

SUPPLEMENTAL MATERIAL: DYNAMICAL COMPUTATION OF THE DENSITY OF STATES AND BAYES FACTORS USING NONEQUILIBRIUM IMPORTANCE SAMPLING

GRANT M. ROTSKOFF AND ERIC VANDEN-EIJNDEN

1. **Derivation of Eqs. (6) and (7).** To derive Eq. (6) for the density $\rho_{\text{ne}}(\mathbf{x})$, consider first the forward time trajectories alone and let us define the forward density $\rho_{\text{ne}}^+(\mathbf{x})$ via

$$\begin{aligned}\langle \phi \rangle_{\text{ne}}^+ &= \frac{1}{\langle \tau^+ \rangle} \int_{\Omega} \int_0^{\tau^+(\mathbf{x})} \phi(\mathbf{X}(t, \mathbf{x})) dt \rho(\mathbf{x}) d\mathbf{x} \\ &\equiv \int_{\Omega} \phi(\mathbf{x}) \rho_{\text{ne}}^+(\mathbf{x}) d\mathbf{x}\end{aligned}\tag{SM.1}$$

The nonequilibrium density $\rho_{\text{ne}}^+(\mathbf{x})$ satisfies the stationary Liouville equation

$$\langle \tau^+ \rangle^{-1} \rho(\mathbf{x}) = \nabla \cdot (\mathbf{b}(\mathbf{x}) \rho_{\text{ne}}^+(\mathbf{x})), \tag{SM.2}$$

with the boundary condition $\rho_{\text{ne}}^+(\mathbf{x}) = 0$ on the regions of $\partial\Omega$ where $\mathbf{b}(\mathbf{x})$ points inward (since, by construction, no mass can be transported there forward in time). Physically, this equation asserts that, at stationarity, the nonequilibrium probability flux out of a small volume in the vicinity of \mathbf{x} is balanced by the rate of reinjection. To solve (SM.2) notice that if we differentiate $\rho_{\text{ne}}^+(\mathbf{X}(t, \mathbf{x}))$ with respect to time, by the chain rule we have

$$\frac{d}{dt} \rho_{\text{ne}}^+ = \mathbf{b} \cdot \nabla \rho_{\text{ne}}^+ = \langle \tau^+ \rangle^{-1} \rho - (\nabla \cdot \mathbf{b}) \rho_{\text{ne}}^+ \tag{SM.3}$$

where all functions are evaluated at $\mathbf{X}(t, \mathbf{x})$ and we used Eq. (1) for $\mathbf{X}(t, \mathbf{x})$ to derive the first equality and (SM.2) to derive the second. Using the Jacobian (Eq. (5)),

$$J(t, \mathbf{x}) = \exp \left(\int_0^t \nabla \cdot \mathbf{b}(\mathbf{X}(s, \mathbf{x})) ds \right), \tag{SM.4}$$

we can write (SM.3) as

$$\frac{d}{dt} (\rho_{\text{ne}}^+(\mathbf{X}(t, \mathbf{x})) J(t, \mathbf{x})) = \langle \tau^+ \rangle^{-1} \rho(\mathbf{X}(t, \mathbf{x})) J(t, \mathbf{x}). \tag{SM.5}$$

Integration from $t = \tau^-(\mathbf{x})$ to $t = 0$ using the boundary condition $\rho_{\text{ne}}^+(\mathbf{X}(\tau^-(\mathbf{x}), \mathbf{x})) = 0$ gives

$$\rho_{\text{ne}}^+(\mathbf{x}) = \langle \tau^+ \rangle^{-1} \int_{\tau^-(\mathbf{x})}^0 J(t, \mathbf{x}) \rho(\mathbf{X}(t, \mathbf{x})) dt. \tag{SM.6}$$

A similar calculation gives the nonequilibrium stationary density resulting from the reverse time propagation of the dynamics, $\rho_{\text{ne}}^-(\mathbf{x})$, defined via

$$\begin{aligned}\langle\phi\rangle_{\text{ne}}^- &= \frac{1}{\langle\tau^-\rangle} \int_{\Omega} \int_{\tau^-(\mathbf{x})}^0 \phi(\mathbf{X}(t, \mathbf{x})) dt \rho(\mathbf{x}) d\mathbf{x} \\ &\equiv \int_{\Omega} \phi(\mathbf{x}) \rho_{\text{ne}}^-(\mathbf{x}) d\mathbf{x}\end{aligned}\tag{SM.7}$$

We obtain

$$\rho_{\text{ne}}^-(\mathbf{x}) = \langle\tau^-\rangle^{-1} \int_0^{\tau^+(\mathbf{x})} J(t, \mathbf{x}) \rho(\mathbf{X}(t, \mathbf{x})) dt.\tag{SM.8}$$

The nonequilibrium density defined via Eq. (4) can then be expressed by superposing $\rho_{\text{ne}}^+(\mathbf{x})$ and $\rho_{\text{ne}}^-(\mathbf{x})$:

$$\begin{aligned}\rho_{\text{ne}}(\mathbf{x}) &= \frac{\langle\tau^+\rangle \rho_{\text{ne}}^+(\mathbf{x}) - \langle\tau^-\rangle \rho_{\text{ne}}^-(\mathbf{x})}{\langle\tau\rangle} \\ &= \langle\tau\rangle^{-1} \int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} J(t, \mathbf{x}) \rho(\mathbf{X}(t, \mathbf{x})) dt\end{aligned}\tag{SM.9}$$

This is Eq. (6).

To derive Eq. (7), start from Eq. (2) and write

$$\begin{aligned}\langle\phi\rangle &= \langle\phi\rho/\rho_{\text{ne}}\rangle_{\text{ne}} \\ &= \frac{1}{\langle\tau\rangle} \int_{\mathbb{R}^d} \int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} \frac{\phi(\mathbf{X}(t, \mathbf{x})) \rho(\mathbf{X}(t, \mathbf{x}))}{\rho_{\text{ne}}(\mathbf{X}(t, \mathbf{x}))} dt \rho(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbb{R}^d} \int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} \frac{\phi(\mathbf{X}(t, \mathbf{x})) \rho(\mathbf{X}(t, \mathbf{x}))}{\int_{\tau^-(\mathbf{X}(t, \mathbf{x}))}^{\tau^+(\mathbf{X}(t, \mathbf{x}))} J(s, \mathbf{X}(t, \mathbf{x})) \rho(\mathbf{X}(s, \mathbf{X}(t, \mathbf{x}))) ds} dt \rho(\mathbf{x}) d\mathbf{x}.\end{aligned}\tag{SM.10}$$

Using $\mathbf{X}(s, \mathbf{X}(t, \mathbf{x})) = \mathbf{X}(s+t, \mathbf{x})$ we have

$$\begin{aligned}J(s, \mathbf{X}(t, \mathbf{x})) &= \exp\left(\int_0^s \nabla \cdot \mathbf{b}(\mathbf{X}(u, \mathbf{X}(t, \mathbf{x}))) du\right) \\ &= \exp\left(\int_0^s \nabla \cdot \mathbf{b}(\mathbf{X}(u+t, \mathbf{x})) du\right) \\ &= \exp\left(\int_t^{s+t} \nabla \cdot \mathbf{b}(\mathbf{X}(u', \mathbf{x})) du'\right) \\ &= \frac{J(s+t, \mathbf{x})}{J(t, \mathbf{x})}\end{aligned}\tag{SM.11}$$

which implies that we can transform the denominator in (SM.10) into

$$\begin{aligned}&\int_{\tau^-(\mathbf{X}(t, \mathbf{x}))}^{\tau^+(\mathbf{X}(t, \mathbf{x}))} J(s, \mathbf{X}(t, \mathbf{x})) \rho(\mathbf{X}(s, \mathbf{X}(t, \mathbf{x}))) ds \\ &= \frac{1}{J(t, \mathbf{x})} \int_{\tau^-(\mathbf{X}(t, \mathbf{x}))}^{\tau^+(\mathbf{X}(t, \mathbf{x}))} J(s+t, \mathbf{x}) \rho(\mathbf{X}(s+t, \mathbf{x}))) ds \\ &= \frac{1}{J(t, \mathbf{x})} \int_{\tau^-(\mathbf{X}(t, \mathbf{x}))+t}^{\tau^+(\mathbf{X}(t, \mathbf{x}))+t} J(s', \mathbf{x}) \rho(\mathbf{X}(s', \mathbf{x}))) ds' \\ &= \frac{1}{J(t, \mathbf{x})} \int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} J(s', \mathbf{x}) \rho(\mathbf{X}(s', \mathbf{x}))) ds'\end{aligned}\tag{SM.12}$$

Inserting this expression in (SM.10) gives

$$\begin{aligned}\langle \phi \rangle &= \int_{\mathbb{R}^d} \int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} \frac{\phi(\mathbf{X}(t, \mathbf{x})) J(t, \mathbf{x}) \rho(\mathbf{X}(t, \mathbf{x}))}{\int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} J(s', \mathbf{x}) \rho(\mathbf{X}(s', \mathbf{x})) ds'} dt \rho(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbb{R}^d} \frac{\int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} \phi(\mathbf{X}(t, \mathbf{x})) J(t, \mathbf{x}) \rho(\mathbf{X}(t, \mathbf{x})) dt}{\int_{\tau^-(\mathbf{x})}^{\tau^+(\mathbf{x})} J(t, \mathbf{x}) \rho(\mathbf{X}(t, \mathbf{x})) dt} \rho(\mathbf{x}) d\mathbf{x}\end{aligned}\tag{SM.13}$$

which is Eq. (7).

2. Comparison with Neal's Annealed Importance Sampling (AIS) method. While our method is conceptually similar to AIS, the estimator (13) has important philosophical and practical differences from AIS, which we highlight here.

In AIS, one defines a sequence of distributions in order to estimate an expectation with respect to a target density ρ_1 , usually known only up to a normalization factor, i.e. we know some $\hat{\rho}_1(\mathbf{x})$ can be evaluated pointwise, but the factor $Z_1 = \int_{\Omega} \hat{\rho}_1(\mathbf{x}) d\mathbf{x}$ needed to get the normalized density $\rho_1(\mathbf{x}) = Z_1^{-1} \hat{\rho}_1(\mathbf{x})$ is not known. In practice, the sampling is done by transporting from an initial density ρ using a nonequilibrium dynamics: ρ can be sampled e.g. using Metropolis-Hastings Monte-Carlo, and is also only known up to a normalization factor in general. This transport defines a nonequilibrium density $\rho_{\text{ne}}(\mathbf{x})$, also known only up to a normalization factor, so that the annealed importance sampling scheme for the expectation of an observable ϕ is

$$\langle \phi \rangle_{\rho_1} = \frac{\langle \phi \hat{\rho}_1 / \hat{\rho}_{\text{ne}} \rangle_{\text{ne}}}{\langle \hat{\rho}_1 / \hat{\rho}_{\text{ne}} \rangle_{\text{ne}}},\tag{SM.14}$$

where the ratio $\hat{\rho}_1 / \hat{\rho}_{\text{ne}}$ plays the role of the weights in AIS.

The situation with our estimator is quite different. First, we use $\rho_1 = \rho$. Secondly, we generate the trajectories in a way such that, even if we do not know the normalization of $\rho_1 = \rho$, we have the property that

$$\langle \hat{\rho}_1 / \hat{\rho}_{\text{ne}} \rangle_{\text{ne}} = \langle \hat{\rho} / \hat{\rho}_{\text{ne}} \rangle_{\text{ne}} = 1\tag{SM.15}$$

by construction. Therefore our reweighting scheme does not require that we estimate a ratio of expectations; therefore, it provides us with an unbiased estimator, unlike (SM.14). Finally, this estimator has lower variance than the direct sample mean estimator, as shown by Eq. (9) in the main text.

3. The mean-field Ising model. We next consider a continuous version of the Curie-Weiss magnet, i.e. the mean-field Ising model with d spins and potential (this is Eq. (15))

$$U(\mathbf{q}) = -\frac{1}{2d} \sum_{i,j=1}^d \cos q_i \cos q_j = -\frac{1}{2d} \left(\sum_{i=1}^d \cos q_i \right)^2\tag{SM.16}$$

The Gibbs (canonical) density for this model is

$$\rho_c(\mathbf{q}) = Z^{-1}(\beta) e^{-\beta U(\mathbf{q})} \quad \text{where} \quad Z(\beta) = \int_{[-\pi, \pi]^d} e^{-\beta U(\mathbf{q})} d\mathbf{q}.\tag{SM.17}$$

This system has similar thermodynamics properties as the standard Curie-Weiss magnet with discrete spins, but it is amenable to simulation by Langevin dynamics since the angles q_i vary continuously.

In particular, like the standard Curie-Weiss magnet, the system with potential (SM.16) displays phase-transitions when β is varied. To see why, and also introduce the scaled free energy that we monitor in our numerical experiments, let us marginalize the Gibbs density (SM.17) on the average magnetization m defined as

$$m = \frac{1}{d} \sum_{i=1}^d \cos q_i. \quad (\text{SM.18})$$

This marginalized density is given by

$$\bar{\rho}_c(m) = \int_{[-\pi, \pi]^K} \rho_c(\mathbf{q}) \delta \left(m - \frac{1}{d} \sum_{i=1}^d \cos q_i \right) d\mathbf{q}. \quad (\text{SM.19})$$

A simple calculation shows that

$$\bar{\rho}_c(m) = \bar{Z}^{-1}(\beta) e^{-\beta d F_d(m, \beta)} \quad \text{where} \quad \bar{Z}(\beta) = \int_{-1}^1 e^{-\beta d F_d(m, \beta)} dm \quad (\text{SM.20})$$

Here we introduced the (scaled) free energy $F_d(m, \beta)$ defined as

$$F_d(m, \beta) = V(m) - \beta^{-1} S_d(m) \quad (\text{SM.21})$$

with potential term

$$V(m) = -\frac{1}{2} m^2 \quad (\text{SM.22})$$

and entropic term

$$S_d(m) = d^{-1} \log \int_{[-\pi, \pi]^d} \delta \left(m - \frac{1}{d} \sum_{i=1}^d \cos q_i \right) d\mathbf{q}. \quad (\text{SM.23})$$

The marginalized density (SM.20) and the free energy (SM.21) can be used to analyze the properties of the system in thermodynamic limit when $d \rightarrow \infty$ and map out its phase transition diagram in this limit. In particular, standard results from large deviation theory recalled below can be used to show that $F_d(m, \beta)$ has a limit as $d \rightarrow \infty$ that has a single minimum at high temperature, but two minima at low temperature. Since $F_d(m, \beta)$ is scaled by d in (SM.20), this implies that density can become bimodal at low temperature, indicative of the presence of two strongly metastable states separated by a free energy barrier whose height is proportional to d .

The limiting free energy $F(m, \beta)$ is defined as

$$F(m, \beta) = \lim_{d \rightarrow \infty} F_d(m) = -\frac{1}{2} m^2 - \beta^{-1} \lim_{d \rightarrow \infty} S_d(m). \quad (\text{SM.24})$$

To calculate the limit of the entropic term, let us define $H(\lambda)$ via the Laplace transform of (SM.23) through

$$\begin{aligned} e^{-dH(\lambda)} &= \int_{-1}^1 e^{-d\lambda m + dS_d(m)} dm \\ &= \int_{[-\pi, \pi]^d} e^{-\lambda \sum_{i=1}^d \cos(q_i)} d\mathbf{q} \\ &= \prod_{i=1}^d \int_{-\pi}^{\pi} e^{-\lambda \cos(q_i)} dq_i \\ &= (2\pi I_0(\lambda))^d, \end{aligned} \quad (\text{SM.25})$$

where $I_0(\lambda)$ is a modified Bessel function. In the large d limit, $S(m)$ can be calculated from $H(\lambda)$ by Legendre transform

$$\begin{aligned} S(m) &= \lim_{d \rightarrow \infty} S_d(m) = \min_{\lambda} \{ \lambda m - H(\lambda) \} \\ &= \min_{\lambda} \{ \lambda m + \log I_0(\lambda) \} + \log(2\pi). \end{aligned} \quad (\text{SM.26})$$

The minimizer $\lambda(m)$ of (SM.26) satisfies

$$m = -\frac{I_1(\lambda(m))}{I_0(\lambda(m))} \quad (\text{SM.27})$$

which, upon inversion, offers a way to parametrically represent $S(m)$ using

$$S(m(\lambda)) = \lambda m(\lambda) + \log I_0(\lambda) + \log(2\pi), \quad m(\lambda) = -\frac{I_1(\lambda)}{I_0(\lambda)}, \quad \lambda \in \mathbb{R}. \quad (\text{SM.28})$$

Similarly we can represent $F(m, \beta)$ as

$$F(m(\lambda), \beta) = -\frac{1}{2}m^2(\lambda) - \beta m(\lambda) - \beta^{-1}S(m(\lambda)), \quad m(\lambda) = -\frac{I_1(\lambda)}{I_0(\lambda)} \quad \lambda \in \mathbb{R}. \quad (\text{SM.29})$$

These are the formulae we use to plot the free energy shown in the inset of Fig. 1. We compare this free energy with the one obtained by estimating the following expectation using our estimator with a single ascent / descent trajectory

$$-(d\beta)^{-1} \log \left\langle e^{-\beta U(\mathbf{q})} \delta \left(m - \frac{1}{d} \sum_{i=1}^d \cos q_j \right) \right\rangle \quad (\text{SM.30})$$

A similar calculation can be used to estimate the factor $Z(\beta)$ in (SM.17) using the identity

$$Z(\beta) = \int_{-1}^1 e^{\frac{1}{2}\beta dm^2 + dS_d(m)} dm \quad (\text{SM.31})$$

In the $d \rightarrow \infty$ limit this implies that

$$G(\beta) \equiv \lim_{d \rightarrow \infty} d^{-1} \log Z(\beta) = \max_{m \in [-1, 1]} \left(\frac{1}{2}\beta dm^2 + dS(m) \right) \quad (\text{SM.32})$$

Since the density of states

$$D(E) = \int_{[-\pi, \pi]^d \times \mathbb{R}^d} \delta \left(E - \frac{1}{2}|\mathbf{p}|^2 - U(\mathbf{q}) \right) d\mathbf{q} d\mathbf{p} \quad (\text{SM.33})$$

is related to $Z(\beta)$ as (accounting for the extra factor coming from the momenta)

$$(2\pi\beta^{-1})^{d/2} Z(\beta) = \int_{\mathbb{R}} D(E) e^{-\beta E} dE \quad (\text{SM.34})$$

we have

$$\lim_{d \rightarrow \infty} d^{-1} \log D(d\mathcal{E}) = Q(\mathcal{E}) \quad (\text{SM.35})$$

with

$$Q(\mathcal{E}) = -\min_{\beta} \left(\beta \mathcal{E} + G(\beta) - \frac{1}{2} \log \beta \right) + \frac{1}{2} \log \pi. \quad (\text{SM.36})$$

This is the formula we used to plot the red curve in Fig. 1 used for comparison with the estimate of $V(E)$ we obtained directly using our estimator.