# Scalable Bayesian Uncertainty Quantification for Neural Network Potentials: Promise and Pitfalls

Stephan Thaler,[*,†] Gregor Doehner,[†] and Julija Zavadlav[*,†,‡]

†*Professorship of Multiscale Modeling of Fluid Materials, Department of Engineering Physics and Computation, TUM School of Engineering and Design, Technical University of Munich, Germany*

‡*Munich Data Science Institute, Technical University of Munich, Germany*

E-mail: stephan.thaler@tum.de; julija.zavadlav@tum.de

**Abstract**

Neural network (NN) potentials promise highly accurate molecular dynamics (MD) simulations within the computational complexity of classical MD force fields. However, when applied outside their training domain, NN potential predictions can be inaccurate, increasing the need for Uncertainty Quantification (UQ). Bayesian modeling provides the mathematical framework for UQ, but classical Bayesian methods based on Markov chain Monte Carlo (MCMC) are computationally intractable for NN potentials. By training graph NN potentials for coarse-grained systems of liquid water and alanine dipeptide, we demonstrate here that scalable Bayesian UQ via stochastic gradient MCMC (SG-MCMC) yields reliable uncertainty estimates for MD observables. We show that cold posteriors can reduce the required training data size and that for reliable UQ, multiple Markov chains are needed. Additionally, we find that SG-MCMC and the Deep Ensemble method achieve comparable results, despite shorter training and less hyperparameter tuning of the latter. We show that both methods can capture aleatoric and epistemic uncertainty reliably, but not systematic uncertainty, which

needs to be minimized by adequate modeling to obtain accurate credible intervals for MD observables. Our results represent a step towards accurate UQ that is of vital importance for trustworthy NN potential-based MD simulations required for decision-making in practice.

# 1 Introduction

Molecular dynamics (MD) simulations are the computational tool of choice to describe complex molecular phenomena. Their computational effort and accuracy depend on the chosen potential energy model. Neural network (NN) potentials,[1–7] which model many-body interactions,[8,9] promise MD simulations at ab initio accuracy[4,10] within the computational complexity of classical molecular mechanics force fields.

The quality of NN potentials is limited by the scarcity of suitable training data,[11] given that data generation via computational quantum mechanics simulations and/or experiments is resource intensive. Hence, potentials are commonly applied outside their training domain due to the high-dimensional chemical space. As NN potentials are data-driven black box models, predictions outside the training domain may be inaccurate or even unphysical.[12–14] This may hinder more widespread adoption of NN potentials in practical applications where less powerful but physically more constrained models are preferred.[15]

Uncertainty quantification (UQ) can provide a remedy as it enables practitioners to quantify the trustworthiness of MD simulation predictions.[16–18] Additionally, the availability of a UQ metric enables more efficient training data generation via active learning,[14,19–23] as well as an adaptive combination of NN potentials with established classical force fields.[24] Bayesian statistics provides a mathematically rigorous approach to UQ. However, classical Bayesian inference schemes based on Markov Chain Monte Carlo (MCMC), such as Hamiltonian (or hybrid[25]) Monte Carlo (HMC),[26] require an evaluation of the likelihood over the whole data set for each parameter update. Frequent full likelihood evaluations are prohibitively expensive for computationally demanding NNs and large data sets.[27] Stochastic

gradient MCMC (SG-MCMC) schemes[28–32] enable scalable Bayesian UQ of NNs by computing stochastic estimates of the gradient of the likelihood on a mini-batch of data. Stochastic variational inference[33,34] represents another scalable Bayesian UQ method, while the Deep Ensemble[35,36] method is a popular non-Bayesian[15,36–38] alternative.

In the context of NN potentials, the Deep Ensemble method is in fact the most common UQ scheme,[5,24,27] but Dropout Monte Carlo[39] and last-layer Gaussian Mixture Models[40] have also been applied. In view of the poor performance of the Deep Ensemble method in an active learning context, Kahle and Zipoli[27] recently hypothesized that Bayesian approaches may provide more reliable uncertainty estimates for NN potentials. However, a comprehensive assessment of Bayesian UQ in the context of NN potentials is still outstanding.

In this work, we investigate scalable Bayesian UQ of MD observables for simulations utilizing NN molecular models. To this end, we first compare the UQ quality of a SG-MCMC method to the popular Deep Ensemble method and a gold-standard[15,31,32,41] HMC sampler based on a Lennard Jones (LJ) system with a 2-body toy NN potential. We then extend the comparison by learning graph NN potentials for coarse-grained (CG) systems of water and alanine dipeptide, demonstrating the practical applicability of SG-MCMC methods to fully-Bayesian modeling of graph NN potentials. Additionally, we investigate the influence of so-called cold posteriors[38] and the number of MCMC chains on the quality of Bayesian UQ. Finally, we advocate distinguishing between different sources of uncertainty; in particular, we highlight the importance of minimizing systematic uncertainties to obtain reliable credible intervals of MD observables.

## 2 Methods

In the following, we briefly summarize the employed SG-MCMC sampler as well as the Deep Ensemble method and continue with an outline of the Bayesian molecular modeling problem considered in this work.

## 2.1 Sources of uncertainty

The uncertainty in physical modeling can be divided into aleatoric, epistemic and systematic uncertainty.[15] Aleatoric uncertainty refers to the inherent stochastic nature of the modeled process, which can be interpreted as randomness in the labels $\mathbf{y}$ for a given input $\mathbf{x}$.[41,42] Epistemic uncertainty refers to the uncertainty about the true hypothesis (model) within the considered hypothesis space (model family). In contrast to aleatoric uncertainty, epistemic uncertainty can be reduced by gathering more data. Finally, systematic uncertainty is caused by model misspecification, i.e., when the true data-generating process is not contained within the hypothesis space. Systematic uncertainty manifests itself in an inconsistency between the data and the hypothesis space.[42]

## 2.2 Bayesian Modeling

A probabilistic model $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$ predicts the distribution of $\mathbf{y}$ reflecting the aleatoric uncertainty, given a training data set $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{N}$ of size $N$. Bayesian UQ additionally aims to quantify the epistemic uncertainty resulting from the model fit to a finite amount of data.[15,37] Instead of selecting a single set of model parameters $\boldsymbol{\theta}$, the Bayesian approach promises more robust predictions by marginalizing over $\boldsymbol{\theta}$.[43] Hence, the goal of Bayesian UQ is to compute the posterior predictive distribution

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \ , \tag{1}$$

where $p(\boldsymbol{\theta}|\mathcal{D})$ is the posterior distribution. The integral in eq. (1) is typically approximated by the Monte Carlo method:

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) \approx \frac{1}{N}\sum_{n=1}^{N} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_n) \ , \tag{2}$$

where $\boldsymbol{\theta}_n$ represents the $n^{\text{th}}$ model parameter set drawn from the posterior. Evaluating eq. (2) requires sampling from the posterior distribution

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} \propto \exp\left(\frac{-\mathcal{U}(\boldsymbol{\theta})}{\mathcal{T}}\right) , \tag{3}$$

with likelihood $p(\mathcal{D}|\boldsymbol{\theta})$ and prior $p(\boldsymbol{\theta})$. In analogy to statistical mechanics, the posterior can be re-written to allow sampling from a Boltzmann-type distribution,[26] with posterior potential energy

$$\mathcal{U}(\boldsymbol{\theta}) = -\sum_{i=1}^{N} \log p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) , \tag{4}$$

and posterior temperature $\mathcal{T}$, which is introduced as an additional hyperparameter. $\mathcal{T} = 1$ corresponds to the Bayesian posterior, while $\mathcal{T} < 1$ are sharper[44] cold posteriors.[38]

The gold-standard HMC[25,26] method leverages the gradient $\nabla_{\boldsymbol{\theta}}\mathcal{U}(\boldsymbol{\theta})$ to simulate Hamiltonian dynamics to generate parameter proposals for the Metropolis Hastings[45] (MH) acceptance step, which guarantees that the equilibrium distribution of the Markov chain corresponds to $p(\boldsymbol{\theta}|\mathcal{D})$. The computation of $\nabla_{\boldsymbol{\theta}}\mathcal{U}(\boldsymbol{\theta})$[26] requires an evaluation of the (NN) model for the whole training data set $\mathcal{D}$ (eq. (4)), rendering HMC computationally intractable for training NN potentials.[27,46]

## 2.3   Stochastic gradient MCMC

Stochastic gradient MCMC (SG-MCMC) methods[28–31] achieve enormous computational speed-ups by replacing $\nabla_{\boldsymbol{\theta}}\mathcal{U}(\boldsymbol{\theta})$ by a stochastic estimate over a mini-batch of data

$$\nabla_{\boldsymbol{\theta}}\tilde{\mathcal{U}}(\boldsymbol{\theta}) = -\frac{N}{B}\sum_{i=1}^{B} \nabla_{\boldsymbol{\theta}}\log p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}}\log p(\boldsymbol{\theta}) , \tag{5}$$

where B is the mini-batch size.

The simplest SG-MCMC scheme is the Stochastic Gradient Langevin Dynamics (SGLD)

method,[28] which updates parameters according to

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \frac{\lambda_k}{2}\nabla_{\boldsymbol{\theta}}\tilde{\mathcal{U}}(\boldsymbol{\theta}) + \boldsymbol{\eta}_t \; ; \quad \boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, \lambda_k\mathbf{I}) \; . \tag{6}$$

The learning rate $\lambda_k$ is decreasing as a function of update step $k$ and $\boldsymbol{\eta}_t$ is a learning rate-dependent Gaussian noise vector. $\lambda_k$ typically follows a polynomial schedule:[28,47]

$$\lambda_k = a(k+1)^{-\gamma} \; , \tag{7}$$

where $a$ is the initial learning rate and $\gamma$ is the decay rate. To reduce the bias due to the omitted MH acceptance step, it is necessary to sample only below a certain learning rate threshold, given that the acceptance probability asymptotically converges to 1 for $\lambda \to 0$. Hence, SGLD smoothly transitions from stochastic posterior maximization to asymptotically unbiased sampling from $p(\boldsymbol{\theta}|\mathcal{D})$ during training.[28,47] In our experiments, we employ a preconditioned version of SGLD (pSGLD),[30] which uses a RMSProp[48] preconditioner to simplify sampling the highly non-convex posterior of NNs,[30,49] as implemented in jax-sgmc.[50]

## 2.4   Deep Ensemble method

Analogous to the Monte Carlo approximation in Bayesian UQ (eq. (2)), the Deep Ensemble method[35,36] estimates epistemic uncertainty from the statistics of predictions from an ensemble of NNs. However, instead of sampling models from the posterior, the ensemble of NNs is generated by minimizing a loss function via stochastic gradient descent, starting from different random NN weight initializations. If desired, aleatoric uncertainty can be quantified by additionally predicting standard deviations and minimizing a negative log-likelihood loss.[36] While most authors consider the Deep Ensemble method non-Bayesian,[15,36–38] Wilson and Izmailov[43] compellingly argue that it can also be interpreted as Bayesian model averaging.

## 2.5 Neural Network Posterior Landscape

The posterior distribution of NNs is high-dimensional, non-convex and multi-modal.[43,44,49,51] The NNs of the Deep Ensemble typically converge to different posterior modes due to the strong decorrelation effect of different random weight initializations.[43,51] Hence, the Deep Ensemble method performs a Bayesian model average of NNs corresponding to different approximate maximum a-posteriori (MAP) points on the NN posterior (assuming regularization terms that mimic the prior).[43] The Deep Ensemble method therefore exploits the NN posterior multi-modality to estimate the uncertainty. By contrast, most scalable Bayesian methods, including single-chain SG-MCMC and stochastic variational inference, have been found to typically approximate a single posterior mode only.[37,41,43] However, sampling multiple posterior modes is essential for robust UQ.[43]

## 2.6 Multi-chain SG-MCMC

Sampling the posterior with multiple randomly initialized SG-MCMC chains appears to be a promising approach. It combines Bayesian posterior exploration along the Markov chain with strong decorrelation from different random initializations, the benefits of which have been shown to be complementary.[44,51] Multi-chain SG-MCMC can be interpreted as a custom trade-off between the number of approximated posterior modes and the amount of Bayesian exploration per mode, with single-mode Bayesian methods and the Deep Ensemble method representing the two extreme cases.

The computational training cost of the Deep Ensemble method and SG-MCMC can be estimated as $C * n_{\text{steps}} * n_{\text{chains}}$, where $n_{\text{chains}}$ is the number of ensemble members (chains), $n_{\text{steps}}$ is the number of parameter updates per ensemble member (chain) and $C$ is the cost per update. Training the different ensemble members (chains) can be parallelized trivially, if desired.

## 2.7 Probabilistic Molecular Modeling

### 2.7.1 Maximum Likelihood Molecular Modeling

The most common training scheme for atomistic (AT) NN potentials is to match the potential energy (possibly also forces and virial) of an underlying high-fidelity model, usually a computational quantum mechanics scheme,[52] given a training data set of $N_{\text{box}}$ molecular states.[8] This can be achieved by minimizing the mean squared error loss function

$$L(\boldsymbol{\theta}) = \frac{1}{N_{\text{box}}} \sum_{i=1}^{N_{\text{box}}} [U_i - U_{i,\boldsymbol{\theta}}]^2 \; , \tag{8}$$

where $U_i$ and $U_{i,\boldsymbol{\theta}}$ are the target and predicted potential energies of molecular state $i$, respectively. The predicted potential energy $U_{\boldsymbol{\theta}}(\mathbf{r})$ depends on atom positions $\mathbf{r}$.

Similarly, for CG systems, the NN potential can be trained via force matching (FM),[53–56] i.e., matching the instantaneous force components $F_j$ acting on each CG particle as computed from the AT force field:

$$L(\boldsymbol{\theta}) = \frac{1}{N_{\text{F}}} \sum_{j=1}^{N_{\text{F}}} [F_j - F_{j,\boldsymbol{\theta}}]^2 \; , \tag{9}$$

where $N_{\text{F}}$ is 3 times the number of CG particles in the training data set. The predicted force components are computed from the CG NN potential $\mathbf{F}_{\boldsymbol{\theta}} = -\nabla_{\mathbf{R}} U_{\boldsymbol{\theta}}^{\text{CG}}(\mathbf{R})$, which acts on CG coordinates $\mathbf{R} = \mathbf{M}(\mathbf{r})$. $\mathbf{M}$ is a linear function that maps from AT to CG coordinates. For infinite data and model capacity, $U_{\boldsymbol{\theta}}^{\text{CG}}(\mathbf{R})$ converges to the potential of mean force (PMF). Given that multiple AT configurations map to the same CG configuration, there exists a lower bound of the loss in eq. (9), which corresponds to the loss of the PMF.[12,54]

### 2.7.2 Bayesian Molecular Modeling

Assuming independent Gaussian homoscedastic aleatoric uncertainty with variance $\sigma_{\text{H}}^2$, the probabilistic model of the energy matching task is $p(U|\mathbf{r}, \boldsymbol{\theta}) \sim \mathcal{N}(U_{\boldsymbol{\theta}}(\mathbf{r}), \sigma_{\text{H}}^2)$. In this case,

the likelihood can be written similar to eq. (8) as[27]

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{i=1}^{N_{\text{box}}} \frac{1}{\sqrt{2\pi\sigma_{\text{H}}^2}} \exp\left(-\frac{[U_i - U_{i,\boldsymbol{\theta}}]^2}{2\sigma_{\text{H}}^2}\right)$$
$$= \left(\frac{1}{\sqrt{2\pi\sigma_{\text{H}}^2}}\right)^{N_{\text{box}}} \exp\left(-\frac{\sum_{i=1}^{N_{\text{box}}}[U_i - U_{i,\boldsymbol{\theta}}]^2}{2\sigma_{\text{H}}^2}\right). \tag{10}$$

The probabilistic model and the likelihood for the FM task follow analogously. The loss minima in eq. (8) and (9) correspond to the likelihood maxima in eq. (10).

The aleatoric uncertainty is uncertainty inherent to the data. When learning atomistic models from simulation data, the aleatoric uncertainty stems from the data-generating simulation and is typically small. For CG systems, the non-injective CG mapping contributes significantly to the aleatoric uncertainty. The variance of the aleatoric uncertainty $\sigma_{\text{H}}^2$ is typically unknown a priori and we model it as a learnable parameter. Thus, the prior $p(\boldsymbol{\theta}) = p(\mathbf{w})p(\sigma_{\text{H}})$ is the product of a prior for the NN potential weights and biases $\mathbf{w}$ and a prior for the aleatoric uncertainty scale.

## 2.8  Neural Network Potential

We choose a graph NN potential, which is a state-of-the-art NN architecture that learns to extract predictive features from the molecular configuration in an end-to-end manner instead of relying on hand-crafted descriptors.[2,3] Specifically, we select our previously published implementation[13] of the DimeNet++[4,5] potential. We set all hyperparameters to their default values, including the graph cut-off radius of $r_{\text{cut}} = 0.5$ nm, except for embedding sizes, which we reduce by factor 4 for computational speed-up. We select a Gaussian prior over all learnable weights and biases $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, 10^2\mathbf{I})$, except for the radial Bessel frequencies,[4] which we model by a uniform distribution.

Given that DimeNet++ trained via FM tends to yield unstable MD simulations,[57] we

augment the NN potential with a fixed, physics-informed "prior" potential $U^{\text{prior}}(\mathbf{R})$:[12,58,59]

$$U_{\boldsymbol{\theta}}^{\text{CG}}(\mathbf{R}) = U_{\boldsymbol{\theta}}^{\text{NN}}(\mathbf{R}) + U^{\text{prior}}(\mathbf{R}) \ . \tag{11}$$

Note that $U^{\text{prior}}(\mathbf{R})$ is not a prior in the Bayesian sense, but rather a physics-informed initialization that enforces physically reasonable predictions in phase-space regions unconstrained by the training data[12,13,60] (see supplementary methods 1 for more details).

## 3  Results

We present three examples (fig. 1) to distinguish between different sources of uncertainty: A LJ toy example features epistemic uncertainty only, while the following two CG systems include a significant amount of aleatoric uncertainty. We additionally show the effects of systematic uncertainty for liquid water and for alanine dipeptide.
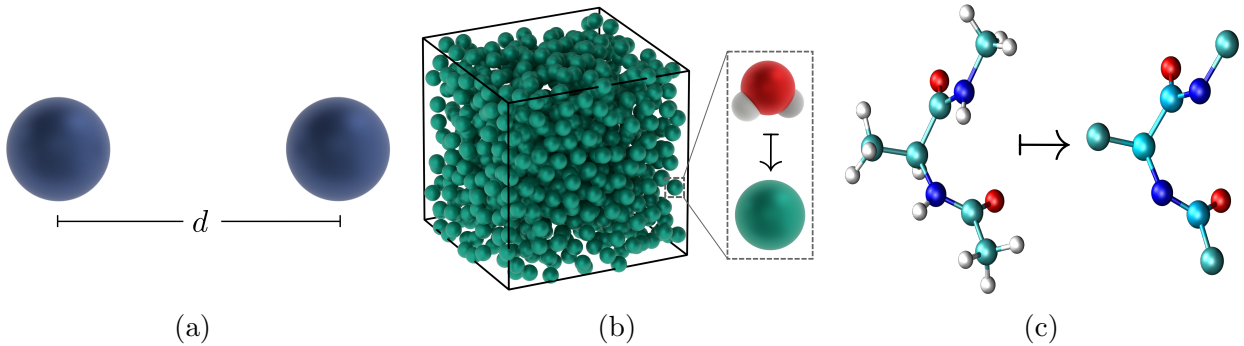


<div align="center">(a)      (b)      (c)</div>

Figure 1: Visualization of numerical test case systems: ($a$) Lennard Jones potential, ($b$) coarse-grained liquid water (adapted from Ref. 13), ($c$) coarse-grained alanine dipeptide.

## 3.1  Lennard Jones Potential

We learn a LJ potential $(\sigma_{\text{LJ}}, \epsilon)$ with a pairwise additive NN potential to benchmark the scalable UQ methods against a HMC scheme. As the reference method, we select the No-U-Turn Sampler (NUTS),[61] which selects the number of HMC integration steps on-the-fly. Additionally, a window adaption warm-up scheme[62,63] automatically selects an appropriate
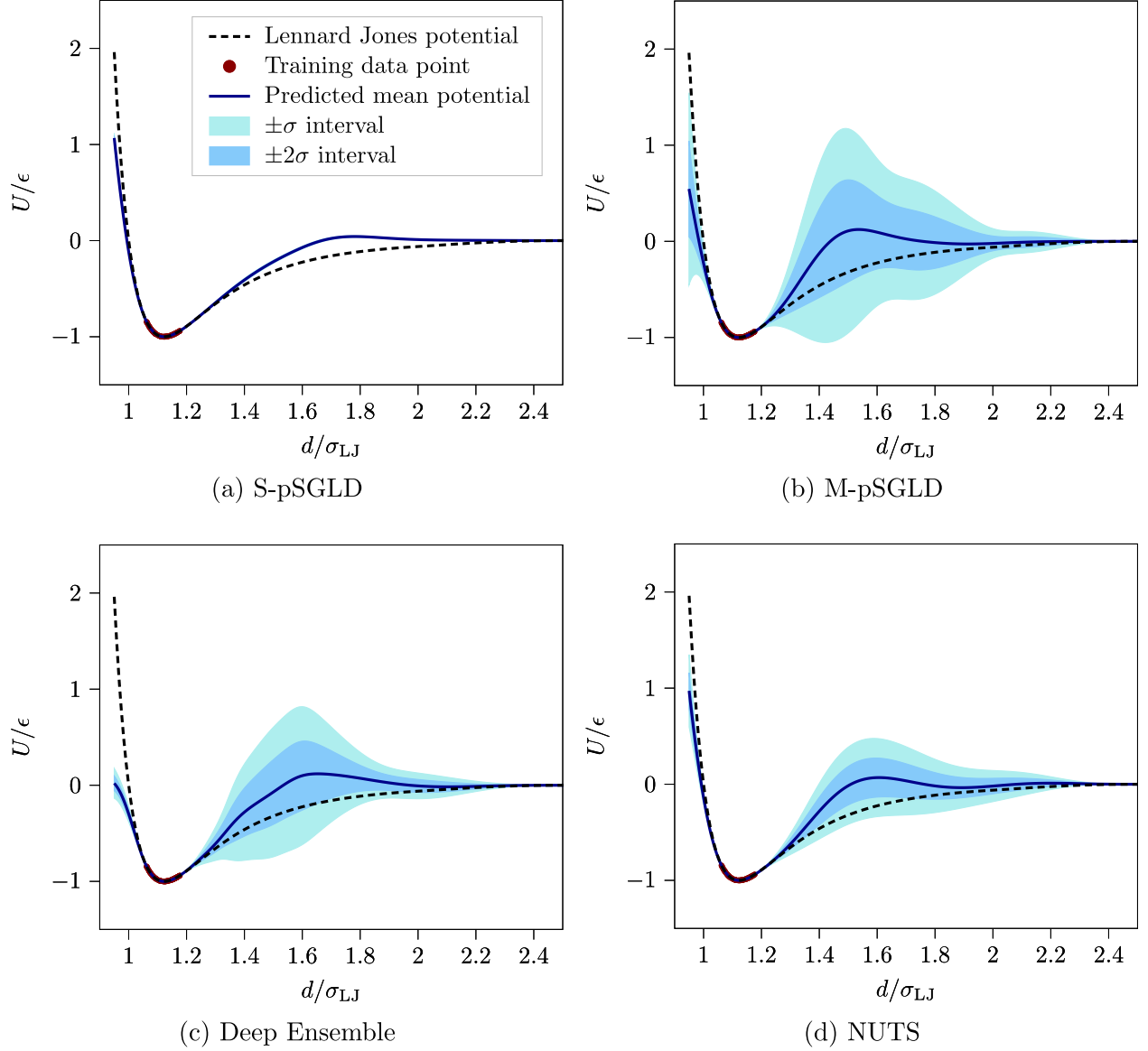
Figure 2: Distribution of NN potentials. Predicted mean potential with $\pm\sigma$ and $\pm2\sigma$ intervals of the single chain pSGLD (*a*), the multi-chain pSGLD (*b*), the Deep Ensemble method (*c*) and the multi-chain No-U-Turn Sampler (NUTS, *d*), compared to the Lennard Jones reference.

mass matrix and step size, such that no hyperparameter tuning is required for the NUTS. The NN potential predicts the pairwise potential energy $U(d)$ given pairwise particle distance $d$ and consists of a single hidden layer with 64 neurons and swish activation, where $d$ is represented by six radial Bessel functions[4] with a cut-off $r_{\mathrm{cut}} = 2.5\sigma_{LJ}$. We choose a Gaussian prior for weights and biases $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and an exponential prior with scale 1 for the aleatoric uncertainty $p(\sigma_{\mathrm{H}})$. For all Bayesian methods, we sample 100 models from the Bayesian posterior ($\mathcal{T} = 1$), evenly distributed over all considered Markov chains. The training data consists of 100 randomly drawn training data points from the well of the LJ potential $d/\sigma_{\mathrm{LJ}} \in [1.0615, 1.1800]$ (supplementary fig. 1). Additional technical details are provided in supplementary methods 2.

In the following, we benchmark pSGLD with a single chain (S-pSGLD), pSGLD with 10 chains (M-pSGLD) and a Deep Ensemble consisting of 10 NNs against a 10 chain NUTS. The obtained mean potentials and corresponding standard deviation intervals are visualized in fig. 2. The mean potentials of all considered methods fit the LJ potential very well where training data were generated: On held-out data within the training interval, we obtained low root-mean-squared errors (RMSE/$\epsilon$) of 0.011 (S-pSGLD), 0.014 (NUTS), 0.023 (M-pSGLD), and 0.025 (Deep Ensemble).

Bayesian methods estimate the scale of the aleatoric uncertainty to $\sigma_{\mathrm{H}}/\epsilon \approx 10^{-3}$. Such low estimates are expected given that the aleatoric uncertainty of the LJ data set is zero and the NN potential has sufficient capacity to interpolate the training data. Hence, we neglect the contribution of the aleatoric uncertainty in the following uncertainty predictions and only show the epistemic uncertainty. The S-pSGLD method samples a single posterior mode only, yielding highly overconfident potential energy predictions outside the training interval. By contrast, the other methods using multiple randomly initialized models sample multiple posterior modes, such that they can capture a significant amount of epistemic uncertainty. Accordingly, the obtained credible intervals include the reference potential across a broad range of distances. Both M-pSGLD and the Deep Ensemble method provide

good approximations to the NUTS reference distribution. However, compared to the NUTS reference, the Deep Ensemble method underestimates uncertainty at short distances and M-pSGLD overestimates uncertainty at medium distances.

All UQ methods sampling multiple modes exhibit a similar shape of the predicted epistemic uncertainty. Local uncertainty maxima are located between $1.4\sigma_{\text{LJ}} < d < 1.8\sigma_{\text{LJ}}$ and the uncertainty significantly increases for $d < 0.9\sigma_{\text{LJ}}$. This uncertainty shape is the result of the NN potential architecture and the location of the training data set: On the one hand, the radial Bessel representation[4] of $d$ smoothly shrinks the NN potential towards 0 at $r_{\text{cut}}$. On the other hand, the training data constrains the potential for $1.06\sigma_{\text{LJ}} < d < 1.18\sigma_{\text{LJ}}$. Hence, these results are consistent with the expectation that the epistemic uncertainty should increase with the distance from points that constrain the potential.
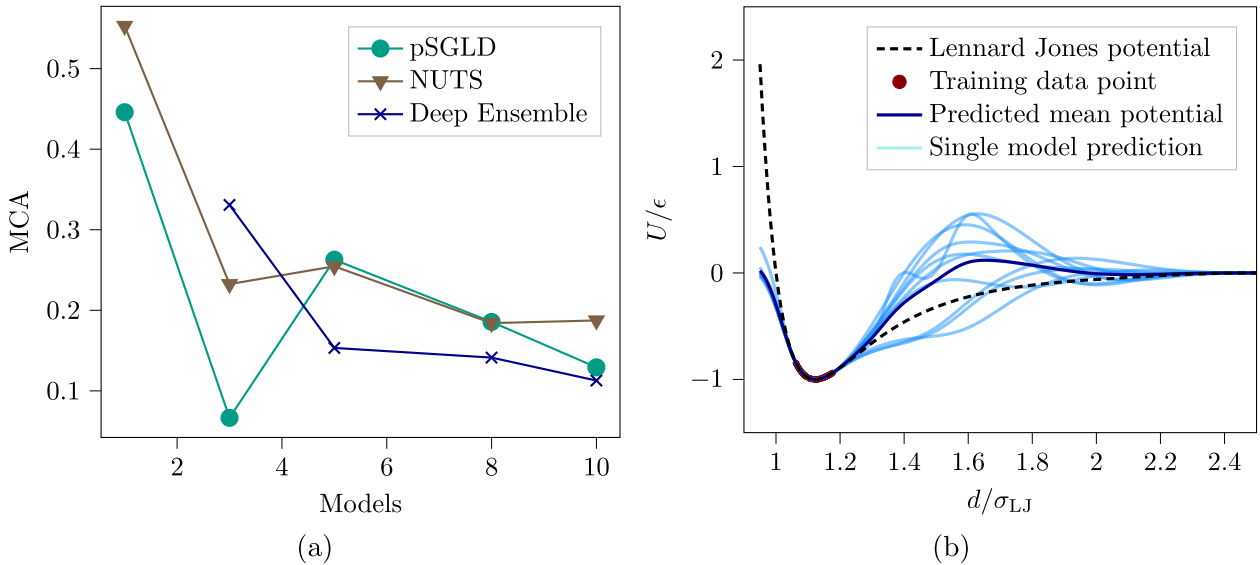


Figure 3: Posterior mode analysis. ($a$) miscalibration area (MCA) of the No-U-Turn Sampler (NUTS), the pSGLD, and the Deep Ensemble methods as a function of the number of randomly initialized models. The MCA includes both within and out-of-distribution test data. ($b$) All predicted potentials of the Deep Ensemble method with resulting mean compared to Lennard Jones reference.

We further investigate the effect of the number of randomly initialized models (number of Markov chains for MCMC) on UQ quality. The miscalibration area (MCA)[64–66] quantifies the agreement between the predicted standard deviation and the true error. For all methods,

13

the MCA shows a decreasing trend with increasing number of randomly initialized models (fig. 3 $a$). This reflects the importance of sampling multiple posterior modes for robust UQ,[43] which can be achieved comparatively easily by exploiting the strong decorrelation effect of random NN initializations.[43,51] Different posterior modes represent different potentials, all of which are consistent with the training data, but differ significantly where there is no data available, thus capturing the epistemic uncertainty (fig. 3 $b$).

The inability to sample multiple posterior modes using a single Markov chain is not unique to pSGLD. A single chain of the NUTS also samples a single posterior mode only and the captured epistemic uncertainty increases with additional chains (see also supplementary fig. 2). This suggests that sampling multiple posterior modes with a single Markov chain is difficult to achieve when training NN potentials, even for sophisticated posterior exploration schemes. Finally, we note that by artificially fixing $\sigma_{\mathrm{H}}$ to a large value as in Ref. 27, the single chain NUTS predicts large epistemic uncertainty outside the training interval (supplementary fig. 3). However, this comes at the cost of a larger error within the training interval (RMSE/$\epsilon = 0.044$ for $\sigma_{\mathrm{H}}/\epsilon = 0.05$) given that models with poorer fit also appear probable due to the allegedly large aleatoric noise in the data.

## 3.2 Coarse-grained Liquid Water

We apply pSGLD and the Deep Ensemble method to CG liquid water, a classic benchmark problem, to test their respective performance both within the training distribution as well as under distribution shift. The reference data consists of 100 cubic boxes of length $l = 3.129$ nm containing 1000 water molecules each, sampled every 1 ps from the TIP4P/2005[67] model at a temperature $T_{\mathrm{ref}} = 298$ K, resulting in a pressure $p_{\mathrm{ref}} = -6.2$ MPa. We divide the data into training, validation and test with a 80%-8%-12% split. Each water molecule is modeled by a CG particle positioned at its center of mass.

We select the repulsive part of the LJ potential as prior potential

$$U^{\mathrm{prior}}(\mathbf{R}) = \sum_{i=1}^{N_{\mathrm{pair}}} \epsilon_{\mathrm{w}} \left( \frac{\sigma_{\mathrm{w}}}{d_i} \right)^{12} \, , \tag{12}$$

with $\epsilon_{\mathrm{w}} = 1$ kJ/mol and $\sigma_{\mathrm{w}} = 0.3165$ nm, where $\sigma_{\mathrm{w}}$ corresponds to the length scale of the SPC water model.[68] This corresponds to the $U^{\mathrm{prior}}$ used in our previous works, where we found DimeNet++ results to be insensitive to the specific prior potential chosen.[13,60] We account for the thermodynamic state point dependency of the PMF[54,69,70] by augmenting the edge embedding of DimeNet++ by two learnable 16 dimensional vector, one multiplied and one divided by $k_{\mathrm{B}}T$. This dual embedding ensures that the temperature dependency effect vanishes for neither high nor low $k_{\mathrm{B}}T$.
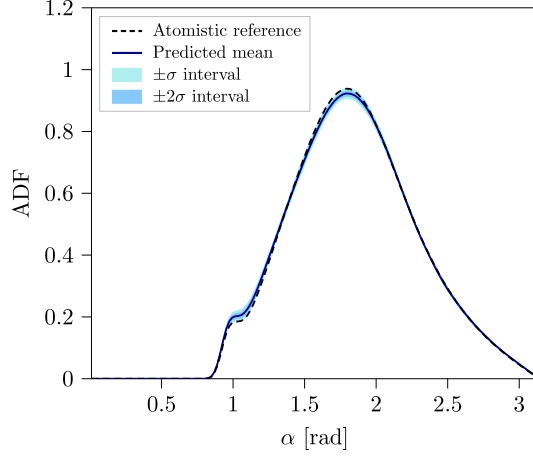
The models are trained with a batch size of 5 boxes, an initial learning rate $a = 5 \cdot 10^{-4}$ and a polynomial learning rate decay schedule (eq. (7)) with $\gamma = 0.55$. We generate a Deep Ensemble of 8 models and train each for 100 epochs with the Adam[71] optimizer with default parameters. For each training trajectory, we select the model parameters with the smallest validation loss, giving the Deep Ensemble method a slight advantage in terms of data usage over the Bayesian methods. For Bayesian modeling, we select a prior distribution $p(\sigma_{\mathrm{H}}) \sim \Gamma(5, 27)$, incorporating the prior knowledge that $\sigma_{\mathrm{H}} > 0$ due to the noise from the non-injective CG mapping.[56] By default, we select a posterior temperature $\mathcal{T} = 0.01$. Each pSGLD chain is run for 10000 epochs, 8000 of which are discarded as burn-in. We randomly subsample the remaining models such that a total of 40 models are selected, evenly distributed over all available chains (8 chains for M-pSGLD, 1 chain for S-pSGLD). One chain of the M-pSGLD method yielded poor potentials and we omitted it for a more balanced comparison.

First, we evaluate the mean force predictions on the test data. The Deep Ensemble method with a RMSE of 135.8 kJ/(mol nm) is more accurate than S-pSGLD and M-pSGLD with RMSE = 137.2 kJ/(mol nm) and RMSE = 136.6 kJ/(mol nm), respectively. We were
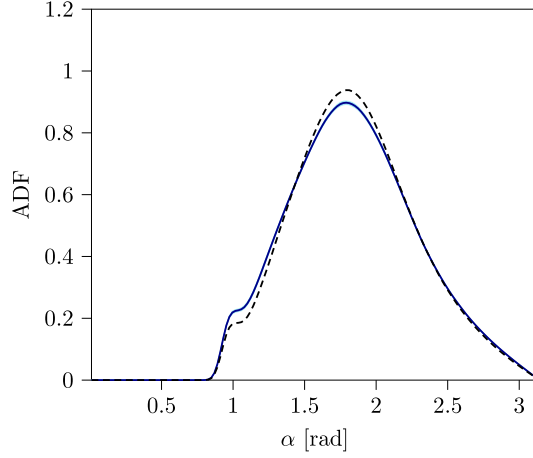
unable to find a set of pSGLD hyperparameters that closed the error differential to the Deep Ensemble method. The force error is dominated by the large aleatoric uncertainty, which is estimated as $\sigma_{\mathrm{H}} = 136.6$ kJ/(mol nm) by S-pSGLD.

We run CG MD simulations at a temperature $T = T_{\mathrm{ref}}$ to investigate the resulting observables without a distribution shift. All CG MD simulations use a time step of 2 fs and are equilibrated for 10 ps, followed by 100 ps of production, where a state is retained every 0.1 ps. The CG MD simulation averages over the aleatoric uncertainty resulting from the CG mapping. Consequently, the predicted standard deviation of observables $\sigma$ includes the epistemic uncertainty as well as a small amount of MD sampling uncertainty due to finite trajectory lengths. Fig. 4 shows the resulting distributions of angular distribution functions (ADF). The mean prediction of the Deep Ensemble method matches the AT reference well and is slightly more accurate than the S-pSGLD and M-pSGLD schemes, reflecting the lower test set RMSE. Additionally, the $2\sigma$ credible interval of the Deep Ensemble method covers the AT reference, and areas with higher uncertainty correspond to areas with larger error. M-pSGLD captures slightly more variance than S-pSGLD, but both schemes are overconfident. The overconfidence of M-pSGLD seems to be primarily attributable to a larger deviation of the predicted mean ADF from the reference curve (in line with the larger test set RMSE) and only secondarily to less captured epistemic uncertainty compared to the Deep Ensemble method. The conclusions drawn from the RDF predictions are identical (supplementary fig. 4).
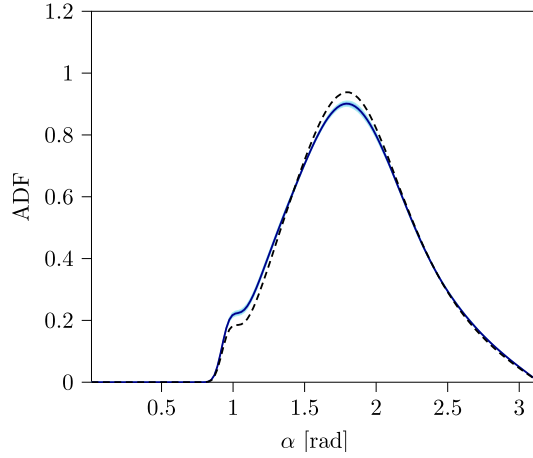
We investigate the impact of the Gaussian prior for weights and biases by retraining the NN potential with an improper uniform distribution. The obtained results for S-pSGLD are largely identical to the Gaussian prior case (supplementary fig. 5). However, the Gaussian prior appears to improve the learning robustness: With the uniform distribution, a total of 4 M-pSGLD Markov chains yielded models with large errors in the mean predictions, compared to only a single Markov chain with the Gaussian prior. Having verified that neither of the priors $p(\mathbf{w})$ and $p(\sigma_{\mathrm{H}})$ are too restrictive, we hypothesize that the higher RMSEs of the

16

(a) Deep Ensemble



(b) S-pSGLD



(c) M-pSGLD

Figure 4: Angular distribution functions (ADF) at $T = T_{\text{ref}}$. Resulting mean ADF with $\pm\sigma$ and $\pm2\sigma$ intervals as predicted by the Deep Ensemble method ($a$), the single chain pSGLD ($b$) and the multi-chain pSGLD ($c$) schemes at a temperature $T = T_{\text{ref}}$, compared to the atomistic reference.

17

pSGLD schemes may be the result of the training, where the coupling of learning rate and additive random noise might impede convergence to models with the highest likelihood.

Next, we investigate the impact of the posterior temperature $\mathcal{T}$ on pSGLD models (fig. 5). With the Bayesian posterior $\mathcal{T} = 1$, S-pSGLD requires a large data set set (1000 boxes)
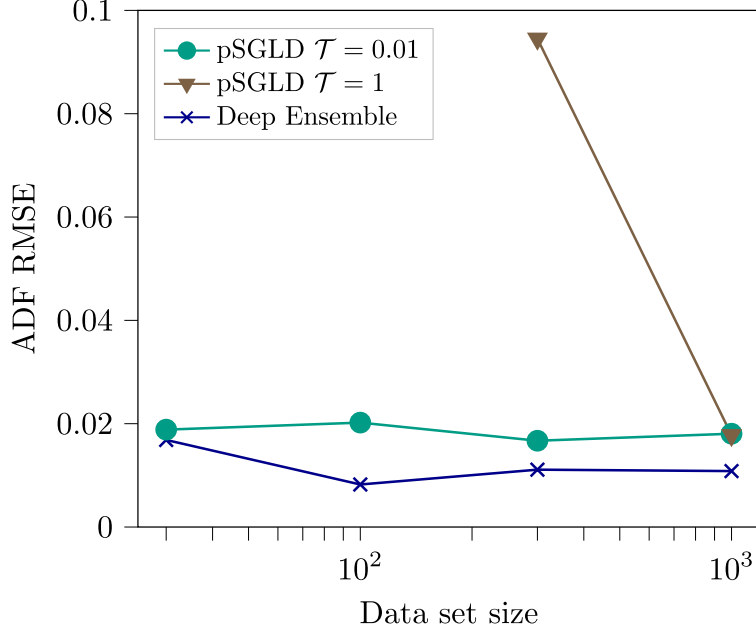


Figure 5: Cold posterior effect. Root mean squared error (RMSE) of the mean predicted angular distribution function (ADF) at $T = T_{\mathrm{ref}}$ of the Deep Ensemble method and the single chain pSGLD scheme with $\mathcal{T} = 1$ and $\mathcal{T} = 0.01$ for different data sizes. Note that pSGLD $\mathcal{T} = 1$ yields unstable MD simulations for data sizes of 30 and 100 boxes.

to sample accurate models. For a medium data set size (300 boxes), the obtained models are highly inaccurate compared to using the cold posterior $\mathcal{T} = 0.01$. For smaller data set sizes, models sampled with the Bayesian posterior result in unstable CG MD simulations. By contrast, the cold posterior allows to sample accurate models with only a fraction of the data (30 boxes). Moreover, the accuracy of pSGLD models hinges on a sufficient amount of burn-in epochs to reduce the learning rate (supplementary fig. 6). Consequently, the pSGLD schemes require significantly more computational training effort in this example than the Deep Ensemble method. Still, the Deep Ensemble method yields more accurate models for all data set sizes considered in fig. 5.
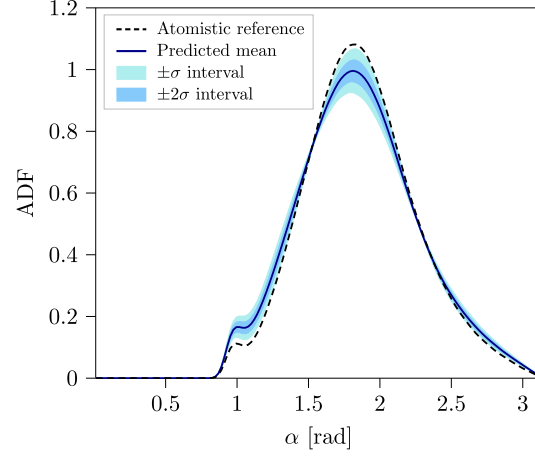
To test the quality of UQ under distribution shift, we apply the obtained models at

a temperature $T = 260$ K. The mean predictions of the considered UQ schemes are very similar to each other and, as expected, deviate from the respective TIP4P results (fig. 6). While S-pSGLD results in highly overconfident predictions, both M-pSGLD and the Deep Ensemble method provide accurate credible intervals, with a slight advantage for the latter. The predicted RDFs allow for identical conclusions (supplementary fig. 7). The accurate ADF credible intervals we obtained with the M-pSGLD and the Deep Ensemble method stand in contrast to previous findings with a 2-body cubic spline model:[72] Given that the 2-body spline cannot model many-body effects, uncertainty with respect to 3-body interactions cannot be captured in the epistemic uncertainty. This model misspecification results in systematic uncertainty not included in the credible interval. This highlights the advantage of the many-body capabilities of NN potentials in a UQ context.
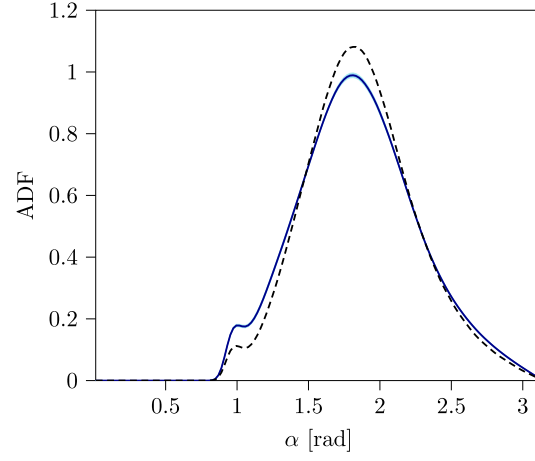
Finally, we study the impact of the $k_\mathrm{B}T$-dependent edge embedding. In the first step, we match the AT reference pressure at $T_\mathrm{ref}$ during the FM training (details in supplementary methods 3). Using the Deep Ensemble method, we then compute the density $\rho$ as a function of temperature with and without the $k_\mathrm{B}T$-dependent embeddings (fig. 7). As desired, the credible interval includes the AT reference and the uncertainty increases with the distance from the training temperature $T_\mathrm{ref}$ for the $k_\mathrm{B}T$-dependent model. By contrast, without $k_\mathrm{B}T$-dependent embedding, the predicted uncertainty barely increases with the distance from $T_\mathrm{ref}$, resulting in overconfident predictions due to model misspecification. Given that the $k_\mathrm{B}T$ dependence enables a broader range of outcomes at $T \neq T_\mathrm{ref}$, the mean predictions also change significantly and yield smaller errors further away from $T_\mathrm{ref}$. These results highlight the potency of scalable UQ methods to quantify errors resulting from applying CG models at different thermodynamic state points than during training.

## 3.3 Coarse-grained Alanine Dipeptide

We consider the benchmark problem of learning the free energy surface (FES) of alanine dipeptide,[73,74] which has recently been shown to be a challenging task for NN potentials

(a) Deep Ensemble

(b) S-pSGLD

(c) M-pSGLD

Figure 6: Out-of-distribution angular distribution functions (ADF) at $T = 260$ K. Resulting mean ADF with $\pm\sigma$ and $\pm 2\sigma$ intervals as predicted by the Deep Ensemble ($a$), the single chain pSGLD ($b$) and the multi-chain pSGLD ($c$) schemes at a temperature $T = 260$ K, compared to the atomistic reference.

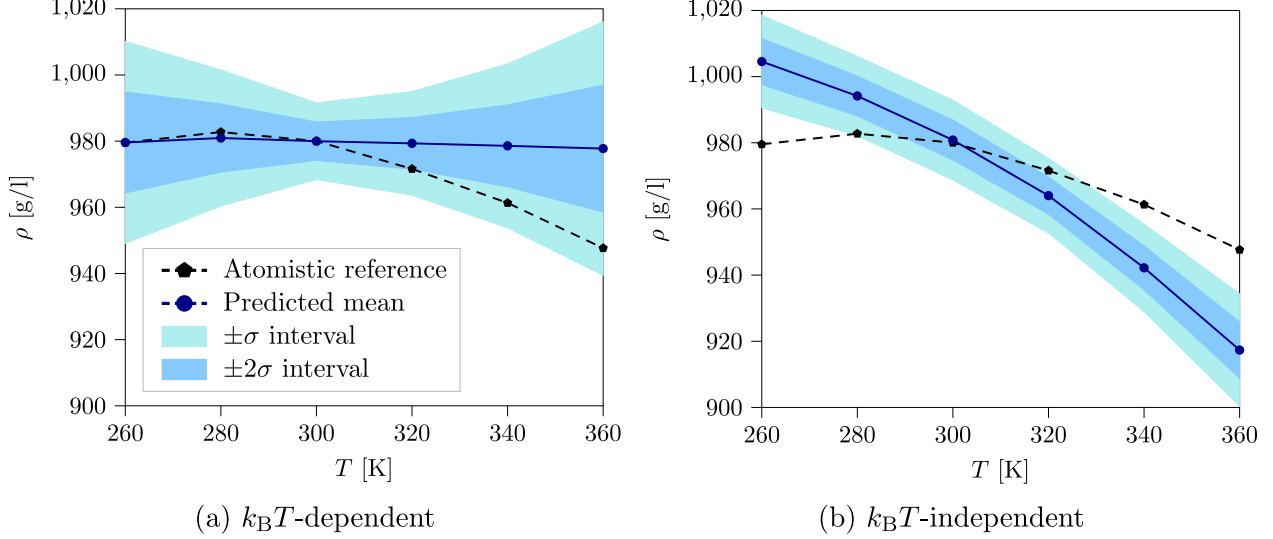| | (a) $k_\mathrm{B}T$-dependent | (b) $k_\mathrm{B}T$-independent |

Figure 7: Water density profile. Resulting mean density $\rho$ at pressure $p = 1$ bar with $\pm\sigma$ and $\pm 2\sigma$ intervals as predicted by the Deep Ensemble method using the $k_\mathrm{B}T$-dependent reference model $(a)$ and the same model without the $k_\mathrm{B}T$-dependent edge embedding $(b)$, compared to the atomistic reference.

trained via FM.[57,60] Here, we investigate the sources of these challenges using the scalable UQ toolbox. We build on the computational setup of our previous study:[60] The CG map retains all 10 heavy atoms of alanine dipeptide, dropping hydrogen atoms and water molecules. The CG particles modeling $CH_3$, $CH$ and $C$ are encoded as different particle types. The training data set consists of a 100 ns AT trajectory at $T_\mathrm{ref} = 300$ K, which is subsampled to $5 \cdot 10^5$ data points by retaining a state every 0.2 ps. The first 80 ns form the training set, the subsequent 8 ns the validation set. To counteract the instability of DimeNet++ in CG MD simulations of alanine dipeptide,[57] we add a prior potential $U^\mathrm{prior}(\mathbf{R})$ (eq. (11)) that consists of harmonic bonds and angles, as well as proper dihedrals. For more technical details on $U^\mathrm{prior}$ and the AT reference data, we refer to our previous work.[60]

We train all models with an initial learning rate $a = 10^{-3}$ and a polynomial learning rate decay schedule (eq. (7)) with $\gamma = 0.55$, as well as a batch size of 512 configurations. The 8 models of the Deep Ensemble method are trained for 1000 epochs, using the Adam[71] optimizer with default parameters. For each training trajectory, we select the model parameters with the smallest validation loss. pSGLD chains are run for 3000 epochs, with the first

2500 epochs discarded as burn-in. We randomly subsample the remaining models such that a total of 40 models are selected, evenly distributed over all available chains (8 chains for M-pSGLD, 1 chain for S-pSGLD). We select a prior distribution $p(\sigma_\mathrm{H}) \sim \Gamma(10, 40)$ and a posterior temperature $\mathcal{T} = 0.05$.

We initially evaluate performance of the considered methods on the test set: The S-pSGLD (RMSE = 414.12 kJ/(mol nm)), M-pSGLD (RMSE = 414.01 kJ/(mol nm)) and Deep Ensemble methods (RMSE = 413.84 kJ/(mol nm)) yield very similar accuracy, with a slight advantage for the Deep Ensemble method. This is in line with the aleatoric uncertainty scale $\sigma_\mathrm{H} = 414.69$, estimated by S-pSGLD.

We perform a 100 ns CG MD production simulation for all sampled models in order to compute the FES. To obtain the same number of trajectories as with the pSGLD schemes, each of the 8 Deep Ensemble models generates 5 trajectories, all starting from different initial states. Despite using a prior potential, some models became stuck in unphysical potential energy "holes",[14] i.e. deep potential energy minima in rarely sampled phase-space regions, which also led to instability in some cases. These potential energy holes might be avoided by employing better prior potentials or by incorporating MD simulations into training, e.g. via active learning[14,19] or alternative training schemes such as relative entropy (RE) minimization.[60,75,76] We note that using the Bayesian posterior $\mathcal{T} = 1$ significantly increased the number of unphysical trajectories (tested for S-pSGLD). For $\mathcal{T} = 1$, we observed results of comparable quality to $\mathcal{T} = 0.05$ only when increasing the data set to $1\mu s$.

First, we investigate UQ results after removing unphysical trajectories that mainly sampled configurations in a potential energy hole. To this end, we removed 1, 7 and 13 trajectories from the S-pSGLD, the M-pSGLD and the Deep Ensemble methods, respectively. The resulting means and standard deviations of the dihedral angles $\phi$ and $\psi$[73,74] are shown in fig. 8. The mean predictions of S-pSGLD and M-pSGLD are very similar, but – consistent with the examples above – S-pSGLD significantly underestimates the epistemic uncertainty. The Deep Ensemble method yields similar mean predictions in $\phi$ and a slightly improved mean

(a) S-pSGLD

(b) S-pSGLD

(c) M-pSGLD
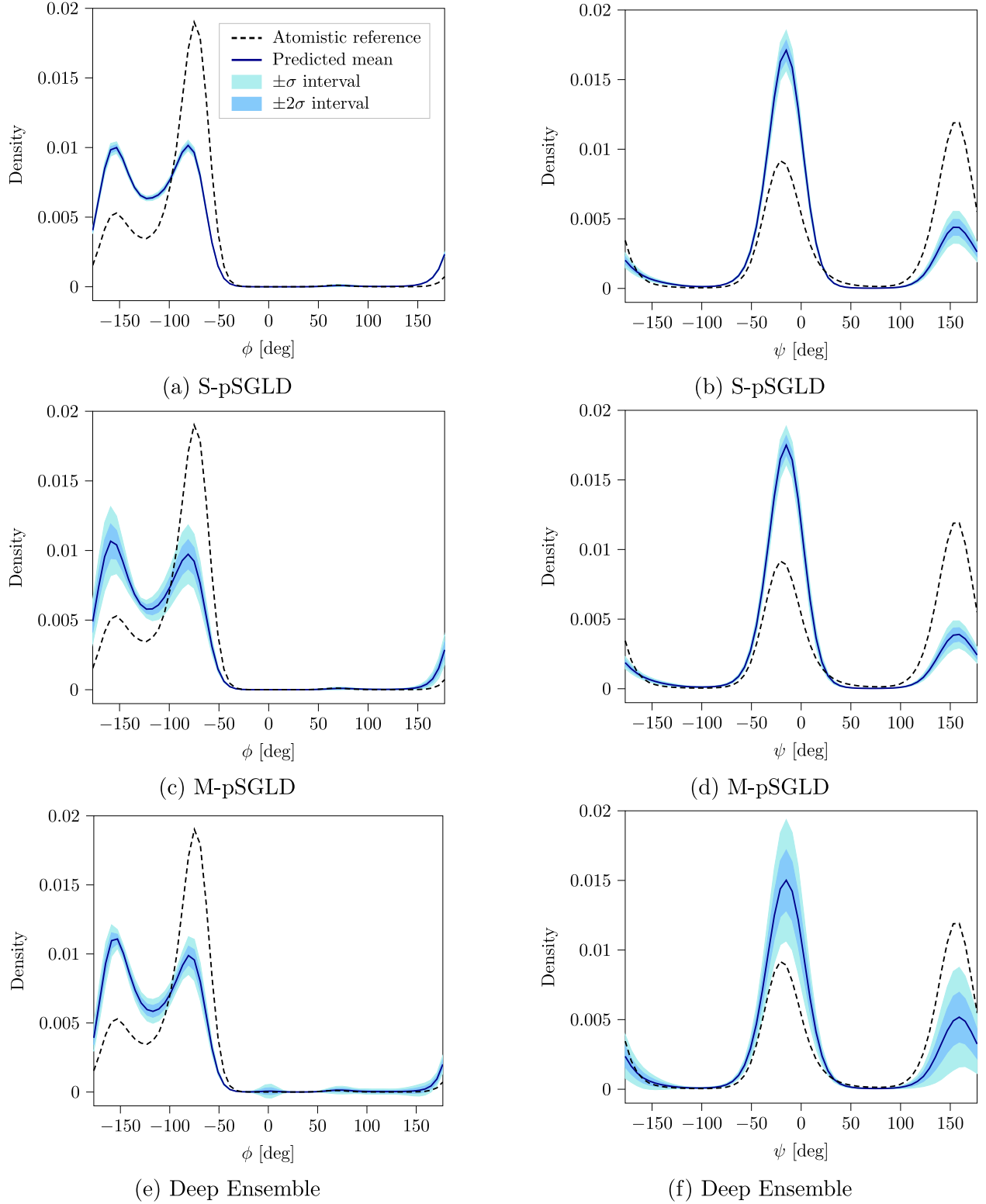
(d) M-pSGLD

(e) Deep Ensemble

(f) Deep Ensemble

Figure 8: Dihedral angle density histograms. Resulting mean distribution of dihedral angles $\phi$ (left column) and $\psi$ (right column) with $\pm\sigma$ and $\pm2\sigma$ intervals as predicted by the single chain pSGLD ($a$, $b$), the multi-chain pSGLD ($c$, $d$) and the Deep Ensemble ($e$, $f$) methods, compared to the atomistic reference.

prediction in $\psi$. The M-pSGLD predicts larger epistemic uncertainty in $\phi$, while the Deep Ensemble method predicts larger uncertainty in $\psi$. However, all considered methods show that the epistemic uncertainty is not sufficiently large to fully account for the deviation from the AT reference in this case.

To contextualize this result, we replace the FM training by RE minimization. Given that the RE model trained on the same data set can match the AT FES accurately,[60] insufficient model capacity is not the main limiting factor. Additionally, a poor approximation of the posterior by the considered UQ methods, which would result in an incorrect size of the predicted credible interval, can also be ruled out: The posterior probability ratio of the RE model and the last model sampled by the S-pSGLD FM scheme is $p(\theta_{\mathrm{RE}}|\mathcal{D})/p(\theta_{\mathrm{FM}}|\mathcal{D}) = \mathrm{e}^{-25872}$. This is in line with previous findings showing that the error on held-out force data is smaller for FM than for RE for alanine dipeptide.[60] Given that the posterior probability ratio of the RE model is numerically zero, UQ schemes with a FM-based posterior cannot sample the RE model.

Multiple mechanisms may contribute to the comparatively weak FES prediction of FM models. First, the FES of a FM model is sensitive to predictions in sparsely resolved transition regions.[60] Second, if a CG MD simulation is able to reach unphysical phase-space regions, sampling such configurations yields an erroneous FES. Both of these mechanisms result in very large epistemic uncertainty. We empirically show this effect when we include trajectories that sampled potential energy holes in the evaluation of the FES distribution (supplementary fig. 8). In particular, the predicted credible intervals of M-pSGLD and the Deep Ensemble method mostly cover the reference FES in this case. Hence, these UQ methods can signal to practitioners that the obtained results are not yet trustworthy.

Rather, more data needs to be generated to further constrain the learned models. We increased the data set size by factor 10 by generating an AT trajectory of 1 $\mu$s length. Interestingly, simply generating more Boltzmann-distributed training data did not solve the potential energy hole problem, nor did it significantly reduce the deviation of the mean

prediction when neglecting trajectories stuck in potential energy holes (supplementary fig. 9). Hence, generating more diverse, non-Boltzmann distributed data sets, e.g. via enhanced sampling schemes[77,78] or active learning,[14,19] seems to be a more promising approach.

The remaining deviation of the mean prediction from the reference FES that is not captured by the predicted epistemic uncertainty (fig. 8) suggests that other, likely systematic sources of error exist. For finite model capacity, a systematic difference between FM and RE minimization is the objective function: RE training minimizes the difference between the potential energy surfaces of the AT and CG models,[79] which is directly related to the FES. In contrast, the optimum of a force-based training objective might trade off accuracy in the FES for improved accuracy in other (e.g. thermodynamic[60]) observables, resulting in systematic uncertainty in the predicted FES. Additionally, numerical errors introduced by the CG MD simulation, similar to the shadow Hamiltonian effect,[80,81] can be corrected for by RE minimization.[60] In FM models, these numerical errors manifest as unquantifiable systematic uncertainty. However, for a comprehensive analysis of the relative impact of each error mechanism, further research is needed.

# 4   Discussion and Conclusion

Our results show that M-pSGLD is well suited to estimate the epistemic uncertainty of MD observables. This method enables fully-Bayesian UQ for NN potentials. All experiments highlight the importance of sampling multiple posterior modes. Exploiting the strong decorrelation effect of multiple random NN initializations via multiple Markov chains is an effective means to this end. In the graph NN examples, cold posteriors proved beneficial to sample both stable and accurate models, reducing the required amount of training data significantly. Hence, we found the number of Markov chains to be the most important additional M-pSGLD hyperparameter, followed by the posterior temperature, the prior distributions and the number of samples per chain.

Both Deep Ensemble and the M-pSGLD methods provided good approximations to the epistemic uncertainty estimated by the NUTS[61] in the LJ example. In addition, the Deep Ensemble method yielded similar UQ quality to M-pSGLD, although it required less training and hyperparameter tuning effort. We found no evidence that the Deep Ensemble method was prone to overconfident predictions, contrasting prior research in an active learning setting.[27] Instead, our results suggest that the Deep Ensemble method quantifies epistemic uncertainty effectively, both within and out of the training distribution.

M-pSGLD promises accurate UQ by leveraging the complementary benefits from sampling multiple posterior modes and additional Bayesian exploration of each mode.[43,44,51] However, further research into SG-MCMC schemes is required before routine application in practice: In our experiments, a single MCMC chain (both pSGLD[28] and NUTS[61]) sampled a single posterior mode only. Hence, the development of methods that sample multiple posterior modes with a single chain, e.g., by leveraging cyclical step size schedules[82] or parallel tempering,[83] is important. Additionally, automatic hyperparameter tuning with a computationally efficient metric could improve the SG-MCMC efficiency; e.g., the popular Stein's discrepancy[84] scales quadratically with the data set size.[31] Finally, recent SG-MCMC samplers such as AMAGOLD[82] or SGGMC[85] include infrequent Metropolis-Hastings acceptance steps to avoid the bias of SGLD.[28,30] Consequently, these samplers use constant learning rates, which may counteract the increased training time of SGLD that results from its small learning rate requirement.[47,82]

We observed a clear cold posterior effect[38] in our experiments with the graph NN potential. For image classification tasks, cold posteriors have demonstrated superior performance in practice.[38,86,87] However, this performance increase is mainly attributed to data augmentation,[44] which increases the effective data size without increasing the data size considered in the likelihood. Analogously, the effective data size might be underestimated by the likelihood in eq. (10), which the cold posterior may correct for: When learning the potential energy of a molecular state, the effective data size is clearly larger than a single data point. For

instance, in the case of a pairwise additive potential, the effective data size corresponds to the total number of particle pairs within a cutoff. For FM, the data size per box considered in the likelihood in eq. (10) equals 3 times the number of CG particles, but whether the effective data size exceeds this value is less clear. More research into the nature of the cold posterior effect is required – ideally resulting in likelihood formulations that better consider the effective data size.

Our results corroborate that for successful UQ, a sufficiently large hypothesis space is necessary: Effects describable by the model can be quantified reliably as epistemic uncertainty, but effects beyond the model capacity become hard to quantify systematic uncertainties.[15,42] For instance, if a potential lacks important many-body interactions or a CG model lacks state-point dependency, the resulting uncertainty estimates are overconfident. Consequently, NN potentials are attractive models in a UQ context, given that they model many-body interactions inherently.

To obtain uncertainty estimates for MD observables, we performed a dedicated MD simulation for every sampled NN potential. This approach is rigorous, as both epistemic uncertainty and MD sampling uncertainty are captured,[18] but also computationally expensive. The computational effort for MD simulations scales linearly with the number of sampled potentials, but the simulations can be parallelized. Distilling the mean potential energy prediction into a single model via student-teacher[30,88,89] learning could improve computational efficiency. With this approach, one could obtain uncertainty estimates for time-averaged observables using a single MD simulation and a reweighting scheme.[24] Concerning the development of computationally more efficient UQ schemes for NN potentials,[40] we have demonstrated that both M-pSGLD and the Deep Ensemble method can serve as reliable baseline schemes. Efficient UQ schemes may pave the way for more reliable MD simulations based on NN potentials to support simulation-based decision-making in health care and material science industries.[18]

# Acknowledgement

# Supporting Information Available

Background information on the prior potential, training details of the Lennard Jones example and training data visualization, pressure correction scheme, NUTS predictions for different number of Markov chains and fixed $\sigma_H$, predicted CG water RDFs, ADFs for S-pSGLD with uniform prior over weights and biases, ADF RMSE for different pSGLD Markov chain lengths, dihedral angle density histograms without removing potential energy holes as well as for the 1 $\mu$s data set. This information is available free of charge via the Internet at https://pubs.acs.org/.

# References

(1) Behler, J.; Parrinello, M. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys. Rev. Lett.* **2007**, *98*, 146401.

(2) Schütt, K. T.; Kindermans, P. J.; Sauceda, H. E.; Chmiela, S.; Tkatchenko, A.; Müller, K. R. SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. Advances in Neural Information Processing Systems. Long Beach, CA, USA, Dec. 4–9, 2017.

(3) Gilmer, J.; Schoenholz, S. S.; Riley, P. F.; Vinyals, O.; Dahl, G. E. Neural message passing for quantum chemistry. Proceedings of the 34th International Conference on Machine Learning. Sydney, Australia, Aug. 6–11, 2017; pp 1263–1272.

(4) Klicpera, J.; Groß, J.; Günnemann, S. Directional Message Passing for Molecular

Graphs. 8th International Conference on Learning Representations. Online, Apr. 26 – May 1, 2020.

(5) Klicpera, J.; Giri, S.; Margraf, J. T.; Günnemann, S. Fast and Uncertainty-Aware Directional Message Passing for Non-Equilibrium Molecules. Machine Learning for Molecules Workshop at NeurIPS. Online, Dec. 12, 2020.

(6) Ko, T. W.; Finkler, J. A.; Goedecker, S.; Behler, J. A fourth-generation high-dimensional neural network potential with accurate electrostatics including non-local charge transfer. *Nat. Commun.* **2021**, *12*, 398.

(7) Batzner, S.; Musaelian, A.; Sun, L.; Geiger, M.; Mailoa, J. P.; Kornbluth, M.; Molinari, N.; Smidt, T. E.; Kozinsky, B. E (3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nat. Commun.* **2022**, *13*, 2453.

(8) Noé, F.; Tkatchenko, A.; Müller, K. R.; Clementi, C. Machine Learning for Molecular Simulation. *Annu. Rev. Phys. Chem.* **2020**, *71*, 361–390.

(9) Schütt, K. T.; Arbabzadah, F.; Chmiela, S.; Müller, K. R.; Tkatchenko, A. Quantum-chemical insights from deep tensor neural networks. *Nat. Commun.* **2017**, *8*, 13890.

(10) Faber, F. A.; Hutchison, L.; Huang, B.; Gilmer, J.; Schoenholz, S. S.; Dahl, G. E.; Vinyals, O.; Kearnes, S.; Riley, P. F.; Von Lilienfeld, O. A. Prediction Errors of Molecular Machine Learning Models Lower than Hybrid DFT Error. *J. Chem. Theory Comput.* **2017**, *13*, 5255–5264.

(11) Stocker, S.; Gasteiger, J.; Becker, F.; Günnemann, S.; Margraf, J. T. How robust are modern graph neural network potentials in long and hot molecular dynamics simulations? *Mach. Learn.: Sci. Technol.* **2022**, *3*, 045010.

(12) Wang, J.; Olsson, S.; Wehmeyer, C.; Pérez, A.; Charron, N. E.; De Fabritiis, G.; Noé, F.;

Clementi, C. Machine Learning of Coarse-Grained Molecular Dynamics Force Fields. *ACS Cent. Sci.* **2019**, *5*, 755–767.

(13) Thaler, S.; Zavadlav, J. Learning neural network potentials from experimental data via Differentiable Trajectory Reweighting. *Nat. Commun.* **2021**, *12*, 6884.

(14) van der Oord, C.; Sachs, M.; Kovács, D. P.; Ortner, C.; Csányi, G. Hyperactive Learning (HAL) for Data-Driven Interatomic Potentials. *arXiv:2210.04225* **2022**,

(15) Gal, Y.; Koumoutsakos, P.; Lanusse, F.; Louppe, G.; Papadimitriou, C. Bayesian uncertainty quantification for machine-learned models in physics. *Nat. Rev. Phys.* **2022**, *4*, 573–577.

(16) Angelikopoulos, P.; Papadimitriou, C.; Koumoutsakos, P. Bayesian uncertainty quantification and propagation in molecular dynamics simulations: a high performance computing framework. *J. Chem. Phys.* **2012**, *137*, 144103.

(17) Zavadlav, J.; Arampatzis, G.; Koumoutsakos, P. Bayesian selection for coarse-grained models of liquid water. *Sci. Rep.* **2019**, *9*, 99.

(18) Wan, S.; Sinclair, R. C.; Coveney, P. V. Uncertainty quantification in classical molecular dynamics. *Philos. Trans. Royal Soc. A* **2021**, *379*, 20200082.

(19) Smith, J. S.; Nebgen, B.; Lubbers, N.; Isayev, O.; Roitberg, A. E. Less is more: Sampling chemical space with active learning. *J. Chem. Phys.* **2018**, *148*, 241733.

(20) Zhang, L.; Lin, D.-Y.; Wang, H.; Car, R.; Weinan, E. Active learning of uniformly accurate interatomic potentials for materials simulation. *Phys. Rev. Mater.* **2019**, *3*, 023804.

(21) Loeffler, T. D.; Patra, T. K.; Chan, H.; Sankaranarayanan, S. K. Active learning a coarse-grained neural network model for bulk water from sparse training data. *Mol. Syst. Des. Eng.* **2020**, *5*, 902–910.

(22) Smith, J. S.; Nebgen, B.; Mathew, N.; Chen, J.; Lubbers, N.; Burakovsky, L.; Tretiak, S.; Nam, H. A.; Germann, T.; Fensin, S., et al. Automated discovery of a robust interatomic potential for aluminum. *Nat. Commun.* **2021**, *12*, 1257.

(23) Xie, S. R.; Rupp, M.; Hennig, R. G. Ultra-fast Force Fields (UF3) framework for machine-learning interatomic potentials. American Physical Society March Meeting. Chicago, IL, USA, Mar. 14–18, 2021.

(24) Imbalzano, G.; Zhuang, Y.; Kapil, V.; Rossi, K.; Engel, E. A.; Grasselli, F.; Ceriotti, M. Uncertainty estimation for molecular dynamics and sampling. *J. Chem. Phys.* **2021**, *154*, 074102.

(25) Duane, S.; Kennedy, A. D.; Pendleton, B. J.; Roweth, D. Hybrid Monte Carlo. *Phys. Lett. B* **1987**, *195*, 216–222.

(26) Neal, R. M. In *Handbook of Markov Chain Monte Carlo*, 1st ed.; Brooks, S., Gelman, A., Jones, G. L., Meng, X.-L., Eds.; Chapman and Hall/CRC: New York, USA, 2011; Chapter MCMC using Hamiltonian Dynamics, pp 139–188.

(27) Kahle, L.; Zipoli, F. Quality of uncertainty estimates from neural network potential ensembles. *Phys. Rev. E* **2022**, *105*, 015311.

(28) Welling, M.; Teh, Y. W. Bayesian learning via stochastic gradient Langevin dynamics. Proceedings of the 28th International Conference on Machine Learning. Bellevue, WA, USA, Jun. 28 – Jul. 2, 2011; pp 681–688.

(29) Chen, T.; Fox, E.; Guestrin, C. Stochastic gradient Hamiltonian Monte Carlo. Proceedings of the 31st International Conference on Machine Learning. Beijing, China, Jun. 21 –26, 2014; pp 1683–1691.

(30) Li, C.; Chen, C.; Carlson, D. E.; Carin, L. Preconditioned Stochastic Gradient Langevin

Dynamics for Deep Neural Networks. Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence. Phoenix, AZ, USA, February 12–17, 2016; pp 1788–1794.

(31) Nemeth, C.; Fearnhead, P. Stochastic gradient Markov chain Monte Carlo. *J. Am. Stat. Assoc.* **2021**, *116*, 433–450.

(32) Lamb, G.; Paige, B. Bayesian Graph Neural Networks for Molecular Property Prediction. Machine Learning for Molecules Workshop at NeurIPS. Online, Dec. 12, 2020.

(33) Graves, A. Practical Variational Inference for Neural Networks. Advances in Neural Information Processing Systems. Granada, Spain, Dec. 12 – 14, 2011.

(34) Hoffman, M. D.; Blei, D. M.; Wang, C.; Paisley, J. Stochastic Variational Inference. *J. Mach. Learn. Res.* **2013**, *14*, 1303–1347.

(35) Hansen, L.; Salamon, P. Neural Network Ensembles. *IEEE Trans. Pattern Anal. Machine Intell.* **1990**, *12*, 993–1001.

(36) Lakshminarayanan, B.; Pritzel, A.; Blundell, C. Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. Advances in Neural Information Processing Systems. Long Beach, CA, USA, Dec. 4–9, 2017.

(37) Ovadia, Y.; Fertig, E.; Ren, J.; Nado, Z.; Sculley, D.; Nowozin, S.; Dillon, J.; Lakshminarayanan, B.; Snoek, J. Can you trust your model's uncertainty? Evaluating predictive uncertainty under dataset shift. Advances in Neural Information Processing Systems. Vancouver, BC, Canada, Dec. 8 – 14, 2019.

(38) Wenzel, F.; Roth, K.; Veeling, B. S.; Swikatkowski, J.; Tran, L.; Mandt, S.; Snoek, J.; Salimans, T.; Jenatton, R.; Nowozin, S. How Good is the Bayes Posterior in Deep Neural Networks Really? Proceedings of the 37th International Conference on Machine Learning. Online, Jul. 13 – 18, 2020; pp 10248–10259.

(39) Wen, M.; Tadmor, E. B. Uncertainty quantification in molecular simulations with dropout neural network potentials. *Npj Comput. Mater.* **2020**, *6*, 124.

(40) Zhu, A.; Batzner, S.; Musaelian, A.; Kozinsky, B. Fast Uncertainty Estimates in Deep Learning Interatomic Potentials. *arXiv preprint arXiv:2211.09866* **2022**,

(41) Gustafsson, F. K.; Danelljan, M.; Schon, T. B. Evaluating Scalable Bayesian Deep Learning Methods for Robust Computer Vision. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops. Seattle, WA, USA, Jun. 14 – 19, 2020; pp 1289–1298.

(42) Hüllermeier, E.; Waegeman, W. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Mach. Learn.* **2021**, *110*, 457–506.

(43) Wilson, A. G.; Izmailov, P. Bayesian Deep Learning and a Probabilistic Perspective of Generalization. Advances in Neural Information Processing Systems. Online, Dec. 6–12, 2020.

(44) Izmailov, P.; Vikram, S.; Hoffman, M. D.; Wilson, A. G. G. What Are Bayesian Neural Network Posteriors Really Like? Proceedings of the 38th International Conference on Machine Learning. Online, Jul. 18–44, 2021; pp 4629–4640.

(45) Hastings, W. K. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **1970**, *57*, 97–109.

(46) Li, D. W.; Brüschweiler, R. Iterative Optimization of Molecular Mechanics Force Fields from NMR Data of Full-Length Proteins. *J. Chem. Theory Comput.* **2011**, *7*, 1773–1782.

(47) Teh, Y. W.; Thiery, A. H.; Vollmer, S. J. Consistency and Fluctuations For Stochastic Gradient Langevin Dynamics. *J. Mach. Learn. Res.* **2016**, *17*, 1–33.

(48) Tieleman, T.; Hinton, G. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning* **2012**, *4*, 26–31.

(49) Dauphin, Y. N.; Pascanu, R.; Gulcehre, C.; Cho, K.; Ganguli, S.; Bengio, Y. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. Advances in Neural Information Processing Systems. Montreal, QC, Canada, Dec. 8–13, 2014.

(50) Thaler, S.; Fuchs, P.; Cukarska, A.; Zavadlav, J. jax-sgmc: Modular Stochastic Gradient MCMC for JAX. 2020; `https://github.com/tummfm/jax-sgmc`.

(51) Fort, S.; Hu, H.; Lakshminarayanan, B. Deep ensembles: A loss landscape perspective. *arXiv preprint arXiv:1912.02757* **2019**,

(52) Smith, J. S.; Zubatyuk, R.; Nebgen, B.; Lubbers, N.; Barros, K.; Roitberg, A. E.; Isayev, O.; Tretiak, S. The ANI-1ccx and ANI-1x data sets, coupled-cluster and density functional theory properties for molecules. *Sci. Data* **2020**, *7*, 134.

(53) Izvekov, S.; Voth, G. A. Multiscale coarse graining of liquid-state systems. *J. Chem. Phys.* **2005**, *123*, 134105.

(54) Noid, W. G.; Chu, J.-W.; Ayton, G. S.; Krishna, V.; Izvekov, S.; Voth, G. A.; Das, A.; Andersen, H. C. The multiscale coarse-graining method. I. A rigorous bridge between atomistic and coarse-grained models. *J. Chem. Phys.* **2008**, *128*, 244114.

(55) Noid, W.; Liu, P.; Wang, Y.; Chu, J.-W.; Ayton, G. S.; Izvekov, S.; Andersen, H. C.; Voth, G. A. The multiscale coarse-graining method. II. Numerical implementation for coarse-grained molecular models. *J. Chem. Phys.* **2008**, *128*, 244115.

(56) Wang, H.; Junghans, C.; Kremer, K. Comparative atomistic and coarse-grained study of water: What do we lose by coarse-graining? *Eur. Phys. J. E* **2009**, *28*, 221–229.

(57) Fu, X.; Wu, Z.; Wang, W.; Xie, T.; Keten, S.; Gomez-Bombarelli, R.; Jaakkola, T. Forces are not Enough: Benchmark and Critical Evaluation for Machine Learning Force Fields with Molecular Simulations. AI for Science: Progress and Promises Workshop at NeurIPS. New Orleans, LA, USA, Dec. 2, 2022.

(58) Das, A.; Andersen, H. C. The multiscale coarse-graining method. III. A test of pairwise additivity of the coarse-grained potential and of new basis functions for the variational calculation. *J. Chem. Phys.* **2009**, *131*, 034102.

(59) Husic, B. E.; Charron, N. E.; Lemm, D.; Wang, J.; Pérez, A.; Krämer, A.; Chen, Y.; Olsson, S.; de Fabritiis, G.; Noé, F.; Clementi, C. Coarse graining molecular dynamics with graph neural networks. *J. Chem. Phys.* **2020**, *153*, 194101.

(60) Thaler, S.; Stupp, M.; Zavadlav, J. Deep coarse-grained potentials via relative entropy minimization. *J. Chem. Phys.* **2022**, *157*, 244103.

(61) Hoffman, M. D.; Gelman, A. The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *J. Mach. Learn. Res.* **2014**, *15*, 1593–1623.

(62) Gelman, A.; Lee, D.; Guo, J. Stan: A probabilistic programming language for Bayesian inference and optimization. *J. Educ. Behav. Stat.* **2015**, *40*, 530–543.

(63) Lao, J.; Louf, R. Blackjax: A sampling library for JAX. 2020; `http://github.com/blackjax-devs/blackjax`.

(64) Tran, K.; Neiswanger, W.; Yoon, J.; Zhang, Q.; Xing, E.; Ulissi, Z. W. Methods for comparing uncertainty quantifications for material property predictions. *Mach. Learn. Sci. Technol.* **2020**, *1*, 025006.

(65) Chung, Y.; Char, I.; Guo, H.; Schneider, J.; Neiswanger, W. Uncertainty Toolbox: an Open-Source Library for Assessing, Visualizing, and Improving Uncertainty Quantification. *arXiv preprint arXiv:2109.10254* **2021**,

(66) Hirschfeld, L.; Swanson, K.; Yang, K.; Barzilay, R.; Coley, C. W. Uncertainty quantification using neural networks for molecular property prediction. *J. Chem. Inf. Model.* **2020**, *60*, 3770–3780.

(67) Abascal, J. L. F.; Vega, C. A general purpose model for the condensed phases of water: TIP4P/2005. *J. Chem. Phys.* **2005**, *123*, 234505.

(68) Berendsen, H. J. C.; Postma, J. P. M.; van Gunsteren, W. F.; Hermans, J. Interaction models for water in relation to protein hydration. Intermolecular forces: Proceedings of the Fourteenth Jerusalem Symposium on Quantum Chemistry and Biochemistry. Jerusalem, Israel, Apr. 13–16, 1981; pp 331–342.

(69) Chaimovich, A.; Shell, M. S. Anomalous waterlike behavior in spherically-symmetric water models optimized with the relative entropy. *Phys. Chem. Chem. Phys.* **2009**, *11*, 1901–1915.

(70) Potestio, R.; Peter, C.; Kremer, K. Computer Simulations of Soft matter: Linking the Scales. *Entropy* **2014**, *16*, 4199–4245.

(71) Kingma, D. P.; Ba, J. L. Adam: A Method for Stochastic Optimization. 3rd International Conference on Learning Representations, ICLR. San Diego, CA, USA, May 7-9, 2015.

(72) Thaler, S.; Zavadlav, J. Uncertainty Quantification for Molecular Models via Stochastic Gradient MCMC. 10th Vienna Conference on Mathematical Modelling. Vienna, Austria, Jul. 27 –29, 2022; pp 19–20.

(73) Pettitt, B. M.; Karplus, M. The potential of mean force surface for the alanine dipeptide in aqueous solution: a theoretical approach. *Chem. Phys. Lett.* **1985**, *121*, 194–201.

(74) Tobias, D. J.; Brooks III, C. L. Conformational equilibrium in the alanine dipeptide

in the gas phase and aqueous solution: A comparison of theoretical results. *J. Phys. Chem.* **1992**, *96*, 3864–3870.

(75) Shell, M. S. The relative entropy is fundamental to multiscale and inverse thermodynamic problems. *J. Chem. Phys.* **2008**, *129*, 144108.

(76) Chaimovich, A.; Shell, M. S. Coarse-graining errors and numerical optimization using a relative entropy framework. *J. Chem. Phys.* **2011**, *134*, 094112.

(77) Schneider, W.; Thiel, W. Anharmonic force fields from analytic second derivatives: Method and application to methyl bromide. *Chem. Phys. Lett.* **1989**, *157*, 367–373.

(78) Barducci, A.; Bonomi, M.; Parrinello, M. Metadynamics. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **2011**, *1*, 826–843.

(79) Chaimovich, A.; Shell, M. S. Relative entropy as a universal metric for multiscale errors. *Physical Review E* **2010**, *81*, 060104.

(80) Toxvaerd, S. Hamiltonians for discrete dynamics. *Phys. Rev. E* **1994**, *50*, 2271.

(81) Toxvaerd, S. Ensemble simulations with discrete classical dynamics. *J. Chem. Phys.* **2013**, *139*, 224106.

(82) Zhang, R.; Cooper, A. F.; De Sa, C. AMAGOLD: Amortized Metropolis adjustment for efficient stochastic gradient MCMC. International Conference on Artificial Intelligence and Statistics. Online, Aug. 26–28, 2020; pp 2142–2152.

(83) Deng, W.; Feng, Q.; Gao, L.; Liang, F.; Lin, G. Non-convex learning via replica exchange stochastic gradient mcmc. Proceedings of the 37th International Conference on Machine Learning. Online, Jul. 13–18, 2020; pp 2474–2483.

(84) Gorham, J.; Mackey, L. Measuring sample quality with Stein's method. Advances in Neural Information Processing Systems. Montreal, QC, Canada, Dec. 7–12, 2015.

(85) Garriga-Alonso, A.; Fortuin, V. Exact Langevin Dynamics with Stochastic Gradients. 3rd Symposium on Advances in Approximate Bayesian Inference. Online, Jan. – Feb., 2021.

(86) Zhang, R.; Li, C.; Zhang, J.; Chen, C.; Wilson, A. G. Cyclical Stochastic Gradient MCMC for Bayesian Deep Learning. 7th International Conference on Learning Representations. New Orleans, LA, USA, May 6–9, 2019.

(87) Heek, J.; Kalchbrenner, N. Bayesian Inference for Large Scale Image Classification. 8th International Conference on Learning Representations. Addis Ababa, Ethiopia, Apr. 26–30, 2020.

(88) Korattikara Balan, A.; Rathod, V.; Murphy, K. P.; Welling, M. Bayesian dark knowledge. Advances in Neural Information Processing Systems. Montreal, QC, Canada, Dec. 7–12, 2015.

(89) Wang, L.; Yoon, K.-J. Knowledge Distillation and Student-Teacher Learning for Visual Intelligence: A Review and New Outlooks. *IEEE Trans. Pattern Anal. Mach. Intell.* **2021**, *44*, 3048–3068.