

Supporting Information – A protocol for identifying accurate collective variables in enhanced molecular dynamics simulations for the description of structural transformations in flexible metal-organic frameworks

Ruben Demuyne¹, Jelle Wieme¹, Sven M.J. Rogge¹, Karen D. Dedecker¹, Louis Vanduyfhuys¹, Michel Waroquier¹ and Veronique Van Speybroeck^{1, a)}
*Center for Molecular Modeling, Ghent University, Technologiepark 903,
B-9052 Zwijnaarde, Belgium*

CONTENTS

I. Free Energy Methods	SI2
A. Basic statistical physics	SI2
B. Umbrella sampling	SI3
1. Weighted histogram analysis method (WHAM)	SI3
C. Variationally enhanced sampling	SI4
D. Thermodynamic integration	SI6
E. Replica Exchange	SI7
1. Temperature Weighted Histogram Analysis Method (TWHAM)	SI7
II. Time-lagged independent component analysis (tICA)	SI8
III. Collective variable transformation	SI8
A. Explicit relation	SI8
1. Theory	SI8
2. Example	SI9
B. Implicit relation	SI9
1. Theory	SI9
IV. Free energy surface in terms of the bending angles	SI13
V. Further analysis inadequate collective variables	SI13

^{a)}Electronic mail: Veronique.VanSpeybroeck@UGent.be.

I. FREE ENERGY METHODS

A. Basic statistical physics

This section is devoted to the theoretical basis of the different free energy methods studied in this work. All methods find their basis in classical statistical physics. More in particular, the computation of the free energy can be related to the classical partition function and hence an integration over the entire phase space, i.e. the space spanned by the N atomic positions \mathbf{r}^N and the N atomic momenta \mathbf{p}^N . Hence, the total partition function Z and free energy F of the global state of a molecular system at a temperature T are defined as:

$$Z = \frac{1}{h^{3N}} \int e^{-\beta\mathcal{H}(\mathbf{r}^N, \mathbf{p}^N)} d\mathbf{r}^N d\mathbf{p}^N \quad (1)$$

$$F = -k_B T \ln Z \quad (2)$$

In this expression, h is the Planck constant, k_B is the Boltzmann constant and $\beta = 1/k_B T$. Herein, we assume the N atoms to be distinguishable, otherwise an additional $1/N!$ needs to be included in our normalization. The Hamiltonian $\mathcal{H}(\mathbf{r}^N, \mathbf{p}^N)$ contains the kinetic energy of all particles and the potential energy $\mathcal{U}(\mathbf{r}^N)$ of all particles.

Suppose we introduce a coordinate $Q(\mathbf{r}^N, \mathbf{p}^N)$ which is a function of all the degrees of freedom of the system. Furthermore, we are interested in the probability that the system is in a state for which $Q(\mathbf{r}^N, \mathbf{p}^N) = q$. In other words, this coordinate $Q(\mathbf{r}^N, \mathbf{p}^N)$ is a way of partitioning all the available microstates of the system (characterized by $\mathbf{r}^N, \mathbf{p}^N$) into a set of macrostates (characterized by q). This results in a partitioning of the partition function and free energy into contributions for every macrostate q :

$$Z(q) = \frac{q_0}{h^{3N}} \int e^{-\beta\mathcal{H}(\mathbf{r}^N, \mathbf{p}^N)} \delta(Q(\mathbf{r}^N, \mathbf{p}^N) - q) d\mathbf{r}^N d\mathbf{p}^N \quad (3)$$

$$F(q) = -k_B T \ln Z(q) \quad (4)$$

In the previous expression, a constant q_0 value is introduced to make $Z(q)$ dimensionless. The latter allows to determine the free energy as a function of q . In essence, the introduction of the constant q_0 results in a constant shift of $F(q)$. Since, we are interested in relative free energy differences, the specific value of q_0 is of no interest. Furthermore, we can also partition the total probability for the system to be in the global state, which is 1, into a probability density $p(q)dq$ for the system to be in a macrostate $[q, q + dq]$ with :

$$p(q) = \frac{1}{h^{3N}} \int p(\mathbf{r}^N, \mathbf{p}^N) \delta(Q(\mathbf{r}^N, \mathbf{p}^N) - q) d\mathbf{r}^N d\mathbf{p}^N \quad (5)$$

$$= \frac{1}{h^{3N} Z} \int e^{-\beta\mathcal{H}(\mathbf{r}^N, \mathbf{p}^N)} \delta(Q(\mathbf{r}^N, \mathbf{p}^N) - q) d\mathbf{r}^N d\mathbf{p}^N = \frac{Z(q)}{q_0 Z} \quad (6)$$

$$= \frac{1}{q_0 Z} e^{-\beta F(q)}$$

From the probability for each macrostate, the corresponding free energy can be determined:

$$F(q) = -k_B T \ln p(q) - k_B T \ln Z = -k_B T \ln p(q) + F \quad (7)$$

In practice, an integration over the entire phase space is unfeasible due to the large number of degrees of freedom. Hence, the partition function Z and corresponding free energy F cannot be computed. Nevertheless, the relative free energy difference between two macrostates can be computed by determination of the probability distribution: $F(q_2) - F(q_1) = -k_B T \ln [p(q_2)/p(q_1)]$. By means of molecular simulations, the relevant parts of the phase space are scanned allowing to determine the probability distribution in terms of the different macrostates. However, these molecular simulations are limited to short time scales, restricting the scan of the phase space to local minima. This results in a non-ergodic sampling if major free energy barriers are present. To overcome the limited sampling,

several enhanced sampling techniques have been proposed: replica exchange, umbrella sampling, variationally enhanced sampling, and thermodynamic integration. Replica exchange enhances the sampling of all degrees of freedom. The other sampling techniques enhance the sampling of the phase space in the direction of the partitioning coordinate, which in literature is often described as the collective variable or the reaction coordinate. Hence, those techniques will help to overcome free energy barriers in the direction of the collective variable, which are sampled poorly in classical molecular dynamics (MD) simulations.

B. Umbrella sampling

As stated before, sampling by regular molecular simulations does not sufficiently sample the regions with low Boltzmann probability, $\exp(-\beta\mathcal{H})$. In order to overcome this issue an external potential is introduced in the umbrella sampling method to enhance the sampling in these regions of low probability.¹ The external potential depends on the partitioning coordinate Q . Introducing this external potential in the simulation results in a biased partition function Z_b , free energy F_b , partitioned partition function $Z_b(q)$ and free energy profile $F_b(q)$:

$$Z_b = \frac{1}{h^{3N}} \int e^{-\beta(\mathcal{H}(\mathbf{r}^N, \mathbf{p}^N) + U_b(Q))} d\mathbf{r}^N d\mathbf{p}^N \quad (8)$$

$$F_b = -k_B T \ln Z_b \quad (9)$$

$$Z_b(q) = \frac{1}{h^{3N}} \int e^{-\beta(\mathcal{H}(\mathbf{r}^N, \mathbf{p}^N) + U_b(Q))} \delta(Q(\mathbf{r}^N, \mathbf{p}^N) - q) d\mathbf{r}^N d\mathbf{p}^N \quad (10)$$

$$F_b(q) = -k_B T \ln [Z_b(q)] \quad (11)$$

$$= -k_B T \ln p_b(q) + F_b \quad (12)$$

The free energy of the unbiased system $F(q)$ can be obtained from the free energy of the biased system $F_b(q)$:

$$F_b(q) = -k_B T \ln [Z_b(q)] \quad (13)$$

$$= -k_B T \ln [Z(q) e^{-\beta U_b(q)}] \quad (14)$$

$$= F(q) + U_b(q) \quad (15)$$

In practice, a set of external potentials $\{U_b^i\}$ is introduced. Here, each external potential focuses on a different region of the partitioning coordinate, such that the total set spans the entire region of interest. A popular choice for the external potential is a set of harmonic functions at uniformly distributed positions q_i over the sampling region: $U_b^i(q) = \frac{k}{2}(q - q_i)^2$. For every external potential U_b^i , a different simulation is performed yielding a set of histograms $H_i(q)$. One finds the biased free energy profiles from the histogram by introducing $p_b^i(q) dq = H_i(q)/M_i$ (with $H_i(q)$ the number of counts in the range $[q, q + dq]$ and M_i the total simulation points in simulation i) in Equation 12. In order to obtain a histogram for the unbiased system from the various biased histograms, different schemes can be considered, e.g. the weighted histogram analysis method (WHAM)^{2,3} and the dynamic histogram analysis method (DHAM)⁴.

1. Weighted histogram analysis method (WHAM)

In principle, the unbiased probability can be written as a function of each of the biased probabilities:

$$p(q) = p_b^i(q) e^{\beta U_b^i(q)} \frac{Z_b^i}{Z} \quad (16)$$

In practice, this will not work because, we can not compute Z_b^i directly. In the weighted histogram analysis method, a weight w_i is assigned to each count H_i in the following way.²

First off, the unbiased probability is written as a function of the weighted biased probabilities:

$$p(q) = \sum_i w_i p_b^i(q) e^{\beta U_b^i(q)} \frac{Z_b^i}{Z} \quad (17)$$

where the weights add up to 1, $\sum_i w_i = 1$. The expression for the unbiased probability contains no information about the weights w_i nor of the values of the ratios $\frac{Z_b^i}{Z}$.

Subsequently, an estimate of the weights is found by choosing them such that the variation of the unbiased probability is minimized. To that end, we assume the biased counts $H_i(q)$ to be distributed according to a Poisson distribution, such that $\text{var}(p_b^i) = \text{var}(H_i)/M_i^2 = E(H_i)/M_i^2 \approx p_b^i/M_i$.

$$\sigma_p^2 = \text{var}(p) = \langle p(q)^2 \rangle - \langle p(q) \rangle^2 \quad (18)$$

$$= \sum_i w_i^2 e^{2\beta U_b^i(q)} \left(\frac{Z_b^i}{Z} \right)^2 [\langle p_b^i(q)^2 \rangle - \langle p_b^i(q) \rangle^2] \quad (19)$$

$$= \sum_i w_i^2 e^{2\beta U_b^i(q)} \left(\frac{Z_b^i}{Z} \right)^2 \frac{p_b^i(q)}{M_i} \quad (20)$$

$$= p(q) \sum_i w_i^2 e^{\beta U_b^i(q)} \left(\frac{Z_b^i}{Z} \right) \frac{1}{M_i} \quad (21)$$

with M_i the number of sampling points obtained in simulation i with external potential U_b^i . Minimizing this expression for the weight functions yields:

$$w_i = \frac{e^{-\beta U_b^i(q)} M_i Z / Z_b^i}{\sum_i e^{-\beta U_b^i(q)} M_i Z / Z_b^i} \quad (22)$$

Substituting these values yields a new expression for the unknown probability density $p(q)$, with only the ratios $\frac{Z_b^i}{Z}$ as unknown values:

$$p(q) = \frac{\sum_i H_i(q)}{\sum_i e^{-\beta U_b^i(q)} M_i Z / Z_b^i} \quad (23)$$

with $H_i(q)$ the counts in bin $[q, q + dq]$ in the simulation with external potential $U_b^i(q)$.

Subsequently, an estimate for the ratios $\frac{Z_b^i}{Z}$ is found by introducing a self-consistent cycle based up the definition of the partition function and the previous expression of the unbiased probability:

$$Z_b^i = \int d\mathbf{r}^N d\mathbf{p}^N e^{-\beta(\mathcal{H} + U_b^i(q))} \quad (24)$$

$$= \int dq Z p(q) e^{-\beta U_b^i(q)} \quad (25)$$

$$= \int dq e^{-\beta U_b^i(q)} \frac{\sum_i H_i(q)}{\sum_i e^{-\beta U_b^i(q)} M_i / Z_b^i} \quad (26)$$

Solving this self-consistent cycle for the unknown partition functions removes the last unknown values in the expression for the unbiased probability, allowing for the determination of the free energy profile (within a constant value). More details on the derivation of the WHAM method can be found in reference 5. In our manuscript, we employ the WHAM script provided in reference 6.

C. Variationally enhanced sampling

In the previous section about umbrella sampling we derived a relation between the free energy of a macrostate $F(q)$, a bias potential $U_t(q)$, and the biased free energy $F_t(q)$ (Eq.

13):

$$F_t(q) = F(q) + U_t(q) \quad (27)$$

Furthermore in Equation (7), we derived the relationship between the probability and free energy of a macrostate. Applying this relation to the biased free energy $F_t(q)$ gives:

$$F_t(q) - F_t = -k_B T \ln p_t(q) \quad (28)$$

$$F_t(q) = -k_B T \ln p_t(q) + cte \quad (29)$$

Combining the Equations (27) and (29) results in :

$$U_t(q) = -F(q) - k_B T \ln p_t(q) + cte \quad (30)$$

This equation describes the relation (up to an irrelevant constant) between the unbiased free energy $F(q)$ of the system – the information we usually want to derive from a simulation – the biased probability $p_t(q)$, and the bias potential $U_t(q)$. In principle, this equation would allow us to determine the bias potential $U_t(q)$ associated with an a priori chosen target probability $p_t(q)$. However, that would require $F(q)$, which is the quantity we aim to calculate. In this enhanced metadynamics method, a functional of a bias potential U_b is introduced⁷:

$$\Omega[U_b] = \frac{1}{\beta} \ln \frac{\int e^{-\beta[F(q)+U_b(q)]} dq}{\int e^{-\beta F(q)} dq} + \int p_t(q) U_b(q) dq \quad (31)$$

This functional has three important properties:

- it has a stationary point at $U_b(q) = -F(q) - k_B T \ln p_t(q) = U_t(q)$
- it is a convex functional of U_b
- the functional value and its derivatives can be calculated from ensemble averages without the explicit value of the free energy:

$$\Omega[U_b] = \frac{1}{\beta} \ln \frac{\int e^{-\beta[F(q)+U_b(q)]} dq}{\int e^{-\beta[F(q)+U_b(q)]} e^{\beta U_b(q)} dq} + \int p_t(q) U_b(q) dq \quad (32)$$

$$= \frac{1}{\beta} \ln \frac{1}{\langle e^{\beta U_b(q)} \rangle_{U_b}} + \langle U_b(q) \rangle_{p_t} \quad (33)$$

$$= -\frac{1}{\beta} \ln \langle e^{\beta U_b(q)} \rangle_{U_b} + \langle U_b(q) \rangle_{p_t} \quad (34)$$

$\langle \cdot \rangle_{U_b}$ represents the ensemble average of the system biased with the potential U_b and can be calculated from a molecular simulation, while $\langle \cdot \rangle_{p_t}$ represents the average according to the target probability distribution $p_t(q)$ and can be calculated by means of numerical or analytical integration. The functional has a global minimum at the bias potential U_t that corresponds to the a priori chosen target probability p_t . After finding the minimum of this functional, one also finds through Equation (30) the free energy $F(q)$ of the unbiased system. The strength of this method is the fact that it also provides a well-defined procedure to construct an efficient bias potential through a variational principle.

A derivation of the functional is discussed in the work of Billionis et al.⁸, where the biasing of the dynamics and the estimation of the free energy profile are unified under the same objective of minimizing the Kullback-Leibler divergence between appropriately selected distributions on the collective variable space. More in particular, the Kullback-Leibler divergence between the target and bias probability,

$$\text{KL}(p_t|p_b) = \int dq p_t(q) \ln \left(\frac{p_t(q)}{p_b(q)} \right), \quad (35)$$

is minimized, where $\text{KL} = 0$ represent matching distributions.

To optimize the functional, the bias potential is expanded into a finite basis set G_k and the variational principle is applied by varying the expansion coefficients α_k until the stationary point Ω is found. Due to the convex nature of Ω , we know that this stationary point is also the global minimum.

$$U_b(q|\boldsymbol{\alpha}) = \sum_k \alpha_k G_k(q) \quad (36)$$

$$\Omega(\boldsymbol{\alpha}) = \frac{1}{\beta} \ln \frac{\int e^{-\beta[F(q)+\sum_k \alpha_k G_k(q)]} dq}{\int e^{-\beta F(q)} dq} + \int p_t(q) \sum_k \alpha_k G_k(q) dq \quad (37)$$

$$\frac{\partial \Omega}{\partial \alpha_k} = \frac{1}{\beta} \left(-\beta \frac{\int G_k(q) e^{-\beta[F(q)+\sum_k \alpha_k G_k(q)]} dq}{\int e^{-\beta[F(q)+\sum_k \alpha_k G_k(q)]} dq} \right) + \int p_t(q) G_k(q) dq \quad (38)$$

$$= -\langle G_k \rangle_{U_b} + \langle G_k \rangle_{p_t} \quad (39)$$

The stationary point is found when the gradient is zero, hence when the average of each basis function with respect to the bias potential U_b on the one hand and the target probability p_t on the other hand is equal. For efficient implementation of the minimizer, the second order derivatives are also required:

$$\frac{\partial^2 \Omega}{\partial \alpha_k \partial \alpha_l} = \beta (\langle G_k G_l \rangle_{U_b} - \langle G_k \rangle_{U_b} \langle G_l \rangle_{U_b}) \quad (40)$$

$$= \beta \text{Cov}[G_k, G_l]_{U_b} \quad (41)$$

The expansion coefficients can then be updated according to a stochastic gradient descent-based algorithm (due to the noise present in the MD simulation, one prefers to use a more robust update scheme than straightforward conjugate gradient method) with update parameter μ :

$$\boldsymbol{\alpha}^{(n+1)} = \boldsymbol{\alpha}^{(n)} - \mu \left[\frac{\partial \Omega}{\partial \boldsymbol{\alpha}} (\bar{\boldsymbol{\alpha}}^{(n)}) + \frac{\partial^2 \Omega}{\partial \boldsymbol{\alpha}^2} (\bar{\boldsymbol{\alpha}}^{(n)}) (\boldsymbol{\alpha}^{(n)} - \bar{\boldsymbol{\alpha}}^{(n)}) \right] \quad (42)$$

Hence, a first order Taylor expansion of the gradient is constructed around the cumulative moving average $\bar{\boldsymbol{\alpha}}^{(n)} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\alpha}^{(i)}$ instead of evaluating the gradient directly at the instantaneous $\boldsymbol{\alpha}^{(i)}$ to minimize the influence of the statistical fluctuations.

In the original work by Valsson and Parrinello exponential basis functions were proposed⁷. In this work, we opt to use Gaussian basis functions instead of the exponential functions.

D. Thermodynamic integration

Suppose the Hamiltonian depends on a partitioning parameter λ : $\mathcal{H}(\mathbf{r}^N, \mathbf{p}^N; \lambda)$ ^{9,10}. The free energy of the system and its derivative with respect to this parameter are given by

$$F(\lambda) = -k_B T \ln \left[C \int e^{-\beta \mathcal{H}(\mathbf{r}^N, \mathbf{p}^N; \lambda)} d\mathbf{r}^N d\mathbf{p}^N \right] \quad (43)$$

$$\frac{\partial F}{\partial \lambda} = -k_B T \frac{C \int -\beta \frac{\partial \mathcal{H}}{\partial \lambda} e^{-\beta \mathcal{H}} d\mathbf{r}^N d\mathbf{p}^N}{C \int e^{-\beta \mathcal{H}} d\mathbf{r}^N d\mathbf{p}^N} \quad (44)$$

$$= \left\langle \frac{\partial \mathcal{H}}{\partial \lambda} \right\rangle_\lambda \quad (45)$$

This results in the following expression for the free energy difference of the system between states λ_2 and λ_1 :

$$F(\lambda_2) - F(\lambda_1) = \int_{\lambda_1}^{\lambda_2} \frac{\partial F}{\partial \lambda} d\lambda \quad (46)$$

$$= \int_{\lambda_1}^{\lambda_2} \left\langle \frac{\partial \mathcal{H}}{\partial \lambda} \right\rangle_\lambda d\lambda \quad (47)$$

The subscript λ indicates that the average has to be computed in the ensemble where the Hamiltonian has a parameter value of λ . Hence, the ensemble changes during the integration, so that multiple ensembles have to be sampled for values of λ in between λ_1 and λ_2 .

E. Replica Exchange

During a molecular dynamics simulation, the sampled phase space is limited to local minima. More in particular, only free energy barriers of the order of $k_B T$ can be surmounted. A valid strategy to extend the sampled phase space to perform a simulation at a higher temperature. This idea is at the basis of replica exchange¹¹. Rather than performing one simulation under certain thermodynamic conditions, multiple simulation are run in parallel at various thermodynamic conditions. Subsequently the sampled phase space is extended by swapping the configuration sampled under those different thermodynamic conditions. To circumvent a non-canonical sampling of the phase space this swapping move is governed by following acceptance rule:

$$P_{acc} = \min \left[1, \exp \left((\beta_n - \beta_m)(\mathcal{U}(\mathbf{r}_i^N) - \mathcal{U}(\mathbf{r}_j^N)) + (\beta_n P_n - \beta_m P_m)(V_i - V_j) \right) \right] \quad (48)$$

where $\{\mathbf{r}_i^N, V_i\}$ and $\{\mathbf{r}_j^N, V_j\}$ represent two different configurations stemming from an MD simulation at two distinct thermodynamic conditions n and m , respectively. In this particular case, the thermodynamic condition depends on a temperature and pressure. In practice, a Monte Carlo simulation is performed to govern the swaps of replicas between the different simulations. Moreover, to suppress the instability due to the swap, velocities are rescaled according to the temperature ratio. In essence, replica exchange is a hybrid scheme of a MD/MC simulation in order to extend the sampled region of the phase space under certain thermodynamic conditions.

Subsequently, a free energy profile can be estimated employing Eq. 7. However, statistical more optimal methods, e.g. temperature weighted histogram analysis method (TWHAM)¹², the multistate Bennett acceptance ratio (MBAR) method¹³ or the transition based reweighting analysis method (TRAM)¹⁴, employ all available data to construct a free energy profile. Hence, also data points corresponding to thermodynamic conditions different from the one of interest are employed for the construction of a free energy profile.

1. Temperature Weighted Histogram Analysis Method (TWHAM)

In this work, we implement and employ the TWHAM method to analyse the data.¹² The proof of WHAM is outlined in section IB for the case of umbrella sampling. TWHAM is based on the same principle, where a global free energy profile is constructed by reweighting local histograms of 'biased' simulations. Rather than various bias potentials, RE simulations are performed at different temperatures $\{T_i\}_{i=1, \dots, N}$. Those various thermodynamic conditions yield a biasing factor, $e^{-(\beta_i - \beta_0)\mathcal{U}}$, which depends on the energy \mathcal{U} of the system and the reference temperature, characterized by the boltzmann β_0 . Hence, the unnormalized probability distribution in terms of both the collective variable and energy can be obtained by solving following equations self-consistently:

$$p(q, E) = \frac{\sum_i H_i(q, E)}{\sum_i e^{-(\beta_i - \beta_0)\mathcal{U}} M_i / Z_b^i} \quad (49)$$

$$Z_b^i = \int \int dq dE e^{-(\beta_i - \beta_0)\mathcal{U}} p(q, E) \quad (50)$$

Iterating those equations until self-consistency is reached yields the bivariate distribution in terms of energy and collective variable. By integrating out the former, a free energy profile in terms of the collective variable q is achieved.

II. TIME-LAGGED INDEPENDENT COMPONENT ANALYSIS (TICA)

The time-lagged independent component analysis (tICA) method is employed for the identification of slow modes in replica exchange data. For the ease of the reader, we summarize the theory of this method in this section. For more details on the theory of tICA, we refer to references 15,16.

Before going into detail on the method, we introduce some notations. The order parameter space has d components $\mathbf{q} = \{q_i(\mathbf{x})\}_{i=1,\dots,d}$ which are a function of the configuration space \mathbf{x} . The independent component space has the same dimensions $\mathbf{z} = \{z_i\}_{i=1,\dots,d}$. The linear transformation from the order parameter space towards the independent component space is denoted with in matrix notation: $\mathbf{z} = \mathbf{U}\mathbf{q}$. The time correlation function of two order parameters is defined as $\mathbf{C}^{\mathbf{q}}_{ij}(\tau) = \langle q_i(t)q_j(t+\tau) \rangle_t$. Furthermore, a time correlation function can be constructed in terms of two independent components, which is defined as $\mathbf{C}^{\mathbf{z}}_{ij} = \langle z_i(t)z_j(t+\tau) \rangle_t$.

In tICA, one searches for the linear transformation \mathbf{U} which maximizes the time correlation function $\mathbf{C}^{\mathbf{z}}(\tau)$ at a fixed lag time τ . Furthermore, it is imposed that the covariance matrix $\mathbf{C}^{\mathbf{z}}(\mathbf{0})$ equals the unity matrix. Introducing the Lagrange multipliers $\boldsymbol{\lambda}$, the Lagrangian (or objective function) can be constructed:

$$\mathbf{F} = \mathbf{C}^{\mathbf{z}}(\tau) - (\mathbf{C}^{\mathbf{z}}(\mathbf{0}) - \mathbf{1})\boldsymbol{\lambda} \quad (51)$$

$$= \mathbf{U}\mathbf{C}^{\mathbf{q}}(\tau)\mathbf{U}^T - (\mathbf{U}\mathbf{C}^{\mathbf{q}}(\mathbf{0})\mathbf{U}^T - \mathbf{1})\boldsymbol{\lambda} \quad (52)$$

Subsequently, by imposing that the derivative of the Lagrangian \mathbf{F} in terms of the linear transformation \mathbf{U} equals zero, one finds a generalized eigenvalue equation of the matrix time correlation matrix:

$$\mathbf{C}^{\mathbf{q}}\mathbf{U}^T = \mathbf{C}^{\mathbf{q}}(\mathbf{0})\mathbf{U}^T\boldsymbol{\lambda} \quad (53)$$

Therein, the Lagrange multipliers equal the eigenvalues of the generalized eigenvalue equation.

III. COLLECTIVE VARIABLE TRANSFORMATION

A. Explicit relation

1. Theory

In this section, we discuss the transformation of collective variables for free energy profiles. In essence, changing the collective variable is the transformation of random variable of a probability distribution. General probability theory learns us that in this case one should take into account the Jacobian to deal with the changing metrics. An easy-to-understand example of such a transformation, is where one is interested in the free energy profile as function of the coordination number $F_s(s)$, given the free energy profile as function of the distance $F_r(r)$. Transforming a coordination number s to a distance r is governed by a (reversible) switching function $r \rightarrow s = f(r)$. Using the relation between probability and free energy, we apply the transformation of random variable:

$$F_r(r) = -k_B T \ln \left[\exp(-\beta F_s(f(r))) \left| \frac{df(r)}{dr} \right| \right] \quad (54)$$

We can extend this transformation to free energy surfaces in multiple dimensions. Suppose we want to transform a free energy surface F_r as function of one set of collective variables (s_1, s_2) to a free energy surface F_s as function of another set of collective variables (r_1, r_2) . Moreover, a transformation exists, which can be inverted, $(r_1, r_2) \rightarrow (s_1, s_2) = (f_1(r_1, r_2), f_2(r_1, r_2))$.

$$F_r(r_1, r_2) = -k_B T \ln [\exp(-\beta F_s(f_1(r_1, r_2), f_2(r_1, r_2))) |J|] \quad (55)$$

with J the Jacobian:

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial r_1} & \frac{\partial f_1}{\partial r_2} \\ \frac{\partial f_2}{\partial r_1} & \frac{\partial f_2}{\partial r_2} \end{bmatrix} \quad (56)$$

2. Example

The latter transformation is only applicable for those cases where a well defined function between the sets of collective variables exists. In the case of a breathing MIL-53(AI), a well defined transition of variables is from the unit cell parameters (a_x, c_z) to the unit cell diagonal D and unit bend cell angle θ :

$$D = g_1(a_x, c_z) = \sqrt{a_x^2 + c_z^2} \quad (57)$$

$$\theta = g_2(a_x, c_z) = \arctan\left(\frac{c_z}{a_x}\right) \quad (58)$$

So suppose, we have a free energy profile in terms of the cell parameters a_x and c_z and we are interested in a free energy profile in terms of the diagonal and angle, we compute the inverse functions and the Jacobian:

$$a_x = f_1(D, \theta) = D \cos(\theta) \quad (59)$$

$$c_z = f_2(D, \theta) = D \sin(\theta) \quad (60)$$

$$|J| = D \quad (61)$$

Hence, we determine the free energy surface in terms of the diagonal and angle using the transformation in eq. 55. The result of such a transformation is provided in Figure S11. Moreover, Figure S11 emphasizes the construction of one dimensional free energy profiles by integrating out one of the variables.

B. Implicit relation

1. Theory

More often one is interested in the transformation of a free energy profile for which no exact function exists. Suppose, we have knowledge of a free energy profile as a function of a first collective variable \mathbf{q}_1 , however, we are interested in a free energy profile as a function of another collective variable \mathbf{q}_2 . In this case, there is no unique relation between these collective variables, we can employ following transformation:

$$F(\mathbf{q}_2) = -k_B T \ln \left[\frac{1}{C} \int_{\mathbf{q}_1} p(\mathbf{q}_2|\mathbf{q}_1) e^{-\beta F(\mathbf{q}_1)} d\mathbf{q}_1 \right] \quad (62)$$

Herein, the conditional probability $p(\mathbf{q}_2|\mathbf{q}_1)$ can stem from an unbiased, constrained, or a biased simulation. Before proving this statement, we introduce some notations. In Eq. 62, the constant C represents the necessary normalization, for the ease of notation. Since, this constant represents a mere vertical shift of the free energy profile, we will drop the constant.

We use following shorthand notations for the phase space: the entire phase space Γ , $\Gamma_{\mathbf{q}_i}$ the phase space associated with constant \mathbf{q}_i points and $\Gamma_{\mathbf{q}_1, \mathbf{q}_2}$ for the phase space associated with constant collective variables \mathbf{q}_1 and \mathbf{q}_2 pair. In essence, the free energy in terms of the second collective variable equals:

$$F(\mathbf{q}_2) = -k_B T \ln \left[\int_{\Gamma_{\mathbf{q}_2}} e^{-\beta H} d\Gamma_{\mathbf{q}_2} \right] \quad (63)$$

The integration over all degrees of freedom other than the second collective variable \mathbf{q}_2 , can be split into two contributions, integrating over the phase space $\Gamma_{\mathbf{q}_1, \mathbf{q}_2}$ and integrating

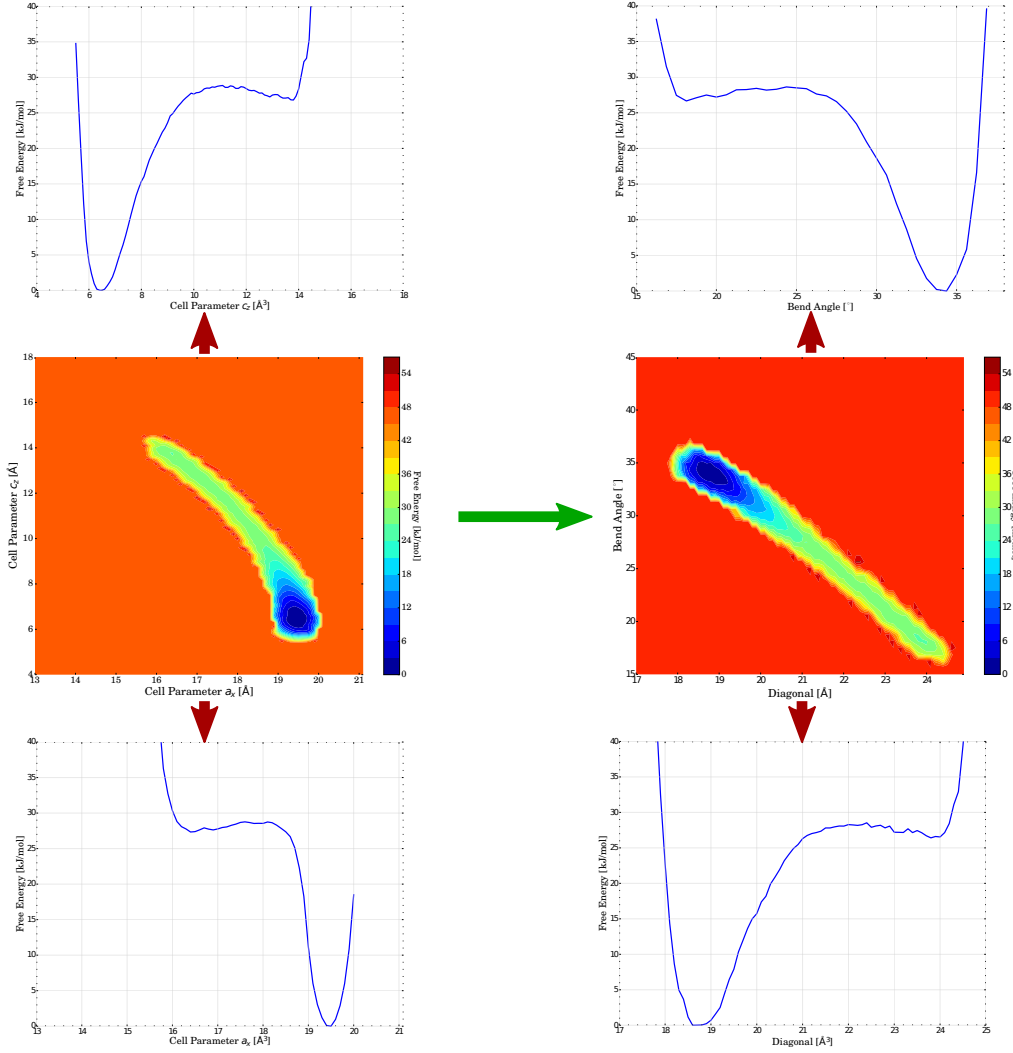


FIG. S11: Using a simple transformation, we change collective variable of a free energy surface (green arrow). Moreover, we construct one dimensional free energy profiles by integrating out one of the variables (red arrows).

over the first collective variable \mathbf{q}_1 . The former can be associated with the multivariate probability $p(\mathbf{q}_1, \mathbf{q}_2)$, which can be further reduced to $p(\mathbf{q}_2|\mathbf{q}_1)p(\mathbf{q}_1)$. Therefore, we find the transformation proposed in Eq. 62.

$$F(\mathbf{q}_2) = -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 \int_{\Gamma_{\mathbf{q}_1, \mathbf{q}_2}} e^{-\beta H} d\Gamma_{\mathbf{q}_1, \mathbf{q}_2} \right] \quad (64)$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} p(\mathbf{q}_1, \mathbf{q}_2) d\mathbf{q}_1 \right] \quad (65)$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} p(\mathbf{q}_2|\mathbf{q}_1)p(\mathbf{q}_1) d\mathbf{q}_1 \right] \quad (66)$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} p(\mathbf{q}_2|\mathbf{q}_1) e^{-\beta F(\mathbf{q}_1)} d\mathbf{q}_1 \right] \quad (67)$$

However, often the construction of a free energy profile is associated with enhanced sampling simulations. Therefore, both the free energy profile $F(\mathbf{q}_1)$ and the conditional probability

$p(\mathbf{q}_2|\mathbf{q}_1)$ needs to be determined from those enhanced sampling simulations. We distinguish between two cases: the enhanced sampling simulation corresponds to constraining in terms of the first collective variable \mathbf{q}_1 or it corresponds to biasing in terms of the first collective variable.

Constrained MD. In the case of a constrained MD simulation, the integration over the phase space $\Gamma_{\mathbf{q}_2}$ in Eq. 63 is again split in two parts, i.e. first integrating over the $\Gamma_{\mathbf{q}_1, \mathbf{q}_2}$ space and subsequently integrating over the \mathbf{q}_1 coordinate.

$$F(\mathbf{q}_2) = -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 \int_{\Gamma_{\mathbf{q}_1, \mathbf{q}_2}} e^{-\beta H} d\Gamma_{\mathbf{q}_1, \mathbf{q}_2} \right] \quad (68)$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 \frac{\int_{\Gamma_{\mathbf{q}_1, \mathbf{q}_2}} e^{-\beta H} d\Gamma_{\mathbf{q}_1, \mathbf{q}_2}}{\int_{\Gamma_{\mathbf{q}_1}} e^{-\beta H} d\Gamma_{\mathbf{q}_1}} \int_{\Gamma_{\mathbf{q}_1}} e^{-\beta H} d\Gamma_{\mathbf{q}_1} \right] \quad (69)$$

Hence, the integration runs over two contributions, the first contribution equals $\langle p(\mathbf{q}_2) \rangle_{\mathbf{q}_1}$ and the second corresponds to the probability in terms of $p(\mathbf{q}_1)$. Therefore, we obtain the relation as proposed in 62:

$$F(\mathbf{q}_2) = -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 \langle p(\mathbf{q}_2) \rangle_{\mathbf{q}_1} e^{-\beta F(\mathbf{q}_1)} \right] \quad (70)$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 p_r(\mathbf{q}_2|\mathbf{q}_1) e^{-\beta F(\mathbf{q}_1)} \right] \quad (71)$$

We stress that the conditional probability stems from constrained MD simulation using the notation p_r .

Biased MD. Finally, we show that the proposed relation also holds in the case the conditional probability stems from a biased MD simulation. If one biases the simulation along the direction of the first collective variable \mathbf{q}_1 , one finds a biased free energy profile along this collective variable:

$$F_b(\mathbf{q}_1) = -k_B T \ln \left[\int_{\Gamma_{\mathbf{q}_1}} e^{-\beta(H+U_b(\mathbf{q}_1))} d\Gamma_{\mathbf{q}_1} \right] = F(\mathbf{q}_1) + U_b(\mathbf{q}_1) \quad (72)$$

Herein, the bias potential $U_b(\mathbf{q}_1)$ is only a function \mathbf{q}_1 and can thus be extracted from integration.

For such a biased simulation, we can proof that Eq. 62 holds in the case of such a biased MD simulation, starting from the definition of the free energy as a function of the the second collective variable \mathbf{q}_2 and splitting up the integration.

$$F(\mathbf{q}_2) = -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 \int_{\Gamma_{\mathbf{q}_1, \mathbf{q}_2}} e^{-\beta H} d\Gamma_{\mathbf{q}_1, \mathbf{q}_2} \right] \quad (73)$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 e^{\beta U_b(\mathbf{q}_1)} \int_{\Gamma_{\mathbf{q}_1, \mathbf{q}_2}} e^{-\beta(H+U_b(\mathbf{q}_1))} d\Gamma_{\mathbf{q}_1, \mathbf{q}_2} \right] \quad (74)$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 e^{\beta U_b(\mathbf{q}_1)} \frac{\int_{\Gamma_{\mathbf{q}_1, \mathbf{q}_2}} e^{-\beta(H+U_b(\mathbf{q}_1))} d\Gamma_{\mathbf{q}_1, \mathbf{q}_2}}{\int_{\Gamma_{\mathbf{q}_1}} e^{-\beta(H+U_b(\mathbf{q}_1))} d\Gamma_{\mathbf{q}_1}} \int_{\Gamma_{\mathbf{q}_1}} e^{-\beta(H+U_b(\mathbf{q}_1))} d\Gamma_{\mathbf{q}_1} \right]$$

$$= -k_B T \ln \left[\int_{\mathbf{q}_1} d\mathbf{q}_1 \frac{\int_{\Gamma_{\mathbf{q}_1, \mathbf{q}_2}} e^{-\beta(H+U_b(\mathbf{q}_1))} d\Gamma_{\mathbf{q}_1, \mathbf{q}_2}}{\int_{\Gamma_{\mathbf{q}_1}} e^{-\beta(H+U_b(\mathbf{q}_1))} d\Gamma_{\mathbf{q}_1}} \int_{\Gamma_{\mathbf{q}_1}} e^{-\beta H} d\Gamma_{\mathbf{q}_1} \right] \quad (75)$$

Hence, we find an integration over two contributions. The last contribution equals the factor $e^{-\beta F(\mathbf{q}_1)}$ (see Eq. 72) and the first contribution equals the conditional probability $p_b(\mathbf{q}_2|\mathbf{q}_1)$. The notation p_b stresses that this probability stems from biased conditions. A special case of biasing occurs when the bias potential exactly compensates the underlying free energy profile, i.e. $U_b(\mathbf{q}_1) = -F(\mathbf{q}_1)$. In that case, the denominator of the first contribution

becomes a constant (i.e. the silently assumed normalization constant). In general, we obtain the proposed transformation in Eq. 62:

$$F(\mathbf{q}_2) = -k_B T \ln \left[\int_{\mathbf{q}_1} p_b(\mathbf{q}_2|\mathbf{q}_1) e^{-\beta F(\mathbf{q}_1)} d\mathbf{q}_1 \right] \quad (76)$$

IV. FREE ENERGY SURFACE IN TERMS OF THE BENDING ANGLES

For CAU-13, we showcase the free energy profiles as function of the averaged bending angle of the CDC linkers in the main text. Since these profiles stem from a projection of the 2D free energy profiles, we include those profiles in this section.

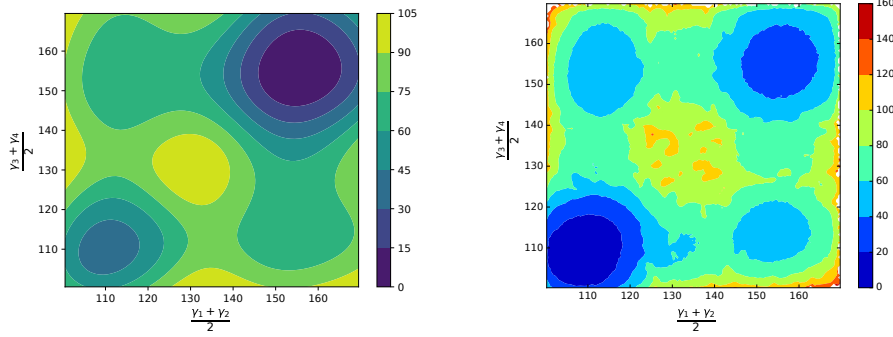


FIG. SI2: Free energy surface for CAU-13 at 300 as a function of the bending angles γ_i obtained with (a) VES and (b) US.

V. FURTHER ANALYSIS INADEQUATE COLLECTIVE VARIABLES

In the main text, the breathing behavior of CAU-13 is studied by means of enhanced sampling simulations employing three different collective variables, i.e. the unit cell volume, bending angles, and dihedral angles. Based on the tICA-RE protocol, the latter are proposed as the slowest varying order parameters and thus the most adequate set of collective variables. In the main text, no qualitative agreement of an enhanced sampling simulation in the direction of the inadequate collective variables with the replica exchange profile is observed, while such an qualitative agreement with for the dihedral angle is present. Hence, a bad choice of biasing coordinate yields an inaccurate free energy profile for techniques such as TI, US, and VES. To further convince the reader of this statement, a further analysis is shown in this section.

To this end, the mean value of volume and average bending angle at each two dimensional dihedral angle space point is shown in Figure SI3. Based on this contour plot, it is clearly visible that the unit cell volume does not characterizes the various stable states. The unit cell volume is not able to distinguish between the two pathways neither allows it to pinpoint whether a stable state or transition state is reached. A fundamental different observation is valid for the bending angle. The average bending angle clearly follows the proposed folding pathways from the cp to lp state, with contour lines orthogonal to the expected transition path. However, in the region of large dihedral angles ($\frac{\phi_1 - \phi_2}{2} \approx 60$, $\frac{\phi_3 - \phi_4}{2} \approx 60$), no uniform increase in the bending angle is observed. Rather than further increasing, the bending angle starts to decrease in the lp state. This leads to an inadequate description of the transformation of the lp state towards the intermediate stable states and thus for the lp-to-cp transformation. In practice, we observe that biased enhanced sampling schemes get stuck in the lp state. This issue is shown in Figure SI4 for a VES simulation.

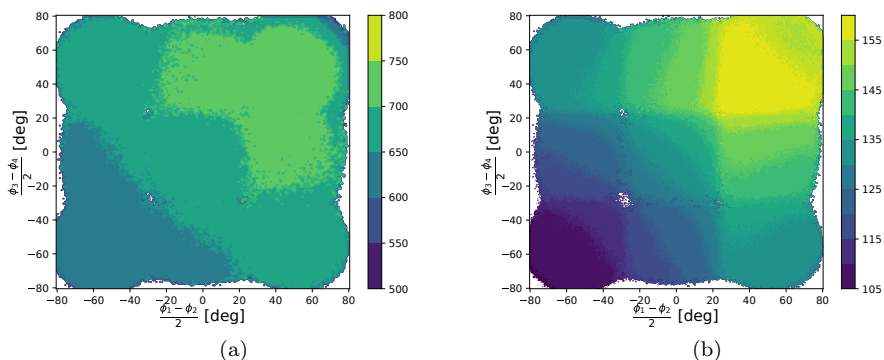


FIG. SI3: Average volume (left in \AA^3) and bending angle (right in degrees) of CAU-13 in terms of the two dimensional collective variable proposed by the tICA-RE protocol. Furthermore, the contours of the free energy surface in terms of the two dimensional collective variable is shown underneath.

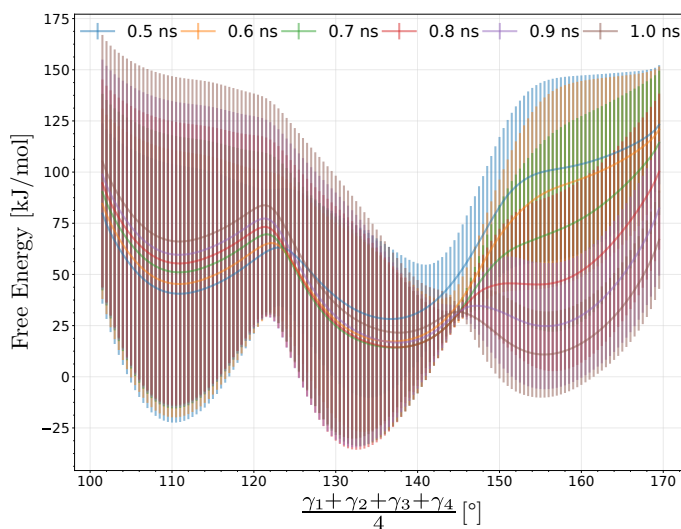


FIG. SI4: One dimensional projection of the free energy surface of CAU-13 constructed in terms of the two dimensional bending angles, $\frac{\gamma_1 + \gamma_2}{2}$ and $\frac{\gamma_3 + \gamma_4}{2}$.

- ¹G. Torrie and J. Valleau, J. Comput. Phys. **23**, 187 (1977).
- ²A. Ferrenberg and R. Swendsen, Phys. Rev. Lett. **61**, 2635 (1988).
- ³S. Kumar, J. M. Rosenberg, D. Bouzida, R. H. Swendsen, and P. A. Kollman, J. Comput. Chem. **13**, 1011 (1992).
- ⁴E. Rosta and G. Hummer, J. Chem. Theory Comput. **11**, 276 (2015).
- ⁵D. Frenkel and B. Smit, *Understanding Molecular Simulation: From Algorithms to Applications*, Vol. 1 (Academic press, 2001).
- ⁶A. Grossfield, “wham: The weighted histogram analysis method,” <http://membrane.urmc.rochester.edu/content/wham>.
- ⁷O. Valsson and M. Parrinello, Phys. Rev. Lett. **113**, 090601 (2014).
- ⁸I. Bilonis and P. S. Koutsourelakis, J. Comput. Phys. **231**, 3849 (2012).
- ⁹J. G. Kirkwood, J. Chem. Phys. **3**, 300 (1935).
- ¹⁰T. P. Straatsma and H. J. C. Berendsen, T.J. Chem. Phys. **89**, 5876 (1988).
- ¹¹Y. Sugita and Y. Okamoto, Chem. Phys. Lett. **314**, 141 (1999).
- ¹²E. Gallicchio, M. Andrec, A. K. Felts, and R. M. Levy, J. Phys. Chem. B **109**, 6722 (2005).
- ¹³M. R. Shirts and J. D. Chodera, J. Chem. Phys. **129**, 124105 (2008).
- ¹⁴H. Wu, F. Paul, C. Wehmeyer, and F. Noé, Proc. Natl. Acad. Sci. **113**, E3221 (2016).
- ¹⁵G. Pérez-Hernández, F. Paul, T. Giorgino, G. De Fabritiis, and F. Noé, J. Chem. Phys. **139**, 015102

(2013).
¹⁶C. R. Schwantes and V. S. Pande, *J. Chem. Theory Comput.* **9**, 2000 (2013).