{(\mathbf{W}_3 \mathbf{h}_a^l)^T (\mathbf{W}_4 \mathbf{h}_b^l + \mathbf{W}_6 \mathbf{c}_{ab})}{\sqrt{d}} \right) \quad (6)$$

$$\beta_a = \text{sigmoid} (\mathbf{W}_5 [\mathbf{W}_1 \mathbf{h}_a^l, m_a, \mathbf{W}_1 \mathbf{h}_a^l - m_a])$$

541 here W_* indicates learnable parameters, d is the hidden size of the attention heads, $[a, b]$ indicates
 542 vector concatenation, $\mathbf{c}_{ab} \in \mathcal{C}$ are the edge features of edge (a, b) , and $\mathcal{N}(a) = \{b | (a, b) \in \mathcal{B} \vee$
 543 $(b, a) \in \mathcal{B}\}$. Between each of the layers, ELU nonlinearities and batch normalization are applied.
 544 After the final message passing layer, we use a learnable set-to-set mapping Vinyals et al. (2016) to
 545 pool the nodes:

$$\begin{aligned} \mathbf{q}_t &= LSTM(\mathbf{q}_{t-1}^*) \\ e_{i,t} &= \mathbf{h}_i^L \cdot \mathbf{q}_t \\ \alpha_{i,t} &= \frac{\exp(e_{i,t})}{\sum_j \exp(e_{j,t})} \\ \mathbf{r}_t &= \sum_{i=1}^N \alpha_{i,t} \mathbf{h}_i^L \\ \mathbf{q}_t^* &= [\mathbf{q}_t, \mathbf{r}_t] \end{aligned} \quad (7)$$

546 where \cdot denotes the dot product and \mathbf{h}_i^L indicates the node embedding after the final message passing
 547 interaction layer. This layer iteratively updates the aggregated set for T processing steps by comput-
 548 ing a weighted sum \mathbf{r}_t of node embeddings, concatenating this sum to the last state \mathbf{q}_t and passing
 549 this concatenated vector \mathbf{q}^* through the LSTM. We found this learnable set-to-set mapping to yield
 550 better results compared to sum or mean reduction. After the set-to-set aggregation, we use a linear
 551 layer Φ to map to the fixed-size latent embedding vector:

$$\mathbf{z} = \Phi(\mathbf{q}_T^*) \quad (8)$$

552 Given this model architecture, we are able to learn a mapping to a latent space, which is by con-
 553 struction of the graph $SE(3)$ -invariant. Moreover, the model is not limited to a fixed-size graph but
 554 can be applied to graphs of distinct molecules.

555 C.2 Decoder

556 The molecular decoder acts as a counterpart to the encoder and reconstructs a molecular state from
 557 a latent representation by predicting the molecule’s internal coordinates for that state. The decoder
 558 architecture was heavily inspired by the work of Winter et al. (2021). As the decoder has to be
 559 applicable to different molecules, we condition the decoder on the time-invariant two-dimensional
 560 molecular graph. Concretely, the decoder predicts a molecular state at time t via:

$$\mathcal{G}_t = \mathcal{D}(z_t, \mathcal{G}_{mol}) \quad (9)$$

561 To do so, we first compute node embeddings for all atoms of $\mathcal{G}_{mol} \in (\mathcal{V}_{mol}, \mathcal{B}_{mol}, \mathcal{X}_{mol}, \mathcal{C}_{mol})$
 562 where nodes represent atoms and edges represent bonds between atoms in the molecule. \mathcal{G}_{mol} is
 563 constant throughout and MD simulation, as only the atom position change. We featurize nodes with
 564 the following attributes: Atomic number, chirality, degree, number of rings the atom is involved
 565 in, implicit valence, formal charge, number of bonded hydrogens, hybridization type, whether or
 566 not it is in an aromatic ring, whether or not it is in a 5 or 6-ring, the residue name and the atom
 567 name. Bonds between atoms are featurized by bond type and a radial basis embedding of the bond
 568 length (Schütt et al., 2017). Since torsion angles are defined by quadruplets of atoms that do not
 569 necessarily have to be direct neighbors, we add additional edges by connecting each node to all its
 570 k -hop neighbors. Concretely, we modify \mathcal{B}_{mol} to be $\mathcal{B}_{mol} := \{(a, b) \mid a \in \mathcal{V}_{mol} \wedge b \in \mathcal{N}^k(a)\}$
 571 where $\mathcal{N}^k(a)$ denotes all nodes that can be reached with a maximum of k hops from a . The
 572 additional edges facilitate the information flow over longer distances during message passing.
 573

574 After an initial node embedding akin to subsection C.1, we apply L message passing layers that
 575 update the node embeddings similar to subsection C.1. With the final node embeddings \mathbf{h}_i^L , we
 576 predict the internal coordinates of the current state by:

$$\begin{aligned} d_{ab}^t &= \Pi_{bond}([\mathbf{h}_a^L, \mathbf{h}_b^L, z_t]) \forall (a, b) \in \mathbb{B} \\ \phi_{abc}^t &= \Pi_{ang}([\mathbf{h}_a^L, \mathbf{h}_b^L, \mathbf{h}_c^L, z_t]) \forall (a, b, c) \in \mathbb{A} \\ \cos\psi_{abcd}^t &= \Pi_{tor_{cos}}([\mathbf{h}_a^L, \mathbf{h}_b^L, \mathbf{h}_c^L, \mathbf{h}_d^L, z_t]) \forall (a, b, c, d) \in \mathbb{T} \\ \sin\psi_{abcd}^t &= \Pi_{tor_{sin}}([\mathbf{h}_a^L, \mathbf{h}_b^L, \mathbf{h}_c^L, \mathbf{h}_d^L, z_t]) \forall (a, b, c, d) \in \mathbb{T} \end{aligned} \quad (10)$$

577 where Π_* are two-layer MLPs with ELU activations and dropout that map from the concatenated
 578 node embeddings and latent state to the single scalar of interest. \mathbb{B} denotes the set of all pairs of
 579 atoms defining a bond, \mathbb{A} is the set of all triplets of atoms defining a bond angle, and \mathbb{T} is the set of
 580 all quadruplets of atoms defining a torsion angle. Note that the decoder outputs a prediction for the
 581 bond angles directly, while for the torsion angles, sin and cos are predicted. This design choice is
 582 grounded on the fact that the models’ parameters could not be optimized to decode the full space of
 583 torsion angles when predicting them directly.
 584

585 D Training and Inference

586 We define the reconstruction loss in terms of internal coordinates by:

$$\begin{aligned} \mathcal{L}_{rec}(\mathcal{G}_i, \hat{\mathcal{G}}_i) &= \xi_b \frac{1}{|\mathbb{B}|} \sum_{(a,b) \in \mathbb{B}} \|d_{ab} - \hat{d}_{ab}\| \\ &+ \xi_a \frac{1}{|\mathbb{A}|} \sum_{(a,b,c) \in \mathbb{A}} \cos(\phi_{abc} - \hat{\phi}_{abc}) \\ &+ \xi_t \frac{1}{2|\mathbb{T}|} \sum_{(a,b,c,d) \in \mathbb{T}} \left(\cos(\psi_{abcd}) - \cos\hat{\psi}_{abcd} \right)^2 + \left(\sin(\psi_{abcd}) - \sin\hat{\psi}_{abcd} \right)^2 \end{aligned} \quad (11)$$

587 where ξ_b, ξ_a, ξ_t are hyperparameters, \mathbb{B} denotes the set of all pairs of atoms defining a bond, \mathbb{A} is
 588 the set of all triplets of atoms defining a bond angle, and \mathbb{T} is the set of all quadruplets of atoms
 589 defining a torsion angle. Note that as described in subsection C.2, the model predicts the bond
 590 angles directly, whereas, for the torsion angles, it predicts $\sin(\psi)$ and $\cos(\psi)$.

591

592 To infer σ^2 , i.e. the amount of noise added during inference, we found that the required noise level
 593 strongly correlates with the variance of the (normalized) torsion angles in the warmup trajectory.
 594 We identified a relationship of

$$\sigma^2 = \frac{1}{|\mathbb{T}|} \sum_{i=1}^{|\mathbb{T}|} \text{Var}(\psi_i) \quad (12)$$

595 to reliably give a good estimate of the noise level with $|\mathbb{T}|$ being the number of torsion for the
 596 respective molecule. While this relationship holds across molecules, we used a noise level of
 597 $\sigma_i^2 = 6 * \text{Var}(\psi_i)$ for the alanine dipeptide model where the factor of six was inferred from the
 598 norm of the latent space.

599

600 E Dataset Details

601 All datasets were created by performing MD simulations using the *openmm* library (Eastman et al.,
 602 2017).

603 The simulation was performed with the parameters shown in Table 2 and Figure 13 shows the free
 604 energy surface based on the two backbone dihedral angles (ϕ, ψ) of alanine dipeptide in implicit
 605 solvation. Given the distribution of (ϕ, ψ) , the free energy surface can be computed by:

$$FES_i = -k_B T \ln [p(\phi_i, \psi_i)] \quad (13)$$

606 where k_B is the Boltzmann constant and T is the temperature of the system. We can ob-
 607 serve five energetically favorable metastable states $\{P_{II}, C_7^{ax}, C_5, \alpha_R, \alpha_L\}$ which we also refer to
 608 as modes of the Boltzmann distribution. Note that the metastable states $\{C_7^{ax}, \alpha_L\}$ are visited rarely.

609

Table 2: Alanine dipeptide dataset properties.

Property	Value
Simulation time	100ns
Integrator	Langevin
Integrator time step	1fs
Forcefield	AMBER ff96
Solvation	OBC GBSA implicit
Frame Spacing	100fs
Temperature	300K

The dipeptide dataset was created with the simulation parameters given in Table 3.

Table 3: Dipeptide dataset properties.

Property	Value
# Peptides	216
Simulation time (each)	12ns
Integrator	Langevin
Integrator time step	1fs
Forcefield	AMBER 14-all
Solvation	implicit GBn
Frame Spacing	120fs
Temperature	300K

610

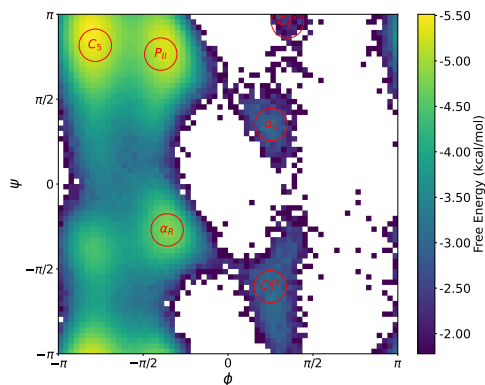


Figure 13: Ramachandran plot of the two backbone dihedral angles of the alanine dipeptide dataset with parameters from Table 2 and metastable states $\{P_{II}, C_7^{ax}, C_5, \alpha_R, \alpha_L\}$ as defined by Vlachas et al. (2022).

611 F Implementation details

612 All experiments were implemented in *PyTorch* (Paszke et al., 2019) using the extension for deep
 613 learning on graphs *Pytorch Geometric* (Fey & Lenssen, 2019). Furthermore, the *scipy* library (Vir-
 614 tanen et al., 2020) is extensively used throughout our implementation and we utilized the *stateinter-*
 615 *preter* package (Novelli et al., 2022) to automatically identify metastable states.

616 The experiments were run on two different machines. All training was run on a machine with two
 617 AMD EPYC 7513 CPU @ 2.60GHz with 32/64 cores each, 504GB of RAM, and eight NVIDIA
 618 RTX A6000 GPUs with 48GB vRam of which only a single one was used at a time. All inference
 619 experiments were performed on a machine with two Intel(R) Xeon(R) Gold 6230 CPU @ 2.10GHz
 620 with 20/40 cores each, 504GB of RAM, and eight NVIDIA Tesla V100 GPUs with 32GB vRam
 621 where again only a single GPU was used at a time.

622 G Additional Model Variations

623 **Dynamical Propagator** We found the LSTM architecture to consistently achieve the best simulation
 624 metrics outperforming the following architectures: Gated Recurrent Unit (GRU) (Cho et al., 2014);
 625 MLP; Mixture Density Network (Bishop, 1994); Transformer for time series forecasting (Wu et al.,
 626 2020). Besides the different architectures, we evaluated if conditioning the dynamical model onto
 627 the molecule it currently works with improves the generalization capabilities of our model. To do
 628 so, we employed another GNN that computes a fixed-size embedding based on the two-dimensional
 629 molecular graph, essentially constructing a learned representation of a certain molecule. This rep-
 630 resentation was then appended to the latent space to facilitate the prediction of correct dynamics for
 631 the propagator. However, we did not encounter any benefits of using this approach in terms of the
 632 quality of predicted trajectories for varying molecules.

633 **Training Schemes** Besides the training scheme described in subsection 3.2, we explored various
 634 methods of improving the robustness of the dynamical model mainly inspired by the approaches
 635 of Brandstetter et al. (2022). The model always gets correct latent states as input at training time
 636 whereas at inference time the propagator gets its own previous prediction as input which introduces
 637 a distribution shift between training and inference time. To mitigate this error, Brandstetter et al.
 638 (2022) suggest the "pushforward trick" which means to instead of using the correct latent state as
 639 input, the previous prediction of the dynamical model is used with a certain probability. Addition-
 640 ally, we tested whether infusing noise at different stages of our pipeline (in cartesian space; in in-
 641 ternal coordinate space; in the latent space) improves the test performance of our dynamical model.
 642 While the above two approaches did not improve the simulation results, we found the approach
 643 of unrolling the LSTM for multiples of its sequence length and cutting the gradients between the
 644 steps as described in subsection 3.2 to be absolutely crucial for the model to learn correct long-term
 645 dynamics.

646 **Pretraining the autoencoder** In contrast to the results of Sidky et al. (2020), we found that pre-
 647 training the autoencoder did not improve simulation results but in fact significantly constrained the
 648 latent space such that dynamical properties could not be modeled precisely anymore.

649 H Hyperparameters

650 For all training, we use the *Adam*⁸ optimizer and the *ReduceLRonPlateau*⁹ learning rate scheduler
 651 with reduction parameter 0.7 and patience 5 epochs. We define an epoch to consist of 12 batches
 652 of trajectories with length T for alanine dipeptide and 16 batches for the peptide models and train
 653 each model for 100 epochs, as we found all training metrics to have fully converged after that time.
 654

655 Training the smaller model on alanine dipeptide took 14.6 hours with a memory consumption of
 656 8.9GB. During inference, the memory consumption was 6B, which is mainly caused by the batched
 657 decoding of structures where we used batches of size $1e5$ and which could be adapted to other hard-
 658 ware limitations. For the dipeptide models, training took approximately three days with a memory
 659 consumption of 43GB. For decoding, we used a batch size of $1e4$, which led to 14GB of used GPU
 660 memory.

661 H.1 Alanine Dipeptide Hyperparameters

662 The parameters were tuned in the order in which they appear in the table from top to bottom. The
 663 final parameters are marked in **bold**.

664 We found the batch size to have a significant impact on the performance of our model, as batches
 665 larger than 8 independent trajectories prevented the models to produce reasonable inference results.
 666 While we do not have concrete evidence, we suspect this to be the case because batches larger than
 667 8 contain too diverse trajectories, essentially impeding the computation of meaningful gradients.

Table 4: Search space for the general hyperparameters, spanning across encoder, decoder and prop-
 agator.

Parameter	Search Space
latent embedding dimension	[5, 10, 32, 64, 75, 100, 128, 256 , 512]
data normalization	[min-max , z-score]
batch size	[2, 4, 8 , 16, 32, 64]
starting learning rate	[1e-3, 5e-4, 1e-4 , 1e-5, 1e-6]
c	[1, 2, 5, 10, 25, 50, 100, 120 , 150, 200]
$\delta_{rec}, \delta_{lat}, \delta_{e2e}, \xi_b, \xi_a, \xi_t$	[0.33, 1 , 2] (independently altered)

668

Table 5: Search space for the hyperparameters of the encoder network.

Parameter	Search Space
# layers	[2, 3, 4, 5 , 6, 7, 8, 10]
# final MLP layers	[1 , 2, 3, 4]
# attention heads	[2, 4, 8 , 16]
node embedding size	[5, 10 , 15, 25]
edge embedding size	[2 , 4, 8, 12]
# readout function	[Set2Set , Sum, Mean]
dropout	[0 , 0.1, 0.15, 0.2]

⁸<https://pytorch.org/docs/stable/generated/torch.optim.Adam.html>

⁹https://pytorch.org/docs/stable/generated/torch.optim.lr_scheduler.ReduceLRonPlateau.html

Table 6: Search space for the hyperparameters of the decoder network.

Parameter	Search Space
# MP layers	[1, 2, 3, 4, 5 , 6, 7, 8, 10]
k-hop edge concatenation	[2 , 3, 4]
# attention heads	[2, 4 , 8, 16]
input node embedding size	[5, 10 , 15, 25]
output node embedding size	[10, 15, 25 , 50, 100]
# final MLP layers	[1, 2, 3 , 4]
dropout MP layers	[0 , 0.1, 0.15]
dropout MLP layers	[0, 0.1 , 0.15]

Table 7: Search space for the hyperparameters of the LSTM propagator.

Parameter	Search Space
k (sequence length)	[1, 3, 5 , 10, 25, 50, 100, 250]
# LSTM layers	[1, 2, 3 , 4, 5, 6]
# MLP layers	[1, 2 , 3]
LSTM dropout	[0, 0.1 , 0.2]
β	0.15

669 **H.2 Dipeptide Hyperparameters**

670 For the training of the peptide models, we identified a batch size of 64 to yield the best results.

Table 8: Search space for the hyperparameters of the dipeptides model. All hyperparameters that are not explicitly listed are the same as for the alanine dipeptide model.

Parameter	Search Space
latent embedding dimension	[128, 256, 512, 1024 , 2048]
# num encoder layers	[4, 5, 6, 8 , 10]
# num decoder layers	[4, 5, 6, 8 , 10]
# LSTM layers	[4, 5 , 6, 8]
c	[1, 2, 5, 10, 25, 50, 100 , 120, 150, 200]
decoder output node embedding size	[10, 15, 25, 50 , 100]
β	0.9