

Enhanced Diffusion Sampling: Efficient Rare Event Sampling and Free Energy Calculation with Diffusion Models

Yu Xie[†], Ludwig Winkler[†], Lixin Sun[†], Sarah Lewis, Adam E. Foster,
 José Jiménez Luna, Tim Hempel, Michael Gastegger, Yaoyi Chen,
 Iryna Zaporozhets, Cecilia Clementi, Christopher M. Bishop, Frank Noé
Microsoft Research AI for Science

[†] Equal contribution.

Abstract

The rare-event sampling problem has long been the central limiting factor in molecular dynamics (MD), especially in biomolecular simulation. Recently, diffusion models such as BioEmu have emerged as powerful equilibrium samplers that generate independent samples from complex molecular distributions, eliminating the cost of sampling rare transition events. However, a sampling problem remains when computing observables that rely on states which are rare in equilibrium, for example folding free energies. Here, we introduce enhanced diffusion sampling, enabling efficient exploration of rare-event regions while preserving unbiased thermodynamic estimators. The key idea is to perform quantitatively accurate steering protocols to generate biased ensembles and subsequently recover equilibrium statistics via exact reweighting. We instantiate our framework in three algorithms: UmbrellaDiff (umbrella sampling with diffusion models), MetaDiff (a batchwise analogue for metadynamics), and ΔG -Diff (free-energy differences via tilted ensembles). Across toy systems, protein folding landscapes and folding free energies, our methods achieve fast, accurate, and scalable estimation of equilibrium properties within GPU-minutes to hours per system—closing the rare-event sampling gap that remained after the advent of diffusion-model equilibrium samplers.

1 Introduction

Molecular dynamics (MD) simulation is a widely used computational approach for generating molecular equilibrium ensembles $p(x)$ and predicting experimental observables $O = \mathbb{E}_{p(x)}[o(x)]$, but its effectiveness is limited by the sampling problem, which consists of two distinct components (Table 1).

1. *Slow mixing problem*—MD produces time-correlated trajectories x_t . Long-lived states or phases lead to trapping the simulation trajectory for long times, resulting in slow exploration and slow convergence of expectation values.
2. *Rare state problem*—even with independent draws from $p(x)$, it can be prohibitive to sample states with small equilibrium probabilities. For example, the probability ratio of unfolded and folded protein states depends exponentially on the folding free energy: $p_u/p_f = \exp(\Delta G_{\text{fold}}/k_B T)$. At 300K, $\Delta G_{\text{fold}} = -5 \text{ kcal/mol}$ implies that ~ 1 in 4.4×10^3 equilibrium samples is unfolded. For a moderately stable protein ($\Delta G_{\text{fold}} = -10 \text{ kcal/mol}$) only ~ 1 in 1.9×10^7 samples is unfolded.

These limitations have motivated enhanced sampling methods over the last 70 years [1], which address rare states by sampling from a biased distribution and then reweighting to recover equilibrium statistics; however, when implemented on top of MD they can remain limited by slow mixing of the unbiased degrees of freedom (Table 1). Recently, generative equilibrium samplers based on normalizing flows and diffusion models [2, 3] have emerged that generate approximately independent equilibrium configurations, removing the slow-mixing bottleneck, but rare-state estimation remains when observables depend on low-probability

	Equilibrium sampling	Enhanced sampling
Molecular Dynamics	Slow mixing problem Rare state problem	Slow mixing problem Rare states sampled
Diffusion samplers	Independent samples Rare state problem	Independent samples Rare states sampled

Table 1: The MD sampling problem consists of a slow mixing problem due to rare interconversion between long-lived states, and a rare state problem as low-probability states are infrequently visited. Diffusion equilibrium samplers tackle the slow mixing problem, enhanced sampling methods tackle the rare state problem. In this paper we explore the combination of both: enhanced diffusion samplers.

regions of $p(x)$ (Table 1). This paper develops a framework for enhanced sampling with diffusion-model samplers, addressing both bottlenecks within a single approach.

Traditional enhanced sampling methods include thermodynamic integration [4], free energy perturbation (FEP) [5], umbrella sampling [6, 7], parallel or simulated tempering [8–10] and metadynamics [11]—see [1] for an extensive review. All these methods sample from a biased distribution, and then later remove this bias from the sampled statistics in order to recover equilibrium statistics [12–14]. They can accelerate sampling by orders of magnitude when suitable collective variables or thermodynamic controls are available. Representative successes include: (i) *Free energy profiles of reactions and of ion permeations through channels* [15–18], where umbrella sampling can restrain sampling along the well-defined reaction coordinate, while the other degrees of freedom relax quickly; (ii) *small-molecule solvation and protein–ligand binding free energies* [19–21], where alchemical methods relying on free energy perturbation (FEP) connect nearby thermodynamic states with feasible local sampling; (iii) *small protein folding in implicit solvent* [22], where replica exchange remains tractable; and (iv) *mutation series in coarse-grained models* [23], where reduced resolution makes otherwise prohibitive transitions amenable to free-energy calculations.

For high-dimensional biomolecular transitions—including protein folding, binding, and conformational changes in explicit solvent—enhanced sampling is often limited by two coupled issues. First, suitable low-dimensional bias coordinates are frequently unknown *a priori* and may only become apparent after substantial sampling. This has motivated adaptive approaches that iteratively discover reaction coordinates and enhance sampling along them [24–32]. Second, these systems typically exhibit a *spectrum* of slow relaxation processes rather than a single dominant timescale, as explored in the MSM literature [33, 34]; when slow modes are weakly separated, long simulations remain necessary to equilibrate degrees of freedom not directly controlled by the bias.

Consequently, successes on explicit-solvent biomolecular problems have often required specialized combinations of methods and/or massive compute. For *all-atom protein folding* free-energy landscapes, temperature replica exchange becomes increasingly inefficient in explicit solvent because the number of replicas grows with system size, and practical studies have relied on hybrids such as REMD+metadynamics [35, 36], bias-exchange metadynamics [37], or multitemperature MD strategies [38], as well as special-purpose hardware or massively distributed simulations plus MSM analysis to obtain quantitative folding landscapes for proteins up to ~ 100 residues [39–42]. For *complex conformational changes*, well-characterized free-energy landscapes have been obtained using metadynamics and its variants [43, 44], string-based approaches combined with umbrella sampling or swarms of trajectories [45, 46], and temperature-accelerated MD in collective variables [47]. Large-scale unbiased simulation with MSMs has also enabled quantitative conformational landscapes in challenging systems [48–52], and brute-force simulations on specialized hardware have resolved long-timescale protein dynamics and GPCR activation/binding mechanisms [53–55]. Finally, while alchemical FEP is well established for relative protein–ligand binding or closely related mutants, it becomes impractical when the alchemical changes induce large-scale conformational rearrangements that are slow and hard to overlap.

Recently, a complementary line of work has emerged from generative deep learning: *Boltzmann generators and Boltzmann emulators* that aim to generate approximately independent configurations from an equilibrium distribution without integrating long MD trajectories. Boltzmann Generators are flow-based approaches that aim to generate independent samples from a Boltzmann distribution defined via an energy function [2, 56]. Boltzmann emulators approximate the equilibrium distribution by training on MD simulation data and

fine-tuning on experimental observables. This approach has recently become popular for sampling protein ensembles, e.g., BioEmu [3, 57–60]. The concept is also emerging in other application areas, such as material sciences [59]. By producing effectively iid samples, these models address the *slow-mixing* bottleneck in Table 1.

However, iid sampling does not remove the *rare-state* bottleneck: estimating observables controlled by low-probability regions still requires a number of samples that scales exponentially in free-energy differences. For example, BioEmu-1 estimates protein fold stabilities by directly sampling the equilibrium ensemble of protein structures and calculating the fraction of unfolded states [3]. For stabilities of up to $\Delta G_{\text{fold}} = -5$ kcal/mol this is still feasible within one GPU hour, but this approach quickly becomes intractable: sampling ~ 1 in 1.9×10^7 unfolded samples for a protein with $\Delta G_{\text{fold}} = -10$ kcal/mol would take on the order of one GPU year [3].

Here we provide a solution for the remaining sampling problems of diffusion-based equilibrium samplers by integrating enhanced sampling models into the diffusion model framework. The key component is a steering algorithm that allows to apply the desired bias potentials to a pretrained diffusion model at inference time. This enables the implementation of several classical enhanced sampling methods in the diffusion model framework, as well as the application of unbiasing methods such as WHAM [13] or MBAR [14]. Furthermore, several technical improvements are presented that provide low-variance estimates of the desired observables even if only a few thousand diffusion-model denoising trajectories can be afforded. Overall, this framework enables converged sampling of equilibrium properties of complex biomolecular processes with large energy differences within GPU minutes to hours per calculation, given a suitable pretrained diffusion model. We demonstrate our framework using the BioEmu model on a biomolecular rare-event problem that is extremely difficult or currently impossible with all-atom MD simulations: efficient calculation of folding free energies for a variety of proteins ranging between 50 and 200 amino acids.

Related ideas of combining enhanced sampling methods with diffusion models have been presented in several recent publications: Ref. [61] introduces a related approach employing adjoint sampling to fine-tune a pre-trained diffusion model towards the sampling of rare states. Ref. [62] introduces a steering approach for exploring rare events with pretrained diffusion models with an approach similar to umbrella-sampling that can be viewed as special case of the present approach. Ref. [63] have presented a steering method which changes the sampled ensemble to be consistent with experimental data or enhance sample diversity.

2 Biasing and unbiasing diffusion model ensembles

We assume that we have a pretrained diffusion model whose output distribution for a given molecule of interest is $p(x)$ — i.e., the model distribution $p(x)$ shall be our *unbiased equilibrium distribution*. We will equivalently operate with probabilities and dimensionless energies, for example:

$$p(x) = e^{-u(x)}, \quad (1)$$

with dimensionless energy $u(x)$. This representation is independent of the thermodynamic ensemble, e.g., in the canonical ensemble, $p(x)$ is the Boltzmann distribution and $u(x) = U(x)/k_B T$ with potential energy $U(x)$, Boltzmann constant k_B and temperature T . Subsequently we will only assume that $p(x)$ and $u(x)$ exist, but not that they can be explicitly evaluated. Energy-based diffusion models, for which $p(x)$ and $u(x)$ can be efficiently evaluated, have been explored recently [64, 65].

Our aim is to compute unbiased expectation values under the equilibrium distribution of the form

$$O = \mathbb{E}_p[o(x)]. \quad (2)$$

For example, the probability of being in the unfolded state, p_u , can be cast in this form, but accurately estimating it—and thus the folding free energy, relies on sampling rare events. As in the traditional enhanced sampling literature, we proceed in two steps:

1. **Generate biased ensembles:** given one or multiple biasing potentials $b_k(x)$, $k = 1, \dots, K$, draw samples from the biased ensembles with energies $u(x) + b_k(x)$. This is achieved by steering the diffusion model at inference time (Section 2.1).

2. **Recover unbiased expectations:** reweight the biased samples given the known bias potentials in order to recover the unbiased expectation values. For $K = 1$, reweighting is trivial, whereas for $K > 1$ it can be performed with the multistate Bennett acceptance ratio (MBAR) method (Section 2.2).

The choice of biasing potentials and strategies will be discussed in Sections 3-5.

2.1 Biased sampling from diffusion models

Given a pretrained diffusion model which samples from $p(x)$, we seek a method to draw samples from a single biased ensemble defined via biasing potential $b(x)$:

$$q(x) := \frac{1}{Z} p(x) e^{-b(x)}, \quad Z := \int p(x) e^{-b(x)} dx, \quad (3)$$

where Z is an unknown normalization constant (partition function) that we will not need to evaluate explicitly.

Several established families of steering methods exist: (i) *Score guidance*, which alters the reverse-time score by adding the gradient of a classifier, constraint, or reward to tilt generations toward desired properties [66–68]. (ii) *Path-integral reweighting* via Sequential Monte Carlo (SMC) with Feynman–Kac (FK) potentials—keeping the proposal dynamics unchanged and accounting for the bias with incremental importance weights along the reverse trajectory; annealed importance sampling (AIS) appears as the no-resampling special case [69–72]. (iii) *Invariant correctors* (e.g., Metropolis-adjusted Langevin, MALA) that are interleaved with the reverse dynamics to mix within the current biased marginal without changing the FK weights [71, 73]. Score-guidance methods are simple and fast but generally do not guarantee unbiased sampling of the target density. FK, AIS and SMC give explicit importance weights and, with resampling, can target the desired biased marginals exactly.

Subsequently, we employ the Feynman-Kac Corrector (FKC) biased sampling methodology from [74] and [75], described in the following paragraphs.

Unbiased diffusion model. We assume our diffusion model has been trained to reverse the corruption process defined by a stochastic differential equation,

$$dx_\tau = f_{1-\tau}(x_\tau) d\tau + \sigma_{1-\tau} dW_\tau, \quad x_{\tau=0} \sim p_{\text{data}} \quad (4)$$

Here W is a standard Wiener process. The drift f and diffusion σ are chosen at training time such that $x \sim p_{\text{noise}} = N(0, I)$ at $\tau = 1$. We choose to write f and σ with $1 - \tau$ as their argument rather than τ for notational convenience when writing about the reverse (denoising) process. We parameterize reverse (denoising) time by $t := 1 - \tau$ and write p_t for the marginal density of x_τ at $\tau = 1 - t$. Thus, $p_0 = p_{\text{noise}}$ and $p_1 = p_{\text{data}}$.

Training the diffusion model consists of learning to compute the score $\nabla \log p_t(x)$ as a function of x and t . Sampling from the diffusion model is performed by first sampling x from p_{noise} and then simulating a reverse process of the form

$$dx_t = g_t(x_t) dt + \tilde{\sigma}_t dW_t \quad (5)$$

where

$$g_t(x) := -f_t(x) + \frac{\sigma_t^2 + \tilde{\sigma}_t^2}{2} \nabla \log p_t(x), \quad (6)$$

and $\tilde{\sigma}_t$ is a hyperparameter we can choose, with $\tilde{\sigma}_t = \sigma_t$ being the standard choice used in [74].

Biased marginals q_t . In order to sample from $q(x)$, we define a family of biased marginal densities q_t :

$$q_t(x) := \frac{p_t(x) e^{-b_t(x)}}{Z_t}, \quad Z_t := \int p_t(x) e^{-b_t(x)} dx \quad (7)$$

where $b_0(x)$ is constant with respect to x , $b_1(x) \equiv b(x)$, and b_t smoothly interpolates between the two.

Generating weighted samples from q_t . Ref. [74] provides methods to modify the dynamics in Eq. 5 so that instead of sampling from p_t , we get weighted samples $(x_1, w_1), \dots, (x_N, w_N)$ from biased marginals q_t . The weights w_1, \dots, w_N are importance weights, encoding the discrepancy between proposal and target biased marginals.

In order to sample from the biased marginals in Eq. 7, we need a slightly more general version of the reward-tilted SDE of [74]. We get this by using control drift $\frac{\tilde{\sigma}_t^2}{2} \nabla b_t$ in the framework of [75], resulting in the following system of SDEs:

$$x_0 \sim p_{\text{noise}} \quad (8)$$

$$dx_t = \left(-f_t + \frac{\sigma_t^2 + \tilde{\sigma}_t^2}{2} \nabla \log p_t(x_t) + \frac{\tilde{\sigma}_t^2}{2} \nabla b_t(x_t) \right) dt + \tilde{\sigma}_t dW_t \quad (9)$$

$$d \log w_t = F_t(x_t) dt - \mathbb{E}_{x \sim q_t}[F_t(x)] \quad (10)$$

$$F_t(x) := -\frac{\partial b_t(x)}{\partial t} + \nabla b_t(x) \cdot \left(f_t(x) - \frac{\sigma_t^2}{2} \nabla \log p_t(x) \right) \quad (11)$$

If we choose $\sigma_t \equiv \tilde{\sigma}$ and $b_t(x) \equiv \lambda_t b(x)$ we recover the reward-tilted SDE from [74].

Effective sample size We track the statistical efficiency of a weighted batch of n samples using the Kish effective sample size (ESS) [76]:

$$\text{ESS} = \frac{(\sum_i w_i)^2}{\sum_i w_i^2} \quad (12)$$

ESS ranges between 1 (all weight on one sample) and n (all weights equal). If the weights are normalized as $\sum_i w_i = 1$, it is $\text{ESS} = 1 / \sum_i w_i^2$.

Resampling During denoising we periodically resample the particles using stratified sampling. Resampling can be done periodically, but should be done if the ESS falls below a predefined threshold (e.g., $n/2$ [71]). After each resampling step, the weights are refreshed, i.e.,

$$(x_t^1, w_t^1), \dots, (x_t^K, w_t^K) \rightarrow (\tilde{x}_t^1, 1), \dots, (\tilde{x}_t^K, 1). \quad (13)$$

Optionally, we can enforce a resampling step after the steering procedure at $t = 0$ per biased ensemble. In that case, all steered samples have uniform weights 1 and the subsequent expressions simplify.

2.2 Unbiasing to the equilibrium distribution

After producing samples from K biased ensembles (Sec. 2.1), we have a set of weighted samples $\{(x_{k,i}, w_{k,i})\}_{i=1}^{N_k}$ from each biased density q_k . Next we combine all biased samples into a single sample with weights W_n :

$$\{(x_{k,i}, w_{k,i})\}_{i=1}^{N_k} \rightarrow \{(x_n, W_n)\}_{n=1}^N. \quad (14)$$

The reweighting coefficients W_n combine (i) the steering importance weights $w_{k,i}$ within each biased density q_k produced via (8-11) and (ii) an additional unbiasing factor that maps each biased ensemble q_k back to p . If we have chosen to do a terminal resampling step at the end of steering, $w_{k,i} \equiv 1$ and W_n only contains the unbiasing from q_k to p .

Equilibrium expectation values The weights W_n must be such that we can express equilibrium expectation values of arbitrary observables $o(x)$ as:

$$\widehat{\mathbb{E}}_p[o] = \frac{\sum_{n=1}^N W_n o(x_n)}{\sum_{n=1}^N W_n}. \quad (15)$$

in a way that is asymptotically unbiased in the infinite sample limit, i.e., $\lim_{N \rightarrow \infty} \widehat{\mathbb{E}}_p[o] = \mathbb{E}_p[o]$. Expression (15) is sufficiently general that all quantities of the unbiased ensemble can be expressed with it. For example,

to compute the folding free energy we define a membership function which assigns folded states to 1 and unfolded states to 0:

$$\chi(x) = \begin{cases} 0 & x \text{ unfolded} \\ 1 & x \text{ folded} \end{cases}, \quad (16)$$

$\chi(x)$ can itself be defined in terms of suitable order parameters, such as the fraction of native contacts from the folded structure, and it can optionally take intermediate values between 0 and 1. In terms of χ , the folding free energy is defined by:

$$\Delta G_{\text{fold}} = -k_B T \ln \frac{p_{\text{folded}}}{p_{\text{unfolded}}} = -k_B T \ln \frac{\hat{\mathbb{E}}_p[\chi]}{1 - \hat{\mathbb{E}}_p[\chi]} \quad (17)$$

using (15). Likewise, a potential of mean force (PMF) along a reaction coordinate $\xi(x)$ can be formulated by, e.g., binning the value range of ξ and expressing the probability of each bin in terms of (15). Practically, we obtain PMFs by constructing a weighted histogram or kernel density estimator from the weighted sample $\{(x_n, W_n)\}_{n=1}^N$.

Single biased ensemble For the special case $K = 1$, i.e., using a single biased ensemble with bias potential $b(x)$, we have generated samples $\{(x_n, w_n)\}_{n=1}^N$ (dropping the window index). The combined importance weights are then trivially given by direct reweighting:

$$W_n = w_n e^{b(x_n)} \quad (18)$$

Unbiasing multiple biased ensembles with MBAR A statistically optimal method to combine samples from multiple biased ensembles towards unbiased ensemble estimators is the multistate Bennett acceptance ratio (MBAR), also known as binless weighted histogram analysis (WHAM) [14, 77, 78]. For the special case $K = 2$, MBAR becomes traditional BAR [12]. MBAR yields minimum-variance, asymptotically unbiased estimates without requiring us to evaluate $\log p(x)$. If the bias energy can be discretized we can instead work with histograms and binned WHAM [7, 13, 79].

Here we describe a weighted MBAR variant which takes the weighted sample terminals from steering, $\{(x_{k,i}, w_{k,i})\}_{i=1}^{N_k}$, and rewrites them into optimal reweighting coefficients for estimating expectations under p (Eq. 14). By convention, MBAR operates on $K + 1$ ensembles, consisting of the unbiased ensemble $b_0(x) \equiv 0$ and the biased ensembles $b_k(x)$ for $k = 1, \dots, K$. Note that we do not draw samples from ensemble 0, so sums run over sampled windows $\ell = 1, \dots, K$. To deal with weighted samples, we define the per-sample effective mass $\alpha_{k,i} \propto w_{k,i}$ and normalize so the window mass sums to the effective sample size of ensemble k , ESS_k (Eq. 12):

$$\alpha_{k,i} = \text{ESS}_k \frac{w_{k,i}}{\sum_{j=1}^{N_k} w_{k,j}} \implies \sum_{i=1}^{N_k} \alpha_{k,i} = \text{ESS}_k. \quad (19)$$

If we have chosen to do terminal resampling at the end of steering, the input samples are already unweighted, $w_{k,i} \equiv 1$, $\text{ESS}_k = N_k$, hence $\alpha_{k,i} = 1$ and we recover standard MBAR.

We then aggregate all samples as $\{x_n\}_{n=1}^N$ with associated α_n (from the bias potential with which they were generated). Let $M_\ell = \sum_{n: \text{orig}(x_n)=\ell} \alpha_n$ denote the window mass, where $M_\ell = \text{ESS}_\ell$ under the default above, and $\text{orig}(x_n)$ denotes the index of the window in which sample x_n was generated. We solve for the reduced free energies of the sampled biased states $k = 1, \dots, K$ via

$$\exp(-\hat{f}_k) = \sum_{n=1}^N \frac{\alpha_n \exp(-b_k(x_n))}{\sum_{\ell=1}^K M_\ell \exp(\hat{f}_\ell - b_\ell(x_n))}, \quad k = 1, \dots, K, \quad (20)$$

and forms target-state weights ($a \in \{0, 1, \dots, K\}$), with the unbiased state $a = 0$ included as the target for reweighting, as

$$W_n^{(a)} = \frac{\alpha_n \exp(-b_a(x_n))}{\sum_{\ell=1}^K M_\ell \exp(\hat{f}_\ell - b_\ell(x_n))}. \quad (21)$$

We then compute unbiased expectations in (15) using

$$W_n \equiv W_n^{(0)}. \quad (22)$$

Using $\alpha_{k,i} \propto w_{k,i}$ with the normalization above lets MBAR automatically down-weight windows with highly skewed particle weights via their ESS_k . If desired, one can replace ESS_k by an external *effective count* (e.g. number of denoising trajectories used for window k , or a correlation-corrected N_k/g_k); equations (20)–(22) are unchanged.

To compute uncertainties, we can use MBAR’s covariance expressions [14], or apply a cluster bootstrap that resamples at the level of independent denoising trajectories to account for within-window correlations, e.g., from late branching.

2.3 Illustration of enhanced diffusion sampling

Before moving on to more complex enhanced sampling protocols, we first demonstrate that enhanced diffusion sampling is efficient in an idealized setting (Fig. 1). We define a family of double-well potentials $u(x)$ where the energy difference between the two minima is given by ΔG ranging from -2 to -14 $k_B T$. The corresponding equilibrium densities are given by $p(x) \propto e^{-u(x)}$, resulting in an equilibrium probabilities between 10^{-1} and 10^{-6} for the high-energy state (Fig. 1a). For each ΔG value, a diffusion model with an analytical score function is defined that generates samples from $p(x)$ exactly. To evaluate the effect of enhanced diffusion sampling on sampling efficiency directly, we defined a single linear bias potential for each double-well potential, $b_{\Delta G} = a(\Delta G)x$ where $a(\Delta G)$ is chosen such that the two energy minima are equal in the biased potential (Fig. 1b).

For each ΔG value we then generate N samples with each of two protocols: (i) direct sampling from the unsteered diffusion model with density $p(x)$, and (ii) steered sampling as described in Sec. 2.1 using $b_{\Delta G}$ and direct reweighting of the biased ensemble with (18). We then estimate ΔG from the samples and compare with the exact value. Both protocols converge to the correct ΔG value in the large N limit (Fig. 1c). Enhanced diffusion sampling converges faster than equilibrium sampling for all inspected ΔG values and the performance gap increases with ΔG . To compare sampling efficiencies between unbiased and enhanced sampling protocols, we run five independent repeats for each N and each protocol and define a sample as being converged when both the absolute difference of the empirical mean of ΔG from the exact value, and the empirical standard deviation of ΔG are within 1 kcal/mol (at 300K) across repetitions. As expected, the number of samples needed to reach convergence with unbiased sampling increases approximately exponentially with ΔG (Fig. 1d, blue). In contrast, the number of samples required by enhanced diffusion sampling only increases mildly with ΔG and ranges between 10 to 100 samples for all ΔG values inspected here (Fig. 1d, orange).

We note that this setup describes the best-case scenario in which only a single biasing potential is needed and the perfect bias is known in advance. In a realistic scenario we will have to try different biasing potentials until all desired states have been sampled, and then combine them with MBAR, which will lead to a greater number of samples needed to convergence and will complicate the relationship between free energy differences and the number of samples needed.

3 UmbrellaDiff — umbrella sampling with diffusion models

Here we describe how the classical umbrella sampling method [6] can be adapted to steered diffusion models. The goal of umbrella sampling is to compute the marginal probability density $p_{\Xi}(\xi) = \int p(x)\delta(\xi(x) - \xi)dx$, or free energy profile—also known as potential of mean force (PMF),

$$u_{\Xi}(\xi) = -\ln p_{\Xi}(\xi), \quad U_{\Xi}(\xi) = k_B T u_{\Xi}(\xi), \quad (23)$$

along a user-chosen reaction coordinate $\xi : \mathcal{X} \rightarrow \mathbb{R}$, e.g., the end-to-end distance of a protein, or the fraction of native contacts. We will formulate umbrella sampling for the simplest case of one-dimensional coordinates, but the extension to multidimensional umbrella sampling is straightforward [7].

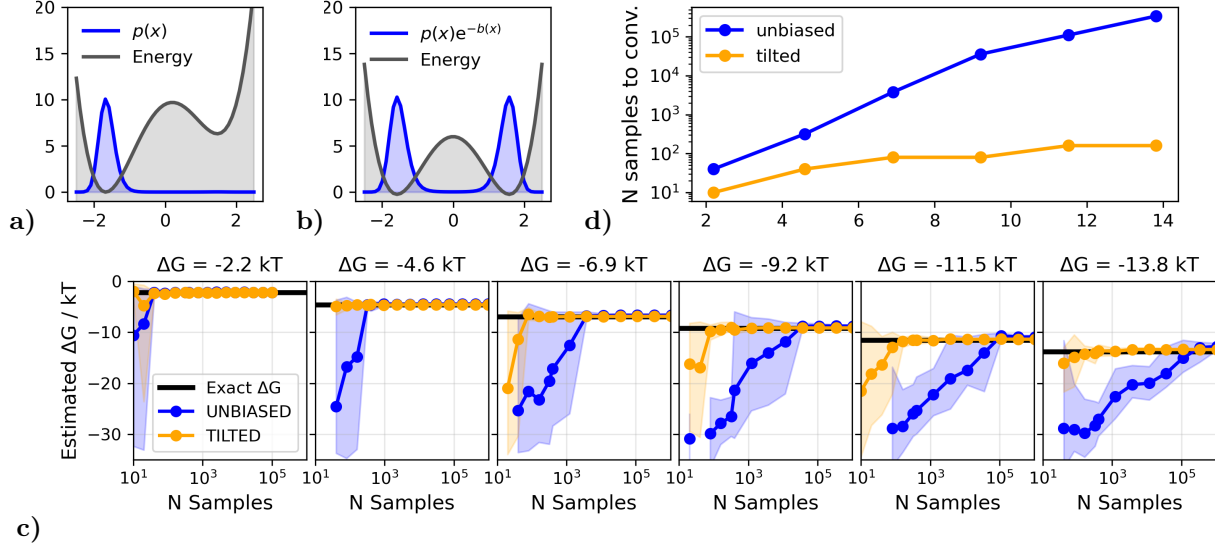


Figure 1: **Enhanced sampling with diffusion models efficiently computes probabilities of rare events** **a)** Unbiased diffusion model probability density $p(x)$ and corresponding energy $-\ln p(x)$. In this example the free energy difference between the two minima is set to $\Delta G = -7 k_B T$. **b)** Biased density $p(x)e^{-b(x)}$ using a linear tilt $b(x) = \beta x$ with β chosen to achieve $\Delta G_{\text{biased}} = 0 k_B T$. **c)** Estimates of ΔG with an unbiased diffusion model and enhanced sampling using a tilting potential that sets $\Delta G_{\text{biased}} = 0 k_B T$. Mean (solid line) and 95% confidence interval (shaded area) are shown. **d)** Number of samples required to get an estimate within 1 kcal/mol of the exact ΔG value and with a standard deviation of 1 kcal/mol, as a function of ΔG .

Umbrella sampling introduces a set of K biased ensembles, each of which restrains sampling near prescribed centers c_k in ξ -space. The most common bias potential is harmonic with stiffness $\kappa_k > 0$:

$$b_k(x) = \frac{1}{2} \kappa_k \|\xi(x) - c_k\|^2, \quad (24)$$

although other choices are possible, e.g., flat-bottom with harmonic tails. For each bias potential, we generate an ensemble using the steering procedure of Sec. 2.1 and remove bias with weighted MBAR (Sec. 2.2). MBAR returns a set of equilibrium weights for each sample, (x_n, W_n) which can then be used in a histogram or kernel density estimate to obtain $p_{\Xi}(\xi)$. Finally we can obtain the PMF via Eq. (23). In order to obtain smooth estimates of $p_{\Xi}(\xi)$ it can be beneficial to increase the number of samples using late branching near $t \approx 0$ (Sec. 2.1).

Designing bias potentials to ensure coverage and overlap. We choose centers $\{c_k\}_{k=1}^K$ to cover the ξ -range of interest. Our prior model of the biased density is that induced by the bias potential alone, i.e., a Gaussian with center c_k and standard deviation $\sigma_k \approx 1/\sqrt{\kappa_k}$. As a rule of thumb, space neighboring centers within one or two σ to obtain ~ 10 – 30% overlap along ξ .

As the effect of $p(x)$ will change the resulting distribution of $q_k(x)$, the bias parameters, may be chosen adaptively, e.g., with a preconditioning such as: (i) start with evenly spaced c_k and $\kappa_k = 1/(c_k - c_{k-1})^2$, (ii) run a small initial sample and estimate empirical $\text{Var}_{q_k}[\xi]$, (iii) update $\kappa_k \leftarrow 1/\widehat{\text{Var}}_{q_k}[\xi]$ to equalize widths.

Diagnostics We recommend to monitor standard diagnostics from umbrella sampling / free-energy analysis to assess (a) whether each biased ensemble is well sampled and (b) whether the set of windows forms a connected path that permits low-variance reweighting [1, 80, 81]. In particular the following diagnostics are suitable:

1. Per-window effective sample size (ESS) of the steering terminal weights $w_{k,i}$ (Eq. 12); when terminals may be correlated (e.g. due to late branching or MCMC rejuvenation), we additionally report

a correlation-corrected effective count $N_k^{\text{eff}} \approx N_k/g_k$ using the statistical inefficiency g_k and use a cluster/bootstrap analysis for uncertainties [82].

2. Window overlap. Along the umbrella coordinate ξ , we quantify overlap between neighboring windows using histogram/KDE-based overlap coefficients (e.g. $O_{k,k+1} = \int \sqrt{q_k(\xi)q_{k+1}(\xi)}d\xi$ or related normalized measures). More generally, we compute the *state overlap matrix* O_{ij} from MBAR/WHAM, where O_{ij} estimates the probability that a configuration sampled from state j could be observed in state i . As a practical rule, require neighboring-window overlap area in the range 0.1 to 0.3 and a well-connected overlap matrix (no near-block-diagonal structure) by inserting additional windows or revising bias parameters [14, 81, 83]
3. MBAR/WHAM uncertainty estimates, including the asymptotic covariance/standard errors for free energies and target observables [14, 79], and stability of the PMF/observables under increasing sample size (and optionally leave-one-window-out checks). For umbrella sampling in particular, we optionally complement MBAR with EMUS error analysis to identifying windows that dominate variance and guide adaptive allocation of sampling effort [84].

These diagnostics can be used to adapt window centers and force constants, increase sampling in selected windows, or insert bridging windows to restore overlap [80, 81].

UmbrellaDiff versus traditional umbrella sampling with MD Compared to traditional umbrella sampling with MD, UmbrellaDiff has several advantages. First, in conventional umbrella sampling one must prepare initial configurations for each window (often via pulling/steered protocols). Then one must equilibrate sufficiently so that neighboring windows overlap not only in the restrained coordinate ξ but also in orthogonal, slow degrees of freedom [80, 81, 85]. This is not an issue with UmbrellaDiff, as each steered sample is independently generated from the noisy state using the procedure in Sec. 2.1.

Secondly, when hidden barriers or multiple channels exist in directions perpendicular to ξ , traditional umbrella simulations can remain trapped in distinct metastable “orthogonal states”, leading to hysteresis and very slow convergence of the PMF unless additional sampling machinery is introduced (Fig. 2) [86, 87]. This is not an issue for UmbrellaDiff, where the only role of the umbrellas is to ensure that low probability values in $p_{\Xi}(\xi)$ are sampled, while the iid sampling property of the diffusion model avoid problems with kinetic trapping and slow intra-window relaxation. Fig. 2 demonstrates that UmbrellaDiff can compute unbiased free energy profiles even if metastable states that are off of the main path of umbrellas in ξ contribute significantly to the free energy, whereas traditional umbrella sampling struggles with these cases.

4 MetaDiff — metadynamics with diffusion models

Metadynamics [11, 88, 89] is a widely used enhanced sampling method with many follow-up variants, including well-tempered metadynamics [90, 91] and associated estimators [92]. The original idea of metadynamics is to run an MD simulation and to periodically deposit repulsive *hills* (typically Gaussians) in the space of chosen collective variables (CVs) $\xi(x)$ at the current CV value. The resulting history-dependent bias progressively fills free-energy minima and drives exploration of new regions [11, 88]. In standard metadynamics, the bias keeps growing, while in well-tempered metadynamics the hill heights are tempered so that the bias converges (up to a known factor) to the negative free energy [90, 91]. Because the bias is explicitly time-dependent, metadynamics is generally analyzed via the accumulated bias and dedicated reweighting/estimators rather than by treating the trajectory as equilibrium sampling at a fixed thermodynamic state; in particular, during the initial “fill-up” period it is a nonequilibrium process, so equilibrium multi-state estimators such as MBAR are not directly applicable [88, 92].

Here we formulate a metadynamics variant for steered diffusion models. In contrast to biased MD trajectories, each invocation of the steering algorithm (Sec. 2.1) targets the equilibrium distribution of the *current* bias condition and returns iid (weighted) samples from that biased ensemble. Thus, each bias update defines a well-posed thermodynamic state, and we *can* apply MBAR online to combine samples across bias conditions and obtain an unbiased free-energy estimate (and diagnostics) at any time, without waiting for a fully “filled” surface. Finally, since diffusion inference is naturally batched, MetaDiff updates the bias from a batch of samples: instead of a single Gaussian hill, we add a batchwise kernel (a mixture of Gaussians in ξ -space).

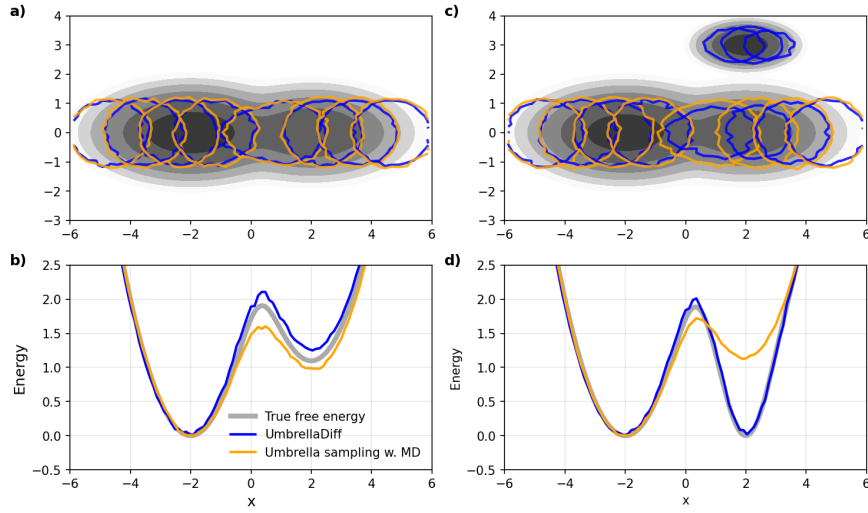


Figure 2: **UmbrellaDiff bypasses kinetic trap problem in traditional umbrella sampling.** a) Simple two-state 2D potential. Umbrella sampling is performed using 8 umbrellas that bias $\xi = x$ using regular umbrella sampling with Langevin dynamics (orange) and UmbrellaDiff (blue). Contours show 95% of the sample density of individual umbrellas. b) Free energy as a function of x using traditional umbrella sampling (orange) and UmbrellaDiff (blue). Both methods approximate the true potential of mean force (grey). c) As a) but with a third high-probability state (top right). Traditional umbrella sampling (orange) does not sample the third state as it is off-path and long simulation trajectories in each umbrella would be required to sample this state. UmbrellaDiff (blue) samples the equilibrium distribution conditioned on the umbrella potential and therefore samples the third state without issues. d) As b). Traditional umbrella sampling estimates a wrong free energy due to failing to sample the third state (orange), UmbrellaDiff estimates the correct free energy (blue).

In the limit of batch size 1, the procedure reduces to standard metadynamics applied to steered diffusion sampling.

We will again formulate the method on a one-dimensional reaction coordinate $\xi : \mathcal{X} \rightarrow \mathbb{R}$, the generalization to multiple dimensions is conceptually straightforward. We represent the bias at outer iteration k as a sum of Gaussians:

$$b_k(\xi) = \sum_{m=1}^{M_k} a_m K_\sigma(\xi; \mu_m), \quad K_\sigma(\xi; \mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\xi - \mu)^2}{2\sigma^2}\right),$$

with $b_0 \equiv 0$. For each outer iteration k , we produce a batch of samples with normalized weights using the steering algorithm (Sec. 2.1): $\{(x_{k,i}, w_{k,i})\}_{i=1}^{N_k}$ with $\xi_{k,i} = \xi(x_{k,i})$. We then deposit one Gaussian per sample sequentially so that b_k remains a sum of Gaussians at all times, in the spirit of classical metadynamics [11] and well-tempered metadynamics [90]. To control the total amount of bias added at iteration k , we choose a per-iteration mass $h_k > 0$ and split it across the N_k samples using per-Gaussian masses $\tilde{h}_{k,i}$ that satisfy $\sum_{i=1}^{N_k} \tilde{h}_{k,i} = h_k$. A simple choice is

$$\tilde{h}_{k,i} = \begin{cases} h_k/N_k, & \text{for unweighted samples,} \\ h_k \frac{w_{k,i}}{\sum_j w_{k,j}}, & \text{for importance-weighted samples,} \end{cases} \quad (25)$$

so each iteration deposits total mass h_k regardless of batch size. In practice, standard metadynamics often uses a *constant* base height ($h_k \equiv h_0$) and relies on exploration to distribute hills, whereas well-tempered metadynamics also works robustly with a constant h_k because the effective deposited height decays automatically via tempering [90]; optionally, a gently decaying schedule $h_k = h_0 \rho^k$ with $\rho \in (0, 1)$ reduces the noise during late iterations.

The final amplitude for each Gaussian i is determined by whether we use a standard or well-tempered Metadynamics scheme:

$$\text{(standard)} \quad a_{k,i} = \tilde{h}_{k,i}, \quad (26)$$

$$\text{(well-tempered, } \gamma > 1) \quad a_{k,i} = \tilde{h}_{k,i} \exp\left(-\frac{b_k^{(i-1)}(\xi_{k,i})}{\gamma - 1}\right). \quad (27)$$

Setting $\gamma = \infty$ in (27) recovers the standard rule (26). Finally, we update the bias. Starting from $b_k^{(0)} \equiv b_k$ we sequentially iterate through the batch for $i = 1, \dots, N_k$:

$$b_k^{(i)}(\xi) = b_k^{(i-1)}(\xi) + a_{k,i} K_\sigma(\xi; \xi_{k,i}), \quad \xi_{k,i} = \xi(x_{k,i}). \quad (28)$$

After processing all N_k terminals, we set $b_{k+1} \equiv b_k^{(N_k)}$ as the new bias potential.

Algorithm (MetaDiff). Given kernel width σ , tempering factor $\gamma \geq 1$, and a base-height schedule $\{h_k\}_{k \geq 0}$:

1. Initialize $b_0 \equiv 0$.
2. For $k = 0, 1, 2, \dots$:
 - (a) Sample from q_k via Steering (Sec. 2.1) and collect weighted samples $\{(x_{k,i}, w_{k,i})\}_{i=1}^{N_k}$.
 - (b) Choose $\tilde{h}_{k,i}$ using Eq. (25).
 - (c) For $i = 1, \dots, N_k$: compute $a_{k,i}$ by (26) (standard) or (27) (well-tempered) using the current $b_k^{(i-1)}$, then update $b_k^{(i)}$ via (28).
 - (d) Set $b_{k+1} \leftarrow b_k^{(N_k)}$. Stop when $\sup_\xi |b_{k+1}(\xi) - b_k(\xi)| < \varepsilon_b$ or when an MBAR-based PMF change across selected checkpoints falls below ε_F .

Mean-field consistency and stationary PMF. As $\sigma \rightarrow 0$ and $h_k \rightarrow 0$ (small hills), the expected increment at ξ under (27)–(28) satisfies

$$\mathbb{E}[b_{k+1}(\xi) - b_k(\xi)] \propto q_{\Xi,k}(\xi) e^{-b_k(\xi)/(\gamma-1)},$$

which is the well-tempered mean-field update [90]. At stationarity, $u_{\Xi}(\xi) = -(\gamma/(\gamma-1)) b_{\infty}(\xi) + C$ (with arbitrary constant C); For $\gamma = \infty$ the update reduces to standard metadynamics; in the idealized fully-filled limit the accumulated bias approximates $-u_{\Xi}(\xi) + C$ [11]. The typical practice is to keep h_k constant and choose γ to tune smoothness/convergence ($\gamma \approx 5$ –20 is common). A modest decay of h_k can further stabilize late iterations.

Illustration on a double-well potential Fig. 3 demonstrates MetaDiff on a one-dimensional double-well potential whose energy minima have a free energy difference of $-13.8 k_B T$. using $1.11 \cdot 10^5$ samples, the unbiased diffusion sampler only samples the lower energy minimum but fails to sample the second minimum (Fig. 3, first column). MetaDiff is run in batches of $1.77 \cdot 10^4$ samples. The second minimum is explored after the second batch and an accurate estimate of the free energy difference is immediately achieved, and further improved in subsequent batches (Fig. 3, columns 2-5).

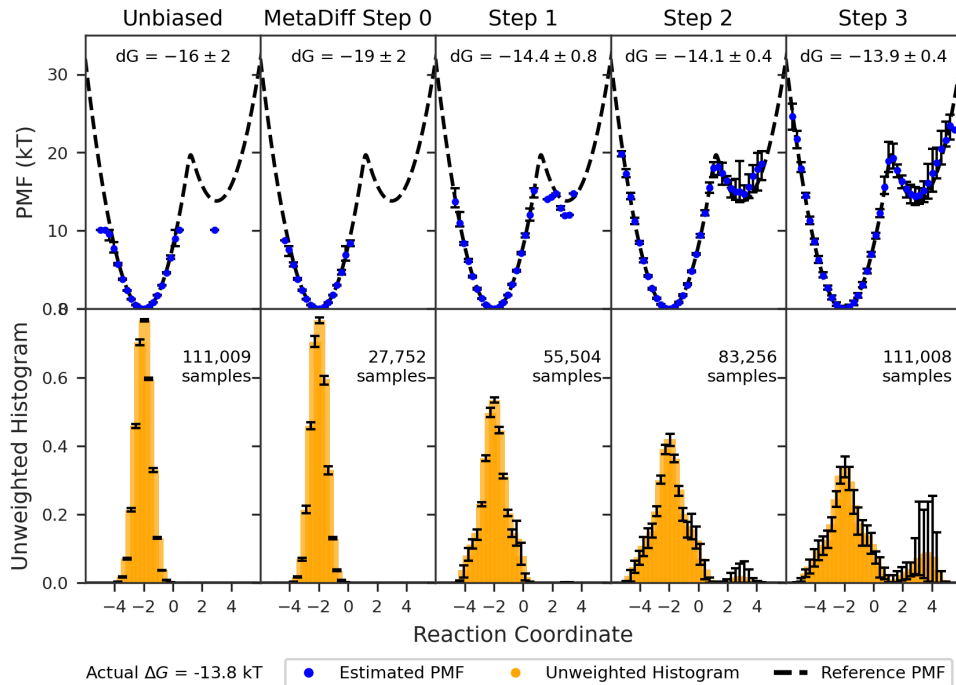


Figure 3: MetaDiff illustration on 1-dimensional double-well potential. Estimated PMF using iterations of metadynamics at a free energy difference of $\Delta G = -13.8 k_B T$.

5 ΔG -Diff: Free energy differences between two states

Here we aim to compute the free energy difference between two configurational states A and B , ΔG_{AB} . Examples include a folding free energy, binding free energy, the free energy change of a reaction or a conformational change. Let $\xi(x) \in [0, 1]$ be a smooth progress coordinate that indicates how far a configuration x has advanced from state A to state B :

$$\xi(x) = \begin{cases} 0 & x \in A \\ (0, 1) & x \notin \{A, B\} \\ 1 & x \in B \end{cases} \quad (29)$$

In practice, ξ can be a sigmoid function of a differentiable fraction-of-native-contacts (FNC), a normalized root-mean-squared-deviation (RMSD) or a committor function defined between sets of configurations A and B [93–96].

To ensure that both states can be sampled, we bias with a centered linear tilt in reduced units

$$b_a(x) = a \left(\xi(x) - \frac{1}{2} \right), \quad a \in \mathbb{R},$$

so $a > 0$ favors A ($\xi \approx 0$) and $a < 0$ favors B ($\xi \approx 1$).

Dominance and overlap diagnostics. We define core sets in the ξ coordinate as $C_A = [0, c]$, $C_B = [1 - c, 1]$ with, e.g., $c = 0.4$. From a weighted batch with DM steering (Sec. 2.1) with tilt a we have samples and normalized weights $(\xi(x_i), \bar{w}_i)$. We define the core probability masses as

$$\hat{\pi}_A(a) = \sum_i \bar{w}_i \mathbf{1}\{\xi(x_i) \in C_A\}, \quad \hat{\pi}_B(a) = \sum_i \bar{w}_i \mathbf{1}\{\xi(x_i) \in C_B\}.$$

Say a has *A-dominance* if $\hat{\pi}_A(a) \geq \vartheta$ and *B-dominance* if $\hat{\pi}_B(a) \geq \vartheta$ with default $\vartheta = 0.5$ (can be increased on demand).

For adjacent tilts $a < a'$, we measure the overlap as follows

$$\hat{\mathcal{O}}_{a \rightarrow a'} \propto \sum_{i \in a} \exp\left(- (a' - a) \left[\xi(x_i) - \frac{1}{2} \right] \right),$$

and require $\hat{\mathcal{O}}_{a \rightarrow a'}, \hat{\mathcal{O}}_{a' \rightarrow a} \geq \varepsilon_{\text{ov}}$ (e.g. 10–30%).

Algorithm (ΔG -Diff). Inputs: tilt step $s > 0$ (default $s \approx 2$ in $k_B T$), initial sample size n (default 100), dominance threshold ϑ (default 0.5), overlap threshold ε_{ov} (default 0.25), $\text{ESS}_{\text{target}}$ (default 100).

1. **Initialize:** $\mathcal{A} \leftarrow \{0\}$. Sample n terminals with $a=0$. Set $a_L \leftarrow 0$ (leftmost), $a_R \leftarrow 0$ (rightmost).
2. **Increase tilts** until both endpoints dominate:
 - (a) If $\hat{\pi}_A(a_R) < \vartheta$: set $a_R \leftarrow a_R + s$, sample n terminals at a_R , insert a_R into \mathcal{A} .
 - (b) If $\hat{\pi}_B(a_L) < \vartheta$: set $a_L \leftarrow a_L - s$, sample n terminals at a_L , insert a_L into \mathcal{A} .
3. **Ensure adjacent overlap:** Sort \mathcal{A} . While any neighboring pair (a, a') violates $\hat{\mathcal{O}}_{a \rightarrow a'} \geq \varepsilon_{\text{ov}}$ or $\hat{\mathcal{O}}_{a' \rightarrow a} \geq \varepsilon_{\text{ov}}$, insert the midpoint $\tilde{a} = (a + a')/2$, sample n terminals at \tilde{a} , insert \tilde{a} into \mathcal{A} , and update overlaps.
4. **Refine for accuracy:** Add more samples to all $a \in \mathcal{A}$ to reach $\text{ESS}_{\text{target}}$.
5. **Unbias and estimate:** Solve weighted MBAR (Sec. 2.2) over all tilts using bias potentials $b_a(x) = a(\xi(x) - \frac{1}{2})$ to obtain target ($a=0$) weights $W_n^{(0)}$, then compute

$$\widehat{\Delta g}_{\text{AB}} = -\log \frac{\sum_n W_n^{(0)} \xi(x_n)}{\sum_n W_n^{(0)} [1 - \xi(x_n)]}, \quad \widehat{\Delta G}_{\text{AB}} = k_B T \widehat{\Delta g}_{\text{AB}}.$$

Uncertainties follow from the MBAR covariance or a cluster bootstrap over denoising-trajectory IDs.

Illustration and results for protein folding with BioEmu Fig. 4a illustrates how ΔG -Diff works for a two-state potential, which is representative of typical protein folding and protein-ligand binding scenarios. When the free energy difference between the two states ΔG is large, almost all equilibrium density is in the more stable state, e.g., the folded or bound state (Fig. 4a, top left). ΔG -Diff defines a series of tilt potentials interpolating between an ensemble in which one state is dominant (top left) and an ensemble in which the other state is dominant (bottom right). If a good guess of ΔG is available, a single tilted ensemble which samples from both states is sufficient. The ΔG -Diff ensembles can be generated with steering and

combined with MBAR to the estimated free energy value. Note that in contrast to classical enhanced sampling methods that use MD, there is no need for having overlapping distributions in coordinate space. A single steered diffusion model ensemble that is able to sample both states (bottom right) is sufficient to obtain an unbiased estimate of ΔG . Multiple ensembles are useful to find the optimal tilt, and combining their samples with MBAR improves statistics.

To connect these trends to practical compute, we implemented FKC reward tilting algorithm [74] in combination with the SDE DPM-Solver++ [97], a second order denoiser, based on the open source BioEmu [3] code. We used BioEmu model version 1.1 to generate samples of a set of relatively stable proteins selected from ProThermDB [98] database, with different predicted unbiased stabilities ranging from 1 to 6 kcal/mol. At room temperature, folding free energy of 5 kcal/mol corresponds to one unfolded sample in 4160 samples, which requires on the order of several GPU hours, with the exact timing depending on GPU and protein size. We have chosen estimated optimal tilts to illustrate the correctness and enhanced sampling efficiency for those proteins. For each protein, we generated 25,000 steered samples in total with small batches, using per-system optimized steering slopes along the collective variable defined as the Root Mean Squared Deviation (RMSD) to the native PDB structure. Unbiased statistics were recovered by direct reweighting (18) since it a single bias potential was used, and folding free energies were computed from the fraction of native contacts (FNC) using a two-state model. To quantify sample efficiency, since different batches are independently sampled, we repeatedly subsampled a given number of batches (and thus a fixed number of particles) without replacement, and evaluated the resulting ΔG estimates as a function of total sample count.

We excluded 8 systems where steering did not yield usable overlap, including 5 where the RMSD range is too large (beyond 2nm), and 3 where the “unbiased” reference estimates were not reliable due to large unbiased ΔG and potentially much larger number of unbiased samples are needed for convergence. After filtering, 18 proteins remain for quantitative comparison. In Fig. 4b, we plot the enhanced sampling ΔG from reweighted steered samples at approximately 1,000 particles against the converged unbiased ΔG , showing close agreement across the full stability range. We next measured the minimum sample count required for convergence, defining convergence as a self-reference mean absolute error (MAE) below 1 kcal/mol, as shown in Fig. 4c. Ubiquitin (1UBQ) is a notable outlier, likely due to the fact that it is not a two-state folder. For unbiased sampling, this required sample count grows exponentially with ΔG , consistent with the fact that the probability of observing the thermodynamically rare state decreases exponentially in free energy difference. In contrast, steered sampling shows a much weaker scaling over the same range, demonstrating a substantial reduction in sample complexity for stable proteins.

6 Discussion

Diffusion-model-based equilibrium samplers remove one of the two fundamental bottlenecks of molecular simulation: slow mixing due to time-correlated trajectories. However, iid sampling alone does not resolve the second bottleneck—the exponential sample complexity of estimating observables dominated by low-probability regions of the equilibrium distribution, such as in protein folding or protein-ligand binding. In this work, we show that the classical bias-and-reweighting paradigm of enhanced sampling can be transferred to diffusion samplers by steering pretrained models to generate biased ensembles and recovering unbiased thermodynamics with exact reweighting (direct reweighting, WHAM/MBAR). In this sense, the present framework addresses both aspects of the sampling problem: diffusion models address slow mixing, and enhanced sampling addresses rarity, together enabling converged free-energy estimation on GPU timescales that were previously dominated by long MD trajectories.

A key consequence of replacing biased MD trajectories by independent (possibly weighted) samples is that several enhanced sampling procedures change character. In UmbrellaDiff, windows no longer need to be dynamically connected by slow transitions, and hidden barriers orthogonal to ξ do not cause kinetic trapping within a window; overlap requirements are purely statistical and can be diagnosed and repaired using standard MBAR/WHAM tools (ESS, overlap matrices, uncertainty estimates). Likewise, in MetaDiff, history-dependent biasing becomes an *equilibrium* sequence of thermodynamic states: each bias update defines a well-posed biased distribution that can be incorporated immediately into MBAR, allowing online diagnostics and early stopping without waiting for a fully “filled” surface. More broadly, once mixing is removed from the underlying sampling engine, bias variables need not correspond to slow dynamical modes, suggesting that reaction-coordinate learning and adaptive window placement may become simpler and more

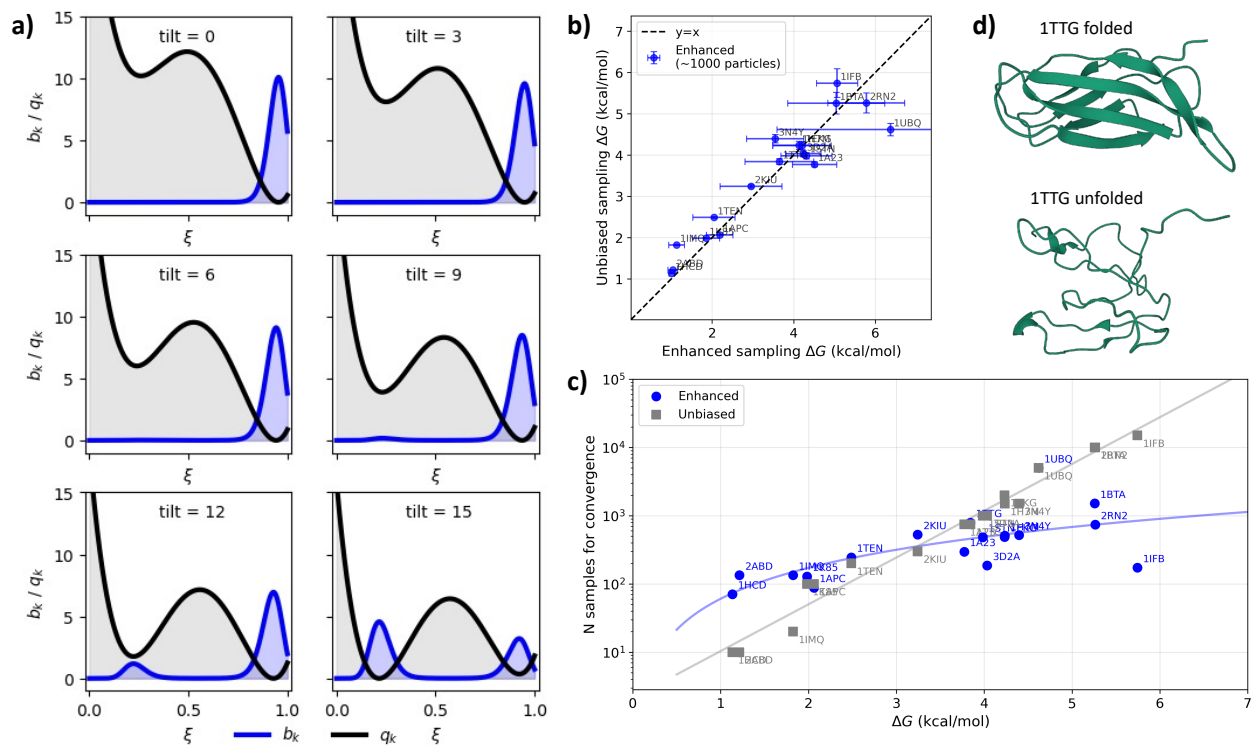


Figure 4: ΔG -Diff illustration and protein folding results with BioEmu: a) Processes with large free energy differences (top left, black), e.g., protein folding, sample almost exclusively one state in equilibrium (top left, blue). ΔG -Diff creates a series of tilted biased potentials until both states are dominantly sampled in at least one ensemble each (top left, bottom right), and these ensembles can be combined with MBAR to yield an estimate of ΔG . b) Comparison of folding free energies ΔG sampled by unsteered BioEmu and ΔG -Diff with single optimal tilt. c) Sampling efficiency as a function of ΔG : unsteered sampling (grey) requires exponentially many samples in ΔG , while the number of samples required by ΔG -Diff shows much weaker scaling with ΔG . d) Sample folded and unfolded state of a stable protein domain of fibronectin (PDB code 1ttg).

robust than in MD-based enhanced sampling.

At the same time, diffusion-based enhanced sampling inherits—and in some cases amplifies—several familiar limitations. First, the approach relies on having a sufficiently accurate pretrained equilibrium model for the molecular system and thermodynamic conditions of interest; any model mismatch propagates into reweighted estimates. Second, as in all importance-sampling schemes, weight degeneracy can arise when biases are too aggressive or when the biased ensembles (and/or the unbiased target) have insufficient overlap, necessitating careful bias design, overlap diagnostics, and potentially additional intermediate states. Finally, the present work focuses on equilibrium properties; extending these ideas to dynamical observables will require additional structure, for example through path reweighting or consistent generative dynamical models.

Looking forward, several directions appear particularly promising. Methodologically, the steering framework can be combined with adaptive schemes that learn reaction coordinates or optimal bias potentials on the fly, and with automated window/tilt construction guided by MBAR diagnostics. Architecturally, energy-based or hybrid diffusion models that allow direct evaluation of $u(x)$ could enable tighter estimator control and new forms of biasing.

Beyond equilibrium ensemble generators, an active line of research learns accelerated *dynamics* or *transport operators* from MD data, including Timewarp [99], implicit transfer operator learning [100], and trajectory generators such as MDGen [101]. Because these models are trained to reproduce finite-lag transitions rather than the equilibrium distribution, their stationary distribution need not coincide with the desired $p(x)$ unless additional constraints or corrections are imposed. One promising route is to wrap learned transport models into Metropolis–Hastings or related accept/reject schemes that target a known equilibrium distribution (as done in Timewarp), in which case bias potentials can be incorporated by changing the target and standard reweighting estimators apply. Another route is to extend biasing and reweighting to *path ensembles*, where one biases trajectory functionals (e.g., for transition-path sampling) and reweights by likelihood ratios of paths; MDGen already highlights transition-path sampling as a target application. Recent work on incorporating equilibrium priors into transfer-operator learning (e.g., BoPITO [102]) further suggests opportunities to unify learned transport with thermodynamically consistent sampling.

Finally, while we have emphasized biomolecular folding and conformational change, the formulation is general and applies to any system where iid equilibrium samplers exist but rare-event statistics are the remaining bottleneck—including materials, soft matter, and condensed-phase chemistry. We anticipate that integrating enhanced sampling principles into diffusion samplers will become a central building block for next-generation molecular simulation workflows, enabling routine computation of free energies and rare-state observables at scales previously accessible only with specialized hardware or massive distributed MD.

7 Acknowledgement

We acknowledge Carles Domingo-Enrich and Akshay Krishnamurthy for helpful discussions. GitHub Copilot has been used in part of the implementation and analysis of this work.

8 Code availability

The BioEmu inference code is in <https://github.com/microsoft/bioemu>. The enhanced sampling feature will be released soon.

References

- [1] J. Hénin, T. Lelièvre, M. R. Shirts, O. Valsson, and L. Delemotte. Enhanced sampling methods for molecular dynamics simulations [article v1.0]. *LiveCoMS*, 4:1583, 2022.
- [2] F. Noé, S. Olsson, J. Köhler, and H. Wu. Boltzmann generators - sampling equilibrium states of many-body systems with deep learning. *Science*, 365:eaaw1147, 2019.
- [3] S. Lewis, t. Hempel, J. Jiménez-Luna, M. Gastegger, Y. Xie, A. Y. K. Foong, V. García Satorras, O. Abdin, B. S. Veeling, I. Zaporozhets, Y. Chen, S. Yang, A. E. Foster, A. Schneuing, J. Nigam, F. Barbero, V. Stimper, A. Campbell, J. Yim, M. Lienen, Y. Shi, S. Zheng, H. Schulz, U. Munir, R. Sordillo, R. Tomioka, C. Clementi, and F. Noé. Scalable emulation of protein equilibrium ensembles with generative deep learning. *Science*, page eadv9817, 2025. doi: 10.1126/science.adv9817.
- [4] J. G. Kirkwood. Statistical mechanics of fluid mixtures. *J. Chem. Phys.*, 3:300–313, 1935.
- [5] R. W. Zwanzig. High-temperature equation of state by a perturbation method. i. nonpolar gases. *J. Chem. Phys.*, 22:1420–1426, 1954.
- [6] G. M. Torrie and J. P. Valleau. Nonphysical Sampling Distributions in Monte Carlo Free-Energy Estimation: Umbrella Sampling. *J. Comp. Phys.*, 23:187–199, 1977.
- [7] M. Souaille and B. Roux. Extension to the weighted histogram analysis method: combining umbrella sampling with free energy calculations. *Comput. Phys. Commun.*, 135:40–57, 2001.
- [8] E. Marinari and G. Parisi. Simulated tempering: A new monte carlo scheme. *Euro. Phys. Lett.*, 19: 451–458, 1992.
- [9] U. H. E. Hansmann. Parallel tempering algorithm for conformational studies of biological molecules. *Chem. Phys. Lett.*, 281:140 – 150, 1997.
- [10] Y. Sugita and Y. Okamoto. Replica-exchange molecular dynamics method for protein folding. *Chem. Phys. Lett.*, 314:141–151, 1999.
- [11] A. Laio and M. Parrinello. Escaping free energy minima. *Proc. Natl. Acad. Sci. USA*, 99:12562–12566, 2002.
- [12] C. H. Bennett. Efficient estimation of free energy differences from monte carlo data. *J. Chem. Phys.*, 22:245–268, 1976.
- [13] A. M. Ferrenberg and R. H. Swendsen. Optimized monte carlo data analysis. *Phys. Rev. Lett.*, 63: 1195–1198, 1989.
- [14] M. R. Shirts and J. D. Chodera. Statistically optimal analysis of samples from multiple equilibrium states. *J. Chem. Phys.*, 129:124105, 2008.
- [15] A. M. Tokita, T. Devergne, A. M. Saitta, and J. Behler. Free energy profiles for chemical reactions in solution from high-dimensional neural network potentials: The case of the strecker synthesis. *arXiv preprint arXiv:2503.05370*, 2025.
- [16] S. Stocker, H. Jung, G. Csányi, C. F. Goldsmith, K. Reuter, and J. T. Margraf. Estimating free energy barriers for heterogeneous catalytic reactions with machine learning potentials and umbrella integration. *J. Chem. Theory Comput.*, 19:6796–6804, 2025.
- [17] J. Hub and B. L. de Groot. Mechanism of selectivity in aquaporins and aquaglyceroporins. *Proc. Natl. Acad. Sci. USA*, 105:1198–1203, 2008.
- [18] F. T. Heer, D. J. Posson, W. Wojtas-Niziurski, C. M. Nimigean, and S. Bernèche. Mechanism of activation at the selectivity filter of the KcsA K⁺ channel. *eLife*, 6:e25844, 2017.

- [19] Z. Cournia, B. K. Allen, and W. Sherman. Relative binding free energy calculations in drug discovery: Recent advances and practical considerations. *J. Chem. Inf. Model.*, 57:2911–2937, 2017.
- [20] J. H. Moore, D. J. Cole, and G. Csányi. Computing solvation free energies of small molecules with first principles accuracy. *arXiv preprint arXiv:2405.18171*, 2025.
- [21] J. D. Chodera, D. L. Mobley, M. R. Shirts, R. W. Dixon, K. Branson, and V. S. Pande. Alchemical free energy methods for drug discovery: Progress and challenges. *Curr. Opin. Struct. Biol.*, 21:150–160, 2011.
- [22] J. W. Pitera and W. Swope. Understanding folding and design: Replica-exchange simulations of “trp-cage” miniproteins. *Proc. Natl. Acad. Sci. USA*, 13:7587–7592, 2003.
- [23] N. E. Charron, K. Bonneau, A. S. Pasos-Trejo, A. Guljas, Y. Chen, F. Musil, J. Venturin, D. Gusew, I. Zaporozhets, A. Krämer, C. Templeton, A. Kelkar, A. E. P. Durumeric, S. Olsson, A. Pérez, M. Majewski, B. E. Husic, A. Patel, G. De Fabritiis, F. Noé, and C. Clementi. Navigating protein landscapes with a machine-learned transferable coarse-grained model. *Nature Chemistry*, 17:1284–1292, 2025.
- [24] J. M. L. Ribeiro, P. Bravo, Y. Wang, and P. Tiwary. Reweighted autoencoded variational Bayes for enhanced sampling (RAVE). *J. Chem. Phys.*, 149:072301, 2018.
- [25] Y. Wang, J. M. L. Ribeiro, and P. Tiwary. Past-future information bottleneck for sampling molecular reaction coordinate simultaneously with thermodynamics and kinetics. *Nature Commun.*, 10:3573, 2019.
- [26] W. Chen, H. Sidky, and A. L. Ferguson. Nonlinear discovery of slow molecular modes using state-free reversible vampnets. *J. Chem. Phys.*, 150:214114, 2019.
- [27] A. Mardt, L. Pasquali, H. Wu, and F. Noé. Vampnets: Deep learning of molecular kinetics. *Nat. Commun.*, 9:5, 2018.
- [28] L. Bonati, G. Piccini, and M. Parrinello. Deep learning the slow modes for rare events sampling. *Proc. Natl. Acad. Sci. USA*, 118:e2113533118, 2021.
- [29] L. Bonati, V. Rizzi, and M. Parrinello. Data-driven collective variables for enhanced sampling. *J. Phys. Chem. Lett.*, 11:2998–3004, 2020.
- [30] D. E. Kleiman and D. Shukla. Active learning of the conformational ensemble of proteins using maximum entropy vampnets. *J. Chem. Theory Comput.*, 19:4377–4388, 2023.
- [31] J. Preto and C. Clementi. Fast recovery of free energy landscapes via diffusion-map-directed molecular dynamics. *Phys. Chem. Chem. Phys.*, 16:19181–19191, 2014.
- [32] W. Zheng, M. A. Rohrdanz, and C. Clementi. Rapid exploration of configuration space with diffusion-map-directed molecular dynamics. *J. Phys. Chem. B*, 117:12769–12776, 2013.
- [33] J.-H. Prinz, H. Wu, M. Sarich, B. G. Keller, M. Senne, M. Held, J. D. Chodera, C. Schütte, and F. Noé. Markov models of molecular kinetics: Generation and validation. *J. Chem. Phys.*, 134:174105, 2011.
- [34] F. Noé, S. Doose, I. Daidone, M. Löllmann, J. D. Chodera, M. Sauer, and J. C. Smith. Dynamical fingerprints for probing individual relaxation processes in biomolecular dynamics with simulations and kinetic experiments. *Proc. Natl. Acad. Sci. USA*, 108:4822–4827, 2011.
- [35] J. Juraszek and P. G. Bolhuis. (Un)Folding mechanisms of the FBP28 WW domain in explicit solvent revealed by multiple rare event simulation methods. *Biophys. J.*, 98:646–656, 2010.
- [36] S. Bajpai, C. R. A. Abreu, N. N. Nair, and M. E. Tuckerman. Solute tempered adiabatic free energy dynamics for enhancing conformational space sampling. *J. Chem. Theory Comput.*, 21:5928–5940, 2025.
- [37] F. Marinelli, F. Pietrucci, A. Laio, and S. Piana. A kinetic model of trp-cage folding from multiple biased molecular dynamics simulations. *PLoS Comput. Biol.*, 5:e1000452, 2009.

- [38] Y. Liu, J. Strümpfer, L. J. Freddolino, M. Gruebele, and K. Schulten. Structural characterization of λ -repressor folding from all-atom molecular dynamics simulations. *J. Phys. Chem. Lett.*, 3:1117–1123, 2012.
- [39] K. Lindorff-Larsen, S. Piana, R. O. Dror, and D. E. Shaw. How fast-folding proteins fold. *Science*, 334:517–520, 2011.
- [40] S. Piana, K. Lindorff-Larsen, and D. E. Shaw. Atomic-level description of ubiquitin folding. *Proc. Natl. Acad. Sci. USA*, 110:5915–5920, 2013.
- [41] G. R. Bowman, V. A. Voelz, and V. S. Pande. Atomistic Folding Simulations of the Five-Helix Bundle Protein Lambda 6-85. *J. Am. Chem. Soc.*, 133:664–667, 2011.
- [42] V. A. Voelz, M. Jäger, S. Yao, Y. Chen, L. Zhu, S. A. Waldauer, G. R. Bowman, M. Friedrichs, O. Bakajin, L. J. Lapidus, S. Weiss, and V. S. Pande. Slow unfolded-state structuring in acyl-CoA binding protein folding revealed by simulation and experiment. *J. Am. Chem. Soc.*, 134:12565–12577, 2012.
- [43] D. Provasi and M. Filizola. Towards ligand-specific conformations of g-protein coupled receptors: A study of the β_2 -adrenergic receptor. *PLoS Comput. Biol.*, 7:e1002193, 2011.
- [44] Y. Wang, E. Papaleo, and K. Lindorff-Larsen. Mapping transiently formed and sparsely populated conformations on a complex energy landscape. *eLife*, 5:e17505, 2016.
- [45] Y. Meng, Y.-L. Lin, and B. Roux. Computational study of the “DFG-flip” conformational transition in c-abl and c-src tyrosine kinases. *J. Phys. Chem. B*, 119:1443–1456, 2015.
- [46] O. Fleetwood, P. Matricon, J. Carlsson, and L. Delemotte. Energy landscapes reveal agonist control of g protein-coupled receptor activation via microswitches. *Biochemistry*, 59:880–891, 2020.
- [47] C. F. Abrams and E. Vanden-Eijnden. Large-scale conformational sampling of proteins using temperature-accelerated molecular dynamics. *Proc. Natl. Acad. Sci. USA*, 107:4961–4966, 2010.
- [48] K. J. Kohlhoff, D. Shukla, M. Lawrenz, G. R. Bowman, D. E. Konerding, D. Belov, R. B. Altman, and V. S. Pande. Cloud-based simulations on google exacycle reveal ligand modulation of gpcr activation pathways. *Nat. Chem.*, 6:15–21, 2014.
- [49] M. M. Sultan, R. A. Denny, R. Unwalla, F. Lovering, and V. S. Pande. Millisecond dynamics of btk reveal kinome-wide conformational plasticity within the apo kinase domain. *Sci. Rep.*, 7:15604, 2017.
- [50] G. R. Bowman and P. L. Geissler. Equilibrium fluctuations of a single folded protein reveal a multitude of potential cryptic allosteric sites. *Proc. Natl. Acad. Sci. USA*, 109:11681–11686, 2012.
- [51] D. Shukla, Y. Meng, B. Roux, and V. S. Pande. Activation pathway of src kinase reveals intermediate states as targets for drug design. *Nat. Commun.*, 5:3397, 2014.
- [52] Y. Meng, D. Shukla, V. S. Pande, and B. Roux. Transition path theory analysis of c-src kinase activation. *Proc. Natl. Acad. Sci. USA*, 113:9193–9198, 2016.
- [53] D. E. Shaw, P. Maragakis, K. Lindorff-Larsen, S. Piana, R.O. Dror, M.P. Eastwood, J.A. Bank, J.M. Jumper, J.K. Salmon, Y. Shan, and W. Wriggers. Atomic-Level Characterization of the Structural Dynamics of Proteins. *Science*, 330:341–346, 2010.
- [54] R. O. Dror, D. H. Arlow, P. Maragakis, T. J. Mildorf, A. C. Pan, H. Xu, D. W. Borhani, and D. E. Shaw. Activation mechanism of the β_2 -adrenergic receptor. *Proc. Natl. Acad. Sci. USA*, 108:18684–18689, 2011.
- [55] R. O. Dror, A. C. Pan, D. H. Arlow, D. W. Borhani, P. Maragakis, Y. Shan, H. Xu, and D. E. Shaw. Pathway and mechanism of drug binding to G-protein-coupled receptors. *Proc. Natl. Acad. Sci. USA*, 108:13118–13123, 2011.

- [56] A. J Havens, B. K. Miller, B. Yan, C. Domingo-Enrich, A. Sriram, D. S. Levine, B. M Wood, B. Hu, B. Amos, B. Karrer, X. Fu, G.-H. Liu, and R. T. Q. Chen. Adjoint sampling: Highly scalable diffusion samplers via adjoint matching. In A. Singh, M. Fazel, D. Hsu, S. Lacoste-Julien, F. Berkenkamp, T. Maharaj, K. Wagstaff, and J. Zhu, editors, *ICML 42*, volume 267 of *Proc. Mach. Learn. Res.*, pages 22204–22237. PMLR, 2025.
- [57] B. Jing, B. Berger, and T. Jaakkola. AlphaFold meets flow matching for generating protein ensembles. In R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp, editors, *ICML 41*, volume 235 of *Proc. Mach. Learn. Res.*, pages 22277–22303. PMLR, 2024.
- [58] G. Janson, A. Jussupow, and M. Feig. Deep generative modeling of temperature-dependent structural ensembles of proteins. *Commun. Chem.*, 8(354), 2025. doi: 10.1038/s42004-025-01737-2.
- [59] S. Zheng, J. He, C. Liu, Y. Shi, Z. Lu, W. Feng, F. Ju, J. Wang, J. Zhu, Y. Min, H. Zhang, S. Tang, H. Hao, P. Jin, C. Chen, F. Noé, H. Liu, and T.-Y. Liu. Predicting equilibrium distributions for molecular systems with deep learning. *Nat. Mach. Intell.*, 6:558–567, 2024.
- [60] Yikai Liu, Zongxin Yu, Richard J. Lindsay, Guang Lin, Ming Chen, Abhilash Sahoo, and Sonya M. Hanson. Exendiff: An experiment-guided diffusion model for protein conformational ensemble generation. *PRX Life*, 3(2):023013, 2025. doi: 10.1103/PRXLife.3.023013. URL <https://link.aps.org/doi/10.1103/PRXLife.3.023013>.
- [61] J. Nam, B. Máté, A. P. Toshev, M. Kaniselman, R. Gómez-Bombarelli, R. T. Q. Chen, B. Wood, G.-H. Liu, and B. Kurt Miller. Enhancing diffusion-based sampling with molecular collective variables. *arXiv:2510.11923*, 2025.
- [62] D. D. Richman, J. Karagunesian, C. M. Suomivuori, and R. O. Dror. Unlocking hidden biomolecular conformational landscapes in diffusion models at inference time. *arXiv:2512.03312*, 2025.
- [63] H. Y. I. Lam, S. P. Ojeda, M. Brezinova, J. Hanke, X. E. Ong, Y. Mu, and M. Vendruscolo. Metadiffusion: inference-time meta-energy biasing of biomolecular diffusion models. <https://doi.org/10.64898/2026.02.10.704873>, 2026.
- [64] M. Arts, V. Garcia Satorras, C.-W. Huang, D. Zügner, M. Federici, C. Clementi, F. Noé, R. Pinsler, and R. van den Berg. Two for one: Diffusion models and force fields for coarse-grained molecular dynamics. *J. Chem. Theory Comput.*, 19:6151–6159, 2023.
- [65] M. Plainier, H. Wu, L. Klein, S. Günnemann, and F. Noé. Consistent sampling and simulation: Molecular dynamics with energy-based diffusion models. In *Adv. Neural Inf. Proc. Syst.*, 2025. arXiv:2506.17139.
- [66] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. In *Adv. Neural Inf. Proc. Syst.*, volume 34, pages 8780–8794, 2021.
- [67] J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [68] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In H. Daumé III and A. Singh, editors, *ICML 37*, volume 119 of *Proc. Mach. Learn. Res.*, pages 8960–8970. PMLR, 2020.
- [69] R. M. Neal. Annealed importance sampling. *Statistics and Computing*, 11:125–139, 2001.
- [70] P. Del Moral. *Feynman-Kac Formulae: Genealogical and Interacting Particle Systems with Applications*. Springer, 2004.
- [71] A. Doucet and A. M. Johansen. A tutorial on particle filtering and smoothing: Fifteen years later. In *The Oxford Handbook of Nonlinear Filtering*. 2009.
- [72] R. Singhal, Z. Horvitz, R. Teehan, M. Ren, Z. Yu, K. McKeown, and R. Ranganath. A general framework for inference-time scaling and steering of diffusion models. *arXiv preprint arXiv:2501.06848*, 2025.

- [73] G. O. Roberts and R. L. Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, 2:341–363, 1996.
- [74] M. Skreta, T. Akhond-Sadeh, V. Ohanesian, R. Bondesan, A. Aspuru-Guzik, A. Doucet, R. Brekelmans, A. Tong, and K. Neklyudov. Feynman-kac correctors in diffusion: Annealing, guidance, and product of experts. *arXiv preprint arXiv:2503.02819*, 2025.
- [75] Y. Ren, W. Gao, L. Ying, G. M Rotskoff, and J. Han. Driftlite: Lightweight drift control for inference-time scaling of diffusion models. *arXiv preprint arXiv:2509.21655*, 2025.
- [76] L. Kish. *Survey sampling*. John Wiley & Sons, 1995.
- [77] A. Kong, P. McCullagh, X.-L. Meng, D. Nicolae, and Z. Tan. A theory of statistical models for monte carlo integration. *J. R. Statist. Soc. B*, 65:585–618, 2003.
- [78] C. Bartels. Analyzing biased monte carlo and molecular dynamics simulations. *Chem. Phys Lett.*, 331: 446–454, 2000.
- [79] S. Kumar, D. Bouzida, R. H. Swendsen, P. A. Kollman, and J. M. Rosenberg. The weighted histogram analysis method for free-energy calculations on biomolecules .1. the method. *J. Comput. Chem.*, 13: 1011 – 1021, 1992.
- [80] J. Kästner. Umbrella sampling. *WIRE Comput. Mol. Sci.*, 1:932–942, 2011.
- [81] P. V. Klimovich, M. R. Shirts, and D. L. Mobley. Guidelines for the analysis of free energy calculations. *J. Comput. Aided Mol. Des.*, 29:397–411, 2015.
- [82] J. D. Chodera, P. J. Elms, W. C. Swope, J. H. Prinz, S. Marqusee, C. Bustamante, F. Noé, and V. S. Pande. A robust approach to estimating rates from time-correlation functions. <http://arxiv.org/abs/1108.2304>, 2011.
- [83] M. Shirts, K. Beauchamp, L. Naden, J. Chodera, J. Rodríguez-Guerra, S. Martiniani, C. Stern, M. Henry, J. Fass, R. Gowers, R. T. McGibbon, B. Dice, C. Jones, D. Dotson, and T. Burgin. choderalab/pymbar: 4.0.0, 2022. URL <https://doi.org/10.5281/zenodo.6872022>.
- [84] E. H. Thiede, B. Van Koten, J. Weare, and A. R. Dinner. Eigenvector method for umbrella sampling enables error analysis. *J. Chem. Phys.*, 145:084115, 2016.
- [85] M. Das and R. Venkatramani. AutoSIM: Redesigning and automating umbrella sampling for biomolecular conformational transitions. *J. Chem. Phys.*, 163:014111, 2025.
- [86] M. Yang, L. Yang, Y. Gao, and H. Hu. Combine umbrella sampling with integrated tempering method for efficient and accurate calculation of free energy changes of complex energy surface. *J. Chem. Phys.*, 141:044108, 2014.
- [87] C. F. Sousa, R. A. Becker, C.-M. Lehr, O. V. Kalinina, and J. S. Hub. Simulated tempering-enhanced umbrella sampling improves convergence of free energy calculations of drug membrane permeation. *J. Chem. Theory Comput.*, 19:1898–1907, 2023.
- [88] A. Laio and F. L. Gervasio. Metadynamics: a method to simulate rare events and reconstruct the free energy in biophysics, chemistry and material science. *Rep. Prog. Phys.*, 71:126601, 2008.
- [89] A. Barducci, M. Bonomi, and M. Parrinello. Metadynamics. *WIREs Comput. Mol. Sci.*, 1:826–843, 2011. doi: 10.1002/wcms.31.
- [90] A. Barducci, G. Bussi, and M. Parrinello. Well-tempered metadynamics: A smoothly converging and tunable free-energy method. *Phys. Rev. Lett.*, 100:020603, 2008.
- [91] J. F. Dama, M. Parrinello, and G. A. Voth. Well-tempered metadynamics converges asymptotically. *Phys. Rev. Lett.*, 112:240602, 2014.

- [92] P. Tiwary and M. Parrinello. A time-independent free energy estimator for metadynamics. *J. Phys. Chem. B*, 119:736–742, 2015. doi: 10.1021/jp504920s.
- [93] L. Onsager. Initial recombination of ions. *Phys. Rev.*, 54:554+, 1938.
- [94] P. G. Bolhuis, C. Dellago, and D. Chandler. Reaction coordinates of biomolecular isomerization. *Proc. Natl. Acad. Sci. USA*, 97:5877–5882, 2000.
- [95] W. E and E. Vanden-Eijnden. Towards a Theory of Transition Paths. *J. Stat. Phys.*, 123:503–523, 2006.
- [96] B. Peters and B. L. Trout. Obtaining reaction coordinates by likelihood maximization. *J. Chem. Phys.*, 125:054108, 2006.
- [97] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *Machine Intelligence Research*, 22(4): 730–751, 2025.
- [98] Rahul Nikam, A Kulandaisamy, K Harini, Divya Sharma, and M Michael Gromiha. ProThermDB: Thermodynamic database for proteins and mutants revisited after 15 years. *Nucleic Acids Res.*, 49: D420–D424, 2021.
- [99] L. Klein, A. Y. K. Foong, T. E. Fjelde, B. Mlodozieniec, M. Brockschmidt, S. Nowozin, F. Noé, and R. Tomioka. Timewarp: Transferable acceleration of molecular dynamics by learning time-coarsened dynamics. *arXiv preprint arXiv:2302.01170*, 2023.
- [100] M. Schreiner, O. Winther, and S. Olsson. Implicit transfer operator learning: Multiple time-resolution surrogates for molecular dynamics. *arXiv preprint arXiv:2305.18046*, 2023.
- [101] B. Jing, H. Stärk, T. Jaakkola, and B. Berger. Generative modeling of molecular dynamics trajectories. *arXiv preprint arXiv:2409.17808*, 2024.
- [102] J. Viguera Diez, M. Schreiner, O. Engkvist, and S. Olsson. Boltzmann priors for implicit transfer operators. *arXiv preprint arXiv:2410.10605*, 2025.

A BioEmu Folding Free Energies Details

We selected 26 proteins from the ProThermDB database, ranging from 76 to 372 residues. The steering potential is a clamped linear function of the collective variable, $V = s \cdot \text{clamp}(CV, c_{\min}, c_{\max})$, where s is a per-system estimated optimal slope, and $c_{\min/\max}$ are the clipping bounds to prevent the potential going unbounded. Within each run, the steered samples are generated in small independent batches with batch size dependent of the protein size, allowing us to pool and subsample batches across runs to study convergence at varying sample sizes. Specifically, we randomly draw (without replacement) subsets of batches from the combined pool and compute ΔG for each draw, repeating 200 times per sample size to obtain statistics. We note the subsampling is performed on the batch level because different batches are independent, but not within each batch because samples in the same batch have already been resampled through the steering algorithm. Folding free energies are estimated from the fraction of native contacts (FNC) using a two-state model with inverse Boltzmann reweighting to correct for the steering bias.

We exclude 8 of the 26 systems from the main analysis: 5 proteins (1QLP, 1CEC, 2LZM, 1RHD, 1RIL) whose RMSD ranges go beyond 2nm, leading to difficult steering and extreme reweighting factors, and 3 proteins (1HK0, 1IOB, 1Y9O) for which the unbiased reference ΔG is unreliable due to unconverged sampling of the unfolded state. The remaining 18 systems (1A23, 1APC, 1BTA, 1EKG, 1H7M, 1HCD, 1IFB, 1IMQ, 1K85, 1STN, 1TEN, 1TTG, 1UBQ, 2ABD, 2KIU, 2RN2, 3D2A, 3N4Y) span predicted unbiased ΔG values from 1.1 to 5.7 kcal/mol.

Fig. 5 shows the catastrophic failure rate—defined as the probability that a random draw of n particles contains no unfolded sample ($\text{FNC} < 0.5$), as a function of sample size for each system. For steered sampling, this rate drops to zero at modest sample sizes across nearly all systems, whereas unbiased sampling requires substantially more particles, particularly for thermodynamically stable proteins. Fig. 6 shows the corresponding mean absolute error (MAE) of ΔG relative to each method’s own converged estimate, as a function of sample size. Steered sampling achieves sub-kcal/mol accuracy at $n \approx 100$ –1,000 depending on stability. While unbiased sampling converges fast on systems with relatively small ΔG , the required number of samples grow exponentially with ΔG . The steered sampling supersedes the unsteered sampling at $\Delta G \approx 3.2$ kcal/mol, and the advantage becomes more significant with larger ΔG .

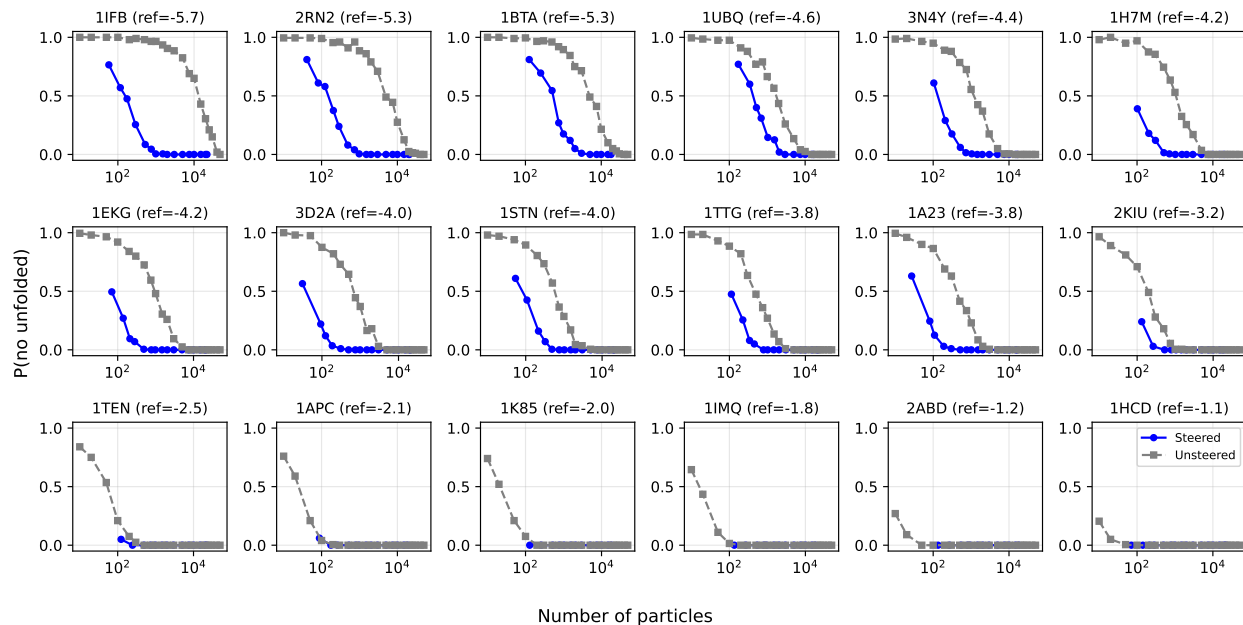


Figure 5: Per-system catastrophic failure rate as a function of sample size for 18 ProThermDB proteins, sorted by reference ΔG . Catastrophic failure is defined as the probability that a random draw of n particles contains no unfolded sample ($FNC < 0.5$). Blue: steered sampling; grey: unbiased sampling. Steered sampling eliminates catastrophic failures at small sample sizes across all systems, while unbiased sampling requires orders of magnitude more particles for very stable proteins.

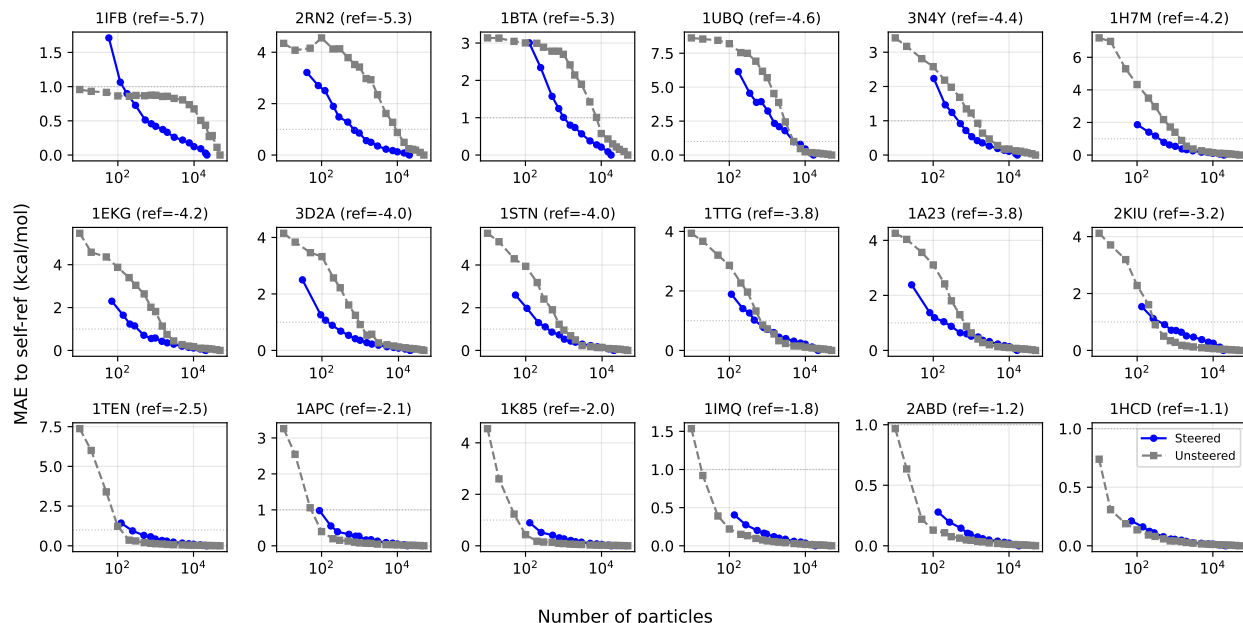


Figure 6: Per-system mean absolute error (MAE) of ΔG as a function of sample size, with each method referenced to its own converged pooled estimate. Proteins are sorted by reference ΔG (most stable at top-left). Steered sampling (blue) converges substantially faster than unbiased sampling (grey), with the advantage most pronounced for more stable systems.